A.KIEŁBASIŃSKI H.SCHWETLICK INTERIORATION INTERIORITIES INTERIORITARI INTERIORI INTERIORITARI INTERIORITARI INTERIORITARI INTERIORITARI INTERIORITARI INTERIORI INTERIORITARI INTERIORI IN

VEB DEUTSCHER VERLAG DER WISSENSCHAFTEN

Mathematik für Naturwissenschaft und Technik

Herausgegeben von H. Heinrich und H. Schubert

Band 18

Numerische lineare Algebra

Eine computerorientierte Einführung

von A. Kiełbasiński und H. Schwetlick

Mit 19 Abbildungen



VEB Deutscher Verlag der Wissenschaften Berlin 1988

ISBN 3-326-00194-0 ISSN 0543-100X

Verlagslektor: Erika Arndt Verlagshersteller: Norma Braun Umschlaggestaltung: Hartwig Hoeftmann © 1988 VEB Deutscher Verlag der Wissenschaften, DDR-1080 Berlin, Postfach 1216 Lizenz-Nr. 206 · 435/127/88 Printed in the German Democratic Republic Gesamtherstellung: VEB Druckhaus "Maxim Gorki", DDR - 7400 Altenburg LSV 1024 Bestellnummer: 571 555 2 03980

Vorwort

Die numerische lineare Algebra nimmt zweifellos eine zentrale Stellung innerhalb der Numerischen Mathematik ein. Von der Theorie her ist sie eines der am besten aufgearbeiteten Gebiete. Es gibt nur wenige Anwendungen der Mathematik, wo sie nicht direkt oder mittelbar angesprochen wird. Ihre Geschichte läßt sich über mehr als 2000 Jahre zurückverfolgen, und seit Beginn des Computerzeitalters ist ihre Entwicklung stürmischer als je zuvor vorangetrieben worden.

Der Zwang, Aufgaben der linearen Algebra auf verfügbaren Computern schnell, zuverlässig und in großer Anzahl lösen zu müssen, hat frühzeitig zur Auswahl, Entwicklung und Analyse von dafür geeigneten Algorithmen und zu deren sachgemäßer Implementierung - meist in Form von Programmpaketen - geführt. Es zeigte sich, daß die bei Computerrechnung mit endlicher Stellenzahl notwendig auftretenden Rundungsfehler einen wesentlichen Einfluß auf die berechneten Resultate haben können. Die Untersuchung des Verhaltens von Algorithmen bei Realisierung in Computerarithmetik wurde daher zu einem zentralen Gegenstand der numerischen linearen Algebra. Entscheidende Fortschritte brachte um die sechziger Jahre die im wesentlichen von WILKINSON entwickelte und perfekt angewendete Methode der "backward error analysis", die eine theoretische Erfassung der Rundungsfehler und daher gesicherte Aussagen über das Stabilitätsverhalten von Algorithmen erlaubt. In seinem Buch "The Algebraic Eigenvalue Problem" - gemeinhin "Bibel" der numerischen linearen Algebra genannt - hat WILKINSON 1965 eine auch heute noch gültige Zusammenfassung der wesentlichen Ergebnisse gegeben und damit die spätere Entwicklung entscheidend beeinflußt.

In den darauffolgenden Jahren entstanden weiterführende Darstellungen zu Teilgebieten der linearen Algebra, etwa LAWSON/HANSONS "Solving Least Squares Problems" aus dem Jahre 1974 und PARLETTS "The Symmetric Eigenvalue Problem" aus dem Jahre 1980. Alle diese anspruchsvollen Monographien erschließen sich jedoch dem Lernenden nicht leicht; sie sprechen mehr den erfahrenen, bereits mit speziellen Kenntnissen der numerischen linearen Algebra vertrauten Leser an. Obwohl die grundlegenden Ideen in entsprechend vereinfachter Form bereits Eingang in neuere einführende Gesamtdarstellungen der Numerischen Mathematik wie etwa die von STOER oder MAESS gefunden haben, fehlten Lehrbücher, die einen breiten Interessentenkreis systematisch zu einem tieferen Verständnis dieser neuen Entwicklung hätten führen können. Dieser Mangel wurde deutlicher, als exzellente Implementierungen der wichtigsten Algorithmen weltweit verfügbar wurden; wir erwähnen das "Handbook for Automatic Computation. Vol. II: Linear Algebra" aus dem Jahre 6

1971, die 1974 bzw. 1977 erschienenen beiden EISPACK-Versionen und das 1979 publizierte LINPACK.

Das vorliegende Lehrbuch ist — nach GOLUB/VAN LOANS "Matrix Computations" aus dem Jahre 1983 und der gleichnamigen "Numerischen linearen Algebra" von BUNSE/BUNSE-GERSTNER aus dem Jahre 1985 — nach Wissen der Autoren das dritte, das die erwähnte Lücke zwischen den anspruchsvollen Monographien und den einführenden Darstellungen zu schließen versucht. Es wendet sich an Mathematiker, aber auch an mathematisch interessierte Naturwissenschaftler und Ingenieure, die sich intensiver mit numerischer linearer Algebra vertraut machen wollen und nicht nur an einer kurzen, übersichtsartigen Zusammenfassung interessiert sind.

Um dieser Zielstellung gerecht zu werden, haben wir uns auf die Behandlung der direkten Verfahren für lineare Gleichungssysteme und Quadratmittelprobleme und die des Eigenwertproblems symmetrischer Matrizen beschränkt. Diese Gebiete sind jedoch vergleichsweise gründlich unter Einbeziehung verschiedenster Gesichtspunkte bis hin zu Fragen der Implementierung erörtert worden. Es wurde eine weitgehend abgeschlossene Darstellung gewählt, die im wesentlichen nur die aus der mathematischen Grundausbildung der Ingenieure bekannten Vorkenntnisse voraussetzt. Um mögliche Unterschiede in diesen Kenntnissen auszugleichen, sind die später benötigten Begriffe und Fakten in einem einführenden ersten Kapitel nochmals übersichtsartig zusammengestellt worden. Im übrigen Text wurden fast alle Beweise im Detail ausgeführt - anfangs ausführlicher, zum Schluß dann knapper -, so daß das Buch auch zum Selbststudium geeignet sein dürfte. In dieser Hinsicht unterscheidet sich die vorliegende Darstellung von der GOLUB/VAN LOANS, wo in den Beweisen zumeist auf die Originalliteratur oder andere Monographien verwiesen wird; dafür ist allerdings dort ein wesentlich umfangreicheres Gebiet behandelt worden. Verzichtet wurde lediglich auf die Herleitung von gewissen Ergebnissen der Rundungsfehleranalyse, wenn der Aufwand dafür ein vertretbares Maß überschritten hätte. Viel Wert wurde auf präzise Begriffsbildungen, die Herausarbeitung der algorithmischen Grundideen und - stärker als etwa bei BUNSE/BUNSE-GERSTNER auf die Untersuchung des Einflusses der Rundungsfehler und besonders auf die sachgemäße Interpretation der Ergebnisse der Rundungsfehleranalyse gelegt. Nach den Erfahrungen der Autoren macht das letztgenannte Problem nicht nur Anfängern Schwierigkeiten, so daß für den Komplex "Aufgaben, Computer, Algorithmen" ein separates Kapitel vorgesehen wurde, in dem die Grundbegriffe der Fehleranalyse und verwandte Begriffe in systematischer Weise entwickelt werden. Der den Grundlagen gewidmete erste Teil schließt mit einem Kapitel über elementare Transformationsmatrizen. Die weiteren drei Teile haben die bereits genannten Problemklassen zum Gegenstand. Sie beginnen jeweils mit der Darstellung der wesentlichen theoretischen Eigenschaften, wobei insbesondere auf Ergebnisse der Störungstheorie und auf Residualkriterien eingegangen wird. In den algorithmisch orientierten Abschnitten wird zunächst der "mathematische" Algorithmus motiviert und formuliert. Es folgen Bemerkungen zur Implementierung, die Formulierung des "Computeralgorithmus" in einer sich selbst erklärenden, einfachen Pseudo-Programmiersprache, Aufwandsbetrachtungen, Rundungsfehleranalyse und deren Interpretation sowie ergänzende und weiterführende Bemerkungen. Da die Algorithmen zum Schluß

immer komplexer werden und eine detaillierte computernahe Formulierung mit akzeptablem Aufwand kaum noch möglich ist, wird ab Kapitel 11 auf die gesonderte Formulierung des Computeralgorithmus verzichtet.

Jedem Kapitel sind Bemerkungen nachgestellt, in denen die historische Entwicklung unter Angabe grundlegender Arbeiten skizziert wird und Hinweise auf weiterführende Resultate zu finden sind. Im Text selbst wird bis auf wenige Ausnahmen keine Literatur zitiert. Jeder der den drei behandelten Aufgabenklassen gewidmeten Teile schließt mit einer kurzen Zusammenfassung, in der die vorgestellten Algorithmen verglichen und Empfehlungen zu ihrer Anwendung gegeben werden. Außerdem wird auch auf im Text nicht behandelte Aufgaben- und Algorithmenklassen eingegangen. Die ins Literaturverzeichnis aufgenommenen Arbeiten können selbstverständlich nur als eine Auswahl angesehen werden; eine Zusammenstellung der fast unübersehbar gewordenen Literatur zur linearen Algebra geben die beiden von FADDEEVA und Mitarbeitern erarbeiteten Bibliographien.

Es ist den Autoren ein Bedürfnis, an dieser Stelle all denen herzlich zu danken, die ihnen im Laufe der mehr als vier Jahre dauernden Erarbeitung des Manuskriptes zur Seite gestanden haben: Die Universität Warschau und die Martin-Luther-Universität Halle ermöglichten uns auch unter schwierigen Bedingungen die für ein solches Projekt erforderlichen wechselseitigen Besuche. Förderlich wirkten sich auch Einladungen des erstgenannten Autors zu Gastvorlesungen an die Technische Universität Dresden aus, wo der zweitgenannte bis 1979 tätig gewesen ist. Herr Prof. Dr. Dr. h. c. H. HEINRICH hat als Mitherausgeber der Reihe, in der dieses Buch erscheint, zahlreiche Anregungen für die Formulierung der Schlußfassung gegeben. Unsere Kollegen Prof. Dr. G. ZIELKE und Dr. O. KNOTH waren uns bei der Beschaffung von Literatur behilflich; Prof. ZIELKE hat auch Teile des Manuskriptes kritisch gelesen. Viele andere Kollegen - stellvertretend sei Dr. V. TILLER genannt unterstützten uns bei der Bewältigung der technischen Arbeiten tatkräftig. Der Deutsche Verlag der Wissenschaften und insbesondere Frau Dipl.-Math. E. ARNDT haben die Entstehung dieses Buches mit Verständnis und Geduld gefördert. Zu herzlichem Dank sind die Autoren Frau E. SCHÄFER verpflichtet, die den größten Teil des Manuskriptes sorgfältig geschrieben hat.

Unser besonderer Dank gilt unseren Familien, die während einer längeren Zeit mancher Belastung ausgesetzt waren und die durch ihr Verständnis, aber auch durch ihre aktive Unterstützung zum Gelingen des Projektes beigetragen haben.

Die Autoren hoffen, daß sich die Arbeit gelohnt hat und der Leser mit Hilfe dieses Buches ein tieferes Verständnis der Probleme der numerischen linearen Algebra gewinnt. Vor allem möge er in die Lage versetzt werden, Möglichkeiten und Grenzen der verschiedenen Algorithmen richtig einzuschätzen, die verfügbare Software effektiv zu nutzen und diese Software gegebenenfalls anzupassen, zu modifizieren oder zu ergänzen, wenn spezielle, von Standardaufgaben abweichende Probleme dies erfordern.

Warschau/Halle, Frühjahr 1986

Inhalt

Те	il I. Grundlagen	•	11
1.	Grundbegriffe der linearen Algebra		11 11 31
2.	Aufgaben, Computer, Algorithmen		45 45 63 73
3.	Elementare Transformationsmatrizen		96 97 100 106
Te	il II. Reguläre lineare Gleichungssysteme		127
4.	Grundlegende Fakten über reguläre Gleichungssysteme		127 127 142 143
5.	Der Gaußsche Algorithmus		153 153 161 168 174
6.	Modifikationen des Gaußschen Algorithmus		184 184 201 207 220 229
7.	Zusammenfassung zum Teil II		237
Те	eil III. Lineare Quadratmittelprobleme		240
8.	Rechteckige Gleichungssysteme und Quadratmittelprobleme		240 240

	8.2. Störungstheorie		•	$\begin{array}{c} 251 \\ 264 \end{array}$
9.	Normalgleichungsverfahren	•		270 270 275
10.	Orthogonalisierungsverfahren			277 277 285 301
11.	Rangdefiziente Quadratmittelprobleme			315 316 325 338
12.	Zusammenfassung zum Teil III	•		350
Tei	il IV. Eigenwertprobleme			353
13.	Das spezielle symmetrische Eigenwertproblem 13.1. 13.1. Grundlegende Eigenschaften, Störungstheorie, Residualkriterien 13.1. 13.2. Das Jacobi-Verfahren 13.2. 13.3. Vektor- und Teilraumiteration 13.3. 13.4. Inverse Iteration 13.4. 13.5. Orthogonale Ähnlichkeitstransformation auf Tridiagonalform 13.6. 13.6. Schneiden des Spektrums und Bisektionsverfahren 13.7. 13.7. Der QR-Algorithmus 13.7.	•	· · · ·	353 354 365 373 389 397 406 415
14.	Das allgemeine symmetrische Eigenwertproblem14.1. Grundlegende Eigenschaften14.2. Explizite Reduktion auf ein spezielles Eigenwertproblem14.3. Vektor- und Teilraumiteration	• • •	• • •	434 434 438 441
15.	Zusammenfassung zum Teil IV			448
16.	Software für Aufgabenklassen der numerischen linearen Algebra	•		450
Lite	eratur	•	•	457
Sac	hverzeichnis			468

I. Grundlagen

1. Grundbegriffe der linearen Algebra

Dieses Kapitel dient der Einführung der später benötigten grundlegenden Begriffe und Aussagen der linearen Algebra. Ein großer Teil des behandelten Stoffes — speziell im Abschnitt 1.1 — ist Gegenstand der mathematischen Grundausbildung an den Hochschulen und wird deshalb nur übersichtsartig zusammengestellt. Weiterführende, dem Naturwissenschaftler und Ingenieur in der Regel nicht vertraute Ergebnisse werden in Form von Aussagen und Sätzen dargeboten. Allerdings muß auch hier auf detaillierte Beweise verzichtet werden; für diese sei auf die Spezialliteratur verwiesen, siehe Bemerkung B 1.1.

1.1. Vektoren, Matrizen, Normen

A. Vektoren und elementare Vektoroperationen

Der Vektorbegriff ist ein zentraler Begriff der linearen Algebra. Unter einem Vektor – genauer: einem *m-dimensionalen Spaltenvektor* – verstehen wir ein *m*-Tupel

$$\boldsymbol{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{pmatrix}$$

von reellen Zahlen $x_1, x_2, ..., x_m$, den Komponenten des Vektors x. Die Menge aller *m*-dimensionalen Vektoren bezeichnen wir mit \mathbb{R}^m ; für \mathbb{R}^1 — die Menge der reellen Zahlen — schreiben wir einfach \mathbb{R} . Für die Bezeichnung von Vektoren verwenden wir i. allg. kleine lateinische Buchstaben x, y, ... Die zugehörigen unten indizierten Größen $x_i, y_i, ...$ stellen die Komponenten des jeweiligen Vektors dar; gelegentlich benutzen wir auch die Schreibweise $(x)_i$. Für reelle Zahlen, die nicht Komponenten von Vektoren sind, werden i. allg. kleine griechische Buchstaben gewählt.

Zwei Vektoren $x, y \in \mathbb{R}^m$ sind gleich, wenn $x_i = y_i$ (i = 1, ..., m) ist. Der Vektor, dessen sämtliche Komponenten 0 sind, heißt Nullvektor $o \in \mathbb{R}^m$.

In der Menge \mathbb{R}^m ist für zwei Vektoren x, y gemäß $(x + y)_i := x_i + y_i$ die Summe x + y definiert. Analog ist das Produkt λx bzw. $x\lambda$ von $x \in \mathbb{R}^m$ mit einer Zahl $\lambda \in \mathbb{R}$ durch $(\lambda x)_i := (x\lambda)_i := \lambda \cdot x_i$ erklärt. Das Zeichen ":=" wird dabei zur Definition der auf der Seite des Doppelpunktes stehenden Größe verwendet. Es gilt $x + y \in \mathbb{R}^m$ sowie $\lambda x \in \mathbb{R}^m$ für alle $x, y \in \mathbb{R}^m$ und $\lambda \in \mathbb{R}$, und aus später klar werdenden Gründen wird \mathbb{R}^m daher *m*-dimensionaler linearer Raum genannt.

Die soeben eingeführten Operationen genügen den Rechenregeln

$$x + y = y + x$$
, $(x + y) + z = x + (y + z) = x + y + z$, (1)

$$(\lambda \mu) \boldsymbol{x} = \lambda(\mu \boldsymbol{x}), \qquad (2)$$

$$(\lambda + \mu) \boldsymbol{x} = \lambda \boldsymbol{x} + \mu \boldsymbol{x}, \qquad \lambda (\boldsymbol{x} + \boldsymbol{y}) = \lambda \boldsymbol{x} + \lambda \boldsymbol{y}$$
 (3)

für alle $x, y, z \in \mathbb{R}^m$ und alle $\lambda, \mu \in \mathbb{R}$.

Das folgende Beispiel zeigt, daß die auf der rechten und linken Seite der Gleichheitszeichen in (2), (3) stehenden Ausdrücke zwar dieselben Größen definieren, ihre Berechnung in der durch die Klammersetzung festgelegten Reihenfolge jedoch i. allg. einen unterschiedlichen Aufwand erfordert:

1.1.1. Beispiel. Man betrachte $\boldsymbol{u} := (\lambda + \mu) \boldsymbol{x}$ und $\boldsymbol{v} := \lambda \boldsymbol{x} + \mu \boldsymbol{x}$ für gegebenes $\boldsymbol{x} \in \mathbf{R}^m$, $\lambda, \mu \in \mathbf{R}$, sowie die zugehörigen Berechnungsvorschriften

 $V_1: \quad \alpha := \lambda + \mu$

for i := 1(1)m do $u_i := \alpha * x_i$,

V₂: for i := 1(1)m do $v_i := (\lambda * x_i) + (\mu * x_i)$.

 V_1 erfordert eine Addition und *m* Multiplikationen, während V_2 *m* Additionen und 2*m* Multiplikationen erfordert, obwohl beide Vorschriften denselben Vektor u = v definieren. \Box

In Berechnungsvorschriften bedeutet das Zeichen ":=", daß der links stehenden Variablen der Wert des rechtsstehenden Ausdrucks zugewiesen wird.

Das Beispiel zeigt, daß in der numerischen linearen Algebra im Gegensatz zur theoretischen linearen Algebra zwischen der *Definition* einer Größe durch einen Ausdruck und der geeignet vorzunehmenden *numerischen Auswertung* dieses Ausdruckes durch eine Berechnungsvorschrift unterschieden werden muß. Neben dem Aufwand ist dabei auch der Einfluß der bei Rechnung mit endlicher Stellenzahl notwendig auftretenden Rundungsfehler zu beachten, siehe Abschnitt 2.2 und 2.3 für eine ausführliche Diskussion dieses Problemkreises.

B. Linearkombinationen und Teilräume

Wir wählen jetzt *n* Vektoren $a^1, \ldots, a^n \in \mathbb{R}^m$ und Zahlen $\lambda_1, \ldots, \lambda_n \in \mathbb{R}$ aus; zur Unterscheidung von Vektoren werden stets obere Indizes verwendet, während reelle Zahlen unten indiziert werden. Der Vektor

$$oldsymbol{y} = \lambda_1 oldsymbol{a}^1 + \lambda_2 oldsymbol{a}^2 + \cdots + \lambda_n oldsymbol{a}^n$$

ist dann wieder ein Vektor aus \mathbb{R}^n und heißt *Linearkombination* der Vektoren a^1, \ldots, a^n mit den *Koeffizienten* $\lambda_1, \ldots, \lambda_n$. Die Menge aller Linearkombinationen von

 a^1, \ldots, a^n wird mit

$$ext{span} \left\{ oldsymbol{a}^1, \, ..., \, oldsymbol{a}^n
ight\} := \left\{ oldsymbol{y} \in oldsymbol{R}^m \colon oldsymbol{y} = \lambda_1 oldsymbol{a}^1 + \dots + \lambda_n oldsymbol{a}^n
ight; \lambda_1, \, ..., \, \lambda_n \in oldsymbol{R}
ight\}$$

bezeichnet.

Für $\boldsymbol{y} = \sum_{j=1}^{n} \lambda_j \boldsymbol{a}^j, \, \boldsymbol{y}' = \sum_{j=1}^{n} \lambda'_j \boldsymbol{a}^j$ und $\alpha \in \mathbb{R}$ gilt $\boldsymbol{y} + \boldsymbol{y}' = \sum_{j=1}^{n} (\lambda_j + \lambda'_j) \, \boldsymbol{a}^j$ sowie $\alpha \boldsymbol{y}$ = $\sum_{j=1}^{n} (\alpha \lambda_j) \, \boldsymbol{a}^j, \, \mathrm{d.} \, \mathrm{h.}, \, \mathrm{auch} \, \boldsymbol{y} + \boldsymbol{y}'$ und $\alpha \boldsymbol{y}$ sind Linearkombinationen $\mathbf{v} \ge \mathbf{n} \, \boldsymbol{a}^1, \ldots, \boldsymbol{a}^n$. Daher hat $X := \mathrm{span} \{ \boldsymbol{a}^1, \ldots, \boldsymbol{a}^n \} \subset \mathbb{R}^m$ folgende Eigenschaften:

Aus $y, y' \in X$ und $\alpha \in \mathbf{R}$ folgt $y + y' \in X$ sowie $\alpha y \in X$.

Eine Teilmenge X von \mathbb{R}^m mit diesen Eigenschaften heißt (linearer) Teilraum von \mathbb{R}^m . Speziell heißt span $\{a^1, \ldots, a^n\}$ der durch die Vektoren a^1, \ldots, a^n aufgespannte Teilraum. Im Dreidimensionalen (m = 3) spannt ein Vektor $a^1 \pm o$ eine Gerade durch den Ursprung auf, und zwei Vektoren a^1, a^2 , die nicht Yielfache voneinander sind, spannen eine durch den Ursprung gehende Ebene auf.

Speziell sind $X = \{o\}$ — die nur aus dem Nullvektor bestehende Menge — und $X = \mathbf{R}^m$ selbst Teilräume von \mathbf{R}^m , so daß insbesondere die Bezeichnung von \mathbf{R}^m als linearer Raum berechtigt ist.

Die m-dimensionalen Koordinatenvektoren

$$\boldsymbol{e}^{i} := \begin{bmatrix} 0\\ \vdots\\ 0\\ 1\\ 0\\ \vdots\\ 0 \end{bmatrix} \quad \leftarrow i \text{-te Komponente } (i-1,...,m)$$

spannen den Raum \mathbf{R}^m auf, denn es gilt $\boldsymbol{x} = \sum_{j=1}^m x_j \boldsymbol{e}^j$ für jedes $\boldsymbol{x} \in \mathbf{R}^m$. Andererseits brauchen *m* beliebige Vektoren nicht notwendig den ganzen Raum aufzuspannen, wie das folgende Beispiel zeigt.

1.1.2. Beispiel. Für

$$\boldsymbol{a}^1 = \begin{pmatrix} 1\\0\\0 \end{pmatrix}, \quad \boldsymbol{a}^2 = \begin{pmatrix} 0\\1\\0 \end{pmatrix}, \quad \boldsymbol{a}^3 = \begin{pmatrix} 1\\1\\0 \end{pmatrix}$$

gilt $a^3 = a^1 + a^2$. Jede Linearkombination $y = \lambda_1 a^1 + \lambda_2 a^2 + \lambda_3 a^3$ läßt sich daher als $y = \lambda_1 a^1 + \lambda_2 a^2 + \lambda_3 (a^1 + a^2) = (\lambda_1 + \lambda_3) a^1 + (\lambda_2 + \lambda_3) a^2$ schreiben, woraus

$$\operatorname{span} \{ \boldsymbol{a^1}, \, \boldsymbol{a^2}, \, \boldsymbol{a^3} \} = \operatorname{span} \{ \boldsymbol{a^1}, \, \boldsymbol{a^2} \}$$

folgt.

Es ist daher sinnvoll, zu fragen, wieviel Vektoren höchstens benötigt werden,

um einen vorgegebenen Teilraum aufzuspannen. Zur Erörterung dieser Frage wird der fundamentale Begriff der linearen Unabhängigkeit eingeführt: Die Vektoren $a^1, \ldots, a^n \in \mathbb{R}^m$ heißen *linear unabhängig*, wenn aus $\lambda_1 a^1 + \cdots + \lambda_n a^n = o$ folgt, daß $\lambda_1 = \cdots = \lambda_n = 0$ ist, d. h., wenn nur die Linearkombination mit sämtlich verschwindenden Koeffizienten den Nullvektor ergibt. Äquivalent dazu ist die mehr geometrische Charakterisierung, daß keiner der Vektoren a^1, \ldots, a^n beim Aufspannen von span $\{a^1, \ldots, a^n\}$ weggelassen werden kann. Im Beispiel 1.1.2 gilt $a^1 + a^2 - a^3$ = o mit nicht identisch verschwindenden Koeffizienten, und a^3 kann beim Aufspannen von span $\{a^1, a^2, a^3\}$ weggelassen werden.

Die Vektoren a^1, \ldots, a^n heißen *linear abhängig*, wenn sie nicht linear unabhängig sind. Dann kann wie im Beispiel mindestens einer dieser Vektoren als Linearkombination der übrigen ausgedrückt und daher beim Aufspannen von span $\{a^1, \ldots, a^n\}$ weggelassen werden. Dieser Prozeß ist so lange fortsetzbar, bis $r \leq n$ linear unabhängige Vektoren erhalten werden — diese seien ohne Einschränkung der Allgemeinheit a^1, \ldots, a^r —, die span $\{a^1, \ldots, a^n\}$ = span $\{a^1, \ldots, a^r\}$ aufspannen.

Eine Menge von linear unabhängigen Vektoren a^1, \ldots, a^r , die einen Teilraum X aufspannen, wird *Basis* von X genannt. Zu jedem Teilraum $X \neq \{o\}$ gibt es unendlich viele verschiedene Basen, aber jede von diesen besteht aus derselben Anzahl rlinear unabhängiger Vektoren, die *Dimension* von X — in Zeichen: $r = \dim(X)$ genannt wird.

Die beiden extremalen Teilräume $X = \{o\}$ bzw. $X = \mathbb{R}^m = \text{span} \{e^1, \dots, e^m\}$ von \mathbb{R}^m haben die Dimension 0 bzw. m, und für jeden Teilraum $X \subset \mathbb{R}^m$ gilt $0 \leq r$ = dim $(X) \leq m$. Insbesondere sind also mehr als m Vektoren aus \mathbb{R}^m stets linear abhängig.

C. Matrizen und elementare Matrixoperationen

Neben dem Vektorbegriff spielt der Begriff einer Matrix in der linearen Algebra eine entscheidende Rolle. Unter einer Matrix – genauer: einer (m,n)-Matrix bzw. Matrix vom Format (m, n) – verstehen wir ein Schema

$$A = egin{pmatrix} a_{11} & a_{12} & \ldots & a_{1n} \ a_{21} & a_{22} & \ldots & a_{2n} \ dots & dots & dots \ a_{m1} & a_{m2} & \ldots & a_{mn} \end{pmatrix} = (a_{ij})$$

von $m \cdot n$ reellen Zahlen a_{ij} (i = 1, ..., m; j = 1, ..., n), den Elementen der Matrix A. Matrizen werden stets mit großen lateinischen, gelegentlich auch griechischen Buchstaben und ihre Elemente mit den zugehörigen kleinen, unten doppelt indizierten Buchstaben bezeichnet. Für a_{ij} schreiben wir auch $(A)_{ij}$. Zur Unterscheidung von Matrizen und deren Elementen werden obere, in Klammern gesetzte Indizes verwendet wie $A^{(0)} = (a_{ij}^{(0)}), A^{(1)} = (a_{ij}^{(1)})$ usw., sofern auf die Elemente Bezug genommen wird. Andernfalls kann auch unten indiziert werden, etwa A_1, A_2 usw.

Der erste Index *i* heißt Zeilenindex, der zweite Index *j* Spaltenindex des Elementes a_{ij} . Die Menge aller (m,n)-Matrizen wird mit $\mathbf{R}^{m,n}$ bezeichnet. In natürlicher Weise

können (m,1)-Matrizen als *m*-dimensionale *Spaltenvektoren*, (1,n)-Matrizen als *n*-dimensionale *Zeilenvektoren* aufgefaßt werden.

Der Spaltenvektor

$$\begin{pmatrix} a_{1j} \\ \vdots \\ a_{mj} \end{pmatrix} \in \mathbf{R}^{m,1}$$

heißt j-te Spalte von A, der Zeilenvektor

$$(a_{i1}, ..., a_{in}) \in \mathbf{R}^{1,n}$$

i-te Zeile von A. Im Fall m = n heißt *A quadratisch* und von der *Ordnung n.* Die Elemente $a_{11}, a_{22}, \ldots, a_{ll}, l := \min \{m, n\}$, einer Matrix *A* werden *Diagonalelemente* von *A*, die zugehörigen Positionen im Rechteckschema *Diagonale* bzw. *Hauptdiagonale* genannt.

Eine Matrix A heißt Diagonalmatrix, wenn $a_{ij} = 0$ für $i \neq j$ ist, d. h., wenn alle Nichtdiagonalelemente verschwinden. Wir verzichten dann auch auf die doppelte Indizierung und schreiben einfach $A = \text{diag}(a_1, \ldots, a_l) = \text{diag}(a_i)$, wobei das Format angegeben werden muß bzw. aus dem Zusammenhang ersichtlich ist. Zum Beispiel sind

$$\boldsymbol{A}_{1} = \begin{pmatrix} 5 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 9 \end{pmatrix}, \quad \boldsymbol{A}_{2} = \begin{pmatrix} 5 & 0 & 0 & 0 \\ 0 & 7 & 0 & 0 \\ 0 & 0 & 9 & 0 \end{pmatrix}, \quad \boldsymbol{A}_{3} = \begin{pmatrix} 5 & 0 & 0 \\ 0 & 7 & 0 \\ 0 & 0 & 9 \\ 0 & 0 & 0 \end{pmatrix}$$

Diagonalmatrizen diag (5, 7, 9) der Formate (3, 3), (3, 4) und (4, 3).

Eine Matrix A wird Bandmatrix der Bandbreite l := p + q + 1 genannt, wenn $a_{ij} = 0$ für j < i - p und i + q < j mit Zahlen $p, q \ge 0$ gilt. Dies bedeutet, daß neben der Hauptdiagonalen nur p untere und q obere Nebendiagonalen von 0 verschiedene Elemente besitzen dürfen. Zum Beispiel sind

	(×	Х	0	0)			(×	Х	0	0	0)
	0	Х	Х	0			×	Х	×	0	0
$A_1 =$	0	0	Х	X	und	$A_2 =$	0	Х	\times	\times	0
	0	0	0	X			0	0	×	X	\times
	0	0	0	0 j			0	0	0	Х	×j

Bandmatrizen der Bandbreite 2 bzw. 3; A_1 heißt (obere) Bidiagonalmatrix, A_2 Tridiagonalmatrix. Hierbei kennzeichnet "ד das Besetztheitsmuster von A, d. h. diejenigen Positionen, wo von 0 verschiedene Elemente auftreten dürfen.

Eine Matrix A heißt obere bzw. untere Trapezmatrix, wenn $a_{ij} = 0$ ist für i > j bzw. i < j. Eine obere Trapezmatrix mit $m \ge n$ wird obere Dreiecksmatrix genannt,

eine untere Trapezmatrix mit $m \leq n$ heißt untere Dreiecksmatrix. Beispiele sind

$$A_{1} = \begin{pmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \\ 0 & 0 & 0 \end{pmatrix}, \qquad A_{2} = \begin{pmatrix} \times & \times & \times & \times \\ 0 & \times & \times & \times \\ 0 & 0 & \times & \times \end{pmatrix},$$
$$A_{3} = \begin{pmatrix} \times & \times & \times \\ 0 & \times & \times \\ 0 & 0 & \times \end{pmatrix}, \qquad A_{4} = \begin{pmatrix} \times & 0 & 0 \\ \times & \times & 0 \\ \times & \times & \times \\ \times & \times & \times \end{pmatrix},$$

wobei A_2 eine obere Trapezmatrix ist, A_1 und A_3 obere Dreiecksmatrizen sind und A_4 eine untere Trapezmatrix ist. Schließlich heißt A obere Hessenbergmatrix, wenn $a_{ij} = 0$ für i > j + 1 gilt; für m = n = 4 ergibt sich das Besetztheitsmuster

$$m{A} = egin{pmatrix} imes & imes & imes \ imes & imes & imes \ imes & imes & imes \ 0 & imes & imes & imes \ 0 & 0 & imes & imes \end{pmatrix}$$

Jede Matrix $A \in \mathbf{R}^{m,n}$ läßt sich als Spaltenvektor $a \in \mathbf{R}^{m,n}$ auffassen, indem die Elemente von A nacheinander spaltenweise als Komponenten des Vektors a angeordnet werden. Dem Element a_{ij} entspricht dann die Komponente a_l von a, wobei durch die Beziehung

$$l = (j - 1) m + i$$
 $(i = 1, ..., m; j = 1, ..., n)$

eine eineindeutige Zuordnung zwischen (i, j) und l gegeben ist. Bei Verwendung der Programmiersprache FORTRAN wird übrigens eine Matrix in dieser Form als einfachindiziertes Feld abgespeichert, wobei die Zeilenzahl m für die Indexzuordnung bekannt sein muß, vgl. Kapitel 16.

Für $A, B \in \mathbb{R}^{m,n}$ und $\lambda \in \mathbb{R}$ kann A + B, $\lambda A \in \mathbb{R}^{m,n}$ definiert werden, indem die zugehörigen "langen" Vektoren denselben Operationen unterworfen werden. Dies führt auf $(A + B)_{ij} = (A)_{ij} - (B)_{ij}$, $(\lambda A)_{ij} = (A\lambda)_{ij} = \lambda(A)_{ij}$, (i = 1, ..., m, j = 1, ..., n), d. h., Summe und Produkt mit einer Zahl sind elementweise zu bilden. Damit wird $\mathbb{R}^{m,n}$ ein linearer Raum der Dimension $m \cdot n$. Die Gleichheit A = B bedeutet $(A)_{ij} = (B)_{ij}$ (i = 1, ..., m; j = 1, ..., n), die Nullmatrix $O \in \mathbb{R}^{m,n}$ besteht aus den $m \cdot n$ Elementen Null, und die Regeln (1), (2), (3) gelten sinngemäß.

Jeder Matrix $A \in \mathbb{R}^{m,n}$ kann die transponierte Matrix $A^{\mathsf{T}} \in \mathbb{R}^{n,m}$ durch die Festlegung $(A^{\mathsf{T}})_{ji} := (A)_{ij}$ (i = 1, ..., m; j = 1, ..., n) zugeordnet werden. Sie entsteht aus A durch Spiegelung an der Diagonalen, im Beispiel

$$A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix}, \quad A^{\intercal} = \begin{pmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{pmatrix}.$$

Durch Transposition geht der Spaltenvektor $\boldsymbol{x} = (x_i) \in \mathbf{R}^m = \mathbf{R}^{m,1}$ in den *m*-dimensionalen Zeilenvektor $\boldsymbol{x}^{\mathsf{T}} = (x_1, \ldots, x_m) \in \mathbf{R}^{1,m}$ über, so daß \boldsymbol{x} platzsparend als $\boldsymbol{x} = (x_1, \ldots, x_m)^{\mathsf{T}} \in \mathbf{R}^m$ geschrieben werden kann. Die Transposition genügt den Regeln

$$(A^{\mathsf{T}})^{\mathsf{T}} = A, \qquad (A + B)^{\mathsf{T}} = A^{\mathsf{T}} + B^{\mathsf{T}}, \qquad (\lambda A)^{\mathsf{T}} = \lambda (A^{\mathsf{T}})$$
(4)

für alle $A, B \in \mathbb{R}^{m,n}$ und alle $\lambda \in \mathbb{R}$.

Eine quadratische Matrix A wird symmetrisch genannt, wenn

$$A = A^{\intercal}$$

gilt; z. B. ist

$$A = egin{pmatrix} 1 & 2 & 3 \ 2 & 4 & 5 \ 3 & 5 & 6 \end{pmatrix}$$

eine symmetrische Matrix. Wegen (4) sind Summe und skalare Vielfache symmetrischer Matrizen wieder symmetrisch. Die Menge $S^{n,n}$ aller reellen symmetrischen Matrizen der Ordnung *n* bildet daher einen Teilraum von $\mathbb{R}^{n,n}$. Jede symmetrische Matrix ist durch die Elemente a_{ij} $(i \leq j)$ im oberen Dreieck eindeutig festgelegt, so daß $S^{n,n}$ von der Dimension n(n + 1)/2 ist.

In vielen Anwendungen ist es zweckmäßig, eine Matrix nicht durch ihre Elemente, sondern durch regelmäßig angeordnete Teilmatrizen niedrigerer Dimension — sog. *Blöcke* — zu charakterisieren. Man spricht dann von Hyper- oder *Blockmatrizen*; Beispiele sind

$$A = egin{pmatrix} A_{11} & A_{12} & A_{13} \ A_{21} & A_{22} & A_{23} \end{pmatrix}, \qquad B = egin{pmatrix} B_{11} & B_{12} & B_{13} \ B_{21} & B_{22} & B_{23} \end{pmatrix}.$$

Dabei wird vorausgesetzt, daß nebeneinander stehende Blöcke dieselbe Zeilenzahl, übereinander stehende dieselbe Spaltenzahl besitzen. Falls A und B dieselbe Blockstruktur haben, ergibt sich die Summe offensichtlich blockweise zu

$$m{A} + m{B} = egin{pmatrix} m{A_{11}} + m{B_{11}} & m{A_{12}} + m{B_{12}} & m{A_{13}} + m{B_{13}} \ m{A_{21}} + m{B_{21}} & m{A_{22}} + m{B_{22}} & m{A_{23}} + m{B_{23}} \end{pmatrix}$$

Durch Transposition von A entsteht

$$A^{\intercal} = egin{pmatrix} A^{\intercal}_{11} & A^{\intercal}_{21} \ A^{\intercal}_{12} & A^{\intercal}_{22} \ A^{\intercal}_{13} & A^{\intercal}_{23} \end{pmatrix}.$$

D. Multiplikative Matrixoperationen

Wir betrachten wieder *n* Vektoren $a^1, \ldots, a^n \in \mathbb{R}^m$, die wir als Spalten der Matrix

$$A = (a_{ij}) = (a^1, ..., a^n) \in \mathbb{R}^{m,n}$$

auffassen. Jede Linearkombination $y = x_1 a^1 + \cdots + x_n a^n$ mit den im Vektor

 $\boldsymbol{x} \in \mathbf{R}^n$ zusammengefaßten Koeffizienten x_1, \ldots, x_n hat dann die Komponenten

$$y_i = \sum_{j=1}^n x_j (\boldsymbol{a}^j)_i = \sum_{j=1}^n a_{ij} x_j \qquad (i = 1, ..., m).$$
(5)

Wir nennen den Vektor $\boldsymbol{y} \in \mathbf{R}^m$ mit den durch (5) definierten Komponenten das Produkt der Matrix $\boldsymbol{A} \in \mathbf{R}^{m,n}$ mit dem Vektor $\boldsymbol{x} \in \mathbf{R}^n$ und schreiben dafür

$$\boldsymbol{y} = A\boldsymbol{x}.\tag{6}$$

Wenn \boldsymbol{x} seinerseits aus einem Vektor $\boldsymbol{w} \in \mathbf{R}^{l}$ durch Multiplikation mit der Matrix $\boldsymbol{B} \in \mathbf{R}^{n,l}$ gemäß

$$\boldsymbol{x} = \boldsymbol{B}\boldsymbol{w} \tag{7}$$

entstanden ist, erhebt sich die Frage, wie die Komponenten von $\boldsymbol{y} = A\boldsymbol{x} = A(\boldsymbol{B}\boldsymbol{w})$ mit denen von \boldsymbol{w} zusammenhängen. Indem (7) komponentenweise als $x_j = \sum_{k=1}^{l} b_{jk} w_k$ geschrieben und in (5) eingesetzt wird, ergibt sich

$$y_{i} = \sum_{j=1}^{n} a_{ij} x_{j} = \sum_{j=1}^{n} a_{ij} \left(\sum_{k=1}^{l} b_{jk} w_{k} \right) = \sum_{k=1}^{l} \left(\sum_{j=1}^{n} a_{ij} b_{jk} \right) w_{k} = \sum_{k=1}^{l} c_{ik} w_{k}$$

 mit

$$c_{ik} := \sum_{j=1}^{n} a_{ij} b_{jk} \qquad (i = 1, ..., m; \ k = 1, ..., l).$$
(8)

Fassen wir die c_{ik} als Elemente einer Matrix $C = (c_{ik}) \in \mathbb{R}^{m,l}$ auf, so gilt folglich y = A(Bw) = Cw. Dieses Ergebnis legt es nahe, für zwei Matrizen $A \in \mathbb{R}^{m,n}$, $B \in \mathbb{R}^{n,l}$ das Produkt

$$\boldsymbol{C} = \boldsymbol{A}\boldsymbol{B} \tag{9}$$

von A mit B als (m, l)-Matrix mit den durch (8) erklärten Elementen zu definieren. Für den Sonderfall eines Spaltenvektors B, d. h. für $B = x \in \mathbb{R}^{n,1}$, ist diese Definition mit der des Produktes Ax identisch.

Mit etwas Schreibarbeit kann nachgerechnet werden, daß die Matrixmultiplikation den Regeln

$$A(\lambda B) = (\lambda A) B = \lambda(AB), \qquad (10)$$

$$(\mathbf{A} + \mathbf{B}) \mathbf{C} = \mathbf{A}\mathbf{C} + \mathbf{B}\mathbf{C}, \qquad \mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{A}\mathbf{B} + \mathbf{A}\mathbf{C}, \tag{11}$$

$$(\boldsymbol{A}\boldsymbol{B})^{\mathsf{T}} = \boldsymbol{B}^{\mathsf{T}}\boldsymbol{A}^{\mathsf{T}} \tag{12}$$

für alle Matrizen A, B, C passenden Formats und alle $\lambda \in \mathbf{R}$ genügt.

Für zwei Matrizen, A, B können beide Produkte AB und BA genau dann gebildet werden, wenn $A \in \mathbb{R}^{m,n}$ und $B \in \mathbb{R}^{n,m}$ ist, aber selbst im Fall m = n gilt i. allg.

$$AB \neq BA$$

wie das Beispiel

$$oldsymbol{A}=egin{pmatrix} 1&2\2&3 \end{pmatrix}, \quad oldsymbol{B}=egin{pmatrix} 2&1\1&4 \end{pmatrix}, \quad oldsymbol{A}B=egin{pmatrix} 4&9\7&14 \end{pmatrix} =egin{pmatrix} 4&7\9&14 \end{pmatrix} =oldsymbol{B}A$$

zeigt. Falls AB = BA gilt, heißen A und B vertauschbar.

Für jede Matrix $A \in \mathbb{R}^{m,n}$ können die beiden Produkte $A^{\mathsf{T}}A \in \mathbb{R}^{n,n}$ sowie $AA^{\mathsf{T}} \in \mathbb{R}^{m,m}$ stets gebildet werden, und beide sind symmetrisch. Das Produkt von oberen bzw. unteren Dreiecksmatrizen ist wieder eine obere bzw. untere Dreiecksmatrix, insbesondere ist das Produkt von Diagonalmatrizen wieder diagonal. Wie das obige Beispiel zeigt, ist dagegen das Produkt symmetrischer Matrizen i. allg. nicht symmetrisch.

Für $A \in \mathbb{R}^{m,n}$, $B \in \mathbb{R}^{n,l}$ und $C \in \mathbb{R}^{l,k}$ kann (AB) C gebildet werden und ist vom Format (m, k). Solche Dreierprodukte genügen dem Assoziativgesetz

$$(AB) C = A(BC) =: ABC, \tag{13}$$

allerdings hat gerade hier die Klammersetzung häufig einen entscheidenden Einfluß auf die Anzahl der zur Berechnung erforderlichen Rechenoperationen und Zwischenspeicher, siehe \ddot{U} 1.1.4.

Eine besondere Rolle spielen die multiplikativen Verknüpfungen zweier Vektoren $a \in \mathbb{R}^m$, $b \in \mathbb{R}^n$. Während die Produkte ab und ba im Fall m, n > 1 nicht gebildet werden können, ist ab^{\intercal} definiert und stellt die Matrix

$$\boldsymbol{a}\boldsymbol{b}^{\mathsf{T}} = (a_i b_j) = \begin{pmatrix} a_1 b_1 & \dots & a_1 b_n \\ \vdots & & \vdots \\ a_m b_1 & \dots & a_m b_n \end{pmatrix} \in \mathbf{R}^{m,n}$$
(14)

dar, die *dyadisches Produkt* von \boldsymbol{a} mit \boldsymbol{b} genannt wird. Ebenso kann $\boldsymbol{b}\boldsymbol{a}^{\mathsf{T}} \in \mathbb{R}^{n,m}$ gebildet werden und ist gleich $(\boldsymbol{a}\boldsymbol{b}^{\mathsf{T}})^{\mathsf{T}}$.

Das Produkt $b^{\mathsf{T}}a$ ist für Vektoren $a, b \in \mathbb{R}^n$ gleicher Dimension stets erklärt und ergibt

$$\boldsymbol{b}^{\mathsf{T}}\boldsymbol{a} = b_1 a_1 + \dots + b_n a_n \in \mathsf{R}. \tag{15}$$

Es genügt dem Kommutativgesetz $\mathbf{a}^{\mathsf{T}}\mathbf{b} = \mathbf{b}^{\mathsf{T}}\mathbf{a}$ und heißt *Skalarprodukt* von \mathbf{a} und \mathbf{b} .

Unter Verwendung des Skalarproduktes läßt sich die folgende einprägsame Regel für die Berechnung der Elemente c_{ik} des Produktes C = AB, $A \in \mathbb{R}^{m,n}$, $B \in \mathbb{R}^{n,l}$ angeben: Falls $a^{i^{\dagger}}$ die Zeilen von A und b^k die Spalten von B gemäß

$$A = \begin{pmatrix} -a^{1\mathsf{T}} & -\\ \vdots \\ -a^{i\mathsf{T}} & -\\ \vdots \\ -a^{m\mathsf{T}} & - \end{pmatrix}, \quad B = \begin{pmatrix} | & | & |\\ b^1 & \dots & b^k & \dots & b^l \\ | & | & | & | \end{pmatrix}$$

bezeichnen, gilt $c_{ik} = \mathbf{a}^{i_{\mathsf{T}}} \mathbf{b}^{k}$ (i = 1, ..., m; j = 1, ..., l), d. h., c_{ik} ist das Skalarprodukt von \mathbf{a}^{i} mit \mathbf{b}^{k} , oder kurz: $c_{ik} = (i$ -te Zeile \mathbf{A}) mal (k-te Spalte \mathbf{B}). In analoger Weise wird bei der Produktbildung von Blockmatrizen vorgegangen. Als Beispiel betrachten wir die Matrizen

$$A = egin{pmatrix} A_{11} & A_{12} & A_{13} \ A_{21} & A_{22} & A_{23} \end{pmatrix} \quad ext{und} \quad B = egin{pmatrix} B_{11} & B_{12} \ B_{21} & B_{22} \ B_{21} & B_{22} \ B_{31} & B_{32} \end{pmatrix}$$

Das Produkt hat dann die Blockstruktur

$$AB = \begin{pmatrix} A_{11}B_{11} + A_{12}B_{21} + A_{13}B_{31} \mid A_{11}B_{12} + A_{12}B_{22} + A_{13}B_{32} \\ \hline A_{21}B_{11} + A_{22}B_{21} + A_{23}B_{31} \mid A_{21}B_{12} + A_{22}B_{22} + A_{23}B_{32} \end{pmatrix}$$

Dabei muß die Blockeinteilung selbstverständlich so sein, daß die Produkte $A_{ii}B_{ik}$ gebildet werden können. Allgemein ist das Produkt zweier Blockmatrizen $A = (A_{ij})$, $B = (B_{ik})$ eine Blockmatrix $C = AB = (C_{ik})$ mit den Blöcken

$$C_{ik} = \sum_{j=1}^m A_{ij} B_{jk},$$

wobei m die Anzahl der Blöcke in den "Zeilen" von A wie in den "Spalten" von Bbezeichnet. Speziell ergibt sich aus den obigen Regeln im Fall $B = (b^1, ..., b^l) \in \mathbf{R}^{n,l}$ die Darstellung des Produktes AB in der Form

$$\boldsymbol{AB} = \begin{pmatrix} | & | \\ \boldsymbol{Ab^1} & \dots & \boldsymbol{Ab^l} \\ | & | \end{pmatrix}.$$

Wir bemerken abschließend, daß sich die Multiplikation eines Vektors x mit einer Zahl λ den Regeln der Matrixmultiplikation unterordnet, wenn $x\lambda$ statt λx geschrieben wird, da dann die Formate (n, 1) von x und (1, 1) von λ im Sinne der Multiplikation zueinander passen.

E. Lineare Abbildungen

Wir betrachten wieder bei gegebener Matrix $A \in \mathbb{R}^{m,n}$ das durch (5), (6) erklärte Produkt

$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x}, \qquad \boldsymbol{x} \in \mathbf{R}^{\boldsymbol{n}}.$$
 (16)

Diese Beziehung kann so gedeutet werden, daß mittels der Matrix A jedem $x \in \mathbf{R}^n$ in eindeutiger Weise ein $y = Ax \in \mathbb{R}^m$ zugeordnet wird. Wir sagen auch, daß die *Matrix A* durch (16) eine *Abbildung* von \mathbb{R}^n in \mathbb{R}^m erklärt, und schreiben $A: \mathbb{R}^n \to \mathbb{R}^m$. Der Vektor y = Ax wird Bild von x genannt; umgekehrt heißt bei gegebenem $y \in \mathbf{R}^{n}$ jedes $x \in \mathbf{R}^{n}$ mit Ax = y ein Urbild von y. Es gibt Matrizen A, so daß

nicht jedes \boldsymbol{y} ein Urbild besitzt, z. B. hat im Fall $\boldsymbol{A} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{pmatrix}$ der Vektor $\boldsymbol{y} = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$ kein Urbild. Andererseits kann es zu einem Bild mehrere Urbilder geben;

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$$
für obige Matrix hat $\boldsymbol{y} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}$ die Urbilder $\boldsymbol{x} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} + \boldsymbol{\xi} \begin{pmatrix} 1 \\ 1 \\ -1 \end{pmatrix}, \boldsymbol{\xi} \in \mathbf{R}$ beliebig.

Die durch (16) festgelegte Abbildung $A : \mathbb{R}^n \to \mathbb{R}^m$ hat die Eigenschaften

$$A(\boldsymbol{x} + \boldsymbol{x}') = A\boldsymbol{x} + A\boldsymbol{x}', \qquad A(\lambda \boldsymbol{x}) = \lambda(A\boldsymbol{x})$$
(17)

und wird deshalb auch lineare Abbildung genannt. Die Menge

$$\mathcal{R}(A) := \{ \boldsymbol{y} \in \mathsf{R}^m \colon \boldsymbol{y} = A\boldsymbol{x}, \, \boldsymbol{x} \in \mathsf{R}^n \}$$
(18)

aller möglichen Bilder heißt Wertebereich von A. Wenn a^1, \ldots, a^n die Spalten von A bezeichnen, ist $\mathcal{R}(A) = \text{span} \{a^1, \ldots, a^n\}$, so daß $\mathcal{R}(A)$ auch Spaltenraum von A genannt wird und insbesondere einen Teilraum von \mathbb{R}^m darstellt. Unter dem Rang von A — in Zeichen: rang (A) — verstehen wir die Dimension von $\mathcal{R}(A)$, d. h. die Anzahl der linear unabhängigen Spalten von A. Es zeigt sich, daß die Anzahl der linear unabhängigen Spalten gleich der Anzahl der linear unabhängigen Zeilen von A ist, d. h., es gilt rang (A) = rang (A^T) . Hieraus folgt $r = \text{rang}(A) \leq \min \{m, n\}$. Im Fall $r = \min \{m, n\}$ sagt man, daß A Vollrang hat, andernfalls wird A rangdefizient genannt.

Wir betrachten jetzt eine Vollrangmatrix $A \in \mathbb{R}^{m,n}$. Im Fall $m \ge n$ muß dann r = n gelten, d. h., alle *n* Spalten von A müssen linear unabhängig sein, und A wird *spaltenregulär* genannt. Entsprechend müssen im Fall $m \le n$ alle *m* Zeilen linear unabhängig sein, und A heißt zeilenregulär. Genau in diesem Fall ist $\mathcal{R}(A) = \mathbb{R}^m$, d. h., jedes $y \in \mathbb{R}^m$ hat mindestens ein Urbild $x \in \mathbb{R}^n$, so daß man auch von einer durch A vermittelten Abbildung von \mathbb{R}^n auf \mathbb{R}^m spricht.

Eine quadratische Matrix $A \in \mathbf{R}^{n,n}$ von vollem Rang r = n ist sowohl spaltenals auch zeilenregulär und wird schlechthin *regulär* genannt; andernfalls heißt sie singulär.

Von besonderer Bedeutung sind diejenigen Vektoren, die durch eine Matrix $A \in \mathbf{R}^{m,n}$ in den Nullvektor abgebildet werden. Sie bilden den Nullraum

$$\mathcal{N}(A) := \{ \boldsymbol{x} \in \mathbf{R}^n \colon A\boldsymbol{x} = \boldsymbol{o} \}$$
⁽¹⁹⁾

von A, der wegen (17) in der Tat ein Teilraum von \mathbb{R}^n ist. Es gilt

$$\dim \left(\mathscr{N}(A) \right) + \dim \left(\mathscr{R}(A) \right) = \dim \left(\mathscr{N}(A) \right) + r = n,$$
⁽²⁰⁾

insbesondere besteht $\mathcal{N}(A)$ nur aus dem Nullvektor genau dann, wenn A spaltenregulär ist. Es gibt dann zu jedem $\mathbf{y} \in \mathcal{R}(A)$ genau ein Urbild $\mathbf{x} \in \mathbb{R}^n$. Ist A zudem quadratisch, also regulär, so existiert folglich zu jedem $\mathbf{y} \in \mathbb{R}^n$ genau ein Urbild $\mathbf{x} \in \mathbb{R}^n$. Die so definierte Abbildung, die jedem $\mathbf{y} \in \mathbb{R}^n$ eindeutig das $\mathbf{x} \in \mathbb{R}^n$ mit $\mathbf{y} = A\mathbf{x}$ zuordnet, heißt die zu $A: \mathbb{R}^n \to \mathbb{R}^n$ inverse Abbildung. Sie ist ebenfalls linear und kann mit einer durch $A \in \mathbb{R}^{n,n}$ eindeutig festgelegten Matrix $A^{-1} \in \mathbb{R}^{n,n}$ in der Form

$$\boldsymbol{x} = \boldsymbol{A}^{-1} \boldsymbol{y}, \qquad \boldsymbol{y} \in \mathbf{R}^n,$$
 (21)

geschrieben werden. Die Matrix A^{-1} heißt die zu A inverse Matrix und ist durch die äquivalenten Bedingungen

$$AA^{-1} = I \quad \text{und} \quad A^{-1}A = I \tag{22}$$

eindeutig charakterisiert, wobei

$$I := I_n := \begin{bmatrix} 1 & & & \\ & 1 & 0 & \\ & & \ddots & \\ 0 & & \ddots & \\ & & & 1 \end{bmatrix} \in \mathbf{R}^{n,n}$$

die *Einheitsmatrix* der Ordnung *n* bezeichnet. Für diese gilt Ix = x für alle $x \in \mathbb{R}^n$, d. h., sie realisiert die *identische Abbildung* von \mathbb{R}^n auf sich. Überdies ist

$$\boldsymbol{I}_{\boldsymbol{m}}\boldsymbol{A} = \boldsymbol{A}\boldsymbol{I}_{\boldsymbol{n}} = \boldsymbol{A} \tag{23}$$

für alle $A \in \mathbf{R}^{m,n}$.

Die Matrixinversion genügt den Regeln

$$(A^{-1})^{-1} = A$$
, $(\lambda A)^{-1} = \frac{1}{\lambda} A^{-1}$, $(A^{\mathsf{T}})^{-1} = (A^{-1})^{\mathsf{T}} = :A^{-\mathsf{T}}$ (24)

für alle regulären $A \in \mathbf{R}^{n,n}$ und alle $\lambda \in \mathbf{R}$ mit $\lambda \neq 0$.

Das Produkt zweier quadratischer Matrizen ist genau dann regulär, wenn jeder Faktor regulär ist. Ist dies der Fall, so gilt

$$(AB)^{-1} = B^{-1}A^{-1}.$$
(25)

F. Quadratische Formen

Es sei jetzt $A \in \mathbb{R}^{n,n}$ eine quadratische Matrix. Dann heißt die Abbildung, die jedem $x \in \mathbb{R}^n$ die reelle Zahl

$$\boldsymbol{x}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x} = \sum_{i,j=1}^{n} a_{ij} x_i x_j \tag{26}$$

zuordnet, die durch A erzeugte quadratische Form. Die Matrix A heißt positiv semidefinit, wenn $x^{\mathsf{T}}Ax \geq 0$ für alle $x \in \mathbb{R}^n$ gilt, d. h., wenn die quadratische Form nur nichtnegative Werte annimmt. Gilt sogar $x^{\mathsf{T}}Ax > 0$ für alle $x \neq o$, so heißt A positiv definit. Wenn die quadratische Form sowohl positive als auch negative Werte annehmen kann, heißt A indefinit. Häufig werden die eingeführten Definitheitsbegriffe nur für symmetrische Matrizen erklärt, vgl. dazu Ü 1.1.5.

G. Determinanten

Jeder quadratischen Matrix $A \in \mathbb{R}^{n,n}$ kann die als *Determinante* von A — in Zeichen: det (A) — bezeichnete reelle Zahl durch die folgende rekursive Definition zugeordnet werden:

(i) Für n = 1 ist det $(A) = a_{11}$.

(ii) Für n > 1 ist

$$\det (A) := \sum_{i=1}^{n} (-1)^{i+j} a_{ij} \det (A^{(i,j)}),$$

wobei $j \in \{1, ..., n\}$ ein beliebiger Spaltenindex ist und $A^{(i,j)}$ diejenige quadratische Untermatrix von A bezeichnet, die durch Streichen der *i*-ten Zeile und *j*-ten Spalte aus A entsteht.

Die Vorschrift (ii) wird auch *Entwicklungssatz* genannt; man spricht auch von der Entwicklung von det (A) nach den Elementen der *j*-ten Spalte. Eine analoge Formel gilt für die Entwicklung nach der *i*-ten Zeile. Man kann zeigen, daß det (A) nicht vom verwendeten Spaltenindex *j* (bzw. analog vom Zeilenindex *i*) abhängt, so daß det (A) durch die obige Definition in der Tat eindeutig festgelegt ist.

Für n = 2 ergibt sich speziell mit j = 1

$$\detegin{pmatrix} a_{11} & a_{12} \ a_{21} & a_{22} \end{pmatrix} = (-1)^2 \, a_{11} \det \, (a_{22}) + (-1)^3 \, a_{21} \det \, (a_{12}) = a_{11}a_{22} - a_{12}a_{21}.$$

Prinzipiell kann det (A) nach der angegebenen Definition für beliebiges n berechnet werden. Aus Aufwandsgründen scheidet dieser Weg für praktische Zwecke allerdings aus, siehe Ü 1.1.7. Aufwandsgünstigere Vorschriften zur Berechnung von det (A) werden in Kapitel 5 beschrieben.

Für die Anwendungen sind die folgenden Eigenschaften von det (A) von Bedeutung:

- (D₁) Bei Vertauschung zweier Zeilen bzw. Spalten ändert sich das Vorzeichen der Determinante.
- (D_2) Die Determinante ändert sich nicht, wenn ein Vielfaches einer Zeile bzw. Spalte zu einer anderen Zeile bzw. Spalte addiert wird.
- (D₃) Bei Multiplikation einer Zeile bzw. Spalte mit einem Faktor $\lambda \in \mathbb{R}$ ändert sich die Determinante um denselben Faktor λ , insbesondere gilt det $(\lambda A) = \lambda^n \det (A)$ für $A \in \mathbb{R}^{n,n}$.
- $(\mathbf{D}_4) \det (\mathbf{A}) = \det (\mathbf{A}^{\mathsf{T}}).$
- $(D_5) \det (AB) = \det (A) \det (B).$
- $(D_6) \det (A) \neq 0$ genau dann, wenn A regulär ist.
- $(\mathbf{D}_7) \det (\boldsymbol{R}) = \det (\boldsymbol{R}_{11}) \det (\boldsymbol{R}_{22}) \cdots \det (\boldsymbol{R}_{ss})$

für jede obere Blockdreiecksmatrix

$$\boldsymbol{R} = \begin{pmatrix} \boldsymbol{R}_{11} & \boldsymbol{R}_{12} & \dots & \boldsymbol{R}_{1s} \\ & \boldsymbol{R}_{22} & \dots & \boldsymbol{R}_{2s} \\ & & \ddots & & \\ & & & \ddots & \\ & & & & \cdot & \boldsymbol{R}_{ss} \end{pmatrix}$$

mit quadratischen Diagonalblöcken $R_{11}, R_{22}, ..., R_{ss}$; insbesondere ist die Determinante einer Dreiecksmatrix gleich dem Produkt ihrer Diagonalelemente.

H. Orthogonalität von Vektoren und orthogonale Matrizen

Zwei Vektoren $p, q \in \mathbb{R}^m$ heißen orthogonal, wenn

$$\boldsymbol{p}^{\mathsf{T}}\boldsymbol{q}=0$$

gilt. Offensichtlich ist der Nullvektor orthogonal zu jedem Vektor, und für nichtverschwindende Vektoren stimmt der Orthogonalitätsbegriff im Fall m = 2, 3 mit dem aus der analytischen Geometrie bekannten überein. Allgemeiner heißt ein System $\{p^1, ..., p^n\}$ von Vektoren $p^1, ..., p^n \in \mathbb{R}^m$ orthogonal, wenn

$$oldsymbol{p^i}^{\intercal}oldsymbol{p^j}=0 \qquad (i \pm j; \, i, \, j=1, \, ..., \, n)$$

ist. Gilt zusätzlich $p^{i^{\intercal}}p^{i} = 1$ (i = 1, ..., n), so heißt das System orthonormal. Nichtverschwindende orthogonale, also insbesondere orthonormale Vektoren sind stets linear unabhängig. Ein Beispiel eines orthonormalen Systems bilden die Einheitsvektoren $\{e^{1}, ..., e^{m}\}$. Eine Matrix $Q \in \mathbb{R}^{m,n}$, $m \ge n$, heißt spaltenorthogonal bzw. spaltenorthonormal, wenn die Spalten $q^{1}, ..., q^{n}$ von Q ein orthogonales bzw. orthonormales System bilden. Dabei ist Q spaltenorthonormal genau dann, wenn

$$\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} = \boldsymbol{I}_{\mathsf{n}} \tag{27}$$

gilt. Entsprechend heißt $P \in \mathbb{R}^{m,n}$, $m \leq n$, zeilenorthogonal bzw. zeilenorthonormal, wenn P^{\intercal} spaltenorthogonal bzw. spaltenorthonormal ist. Zeilenorthonormale Matrizen P sind durch die Bedingung

$$\boldsymbol{P}\boldsymbol{P}^{\mathsf{T}} = \boldsymbol{I}_{\boldsymbol{m}} \tag{28}$$

charakterisiert.

ŀ

Eine quadratische spaltenorthonormale Matrix $Q \in \mathbb{R}^{n,n}$ heißt schlechthin orthogonal. Sie ist dann auch zeilenorthonormal und umgekehrt, d. h., eine orthogonale Matrix ist durch die gleichwertigen Bedingungen

$$Q^{\mathsf{T}}Q = I$$
 bzw. $QQ^{\mathsf{T}} = I$ (29)

charakterisiert. Jede orthogonale Matrix ist wegen (29) regulär, und es gilt

$$\boldsymbol{Q}^{-1} = \boldsymbol{Q}^{\mathsf{T}}.\tag{30}$$

Aus (29) folgt unter Beachtung der Determinanteneigenschaften (D₅), (D₄) und (D₇) sofort $1 = \det(I) = \det(Q) \det(Q^{T}) = [\det(Q)]^2$, also

$$|\det\left(\boldsymbol{Q}\right)| = 1 \tag{31}$$

für jedes orthogonale Q.

Das Produkt orthogonaler Matrizen ist wieder orthogonal. Wir bemerken abschließend, daß eine Matrix $Q \in \mathbb{R}^{n,n}$ genau dann orthogonal ist, wenn für alle $u, v \in \mathbb{R}^n$

$$(\mathbf{Q}\mathbf{u})^{\mathsf{T}} \mathbf{Q}\mathbf{v} = \mathbf{u}^{\mathsf{T}}\mathbf{v} \tag{32}$$

gilt, d. h., wenn die durch Q vermittelte Abbildung y = Qx das Skalarprodukt zweier Vektoren invariant läßt.

I. Vektor- und Matrixnormen

Für jedes $\boldsymbol{x} \in \mathbf{R}^n$ ist $\boldsymbol{x}^\intercal \boldsymbol{x} = \sum_{i=1}^n (x_i)^2 \ge 0$, so daß die Größe

$$\|\boldsymbol{x}\|_{2} := \sqrt{\boldsymbol{x}^{\mathsf{T}}\boldsymbol{x}} = \sqrt{\sum_{i=1}^{n} |x_{i}|^{2}}$$
(33)

stets definiert ist. Durch (33) wird jedem $x \in \mathbb{R}^n$ die nichtnegative Zahl $||x||_2$ zugeordnet, die *Euklidische Norm* von x. Sie gibt für n = 2 und n = 3 die aus der analytischen Geometrie bekannte Länge des zu x gehörenden Ortsvektors wieder, siehe Abb. 1.1.1.



Die Euklidische Norm hat die Eigenschaften

- $(\mathbf{N}_1) \|\boldsymbol{x}\|_2 \geq 0, \|\boldsymbol{x}\|_2 = 0 \text{ genau für } \boldsymbol{x} = \boldsymbol{o},$
- (N₂) $\|\lambda x\|_2 = |\lambda| \|x\|_2$,
- (N₃) $\|\boldsymbol{x} + \boldsymbol{y}\|_2 \le \|\boldsymbol{x}\|_2 + \|\boldsymbol{y}\|_2$

für alle $x, y \in \mathbb{R}^n$ und alle $\lambda \in \mathbb{R}$. Die Eigenschaft (N₃) wird *Dreiecksungleichung* genannt; sie ist gleichwertig mit

 $(\mathbf{N}_3') |||\mathbf{x}||_2 - ||\mathbf{y}||_2| \leq ||\mathbf{x} - \mathbf{y}||_2$ für alle $\mathbf{x}, \mathbf{y} \in \mathbf{R}^n$.

Ihre Gültigkeit folgt aus der sog. Schwarzschen Ungleichung

 $|\boldsymbol{x}^{\mathsf{T}}\boldsymbol{y}| \leq \|\boldsymbol{x}\|_{2} \|\boldsymbol{y}\|_{2} \quad \text{für alle } \boldsymbol{x}, \boldsymbol{y} \in \mathbf{R}^{n},$ (34)

durch die Skalarprodukt und Euklidische Norm verknüpft sind. Dabei gilt in (34) das Gleichheitszeichen genau dann, wenn x und y linear abhängig sind. Man kann elementar nachrechnen, daß im Fall n = 2 und n = 3 der cos des durch zwei nichtverschwindende Vektoren x, y gebildeten Winkels $\varphi = \measuredangle (x, y)$ durch

$$\cos\varphi = \frac{\boldsymbol{x}^{\mathsf{T}}\boldsymbol{y}}{\|\boldsymbol{x}\|_{2} \|\boldsymbol{y}\|_{2}} \tag{35}$$



Abb. 1.1.2. Winkel zweier Vektoren

ausgedrückt werden kann, siehe Abb. 1.1.2. Da die rechte Seite von (35) für beliebiges *n* wegen (34) betragsmäßig höchstens gleich 1 ist, kann durch (35) auch für n > 3 zwischen x und y ein Winkel φ mit $0 \le \varphi \le \pi$ definiert werden.

Wir kehren jetzt wieder zur Euklidischen Norm zurück. Obwohl diese ein natürliches Maß für die "Größe" eines Vektors ist, können auch andere Vorschriften angegeben werden, die sich ebenfalls in dieser Weise interpretieren lassen und wie $||\boldsymbol{x}||_2$ die Eigenschaften (N_1) , (N_2) , (N_3) erfüllen. Beispiele dafür sind die

Summennorm

$$\|\boldsymbol{x}\|_{1} := \sum_{i=1}^{n} |x_{i}| \tag{36}$$

und die Maximumnorm

$$\|\boldsymbol{x}\|_{\infty} := \max_{i=1,\dots,n} |x_i|.$$
(37)

Der Unterschied zwischen den angegebenen Normen wird deutlich, wenn man sich die jeweiligen sog. *Einheitskugeln* $\mathscr{S} = \{x \in \mathbb{R}^n : ||x|| \leq 1\}$ etwa für n = 2 veranschaulicht, siehe Abb. 1.1.3.



Abb. 1.1.3. Einheitskugeln für verschiedene Vektornormen im Fall n = 2

Jede der eingeführten Vektornormen läßt sich durch die übrigen abschätzen, siehe Ü1.1.8.

Wir betrachten jetzt die durch eine Matrix $A \in \mathbb{R}^{m,n}$ vermittelte Abbildung y = Ax und fragen, um welchen Faktor sich $||x||_p$ beim Übergang zu $||y||_p = ||Ax||_p$ maximal verändert, d. h., wir suchen den maximalen Wert von $||Ax||_p/||x||_p$, wenn x alle nichtverschwindenden Vektoren aus \mathbb{R}^n durchläuft. Dabei bezeichnet $||x||_p$

bzw. $\|y\|_p$ eine Norm des Typs $p \in \{1, 2, \infty\}$ im Urbildraum \mathbb{R}^n bzw. im Bildraum \mathbb{R}^m . Der gesuchte Wert

$$\|\boldsymbol{A}\|_{p} := \max\left\{\frac{\|\boldsymbol{A}\boldsymbol{x}\|_{p}}{\|\boldsymbol{x}\|_{p}} : \boldsymbol{x} \neq \boldsymbol{o}\right\} = \max\left\{\|\boldsymbol{A}\boldsymbol{x}\|_{p} : \|\boldsymbol{x}\|_{p} \leq 1\right\}$$
(38)

existiert für jede Matrix $A \in \mathbb{R}^{m,n}$. Er genügt den Normeigenschaften (N_1) , (N_2) , (N_3) und ist daher eine Norm auf $\mathbb{R}^{m,n}$, die die der Vektornorm $||x||_p$ zugeordnete Matrixnorm genannt wird. Aus (38) folgt sofort die Gültigkeit der Abschätzung

 $\|A\boldsymbol{x}\|_{p} \leq \|A\|_{p} \|\boldsymbol{x}\|_{p} \quad \text{für alle } \boldsymbol{x} \in \mathbf{R}^{n}.$ (39)

Für mindestens ein $x \neq o$ steht in (39) das Gleichheitszeichen, so daß die Zahl $||A||_p$ dort nicht durch eine kleinere ersetzt werden kann.

Für die Indizes p = 1 und $p = \infty$ lassen sich die den Vektornormen $||x||_1$ und $||x||_{\infty}$ zugeordneten Matrixnormen in einfacher Weise aus den Elementen von A gemäß

$$\|\boldsymbol{A}\|_{1} = \max_{j=1,\dots,n} \sum_{i=1}^{m} |a_{ij}|$$
(40)

bzw.

$$\|A\|_{\infty} = \max_{i=1,\dots,m} \sum_{j=1}^{n} |a_{ij}|$$
(41)

berechnen. Diese Formeln legen die Benennungen Spaltensummennorm für $||A||_1$ und Zeilensummennorm für $||A||_{\infty}$ nahe. Dagegen kann $||A||_2$ nicht in solch einfacher Weise bestimmt werden. Wir werden im Abschnitt 1.2 zeigen, daß

$$\|\boldsymbol{A}\|_{2} = \sqrt{\lambda_{\max}(\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A})} \tag{42}$$

gilt, wobei $\lambda_{\max}(A^{\mathsf{T}}A)$ den größten Eigenwert der symmetrischen und positiv semidefiniten Matrix $A^{\mathsf{T}}A$ bezeichnet, der stets nichtnegativ ist. Die Norm $||A||_2$ heißt Spektralnorm von A. Die in Analogie zur Euklidischen Vektornorm aus den Matrixelementen gebildete Norm

$$\|A\|_{F} := \sqrt{\sum_{i=1}^{m} \sum_{j=1}^{n} |a_{ij}|^{2}}$$
(43)

ist natürlich auch eine Matrixnorm – d. h. eine Norm auf $\mathbb{R}^{m,n}$ – mit den Eigenschaften (N₁), (N₂), (N₃); sie wird *Frobeniusnorm* genannt. Für diese gilt

$$\|A\|_{F}/\sqrt{r} \le \|A\|_{2} \le \|A\|_{F} \tag{44}$$

mit $r = \operatorname{rang} (A) \leq \min \{m, n\}$. Insbesondere kann daher in der Ungleichung (39) für p = 2 die Zahl $||A||_2$ durch die leicht berechenbare obere Schranke $||A||_F$ ersetzt werden: Für alle $x \in \mathbb{R}^n$ gilt

$$\|Ax\|_{2} \leq \|A\|_{F} \|x\|_{2}. \tag{45}$$

Allerdings gibt es in (45) i. allg. kein $x \neq o$, so daß das Gleichheitszeichen steht.

Für den Fall einer Diagonalmatrix $D = \text{diag}(d_1, ..., d_l) \in \mathbb{R}^{m,n}$ mit $l = \min\{m, n\}$ gilt

$$\|\boldsymbol{D}\|_{1} = \|\boldsymbol{D}\|_{2} = \|\boldsymbol{D}\|_{\infty} = \max_{j=1,\dots,l} |d_{j}|, \qquad (46)$$

aber

$$\|\boldsymbol{D}\|_F = \sqrt{\sum_{j=1}^{l} (d_j)^2}.$$
 (47)

Aus der in (32) festgestellten Invarianz des Skalarproduktes folgt eine wichtige Eigenschaft der mit dem Index 2 bzw. F gekennzeichneten Vektor- bzw. Matrixnormen, nämlich deren Invarianz unter orthogonalen Transformationen: Sind $Q \in \mathbf{R}^{n,n}, P \in \mathbf{R}^{m,m}$ orthogonale Matrizen, so gilt

$$\|Qx\|_2 = \|x\|_2 \tag{48}$$

für alle $x \in \mathbf{R}^n$ und

$$\|PAQ\|_{2} = \|A\|_{2}, \qquad \|PAQ\|_{F} = \|A\|_{F}$$
(49)

für alle $A \in \mathbb{R}^{m,n}$. Man beachte, daß die Einheitsmatrix auch orthogonal ist, so daß in (49) einer der Faktoren P oder Q auch fehlen darf. Insbesondere folgt aus (47), (48) und (49)

$$\|Q\|_2 = 1, \qquad \|Q\|_F = \sqrt{n}$$
 (50)

für jedes orthogonale Q.

Alle bisher betrachteten Matrixnormen $||A||_p$, $p \in \{1, 2, \infty, F\}$, haben die Eigenschaft

$$|a_{ij}| \leq ||A|| \qquad (i = 1, ..., m; j = 1, ..., n).$$
 (51)

Für die Indizes 1, ∞ und F ist dies offensichtlich, für die Spektralnorm folgt (51) mittels (34) aus $|a_{ij}| = |e^{i\intercal} (Ae^j)| \leq ||e^i||_2 ||Ae^j||_2 \leq ||e^i||_2 ||A||_2 ||e^j||_2 = ||A||_2$. Überdies gilt für die erwähnten Normen stets

$$\|\boldsymbol{A}\boldsymbol{B}\| \le \|\boldsymbol{A}\| \, \|\boldsymbol{B}\|,\tag{52}$$

sofern das Produkt **AB** gebildet werden kann. Im Fall $||A||_p$, $p \in \{1, 2, \infty\}$, ist dies eine Folge der Eigenschaft (39), denn für jedes x gilt

$$\|(AB) x\|_p = \|A(Bx)\|_p \le \|A\|_p \|Bx\|_p \le \|A\|_p \|B\|_p \|x\|_p$$

so da $\beta \|A\|_p \|B\|_p$ in der Tat eine obere Schranke für $\|AB\|_p$ ist. Für die Frobeniusnorm ergibt sich (52) unter Beachtung von (44) aus der Ungleichung

$$\|AB\|_{F} \leq \|A\|_{2} \|B\|_{F}. \tag{53}$$

Zum Nachweis von (53) setzen wir $C := AB = (c^1, ..., c^l), B = (b^1, ..., b^l)$. Dann gilt $c^j = Ab^j$ (j = 1, ..., l), wobei c^j bzw. b^j die Spalten von C bzw. B bezeichnen. Wegen (39) folgt hieraus $\|c^j\|_2 \leq \|A\|_2 \|b^j\|_2$, also

$$\|C\|_F^2 = \sum_{j=1}^l (\|c^j\|_2)^2 \le \|A\|_2^2 \sum_{j=1}^l (\|b^j\|_2)^2 = \|A\|_2^2 \|B\|_F^2$$

wie behauptet. Die Eigenschaft (52) wird *Multiplikativität* der Matrixnorm ||A|| genannt.

J. Ordnungen und Absolutbeträge

Bei der Analyse von Algorithmen der linearen Algebra treten häufig elementweise Abschätzungen von Vektoren und Matrizen auf. Um solche Abschätzungen einfach beschreiben zu können, führen wir im \mathbf{R}^m eine *natürliche Halbordnung* genannte Ordnungsrelation " \leq " durch die Festlegung

$$\boldsymbol{x} \leq \boldsymbol{y}, \quad \text{wenn } x_i \leq y_i \quad (i = 1, ..., m)$$

$$(54)$$

für $x, y \in \mathbb{R}^m$ ein. Analog wird für $A, B \in \mathbb{R}^{m,n}$

$$A \leq B$$
, wenn $a_{ij} \leq b_{ij}$ $(i = 1, ..., m; j = 1, ..., n)$ (55)

gesetzt. Im Gegensatz zur Situation bei reellen Zahlen gibt es Vektoren x, y, für die weder $x \ge y$ noch $y \ge x$ gilt, so daß die Benennung Halbordnung berechtigt ist. In entsprechender Weise wird für $x \in \mathbb{R}^m$ ein *Betragsvektor* $|x| \in \mathbb{R}^m$ als Vektor mit den Komponenten

$$(|\boldsymbol{x}|)_i = |x_i| \qquad (i = 1, ..., m)$$
 (56)

und für $A \in \mathbf{R}^{m,n}$ eine *Betragsmatrix* $|A| \in \mathbf{R}^{m,n}$ als Matrix mit den Elementen

$$(|A|)_{ij} = |a_{ij}|$$
 $(i = 1, ..., m; j = 1, ..., n)$ (57)

eingeführt. Ordnungsbeziehungen und Betragsbildung genügen dabei denselben Beziehungen wie im Bereich der reellen Zahlen, insbesondere dürfen Ungleichungen addiert und mit nichtnegativen Faktoren multipliziert werden, und die Beträge erfüllen die Beziehungen

$$|\boldsymbol{x}| = 0$$
 genau für $\boldsymbol{x} = \boldsymbol{o}, \quad |\lambda \boldsymbol{x}| = |\lambda| |\boldsymbol{x}|$ sowie $|\boldsymbol{x} + \boldsymbol{y}| \le |\boldsymbol{x}| + |\boldsymbol{y}|$
(58)

für alle $x, y \in \mathbb{R}^m$ und $\lambda \in \mathbb{R}$. Dieselben Beziehungen gelten auch, wenn x, y durch Matrizen $A, B \in \mathbb{R}^{m,n}$ ersetzt werden. Falls das Produkt AB gebildet werden kann, ist überdies

$$|\boldsymbol{A}\boldsymbol{B}| \leq |\boldsymbol{A}| |\boldsymbol{B}|, \tag{59}$$

insbesondere gilt

$$|A\boldsymbol{x}| \leq |A| \, |\boldsymbol{x}| \tag{60}$$

für alle $A \in \mathbb{R}^{m,n}$ und $x \in \mathbb{R}^n$.

Die eingeführten Ordnungsbeziehungen und Beträge sind mit den Normen in dem Sinne verträglich, daß

$$\||\boldsymbol{x}|\| = \|\boldsymbol{x}\| \tag{61}$$

sowie

$$\operatorname{aus} |\boldsymbol{x}| \leq |\boldsymbol{y}| \text{ folgt } \|\boldsymbol{x}\| \leq \|\boldsymbol{y}\| \tag{62}$$

gilt, wobei ||x|| eine der Vektornormen $||x||_p$, $p \in \{1, 2, \infty\}$, ist.

Eine Norm mit der Eigenschaft (61) bzw. (62) heißt absolut bzw. monoton. Es ist sofort zu sehen, daß die analogen Eigenschaften

$$|||A||| = ||A|| \tag{63}$$

bzw.

$$aus |\mathbf{A}| \leq |\mathbf{B}| \text{ folgt } \|\mathbf{A}\| \leq \|\mathbf{B}\| \tag{64}$$

auch für die Matrixnormen $||\mathbf{A}||_1$, $||\mathbf{A}||_{\infty}$ und $||\mathbf{A}||_F$ gelten, während $||\mathbf{A}||_2$ weder absolut noch monoton ist, siehe Ü 1.1.11. Es gilt jedoch

$$\|A\|_{2} \leq \||A|\|_{2} \tag{65}$$

für jedes $A \in \mathbf{R}^{m,n}$ sowie

aus
$$|A| \leq |B|$$
 folgt $|||A|||_2 \leq |||B|||_2$. (66)

Übungsaufgaben

Ü 1.1.1. Für $A \in \mathbb{R}^{m,n}$, $x \in \mathbb{R}^n$ kann $y = Ax \in \mathbb{R}^m$ wie folgt berechnet werden:

oder

 $\begin{array}{l} V_2: \mbox{ for } i:=1(1)m \mbox{ do } y_i:=0 \\ \mbox{ for } j:=1(1)n \mbox{ do } \\ \mbox{ for } i:=1(1)m \mbox{ do } y_i:=y_i+a_{ii}*x_i. \end{array}$

Man mache sich die Unterschiede zwischen beiden Varianten klar und überlege sich, daß V_2 z. B. dann vorzuziehen ist, wenn A spaltenweise auf einem Magnetband abgespeichert ist.

U1.1.2. Es sei $S \in \mathbb{R}^{m,n}$ eine reguläre Matrix und $X \subset \mathbb{R}^n$ ein Teilraum von \mathbb{R}^n . Man beweise, daß dann $SX := \{y = Sx : x \in X\}$ auch ein Teilraum von \mathbb{R}^n ist und daß dim $(SX) = \dim(X)$ gilt. Außerdem ist $\{a^1, \ldots, a^r\}$ eine Basis von X genau dann, wenn $\{Sa^1, \ldots, Sa^r\}$ eine Basis von SX ist.

Ü 1.1.3. Es sei $A := ab^{\intercal}$ das dyadische Produkt der Vektoren $a, b \in \mathsf{R}^n, a \neq o, b \neq o$. Man zeige

(i) $\mathcal{R}(A) = \text{span } \{a\}, \text{ rang } (A) = 1 \text{ und } \mathcal{N}(A) = \{x \in \mathbb{R}^n : b^{\mathsf{T}}x = 0\}, \text{ dim } \mathcal{N}(A) = n - 1.$

(ii) $A = A^{\dagger}$ genau dann, wenn a und b linear abhängig sind, d. h., wenn $a = \lambda b$ ist mit einer Zahl $\lambda = 0$.

(iii) $||A||_2 = ||A||_F = ||a||_2 ||b||_2$, $||A||_{\infty} = ||a||_{\infty} ||b||_1$, $||A||_1 = ||a||_1 ||b||_{\infty}$.

Ü 1.1.4. Es sei $A = ab^{\mathsf{T}}$ wie in **Ü** 1.1.3 erklärt. Dann kann $y = Ax = ab^{\mathsf{T}}x$ nach den beiden Varianten $y = (ab^{\mathsf{T}}) x$ bzw. $y = a(b^{\mathsf{T}}x)$ berechnet werden. Man vergleiche beide Vorschriften hinsichtlich der Anzahl der erforderlichen Rechenoperationen und des – auch für Zwischenergebnisse – benötigten Speicherplatzes.

Ü 1.1.5. Man beweise: (i) Jede Matrix $A \in \mathbb{R}^{n,n}$ kann in der Form A = B + C als Summe der symmetrischen Matrix $B = (A + A^{\mathsf{T}})/2$ und der schiefsymmetrischen Matrix $C = (A - A^{\mathsf{T}})/2$ dargestellt werden. Dabei heißt $C \in \mathbb{R}^{n,n}$ schiefsymmetrisch, wenn $C = -C^{\mathsf{T}}$ gilt. Die Matrix B wird symmetrischer Anteil von A genannt.

(ii) Die quadratische Form einer schiefsymmetrischen Matrix C verschwindet identisch, d. h., es gilt $\mathbf{x}^{\mathsf{T}}C\mathbf{x} = 0$ für alle $\mathbf{x} \in \mathbf{R}^n$. Hieraus folgt $\mathbf{x}^{\mathsf{T}}A\mathbf{x} = \mathbf{x}^{\mathsf{T}}B\mathbf{x}$ für alle $\mathbf{x} \in \mathbf{R}^n$, d. h., die quadratische Form einer beliebigen Matrix wird allein durch deren symmetrischen Anteil bestimmt. **U** 1.1.6. Es sei $A = (a_{ij}) \in \mathbb{R}^{n,n}$ eine symmetrische und positiv definite Matrix. Man beweise: (i) $a_{ii} > 0$ (i = 1, ..., n),

(ii)
$$|a_{ij}| < (a_{ii} + a_{jj})/2$$
 (*i*, *j* = 1, ..., *n*; *i* + *j*),

(iii) $|a_{ij}| < \sqrt{a_{ii}a_{jj}}$ (*i*, *j* = 1, ..., *n*; *i* = *j*).

Hinweis: Man betrachte $x^{T}Ax$ für $x = \xi e^{i} + \eta e^{j}$, wobei $\xi = 1$, $\eta = 0$ ist für (i), $\xi = 1$, $\eta = \pm 1$ für (ii) und ξ , η beliebig für (iii).

Ü 1.1.7. Es ist zu zeigen, daß die Berechnung von det (A) für $A \in \mathbb{R}^{n,n}$ nach dem Entwicklungssatz n! Multiplikationen und (n-1)! Additionen erfordert.

Ü 1.1.8. Man beweise, daß für die Vektornormen $||x||_p$, $p \in \{1, 2, \infty\}$, die Ungleichungen

$$\|\boldsymbol{x}\|_{\infty} \leq \|\boldsymbol{x}\|_{1} \leq n \|\boldsymbol{x}\|_{\infty}, \qquad \|\boldsymbol{x}\|_{\infty} \leq \|\boldsymbol{x}\|_{2} \leq \sqrt{n} \|\boldsymbol{x}\|_{\infty}, \qquad \|\boldsymbol{x}\|_{2} \leq \|\boldsymbol{x}\|_{1} \leq \sqrt{n} \|\boldsymbol{x}\|_{2}$$

für jedes $x \in \mathbb{R}^n$ gelten, und zeige, daß die hier auftretenden Konstanten nicht verbessert werden können.

Ü 1.1.9. Man zeige, daß für jedes $x \in \mathbf{R}^n$

$$\|m{x}\|_2^2 \leq \|m{x}\|_1 \, \|m{x}\|_\infty \leq 0.5 ig(
ot \! n+1 ig) \, \|m{x}\|_2^2$$

gilt.

Ü 1.1.10. Man beweise: Wenn die Matrix $A \in \mathbf{R}^{m,n}$ die Blockstruktur

$$A = \left(\frac{A_{11} \mid A_{12}}{A_{21} \mid A_{22}} \right)$$

besitzt, gilt

$$||A_{ij}||_p \leq ||A||_p$$
 (*i*, *j* = 1, 2)

für jede der Matrixnormen $p \in \{1, 2, \infty, F\}$.

Ü1.1.11. Am Beispiel der Matrizen

 $A = egin{pmatrix} 0.6 & 0.6 \ 0.6 & 0.6 \end{pmatrix}, \quad B = rac{1}{\sqrt{2}} egin{pmatrix} 1 & 1 \ -1 & 1 \end{pmatrix}$

ist zu zeigen, daß $||A||_2$ weder eine absolute noch eine monotone Norm ist.

Hinweis: Man beachte, daß **B** orthogonal ist, und stelle **A** in der Form $\mathbf{A} = 0.6ee^{\intercal}$ mit $\mathbf{e} = (1, 1)^{\intercal}$ dar. Die Norm $||\mathbf{A}||_2$ kann dann nach Ü 1.1.3 berechnet werden. Analog ist $|||\mathbf{B}|||_2$ bestimmbar.

Ü 1.1.12. Man beweise, daß für jede orthogonale obere Dreiecksmatrix A

|A| = I

gilt, d. h., A ist diagonal mit Diagonalelementen ± 1 .

1.2. Eigenwerte, Singulärwerte

A. Eigenwerte und charakteristisches Polynom

Eine Zahl $\lambda \in \mathbb{R}$ wird *Eigenwert* der quadratischen Matrix $A \in \mathbb{R}^{n,n}$ genannt, wenn es einen Vektor $x \in \mathbb{R}^n$, $x \neq o$, gibt, für den

$$A\boldsymbol{x} = \lambda \boldsymbol{x} \tag{1}$$

gilt. Jeder derartige Vektor x heißt ein zu λ gehörender *Eigenvektor*; wir nennen $\{\lambda, x\}$ auch ein *Eigenpaar* der Matrix A. Die Eigenvektoren sind also diejenigen Vektoren, die durch A in ein Vielfaches von sich selbst abgebildet werden, und der Proportionalitätsfaktor stellt gerade den Eigenwert dar. Die Bedingung $x \pm o$ ist dabei erforderlich, denn für x = o ist (1) trivialerweise für beliebiges λ erfüllt. Mit x ist auch jedes nichtverschwindende Vielfache Eigenvektor von A.

Die Gleichung (1) läßt sich äquivalent in der Form

$$(\boldsymbol{A} - \lambda \boldsymbol{I}) \boldsymbol{x} = \boldsymbol{o} \tag{2}$$

schreiben, d. h. als homogenes Gleichungssystem $M(\lambda) x = o$ mit der von λ abhängenden Koeffizientenmatrix $M(\lambda) := A - \lambda I$. Die Eigenwerte von A sind diejenigen λ , für die (2) von o verschiedene Lösungen besitzt, für die also $A - \lambda I$ singulär ist. Unter Beachtung der Determinanteneigenschaft (D₆) aus Abschnitt 1.1.G ergibt sich hieraus

$$\det\left(\boldsymbol{A}-\boldsymbol{\lambda}\boldsymbol{I}\right)=0\tag{3}$$

als notwendige und hinreichende Bedingung dafür, daß λ ein Eigenwert von A ist. Für n = 1 bedeutet dies

$$\det (a_{11} - \lambda) = (-\lambda) + a_{11} = 0,$$

für n = 2 geht (3) in

$$\det \begin{pmatrix} a_{11} - \lambda & a_{12} \\ a_{21} & a_{22} - \lambda \end{pmatrix} = (-\lambda)^2 - (a_{11} + a_{22}) \lambda + a_{11}a_{22} - a_{12}a_{21} = 0$$

über, und allgemein erhält man

$$\det \begin{pmatrix} a_{11} - \lambda & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} - \lambda \end{pmatrix}$$

= $(-1)^n \lambda^n + c_{n-1} \lambda^{n-1} + \dots + c_1 \lambda + c_0 = 0$

mit gewissen reellen Koeffizienten c_0, \ldots, c_{n-1} , die Produktsummen der Elemente von A sind. Der Ausdruck det $(A - \lambda I) =: p(\lambda)$ stellt ein Polynom vom Grad nmit dem Führungskoeffizienten $(-1)^n$ dar und wird *charakteristisches Polynom* der Matrix A genannt. Die Eigenwerte von A im Sinne der obigen Definition sind also genau die reellen Nullstellen des charakteristischen Polynoms.

Obwohl das charakteristische Polynom einer reellen Matrix nur reelle Koeffizienten hat, kann es auch komplexe Nullstellen besitzen, wie das Beispiel

$$oldsymbol{A} = egin{pmatrix} 2 & 1 \ -2 & 4 \end{pmatrix}, \quad \det \left(oldsymbol{A} - \lambda oldsymbol{I}
ight) = \lambda^2 - 6\lambda + 10 = 0$$

mit den Nullstellen $\lambda_{1,2} = 3 \pm \sqrt{-1} = 3 \pm i$ zeigt. Es ist daher angebracht, auch komplexe Eigenwerte $\lambda = \alpha + i\beta$ zuzulassen. Für solche hat (2) die Gestalt

$$M(\lambda) \boldsymbol{x} := (\boldsymbol{A} - \lambda \boldsymbol{I}) \boldsymbol{x} = [(\boldsymbol{A} - \alpha \boldsymbol{I}) - \mathrm{i}\beta \boldsymbol{I}] \boldsymbol{x} = \boldsymbol{o},$$

d. h., die Koeffizientenmatrix $M(\lambda)$ wird ebenfalls komplex, und es müssen dann auch komplexe Eigenvektoren x = u + iv als Lösungen von (2) betrachtet werden. Wir erinnern an dieser Stelle daran, daß jedes λ aus der Menge **C** der komplexen Zahlen die Gestalt $\lambda = \alpha + i\beta$ mit $\alpha, \beta \in \mathbf{R}$ hat. Dabei ist α der Realteil, β der Imaginärteil von λ , und i bezeichnet die imaginäre Einheit, für die $i^2 = -1$ gilt. Die Addition bzw. Subtraktion komplexer Zahlen geschieht durch Addition bzw. Subtraktion der Real- und Imaginärteile, während die Multiplikation durch

$$\lambda_1 \lambda_2 = (\alpha_1 + \mathrm{i}\beta_1) (\alpha_2 + \mathrm{i}\beta_2) = (\alpha_1 \alpha_2 - \beta_1 \beta_2) + \mathrm{i}(\alpha_1 \beta_2 + \alpha_2 \beta_1) \tag{4}$$

erklärt ist. Die zu $\lambda = \alpha + i\beta$ konjugiert komplexe Zahl ist $\overline{\lambda} = \alpha - i\beta$, und $|\lambda| = \sqrt{\lambda \overline{\lambda}} = \sqrt{\alpha^2 + \beta^2} \ge 0$ ist der absolute Betrag von λ bzw. $\overline{\lambda}$. Im Fall $\lambda_2 \neq 0$ wird schließlich die komplexe Division gemäß

$$\frac{\lambda_1}{\lambda_2} = \frac{\lambda_1 \bar{\lambda}_2}{\lambda_2 \bar{\lambda}_2} = \frac{\alpha_1 \alpha_2 + \beta_1 \beta_2}{\alpha_2^2 + \beta_2^2} + i \frac{-\alpha_1 \beta_2 + \alpha_2 \beta_1}{\alpha_2^2 + \beta_2^2}$$
(5)

eingeführt. Analog zum Vorgehen im Abschnitt 1.1 wird dann der *m*-dimensionale komplexe Vektorraum \mathbb{C}^m als Menge aller *m*-Tupel $\boldsymbol{x} = \boldsymbol{u} + i\boldsymbol{v}$ mit $\boldsymbol{u}, \boldsymbol{v} \in \mathbb{R}^m$ definiert, in dem eine Addition und eine Multiplikation mit einem komplexen Skalar $\lambda \in \mathbb{C}$ komponentenweise erklärt ist. Entsprechend ist $\mathbb{C}^{m,n}$ die Menge aller komplexen (m,n)-Matrizen $\boldsymbol{A} = \boldsymbol{U} + i\boldsymbol{V}$ mit $\boldsymbol{U}, \boldsymbol{V} \in \mathbb{R}^{m,n}$, wobei wieder alle Operationen im komplexen Sinne zu verstehen sind. Für $\boldsymbol{A} \in \mathbb{C}^{m,n}$ ist $\boldsymbol{A}^{\mathsf{H}} \in \mathbb{C}^{n,m}$ diejenige Matrix, die durch Transposition und Übergang zum Konjugiertkomplexen aus \boldsymbol{A} entsteht, d. h. $\boldsymbol{A}^{\mathsf{H}} := \bar{\boldsymbol{A}}^{\mathsf{T}}$. Das Skalarprodukt von $\boldsymbol{x}, \boldsymbol{y} \in \mathbb{C}^m$ wird als komplexe Zahl $\boldsymbol{y}^{\mathsf{H}}\boldsymbol{x}$ $= \sum_{j} \bar{\boldsymbol{y}}_j \boldsymbol{x}_j$ eingeführt. Dann gilt $\boldsymbol{y}^{\mathsf{H}}\boldsymbol{x} = \bar{\boldsymbol{x}}^{\mathsf{H}}\boldsymbol{y}$, d. h., das Skalarprodukt ist nicht mehr kommutativ. Jedoch ist $\boldsymbol{x}^{\mathsf{H}}\boldsymbol{x} = \sum_{j} |\boldsymbol{x}_j|^2$ stets reell und nichtnegativ, so daß $||\boldsymbol{x}||_2$ $= \sqrt{\boldsymbol{x}^{\mathsf{H}}\boldsymbol{x}}$ für alle $\boldsymbol{x} \in \mathbb{C}^m$ erklärt ist und die Normeigenschaften besitzt, wobei (N_2) als $||\lambda x|| = |\lambda| \cdot ||\boldsymbol{x}||$ für alle $\boldsymbol{x} \in \mathbb{C}^m$ und alle $\lambda \in \mathbb{C}$ mit dem oben eingeführten Absolutbetrag $|\lambda|$ zu verstehen ist. Indem die Absolutbeträge überall durch ihre komplexe Version und der obere Index ,, $\mathbb{T}^{"}$ durch , $\mathbb{H}^{"}$ ersetzt werden, übertragen sich die Begriffe und Aussagen des Abschnittes 1.1 sinngemäß auf den komplexen Fall.

Nach diesen Vorbereitungen können wir nun auch eine komplexe Matrix $A \in \mathbb{C}^{n,n}$ betrachten und in Erweiterung der eingangs gegebenen Definitionen komplexe Eigenwerte und Eigenvektoren zulassen. Die Eigenwerte von A sind dann genau die reellen und komplexen Nullstellen des charakteristischen Polynoms. Da die Nullstellen eines Polynoms stetige Funktionen der Koeffizienten, die Koeffizienten des charakteristischen Polynoms wieder stetige Funktionen der Matrixelemente sind, ergibt sich hieraus: Die Eigenwerte einer Matrix sind stetige Funktionen der Matrixelemente. Für reelle Matrizen $A \in \mathbb{R}^{n,n} \subset \mathbb{C}^{n,n}$ treten komplexe Eigenwerte stets als Paare konjugiert komplexer Eigenwerte auf, und die Eigenvektoren können entsprechend gewählt werden, siehe Ü 1.2.2.

Ein Eigenwert $\lambda_j \in \mathbf{C}$ von A kann u. U. mehrfache Nullstelle des charakteristischen Polynoms $p(\lambda)$ sein. Wir sagen, daß λ_j die algebraische Vielfachheit $\alpha_j \geq 1$ besitzt, wenn λ_i eine α_i -fache Nullstelle von $p(\lambda)$ ist, d. h., wenn

$$p(\lambda_j)=p'(\lambda_j)=\dots=p^{(lpha_j-1)}(\lambda_j)=0\,,\qquad p^{(lpha_j)}(\lambda_j)=0$$

gilt. Jede Matrix $A \in \mathbb{C}^{n,n}$ hat also genau *n* entsprechend ihren algebraischen Vieltachheiten gezählte Eigenwerte.

Die Anzahl γ_j der zu einem Eigenwert λ_j gehörenden linear unabhängigen Eigenvektoren wird geometrische Vielfachheit von λ_j genannt, d. h. $\gamma_j := \dim \mathcal{N}(\boldsymbol{A} - \lambda_j \boldsymbol{I})$. Dabei gilt stets

$$1 \leq \gamma_j \leq \alpha_j \leq n, \tag{6}$$

und es kann Eigenwerte λ_j mit $\gamma_j < \alpha_j$ geben. Solche Eigenwerte sind notwendig algebraisch mehrfach. Ein Beispiel gibt die Matrix

$$oldsymbol{A} = egin{pmatrix} 3 & 1 \ -1 & 1 \end{pmatrix} \quad ext{mit} \quad ext{det} \left(oldsymbol{A} - \lambda oldsymbol{I}
ight) = (\lambda - 2)^2.$$

Sie besitzt den algebraischen zweifachen Eigenwert $\lambda_1 = 2$, aber die zugehörigen Eigenvektoren als Lösungen von (A - 2I) x = o bilden nur den eindimensionalen, durch $x = (1, -1)^{\mathsf{T}}$ aufgespannten Teilraum. Eine Matrix A, die mindestens einen Eigenwert λ_i mit $\gamma_i < \alpha_i$ besitzt, heißt *defektiv*.

Eigenvektoren, die zu verschiedenen Eigenwerten gehören, sind stets linear unabhängig. Die Anzahl der linear unabhängigen Eigenvektoren einer Matrix ist gleich der Summe der geometrischen Vielfachheiten ihrer Eigenwerte.

Wir bemerken abschließend, daß A und A^{\intercal} dieselben Eigenwerte besitzen, denn wegen der Determinanteneigenschaft (D₄) haben A und A^{\intercal} dasselbe charakteristische Polynom. Die Eigenvektoren sind jedoch i. allg. verschieden.

B. *Ähnlichkeitstransformationen und Eigenwertzerlegung*

Wir wollen das Verhalten der Eigenwertgleichung (1) gegenüber Transformationen

$$\boldsymbol{x} \to \boldsymbol{y} = \boldsymbol{T}^{-1} \boldsymbol{x}, \qquad \boldsymbol{T} \in \mathbf{C}^{n,n} ext{ regular}, ext{(7)}$$

untersuchen, d. h., wir ersetzen x durch Ty. Einsetzen von (7) in (1) führt auf das transformierte Problem

$$By = \lambda y \quad \text{mit} \quad B := T^{-1}AT. \tag{8}$$

1.2.1. Aussage. Die Matrix $A \in \mathbb{C}^{n,n}$ werde mittels der regulären Matrix $T \in \mathbb{C}^{n,n}$ in $B = T^{-1}AT$ transformiert. Dann haben A und B dieselben Eigenwerte, und sowohl ihre algebraischen als auch ihre geometrischen Vielfachheiten stimmen überein. Dabei ist x Eigenvektor von A zum Eigenwert λ genau dann, wenn $y = T^{-1}x$ Eigenvektor von B zum selben Eigenwert ist.

Beweis. Aus $\boldsymbol{B} - \lambda \boldsymbol{I} = \boldsymbol{T}^{-1}\boldsymbol{A}\boldsymbol{T} - \lambda \boldsymbol{T}^{-1}\boldsymbol{T} = \boldsymbol{T}^{-1}(\boldsymbol{A} - \lambda \boldsymbol{I}) \boldsymbol{T}$ folgt wegen (D₅) sofort

$$\det (\boldsymbol{B} - \lambda \boldsymbol{I}) = \det (\boldsymbol{T}^{-1}) \det (\boldsymbol{A} - \lambda \boldsymbol{I}) \det (\boldsymbol{T}).$$

Analog ergibt sich aus $T^{-1}T = I$ die Beziehung det (T^{-1}) det $(T) = \det(I) = 1$, was auf det $(B - \lambda I) = \det(A - \lambda I)$, also die Gleichheit der charakteristischen Polynome von B und

A führt. Daher haben **B** und **A** dieselben Eigenwerte einschließlich ihrer algebraischen Vielfachheiten. Die Aussage über die Eigenvektoren folgt aus der Äquivalenz von (1) und (8) und zieht $\mathcal{N}(\boldsymbol{B} - \lambda \boldsymbol{I}) = \boldsymbol{T}^{-1}\mathcal{N}(\boldsymbol{A} - \lambda \boldsymbol{I})$ nach sich. Nach Ü 1.1.2 haben dann $\mathcal{N}(\boldsymbol{B} - \lambda \boldsymbol{I})$ und $\mathcal{N}(\mathbf{A} - \lambda \mathbf{I})$ dieselbe Dimension, d. h., auch die geometrischen Vielfachheiten sind dieselben.

Die Transformation $A \rightarrow B = T^{-1}AT$ mittels einer regulären Matrix T wird \ddot{A} hnlichkeitstransformation genannt; A und B heißen \ddot{a} hnlich. Aussage 1.2.1 bedeutet: Die Eigenwerte einer Matrix einschließlich ihrer algebraischen und geometrischen Vielfachheiten sind invariant unter Ähnlichkeitstransformationen.

Für theoretische wie praktische Zwecke sind diejenigen Ähnlichkeitstransformationen von besonderem Interesse, die A in eine Matrix B überführen, deren Eigenwerte sofort abgelesen bzw. einfach berechnet werden können. Als Zielmatrizen mit dieser Eigenschaft bieten sich speziell Blockdreiecksmatrizen an.

1.2.2. Aussage. Die Eigenwerte der Blockdreiecksmatrix

$$\boldsymbol{R} = \begin{pmatrix} \boldsymbol{R}_{11} & \boldsymbol{R}_{12} & \dots & \boldsymbol{R}_{1s} \\ & \boldsymbol{R}_{22} & \dots & \boldsymbol{R}_{2s} \\ & & \ddots & \vdots \\ & \boldsymbol{O} & & & \boldsymbol{R}_{ss} \end{pmatrix} \in \boldsymbol{C}^{n,n}$$

mit quadratischen Diagonalblöcken \mathbf{R}_{jj} der Dimension n_j , $n_1 + \cdots + n_s = n$, sind gerade die Eigenwerte aller dieser Diagonalblöcke.

Beweis. Mit der Determinanteneigenschaft (D₇) erhält man

 $\det \left(\boldsymbol{R} - \lambda \boldsymbol{I}\right) = \det \left(\boldsymbol{R}_{11} - \lambda \boldsymbol{I}\right) \det \left(\boldsymbol{R}_{22} - \lambda \boldsymbol{I}\right) \cdots \det \left(\boldsymbol{R}_{ss} - \lambda \boldsymbol{I}\right),$

denn $\mathbf{R} = \lambda \mathbf{I}$ ist eine Blockdreiecksmatrix derselben Struktur wie \mathbf{R} .

Aus 1.2.2 ergibt sich als wichtiger Sonderfall: Die Eigenwerte einer Dreiecksmatrix und speziell die einer Diagonalmatrix sind gleich den Diagonalelementen.

Besonders einfach in bezug auf die Struktur ihrer Eigenwerte und Eigenvektoren sind reelle symmetrische Matrizen.

1.2.3. Aussage. Für jede reelle symmetrische Matrix A gilt:

- \boldsymbol{A} besitzt nur reelle Eigenwerte, und die Eigenvektoren können reell gewählt werden.
- (ii) Eigenvektoren zu verschiedenen Eigenwerten sind orthogonal.
 (iii) Algebraische und geometrische Vielfachheiten der Eigenwerte stimmen überein, d. h., zu einem Eigenwert λ_i der algebraischen Vielfachheit α_i gehören auch α_i linear unabhängige Eigenvektoren.

Beweis. (i) Es sei $\lambda \in \mathbf{C}$ ein Eigenwert und $x \in \mathbf{C}^n$ ein zugehöriger Eigenvektor, d. h., es gelte $y := Ax = \lambda x$. Durch Linksmultiplikation mit x^{H} folgt $x^{\mathsf{H}}y = \lambda x^{\mathsf{H}}x$, also $\lambda = x^{\mathsf{H}}y/x^{\mathsf{H}}x$. Nun ist $x^{\mathsf{H}}x$ reell und positiv, und für den Zähler gilt $\overline{x^{\mathsf{H}}y} = x^{\mathsf{T}}\overline{y} = \overline{y}^{\mathsf{T}}x = y^{\mathsf{H}}x = (\overline{Ax})^{\mathsf{T}}x$ $= \bar{x}^{\mathsf{T}} \bar{A}^{\mathsf{T}} x = x^{\mathsf{H}} A x = x^{\mathsf{H}} y$, da A reell und symmetrisch ist. Hieraus folgt $\bar{\lambda} = \overline{x^{\mathsf{H}} y} / \overline{x^{\mathsf{H}} x}$ $= x^{H}y/x^{H}x = \lambda$, d. h., λ ist reell. Die Eigenvektoren als Lösungen des reellen Systems $(A - \lambda I) x = o$ können dann stets reell gewählt werden.

(ii) Es seien $\lambda_i = \lambda_j$ zwei Eigenwerte von A mit zugehörigen reellen Eigenvektoren x^i, x^j . Aus $Ax^i = \lambda_i x^i$ folgt $(x^j)^{\mathsf{T}} Ax^i = \lambda_i (x^j)^{\mathsf{T}} x^i$, und ebenso erhält man $(x^i)^{\mathsf{T}} Ax^j = \lambda_i (x^i)^{\mathsf{T}} x^j$.

Subtraktion beider Gleichungen liefert unter Beachtung der Symmetrie von A die Beziehung $0 = (\lambda_i - \lambda_j) [x^{i \intercal} x^j]$. Wegen $\lambda_i \neq \lambda_j$ zieht dies $x^{i \intercal} x^j = 0$, also die Orthogonalität der Eigenvektoren nach sich.

Auf den Beweis von (iii) soll an dieser Stelle verzichtet werden; ein konstruktiver Beweis wird später in Abschnitt 13.2 mittels des Jacobi-Verfahrens gegeben werden.

Aus (iii) folgt, daß im Fall eines Eigenwertes λ_i der algebraischen Vielfachheit α_i jeweils beliebige α_i linear unabhängige Vektoren aus dem α_i -dimensionalen Teilraum $\mathcal{N}(\mathbf{A} - \lambda_i \mathbf{I})$ als Eigenvektoren zu λ_i genommen werden können. Insbesondere können diese Vektoren als orthogonale Basis von $\mathcal{N}(\boldsymbol{A} - \lambda_i \boldsymbol{I})$ ausgewählt werden. Es gibt daher stets n orthogonale Eigenvektoren u^1, \ldots, u^n zu den Eigenwerten $\lambda_1, \ldots, \lambda_n$, von denen gewisse entsprechend ihrer Vielfachheit mehrfach auftreten können. Ohne Einschränkung der Allgemeinheit können die Eigenvektoren als normiert bezüglich $||\mathbf{x}||_2$ vorausgesetzt werden. Die Gleichungen $A\mathbf{u}^j = \lambda_j \mathbf{u}^j$ (j = 1, ..., n)lassen sich dann mit der orthogonalen Matrix $U := (u^1, ..., u^n)$ und der Diagonalmatrix $\Lambda := \text{diag}(\lambda_1, \ldots, \lambda_n)$ in Matrixform $AU = U\Lambda$ schreiben. Durch Linksmultiplikation mit U^{-1} folgt schließlich $U^{-1}AU = U^{\dagger}AU = \Lambda$. Damit ist der folgende Satz bewiesen:

1.2.4. Eigenwertzerlegung reeller symmetrischer Matrizen. Zu jeder reellen symmetrischen Matrix $A \in \mathbb{R}^{n,n}$ gibt es eine orthogonale Matrix $U \in \mathbb{R}^{n,n}$ mit

$$\boldsymbol{U}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{U} = \boldsymbol{\Lambda} = \operatorname{diag}\left(\lambda_1, \dots, \lambda_n\right) \in \mathbf{R}^{n,n},\tag{9}$$

wobei $\lambda_1, \ldots, \lambda_n$ die Eigenwerte von A sind.

Satz 1.2.4 besagt: Jede reelle symmetrische Matrix kann orthogonal ähnlich auf reelle Diagonalform transformiert werden. In der nach A aufgelösten Form lautet (9)

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{A}\boldsymbol{U}^{\mathsf{T}} = (\boldsymbol{u}^{1}, ..., \boldsymbol{u}^{n}) \begin{pmatrix} \lambda_{1}\boldsymbol{u}^{1\mathsf{T}} \\ \vdots \\ \lambda_{n}\boldsymbol{u}^{n\mathsf{T}} \end{pmatrix} = \sum_{j=1}^{n} \lambda_{j}\boldsymbol{u}^{j}\boldsymbol{u}^{j\mathsf{T}}, \qquad (10)$$

was die Benennung Eigenwertzerlegung motiviert.

Bei nichtsymmetrischen Matrizen ist die Situation weitaus komplizierter, da dann sowohl komplexe Eigenpaare als auch weniger als n linear unabhängige Eigenvektoren auftreten können. Wir schließen zunächst den letzten Fall aus, d. h., wir betrachten eine nichtdefektive Matrix A, und $x^1, \ldots, x^n \in C^n$ seien linear unabhängige Eigenvektoren zu den Eigenwerten $\lambda_1, \ldots, \lambda_n \in \mathbf{C}$. Mit der regulären Matrix $X := (x^1, \ldots, x^n)$ $\in \mathbf{C}^{n,n}$ und $\Lambda := \operatorname{diag}(\lambda_1, \ldots, \lambda_n) \in \mathbf{C}^{n,n}$ folgt dann $AX = X\Lambda$, also $X^{-1}AX = \Lambda$.

1.2.5. Aussage. Zu jeder nichtdefektiven Matrix $A \in \mathbb{R}^{n,n}$ gibt es eine reguläre Matrix $X \in \mathbf{C}^{n,n}$ mit

$$\boldsymbol{X}^{-1}\boldsymbol{A}\boldsymbol{X} = \boldsymbol{\Lambda} = \operatorname{diag}\left(\lambda_1, \dots, \lambda_n\right) \in \mathbf{C}^{n,n}.$$
(11)

Falls A nur reelle Eigenwerte hat, kann X reell gewählt werden.

Aussage 1.2.5 begründet, warum nichtdefektive Matrizen — d. h. Matrizen, bei denen die algebraische Vielfachheit der Eigenwerte mit der geometrischen übereinstimmt — auch diagonalähnlich genannt werden. Speziell sind Matrizen mit n verschiedenen Eigenwerten wegen (6) diagonalähnlich.

Der Fall defektiver, also nichtdiagonalähnlicher Matrizen wird durch den folgenden Satz beschrieben:

1.2.6. Satz. Zu jeder Matrix $A \in \mathbb{R}^{n,n}$ gibt es eine reguläre Matrix $T \in \mathbb{C}^{n,n}$, so daß

$$\boldsymbol{T}^{-1}\boldsymbol{A}\boldsymbol{T} = \boldsymbol{J} = \begin{pmatrix} \boldsymbol{J}(\lambda_1, n_1) & & \\ & \boldsymbol{J}(\lambda_2, n_2) & & \\ & & \ddots & \\ & & & \boldsymbol{J}(\lambda_s, n_s) \end{pmatrix}$$
(12)

mit Diagonalblöcken

$$oldsymbol{J}(\lambda_j,\,n_j)\,=\,egin{pmatrix} \lambda_j & 1 & 0\ \lambda_j & 1\ & \ddots\ & \ddots\ 0 & & \ddots\ & 1\ 0 & & \lambda_j \end{bmatrix}\in oldsymbol{C}^{n_j,n_j}$$

und $n_1 + \cdots + n_s = n$ gilt. Falls *A* nur reelle Eigenwerte hat, kann *T* reell gewählt werden.

Die Matrix J wird Jordansche Normalform von A genannt; die $J(\lambda_j, n_j)$ heißen Jordanblöcke, für $n_j = 1$ ist $J(\lambda_j, 1) = (\lambda_j)$. Die Matrix J ist bis auf die Anordnung der Diagonalblöcke eindeutig durch A bestimmt. Die Anzahl s der Diagonalblöcke ist gleich der Anzahl der linear unabhängigen Eigenvektoren von A. Diese stehen in denjenigen Spalten von T, die den ersten Spalten der Jordanblöcke' entsprechen, vgl. Ü 1.2.3. Sie werden durch sog. Hauptvektoren zu einer Basis von \mathbb{C}^n ergänzt. Eine diagonalähnliche Matrix hat genau n Jordanblöcke der Ordnung 1.

Ein Beispiel für die Jordansche Normalform einer Matrix ist



Die Matrix J und damit auch A besitzt den algebraisch sechsfachen Eigenwert $\lambda = 4$ mit der geometrischen Vielfachheit 3 sowie den algebraisch doppelten Eigenwert $\lambda = 2$ mit der geometrischen Vielfachheit 1. Die drei linear unabhängigen Eigenvektoren zu $\lambda_1 = 4$ stehen in den Spalten 1, 3 und 4 von T, der zu $\lambda_2 = 2$ gehörende in Spalte 7. Das Beispiel weist darauf hin, daß mehrere Jordanblöcke zum selben Eigenwert gehören können.
In Ü 1.2.4 wird gezeigt, daß die Einsen in der Nebendiagonale der Jordanschen Normalform durch beliebige Zahlen $\varepsilon > 0$ ersetzt werden können.

Die Ähnlichkeitstransformation aus 1.2.5 und 1.2.6 zielen wie die aus 1.2.4 auf eine Diagonal- bzw. Fast-Diagonalform hin, allerdings ist dann die Transformationsmatrix nicht mehr orthogonal. Ein alternatives, für die Algorithmenentwicklung günstigeres Vorgehen besteht darin, als Zielmatrix lediglich eine Dreiecksform anzustreben, dafür aber mit orthogonalen bzw. unitären Transformationen als deren komplexe Entsprechung zu arbeiten. Eine Matrix $U \in \mathbb{C}^{n,n}$ heißt dabei *unitär*, wenn

 $U^{\mathsf{H}}U = I$ bzw. äquivalent $UU^{\mathsf{H}} = I$

gilt, und für solche Matrizen ist

 $U^{\mathsf{H}} = U^{-1},$

vgl. (1.1.29) und (1.1.30).

1.2.7. Satz von Schur. Zu jeder Matrix $A \in \mathbb{R}^{n,n}$ gibt es eine unitäre Matrix $U \in \mathbb{C}^{n,n}$, so da β

$$\boldsymbol{U}^{\mathsf{H}}\boldsymbol{A}\boldsymbol{U} = \boldsymbol{R} \tag{13}$$

mit einer oberen Dreiecksmatrix $R \in \mathbb{C}^{n,n}$ gilt. Falls A nur reelle Eigenwerte besitzt, kann U reell, d. h. orthogonal gewählt werden.

Jede Matrix läßt sich also unitär ähnlich auf obere Dreiecksform transformieren, und die Diagonalelemente von \mathbf{R} sind gerade die Eigenwerte von \mathbf{A} .

Für die Anwendungen hat Satz 1.2.7 den Nachteil, daß im Fall einer Matrix mit komplexen Eigenwerten sowohl \mathbf{R} als auch U komplex werden und gegebenenfalls in komplexer Arithmetik berechnet werden müssen. Andererseits treten für reelles \mathbf{A} stets Paare von konjugiert komplexen Eigenwerten $\alpha \pm i\beta$ auf, die selbst Eigenwerte von reellen (2, 2)-Matrizen sind, siehe Ü 1.2.2. Es liegt daher nahe, auf das explizite Auftreten der komplexen Eigenwerte von \mathbf{A} in der Diagonalen von \mathbf{R} zu verzichten und dort stattdessen reelle Zweierblöcke anzusiedeln, deren zwei Eigenwerte gerade ein Paar konjugiert komplexer oder zwei reelle Eigenwerte von \mathbf{A} sind. Dann kann die gesamte Transformation sogar im Reellen durchgeführt werden, wie der folgende Satz als Modifikation von 1.2.7 zeigt.

1.2.8. Satz. Zu jeder Matrix $A \in \mathbf{R}^{n,n}$ gibt es eine orthogonale Matrix $U \in \mathbf{R}^{n,n}$ mit

$$U^{\mathsf{T}}AU = R, \tag{14}$$

wobe
i ${\pmb R}$ eine obere Blockdreiecksmatrix

$$m{R} = egin{pmatrix} m{R}_{11} & m{R}_{12} & \dots & m{R}_{1s} \ m{R}_{22} & \dots & m{R}_{2s} \ & & \ddots & \vdots \ & & & m{R}_{ss} \end{pmatrix} \in m{R}^{n,n}$$

mit quadratischen Diagonalblöcken $\mathbf{R}_{ij} \in \mathbf{R}^{n_j, n_j}$ der Dimension $n_j = 1$ oder $n_j = 2$ ist.

Nach 1.2.2 sind die Eigenwerte von \mathbf{R} gerade die Eigenwerte der Diagonalblöcke \mathbf{R}_{jj} . Für Blöcke der Ordnung 1 sind dies die Diagonalelemente selbst, für (2,2)-Blöcke ergeben sie sich aus dem zugehörigen charakteristischen Polynom, d. h. aus Wurzeln einer quadratischen Gleichung. Die Sätze 1.2.7 und 1.2.8 bilden die theoretische Grundlage für den QR-Algorithmus zur Eigenwertbestimmung nichtsymmetrischer Matrizen, auf den wir allerdings in diesem Buch nicht eingehen.

Wir geben abschließend zu diesem Komplex noch einige einfach berechenbare Schranken für die Eigenwerte einer Matrix an.

1.2.9. Aussage. Für jeden Eigenwert $\lambda_i \in \mathbf{C}$ einer Matrix $A \in \mathbf{R}^{n,n}$ gilt

$$|\lambda_j| \le \|A\|,\tag{15}$$

wobei ||A|| irgendeine multiplikative Matrixnorm, speziell also eine der Normen $||A||_p$, $p \in \{1, 2, \infty, F\}$, sei.

Beweis. Es sei $\{\lambda, x\}$ ein Eigenpaar von A. Aus $Ax = \lambda x$ folgt dann auf Grund der vorausgesetzten Multiplikativität der Matrixnorm sofort $|\lambda| ||x|| = ||\lambda x|| = ||Ax|| \le ||A|| ||x||$, wegen $x \neq o$ also die Behauptung. \Box

Bessere Schranken liefert in der Regel der folgende Satz.

1.2.10. Kreisesatz von Geršgorin. Jeder Eigenwert $\lambda_j \in \mathbf{C}$ einer Matrix $A \in \mathbf{R}^{n,n}$ liegt in mindestens einem der Kreise

$$\mathscr{K}_i := \left\{ \lambda \in \mathbf{C} \colon |\lambda - a_{ii}| \leq \sum_{\substack{j=1\\i\neq i}}^n |a_{ij}| = :r_i \right\} \quad (i = 1, ..., n).$$
(16)

Wir bemerken zunächst, daß \mathcal{K}_i einen Kreis der komplexen Ebene mit dem Mittelpunkt a_{ii} und dem Radius r_i darstellt, und beweisen jetzt die Aussage des Satzes. Es sei dazu $\{\lambda, x\}$ ein Eigenpaar von A. Dann wird i so gewählt, daß $|x_i| = ||x||_{\infty}$ gilt. Aus $(Ax)_i = \lambda x_i$ ergibt sich dann sofort $(Ax)_i - a_{ii}x_i = (\lambda - a_{ii}) x_i$, also

$$|\lambda - a_{ii}| |x_i| = |(Ax)_i - a_{ii}x_i| \le \sum_{\substack{j=1\j \neq i}}^n |a_{ij}x_j| = \sum_{\substack{j=1\j \neq i}}^n |a_{ij}| |x_j| \le \left\{\sum_{\substack{j=1\j \neq i}}^n |a_{ij}|
ight\} |x_i|,$$

d. h. $\lambda \in \mathcal{K}_i$. \Box

1.2.11. Verschärfung. Hat die Vereinigung von r Kreisen \mathcal{K}_i keine gemeinsamen Punkte mit den restlichen Kreisen, so liegen in ihr genau r Eigenwerte von A. Insbesondere enthält jeder Kreis \mathcal{K}_i genau einen Eigenwert, wenn alle Kreise disjunkt sind.

Dies ergibt sich aus 1.2.10, indem die Nichtdiagonalelemente stetig und monoton von 0 auf ihren vorgegebenen Wert transformiert werden und die stetige Abhängigkeit der Eigenwerte von den Matrixelementen berücksichtigt wird. □

C. Kongruenztransformationen und quadratische Formen

Wir betrachten jetzt das Verhalten der durch die symmetrische Matrix $A \in S^{n,n}$ erzeugten quadratischen Form $x^{T}Ax$ gegenüber der Transformation (7). Einsetzen von $\boldsymbol{x} = \boldsymbol{T} \boldsymbol{y}$ führt auf

$$\boldsymbol{x}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x} = \boldsymbol{y}^{\mathsf{T}} \boldsymbol{B} \boldsymbol{y} \quad \text{mit} \quad \boldsymbol{B} := \boldsymbol{T}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{T}.$$
 (17)

Die Transformation $A \rightarrow B = T^{\mathsf{T}}AT$ mittels einer regulären Matrix T wird Kongruenztransformation genannt; A und B heißen kongruent. Dabei bleibt die Symmetrie der Matrix erhalten. Wir suchen wieder solche Transformationen, für die B und damit die zugehörige Form $y^{\mathsf{T}}By$ möglichst einfache Gestalt annehmen. Dies ist der Fall, wenn B diagonal ist. Hierfür bietet sich die Ähnlichkeitstransformation (9) an, die wegen der Orthogonalität der Transformationsmatrix U gleichzeitig eine Kongruenztransformation ist.

1.2.12. Aussage. Es sei $U^{\mathsf{T}}AU = \Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_n)$ mit $U^{\mathsf{T}}U = I$ die Eigenwertzerlegung der reellen symmetrischen Matrix $A \in S^{n,n}$. Dann gilt für jedes $x \in \mathbb{R}^n$

$$\boldsymbol{x}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x} = \boldsymbol{y}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{y} = \sum_{j=1}^{n} \lambda_j (\boldsymbol{y}_j)^2 \tag{18}$$

mit $\boldsymbol{y} := \boldsymbol{U}^{\mathsf{T}} \boldsymbol{x}$.

Die hier gegebene Transformation der quadratischen Form auf eine Summe von Quadraten mit den Eigenwerten λ_i von A als Koeffizienten heißt Hauptachsentransformation. Aus 1.2.12 folgt insbesondere: A ist positiv definit bzw. semidefinit genau dann, wenn alle Eigenwerte positiv bzw. nichtnegativ sind.

Wir werden später im Abschnitt 6.1.A sehen, daß es auch einfachere, allerdings nichtorthogonale Kongruenztransformationen von A auf Diagonalform $\Delta = \text{diag}(\delta_j)$ gibt. Die Zahlen δ_j sind dann nicht mehr die Eigenwerte von A, aber ihre Vorzeichenverteilung ist dieselbe wie die der λ_j .

1.2.13. Trägheitsgesetz für quadratische Formen. Für die reelle symmetrische Matrix $A \in S^{n,n}$ und reguläres $T \in \mathbb{R}^{n,n}$ gelte

 $T^{\mathsf{T}}AT = \varDelta = \operatorname{diag}(\delta_i).$

Dann ist die Anzahl der positiven, negativen und verschwindenden Diagonalelemente von \varDelta unabhängig von der Transformationsmatrix T und identisch mit der Anzahl der positiven, negativen und verschwindenden Eigenwerte von A.

Satz 1.2.13 läßt sich wie folgt lesen: Die Vorzeichenverteilung der Eigenwerte reeller symmetrischer Matrizen ist invariant unter Kongruenztransformationen.

D. Äquivalenztransformationen und Singulärwertzerlegung

Wir betrachten abschließend eine beliebige, nicht notwendig quadratische Matrix $A \in \mathbf{R}^{m,n}$ und wollen untersuchen, wie sich die durch A erzeugte Abbildung

$$\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} \tag{19}$$

gegenüber Transformationen verhält. Da Urbild- und Bildraum hier i. allg. verschieden sind, müssen x und y unterschiedlichen Transformationen unterworfen werden,

d. h., wir betrachten Transformationen der Form

$$x \to x' = S^{-1}x, \ S \in \mathbb{R}^{n,n}$$
 regulär, $y \to y' = T^{-1}y, \ T \in \mathbb{R}^{m,m}$ regulär.
$$(20)$$

Einsetzen von (20) in (19) führt auf

$$\mathbf{y}' = \mathbf{B}\mathbf{x}' \quad \text{mit} \quad \mathbf{B} := \mathbf{T}^{-1}\mathbf{A}\mathbf{S} \in \mathbf{R}^{m,m}. \tag{21}$$

Die Transformation $A \rightarrow B = T^{-1}AS$ mit regulären Matrizen S, T wird Äquivalenztransformation genannt; A und B heißen äquivalent. Dabei gilt

$$\mathscr{N}(\mathbf{B}) = \mathbf{S}^{-1}\mathscr{N}(\mathbf{A}), \qquad \mathscr{R}(\mathbf{B}) = \mathbf{T}^{-1}\mathscr{R}(\mathbf{A}). \tag{22}$$

Hieraus folgt: Der Rang einer Matrix ist invariant unter Äquivalenztransformationen; man beachte Ü 1.1.2.

Wie bei den bisher betrachteten Transformationen ist es auch hier möglich, die Matrix A auf Diagonalform zu transformieren, und dies sogar mittels orthogonaler Matrizen.

1.2.14. Singulärwertzerlegung. Zu jeder Matrix $A \in \mathbb{R}^{m,n}$ gibt es orthogonale Matrizen $U \in \mathbb{R}^{m,m}$ und $V \in \mathbb{R}^{n,n}$, so daß

$$\boldsymbol{U}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{V} = \boldsymbol{\Sigma} = \text{diag}\;(\boldsymbol{\sigma}_1,\;..,\boldsymbol{\sigma}_l) \in \mathbf{R}^{m,n} \tag{23}$$

wit $l := \min \{m, n\}$ und $\sigma_1 \ge \sigma_2 \ge \cdots \ge \sigma_r > 0, \quad \sigma_{r+1} = \cdots = \sigma_l = 0$ gilt. Die Zahlen $\sigma_1, \ldots, \sigma_l$ sind durch A eindeutig festgelegt und heißen Singulärwerte von A.

Die Diagonalmatrix Σ ist also von der Gestalt

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1 & & \\ & \cdot & \\ & & \cdot & \\ & & \cdot & \\ & & & \sigma_r \\ \hline & & & \sigma_r \\ \hline & & & 0 \end{pmatrix} \in \mathbf{R}^{m,n};$$

im Fall r = l = n fehlen die rechten, im Fall r = l = m die unteren Nullblöcke in Σ.

Bevor wir 1.2.14 beweisen, sollen einige Folgerungen angegeben werden:

Unter Beachtung von (22) ergibt sich sofort

$$\mathcal{R}(A) = U\mathcal{R}(\Sigma), \qquad \mathcal{N}(A) = V\mathcal{N}(\Sigma),$$
(24)

insbesondere ist also rang $(A) = rang (\Sigma) = r$, d. h., die Anzahl r der positiven Singulärwerte ist gleich dem Rang der Matrix.

Auflösen von (23) nach A liefert

$$A = U\Sigma V^{\mathsf{T}} = \sum_{j=1}^{r} \sigma_j u^j v^{j\mathsf{T}}, \qquad (25)$$

was die Benennung Singulärwertzerlegung motiviert. Die Spalten von V bzw. U werden gelegentlich *Singulärvektoren* genannt; für diese gilt

$$Av^{j} = \sigma_{i}u^{j} \qquad (j = 1, ..., l).$$

Einsetzen von (25) liefert

$$(A^{\mathsf{T}}A) V = V(\Sigma^{\mathsf{T}}\Sigma) = V \operatorname{diag} (\sigma_1^2, \dots, \sigma_r^2, 0, \dots, 0),$$
(26)

und ebenso folgt

$$(AA^{\mathsf{T}}) U = U(\Sigma\Sigma^{\mathsf{T}}) = U \operatorname{diag} (\sigma_1^2, \dots, \sigma_r^2, 0, \dots, 0).$$
(27)

Die Zahlen $\sigma_1^2, \ldots, \sigma_r^2$ sind also die positiven Eigenwerte sowohl von $A^{\intercal}A$ als auch von AA^{\intercal} , und die restlichen n - r bzw. m - r Eigenwerte von $A^{\intercal}A$ bzw. AA^{\intercal} sind 0. Insbesondere gilt

$$r = \operatorname{rang} \left(A \right) = \operatorname{rang} \left(A^{\mathsf{T}} A \right) = \operatorname{rang} \left(A A^{\mathsf{T}} \right), \tag{28}$$

und $A^{\mathsf{T}}A$ wie AA^{T} sind positiv semidefinit. Die Spalten von V sind die Eigenvektoren von $A^{\mathsf{T}}A$, die Spalten von U die Eigenvektoren von AA^{T} .

Wir kommen jetzt zum Beweis von 1.2.14. Die Matrix $A^{\mathsf{T}}A$ ist symmetrisch und wegen $\mathbf{x}^{\mathsf{T}}A^{\mathsf{T}}A\mathbf{x} = ||A\mathbf{x}||_2^2 \ge 0$ positiv semidefinit. Nach 1.2.4 existiert dann ein orthogonales $V = (v^1, \ldots, v^n) \in \mathbb{R}^{n,n}$ mit $V^{\mathsf{T}}(A^{\mathsf{T}}A) V = A = \text{diag}(\lambda_i)$ und $\lambda_i \ge 0$. Ohne Einschränkung der Allgemeinheit kann $\lambda_1 \ge \cdots \ge \lambda_r > 0$, $\lambda_{r+1} = \cdots = \lambda_n = 0$ mit einem gewissen $r \in \{0, \ldots, n\}$ angenommen werden. Aus $A^{\mathsf{T}}Av^j = \lambda_j v^j$ folgt $\lambda_j = v^{\mathsf{T}}A^{\mathsf{T}}Av^j = ||Av^j||_2^2$, also $Av^j \neq o$ für $j = 1, \ldots, r$ und $Av^j = o$ für $j = r + 1, \ldots, n$. Mit den Festlegungen $\sigma_j := \sqrt{\lambda_j}$ und $u^j := Av^j/\sigma_j$ $(j = 1, \ldots, r)$ ergibt sich dann für $i, j = 1, \ldots, r$

$$\boldsymbol{u}^{i\mathsf{T}}\boldsymbol{u}^{j} = \frac{1}{\sigma_{i}\sigma_{j}}\boldsymbol{v}^{i\mathsf{T}}\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{v}^{j} = \frac{\lambda_{h}}{\sigma_{i}\sigma_{j}}\boldsymbol{v}^{i\mathsf{T}}\boldsymbol{v}^{j} = \begin{cases} 1 & \text{für } i=j, \\ 0 & \text{für } i\neq j. \end{cases}$$

Die Vektoren $\{u^1, \ldots, u^r\}$ bilden daher ein orthonormiertes System und können durch Hinzunahme weiterer Vektoren u^{r+1}, \ldots, u^m zu einer orthonormalen Basis von \mathbb{R}^m ergänzt werden. Dann ist $U = (u^1, \ldots, u^m)$ orthogonal. Mit der Festlegung $\sigma_{r+1} = \cdots = \sigma_l = 0$ folgt wegen

$$(\boldsymbol{U}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{V})_{ij} = \boldsymbol{u}^{i\mathsf{T}}\boldsymbol{A}\boldsymbol{v}^{j} = \begin{cases} \sigma_{j}\boldsymbol{u}^{i\mathsf{T}}\boldsymbol{u}^{j} & \text{für } j = 1, ..., r, \\ 0 & \text{für } j = r+1, ..., n \end{cases}$$

die gesuchte Beziehung $U^{\mathsf{T}}AV = \Sigma$.

Die in Abschnitt 1.1 eingeführten Matrixnormen $||A||_2$ und $||A||_F$ können in einfacher Weise durch die Singulärwerte ausgedrückt werden.

1.2.15. Aussage. Für jedes $A \in \mathbb{R}^{m,n}$ gilt

$$\|A\|_{2} = \sigma_{1}, \qquad \|A\|_{F} = \sqrt{(\sigma_{1})^{2} + \dots + (\sigma_{r})^{2}}.$$
 (29)

Beweis. Wegen der Orthogonalinvarianz beider Normen ist $||A|| = ||\Sigma||$. Aus (1.1.46) und (1.1.47) ergibt sich dann die Behauptung.

Der erste Teil von (29) ist wegen (26) überdies identisch mit (1.1.42), womit diese Formel nachträglich bewiesen ist; außerdem folgt (1.1.44) ebenfalls aus (29).

Abschließend geben wir einige einfache Beziehungen zwischen Eigen- und Singulärwerten quadratischer Matrizen an. **1.2.16.** Aussage. Die quadratische Matrix $A \in \mathbb{R}^{n,n}$ besitze die Singulärwerte $\sigma_1 \geq \cdots \geq \sigma_n$ und die gemäß $|\lambda_1| \geq \cdots \geq |\lambda_n|$ geordneten Eigenwerte $\lambda_1, \ldots, \lambda_n \in \mathbb{C}$. Dann gelten folgende Beziehungen:

(i)
$$\det(A) = \prod_{j=1}^{n} \lambda_{j}, \quad |\det(A)| = \prod_{j=1}^{n} \sigma_{j}.$$
 (30)

(ii)
$$\sigma_n \leq |\lambda_j| \leq \sigma_1$$
 $(j = 1, ..., n).$ (31)

(iii)
$$\sigma_j = |\lambda_j|$$
 $(j = 1, ..., n),$ falls $A = A^{\intercal}$ ist. (32)

Beweis. (i) Unter Beachtung von (1.1.31) folgt der erste Teil direkt aus (13), während sich der zweite analog aus (23) ergibt.

(ii) Die rechte Seite der Ungleichung ist der Sonderfall von 1.2.9 für $||A||_2$. Die linke Seite ist im Fall $\sigma_n = 0$ trivial. Es sei also $\sigma_n > 0$, d. h., A sei regulär. Dann ist A^{-1} vorhanden und hat die Singulärwerte $1/\sigma_j$ sowie die Eigenwerte $1/\lambda_j$. Anwendung der eben bewiesenen Ungleichung auf A^{-1} liefert dann die linke Seite.

(iii) Aus (26) folgt für symmetrisches A unmittelbar $A^2 = A^{\mathsf{T}}A = V\Sigma^2 V^{\mathsf{T}}$, d. h. $\lambda_j^2 = \sigma_j^2$ und damit $|\lambda_j| = \sigma_j$. Man beachte dabei, daß A^2 die Eigenwerte λ_j^2 hat. \Box

Für nichtsymmetrisches A können die Singulärwerte beliebig stark von den Eigenwerten abweichen, vgl. Ü 1.2.8.

Übungsaufgaben

Ü 1.2.1. Man zeige, daß für A = U + iV mit $U, V \in \mathbb{R}^{m,n}$, a = b + ic mit $b, c \in \mathbb{R}^m$ und x = u + iv mit $u, v \in \mathbb{R}^n$ die folgenden Beziehungen gleichwertig sind:

(i)
$$Ax = a$$
, (ii) $\overline{A}\overline{x} = \overline{a}$, (iii) $\left(\frac{U}{V} - \frac{V}{U}\right) \left(\frac{u}{v}\right) = \left(\frac{b}{c}\right)$.

Überdies gilt

$$Q_m^{\mathsf{H}} \left(\begin{array}{c|c} U + \mathrm{i}V & O \\ \hline O & U - \mathrm{i}V \end{array} \right) Q_n = \left(\begin{array}{c|c} U & -V \\ \hline V & U \end{array} \right) \quad \mathrm{mit} \quad Q_k := rac{1}{\sqrt{2}} \left(\begin{array}{c|c} I_k & -\mathrm{i}I_k \\ \mathrm{i}I_k & I_k \end{array} \right),$$

und Q ist unitär.

Ü 1.2.2. Man beweise, daß für $A \in \mathbb{R}^{n,n}$ die folgenden Aussagen äquivalent sind:

(i) $(\alpha + i\beta, u + iv)$ ist Eigenpaar von A.

(ii) $(\alpha - i\beta, \boldsymbol{u} - i\boldsymbol{v})$ ist Eigenpaar von A.

(iii) Mit
$$\Gamma := \begin{pmatrix} \alpha & \beta \\ -\beta & \alpha \end{pmatrix} \in \mathbf{R}^{2,2}$$
 und $X := \begin{pmatrix} | & | \\ u, v \\ | & | \end{pmatrix} \in \mathbf{R}^{n,2}, X \neq O$, gilt $AX = X\Gamma$. Man be-

achte dabei, daß Γ gerade die Eigenwerte $\lambda_{1,2} = \alpha \pm i\beta$ besitzt.

Ü 1.2.3. Man überlege sich, daß der Jordanblock $J(\lambda, k) \in \mathbb{C}^{k,k}$ den Eigenwert λ der algebraischen Vielfachheit k und der geometrischen Vielfachheit 1 besitzt. Wie lautet ein zugehöriger Eigenvektor?

Ü 1.2.4. Es sei $J = T^{-1}AT$ die Jordansche Normalform von $A \in \mathbb{R}^{n,n}$. Mit $\varepsilon > 0$ werde $D := \text{diag}(1, \varepsilon, \dots, \varepsilon^{n-1}), T_{\varepsilon} := TD_{\varepsilon}$ und $J_{\varepsilon} := T_{\varepsilon}^{-1}AT_{\varepsilon}$ gebildet. Man zeige, daß J_{ε} diejenige Matrix ist, die aus J entsteht, indem dort die Einsen in der oberen Nebendiagonalen durch ε ersetzt werden. Dabei gilt $||T_{\varepsilon}^{-1}|| \cdot ||T_{\varepsilon}|| \to \infty$ für $\varepsilon \to 0$, d. h., die Transformationsmatrix T_{ε} wird für $\varepsilon \to 0$ beliebig schlecht konditioniert, siehe Abschnitt 4.1. Ü 1.2.5. Man überlege sich, daß Satz 1.2.7 für den Spezialfall einer reellen symmetrischen Matrix in Satz 1.2.4 übergeht.

Hinweis: Man zeige $R^{\tilde{T}} = R$ und beachte, daß eine symmetrische obere Dreiecksmatrix notwendig diagonal ist.

Ü 1.2.6. Es ist zu zeigen, daß genau dann $||A||_F = ||A||_2$ gilt, wenn rang $(A) \leq 1$ ist.

Ü 1.2.7. Betrachtet werde eine Matrix $A \in \mathbb{R}^{m,n}$ vom Rang r mit der Singulärwertzerlegung (25). Man zeige, daß für festes $s \in \{1, ..., r\}$ die Matrix $A_s := \sum_{j=1}^{s} \sigma_j u^j v^{j\top}$ diejenige Matrix aus der Menge $\{B \in \mathbb{R}^{m,n}: \operatorname{rang}(B) \leq s\}$ ist, die $||B - A||_F$ minimiert. Dies gilt auch für die Spektralnorm, allerdings ist der Beweis komplizierter.

Ü 1.2.8. Man zeige, daß sich die Singulärwerte der Matrix

$$A(\alpha) := \begin{pmatrix} 1 & \alpha \\ 0 & 1 \end{pmatrix}$$

für großes α wie $\sigma_1(\alpha) \sim \sqrt{2 + \alpha^2}$, $\sigma_2(\alpha) \sim 1/\sqrt{2 + \alpha^2}$ verhalten, während die Eigenwerte durch $\lambda_1(\alpha) = \lambda_2(\alpha) = 1$ für alle α gegeben sind.

Ü 1.2.9. Man zeige, daß für jede Matrix $A \in \mathbb{R}^{m,n}$

$$(||A||_2)^2 = (||A^{\mathsf{T}}||_2)^2 = ||A^{\mathsf{T}}A||_2 = ||AA^{\mathsf{T}}||_2 \le ||A||_1 ||A||_{\infty}$$
(33)

gilt, und folgere hieraus für symmetrisches $A = A^{\intercal} \in \mathbf{R}^{n,n}$

$$\|A\|_{2} \leq \||A|\|_{2} \leq \|A\|_{1} = \|A\|_{\infty}.$$
(34)

 $\text{Hinweis: Man beachte } \sigma_1^2 = \lambda_{\max}(A^{\mathsf{T}}A) \leq \|A^{\mathsf{T}}A\|_{\infty} \leq \|A^{\mathsf{T}}\|_{\infty} \|A\|_{\infty} = \|A\|_1 \|A\|_{\infty}.$

Ü 1.2.10. Man beweise die Gültigkeit von

$$\|\boldsymbol{A}\|_{2} = \max\{\|\boldsymbol{y}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x}\|:\|\boldsymbol{y}\|_{2} \leq 1, \|\boldsymbol{x}\|_{2} \leq 1\}$$
(35)

für jede Matrix $A \in \mathbf{R}^{m,n}$ und von

$$\|\boldsymbol{A}\|_{2} = \max\left\{ |\boldsymbol{x}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{x}| : \|\boldsymbol{x}\|_{2} \leq 1 \right\}$$

$$\tag{36}$$

für symmetrisches $A = A^{\mathsf{T}} \in \mathbb{R}^{n,n}$.

Bemerkungen zum Kapitel 1

B 1.1. Der in diesem Kapitel behandelte Stoff ist Gegenstand der meisten Lehrbücher der linearen Algebra und Matrixtheorie. Wir verweisen z. B. auf GANTMAHER [86] als theoretisch orientierte sowie FADDEEV/FADDEEVA [63] als numerisch orientierte detaillierte Darstellungen.

B 1.2. Der größte Teil der in diesem Buch dargestellten Gebiete läßt sich im Reellen behandeln, lediglich bei Eigenwertaufgaben nichtsymmetrischer Matrizen müssen komplexe Eigenwerte und damit komplexe Vektoren und Matrizen zugelassen werden. Aus diesem Grund wählen wir die Darstellung im Reellen als Basis und nehmen nur dann, wenn es der Gegenstand erfordert, eine Erweiterung ins Komplexe vor.

B 1.3. Die reelle Blockdreiecksversion 1.2.8 des klassischen Schurschen Satzes 1.2.7 ist Grundlage für den im Reellen ausführbaren doppelten **QR**-Algorithmus zur Eigenwertbestimmung bei nichtsymmetrischen Matrizen. Sie ist mit einer Beweisskizze bei STEWART [73] nachzulesen. Das Trägheitsgesetz 1.2.13 wird meist nach SYLVESTER benannt. **B 1.4.** Obwohl die Singulärwertzerlegung 1.2.14 für quadratische Matrizen bereits 1889 durch SYLVESTER und für allgemeine Matrizen 1939 durch ECKART/YOUNG gefunden wurde, hat sie erst in jüngster Zeit Eingang in die numerische lineare Algebra und die zugehörige Spezialliteratur gefunden, siehe etwa FORSYTHE/MOLER [67], STEWART [73], LAWSON/HANSON [74] und VOEVODIN [77].

2. Aufgaben, Computer, Algorithmen

Dieses Kapitel spielt für das Verständnis des vorliegenden Buches eine grundlegende Rolle. Es dient

- der Vorstellung der hier behandelten konkreten Aufgabenklassen
- der Diskussion des Begriffes "Numerisches Problem" und dessen wesentlichen Eigenschaften, speziell der Abhängigkeit der Resultate von den Eingangsdaten und deren Fehlern
- der Untersuchung der auf derzeitigen Computern realisierten Arithmetik
- der Einführung und Erläuterung des Begriffs "Numerischer Algorithmus" und damit zusammenhängender Probleme wie Auswirkung von Rundungsfehlern während der Computerrechnung.

In der klassischen Mathematikausbildung werden diese Probleme meist nur kurz bzw. überhaupt nicht behandelt. Der Leser sei deshalb aufgefordert, die folgenden Seiten besonders sorgfältig zu studieren, die Beispiele gründlich nachzuvollziehen und sich um das aktive Verständnis der z. T. nicht ganz einfachen Begriffe und Denkweisen zu bemühen. Auch in den späteren Kapiteln wird der hier angesprochene Problemkreis immer wieder an konkreten Aufgaben erläutert werden, so daß sich nach und nach der erwünschte Zustand der Vertrautheit und aktiven Beherrschung einstellen wird.

2.1. Aufgabenklassen der linearen Algebra

Die Aufgaben der linearen Algebra treten in der Regel nicht eigenständig, sondern einzeln bzw. mehrfach als Teilaufgaben bei der Lösung komplexer Probleme aus den verschiedensten Anwendungsgebieten auf.

A.1. Lösung quadratischer linearer Gleichungssysteme

Gegeben sind eine reguläre Koeffizientenmatrix $A \in \mathbb{R}^{n,n}$ und eine rechte Seite $b \in \mathbb{R}^n$. Gesucht ist die in diesem Fall eindeutige Lösung $x \in \mathbb{R}^n$ des linearen Gleichungssystems

$$A\boldsymbol{x} = \boldsymbol{b}.$$
 (1)

Beispiele für das Auftreten solcher Gleichungssysteme sind

- Bilanzprobleme in Technik und Ökonomie
- Analyse linearer elektrischer und mechanischer Netzwerke (Systeme hoher Dimension mit vielen unregelmäßig verteilten Nullen in A)

- Lösung nichtlinearer Gleichungssysteme und Minimierungsaufgaben mittels Newton-ähnlicher Verfahren (pro Iterationsschritt ein System $A_k x^k = b^k$ zu lösen)
- Interpolation und Approximation von Kurven und Flächen mittels Spline- und anderer Funktionen (Systeme mit Bandmatrizen)
- Integration von Anfangswertaufgaben bei Systemen gewöhnlicher Differentialgleichungen mittels impliziter Verfahren (Systeme hoher Dimension mit vielen unregelmäßig verteilten Nullen, siehe Ü 2.1.1)
- Diskretisierung von Randwertaufgaben bei gewöhnlichen und partiellen Differentialgleichungen mittels Differenzenverfahren oder finiter Elemente (Systeme hoher Dimension mit vielen regelmäßig verteilten Nullen, Systeme mit Bandmatrizen, siehe Ü 2.1.2)
- Lösung von Anfangs-Randwert-Aufgaben bei partiellen Differentialgleichungen.

Bei den angegebenen Beispielen treten u. a. folgende Sonderfälle und Modifikationen der Aufgabenstellung auf:

- Lösung von Systemen mit Matrizen spezieller Struktur:
 A symmetrisch bzw. symmetrisch und positiv definit und/oder
 A von hoher Dimension mit vielen regelmäßig bzw. unregelmäßig verteilten Nullen, speziell A Bandmatrix (Diskretisierung selbstadjungierter Differentialgleichungen, Optimierungsalgorithmen)
- Lösung mehrerer Systeme $Ax^k = b^k$ (k = 1, ..., q) mit derselben Koeffizientenmatrix (iterative Verbesserung, Lösung von Matrixgleichungen AX = B, $X = (x^1, ..., x^q), B = (b^1, ..., b^q) \in \mathbb{R}^{n,q}$)
- Lösung mehrerer Systeme $A_k x^k = b^k$ (k = 1, 2, ...) mit verschiedenen Koeffizientenmatrizen, wobei

rang $(A_{k+1} - A_k) = 1$ oder 2 (Quasi-Newton-Verfahren)

oder

rang $(A_{k+1} - A_k) = n$, aber *n* groß und Besetztheitsmuster aller A_k identisch (Anfangswertaufgaben bei Differentialgleichungen).

Eine Matrix, die viele Nullelemente enthält, wird schwach besetzte Matrix (engl. "sparse matrix", russ. "разреженная матрица") genannt. Man spricht von schwacher Besetztheit, wenn die Anzahl der Nichtnullelemente kleiner als etwa 10% der Gesamtzahl aller Elemente ist; meist ist sie sogar wesentlich geringer. Schwach besetzte lineare Gleichungssysteme entstehen z. B. bei der Analyse von Netzwerken und bei der Diskretisierung von Differentialgleichungen und erfordern i. allg. spezielle Lösungsverfahren, deren ausführliche Behandlung über den Rahmen dieses Buches hinausgeht. Einige Basistechniken werden im Abschnitt 6.4 vorgestellt.

A 2. Berechnung inverser Matrizen

Gegeben ist eine reguläre Matrix $A \in \mathbb{R}^{n,n}$. Gesucht ist die inverse Matrix

 A^{-1} .

Ein wichtiger *Sonderfall* ist die Inversion einer symmetrischen, positiv definiten Matrix.

Beispiele für Probleme, bei denen inverse Matrizen explizit berechnet werden müssen, kommen aus der Statistik.

2.1.1. Warnung. Ausdrücke der Form $X = A^{-1}B$, $B \in \mathbb{R}^{n,q}$, sollten i. allg. nicht durch Berechnung von A^{-1} und anschließende Multiplikation mit B berechnet werden, sondern billiger und genauer durch direkte Berechnung von X aus der Matrixgleichung AX = B.

Insbesondere sollte die Lösung x des regulären Systems Ax = b i. allg. nicht durch numerische Auswertung der Formel $x = A^{-1}b$ berechnet werden.

Eine *Modifikation* der Aufgabenstellung ist die Berechnung mehrerer inverser Matrizen A_k , wobei A_{k+1} aus A_k durch

- eine Rang-1-Änderung, speziell Ändern einer Zeile oder Spalte
- Streichen einer Zeile und Spalte
- Hinzufügen einer Zeile und Spalte

entsteht, siehe Abschnitt 6.5.

A 3. Berechnung von Determinanten

Gegeben ist die quadratische, nicht notwendig reguläre Matrix $A \in \mathbb{R}^{n,n}$. Gesucht ist die Zahl

 $\det(A)$.

Beispiele für das Auftreten von Determinantenberechnungen sind

- Verfahren zur Lösung linearer bzw. nichtlinearer Eigenwertprobleme $M(\lambda) x = o$, die auf der Bestimmung der Wurzeln der skalaren Gleichung $p(\lambda) := \det (M(\lambda))$ = 0 beruhen
- Verfahren zur Berechnung von Rückkehrpunkten und Verzweigungspunkten bei parameterabhängigen nichtlinearen Gleichungssystemen.

Dabei heißt das Eigenwertproblem $M(\lambda) x = o$ linear, wenn $M(\lambda) = A_0 + \lambda A_1$ mit Matrizen $A_0, A_1 \in \mathbb{R}^{n,n}$ gilt; andernfalls heißt es nichtlinear. Zum Beispiel legt $M(\lambda) := A_0 + \lambda A_1 + \lambda^2 A_2$ ein nichtlineares — nämlich quadratisches — Eigenwertproblem fest.

B 1. Lösung linearer Quadratmittelprobleme

Gegeben sind eine Koeffizientenmatrix $A \in \mathbb{R}^{m,n}$ mit $m \ge n$ und eine rechte Seite $b \in \mathbb{R}^m$. Dann hat das überbestimmte lineare Gleichungssystem Ax = b i. allg. keine Lösung. Gesucht ist deshalb eine Lösung $x \in \mathbb{R}^n$ des linearen Quadratmittelproblems

$$\|A\boldsymbol{x} - \boldsymbol{b}\|_{2} \to \underset{\boldsymbol{x} \in \boldsymbol{R}^{n}}{\operatorname{Minimum}!}, \tag{2}$$

d. h. ein Vektor x, der das *Residuum* r(x) := b - Ax des überbestimmten Systems in der Euklidischen Norm minimiert. Für (2) verwenden wir äquivalent die abkürzende Schreibweise

 $Ax \cong b$.

Beispiele für das Auftreten linearer Quadratmittelprobleme sind

- lineare Regressionsprobleme (siehe Ü 2.1.3)
- Verfahren vom Gauß-Newton-Typ zur Lösung nichtlinearer Regressionsprobleme.

Sonderfälle und Modifikationen der Aufgabenstellung sind

- Lösung von Quadratmittelproblemen mit schwach besetzten Matrizen hoher Dimension
- Lösung mehrerer Probleme $Ax^k \cong b^k$ (k = 1, ..., q) mit derselben Koeffizientenmatrix (iterative Verbesserung von Quadratmittellösungen, Matrixprobleme $||AX - B||_F \to \text{Min}!, B, X \in \mathbb{R}^{m,q}$, siehe Ü 2.1.4)
- Lösung mehrerer Probleme $A_k x^k \cong b^k$ mit unterschiedlichen Matrizen A_k , wobei A_{k+1} aus A_k entsteht durch Rang-1-Änderung, speziell Ändern einer Zeile oder Spalte, Hinzufügen oder Streichen einer Zeile oder Spalte (Lineare Regressionsprobleme mit Änderungen in Ansatzfunktionen und Meßwertsätzen, siehe Ü 2.1.3)
- Lösung regularisierter Quadratmittelprobleme im Fall schlecht konditionierter bzw. inkorrekt gestellter Aufgaben $Ax \cong b$. Dieses Problem hängt eng mit der numerischen Bestimmung des Ranges der Matrix A zusammen, siehe Kapitel 11
- Lösung linearer Quadratmittelprobleme mit linearen Gleichheitsneben bedingungen :

 $\|Ax - b\|_2 \rightarrow \text{Minimum bei } Cx = d$

mit $C \in \mathbb{R}^{p,n}$, $d \in \mathbb{R}^p$ und $m \ge n - p$.

B 2. Berechnung von Pseudoinversen

Gegeben ist eine Matrix $A \in \mathbb{R}^{m,n}$ beliebigen Formats (m, n). Gesucht ist ihre Pseudoinverse

 $A^+.$

Dabei ist $A^+ \in \mathbb{R}^{n,m}$ gerade diejenige, durch A eindeutig festgelegte Matrix, die jedem $b \in \mathbb{R}^m$ gemäß $x^+ := A^+b$ die bezüglich $||x||_2$ normkleinste Lösung x^+ von $Ax \cong b$ zuordnet, siehe Abschnitt 8.1.

Beispiele für Aufgaben, bei denen Pseudoinverse explizit benötigt werden, kommen aus der Statistik.

2.1.2. Warnung. Ausdrücke der Form $Y = A^+B$, $B \in \mathbb{R}^{m,q}$, sollten i. allg. nicht durch Berechnung von A^+ und anschließende Multiplikation mit B, sondern direkt durch Berechnung der Lösungen y^k von $Ay^k \cong b^k$ (k = 1, ..., q) mit minimalem $||y^k||_2$ ermittelt werden, wobei $Y = (y^1, ..., y^q)$ und $B = (b^1, ..., b^q)$.

Modifikationen der Aufgabe bestehen in Änderungen und Regularisierungen von A analog zu B 1.

(3)

C1. Das spezielle symmetrische Eigenwertproblem

Gegeben ist eine symmetrische Matrix $A \in S^{n,n}$. Gesucht sind Eigenwerte und gegebenenfalls zugehörige orthonormierte Eigenvektoren des Eigenwertproblems

$$A\boldsymbol{x} = \lambda \boldsymbol{x}.\tag{4}$$

Beispiele für die Herkunft solcher Probleme sind

- Diskretisierung spezieller Typen von Eigenwertaufgaben bei Differential- und Integralgleichungen
- allgemeine symmetrische Eigenwertprobleme, die auf spezielle reduziert worden sind, siehe Abschnitt 14.2.
- statistische Fragestellungen.

Sonderfälle des Problems stellen dar

- die Berechnung der p größten oder kleinsten Eigenwerte
- die Berechnung der in einem vorgegebenen Intervall [a, b] liegenden Eigenwerte und gegebenenfalls Berechnung der zugehörigen Eigenvektoren
- Aufgaben mit Matrizen spezieller Struktur wie Bandmatrizen o. ä.

C 2. Das allgemeine symmetrische Eigenwertproblem

Gegeben sind die symmetrischen Matrizen $A, B \in S^{n,n}$, wobei B positiv definit sei. Gesucht sind gewisse Eigenwerte und gegebenenfalls zugehörige Eigenvektoren des Matrizenpaares $\{A, B\}$ als Lösungen des sog. allgemeinen symmetrischen Eigenwertproblems

$$A\boldsymbol{x} = \lambda \boldsymbol{B}\boldsymbol{x}.$$
 (5)

Das wichtigste *Beispiel* für das Auftreten allgemeiner symmetrischer Eigenwertprobleme ist die Untersuchung des Eigenschwingungsverhaltens von mechanischen oder elektrischen Systemen, siehe Ü 2.1.6.

C 3. Das spezielle nichtsymmetrische Eigenwertproblem

Gegeben ist die nichtsymmetrische Matrix $A \in \mathbb{R}^{n,n}$. Gesucht sind gewisse Eigenwerte und gegebenenfalls zugehörige Eigenvektoren von A als Lösungen von

$$Ax = \lambda x$$
.

Beispiele für das Auftreten solcher Probleme sind

- Diskretisierung von nicht-selbstadjungierten Eigenwertproblemen bei Differentialgleichungen
- -- Bestimmung der allgemeinen Lösung von linearen Differentialgleichungssystemen $\dot{y} = Ay$, siehe Ü 2.1.5.
- 4 Schwetlick, Numerische Algebra

D. Numerische Probleme

Alle bisher beschriebenen Aufgabenklassen haben die folgende Struktur: Jede einzelne Aufgabe der Klasse ist durch eine endliche Zahl von *Eingangsdaten* mit gewissen Eigenschaften charakterisiert, denen durch die mathematische Aufgabenstellung in eindeutiger Weise eine endliche Zahl von *Ausgangsdaten* als Ergebnis bzw. Lösung zugeordnet wird.

2.1.3. Beispiel. Eingangs- und Ausgangsdaten für ausgewählte Aufgabenklassen

- (i) Lösung beliebig regulärer Gleichungssysteme Ax = b der Ordnung n. Eingangsdaten: a_{ij} , b_i (i, j = 1, ..., n), wobei A regulär ist Ausgangsdaten: x_i (i = 1, ..., n)
- (ii) Lösung regulärer Gleichungssysteme Ax = b mit Tridiagonalmatrizen der Ordnung n. Eingangsdaten: a_{ij}, b_i $(i, j = 1, ..., n, |i - j| \leq 1)$, wobei A regulär ist Ausgangsdaten: x_i (i = 1, ..., n)
- (iii) Bestimmung der drei kleinsten Eigenwerte reeller symmetrischer Matrizen A der Ordnung n.

Eingangsdaten: a_{ij} $(i, j = 1, ..., n, i \leq j)$

```
Ausgangsdaten: \lambda_{n-2}, \lambda_{n-1}, \lambda_n
```

Aufgaben des angeführten Typs werden numerische Problemklassen genannt. Eine solche numerische Problemklasse $\{\mathcal{E}, \boldsymbol{P}\}$ ist also eine Abbildung

$$\boldsymbol{P}\colon \mathscr{E} \subset \mathbf{R}^N \to \mathbf{R}^M,\tag{6}$$

die jedem N-dimensionalen Satz $e = (e_1, \ldots, e_N)^{\mathsf{T}} \in \mathcal{E}$ von Eingangsdaten aus dem zur Klasse gehörenden Definitionsbereich \mathcal{E} in eindeutiger Weise einen *M*-dimensionalen Satz $\boldsymbol{a} = (a_1, \ldots, a_M)^{\mathsf{T}}$ von Ausgangsdaten zugeordnet. Dabei kann \boldsymbol{P} wie im Fall des regulären Gleichungssystems explizit angebbar sein — nämlich als $\boldsymbol{x} = \boldsymbol{A}^{-1}\boldsymbol{b}$ — bzw. wie im Fall der Eigenwertberechnung nur implizit definiert sein. Die Bestimmung von $\boldsymbol{a} = \boldsymbol{P}(\boldsymbol{e})$ für ein einzelnes, festes $\boldsymbol{e} \in \mathcal{E}$ heißt *numerisches Problem* $\{\boldsymbol{e}, \boldsymbol{P}\}$ aus der Klasse $\{\mathcal{E}, \boldsymbol{P}\}$. So ist die Lösung eines konkreten regulären Gleichungssystems $\boldsymbol{Ax} = \boldsymbol{b}$ für festes $\boldsymbol{A}, \boldsymbol{b}$ ein numerisches Problem aus der Klasse "Lösung regulärer Gleichungssysteme".

2.1.4. Merksatz. Numerische Probleme sind dadurch gekennzeichnet, daß die Eingangsdaten als fehlerbehaftet angesehen werden müssen.

Dabei sind drei Fehlerquellen zu unterscheiden, die in der Regel gleichzeitig auftreten:

- Die Eingangsdaten sind Ergebnisse von Messungen bzw. von vorausgegangenen Rechnungen und deshalb notwendig fehlerbehaftet
- Wenn $\boldsymbol{a} = \boldsymbol{P}(\boldsymbol{e})$ durch einen numerischen Algorithmus auf einem Computer berechnet werden soll, müssen die Eingangsdaten e_j als Computerzahlen im Speicher dargestellt werden. Selbst bei fehlerfreien Eingangsdaten wird dabei i. allg. ein - wenn auch kleiner - Darstellungsfehler hervorgerufen, vgl. Abschnitte 2.2 und 2.3

- Für viele Algorithmen läßt sich zeigen, daß die berechnete und durch Rundungsfehler verfälschte Computernäherung \tilde{a} für a = P(e) die exakte Lösung eines benachbarten Problems mit den gestörten Eingangsdaten e + de ist, d. h., es gilt $\tilde{a} = P(e + de)$, wobei für die Störung de eine kleine Schranke Δe mit $||de|| \leq \Delta e$ angegeben werden kann. Das ist das Prinzip der rückwärtigen Fehleranalyse, vgl. Abschnitt 2.3. Die Berechnung und Abschätzung von de ist oft wesentlich einfacher als die direkte Abschätzung von $\tilde{a} - a$ und hat sich speziell in der numerischen linearen Algebra als äußerst tragfähig erwiesen.

Unabhängig von der Art des Eingangsfehlers liegt dabei folgende Situation vor: Die Eingangsdaten e des konkreten Problems $\{e, P\}$ sind nur im Rahmen eines gewissen *Fehlerniveaus* bekannt, das hier der Einfachheit halber durch die kollektive Normschranke $\Delta e > 0$ charakterisiert werden soll: Neben e sind alle Eingangsdaten \tilde{e} mit einem Fehler — genauer: *absolutem Fehler* —

$$\boldsymbol{\delta \boldsymbol{e}} := \boldsymbol{\tilde{e}} - \boldsymbol{e},\tag{7}$$

für den

$$\|\boldsymbol{\delta}\boldsymbol{e}\| = \|\boldsymbol{\tilde{e}} - \boldsymbol{e}\| \leq \Delta \boldsymbol{e} \tag{8}$$

gilt, als gleichberechtigt anzusehen. Selbstverständlich können auch individuelle Fehlerschranken des Typs $|\delta e_i| = |\tilde{e}_i - e_i| \leq \Delta e_i$ verwendet werden.

Das fehlerbehaftete Problem $\{e, P\}$ mit dem Fehlerniveau Δe ist daher durch die Menge

$$\mathscr{E}(\boldsymbol{e}, \Delta \boldsymbol{e}) := \{ \boldsymbol{\tilde{e}} \in \mathscr{E} : \| \boldsymbol{\tilde{e}} - \boldsymbol{e} \| \leq \Delta \boldsymbol{e} \}$$
(9)

von gleichberechtigten Eingangsdaten charakterisiert. Als Lösung muß dann neben a = P(e) jedes $\tilde{a} = P(\tilde{e})$ mit $\tilde{e} \in \mathcal{E}(e, \Delta e)$ akzeptiert werden, d. h., alle Elemente der Menge

$$\mathcal{A}(\boldsymbol{e}, \Delta \boldsymbol{e}) := \boldsymbol{P}(\mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{e})) := \{ \boldsymbol{\tilde{a}} = \boldsymbol{P}(\boldsymbol{\tilde{e}}) : \boldsymbol{\tilde{e}} \in \mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{e}) \}$$
(10)

müssen als gleichberechtigte Lösungen angesehen werden. Die Lösung des fehlerbehafteten Problems $\{e, P\}$ mit dem Fehlerniveau Δe besteht daher eigentlich in der Bestimmung der gesamten zulässigen Lösungsmenge $\mathcal{A}(e, \Delta e)$.

Damit das fehlerbehaftete numerische Problem vernünftig gestellt ist, muß $\mathcal{A}(e, \Delta e)$ zumindest beschränkt sein, d. h., die Norm der Elemente von \mathcal{A} darf eine obere Schranke nicht überschreiten. Als Beispiel betrachten wir eine fehlerbehaftete reguläre Matrix $M \in \mathbb{R}^{n,n}$ mit dem Fehlerniveau ΔM . Bei zu großem ΔM liegen in der Menge $\{\widetilde{M} : \|\widetilde{M} - M\| \leq \Delta M\}$ auch singuläre Matrizen. Das Problem der Inversion von M ist dann zum Fehlerniveau ΔM nicht vernünftig gestellt, denn in $\mathcal{E}(M, \Delta M) = \{\widetilde{M} : \widetilde{M} \text{ regulär und } \|\widetilde{M} - M\| \leq \Delta M\}$ gibt es Matrizen \widetilde{M} mit beliebig großem $\|\widetilde{M}^{-1}\|$, d. h., $\mathcal{A}(M, \Delta M)$ ist nicht beschränkt; siehe Abschnitt 4.1. In diesem Fall kann durch Verkleinerung von ΔM – also genauere Bestimmung der Elemente von M – versucht werden, beschränkte Lösungsmengen zu erhalten. Im Abschnitt 4.1 werden wir sehen, daß dies für das diskutierte Beispiel der Inversion in der Tat möglich ist. Andernfalls muß die Problemstellung z. B. durch Regularisierung o. ä. modifiziert werden, vgl. Abschnitte 11.2 und 11.3. Ob dies sinnvoll und zulässig ist,

kann nicht allein vom Standpunkt der Mathematik, sondern nur vom Standpunkt des im Hintergrund stehenden realen Problems unter Einbeziehung des Naturwissenschaftlers oder Ingenieurs entschieden werden.

Es stellt sich heraus: Die exakte Bestimmung der Lösungsmenge $\mathcal{A}(\mathbf{e}, \Delta \mathbf{e})$ ist für die Aufgabenklassen der linearen Algebra praktisch nicht möglich. Man behilft sich daher mit der Konstruktion von Mengen, die $\mathcal{A}(\mathbf{e}, \Delta \mathbf{e})$ einschließen. Eine Möglichkeit dafür bildet die Intervallanalysis bzw. Intervallmathematik, siehe B 2.8. Wir gehen einen anderen, in der Regel weitaus billigeren Weg und bestimmen neben dem Repräsentanten $\mathbf{a} = \mathbf{P}(\mathbf{e})$ eine Fehlerschranke $\Delta \mathbf{a}$ so, daß

$$\|\boldsymbol{\delta}\boldsymbol{a}\| = \|\boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e})\| \leq \Delta \boldsymbol{a} \quad \text{für alle } \boldsymbol{\tilde{e}} \in \mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{e})$$
(11)

gilt. Dies ist gleichbedeutend mit

$$\mathcal{A}(\boldsymbol{e}, \Delta \boldsymbol{e}) \subset \mathcal{S}(\boldsymbol{a}, \Delta \boldsymbol{a}) := \{ \boldsymbol{\tilde{a}} \in \mathbf{R}^{M} : \| \boldsymbol{\tilde{a}} - \boldsymbol{a} \| \leq \Delta \boldsymbol{a} \},$$
(12)

d. h., die Lösungsmenge wird mit einer Normkugel um a mit dem Radius Δa eingeschlossen, siehe Abb. 2.1.1 für eine Veranschaulichung unter Verwendung der Euklidischen Norm.



Abb. 2.1.1. Lösungsmenge und Fehlerschranken bei fehlerbehafteten Eingangsdaten

Der kleinste in (11) mögliche Wert für Δa wird mit Δa_{opt} bezeichnet und gibt das *optimale* (absolute) *Fehlerniveau* der Ausgangsdaten an:

$$\Delta \boldsymbol{a}_{\text{opt}} := \min \left\{ \Delta \boldsymbol{a} > 0 \colon \| \boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e}) \| \leq \Delta \boldsymbol{a} \text{ für alle } \boldsymbol{\tilde{e}} \in \mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{e}) \right\}.$$
(13)

Eine äquivalente Festlegung ist

$$\Delta \boldsymbol{a}_{\text{opt}} = \max \left\{ \| \boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e}) \| : \boldsymbol{\tilde{e}} \in \mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{e}) \right\}.$$
(14)

Wir bemerken hier, daß "min" bzw. "max" eigentlich durch "inf" bzw. "sup" ersetzt werden müßten. Bei abgeschlossenem $\mathscr{E}(e, \varDelta e)$ und stetigem P — dies ist bei unseren Beispielen der Fall — sind jedoch die einfachen Begriffe "min" bzw. "max" zulässig.

Aus der Definition von $\Delta \boldsymbol{a}_{opt}$ folgt: Bei einem Fehlerniveau $\Delta \boldsymbol{e}$ der Eingangsdaten ist der Fehler der Ausgangsdaten durch $\Delta \boldsymbol{a}_{opt}$ beschränkt, und es gibt zulässige Eingangsdaten $\tilde{\boldsymbol{e}} \in \mathcal{E}(\boldsymbol{e}, \Delta \boldsymbol{a})$ mit $\|\boldsymbol{P}(\tilde{\boldsymbol{e}}) - \boldsymbol{P}(\boldsymbol{e})\| = \Delta \boldsymbol{a}_{opt}, d. h., das optimale Fehlerniveau$ wird angenommen. Man nennt $\Delta \mathbf{a}_{opt}$ aus diesem Grunde auch unvermeidlichen Fehler der Ausgangsdaten. Eine höhere Genauigkeit ist bei der in den Eingangsdaten $\mathcal{E}(\mathbf{e}, \Delta \mathbf{e})$ liegenden Information nicht zu erwarten, siehe Abb. 2.1.1, wo mit $\Delta \mathbf{a} = \Delta \mathbf{a}_{opt}$ gearbeitet wird.

Für die Berechnung von Δa muß die Änderung $\delta a := \tilde{a} - a = P(\tilde{e}) - P(e)$ für kleine Störungen $\tilde{e} - e$ untersucht und abgeschätzt werden. Dies ist Gegenstand der *Störungstheorie*, die exakte oder zumindest in der Größenordnung richtige Ausdrücke für $P(\tilde{e}) - P(e)$ als Funktion von $\tilde{e} - e$ bzw. für deren Normen bereitstellt.

Eine natürliche Forderung für eine vernünftige Problemstellung $\{e, P\}$ ist die Stetigkeit von P im Punkt e: Wenn \tilde{e} gegen e strebt, soll $\tilde{a} = P(\tilde{e})$ gegen a = P(e) streben; d. h., kleine Änderungen der Eingangsdaten sollen kleine Änderungen der Ausgangsdaten bewirken. Ein in diesem Sinne vernünftiges numerisches Problem wird als korrekt gestellt bezeichnet (engl. "properly posed", russ. "корректно поставлено"). Für numerische Zwecke ist die Eigenschaft "korrekt gestellt" zu schwach; die Stetigkeit im Punkt e muß relativ gutartig und quantitativ erfaßbar sein. Zur Motivierung betrachten wir das folgende Beispiel:

2.1.5. Beispiel. Stetigkeitsverhalten der Problemklasse

$$a = P(e) = 1/e, \quad e \in \mathscr{E} := \{e \in \mathbf{R} : e > 0\}.$$
 (15)

Für festes $e \in \mathscr{E}$ – also e > 0 – liegt auch jedes \tilde{e} mit $|\tilde{e} - e| \leq \Delta e$ in \mathscr{E} , sofern $\Delta e \leq \Delta_0 e < e$ ist. In diesem Fall folgt

$$|P(\tilde{e}) - P(e)| = \left|\frac{1}{\tilde{e}} - \frac{1}{e}\right| = \frac{1}{e\tilde{e}} |\tilde{e} - e| \le \frac{1}{e(e - \Delta e)} |\tilde{e} - e|.$$
(16)

Die Änderung $\tilde{e} - e$ der Eingangsdaten kann sich also höchstens mit dem Verstärkungsfaktor $L(e, \Delta e) := 1/[e(e - \Delta e)]$ auf die Ausgangsdaten auswirken. Der Faktor $L(e, \Delta e)$ hängt von e und Δe ab, und für genügend kleines Δe unterscheidet er sich immer weniger von $L(e, 0) = 1/e^2$. \Box

Das obige eindimensionale Beispiel legt es nahe, die folgenden allgemeinen Begriffe einzuführen:

2.1.6. Begriffe der Lipschitzstetigkeit. Betrachtet wird die numerische Problemklasse $\{\mathscr{E}, \mathbf{P}\}$.

(i) Das Problem $\{e, P\}$ der Klasse $\{\mathcal{E}, P\}$ heißt lokal lipschitzstetig, wenn es ein (i. allg. von e abhängiges) Fehlerniveau $\Delta_0 e > 0$ gibt, so daß zu jedem $\Delta e \leq \Delta_0 e$ eine lokale Lipschitzkonstante $L(e, \Delta e) > 0$ existiert mit

$$\|\boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e})\| \leq L(\boldsymbol{e}, \boldsymbol{\varDelta}\boldsymbol{e}) \|\boldsymbol{\tilde{e}} - \boldsymbol{e}\|$$
(17)

für alle $\tilde{e} \in \mathscr{E}$ mit $\|\tilde{e} - e\| \leq \Delta e$.

(ii) Die Klasse $\{\mathcal{E}, P\}$ heißt lokal lipschitzstetig, wenn jedes Problem $\{e, P\}, e \in \mathcal{E}$, der Klasse lokal lipschitzstetig ist.

(iii) Die Klasse $\{\mathcal{E}, \mathbf{P}\}$ heißt (schlechthin) lipschitzstetig, wenn eine (globale) Lipschitzkonstante L > 0 mit

$$\|\boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e})\| \leq L \|\boldsymbol{\tilde{e}} - \boldsymbol{e}\|$$
(18)

für alle $\tilde{e}, e \in \mathcal{E}$ existiert.

Dabei wird stets vorausgesetzt, daß die Lipschitzkonstanten minimal gewählt werden, so daß die Abschätzungen (17) bzw. (18) nicht verbessert werden können.

Wir bemerken noch, daß die Eigenschaft der Lipschitzstetigkeit unabhängig von den verwendeten Normen ist, während die Zahlenwerte der Lipschitzkonstanten natürlich von den Normen abhängen.

Sämtliche der oben eingeführten Stetigkeitsbegriffe treten bei den später behandelten Problemklassen tatsächlich auf. Als Beispiele führen wir an:

- Die Bestimmung der Normallösung eines Quadratmittelproblems $Ax \cong b$ ist ein lokal lipschitzstetiges Problem, falls A Vollrang hat. Andernfalls liegt nicht einmal Stetigkeit vor, so daß die Problemklasse, für beliebige Eingangsdaten $\{A, b\}$ die Normallösung zu bestimmen, nicht lokal lipschitzstetig ist.
- Die Lösung regulärer linearer Gleichungssysteme Ax = b ist eine lokal lipschitzstetige Problemklasse.
- Die Eigenwertbestimmung symmetrischer Matrizen A ist eine schlechthin lipschitzstetige Problemklasse mit der Lipschitzkonstanten L = 1, sofern $\tilde{A} - A$ in der *F*-Norm gemessen wird.

2.1.7. Bemerkung (i) Die lokale Lipschitzstetigkeit ist eine starke Form der Stetigkeit, deren Vorliegen keineswegs selbstverständlich ist. Zur Illustration betrachten wir die Problemklasse

$$a = P(e) = e^{\alpha}, \qquad e \in \mathscr{E} := \{e \in \mathbf{R} : e \ge 0\}, \qquad 0 < \alpha \le 1 ext{ fest.}$$

Für das durch e = 0 festgelegte Problem und beliebiges $\tilde{e} > 0$ gilt

$$|P(ilde{e})-P(e)|=|e^{lpha}|=1\cdot| ilde{e}-e|^{lpha}$$
 ,

d. h., die Änderung $\tilde{e} - e$ wirkt sich nicht in erster, sondern in der Potenz α auf die Ausgangsdaten aus. Man spricht dann von Hölderstetigkeit. Etwa zur Definition der lokalen Hölderstetigkeit müßte (17) durch

$$\|P(\tilde{e}) - P(e)\| \leq L(e, \Delta e) \|\tilde{e} - e\|^{\mathfrak{a}(e)}$$
⁽¹⁹⁾

ersetzt werden, wobei $\alpha(e)$ den größtmöglichen Exponenten bezeichnet, für den (19) gilt; die übrigen Begriffe wären sinngemäß zu übertragen. Je kleiner α ist, um so schwächer ist diese Form der Stetigkeit: Im Beispiel ergibt sich für e = 0, $\tilde{e} - e = 10^{-6}$ etwa im Fall $\alpha = 1/2$ die Änderung $\delta a = P(\tilde{e}) - P(e) = 10^{-3}$, im Fall $\alpha = 1/12$ dagegen $\delta a = 10^{-1/2} = 0.3!$ Man beachte, daß für e > 0 lokale Lipschitzstetigkeit vorliegt, also α in der Tat von e abhängen kann.

(ii) Es gibt auch in der linearen Algebra Probleme, die nur lokal hölderstetig mit $\alpha < 1$ sind. So können die Eigenwerte nichtsymmetrischer Matrizen A lokal hölder-

stetig mit $\alpha \ge 1/n$ von A abhängen, und der Wert $\alpha = 1/n$ kann angenommen werden: Die Matrix

 $\boldsymbol{A}(\varepsilon) := \begin{pmatrix} 1 & 1 \\ & 1 & 1 \\ & & \ddots & \\ & & \ddots & \ddots \\ & & & \ddots & \ddots \\ \varepsilon & & \ddots & \ddots & 1 \\ \varepsilon & & \ddots & \ddots & 1 \\ \varepsilon & & \ddots & \ddots & 1 \end{pmatrix} \in \mathbf{R}^{n,n}$

hat das charakteristische Polynom det $(A(\varepsilon) - \lambda I) = (1 - \lambda)^n + (-1)^{n+1} \varepsilon$, so daß für die Eigenwerte $\lambda_j(\varepsilon)$ die Beziehung

$$|\lambda_j(\varepsilon) - \lambda_j(0)| = \sqrt[n]{\varepsilon} = \varepsilon^{1/n}$$

folgt; man beachte $\lambda_j(0) = 1$ (j = 1, ..., n). Eine Störung $\varepsilon = 10^{-6}$ bewirkt also im Fall n = 12 eine Änderung der Eigenwerte um $|\delta\lambda_j| = 10^{-1/2} = 0.3$. Da wir das nichtsymmetrische Eigenwertproblem nicht behandeln, gehen wir auf die Hölderstetigkeit nicht weiter ein.

(iii) Falls P in der Umgebung von $e \in \mathcal{E}$ stetige partielle Ableitungen erster Ordnung besitzt, ist das Problem $\{e, P\}$ lokal lipschitzstetig, vgl. Ü 2.1.7.

Mittels der lokalen Lipschitzkonstanten läßt sich das optimale Fehlerniveau einfach abschätzen bzw. unter zusätzlichen Glattheitsvoraussetzungen sogar asymptotisch exakt angeben.

2.1.8. Aussage. Betrachtet wird das numerische Problem $\{e, P\}$.

(i) Wenn $\{e, P\}$ lokal lipschitzstetig ist, gilt für $\Delta e \leq \Delta_0 e$

$$\Delta \boldsymbol{a}_{\text{opt}} \leq L(\boldsymbol{e}, \Delta \boldsymbol{e}) \, \Delta \boldsymbol{e}. \tag{20}$$

(ii) Wenn P in der Umgebung von e genügend glatt im Sinne der Existenz stetiger partieller Ableitungen zweiter Ordnung ist, liegt lokale Lipschitzstetigkeit vor, die lokale Lipschitzkonstante $L(e, \Delta e)$ hängt stetig von $\Delta e \ge 0$ ab, der Grenzwert

$$L(\boldsymbol{e}) := \lim_{\Delta \boldsymbol{e} \to +0} L(\boldsymbol{e}, \Delta \boldsymbol{e})$$
⁽²¹⁾

existiert, und es gilt $L(e, \Delta e) = L(e) + O(\Delta e)$ sowie

$$\Delta \boldsymbol{a}_{\text{ont}} = \left[\boldsymbol{L}(\boldsymbol{e}) + \boldsymbol{O}(\Delta \boldsymbol{e}) \right] \Delta \boldsymbol{e} \tag{22}$$

für $\Delta e \to 0$, $\Delta e \leq \Delta_0 e$ mit genügend kleinem $\Delta_0 e > 0$.

Für einen festen Exponenten $\beta > 0$ bezeichnet $O(\xi^{\beta})$ dabei eine Funktion von $\xi \in \mathbf{R}$, die für $\xi \to 0$ mindestens so schnell wie ξ^{β} gegen 0 geht, d. h., es gilt

 $|O(\xi^{\beta})| \leq C |\xi|^{\beta}$

für genügend kleines $|\xi|$ mit einer Konstanten C > 0.

Beweis. Aussage (i) folgt sofort aus (13) und (17). Teil (ii) ergibt sich aus der Taylorformel mit Restglied zweiter Ordnung für P an der Stelle e, vgl. wieder Ü 2.1.7.

Bei den später behandelten Problemen ist die Glattheitsforderung aus (ii) immer erfüllt. Wir nehmen daher im folgenden stets an, daß P dieser Voraussetzung genügt und somit die Darstellung (22) für das optimale Fehlerniveau gültig ist. Aus dieser folgt

$$\Delta \boldsymbol{a}_{\rm out} \approx L(\boldsymbol{e}) \,\Delta \boldsymbol{e} \tag{23}$$

für genügend kleines Ae, und (23) ist asymptotisch exakt im Sinne von

$$\frac{\angle \mathbf{a}_{opt}}{\varDelta \mathbf{e}} \to \mathbf{L}(\mathbf{e}) \quad \text{für} \quad \varDelta \mathbf{e} \to 0.$$
(24)

Die Zahl L(e) wird absolute Konditionszahl des lokal lipschitzstetigen Problems $\{e, P\}$ genannt. Sie gibt den maximal möglichen Verstärkungsfaktor an, mit dem sich eine kleine Störung der Eingangsdaten auf die Ausgangsdaten auswirken kann.

In den meisten Fällen ist es zweckmäßig, mit relativen Fehlern und Fehlerschranken anstelle von absoluten zu arbeiten. Für eine skalare Größe $e_i \in \mathbf{R}$ heißt unter der Voraussetzung $e_i \neq 0$ der Quotient

$$\varepsilon_i := \varepsilon(e_i) := \frac{\tilde{e}_i - e_i}{e_i} = \frac{\delta e_i}{e_i}$$
(25)

relativer Fehler von e_i und wird üblicherweise in Prozent angegeben; gelegentlich wird auch mit $(\tilde{e}_i - e_i)/\tilde{e}_i$ gearbeitet, siehe B 2.1 und Ü 2.1.8. Durch ε_i kann die im Repräsentanten e_i enthaltene Information unabhängig von dessen Größenordnung und damit besser charakterisiert werden; z. B. läßt sich $-\log_{10} |\varepsilon_i|$ als Anzahl der gültigen Dezimalziffern in der Dezimaldarstellung von e_i deuten, vgl. Abschnitt 2.2. Bei Vektoren $e = (e_i) \in \mathbb{R}^N$ können die individuellen relativen Fehler der Komponenten für jedes *i* gemäß (25) definiert werden. Häufig ist es jedoch einfacher und zweckmäßiger, mit kollektiven relativen Fehlern

$$\frac{\|\boldsymbol{\tilde{e}} - \boldsymbol{e}\|}{\|\boldsymbol{e}\|} \quad \text{bzw.} \quad \frac{\|\boldsymbol{P}(\boldsymbol{\tilde{e}}) - \boldsymbol{P}(\boldsymbol{e})\|}{\|\boldsymbol{P}(\boldsymbol{e})\|}$$

zu arbeiten. Dazu muß natürlich $e \neq o$ bzw. $P(e) \neq 0$ vorausgesetzt werden.

Wenn $\{e, P\}$ lokal lipschitzstetig ist, läßt sich (17) mit

$$K(e, \Delta e) := L(e, \Delta e) \frac{\|e\|}{\|P(e)\|}$$

äquivalent in der Form

$$\frac{\|\boldsymbol{\delta a}\|}{\|\boldsymbol{a}\|} = \frac{\|\boldsymbol{P}(\tilde{\boldsymbol{e}}) - \boldsymbol{P}(\boldsymbol{e})\|}{\|\boldsymbol{P}(\boldsymbol{e})\|} \leq K(\boldsymbol{e}, \varDelta \boldsymbol{e}) \frac{\|\tilde{\boldsymbol{e}} - \boldsymbol{e}\|}{\|\boldsymbol{e}\|}$$

schreiben, d. h., die relativen Fehler genügen einer zu (17) analogen Ungleichung mit der relativen lokalen Lipschitzkonstanten $K(e, \Delta e)$ statt $L(e, \Delta e)$. Mit $L(e, \Delta e)$ ist auch $K(e, \Delta e)$ optimal, d. h. nicht durch eine kleinere Konstante ersetzbar. Die $2.1.8(\mathrm{i})$ entsprechende Abschätzung des relativen optimalen Fehlerniveaus lautet dann

$$\frac{\Delta \boldsymbol{a}_{\text{opt}}}{\|\boldsymbol{a}\|} \leq K(\boldsymbol{e}, \Delta \boldsymbol{e}) \frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|} \quad \text{für} \quad \frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|} \leq \frac{\Delta_0 \boldsymbol{e}}{\|\boldsymbol{e}\|},$$
(26)

und analog zu 2.1.8(ii) ist $K(e, \Delta e)$ für genügend glattes P bezüglich Δe stetig, der Grenzwert

$$K(\boldsymbol{e}) := \lim_{\Delta \boldsymbol{e} \to +0} K(\boldsymbol{e}, \Delta \boldsymbol{e}) = L(\boldsymbol{e}) \|\boldsymbol{e}\| / \|\boldsymbol{P}(\boldsymbol{e})\|$$
(27)

existiert, und es gilt $K(e, \Delta e) = K(e) + \theta(\Delta e/||e||)$ sowie

$$\frac{\Delta \boldsymbol{a}_{\text{opt}}}{\|\boldsymbol{a}\|} = \left[K(\boldsymbol{e}) + \theta \left(\frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|} \right) \right] \frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|}$$

für $\Delta e/\|e\| \to 0$, $\Delta e/\|e\| \le \Delta_0 e/\|e\|$, $\Delta_0 e$ genügend klein. Dies bedeutet

$$\frac{\Delta \boldsymbol{a}_{\text{opt}}}{\|\boldsymbol{a}\|} \approx K(\boldsymbol{e}) \frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|}$$
(28)

für genügend kleines $\Delta e/||e||$, und (28) ist im Sinne von

$$\frac{\Delta \boldsymbol{a}_{opt}/\|\boldsymbol{a}\|}{\Delta \boldsymbol{e}/\|\boldsymbol{e}\|} \to K(\boldsymbol{e}) \quad \text{für} \quad \frac{\Delta \boldsymbol{e}}{\|\boldsymbol{e}\|} \to 0$$
(29)

asymptotisch exakt. Das optimale relative Fehlerniveau der Ausgangsdaten ist also das K(e)-fache des relativen Fehlerniveaus der Eingangsdaten, sofern letzteres genügend klein ist. Der relative Fehler der Eingangsdaten kann sich höchstens mit dem Verstärkungsfaktor K(e) als relativer Fehler der Ausgangsdaten auswirken. Die Zahl K(e) wird deshalb relative Konditionszahl des Problems $\{e, P\}$ genannt.

Die folgenden Beispiele sollen die neu eingeführten Begriffe erläutern und zeigen, wie Lipschitzkonstanten und Konditionszahlen in konkreten Fällen berechnet werden können.

2.1.9. Beispiel. Bestimmung der Euklidischen Norm $r(x) := ||x||_2$ für $x \in \mathbb{R}^n$. Wegen der Normeigenschaft (N'_3) aus 1.1.I gilt

$$|\delta r| = |||x + \delta x||_2 - ||x||_2| \le ||(x + \delta x) - x||_2 = ||\delta x||_2$$
(30)

für alle $x, \, \delta x \in \mathbb{R}^n$, d. h., die Aufgabenklasse ist lipschitzstetig mit der Konstanten L = 1. Aus (30) folgt

$$\left|\frac{\delta r}{r}\right| = \frac{|\|\boldsymbol{x} + \boldsymbol{\delta x}\|_2 - \|\boldsymbol{x}\|_2|}{\|\boldsymbol{x}\|_2} \le \frac{\|\boldsymbol{\delta x}\|_2}{\|\boldsymbol{x}\|_2}$$
(31)

für alle $x, \, \delta x \in \mathbb{R}^n$ mit $x \neq o$, d. h., auch die relativen Fehler können sich nicht verstärken. Da in (30) und (31) für $\delta x = \lambda x, \, x \neq o, \, \lambda \neq 0$ beliebig, das Gleichheitszeichen steht, folgt L(x) = 1 für alle x und K(x) = 1 für alle $x \neq o$. 58

Um den Einfluß von dx nicht nur in einer Normabschätzung, sondern direkt erfassen zu können, betrachten wir

$$egin{aligned} \delta r &= \|m{x} + m{d} m{x}\|_2 - \|m{x}\|_2 = (\|m{x}\|_2^2 + 2m{x}^{\intercal}m{d} m{x} + \|m{d} m{x}\|_2^2)^{1/2} - (\|m{x}\|_2^2)^{1/2} \ &= [(\|m{x}\|_2^2 + 2m{x}^{\intercal}m{d} m{x} + \|m{d} m{x}\|_2^2) - (\|m{x}\|_2^2)]/[(\|m{x}\|_2^2 + 2m{x}^{\intercal}m{d} m{x} + \|m{d} m{x}\|_2^2)^{1/2} \ &+ (\|m{x}\|_2^2)^{1/2}]. \end{aligned}$$

Berücksichtigung von Gliedern erster Ordnung in dx führt auf

$$\delta r = \|\boldsymbol{x} + \boldsymbol{\delta x}\|_2 - \|\boldsymbol{x}\|_2 = \frac{\boldsymbol{x}^{\mathsf{T}} \boldsymbol{\delta x}}{\|\boldsymbol{x}\|_2} + O(\|\boldsymbol{\delta x}\|^2) \quad \text{für} \quad \boldsymbol{x} \neq \boldsymbol{o}.$$
(32)

Die Beziehung $v(\varepsilon) = O(\varepsilon^p)$ für $\varepsilon \in \mathbf{R}$ bedeutet dabei $||v(\varepsilon)|| \leq C |\varepsilon|^p$ für genügend kleines $|\varepsilon|$ mit einer Konstanten C > 0. Die Darstellung (32) besagt u. a., daß sich eine zu x orthogonale Störung δx in erster Ordnung nicht auf $||x + \delta x||_2$ auswirkt, während der Einfluß für $\delta x = \lambda x$ maximal wird. Dies ist anschaulich evident.

Zahlenbeispiel. Für n = 2 und $\boldsymbol{x} = (0.9987, 0.05097)^{\mathsf{T}}$ ist $r(\boldsymbol{x}) = \|\boldsymbol{x}\|_2 = 1.000000$. Das zu $\Delta \boldsymbol{x} = 0.01$ gehörende optimale Fehlerniveau ist $\Delta r_{\text{opt}} = \Delta \boldsymbol{x} = 0.01$.

(i) Für $\delta x = (0.01, 0)^{T}$ ergibt sich $r(x + \delta x) = 1.009987$, also

 $\delta r = 0.009987$ und $|\delta r|/\Delta r_{opt} = 99.87\%$.

(ii) Für $\delta x = (0, 0.01)^{T}$ ergibt sich $r(x + \delta x) = 1.000559$, also

 $\delta r = 0.000559$ und $|\delta r|/\Delta r_{opt} = 5.59\%$.

Das Beispiel zeigt, daß ein Fehler δr in der Größenordnung Δr_{opt} nicht für jede im Rahmen des Fehlerniveaus $\Delta x = 0.01$ zulässige Störung δx angenommen wird. Da der "wahre" Fehler δx der Eingangsdaten i. allg. nicht bekannt ist, muß jedoch mit einem Fehler von maximal $\Delta r_{opt} = 0.01$ gerechnet werden.

Falls die Eingangsdaten e aus verschiedenen Gruppen — etwa e = (A, b) bei linearen Gleichungssystemen bzw. e = (x, y) beim Skalarprodukt — bestehen, ist es zweckmäßig, den Einfluß jeder Gruppe getrennt zu berücksichtigen und additiv zu überlagern. Man kommt dann zu sog. *partiellen Konditionszahlen* bezüglich jeder Gruppe. Wir demonstrieren diese Verallgemeinerung der bisher eingeführten Begriffe am Beispiel des Skalarproduktes $x^{T}y$.

2.1.10. Beispiel. Bestimmung des Skalarproduktes $s = s(x, y) = x^{\mathsf{T}}y$ für $x, y \in \mathbb{R}^n$. Hier gilt

$$\delta s = (\boldsymbol{x} + \boldsymbol{d} \boldsymbol{x})^{\mathsf{T}} (\boldsymbol{y} + \boldsymbol{d} \boldsymbol{y})^{\mathsf{T}} - \boldsymbol{x}^{\mathsf{T}} \boldsymbol{y} = \boldsymbol{d} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{y} + \boldsymbol{x}^{\mathsf{T}} \boldsymbol{d} \boldsymbol{y} + \boldsymbol{d} \boldsymbol{x}^{\mathsf{T}} \boldsymbol{d} \boldsymbol{y},$$

also

$$\begin{aligned} |\delta s| &\leq ||\delta x||_2 ||y||_2 + ||\delta y||_2 ||x||_2 + ||\delta x||_2 ||\delta y||_2 \\ &= (||y||_2 + 0.5 ||\delta y||_2) ||\delta x||_2 + (||x||_2 + 0.5 ||\delta x||_2) ||\delta x||_2. \end{aligned}$$
(33)

Für beliebige $\Delta x, \Delta y > 0$ und $\| \delta x \|_2 \leq \Delta x, \| \delta y \|_2 \leq \Delta y$ folgt

 $|\delta s| \leq (\|\boldsymbol{y}\|_2 + 0.5 \Delta \boldsymbol{y}) \, \|\boldsymbol{\delta x}\|_2 + (\|\boldsymbol{x}\|_2 + 0.5 \Delta \boldsymbol{x}) \, \|\boldsymbol{\delta y}\|_2, \tag{34}$

d. h., die Problemklasse ist lokal lipschitzstetig mit den partiellen lokalen Lipschitzkonstanten $L_{\mathbf{x}}(\mathbf{x}, \mathbf{y}, \Delta \mathbf{x}, \Delta \mathbf{y}) := \|\mathbf{y}\|_2 + 0.5\Delta \mathbf{y}, L_{\mathbf{y}}(\mathbf{x}, \mathbf{y}, \Delta \mathbf{x}, \Delta \mathbf{y}) := \|\mathbf{x}\|_2 + 0.5\Delta \mathbf{x}.$ Für $\mathbf{x} \mathbf{y} + 0$ ergibt sich aus (33)

Für $x^{\mathsf{T}}y \neq 0$ ergibt sich aus (33)

$$\frac{|\delta s|}{|s|} \leq \frac{||\bm{x}||_2 ||\bm{y}||_2}{|\bm{x}^{\mathsf{T}}\bm{y}|} \left\{ \frac{||\delta \bm{x}||_2}{||\bm{x}||_2} + \frac{||\delta \bm{y}||_2}{||\bm{y}||_2} + \frac{||\delta \bm{x}||_2}{||\bm{x}||_2} \cdot \frac{||\delta \bm{y}||_2}{||\bm{y}||_2} \right\}.$$

Ist $\|\boldsymbol{\delta x}\|_2/\|\boldsymbol{x}\|_2 \leq \Delta \boldsymbol{x}/\|\boldsymbol{x}\|_2$ und $\|\boldsymbol{\delta y}\|_2/\|\boldsymbol{y}\|_2 \leq \Delta \boldsymbol{y}/\|\boldsymbol{y}\|_2$, so folgt

$$\frac{|\delta s|}{|s|} \leq \frac{\|\boldsymbol{x}\|_{2} \|\boldsymbol{y}\|_{2}}{|\boldsymbol{x}^{\mathsf{T}}\boldsymbol{y}|} \left\{ \left(1 + 0.5 \frac{\Delta \boldsymbol{y}}{\|\boldsymbol{y}\|_{2}}\right) \frac{\|\delta \boldsymbol{x}\|_{2}}{\|\boldsymbol{x}\|_{2}} + \left(1 + 0.5 \frac{\Delta \boldsymbol{x}}{\|\boldsymbol{x}\|_{2}}\right) \frac{\|\delta \boldsymbol{y}\|_{2}}{\|\boldsymbol{y}\|_{2}} \right\}.$$
(35)

Aus (34) bzw. (35) liest man für Δx , $\Delta y \to 0$ die absoluten bzw. relativen partiellen Konditionszahlen

$$L_{x}(x, y) := ||y||_{2}, \quad L_{y}(x, y) := ||x||_{2} \quad \text{bzw.} \quad K_{x}(x, y) = K_{y}(x, y) = \frac{||x||_{2} ||y||_{2}}{|x^{\mathsf{T}}y|} \quad (36)$$

ab. Aus (35) bzw. (36) folgt insbesondere, daß der relative Fehler des Skalarproduktes fast orthogonaler Vektoren beliebig groß werden kann, d. h., die Skalarproduktbestimmung fast orthogonaler Vektoren ist bezüglich des relativen Fehlers ein schlecht konditioniertes Problem.

Zahlenbeispiel. Für n = 2 und $x = (0.9987, 0.05097)^T$, $y = (1.023, -20.16)^T$ ist $s = x^T y$ = -0.005885. Für die im Rahmen des absoluten Fehlerniveaus $\Delta x = \Delta y = 0.01$ zulässigen Störungen $\delta x = (0, -0.01)^T$ und $\delta y = (0.01, 0)^T$ ergibt sich $(x + \delta x)^T (y + \delta y) = 0.205702$, also $\delta s = 0.2116$ und $\delta s/s = -35.95$. Die absoluten Konditionszahlen sind $L_x = ||y||_2 = 20.19$, $L_y = ||x||_2 = 1.000$, die relativen $K_x = K_y = ||x||_2 ||y||_2/|x^Ty| = 3430$, d. h., die absoluten Fehler können sich etwa um den Faktor 20, die relativen Fehler jedoch um den Faktor 3430 verstärken. Die in Analogie zu (24) bzw. (29) berechneten Näherungswerte für Δs_{opt} bzw. $\Delta s_{opt}/|s| \sin \Delta s_{opt} \approx L_x \Delta x + L_y \Delta y = 0.1219$ bzw. $\Delta s_{opt}/|s| \approx K_x \cdot \Delta x/||x||_2 + K_y \cdot \Delta y/||y||_2$ = 36.000, d. h., mit der obigen Störung δx , δy wird das unvermeidliche Fehlerniveau praktisch erreicht. \Box

Wir haben bisher korrekt gestellte Probleme $\{e, P\}$ aus der Klasse $\{\mathcal{E}, P\}$ betrachtet, d. h. Probleme, die eine eindeutige, stetig von e abhängende Lösung a = P(e)besitzen. Ist dies nicht der Fall, so heißt das Problem *inkorrekt* gestellt (engl. "illposed", "not properly posed", russ. "некорректно поставлено"). Es gibt dann evtl. keine Lösung, oder es existieren mehrere Lösungen, oder die Lösung hängt nicht stetig von den Eingangsdaten ab. Ein Beispiel für Mehrdeutigkeit ist die Eigenvektorbestimmung bei mehrfachen Eigenwerten symmetrischer Matrizen, siehe Ü 2.1.9. Hier kann durch eine Präzisierung der Problemstellung bzw. durch Hinzunahme von Bedingungen, welche die Eindeutigkeit garantieren, zu korrekten Aufgabenstellungen gelangt werden.

Im letzten Fall, d. h. bei Unstetigkeit von P in e, muß die Problemstellung durch Veränderung von \mathcal{E} bzw. P modifiziert und damit Stetigkeit erzwungen werden. Man spricht dann von einer *Regularisierung* des inkorrekt gestellten Problems.

2.1.11. Beispiel. Für $x \in \mathscr{E} := \mathbb{R}$ sei $x^+ := P(x)$ definiert als

$$P(x) = egin{cases} 1/x & ext{für} & x \neq 0, \ 0 & ext{für} & x = 0. \end{cases}$$

Für x = 0 und $\delta x \neq 0$ gilt $|P(x + \delta x) - P(x)| = 1/|\delta x|$, d. h., eine Störung $|\delta x| > 0$ von x wirkt sich wie $1/|\delta x|$ auf das Ergebnis aus. Für x = 0 liegt daher eine extrem inkorrekt gestellte Aufgabe vor. Mögliche Regularisierungen sind

$$P(x, \alpha) := rac{x}{lpha^2 + x^2}$$
 bzw. $P(x, \alpha) := egin{cases} 1/x & ext{für} & |x| \ge lpha, \ 0 & ext{für} & |x| < lpha, \ \end{cases}$

wobei x > 0 ein geeignet festzulegender Regularisierungsparameter ist. \Box

Wir werden in Kapitel 8 sehen, daß die Lösung des Quadratmittelproblems $Ax \cong b$ wie auch die Bestimmung der Pseudoinversen A^+ im Sinne von 2.1.11 inkorrekt gestellte Aufgaben sind, sofern A rangdefizient ist. Zur Regularisierung werden beide angegebenen Möglichkeiten benutzt, siehe Abschnitte 11.2 und 11.3.

Übungsaufgaben

Ü 2.1.1. Zur Integration eines steifen Differentialgleichungssystems $\dot{y} = f(t, y)$, $y = y(t) = (y_1(t), ..., y_n(t))^{\mathsf{T}}$, $t \in [0, T]$, kann das als Trapezregel bezeichnete implizite Verfahren

$$y^{k+1} = y^k + 0.5h_k[f(t_{k+1}, y^{k+1}) + f(t_k, y^k)], \quad h_k > 0$$
 Schrittweite,

verwendet werden. Man zeige, daß bei Anwendung auf die lineare Differentialgleichung

$$\dot{\boldsymbol{y}} = \boldsymbol{A}(t) \, \boldsymbol{y} + \boldsymbol{g}(t), \qquad \boldsymbol{A}(.) \colon [0, T] \to \mathbf{R}^{n, n}, \qquad \boldsymbol{g}(.) \colon [0, T] \to \mathbf{R}^{n}$$

in jedem Schritt das Gleichungssystem

$$[I - 0.5h_k A(t_{k+1})] \mathbf{y}^{k+1} = \mathbf{y}^k + 0.5h_k [A(t_k) \mathbf{y}^k + \mathbf{g}(t_{k+1}) + \mathbf{g}(t_k)]$$

zu lösen ist. Dabei ist $\mathbf{y}^k \in \mathbf{R}^n$ die Näherung für $\mathbf{y}(t_k)$.

Ü 2.1.2. Zur Lösung des Randwertproblems

$$-y'' + p(x) y = q(x), \qquad y(0) = y(1) = 0,$$

für die skalare Funktion y = y(t) nach dem gewöhnlichen Differenzenverfahren wird das Erfülltsein der Differentialgleichung in den Gitterpunkten $x_i := i/(n + 1)$ (i = 1, ..., n) gefordert und $y''(x_i)$ durch die Differenzenapproximation

$$y^{\prime\prime}(x_i) pprox rac{1}{h^2} \left[y(x_{i-1}) - 2y(x_i) + y(x_{i+1})
ight]$$

ersetzt, wobei h := 1/(n + 1) ist. Man zeige, daß dadurch das lineare Gleichungssystem

$$\begin{split} -y_{i-1} + (2 + h^2 p_i) \, y_i - y_{i-1} &= h^2 q_i \qquad (i = 1, \, ..., \, n), \\ y_0 &= y_{n+1} = 0, \qquad p_i = p(x_i), \qquad q_i = q(x_i) \end{split}$$

für die Näherungswerte $y_i\approx y(x_i)$ mit einer symmetrischen und tridiagonalen Koeffizientenmatrix entsteht.

Ü 2.1.3. Der Ansatz $y = \sum_{j=1}^{n} x_j \varphi_j(t)$ mit *n* Ansatzfunktionen $\varphi_j(t)$ und freien Parametern x_j soll so an $m \ge n$ Meßwertpaare (t_i, y_i) (i = 1, ..., m) angepaßt werden, daß

$$\sum_{i=1}^{m} \left| y_i - \sum_{j=1}^{n} x_j \varphi_j(t_i) \right|^2 \to \underset{x_1, \dots, x_n}{\operatorname{Minimum}!}$$

Man schreibe diese Aufgabe als lineares Quadratmittelproblem $Ax \simeq b$. Wie ändert sich A bzw. b beim ändern, Streichen oder Hinzufügen einer Ansatzfunktion bzw. eines Meßwertpaares?

Ü 2.1.4. Es seien $A \in \mathbb{R}^{m,n}$, $B = (b^1, ..., b^q) \in \mathbb{R}^{m,q}$ sowie $X = (x^1, ..., x^q) \in \mathbb{R}^{n,q}$ gegeben. Man zeige, daß X genau dann $||AX - B||_F$ minimiert, wenn die Spalten x^k die Quadratmittelprobleme $Ax \simeq b^k$ (k = 1, ..., q) lösen.

Ü 2.1.5. Es sei $A \in \mathbb{R}^{n,n}$ eine diagonalähnliche Matrix mit den Eigenwerten $\{\lambda_j\}$ und den linear unabhängigen Eigenvektoren $\{s^j\}$. Man beweise, daß $y(t) := \sum_{j=1}^n a_j e^{\lambda_j t} s^j$ für beliebige Koeffizienten $a_j \in \mathbb{C}$ die Differentialgleichung

$$\dot{y} = Ay$$

löst und daß $\boldsymbol{y}(t) := e^{At}\boldsymbol{y}^0$ mit $e^{At} := \boldsymbol{S} e^{At}\boldsymbol{S}^{-1}$, $e^{At} := \text{diag} (e^{\lambda jt})$, $\boldsymbol{S} := (\boldsymbol{s}^1, ..., \boldsymbol{s}^n)$ die Lösung zum Anfangswert $\boldsymbol{y}(0) = \boldsymbol{y}^0$ ist.

 $\ddot{\mathbf{U}}$ 2.1.6. Gegeben sei das folgende, unter dem Einfluß der Schwerkraft stehende schwingfähige System.



Die Bewegung der Punktmassen m_1, m_2 sei auf die vertikale z-Achse eingeschränkt, und k_i, l_i bezeichne Federkonstante und Länge der Feder *i*. Man überlege sich:

(i) Die Bewegungsgleichungen des Systems lauten

$$M\ddot{z} = -Kz + p$$

 mit

$$m{z} = inom{z_1}{z_2}, \ \ m{p} = inom{m_1g + k_1l_1 - k_2l_2}{m_2g + k_2l_2}, \ \ m{M} = inom{m_1 \ 0}{0 \ m_2}, \ \ m{K} = inom{k_1 + k_2 \ -k_2}{-k_2 \ k_2},$$

wobei M und K symmetrisch und positiv definit sind.

(ii) Die Gleichgewichtslage z^0 ergibt sich aus dem linearen Gleichungssystem Kz = p. (iii) Die Auslenkungen $u := z - z^0$ aus der Gleichgewichtslage genügen der Differentialgleichung

$$M\ddot{u} = -Ku$$
,

und die Eigenfrequenzen ω_i von periodischen Lösungen $u = e^{i\omega t}x$ ergeben sich aus dem allgemeinen symmetrischen Eigenwertproblem

$$\omega^2 M x = K x.$$

Ü 2.1.7. Es sei $\{e, P\}$ ein numerisches Problem, und $P: \mathbb{R}^N \to \mathbb{R}^M$ besitze in der Umgebung von e stetige partielle Ableitungen $\partial P_i/\partial e_j$. Man überlege sich:

(i) Das Problem $\{e, P\}$ ist lokal lipschitzstetig, und es gilt

$$\|\boldsymbol{P}(\tilde{\boldsymbol{e}}) - \boldsymbol{P}(\boldsymbol{e})\| \leq \sum_{j=1}^{N} L_j(\boldsymbol{e}, \Delta \boldsymbol{e}) |\tilde{e}_j - e_j|$$

 \mathbf{mit}

$$L_j(\boldsymbol{e}, \varDelta \boldsymbol{e}) := \max \left\{ \left\| \frac{\partial \boldsymbol{P}}{\partial e_j} \left(\boldsymbol{e}_j \right) \right\| : \| \tilde{\boldsymbol{e}} - \boldsymbol{e} \| \leq \varDelta \boldsymbol{e}
ight\}.$$

Insbesondere sind die Zahlen $L_j(e) := \left\| \frac{\partial P}{\partial e_j}(e) \right\|$ die partiellen Konditionszahlen von $\{e, P\}$.

(ii) Besitzt **P** stetige partielle Ableitungen zweiter Ordnung, so ist

$$P(\hat{e}) - P(e) = \sum_{j=1}^{N} \left[\frac{\partial P}{\partial e_j}(e) \right] (\tilde{e}_j - e_j) + O(||\tilde{e} - e||^2).$$

Dabei ist $\frac{\partial \boldsymbol{P}}{\partial e_j}(\boldsymbol{e}) = \left(\frac{\partial P_1}{\partial e_j}, \dots, \frac{\partial P_M}{\partial e_j}\right)^{\mathsf{T}}$ der Vektor der partiellen Ableitungen von \boldsymbol{P} nach e_j .

Ü 2.1.8. Die Zahl $e \neq 0$, $e \in \mathbb{R}$, sei im Rahmen des relativen Fehlerniveaus $\sigma > 0$ gegeben, d. h., e repräsentiere alle Zahlen der Menge $\mathscr{E}_r := \mathscr{E}_r(e, \sigma) := \{\tilde{e} \in \mathbb{R} : |(\tilde{e} - e)/e| \leq \sigma\}$. Man zeige, daß alle zulässigen $\tilde{e} \in \mathscr{E}_r$ genau dann von 0 verschieden sind, wenn $\sigma < 1$ gilt.

Im Fall $\sigma < 1$ kann daher auch

$$arepsilon':=arepsilon'(e):=rac{ ilde{e}-e}{ ilde{e}} ext{ statt } arepsilon=arepsilon(e)=rac{ ilde{e}-e}{e}$$

als relativer Fehler von e definiert werden. Man rechne nach, daß

$$arepsilon' = rac{arepsilon}{1+arepsilon} \;\; ext{ sowie } \;\; |arepsilon'| \leq rac{|arepsilon|}{1-|arepsilon|} \leq rac{\sigma}{1-\sigma} =: \sigma'$$

für alle $\tilde{e} \in \mathcal{E}_r$ gilt. Für kleines relatives Fehlerniveau stimmen daher ε' und ε sowie die zugehörigen Schranken σ' und σ praktisch überein.

Ü 2.1.9 (PARLETT [80a]). Gegeben ist die symmetrische Matrix

$$A(arepsilon):=egin{pmatrix} 1+arepsilon\cos(2/arepsilon)&arepsilon\sin(2/arepsilon)\ arepsilon\sin(2/arepsilon)&1-arepsilon\cos(2/arepsilon)\end{pmatrix}$$

für $\varepsilon \geq 0$. Man überlege sich:

(i) $A(\varepsilon)$ hängt stetig von ε ab und besitzt die Eigenwerte $\lambda_1(\varepsilon) = 1 + \varepsilon$, $\lambda_2(\varepsilon) = 1 - \varepsilon$. Die zugehörigen orthonormalen Eigenvektoren sind im Fall $\varepsilon > 0$ bis auf das Vorzeichen eindeutig bestimmt und lauten

$$\boldsymbol{u}^{1}(\varepsilon) = (\cos(1/\varepsilon), \sin(1/\varepsilon))^{\mathsf{T}}, \quad \boldsymbol{u}^{2}(\varepsilon) = (\sin(1/\varepsilon), -\cos(1/\varepsilon))^{\mathsf{T}}.$$

(ii) Im Fall $\varepsilon = 0$ hat $A(\varepsilon)$ den doppelten Eigenwert $\lambda_1(0) = \lambda_2(0) = 1$, aber die Eigenvektoren $u^1(\varepsilon)$, $u^2(\varepsilon)$ haben für $\varepsilon \to +0$ keinen Grenzwert. Jedes $u \in \mathbb{R}^2$ mit $||u||_2 = 1$ ist ein normierter Eigenvektor von A(0) und Häufungspunkt der Menge $\{u = u^1(\varepsilon), u^2(\varepsilon) : \varepsilon > 0\}$.

2.2. Computerarithmetik

A. Numerisches Rechnen

Das derzeitige numerische Rechnen beruht fast ausschließlich auf der seit Jahrhunderten bekannten Entwicklung von reellen Zahlen nach Potenzen einer Basis und der Anordnung der Entwicklungskoeffizienten als Ziffern im zugehörigen Stellenwertsystem. Beispiele für endliche sowie unendliche Potenzentwicklungen zur Basis 10 mit den Koeffizienten bzw. Ziffern 0, 1, ..., 9 sind

$$\begin{array}{ll} 5 \times 10^2 + 0 \times 10^1 + 7 \times 10^0 &= 507, \\ 6 \times 10^{-3} + 6 \times 10^{-4} + 6 \times 10^{-5} + \cdots &= 0.00666..., \\ 3 \times 10^0 + 1 \times 10^{-1} + 4 \times 10^{-2} + 1 \times 10^{-3} + \cdots &= 3.141... \end{array}$$

Neben der übersichtlichen Schreibweise und den einfachen Vorschriften zur Ausführung der arithmetischen Operationen hat diese Zahldarstellung den Vorteil, daß die erste von 0 verschiedene Ziffer und ihre Stellung relativ zum Dezimalpunkt die wesentlichste Information über die Zahl enthält. Alle weiter rechts stehenden Ziffern liefern zusätzliche, jedoch weniger und weniger wesentliche Information.

Bei den meisten Rechnungen treten Zahlen unterschiedlicher Größenordnung auf, so daß sich die halblogarithmische Darstellung der Potenzentwicklung anbietet. Für die obigen Beispiele lautet diese

$0.507 imes10^3$	(= 507.),
$0.666\ldots imes10^{-2}$	(= 0.00666),
$0.3141\ldots imes10^1$	(= 3.141).

Allgemein kann jede reelle Zahl $x \neq 0$ in der Form

$$x = m_{R} \cdot 10^{e} = \pm 0.m_{1}m_{2}m_{3}... \times 10^{e}, \qquad m_{i} \in \{0, 1, ..., 9\}$$
(1)

mit der Mantisse m_R und dem ganzzahligen Exponenten e geschrieben werden, wobei die Normalisierungsbedingung $m_1 \neq 0$ als erfüllt vorausgesetzt wird.

Die numerische Praxis erfordert die Beschränkung auf Potenzentwicklungen mit endlicher Stellenzahl. Da dann i. allg. nicht mehr alle Ziffern berücksichtigt werden können, entstehen Fehler.

2.2.1. Bemerkung. Numerisches Rechnen ist durch die Einschränkung auf endliche Stellenzahl gekennzeichnet. Dadurch werden i. allg. sowohl bei der Darstellung von Zahlen als auch bei der Ausführung arithmetischer Operationen Fehler erzeugt, die Darstellungs- bzw. Rundungsfehler genannt werden.

Mit wachsender Länge der Potenzentwicklung wächst die Genauigkeit, aber auch der Rechenaufwand, d. h. der zum Speichern der Zahlen benötigte Speicherplatz und die Rechenzeit für die Ausführung arithmetischer und anderer Operationen. Die Länge der benutzten Potenzentwicklung soll daher ein vernünftiger Kompromiß zwischen Genauigkeitsanforderungen und Rechenkosten sein. Handrechnung erlaubt die flexible Anpassung der Stellenzahl an den Rechengang. Computerrechnung ist Rechnung mit fester Stellenzahl, die für übliche Genauigkeitsforderungen meist höher als eigentlich nötig ist. Andererseits kann wegen der hohen Zahl der Rechenoperationen der summarische Einfluß der Rundungsfehler wesentlich größer werden als bei Handrechnung.

2.2.2. Bemerkung. Die Einschätzung der Genauigkeit von Zahlen, die Ergebnisse von Computerrechnungen sind, erfordert die Kenntnis der wesentlichen Eigenschaften der Computerarithmetik und eine Fehleranalyse des benutzten Verfahrens.

Die Arithmetiken der derzeitigen Computertypen unterscheiden sich in vielen Einzelheiten. Ihre wesentlichen Eigenschaften sind jedoch dieselben und können durch wenige Kenngrößen charakterisiert werden. Dies ist Gegenstand der nachfolgenden Ausführungen; als Beispiel gehen wir auf die Arithmetik der Rechner der ESER-Reihe und der Reihe IBM 370 und Nachfolger ein. Grundlegende Eigenschaften von numerischen Algorithmen und Elemente der Fehleranalyse werden im nachfolgenden Abschnitt 2.3 behandelt.

B. Computerzahlen

Ausgangspunkt für die Darstellung einer Computerzahl ist die halblogarithmische Darstellung (1), wobei jedoch nur eine feste Anzahl von Mantissenstellen und ein fester Exponentenbereich zugelassen werden. Aus Gründen der effektiven Realisierung der Hardware arbeiten die meisten Computer dabei nicht mit der Basis 10, sondern mit einer Zweierpotenz als Basis. Üblich sind die Basen 2, 8 oder 16; man spricht von binärer, oktaler oder hexadezimaler Arithmetik.

2.2.3. Begriffe. Es seien β , t, E_1 und E_2 positive ganze Zahlen, $\beta \ge 2$. Dann besteht die Menge $\Re = \Re(\beta, t, E_1, E_2)$ der Computerzahlen (auch: Gleitpunktzahlen) zur Basis β mit der Mantissenlänge t und dem Exponentenbereich $[-E_1, E_2]$ aus allen Zahlen der Gestalt

 $x = m \cdot \beta^{e}$.

Dabei ist $e - \det Exponent$ von $x - \operatorname{eine}$ beliebige ganze Zahl mit $-E_1 \leq e \leq E_2$, und m — die Mantisse von x — ist ein beliebiger t-stelliger Bruch zur Basis β der Form

$$m = \pm 0.m_1m_2...m_t = \pm (m_1\beta^{-1} + m_2\beta^{-2} + \cdots + m_t\beta^{-t})$$

 $m = \pm 0.m_1 m_2 \dots m_t - \pm (m_1 p - 1) m_2 p$ mit den Ziffern $m_i \in \{0, 1, \dots, \beta - 1\}$ $(i = 1, \dots, t)$, wobei folgendes gilt: entweder $m_1 \ge 1$ oder $m_1 = m_2 = \dots = m_t = 0$ und $e = -E_1$. (2)

Die Alternative (2) besagt, daß für jede Computerzahl $x \neq 0$ die erste Mantissenstelle m_1 von 0 verschieden ist; man spricht von normalisierter Gleitpunktdarstellung. Die Normalisierungsbedingung ist nur für die Computerzahl x = 0 verletzt; für diese gilt m = 0. Die Zuordnung des kleinsten Exponenten $e = -E_1$ zur Computerzahl x = 0 ist üblich und zweckmäßig, aber nicht zwingend.

Für den Exponentenbereich $[-E_1, E_2]$ gilt meist $E_1 = E_2 + 1$ oder $E_1 = E_2$, d. h., er ist fast symmetrisch.

Zur Darstellung einer Zahl im Computer wird in der Regel ein *Computerwort* benutzt, das aus einer bestimmten Anzahl von Bits besteht; gelegentlich werden auch mehrere Computerworte verwendet. Dabei ist ein *Bit* die Einheit der binären Information und kann die Werte 0 oder 1 repräsentieren; jeweils 8 Bit werden als ein *Byte* bezeichnet.

2.2.4. Beispiel. Computer der ESER- und IBM-Reihen rechnen zur Basis $\beta = 16$ (hexadezimale Darstellung mit den Ziffern

 $m_i \in \{0, 1, \dots, 9, 10, \dots, 15\} = \{0, \dots, 9, A, B, C, D, E, F\}$

Die Wortlänge ist 4 Byte, also 32 Bit.

(i) In einfacher Genauigkeit (REAL*4 bzw. REAL in FORTRAN) wird ein Wort zur Zahldarstellung verwendet mit der folgenden Aufteilung:



Dabei ist v das Vorzeichen von m (v = 0 für $m \ge 0$, v = 1 für m < 0). Die Zahl c := e + 64 ist die Charakteristik von x und wird mit 7 bit binär codiert, so daß $0 \le c \le 2^7 - 1$, also $-64 \le e \le 63$ gilt. Die 6 hexadezimalen Mantissenziffern werden binär in den Positionen 9 bis 32 codiert. Die Menge der einfachgenauen Computerzahlen ist daher $\Re(16, 6, 64, 63)$.

(ii) In doppelter Genauigkeit (REAL*8 bzw. DOUBLE PRECISION in FOR-TRAN) wird ein Doppelwort zur Zahldarstellung verwendet. Das linke entspricht dem bei einfachgenauer Darstellung, das rechte nimmt in den Positionen 1 bis 32 die weiteren 8 hexadezimalen Mantissenziffern m_7, \ldots, m_{14} in binärer Codierung auf. Die Menge der doppeltgenauen Computerzahlen ist daher $\Re(16, 14, 64, 63)$.

Aus 2.2.3 folgt für jedes $x \in \mathbf{R}$, $x \neq 0$, die Beziehung $\beta^{-1} \leq |m| \leq 1 - \beta^{-t}$, also

$$\beta^{e-1} \leq |x| \leq \beta^e (1 - \beta^{-t}) < \beta^e.$$
(3)

Für zwei benachbarte, von 0 verschiedene Computerzahlen x, x' ergibt sich damit

$$\beta^{-t} < |x - x'|/|x| \leq \beta^{1-t}.$$

Schließlich zeigt (3), daß für jedes $x \in \Re$

entweder x = 0 oder $\beta^{-E_1-1} \leq |x| < \beta^{E_2}$

gilt. Zusammenfassend erhalten wir die folgende Aussage:

2.2.5. Eigenschaften der Computerzahlen. Die Computerzahlen aus $\Re(\beta, t, E_1, E_2)$ bilden im Intervall [-MAX, +MAX], MAX := β^{E_1} , ein endliches, symmetrisches und ungleichmäßiges Gitter. Das Intervall [-MIN, +MIN], MIN := β^{-E_1-1} , ent-

hält nur die drei Computerzahlen -MIN, 0, +MIN. Die relative Feinheit des Gitters in den Intervallen [-MAX, -MIN] und [MIN, MAX] genügt den Ungleichungen

$$eta^{-t} < rac{|x-x'|}{|x|} \leq eta^{1-t}$$
 ,

wobei x, x' zwei benachbarte Computerzahlen aus diesen Intervallen sind.

2.2.6. Beispiel. Für ESER- und IBM-Computer gilt

 $MIN = 16^{-65} = 5.4 \times 10^{-79}, MAX = 16^{63} = 7.2 \times 10^{75}$

in einfacher wie doppelter Genauigkeit. Die relative Feinheit ist durch die Ungleichungen

$$5.9 imes 10^{-8} = 16^{-6} < |x - x'|/|x| \le 16^{-5} = 9.5 imes 10^{-7}$$

in einfacher Genauigkeit bzw.

$$1.4 imes 10^{-17} = 16^{-14} < |x - x'|/|x| \le 16^{-13} = 2.2 imes 10^{-16}$$

in doppelter Genauigkeit charakterisiert.

Für reale Computer ist β^{-t} stets eine sehr kleine Größe. In Abschätzungen werden wir daher Summanden der Größenordnung β^{-t} gegenüber solchen der Größenordnung 1 vernachlässigen; z. B. schreiben wir statt

$$|\lambda| \leq 1.5 + 4 \cdot eta^{-t}$$

einfach

$$|\lambda| \leq 1.5$$
,

d. h., $\gamma \leq \delta$ bedeutet $\gamma \leq \delta(1 + \varrho)$ mit $\varrho = O(\beta^{-t})$. Analog wird das Zeichen "=" eingeführt. Später werden wir in solchen Abschätzungen der Einfachheit halber auch den Punkt weglassen.

C. Darstellung reeller Zahlen durch Computerzahlen

Für die Computerrechnung müssen reelle Zahlen $x \in \mathbf{R}$ durch die endlich vielen Computerzahlen aus $\Re(\beta, t, E_1, E_2)$ approximiert werden.

Im Fall |x| > MAX ist dies nicht sinnvoll möglich, weil der Exponentenbereich nicht ausreicht, um wenigstens die Größenordnung von x richtig wiederzugeben. Man spricht von *Exponentenüberlauf*, und die Rechnung wird i. allg. abgebrochen.

Im Fall $|x| \leq MAX$ muß der Zahl $x \in \mathbf{R}$ jedoch in eindeutiger Weise eine Computerzahl rd $(x) \in \Re$ zugeordnet werden. Für |x| < MIN ist dies rd (x) = 0. Im Fall 0 < |x| < MIN spricht man dann von *Exponentenunterlauf*.

Für MIN $\leq |x| \leq$ MAX sind zwei typische Zuordnungen üblich, nämlich

- Rundung (auch: symmetrische Rundung, engl. "rounding", russ. "правильное округление")
- Abbrechen (auch: unsymmetrische Rundung, engl. "chopping").

Im Fall der Rundung ist rd (x) die x am nächsten gelegene Computerzahl. Wenn zwei solche Computerzahlen existieren, wird eine von beiden — z. B. die betragsgrößere — genommen.

Beim Abbrechen ist rd (x) die x nächstgelegene Computerzahl zwischen x und 0.

Man überzeugt sich, daß im Fall MIN $\leq |x| \leq$ MAX die Computerdarstellung rd (x) der reellen Zahl x wie folgt gefunden werden kann: Man bestimmt die halblogarithmische Darstellung $x = m_R \beta^e$ von x mit $\beta^{-1} \leq |m_R| < 1$ und rundet die (normalisierte) Mantisse auf t Stellen bzw. bricht sie nach t Stellen ab. Dazu müssen die ersten t + 1 bzw. t Stellen von m_R berechnet werden.

Wir nennen die Funktion

 $rd: [-MAX, MAX] \subset \mathbf{R} \to \Re$

Rundungs- oder Darstellungsfunktion; die Größen

rd (x) – x bzw.
$$\frac{\operatorname{rd}(x) - x}{x}$$
 im Fall $x \neq 0$

heißen absoluter bzw. relativer Rundungs- oder Darstellungsfehler von x.

Zur Untersuchung dieser Rundungsfehler betrachten wir die folgenden Fälle.

Fall 1: x = 0: Dann ist rd (x) = 0, also rd (x) - x = 0. Die Zahl x = 0 wird exakt dargestellt. Fall 2: 0 < |x| < MIN: Dann ist rd (x) = 0, also

rd (x)
$$-x = -x$$
 sowie $\varepsilon(x) = \frac{\operatorname{rd}(x) - x}{x} = -1$.

Im Fall 0 < |x| < MIN ist der relative Rundungsfehler betragsmäßig gleich 100%. Die gesamte in $x \neq 0$ enthaltene Information geht bei Übergang zur Computerdarstellung rd (x) = 0 verloren. Exponentenunterlauf kann daher — muß aber nicht — eine Quelle für wesentliche Genauigkeitsverluste sein in Abhängigkeit davon, wie die durch Unterlauf entstandenen Nullen in die nachfolgenden Rechnungen eingehen.

Fall 3: MIN $\leq |x| \leq MAX$:

Aus der Definition von rd und der Abschätzung aus 2.2.5 folgt

$$|\varepsilon(x)| = \left|\frac{\operatorname{rd}(x) - x}{x}\right| \leq r := \begin{cases} 0.5\beta^{1-t} & \text{für symmetrische Rundung,} \\ \beta^{1-t} & \text{für unsymmetrische Rundung.} \end{cases}$$
(4)

Die Zahl v heißt relatives Rundungsfehlerniveau der Arithmetik, gelegentlich auch relative Maschinengenauigkeit, und ist i. allg. sehr klein. Abschätzung (4) besagt: Im Fall MIN $\leq |x| \leq MAX$ ist der relative Rundungsfehler klein. Die Rundung bewirkt nur einen geringen Verlust der in x enthaltenen Information.

Die Abschätzung (4) läßt sich äquivalent auch als

$$\operatorname{rd}(x) = x (1 + \varepsilon(x)) \quad \operatorname{mit} \quad |\varepsilon(x)| \leq \nu \tag{5}$$

schreiben und wie folgt deuten: Die Computerdarstellung rd (x) ist im Fall $MIN \leq |x| \leq MAX$ die exakte Darstellung einer wenig gestörten Zahl $x(1 + \varepsilon(x))$. Wenn $\varepsilon(x) = 0$ für x = 0 definiert wird, gilt (5) auch für x = 0.

Zusammenfassend erhalten wir:

2.2.7. Eigenschaften der Rundungsfunktion. Die Rundungsfunktion ordnet jeder reellen Zahl $x \in [-MAX, MAX]$ die Computerdarstellung rd $(x) \in \Re(\beta, t, E_1, E_2)$ zu. Dabei gilt

 $\mathrm{rd}\ (x) = x\big(1 + \varepsilon(x)\big)$

 mit

$$arepsilon(x)=0$$
 für $x=0,$
 $arepsilon(x)=-1$ für $0<|x|< ext{MIN},$

$$|\varepsilon(x)| \leq v$$
 for MIN $\leq |x| \leq MAX$

und dem relativen Rundungsfehlerniveau

$$v = \begin{cases} 0.5\beta^{1-t} & \text{für symmetrische Rundung,} \\ \beta^{1-t} & \text{für unsymmetrische Rundung.} \end{cases}$$

Im Fall $x \neq 0$ stellt $\varepsilon(x)$ den relativen Rundungsfehler von x dar.

2.2.8. Bemerkung. (i) Der Nachteil der unsymmetrischen Rundung besteht weniger in dem verdoppelten relativen Rundungsfehlerniveau, sondern mehr in der Tatsache, daß der relative Rundungsfehler $\epsilon(x)$ stets nichtpositiv ist. Damit ist die Chance einer "stochastischen Auslöschung" relativer Rundungsfehler, z. B. bei Prozessen, in denen sich die relativen Fehler addieren, nicht mehr gegeben, v gl. Abschnitt 2.3.

(ii) Als einfaches Modell für die stochastische Untersuchung des relativen Rundungsfehlers $\varepsilon = \varepsilon(x)$ nehmen wir an, daß x eine auf $[a, b] \subset [MIN, MAX]$ gleichverteilte Zufallsgröße sei, wobei (b - a)/a ein ganzzahliges Vielfaches der Basis β ist. Dann stellt ε auch eine Zufallsgröße dar, und im Fall symmetrischer Rundung gilt

$$E(arepsilon)=0+O(\mathbf{r}^2) \quad ext{sowie} \quad D^2(arepsilon)=rac{\mathbf{r}^2}{3eta}+O(\mathbf{r}^3).$$

Für unsymmetrische Rundung ist dagegen

$$E(arepsilon)=-rac{\lneta}{2(eta-1)}\,
u+O(
u^2) \hspace{1.5cm} ext{sowie} \hspace{1.5cm} D^2(arepsilon)=\left[rac{1}{3eta}-rac{1}{4}\left(rac{\lneta}{eta-1}
ight)^2
ight]
u^2+O(
u^3).$$

Für Hinweise zur stochastischen Rundungsfehleranalyse siehe B 2.4. 🗌

2.2.9. Beispiel. Bei Verwendung von FORTRAN und PL1 arbeiten ESER- und IBM-Computer mit unsymmetrischer Rundung.

In einfacher Genauigkeit ergibt sich mit $v = 9.5 \times 10^{-7}$

$$E(\varepsilon) = -8.8 \times 10^{-8}, \qquad \sqrt{D^2(\varepsilon)} = 1.1 \times 10^{-7}.$$

In doppelter Genauigkeit erhält man mit $v = 2.2 \times 10^{-16}$

$$E(arepsilon)=-2.1 imes10^{-17}$$
, $\sqrt{D^2(arepsilon)}=2.5 imes10^{-17}$,

Die statistischen Werte sind etwa eine Zehnerpotenz günstiger als die maximale Fehlerschranke r.

Hardwaremäßig ist bei den genannten Computern auch eine symmetrische Rundung realisiert, die jedoch nur über Assemblerprogramme genutzt werden kann. □

D. Computerarithmetik

Für zwei Computerzahlen $x, y \in \Re$ ist das exakte Ergebnis

 $z = x \diamondsuit y, \qquad \diamondsuit \in \{+, -, *, /\}$

der arithmetischen Operation \diamondsuit in der Regel keine Computerzahl (bei der Division muß selbstverständlich $y \neq 0$ gefordert werden).

Als Beispiel betrachten wir $\Re = \Re(10, 3, 5, 5)$ und $x = 0.718 \times 10^2$, $y = 0.591 \times 10^{-4}$. Für diese Argumente ist

$$egin{aligned} &z_1 = x + y = 0.718000591 imes 10^2, &z_3 = y * y = 0.349281 imes 10^{-8}, \ &z_2 = x * y = 0.424338 imes 10^{-2}, &z_4 = x/y = 0.121489 \ldots imes 10^7. \end{aligned}$$

Sämtliche Resultate gehören nicht zu \Re , denn in allen treten mehr als drei Mantissenstellen auf, und bei z_3 bzw. z_4 wird zusätzlich der Exponentenbereich nach unten bzw. nach oben überschritten.

Im folgenden bezeichne

fl $(x \diamondsuit y) \in \Re$

für $x, y \in \Re$ das Computerresultat der arithmetischen Operation $x \diamondsuit y$ (nach der englischen Bezeichnung "floating point arithmetic" für Gleitpunktarithmetik).

Von einer guten Computerarithmetik sollte verlangt werden, daß dieses Computerresultat die Computerdarstellung des exakten Resultates ist, sofern kein Überlauf eintritt. Im obigen Beispiel müßte dann gelten:

 ${
m fl}\,(x+y)=0.718 imes\,10^2,~~{
m fl}\,(y*y)=0~~({
m Unterlauf}),$

fl
$$(x * y) = 0.424 \times 19^{-2}$$
, fl (x/y) nicht erklärt (Überlauf).

2.2.10. Postulat. Die Computerarithmetik in $\Re = \Re(\beta, t, E_1, E_2)$ arbeite so, daß für $x, y \in \Re$ im Fall $x \diamondsuit y \in [-MAX, MAX]$

$$fl(x \diamondsuit y) = rd(x \diamondsuit y) \tag{6}$$

gilt.

2.2.11. Aussage. Unter der Voraussetzung 2.2.10 gilt für $x, y \in \Re$ mit $x \diamondsuit y = 0$ oder MIN $\leq |x \diamondsuit y| \leq MAX$

$$\mathrm{fl} \ (x \diamondsuit y) = (x \diamondsuit y) \ (1 + \varepsilon) \quad \mathrm{mit} \quad |\varepsilon| \leq \nu. \tag{7}$$

Dies folgt sofort aus 2.2.7 und besagt: Der relative Rundungsfehler des Resultats ist im Fall fl $(x \diamond y) \neq 0$ höchstens gleich v.

Falls $x' = x(1 + \varepsilon), y' = y(1 + \varepsilon)$ geschrieben wird, läßt sich (7) als

$$\mathrm{fl} \ (x\pm y)=x'\pm y', \quad \mathrm{fl} \ (x*y)=x*y'=x'*y, \quad \mathrm{fl} \ (x/y)=x'/y$$

schreiben und wie folgt interpretieren: Das Computerresultat einer arithmetischen Operation ist das exakte Resultat derselben Operation mit wenig gestörten Operanden. Bei Bedarf können für $\diamond \in \{*, \}$ die Störungen auf beide Operanden x, y verteilt werden, siehe Ü 2.2.3.

2.2.12. Bemerkung. Die Voraussetzung 2.2.10 ist für viele Computer — darunter für solche der ESER- und IBM-Reihen — erfüllt. Es gibt jedoch Computertypen, die eine schwächere Arithmetik besitzen, für die (6) nicht gilt. Dabei treten folgende Abweichungen auf:

(i) Die Darstellung (7) gilt mit der Abschätzung $|\varepsilon| \leq Kv$, wobei $K \leq 3$ ist o. ä.

(ii) Für
$$\diamondsuit \in \{+, -\}$$
 ist (7) durch

$$\mathrm{fl}\ (x\pm y)=x(1+arepsilon_x)\pm y(1+arepsilon_y),\qquad |arepsilon_x|,\,|arepsilon_y|\leq
u,$$

zu ersetzen. Hieraus folgt

$$|fl(x \pm y) - (x \pm y)| \le \nu(|x| + |y|).$$
(8)

Der Fall (ii) tritt auf, wenn die Addition nur mit t Mantissenstellen ohne Schutzstelle ausgeführt wird. Aus (8) folgt für $x \pm y \neq 0$

$$\frac{|\text{fl} (x \pm y) - (x \pm y)|}{|x \pm y|} \leq \frac{|x| + |y|}{|x \pm y|} \, r,$$

d. h., im Fall $|x \pm y| \ll |x| + |y|$ kann der relative Rundungsfehler wesentlich größer als v werden.

Es zeigt sich jedoch, daß diese Abweichungen von (6) bzw. (7) die Genauigkeit der wichtigsten numerischen Prozesse der linearen Algebra nur unwesentlich beeinflussen. Wir nehmen daher im folgenden an, daß 2.2.10 erfüllt ist.

2.2.13. Bemerkung. Für $x, y \in \Re$ gelte $x = m_x \beta^{\mathfrak{e}}, y = m_y \beta^{\mathfrak{e}}, \operatorname{sgn}(x) = \operatorname{sgn}(y)$ und MIN $\leq |x - y|$. Dann ist

$$fl(x-y) = x - y, \qquad (9)$$

d. h., die Differenz zweier Computerzahlen desselben Vorzeichens und mit demselben Exponenten ist exakt. Falls dabei mehrere führende Stellen von m_x und m_y übereinstimmen, spricht man von Auslöschung.

Dies folgt sofort aus dem Postulat 2.2.10, denn unter den angegebenen Voraussetzungen ist $x - y \in \Re$, also fl (x - y) = x - y. Die Arithmetik aller derzeitiger Computer arbeitet jedoch so, daß 2.2.13 auch dann richtig bleibt, wenn 2.2.10 z. B. für Addition und Subtraktion nicht erfüllt ist.

Ein Beispiel für Auslöschung in $\Re(10, 8, 20, 20)$ ist

fl $(0.98765432 \times 10^{1} - 0.98765321 \times 10^{1}) = 0.11100000 \times 10^{-4}$.

Auslöschung kann auch bei um 1 verschiedenen Exponenten auftreten, wie durch

fl $(0.10000000 \times 10^{1} - 0.99990000 \times 10^{0}) = 0.10000000 \times 10^{-3}$

gezeigt wird. Allgemein ist das Auftreten von Auslöschung durch die Bedingung $|x \pm y| \ll |x| + |y|$ charakterisiert.

Bei Auslöschung gehen führende Mantissenstellen verloren. Dies kann — aber muß nicht — zu einem wesentlichen Verlust der in x, y enthaltenen Information führen.

Zur Erläuterung dieser Feststellung fassen wir x, y aus 2.2.13 als Computerdarstellungen der reellen Zahlen $\tilde{x}, \tilde{y} \in \mathbf{R}$ mit $\tilde{x} - \tilde{y} \neq 0$ auf. Dann ist $x = \operatorname{rd}(\tilde{x})$ $= \tilde{x}(1 + \varepsilon_x), y = \operatorname{rd}(\tilde{y}) = \tilde{y}(1 + \varepsilon_y)$, wobei $\varepsilon_x, \varepsilon_y$ mit $|\varepsilon_x|, |\varepsilon_y| \leq v$ die relativen Darstellungsfehler bezeichnen. Unter Beachtung von (9) folgt

$$\frac{|\mathrm{fl}\,(x-y)-(\tilde{x}-\tilde{y})|}{|\tilde{x}-\tilde{y}|} = \frac{|(x-y)-(\tilde{x}-\tilde{y})|}{|\tilde{x}-\tilde{y}|} = \frac{|\varepsilon_x\tilde{x}-\varepsilon_y\tilde{y}|}{|\tilde{x}-\tilde{y}|} \leq \frac{|\tilde{x}|\,+\,|\tilde{y}|}{|\tilde{x}-\tilde{y}|}\,\mathsf{r}$$

Bei Auslöschung kann sich daher die Ungenauigkeit der Daten mit dem großen Faktor $(|\tilde{x}_i + |\tilde{y}|)/|\tilde{x} - \tilde{y}|$ auf die Differenz übertragen, obwohl fl (x - y) fehlerfrei aus x, y berechnet wird!

Auslöschung kann auch bei der Auswertung komplizierterer Formeln auftreten und läßt sich in manchen Fällen durch geschickte Umformungen vermeiden. Beispiele sind

$$1 - \sqrt{1 - x} = rac{(1 - \sqrt{1 - x})(1 + \sqrt{1 - x})}{1 + \sqrt{1 - x}} = rac{x}{1 + \sqrt{1 - x}}, \quad |x| ext{ klein,}$$

oder

$$1 - \cos x = 1 - \left(1 - 2\left(\sin \frac{x}{2}\right)^2\right) = 2\left(\sin \frac{x}{2}\right)^2$$
, |x| klein

Eine weitere fehlerfrei ausgeführte Operation ist die Multiplikation bzw. Division mit einer Potenz von β .

2.2.14. Bemerkung. Für $x, y \in \Re$ und eine Potenz β^k , k ganz, gilt

$$\mathrm{fl}\;(x*\beta^{k})=x*\beta^{k},\;\;\mathrm{fl}\;(y/\beta^{k})=y/\beta^{k},\;\;\mathrm{falls}\;\mathrm{MIN}\;{\leq}\left\{\begin{matrix}|x\beta^{k}|\\|y/\beta^{k}|\end{matrix}\right\}\leq\mathrm{MAX}$$

In manchen numerischen Prozessen der linearen Algebra treten neben den arithmetischen Operationen auch Quadratwurzelberechnungen auf. Die entsprechende Computerfunction bezeichnen wir mit sqrt — gelegentlich auch mit fl (\sqrt{x}) — und fordern:

Für
$$x \in \Re$$
, $x \ge 0$, gelte sqrt $(x) = \operatorname{rd}(\sqrt[]{x})$. (10)

2.2.15. Bemerkung. (i) Auf den meisten Computern ist über eine höhere Programmiersprache wie FORTRAN oder PL1 auch eine einfach- und doppeltgenaue komplexe Arithmetik nutzbar. Da wir ausschließlich mit reellen Eingangsdaten versehene Probleme und im Reellen arbeitende Algorithmen betrachten, gehen wir auf diese Arithmetik nicht ein.

(ii) Fast alle Computer verfügen außer über die beschriebene Gleitpunktarithmetik noch über eine ganzzahlige, rundungsfehlerfreie Arithmetik für Indexrechnungen o. ä., die auf numerische Prozesse keinen wesentlichen Einfluß hat.

Übungsaufgaben

Ü 2.2.1. Man überlege sich, daß die Menge $\Re(\beta, t, E_1, E_2)$ aus $2(E_1 + E_2 + 1)\beta^{t-1}(\beta - 1) + 1$ verschiedenen Zahlen besteht.

Ü2.2.2.Es ist zu zeigen, daß das relative Rundungsfehlernive
au ν unter der Voraussetzung 2.2.10 durch

$$\nu = \min \{\nu' \in \Re: \text{fl} (1 + \nu') > 1\}$$
(11)

charakterisiert werden kann. Man schreibe ein FORTRAN-Programm, das ν auf der Grundlage von (11) mit einem relativen Fehler von höchstens 0.1 berechnet.

Ü 2.2.3. Man zeige, daß im Fall $|\varepsilon| \leq v$

$$1 + \varepsilon = \frac{1}{1 + \varepsilon_1} = \frac{1 + \varepsilon_2}{1 + \varepsilon_3} = (1 + \varepsilon_4)^2$$
(12)

mit $\varepsilon_1 = -\varepsilon$, $\varepsilon_2 = -\varepsilon_3 = \varepsilon/2$, $\varepsilon_4 = \varepsilon/2$, also

$$|\varepsilon_1| \leq \nu, \quad |\varepsilon_2|, |\varepsilon_3|, |\varepsilon_4| \leq \nu/2 \tag{13}$$

gilt. Für die Operationen * und / läßt sich dann (7) in der Gestalt

fl
$$(x * y) = [x(1 + \varepsilon)] * y = x * [y(1 + \varepsilon)] = [x(1 + \varepsilon_4)] * [y(1 + \varepsilon_4)]$$

bzw.

fl
$$(x/y) = [x(1 + \varepsilon)]/y = x/[y(1 + \varepsilon_1)] = [x(1 + \varepsilon_2)]/[y(1 + \varepsilon_3)]$$

mit den Schranken (13) schreiben, d. h., die Störungen können nach Wunsch auf jeden der einzelnen bzw. auf beide Operanden verteilt werden.

Ü 2.2.4. In $\Re(10, 4, 10, 10)$ seien die Zahlen $x = 0.2345 \times 10^{-2}$, $y = 0.6789 \times 10^{9}$, $z = 0.1234 \times 10^{2}$ gegeben. Man überprüfe, daß bei Rechnung mit unsymmetrischer Rundung

$${
m fl}\left({
m fl}\left(x+y
ight)+z
ight)=0.1302 imes10^2\pm{
m fl}\left(x+{
m fl}\left(y+z
ight)
ight)=0.1301 imes10^2$$

sowie

$$fl(x * fl(y + z)) = 0.3050 \times 10^{-1} \neq fl(fl(x * y) + fl(x * z)) = 0.3052 \times 10^{-1}$$

gilt. In Gleitpunktarithmetik gelten also weder das Assoziativgesetz noch das Distributivgesetz. **Ü** 2.2.5. Bei der Berechnung von $||v||_2 - |v_1|$ für $v = (v_1, ..., v_n)^{\mathsf{T}} \in \mathbb{R}^n$ tritt im Fall $|v_1| \gg |v_j|$ (j = 2, ..., n) Auslöschung auf. Man überlege sich, daß dieser Nachteil bei Berechnung gemäß

$$||v||_2 - |v_1| = rac{\sum\limits_{j=2}^n |v_j|^2}{||v||_2 + |v_1|}$$

vermieden wird.

2.3. Numerische Algorithmen und Grundlagen der Fehleranalyse

Die Lösung des numerischen Problems $\{e, P\}$ aus der Klasse $\{\mathcal{E}, P\}$ erfordert die Berechnung der zu e gehörenden Ausgangsdaten $a^* = P(e)$. Dazu wird ein numerischer Algorithmus benötigt, d. h. eine Vorschrift $V_0: \mathcal{E} \subset \mathbb{R}^N \to \mathbb{R}^M$, die jedem Satz $e \in \mathcal{E}$ von Eingangsdaten unter Ausführung endlich vieler arithmetischer Operationen $\{+, -, *, /\}$ in eindeutiger Weise einen Satz $a^* = V_0(e)$ von Ausgangsdaten zuordnet. Ein numerischer Algorithmus ist also eine spezielle Datentransformation von $e = (e_1, \ldots, e_N)^T$ in $a^* = (a_1^*, \ldots, a_M^*)^T$, bei der nur endlich viele arithmetische Operationen erlaubt sind. Die Forderung nach der Endlichkeit der Zahl der Daten in endlicher Zeit nur endlich viele Operationen mit endlich vielen Operanden ausgeführt werden. Es ist zweckmäßig, außer den arithmetischen Operationen auch die Berechnung von elementaren Funktionen — in der linearen Algebra speziell die von \sqrt{x} — als Basisoperation zuzulassen.

Üblicherweise wird ein numerischer Algorithmus als Programm in einer höheren Programmiersprache wie FORTRAN, PL1 usw. angegeben. Da solche Programme i. allg. jedoch unübersichtlich und schwer lesbar sind, beschreiben wir Algorithmen in traditioneller mathematischer Schreibweise oder in einer sich selbst erklärenden Pseudo-Programmiersprache, wodurch der wesentliche mathematische Gehalt deutlicher sichtbar wird. Zwischen einer solchen kompakten Beschreibung und einem zuverlässigen Computerprogramm zur Lösung der entsprechenden Aufgaben liegt jedoch noch ein großer Schritt, vgl. Kapitel 16.

Numerische Algorithmen werden im Bereich der reellen Zahlen und exakten Rechenoperationen entwickelt. Bei der Realisierung auf einem Computer mit dem Zahlenbereich \Re und der zugehörigen Computerarithmetik fl müssen jedoch alle Eingangsdaten und Zwischenergebnisse durch ihre Computerdarstellung und alle arithmetischen Operationen durch ihre Computerrealisierungen ersetzt werden. Bei iterativen Algorithmen, die theoretisch unendlich viele Iterationsschritte erfordern, um die gesuchte Lösung als Grenzwert der Iteriertenfolge zu liefern, muß außerdem ein geeignetes Abbruchkriterium vorhanden sein, das nach endlich vielen Schritten eine genügend genaue Näherung akzeptiert. Der konzeptionelle Algorithmus V_0 geht damit in seine Computerrealisierung $V : \mathfrak{G} \subset \mathfrak{R}^N \to \mathfrak{R}^M$ über. Mit \mathfrak{R}^N ist hier die Menge aller N-dimensionalen Vektoren mit Computerzahlen aus \mathfrak{R} als Komponenten bezeichnet worden, und $\mathfrak{G} \subset \mathfrak{R}^N$ ist der Definitionsbereich der Computerrealisierung V.

Bei der Ausführung von V auf einem Computer entstehen notwendig Rundungsfehler, so daß selbst bei exakt darstellbaren Eingangsdaten $e \in \mathfrak{E} \subset \mathfrak{R}^N$ und einem
endlichen konzeptionellen Algorithmus ein von a^* verschiedener Wert $a = V(e) \in \Re^M$ berechnet wird. Die Erfassung der bei der Realisierung von V_0 – d. h. beim Übergang von V_0 zu V — erzeugten Rundungsfehler ist Gegenstand der Fehleranalyse des Algorithmus V_0 . Das Verhalten gegenüber Rundungsfehlern ist ein wesentliches Charakteristikum eines numerischen Algorithmus.

Neben dem Fehlerverhalten spielt der Aufwand, d. h. die Anzahl der auszuführenden arithmetischen Operationen und die Anzahl der für das Zwischenspeichern von Hilfsgrößen erforderlichen Speicherplätze S bei der Einschätzung von Algorithmen eine entscheidende Rolle. Wir werden später sehen, daß in der numerischen linearen Algebra Operationen des Typs r := s + u * v, also eine Multiplikation, der eine Addition folgt, am häufigsten vorkommen. Diese Kombination soll mit "opms" bezeichnet werden. Für eine einzelne Addition/Subtraktion, Multiplikation/Division bzw. Quadratwurzelberechnung sollen die Bezeichnungen "ops", "opm" bzw. "opr" verwendet werden.

Im folgenden werden wir einige einfache Basisalgorithmen der linearen Algebra genauer untersuchen und dabei wesentliche Begriffe und Prinzipien der Fehleranalyse herausarbeiten. Dabei wird generell vorausgesetzt, da β bei der Realisierung weder Über- noch Unterlauf eintritt.

A. Die Summe

Aufgabe: Für die Eingangsdaten $x_i \in \Re$ (i = 1, ..., n) ist die Summe

$$z^* = \sum_{i=1}^n x_i$$

zu berechnen.

2.3.1. Algorithmus zur Summation von n Zahlen.

z := 0

for i := 1(1)n do $z := z + x_i$

Aufwand: n ops (Wir geben nur solche Aufwandsgrößen an, die mindestens proportional zu n sind.)

Hier und im folgenden soll das Zeichen "=" bedeuten, daß der rechtsstehende Ausdruck in der jeweils betrachteten Computerarithmetik fl in der durch Klammern bzw. die Vorrangsregeln festgelegten Reihenfolge ausgewertet und das Computerresultat der linksstehenden Variablen als Wert zugewiesen wird. Die Anweisung

r := s + u * v

ist daher im Sinne von

 $r = \mathrm{fl}\left(s + \mathrm{fl}\left(u * v\right)\right)$

zu verstehen. Wir schreiben auch kurz

$$r = \mathrm{fl} \left(s + u * v \right)$$

da auf Grund der Vorrangregeln klar ist, daß das Produkt fl (u * v) zuerst berechnet werden muß. Analog wird bei komplizierteren Ausdrücken verfahren.

$$z_0 = 0,$$

 $z_i = \text{fl}(z_{i-1} + x_i) \qquad (i = 1, ..., n)$
(1)

beschreiben, und $z = z_n$ ist der berechnete Wert für die exakte Summe z^* .

Zur Erfassung der Rundungsfehler wenden wir 2.2.11 auf die Vorschrift (1) an und erhalten

$$z_i = (z_{i-1} + x_i) (1 + \varepsilon_i) \quad \text{mit} \quad |\varepsilon_i| \leq \nu,$$
 (2)

wobei die bei der *i*-ten Addition gemäß (2.2.7) auftretende Störung ε mit ε_i bezeichnet worden ist.

Aus (2) ergibt sich $z_1 = x_1(1 + \varepsilon_1)$, $z_2 = x_1(1 + \varepsilon_1)(1 + \varepsilon_2) + x_2(1 + \varepsilon_2)$, allgemein

$$z_{j} = \sum_{i=1}^{j} x_{i} \left\{ \prod_{k=i}^{j} (1 + \varepsilon_{k}) \right\} \qquad (j = 1, ..., n).$$
(3)

Zur Abschätzung der Produkte $\prod (1 + \varepsilon_k)$ stellen wir die folgende Aussage bereit: **2.3.2.** Aussage. Es gelte $|\varepsilon_k| \leq \nu$ (k = 1, ..., m). Dann ist

$$\prod_{k=1}^{m} (1 + \varepsilon_k) = 1 + \varrho^{(m)} \quad \text{mit} \quad \varrho^{(m)} = \sum_{k=1}^{m} \varepsilon_k + O(\nu^2).$$
(4)
Im Fall $m\nu < 2$ gelten die Abschätzungen

$$|\varrho^{(m)}| \leq \frac{m\nu}{1 - 0.5m\nu} = m\nu.$$
⁽⁵⁾

Beweis. Aus

$$\prod_{k=1}^{m} (1+\varepsilon_k) = 1 + \sum_{k=1}^{m} \varepsilon_k + \sum_{\substack{k,l=1\\k(6)$$

ergibt sich sofort (4). Zum Nachweis von (5) beachten wir, daß die erste Summe in (6) aus m Summanden, die zweite aus $m(m-1)/(1\cdot 2)$, die dritte aus $m(m-1)(m-2)/(1\cdot 2\cdot 3)$ besteht usw. Im Fall $q := 0.5m\nu < 1$ ist daher

$$\begin{split} |\varrho^{(m)}| &\leq \frac{m\nu}{1} + \frac{m\nu(m-1)\nu}{1\cdot 2} + \frac{m\nu(m-1)\nu(m-2)\nu}{1\cdot 2\cdot 3} + \dots + \frac{m\nu(m-1)\nu\dots 2\nu 1\nu}{1\cdot 2\cdots (m-1)m} \\ &\leq m\nu(1+q+q^2+\dots+q^{m-1}) \leq m\nu(1+q+q^2+\dots) = m\nu/(1-q). \ \Box \end{split}$$

Im Fall $mv \leq 0.1$ ist $1/(1 - 0.5mv) \leq 1.06$, also $|\varrho^{(m)}| \leq 1.06mv$. Im folgenden werden wir daher stets $mv \leq 0.1$ voraussetzen und bei der Abschätzung von $\varrho^{(m)}$ die in erster Ordnung gültige einfachere Schranke mv statt der exakten Schranke 1.06mv verwenden. Für reale Dimensionen und Computer ist $m\nu$ in der Regel wesentlich kleiner, so daß dieses Vorgehen um so mehr berechtigt ist. Für $m \leq 1000$ und $\nu = 9.5 \times 10^{-7}$ (eińfachgenaue Arithmetik bei ESER- und IBM-Computern) ist z. B. $1/(1-0.5m\nu) \leq 1.0005!$

Aussage 2.3.2 besagt: Das Produkt aus m Faktoren $(1 + \varepsilon_k)$ mit $|\varepsilon_k| \leq \nu$ kann durch einen Faktor $(1 + \varrho^{(m)})$ derselben Form ersetzt werden, wobei sich die Schranke von $\varrho^{(m)}$ auf das m-fache der Schranke der ε_k erhöht. In erster Näherung ist $\varrho^{(m)}$ gleich der Summe der Störungen ε_k .

Für j = n ergibt sich aus (3) unter Beachtung von 2.3.2 und wegen $\varepsilon_1 = 0$ – die erste Addition $z_1 = z_0 + x_1 = 0 + x_1$ wird fehlerfrei ausgeführt – die folgende Fehleraussage.

2.3.3. Aussage. Die zum Fehlerniveau \dot{v} gemäß 2.3.1 berechnete Summe z läßt sich als

$$z = \sum_{i=1}^{n} x_i (1 + \varepsilon_i^{(n)}) \tag{7}$$

mit relativen Störungen

$$\varepsilon_{i}^{(n)} = \prod_{k=i}^{n} (1 + \varepsilon_{k}) - 1 = \sum_{k=i}^{n} \varepsilon_{k} + O(r^{2}), \qquad |\varepsilon_{k}| \leq \nu, \qquad (8)$$

darstellen. Die Störungen sind klein im Sinne von

$$|\varepsilon_i^{(n)}| \stackrel{<}{=} \min(n-i+1,n-1) \nu \leq (n-1) \nu, \qquad (9)$$

und der erzeugte Rundungsfehler $z^* - z$ genügt der Ungleichung

$$|z^* - z| \leq \left\{ \sum_{i=1}^n \min(n - i + 1, n - 1) |x_i| \right\} \nu \leq \left\{ (n - 1) \sum_{i=1}^n |x_i| \right\} \nu.$$
 (10)

Die Abschätzungen (10) stellen sog. *a-priori-Schranken* dar, denn sie enthalten nur die Eingangsdaten x_i der Aufgabe. Eine i. allg. wesentlich günstigere *a-posteriori-Schranke*, die mit erst im Laufe oder nach Beendigung der Rechnung verfügbaren Größen arbeitet, kann wie folgt erhalten werden: Aus (2) folgt

$$z_i = z_{i-1} + x_i + (z_{i-1} + x_i) \, \varepsilon_i = z_{i-1} + x_i + z_i \varepsilon_i / (1 + \varepsilon_i) \, .$$

Für $\delta_i = \varepsilon_i/(1 + \varepsilon_i)$ gilt $|\delta_i| \leq \nu/(1 - \nu) = \nu$, so daß Summation auf

$$z = z_n - z_0 = \sum_{i=1}^n (z_i - z_{i-1}) = \sum_{i=1}^n (x_i + z_i \delta_i) = z^* + \sum_{i=1}^n \delta_i z_i,$$

also

$$|z^* - z| \leq \frac{\nu}{\sum_{i=1}^n} |z_i| \tag{11}$$

führt. Die Schranke (11) kann leicht in derselben Schleife wie z selbst berechnet werden.

2.3.4. Bemerkungen. (i) Die linke Schranke in (10) für den erzeugten Rundungsfehler wird minimal, wenn die x_i in betragsmäßig wachsender Folge summiert werden.

(ii) Der relative Rundungsfehler $|z^* - z|/|z^*|$ kann bei Summation von Zahlen unterschiedlichen Vorzeichens beliebig groß werden.

(iii) Der relative Rundungsfehler ist bei Summation von Zahlen einheitlichen Vorzeichens klein:

$$|z^* - z|/|z^*| \leq (n-1)\nu. \square$$

Dies ergibt sich sofort aus (10) und der Tatsache, daß der Quotient $\sum |x_i|/|\sum x_i|$ bei unterschiedlichem Vorzeichen der x_i beliebig groß werden kann, bei einheitlichem Vorzeichen jedoch den Wert 1 hat.

Aus 2.3.3 folgt ferner, daß die relativen Störungen $\varepsilon_i^{(n)}$ in erster Ordnung gleich der Summe $\varepsilon_i + \varepsilon_{i+1} + \cdots + \varepsilon_n$ der relativen Rundungsfehler der einzelnen Additionen sind. Bei symmetrischer Rundung werden die ε_k unterschiedliches Vorzeichen haben und können sich auslöschen. Unter geeigneten stochastischen Voraussetzungen wird dann statt der im ungünstigsten Fall zu erwartenden Relation $|\varepsilon_i^{(n)}| = (n-i+1)\nu$ im Mittel ein Wert $|\varepsilon_i^{(n)}|$ in der Größenordnung von $\sqrt{(n-i+1)/(3\beta)} \nu$ zu erwarten sein, vgl. 2.2.8. Bei unsymmetrischer Rundung sind alle ε_k vom selben Vorzeichen, so daß der Auslöschungseffekt nicht auftreten kann.

In 2.3.4 wurde festgestellt, daß es Eingangsdaten $\{x_i\}$ gibt, die bei der Summation zu beliebig großen relativen Rundungsfehlern $|z^* - z|/|z^*|$ des Resultats führen, d. h., die berechnete Summe wird sehr ungenau. Es liegt daher die Frage nahe, ob diese Ungenauigkeit durch das verwendete Summationsverfahren oder durch andere Ursachen hervorgerufen wird. Zur Beantwortung dieser Frage erinnern wir daran, daß die Computerzahl $x \in \Re$ nicht nur sich selbst, sondern alle reellen Zahlen \tilde{x} mit $x = \operatorname{rd}(\tilde{x})$ repräsentiert, vgl. Abschnitt 2.2. Die dadurch bedingte maximale relative Unsicherheit in x kann dann für $x \neq 0$ im günstigsten Fall den Wert ν/β , im ungünstigsten den Wert ν erreichen in Abhängigkeit von der Größenordnung der Mantisse, siehe 2.2.5. Da für reale Computer β nicht groß ist und wir auch die ungünstigste Situation erfassen müssen, sehen wir etwas vereinfachend die folgende Annahme als erfüllt an:

2.3.5. Annahme. Jedes Eingangsdatenelement $x \in \Re$, $x \neq 0$, enthält eine natürliche Unsicherheit auf dem relativen Fehlerniveau ν , die durch den Darstellungsfehler hervorgerufen wird. Durch $x \in \Re$, $x \neq 0$, werden alle reellen Zahlen \tilde{x} mit

$$ilde{x} = x(1+\vartheta), \quad |\vartheta| \leq v,$$
(12)

repräsentiert.

Im Fall der Summation repräsentiert also der Eingangsdatensatz $\{x_i\}_{i=1}^n, x_i \in \Re$, alle reellen Eingangsdaten

$$\tilde{x}_i = x_i(1 + \vartheta_i) \quad \text{mit} \quad |\vartheta_i| \leq v \quad (i = 1, ..., n).$$
 (13)

Die zu $\{\tilde{x}_i\}$ gehörende exakte Summe ist $\tilde{z} = \sum_{i=1}^n \tilde{x}_i$.

Nach diesen Vorbereitungen können wir die gestellte Frage beantworten, indem wir die Abschätzungen aus 2.3.3 in zweierlei Weise interpretieren. Erste Interpretation: Unter Beachtung von (13) folgt

$$\tilde{z} = \sum_{i=1}^{n} x_i (1 + \vartheta_i) \quad \text{mit} \quad |\vartheta_i| \le \nu,$$
(14)

aus 2.3.3 liest man dagegen

$$z = \sum_{i=1}^{n} x_i (1 + \varepsilon_i^{(n)}) \quad \text{mit} \quad |\varepsilon_i^{(n)}| \leq (n-1) v \tag{15}$$

ab. Ein Vergleich von (14) und (15) führt auf das folgende Resultat:

2.3.6. Aussage. Die mittels des Verfahrens 2.3.1 berechnete Summe z ist die exakte Summe der gestörten Eingangsdaten $\{x_i(1 + \varepsilon_i^{(n)})\}$, wobei $|\varepsilon_i^{(n)}| \leq (n - 1) r$ ist. Die durch die Rundungsfehler erzeugten Störungen $\varepsilon_i^{(n)}$ vergrößern die in den Eingangsdaten enthaltenen Unsicherheiten also höchstens um den Faktor n - 1.

Zweite Interpretation: Ein durch $\{x_i\}$ repräsentierter Eingangsdatensatz $\{\tilde{x}_i\}$ liefert die exakte Summe \tilde{z} , wobei wegen (14)

$$\left|\tilde{z} - z^*\right| = \left|\sum_{i=1}^n \left(\tilde{x}_i - x_i\right)\right| = \left|\sum_{i=1}^n x_i \vartheta_i\right| \le v \sum_{i=1}^n |x_i| \tag{16}$$

gilt. Diese Abschätzung ist in bezug auf (13) scharf, d. h., es gibt Störungen $|\vartheta_i| = \nu$, für die das Gleichheitszeichen steht. Die rechte Seite von (16) gibt also das unvermeidliche absolute Fehlerniveau $\Delta z_{\min}(x_i, \nu)$ der Summe an, falls jeder Summand im Bahmen des relativen Fehlerniveaus ν gegeben ist und $\sum_{i=1}^{n} |x_i|$ ist die entspre-

im Rahmen des relativen Fehlerniveaus ν gegeben ist, und $\sum_{i=1}^{n} |x_i|$ ist die entsprechende Konditionszahl, vgl. Abschnitt 2.2.

Der Fehler von z kann nach 2.3.3 gemäß

$$|z^* - z| \leq (n-1) \nu \sum_{i=1}^n |x_i|$$
(17)

abgeschätzt werden. Der Vergleich von (16) und (17) führt auf die folgende Aussage.

2.3.7. Aussage. Der im Verfahren 2.3.1 erzeugte Rundungsfehler von z ist höchstens das (n - 1)-fache des unvermeidlichen Fehlerniveaus

$$\varDelta z_{opt}(x_i, v) = v \sum_{i=1}^n |x_i|,$$

das durch die Unsicherheit der Eingangsdaten hervorgerufen wird.

Beide Interpretationen führen uns zu dem Ergebnis, daß die in den Eingangsdaten $\{x_i\}_{i=1}^n$ enthaltene Unsicherheit durch das Verfahren 2.3.1 nur quantitativ und zwar höchstens um den Faktor F = n - 1 — vergrößert wird. Es ist jedoch wesentlich, den Unterschied zwischen beiden Interpretationen zu erkennen.

Erste Interpretation: Das technische Mittel ist die Konstruktion von Störungen $\{\varepsilon_i^{(n)}\}$ der Eingangsdaten $\{x_i\}$ derart, daß das berechnete Ergebnis gerade das exakte Ergebnis

zu den gestörten Eingangsdaten ist, und die Angabe von Schranken für diese Störungen. Dann werden die Schranken der Störungen mit dem Niveau des Darstellungsfehlers verglichen. Die entsprechende Philosophie besteht darin, die während der Rechnung entstehenden Rundungsfehler als Störungen der Eingangsdaten zu interpretieren und mit deren natürlicher Unsicherheit zu vergleichen. Dies ist die Idee der sog. Rückwärtsanalyse (engl. "backward error analysis").

Zweite Interpretation: Das technische Mittel ist die Konstruktion von möglichst realistischen Abschätzungen sowohl für das unvermeidliche Fehlerniveau $\Delta z_{opt}(x_i, v)$ des exakten Resultates z^* als auch für den durch das Verfahren erzeugten Rundungsfehler $|z^* - z|$. Dann werden diese beiden Abschätzungen miteinander verglichen. Die entsprechende Philosophie besteht darin, die durch Rechnung erzeugten Rundungsfehler in den Ausgangsdaten mit derem unvermeidlichen Fehlerniveau zu vergleichen, das durch die Unsicherheit der Eingangsdaten bedingt ist. Das ist die Idee der sog. Vorwärtsanalyse.

Auf Grund von Aussage 2.3.6 (Ergebnis der Rückwärtsanalyse) werden wir das Summationsverfahren 2.3.1 als numerisch gutartig mit der Fehlerkumulationskonstanten F = n - 1 bezeichnen, denn das berechnete Ergebnis ist das exakte Ergebnis zu gestörten Eingangsdaten, deren Störungsniveau höchstens das F-fache der Darstellungsunsicherheit erreicht. Auf Grund von Aussage 2.3.7 (Ergebnis der Vorwärtsanalyse) wird das Verfahren numerisch stabil mit der Fehlerkumulationskonstanten F = n - 1 genannt, weil der durch das Verfahren erzeugte Rundungsfehler von z höchstens das F-fache des durch die Unsicherheit der Eingangsdaten bedingten unvermeidlichen Fehlers $\Delta z_{out}(x_i, v)$ betragen kann.

Aus den bisherigen Ausführungen wird klar: Die numerische Gutartigkeit (mit nicht zu großem F) stellt die bestmögliche Qualität eines Algorithmus dar, weil dann die in den Eingangsdaten enthaltene Information bestmöglich auf die Ausgangsdaten übertragen wird. Unter den gutartigen Algorithmen sind selbstverständlich diejenigen mit den kleinsten Fehlerkumulationskonstanten bezüglich ihrer Rundungsfehlerempfindlichkeit die günstigsten, allerdings muß ein kleineres F in der Regel durch einen höheren Aufwand erkauft werden.

Für genügend kleines r und lokal lipschitzstetige Aufgabenklassen zieht die numerische Gutartigkeit die numerische Stabilität nach sich; die Umkehrung gilt i. allg. natürlich nicht. Numerische Stabilität ist eine Mindestforderung an ein vernünftiges Verfahren. Liegt sie nicht vor, so können für gewisse Probleme beliebig große Genauigkeitsverluste in bezug auf das unvermeidliche Fehlerniveau auftreten. Man beachte dabei, daß auch ein stabiles Verfahren beliebig große Rundungsfehler erzeugen kann, nämlich dann, wenn das unvermeidliche Fehlerniveau entsprechend hoch ist. Es wäre auch unangemessen zu verlangen, daß ein Verfahren die in den Daten bereits enthaltene Unsicherheit eliminiert.

Zur Illustration betrachten wir das folgende Beispiel.

2.3.8. Beispiel. In $\Re(10, 4, 9, 9)$ ist bei Rechnung mit unsymmetrischer Rundung (Abbrechen nach der vierten Mantissenstelle) die Summe $z^* = \sum_{i=1}^{6} x_i$ für die in der Tabelle angegebenen Summanden zu berechnen.

<i>i</i>	x_i	x'i	$\varepsilon_i^{(6)} = (x_i' - x_i)/x_i$
1 2 3 4 5	$\begin{array}{c} 0.1065 \times 10^2 \\ 0.1298 \times 10^0 \\ -0.5102 \times 10^{-1} \\ 0.7982 \times 10^{-1} \\ 0.1399 \times 10^0 \\ 0.5570 \\ 100 \\ 0.100 \\ 0.000 \\ 0$	$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$ \begin{array}{c} -4.4 \times 10^{-3} \\ -4.5 \times 10^{-3} \\ -3.5 \times 10^{-3} \\ -2.8 \times 10^{-3} \\ -1.8 \times 10^{-3} \\ -1.8 \times 10^{-3} \end{array} $

Man erhält z = 10.96, die exakte Summe ist $z^* = 11.00829$. Der erzeugte Rundungsfehler ist $\delta z = z^* - z = 0.04829$, und das durch den Darstellungsfehler $\nu = 10^{-3}$ bedingte unvermeidliche Fehlerniveau von z ist $\Delta z_{opt}(x_i, v) = v \sum_{i=1}^{6} |x_i| = 0.011$. Die tatsächliche Fehler-kumulation ist $|\delta z|/\Delta z_{opt} = 4.35$, also 87% des nach Aussage 2.3.7 maximal möglichen Wertes F = n - 1 = 5. Zum Vorgleich sind ist des nacht als des nacht. F = n - 1 = 5. Zum Vergleich sind in der zweiten Spalte gestörte Summanden $x'_i = x_i(1 + \varepsilon_i^{(6)})$ angegeben worden, deren Störungen $\varepsilon_{i}^{(6)}$ den Abschätzungen (9) genügen und deren exakte Summe $\sum_{i=1}^{n} x'_{i}$ gerade den Wert z = 10.96 ergibt. Wir bemerken schließlich, daß Summation der x, nach wachsenden Beträgen auf den Wert 11.00, also das unsymmetrisch gerundete exakte Ergebnis führt.

Für ein günstigeres, aber auch aufwendigeres Summationsverfahren mit $F = \log_2 n$ siehe Ü 2.3.1.

B. Das Skalarprodukt

Aufgabe: Für die Eingangsdaten $x_i, y_i \in \Re$ (i = 1, ..., n) ist das Skalarprodukt

$$s^* = \sum_{i=1}^n x_i y_i$$

zu berechnen.

Bezeichnet \Re^n wie eingangs die Menge aller Vektoren $\boldsymbol{x} = (x_1, ..., x_n)^{\intercal}$ mit $x_i \in \Re$ (i = 1, ..., n), so läßt sich die Aufgabe wie folgt formulieren:

Für $x, y \in \Re^n$ ist $s^* = x^{\mathsf{T}}y$ zu berechnen.

2.3.9. Algorithmus zur Berechnung des Skalarproduktes.

$$s := 0$$

s := 0for i := 1(1)n do $s := s + x_i * y_i$

Autwand: n opms

Zur Rundungsfehleranalyse schreiben wir 2.3.9 in der äquivalenten Form

$$s_0 = 0$$
,
 $s_i = \text{fl}(s_{i-1} + g_i), \quad g_i = \text{fl}(x_i * y_i) \quad (i = 1, ..., n).$

Analog zum Vorgehen bei der Summation folgt unter Beachtung von 2.2.11

$$s_i = (s_{i-1} + g_i) (1 + \varepsilon_i), \qquad |\varepsilon_i| \leq \nu, \tag{18}$$

$$g_i = x_i y_i (1 + \delta_i),$$
 $|\delta_i| \le v$ $(i = 1, ..., n).$ (19)

Die Rekursion (18) ist bis auf die Bezeichnung identisch mit der Rekursion (2) bei der Summation. Aus 2.3.3 folgt daher

$$s = s_n = \sum_{i=1}^n g_i (1 + \varepsilon_i^{(n)}),$$

und $\varepsilon_i^{(n)}$ genügt den Abschätzungen (9). Einsetzen von (19) führt auf

$$s = \sum_{i=1}^{n} x_i y_i (1 + \delta_i) (1 + \varepsilon_i^{(n)}).$$

Nun gilt

$$(1+\delta_i)(1+\varepsilon_i^{(n)}) = 1+\varphi_i^{(n)} \quad \text{mit} \quad \varphi_i^{(n)} = \delta_i + (1+\delta_i)\varepsilon_i^{(n)}, \tag{20}$$

woraus wegen $1 + \nu = 1$ und (9) sofort

$$\begin{aligned} |\varphi_i^{(n)}| &\leq \nu + (1+\nu) \, |\varepsilon_i^{(n)}| \, \stackrel{<}{=} \, \nu + (1+\nu) \min \left(n-i+1, n-1\right) \nu \\ &= \min \left(n-i+2, n\right) \nu \end{aligned}$$

folgt. Schließlich kann

$$(1 + \varphi_i^{(n)}) = (1 + \psi_i^{(n)})^2 \quad \text{mit} \quad |\psi_i^{(n)}| \leq |\varphi_i^{(n)}|/2 \leq 0.5 \min(n - i + 2, n) \nu$$

geschrieben werden, so daß jeder der Faktoren x_i und y_i mit einer Störung belastet werden kann.

2.3.10. Rückwärtsanalyse der Skalarproduktberechnung. Das gemäß Algorithmus 2.3.9 zum Rundungsfehlerniveau v berechnete Skalarprodukt s ist das exakte Skalarprodukt

$$s = \sum_{i=1}^{n} x_i y_i (1 + \varphi_i^{(n)}) = \sum_{i=1}^{n} [x_i (1 + \psi_i^{(n)})] [y_i (1 + \psi_i^{(n)})]$$
(21)

mit gestörten Eingangsdaten, und die Störungen genügen den Abschätzungen

$$\frac{|\varphi_i^{(n)}|}{2|\psi_i^{(n)}|} \bigg\} \stackrel{\text{def}}{=} \min(n-i+2,n) \, \nu \leq n\nu.$$

$$(22)$$

Das Verfahren 2.3.9 ist also in bezug auf die Eingangsdaten $x, y \in \mathbb{R}^n$ numerisch gutartig mit F = n/2.

Wird x als fester Parameter und nur y als Eingangsdatenelement angesehen, so kann $s = \sum_{i=1}^{n} x_i [y_i(1 + \varphi_i^{(n)})]$ geschrieben werden, d. h., nur y braucht mit Störungen belastet zu werden, die dann allerdings unter Umständen doppelt so groß werden können. In dieser Interpretation — d. h. in bezug auf y — liegt also ebenfalls nume-

6 Schwetlick, Numerische Algebra

rische Gutartigkeit vor, aber mit F = n. In den späteren Anwendungen von 2.3.10 werden wir immer nur einen der Faktoren x, y mit Fehlern belasten, weil dann die Abschätzungen einfacher werden.

Zur Vorwärtsanalyse bestimmen wir zunächst das durch den Darstellungsfehler hervorgerufene unvermeidliche Fehlerniveau Δs_{opt} und nehmen dazu an, daß durch x_i, y_i auch alle \tilde{x}_i, \tilde{y}_i mit

$$ilde{x}_i = x_i(1+arsigma_i), \qquad ilde{y}_i = y_i(1+\eta_i), \qquad |arsigma_i|, |\eta_i| \leq
u \qquad (i=1,...,n)$$

repräsentiert werden. Mit $\tilde{s} = \tilde{y}^{\mathsf{T}} \tilde{x}$ ergibt sich dann

$$egin{aligned} &| ilde{s}-s^{m{st}}| = \left|\sum\limits_{i=1}^n \left\{x_i y_i (1+\xi_i) \; (1+\eta_i) - x_i y_i
ight\}
ight| \ &\leq \sum\limits_{i=1}^n \left\{|\xi_i| + |\eta_i| + |\xi_i| \; |\eta_i|
ight\} \; |x_i| \; |y_i| \ &\leq 2
u \sum\limits_{i=1}^n |x_i y_i| = 2
u \; |m{x}|^{\mathsf{T}} \; |m{y}| \leq 2
u \; ||m{x}||_2 \; ||m{y}||_2 \end{aligned}$$

vgl. Beispiel 2.1.10.

Aus (21) folgt weiter

$$|s^* - s| = \left|\sum_{i=1}^n x_i y_i \varphi_i^{(n)}\right| \leq \sum_{i=1}^n |x_i y_i| |\varphi_i^{(n)}|,$$

mit (22) also das folgende Resultat.

2.3.11. Vorwärtsanalyse der Skalarproduktberechnung. Der durch Algorithmus 2.3.9 erzeugte Rundungsfehler $s^* - s$ genügt der Ungleichung

$$|s^* - s| \leq \left\{ \sum_{i=1}^n \min(n - i + 2, n) |x_i y_i| \right\} v \leq \left\{ n \sum_{i=1}^n |x_i y_i| \right\} v \leq n ||\mathbf{x}||_2 ||\mathbf{y}||_2 v.$$
(23)

Das durch die Unsicherheit der Eingangsdaten hervorgerufene unvermeidliche Fehlerniveau ist

Das Verfahren 2.3.9 ist also numerisch stabil mit der Konstanten F = n/2.

2.3.12. Bemerkung. (i) Aus (20) und (8) folgt

$$\varphi_i^{(n)} = \delta_i + \varepsilon_i^{(n)} + O(\nu^2) = \delta_i + \varepsilon_i + \varepsilon_{i+1} + \dots + \varepsilon_n + O(\nu^2), \qquad (25)$$

wobei δ_i und ε_i durch (18), (19) erklärt sind. Bei symmetrischer Rundung ist daher zu erwarten, daß bei der Summation der n - i + 2 relativen Rundungsfehler δ_i , $\varepsilon_i, \ldots, \varepsilon_n$ stochastische Auslöschung eintritt.

(ii) Aus (23) folgt, daß bei einheitlichem Vorzeichen der Teilprodukte $x_i y_i$ der relative Fehler von *s* klein ist: $|s^* - s|/|s^*| \leq n\nu$. Bei unterschiedlichem Vorzeichen

der Teilprodukte kann der relative Fehler von *s* jedoch beliebig groß werden, vgl. Beispiel 2.1.10.

(iii) Wird x als exakter Parameter betrachtet, so ergibt sich bezüglich der Darstellungsfehler von y die Schranke $\Delta s_{opt}(y_i, v) = v \sum_{i=1}^{n} |x_i| |y_i|$, also der halbe Wert von (24). Bei festem x liegt dann bezüglich y Stabilität mit F = n vor in Übereinstimmung zur Nachbemerkung von 2.3.10. \Box

Eine Möglichkeit zur Abschätzung des erzeugten Rundungsfehlers $|s^* - s|$ gibt die nachfolgende, analog zu (11) abgeleitete a-posteriori-Schranke: Aus (18), (19) folgt

$$s_i = (s_{i-1} + g_i) (1 + \varepsilon_i) = s_{i-1} + g_i + (s_{i-1} + g_i) \varepsilon_i$$

= $s_{i-1} + g_i + s_i \varepsilon_i / (1 + \varepsilon_i)$

sowie

$$g_i = x_i y_i (1+\delta_i) = x_i y_i + x_i y_i \delta_i = x_i y_i + g_i \delta_i / (1+\delta_i)$$
 ,

also insgesamt

$$s_i = s_{i-1} + x_i y_i + g_i \delta_i / (1 + \delta_i) + s_i \varepsilon_i / (1 + \varepsilon_i)$$

Summation führt direkt auf

$$|s^* - s| = \left| \sum_{i=1}^n \left\{ g_i \delta_i / (1 + \delta_i) + s_i \varepsilon_i / (1 + \varepsilon_i) \right\} \right| \stackrel{<}{=} v \sum_{i=1}^n \left\{ |g_i| + |s_i| \right\} =: d^*.$$
(26)

Der folgende Algorithmus ist eine Erweiterung von 2.3.9, in der gleichzeitig mit s die Schranke d durch Auswertung von (26) mit berechnet wird.

2.3.13. Algorithmus zur Berechnung des Skalarproduktes und einer a-posteriori-Fehlerschranke.

$$s := d := 0$$

for $i := 1(1)n$ do
$$\begin{vmatrix} g := x_i * y_i \\ s := s + g \\ d := d + |g| + |s| \end{vmatrix}$$
$$d := v * d$$

Aufwand: n opms + 2n ops, also 2n ops mehr als für Algorithmus 2.3.9.

Auch für schlechtkonditionierte Aufgaben, bei denen der Quotient $\sum |x_iy_i|/|\sum x_iy_i|$ bzw. $||\boldsymbol{x}||_2 ||\boldsymbol{y}||_2/|\boldsymbol{x}^{\mathsf{T}}\boldsymbol{y}|$ sehr groß wird, vgl. 2.3.11 und 2.1.10, kann eine hohe relative Genauigkeit erreicht werden, indem die Teilprodukte x_iy_i in höherer Genauigkeit berechnet und aufsummiert werden. Zusätzlich zur Standard-Gleitpunktarithmetik.fl mit dem Fehlerniveau ν muß dazu eine Arithmetik fl₁ mit dem kleineren Fehlerniveau $v_1 \ll \nu$ verfügbar sein. Zur Illustration betrachten wir das folgende Beispiel.

6*

2.3.14. Beispiel. Berechnung von $s^* = \sum_{i=1}^{3} x_i y_i$ für die unten angegebenen Eingangsdaten x_i , $y_i \in \Re = \Re(10, 3, 9, 9)$. Die Standardarithmetik fl ist die in \Re mit unsymmetrischer Rundung. Produktbildung $g_i = \mathrm{fl}_1 (x_i * y_i)$ und Summation $s_i = \mathrm{fl}_1 (s_{i-1} + g_i)$ erfolgen in der Arithmetik fl₁ von $\Re_1 = \Re(10, 6, 9, 9)$ mit unsymmetrischer Rundung, also mit doppelter Mantissenlänge. Die Fehlerniveaus sind $v = 10^{-2}$ bzw. $v_1 = 10^{-5}$. Zum Vergleich wurden zusätzlich die in der Standardarithmetik berechneten Werte mit angegeben.

i	x_i	y_i	$\mathrm{fl}\left(x_{i} \ast y_{i}\right)$	fl ($s_{i-1} + g_i$)	$\mathrm{fl}_{1}\left(x_{i} \ast y_{i}\right)$	$\mathrm{fl}_{1}\left(s_{i-1}+g_{i}\right)$
1 2	$\begin{array}{r} 0.341 + 00 \\ 0.567 - 01 \end{array}$	$0.107\!+\!02$ $0.164\!+\!00$	$ \begin{smallmatrix} 0.364 + 01 \\ 0.929 - 02 \end{smallmatrix} $	$0.364 + 01 \\ 0.364 + 01$	$\substack{0.364870+01\\0.929880-02}$	$\substack{0.364870+01\\0.365799+01}$
3	-0.518 + 01	$0.706\!+\!00$	-0.365+01	-0.100 - 01	-0.365708+01	0.910000-03

Die Schreibweise -0.518+01 bedeutet hier und im folgenden stets $-0.518 \times 10^{+01}$ usw. Der exakte Wert ist $s^* = 0.9188 \times 10^{-3}$.

Bei Rechnung in fl ergibt sich $s = -0.100 \times 10^{-1}$, d. h., selbst das Vorzeichen ist falsch. Der relative Rundungsfehler $|s^* - s|/|s^*| = 11.9$ überschreitet das Niveau $\nu = 10^{-2}$ des Darstellungsfehlers der Daten um den Faktor 1190! Der durch Akkumulation der Summen in fl₁ berechnete und danach auf einfache Genauigkeit gerundete Wert $s^{ac} = 0.910 \times 10^{-3}$ hat dagegen einen Rundungsfehler $|s^* - s^{ac}|/|s^*| = 0.96 \times 10^{-2}$ in der Größenordnung von v. Man beachte jedoch, daß diese hohe Genauigkeit von sac in bezug auf s* nur dann einen Informationsgewinn bringt, wenn die Eingangsdaten x_i , y_i exakt bzw. sehr genau sind. Sind die Eingangsdaten dagegen nur im Rahmen des Darstellungsfehlers genau, so müssen alle \tilde{x}_i, \tilde{y}_i mit $x_i = \operatorname{rd}(\tilde{x}_i)$ und $y_i = \operatorname{rd}(\tilde{y}_i)$ als gleichwertig angesehen werden. Über der Menge dieser zulässigen Eingangsdaten kann $\tilde{s} = \sum_{1=i}^{3} \tilde{x}_i \tilde{y}_i$ alle Werte zwischen $0.341 \times 10.7 + 0.0567 \times 0.164$ - $5.19 \times 0.707 = -0.011331$ und $0.342 \times 10.8 + 0.0568 \times 0.165 - 5.18 \times 0.706 = 0.045892$ annehmen. Im Sinne des Darstellungsfehlers ist daher jedes \tilde{z} zwischen -0.011... und +0.045... als zulässiges Resultat zu akzeptieren, vgl. Abschnitt 2.1. Speziell ist $s^* = 0.009188$ gegenüber s = -0.01 weder besser noch schlechter. Zum Beispiel ergibt sich s = -0.01 als exakter Wert, wenn x_3 durch $x'_3 = -5.189$ und y_3 durch $y'_3 = 3.6679988/5.189 = 0.70687...$ ersetzt werden; die Computerdarstellungen der modifizierten Werte stimmen mit den ursprünglichen überein.

Beispiel 2.3.14 zeigt, daß die Unsicherheit der Eingangsdaten durch keinen noch so genauen Algorithmus eliminiert werden kann. Skalarproduktberechnung durch Akkumulation in höherer Genauigkeit ist daher nur dann sinnvoll, wenn die Eingangsdaten x, y als exakt angesehen werden können oder wenn n und damit die Fehlerkumulation F = n/2 von 2.3.9 sehr groß ist. Im ersten Fall bringt Akkumulation mit höherer Genauigkeit tatsächlich eine wesentlich größere Genauigkeit, allerdings treten exakte Eingangsdaten sehr selten auf. Im zweiten Fall — n groß — kann die Fehlerkumulation durch Akkumulation in höherer Genauigkeit praktisch auf den Wert F = 1gebracht werden. In diesem Fall läßt sich aber auch in der Standardarithmetik fl durch Verwendung der Kaskadensummation aus Ü 2.3.1 zur Summation der Produkte fl $(x_i * y_i)$ eine Reduktion auf $F = (\log_2 n + 1)/2$ erreichen, also ein wesentlich kleinerer Wert als F = n/2 für die Standardsummation. Für $n = 2^{10} = 1024$ ist z. B. $(\log_2 n + 1)/2 = 5.5$, aber n/2 = 512!

C. Darstellung von Vektoren und Matrizen in der Computerarithmetik

Bei den Aufgabenklassen der linearen Algebra treten in der Regel Vektoren $\boldsymbol{b} = (b_i) \in \mathbf{R}^n$ und Matrizen $\boldsymbol{A} = (a_{ij}) \in \mathbf{R}^{m,n}$ als Eingangsdaten auf. Die zugehörigen Computerdarstellungen sollen mit rd (\boldsymbol{b}) bzw. rd (\boldsymbol{A}) bezeichnet werden; sie entstehen durch elementweise Rundung, d. h.

$$\mathrm{rd}\;(\boldsymbol{b}) := \left(\mathrm{rd}\;(b_i)\right) \in \mathfrak{R}^{\boldsymbol{m}}, \qquad \mathrm{rd}\;(\boldsymbol{A}) := \left(\mathrm{rd}\;(a_{ij})\right) \in \mathfrak{R}^{\boldsymbol{m},\boldsymbol{n}}. \tag{27}$$

Analog zu \Re^m bezeichnet $\Re^{m,n}$ dabei die Menge aller (m,n)-Matrizen, deren Elemente Computerzahlen aus \Re sind. Im Fall

$$b_i = 0 \text{ oder MIN} \le |b_i| \le \text{MAX}$$
 bzw. $a_{ij} = 0 \text{ oder MIN} \le |a_{ij}| \le \text{MAX}$ (28)

folgt aus 2.2.7

$$\begin{aligned} |\mathrm{rd}\ (b_i) - b_i| &\leq \nu \, |b_i| \quad \text{bzw.} \quad |\mathrm{rd}\ (a_{ij}) - a_{ij}| &\leq \nu \, |a_{ij}| \\ (i &= 1, \dots, m; \, j = 1, \dots, n). \end{aligned} \tag{29}$$

Unter Verwendung der in 1.1.J eingeführten Ordnungsbegriffe läßt sich (29) kurz als

$$|\mathrm{rd}\,(oldsymbol{b})-oldsymbol{b}|\leq
u\,|oldsymbol{b}|\quad\mathrm{bzw.}\quad|\mathrm{rd}\,(A)-A|\leq
u\,|A|$$

schreiben. Für die Vektornormen mit dem Index $p \in \{1, 2, \infty\}$ bzw. für die Matrixnormen mit dem Index $p \in \{1, 2, \infty, F\}$ folgt hieraus

$$\|\operatorname{rd}(b) - b\| \leq v \|b\|$$
 bzw. $\|\operatorname{rd}(A) - A\| \leq v \||A|\|$,

vgl. wieder Abschnitt 1.1.J. Auf der rechten Seite kann |||A||| für die Matrixnormen mit $p \in \{1, \infty, F\}$ durch ||A|| ersetzt werden, denn diese Normen sind wie die betrachteten Vektornormen absolut und monoton. Für p = 2 gilt wegen (1.1.44)

$$\||A|\|_{2} \leq \||A|\|_{F} = \|A\|_{F} \leq \sqrt[n]{\min(m, n)} \|A\|_{2}.$$
(30)

Zusammenfassend erhalten wir

2.3.15. Aussage. Für $b \in \mathbb{R}^m$ bzw. $A \in \mathbb{R}^{m,n}$ sei (28) erfüllt, d. h., bei Übergang zu den Computerdarstellungen rd (b) bzw. rd (A) trete weder Über- noch Unterlauf auf. Dann gilt

(i) Die Computerdarstellungen haben eine hohe individuelle Genauigkeit im Sinne von

$$\operatorname{rd}(\boldsymbol{b}) - \boldsymbol{b}| \leq v |\boldsymbol{b}| \quad \operatorname{bzw.} \quad |\operatorname{rd}(\boldsymbol{A}) - \boldsymbol{A}| \leq v |\boldsymbol{A}|. \tag{31}$$

(ii) Die hohe individuelle Genauigkeit überträgt sich auf die Normen im Sinne von

$$\|\mathrm{rd}\,(\boldsymbol{b}) - \boldsymbol{b}\| \leq v \,\|\boldsymbol{b}\| \tag{32}$$

bzw.

$$\|\mathrm{rd}\,(A) - A\| \leq \begin{cases} v \|A\| & \text{für } p = 1, \,\infty, F, \\ v \||A|\| \leq v \,\sqrt{\min(m, n)} \,\|A\| & \text{für } p = 2. \end{cases}$$
(33)

2.3.16. Bemerkung. (i) Eine im Sinne von

$$\|\tilde{\boldsymbol{b}} - \boldsymbol{b}\| \leq \boldsymbol{v} \|\boldsymbol{b}\| \quad \text{bzw.} \quad \|\tilde{\boldsymbol{A}} - \boldsymbol{A}\| \leq \boldsymbol{v} \||\boldsymbol{A}|\| \tag{34}$$

der Norm nach genaue Darstellung von b bzw. A durch \tilde{b} bzw. \tilde{A} ist gleichbedeutend damit, daß die betragsgroßen Elemente von b bzw. A durch die entsprechenden Elemente von \tilde{b} bzw. \tilde{A} mit hoher individueller Genauigkeit dargestellt werden.

(ii) Für viele Aufgaben der numerischen linearen Algebra kann die im Sinne von (31) hohe individuelle Genauigkeit aller Eingangsdaten nicht ausgenutzt werden. Um akzeptable Resultate zu garantieren, genügt eine der Norm nach hohe Genauigkeit der Eingangsdaten.

Zur Illustration betrachten wir die folgenden Beispiele, siehe auch Ü 2.3.2.

2.3.17. Beispiel. Es sei $\Re = \Re(10, 3, 9, 9)$, die Rundung sei unsymmetrisch.

(i) Der Vektor $\mathbf{b} = (6.6666, 0.0033333)^{\mathsf{T}} \in \mathbb{R}^2$ wird durch $\operatorname{rd}(\mathbf{b}) = (0.666 \times 10^4, 0.333 \times 10^{-2})^{\mathsf{T}} \in \mathbb{R}^2$ mit dem kollektiven relativen Rundungsfehler $||\operatorname{rd}(\mathbf{b}) - \mathbf{b}||_2/||\mathbf{b}||_2 = 0.990 \times 10^{-3}$ dargestellt. Wird $\operatorname{rd}(\mathbf{b})$ durch $\tilde{\mathbf{b}} = (0.666 \times 10^4, 0.400 \times 10^{-2})^{\mathsf{T}} \in \mathbb{R}^2$ ersetzt, so ergibt sich praktisch derselbe kollektive Darstellungsfehler $||\tilde{\mathbf{b}} - \mathbf{b}||/||\mathbf{b}|| = 0.995 \times 10^{-3}$, obwohl $\tilde{\mathbf{b}}_2$ in keiner Ziffer mit b_2 übereinstimmt und einen individuellen Darstellungsfehler $||\tilde{\mathbf{b}}_2 - \mathbf{b}_2|/|b_2| = 2 \times 10^{-1}$ von 20ν aufweist.

(ii) Für $\boldsymbol{x} = (0.578, 0.713 \times 10^{-3})^{\mathsf{T}}$, $\boldsymbol{y} = (0.202, 0.513 \times 10^{1})^{\mathsf{T}} \in \Re^2$ soll $s = \mathrm{fl}(\boldsymbol{x}^{\mathsf{T}}\boldsymbol{y})$ berechnet werden. Es ergibt sich $s = \mathrm{fl}(\mathrm{fl}(x_1 \ast y_1) + \mathrm{fl}(x_2 \ast y_2)) = \mathrm{fl}(0.116 + 0.365 \times 10^{-2}) = 0.118$. Sind $\boldsymbol{x}, \boldsymbol{y}$ Computerdarstellungen der Vektoren $\tilde{\boldsymbol{x}}, \tilde{\boldsymbol{y}} \in \mathbb{R}^2$, so haben die Komponenten eine individuelle relative Unsicherheit von $\boldsymbol{v} = 10^{-2}$. Diese bewirkt – unabhängig von den bei der Berechnung auftretenden Rundungsfehlern – im ersten Summanden von s eine absolute Unsicherheit von 10^{-3} , beim zweiten eine von 10^{-5} . Die Unsicherheit in s würde nicht wesentlich größer werden, wenn x_2 einen um den Faktor 100 größeren Fehler hätte. Nur bei spezieller Korrelation der Daten – wenn etwa y_2 in der Größenordnung 10^3 liegen würde – könnte die individuelle hohe Genauigkeit von x_2 ausgenutzt werden.

Aus 2.3.16(i) folgt u. a., daß beim Übergang zu rd (b) bzw. rd (A) auftretender Unterlauf die Gültigkeit von (32) bzw. (33) nicht beeinträchtigt, sofern in b bzw. Amindestens ein Element in vernünftiger Größenordnung vorhanden ist.

D. Einfache Vektor- und Matrixoperationen

Wir betrachten hier nur einige einfache Beispiele, um die Technik und die Sprache der Rundungsfehleranalyse zu demonstrieren.

Aufgabe 1: Für $a, b, x \in \Re^n$ ist

 $\boldsymbol{z^*} = (\boldsymbol{I} - \boldsymbol{a} \boldsymbol{b}^{\mathsf{T}}) \, \boldsymbol{x}$

zu berechnen.

Algorithmus: $\alpha := b^{\mathsf{T}}x, z := x - a * \alpha$ Aufwand: 2n opms Rundungsfehleranalyse: Nach 2.3.10 gilt

$$\alpha = \sum_{j=1}^{n} b_j x_j (1 + \varphi_j^{(n)}), \qquad |\varphi_j^{(n)}| \leq n\nu,$$

und für die Komponenten z_i von z ergibt sich nach 2.2.11

$$z_{i} = fl(x_{i} - a_{i} * \alpha) = [x_{i} - a_{i} * \alpha(1 + \delta_{i})](1 + \xi_{i}), \quad |\delta_{i}|, |\xi_{i}| \leq \nu, \quad (35)$$

also

$$z_i = x_i(1+\xi_i) - a_i(1+\delta_i) (1+\xi_i) \left\{ \sum_{j=1}^n b_j \left[\frac{1+\varphi_j^{(n)}}{1+\xi_j}
ight] x_j(1+\xi_j)
ight\}$$

Nach 2.3.2 ist $(1 + \delta_i) (1 + \xi_i) = 1 + \alpha_i$ mit $|\alpha_i| \leq 2\nu$ und ebenso $(1 + \varphi_j^{(n)})/(1 + \xi_j) = 1 + \beta_j$ mit $|\beta_j| \leq (n + 1)\nu$, vgl. Ü 2.2.3. Mit den durch $\delta x_i = \xi_i x_i, \, \delta a_i = \alpha_i a_i$ und $\delta b_i = \beta_i b_i$ festgelegten Störungen $\delta x, \, \delta a$ und δb folgt dann

$$z = (x + dx) - (a + da) (b + db)^{\mathsf{T}} (x + dx)$$
$$= [I - (a + da) (b + db)^{\mathsf{T}}] (x + dx),$$

und die Störungen sind im Sinne von

$$|\delta \boldsymbol{a}| \leqq 2 \boldsymbol{v} |\boldsymbol{a}|, \quad |\delta \boldsymbol{b}| \leqq (n+1) \boldsymbol{v} |\boldsymbol{b}| \text{ und } |\delta \boldsymbol{x}| \leqq \boldsymbol{v} |\boldsymbol{x}|$$

individuell klein. Der Algorithmus ist also numerisch gutartig mit F = n + 1, und zwar elementweise in bezug auf die Eingangsdaten a, b und x. Daß die Verteilung der Störungen in gewissen Grenzen willkürlich ist, zeigt die folgende gemischte Fehleranalyse, bei der x nicht, dafür aber z mit Störungen belastet wird. Aus (35) folgt auch

$$z_i/(1+\xi_i) = z_i(1+\eta_i) = x_i - a_i(1+\delta_i) \left\{ \sum_{j=1}^n b_j(1+\varphi_j^{(n)}) x_j \right\},$$

wobei η_i mit $|\eta_i| \leq v$ über $1/(1 + \xi_i) = (1 + \eta_i)$ eingeführt worden ist. Mit den durch $\delta z_i = \eta_i z_i$, $\delta a_i = \delta_i a_i$ und $\delta b_i = \varphi_i^{(n)} b_i$ festgelegten Störungen δz , δa und δb , die sämtlich wegen

$$|\delta \boldsymbol{z}| \leq v |\boldsymbol{z}|, \quad |\delta \boldsymbol{a}| \leq v |\boldsymbol{a}|, \quad |\delta \boldsymbol{b}| \leq nv |\boldsymbol{b}|$$

individuell klein sind, ergibt sich

$$z + dz = x - (a + da) (b + db)^{\mathsf{T}} x = [I - (a + da) (b + db)^{\mathsf{T}}] x.$$

Das berechnete Resultat z ist das um dz gestörte exakte Resultat zu den um da und db gestörten Eingangsdaten, während x nicht mit Störungen belastet wird.

Aufgabe 2. Für $A \in \Re^{m,n}$ und $x \in \Re^n$ ist

$$y^* = Ax$$

zu berechnen.

Algorithmus: for
$$i := 1(1)m$$
 do $y_i := \sum_{j=1}^n a_{ij} * x_j$

Aufwand: mn opms

Rundungsfehleranalyse: Die *i*-te Komponente y_i ergibt sich als Skalarprodukt der *i*-ten Zeile von A mit x, so daß wegen 2.3.10

$$y_i = \sum_{j=1}^n a_{ij}(1+\varphi_{ij}) x_j \quad \text{mit} \quad |\varphi_{ij}| \leq nr$$
(36)

gilt. Mit der Störungsmatrix δA , $(\delta A)_{ij} = \varphi_{ij}a_{ij}$ folgt daher:

Das berechnete Produkt y läßt sich in der Form

$$\boldsymbol{y} = (\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}) \, \boldsymbol{x} \quad \text{mit} \quad |\boldsymbol{\delta}\boldsymbol{A}| \, \leqq \, n\boldsymbol{v} \, |\boldsymbol{A}| \tag{37}$$

schreiben, d. h., die Matrix-Vektor-Multiplikation ist numerisch gutartig mit F = n elementweise in bezug auf A.

In (36) lassen sich auch die Komponenten von \boldsymbol{x} mit den Störungen φ_{ij} belasten: Mit der Störung $d\boldsymbol{x}^{(i)}$, $(d\boldsymbol{x}^{(i)})_i = \varphi_{ii}x_i$, gilt

$$y_i = [\mathbf{A}(\mathbf{x} + \mathbf{\delta}\mathbf{x}^{(i)})]_i \quad \text{mit} \quad |\mathbf{\delta}\mathbf{x}^{(i)}| \le n\nu |\mathbf{x}|.$$
(38)

Da die Störungen $dx^{(i)}$ von x für jede Komponente i. allg. verschieden sind, ist nur die Berechnung einer einzelnen Komponente y_i , nicht aber die Berechnung des gesamten Vektors $y = (y_1, \ldots, y_m)^{\mathsf{T}}$ numerisch gutartig in bezug auf x, d. h., es gibt i. allg. keine kleine Störung dx mit y = A(x + dx). Allerdings folgt aus (38) die numerische Stabilität in bezug auf x, d. h. im Vergleich auf das durch den Darstellungsfehler von x hervorgerufene unvermeidliche Fehlerniveau von $y^* = Ax$: Wenn neben x alle $\tilde{x} = x + dx$ mit $|dx| \leq v |x|$ als zulässig angesehen werden, ergibt sich

$$ilde{oldsymbol{y}} - oldsymbol{y}^*| = |A(ilde{oldsymbol{x}} - oldsymbol{x})| \le |A| \ |\delta oldsymbol{x}| \le
u \ |A| \ |oldsymbol{x}| = oldsymbol{y}_{ ext{opt}}(x_i,
u)$$
(39)

als Vektor der unvermeidlichen Fehlerniveaus in jeder Komponente von $\boldsymbol{y^*}$. Andererseits folgt aus (38) sofort

$$|y_i^{st}-y_i|=|(oldsymbol{A}\deltaoldsymbol{x^{(i)}})_i|\leq (|oldsymbol{A}|~|oldsymbol{\delta x^{(i)}}|)_i\leq n
u(|oldsymbol{A}|~|oldsymbol{x}|)_i,$$

also

$$|\boldsymbol{y^*} - \boldsymbol{y}| \le n\nu |\boldsymbol{A}| |\boldsymbol{x}| \tag{40}$$

als Schranke für den erzeugten Rundungsfehler. Es liegt daher (komponentenweise) numerische Stabilität mit F = n in bezug auf x vor. \Box

Aufgabe 3: Für $A \in \Re^{m,n}$, $B \in \Re^{n,l}$ ist

$$C^* = AB$$

zu berechnen.

Algorithmus: for i := 1(1)m do

for
$$j := 1(1)l$$
 do $c_{ij} := \sum_{k=1}^{n} a_{ik} * b_{kj}$

Aufwand: mnl opms

Rundungsfehleranalyse: Da die j-te Spalte $c^{(j)}$ von C als Produkt von A mit der j-ten Spalte $b^{(j)}$ von B berechnet wird, also $c^{(j)} = Ab^{(j)}$ gilt, können die Ergebnisse

von Aufgabe 2 direkt angewendet werden und liefern: Zu jedem j gibt es eine Störung $dA^{(j)} \in \mathbb{R}^{m,n}$, so daß die j-te Spalte der berechneten Ergebnismatrix als

$$\boldsymbol{c}^{(j)} = (\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}^{(j)}) \boldsymbol{b}^{(j)} \quad \text{mit} \quad |\boldsymbol{\delta}\boldsymbol{A}^{(j)}| \leq \nu n |\boldsymbol{A}|$$
(41)

geschrieben werden kann. Im allgemeinen existieren jedoch keine relativ kleinen Störungen δA und δB mit $C = (A + \delta A) (B + \delta B)$, so daß die Matrix-Matrix-Multiplikation nicht numerisch gutartig ist.

Zur Stabilitätsuntersuchung betrachten wir gestörte Faktoren $\tilde{A} = A + \delta A$, $\tilde{B} = B + \delta B$ mit individuell kleinen Störungen $|\delta A| \leq v |A|$, $|\delta B| \leq v |B|$. Dann folgt

$$|C - C^*| = |(A + \delta A) (B + \delta B) - AB| \leq |\delta A| |B| + |A| |\delta B| + |\delta A| |\delta B|$$

$$\leq (2\nu + \nu^2) |A| |B|, \qquad (42)$$

d. h., das elementweise Niveau des unvermeidlichen Fehlers ist $2\nu |\mathbf{A}| |\mathbf{B}|$. Aus (41) folgt

$$|c^{*(j)} - c^{(j)}| = |\delta A^{(j)} b^{(j)}| \leq vn |A| |b^{(j)}|, \text{ also } |C^* - C| \leq vn |A| |B|. (43)$$

Die Matrix-Matrix-Multiplikation ist daher numerisch stabil mit F = n/2 in bezug auf A und B. Wird nur die Darstellungsunsicherheit von B betrachtet und A als fixierter Parameter angesehen, so tritt in (42) der Faktor v statt $2v + v^2$ auf. In diesem Modell ist F = n. \Box

E. Vermeidung von Über- und Unterlauf

Wir haben bisher vorausgesetzt, daß bei der Realisierung der Algorithmen in \Re weder Über- noch Unterlauf eintritt. Bei den heute meist üblichen Werten der Unterlaufschranke MIN und der Überlaufsschranke MAX — vgl. Abschnitt 2.2 — deutet ein tatsächlich vorkommender Über- bzw. Unterlauf fast immer auf einen Programm- bzw. Anwendungsfehler hin. Die Hauptursache ist meist die fortlaufende Multiplikation großer bzw. kleiner Zahlen etwa bei der Determinantenberechnung oder bei der Berechnung der Exponentialfunktion sowie Division durch die Computernull, die meist gesondert angezeigt wird. Insbesondere sollte Über- bzw. Unterlauf bei der Berechnung von Zwischenergebnissen vermieden werden, wenn die Eingangs- und Ausgangsdaten betragsmäßig in [MIN, MAX] liegen. Durch geschickte Programmierung läßt sich das häufig erreichen; als Beispiel gehen wir auf die Berechnung der Euklidischen Norm eines Vektors ein.

Aufgabe: Für
$$x \in \mathfrak{R}^n$$
 ist $\|x\|_2 = \sqrt{\sum\limits_{i=1}^n (x_i)^2}$ zu berechnen.

Die wörtliche Umsetzung dieser Formel führt in $\Re(10, 4, 9, 9)$ z. B. für $x = (0.3 \times 10^6, 0.1 \times 10^{-2})^{\intercal}$ zu Überlauf beim Quadrieren von x_1 , obwohl $||x||_2 = 0.300 \dots \times 10^6$ in \Re gut darstellbar ist. Durch eine geeignete Skalierung läßt

sich der beschriebene Effekt vermeiden, etwa durch Auswertung von

$$\|m{x}\|_2 = |x_\mu| \; \left| \int_{i=1}^n (x_i/|x_\mu|)^2 \;, \qquad |x_\mu| = \max_i \, |x_i|$$

nach dem Algorithmus

 $\begin{array}{l} \operatorname{xmax} := 0\\ \operatorname{for} i := 1(1)n \ \operatorname{do} \ \operatorname{if} |x_i| > \operatorname{xmax} \ \operatorname{then} \ \operatorname{xmax} := |x_i|\\ \operatorname{if} \operatorname{xmax} = 0\\ \operatorname{then} \ \operatorname{normx} := 0\\ \operatorname{else} \ \left| \begin{array}{l} s := 0\\ \operatorname{for} i := 1(1)n \ \operatorname{do} \ s := s + (x_i/\operatorname{xmax}) \uparrow 2\\ \operatorname{normx} := \operatorname{xmax} * \operatorname{sqrt} (s) \end{array} \right| \end{array}$

Die Quotienten $x_i/|x_{\mu}|$ sind betragsmäßig durch 1 beschränkt, so daß kein Überlauf eintreten kann, und mindestens einer ist gleich 1, so daß evtl. eintretender Unterlauf die Genauigkeit des Resultats nicht beeinflußt. Unter der Annahme (2.2.10) und mit 2.3.10 läßt sich zeigen, daß der angegebene Algorithmus numerisch gutartig ist mit F = (n + 6)/2. Da alle Summanden nichtnegativ sind, ist der erzeugte Rundungsfehler klein: $|||x||_2 - \operatorname{normx} |/||x||_2 \leq 0.5(n + 6) \nu$.

Durch raffiniertere Vorschriften kann der oben notwendige zweimalige Durchlauf der x_i — einmal zur Maximumsuche, einmal zur Aufsummierung der Quadrate — auf einen reduziert werden, siehe B 2.6 für Literaturhinweise.

Bei klassischer Auswertung gemäß norm $\mathbf{x} := \operatorname{sqrt} (\mathbf{x}^{\mathsf{T}}\mathbf{x})$ liegt numerische Gutartigkeit mit F = (n + 2)/2 vor.

F. Zusammenfassung

Nachdem die Grundbegriffe der Fehleranalyse an Hand mehrerer Beispiele eingeführt worden sind, sollen sie jetzt in allgemeiner Form fixiert werden.

Wir betrachten die numerische Problemklasse $\{\mathscr{E}, P\}$, d. h. eine Abbildung P: $\mathscr{E} \subset \mathbb{R}^N \to \mathbb{R}^M$, und einen (konzeptionellen) numerischen Algorithmus $V_0: \mathscr{E} \subset \mathbb{R}^N$ $\to \mathbb{R}^M$ zur Lösung von Aufgaben $\{e, P\}$ aus der Klasse $\{\mathscr{E}, P\}$. Die Computerealisierung von V_0 in einer Gleitpunktarithmetik mit genügend kleinem v sei V: $\mathfrak{E} \subset \mathfrak{R}^N \to \mathfrak{R}^M$, wobei \mathfrak{E} eine geeignete Einschränkung von \mathscr{E} auf \mathfrak{R}^N sei, d. h., es gelte $\mathfrak{E} \subset \mathscr{E}$. Durch V wird jedem Eingangsdatensatz $e \in \mathfrak{E}$ die numerische Lösung $a = V(e) \in \mathfrak{R}^M$ zugeordnet; die zu e gehörende exakte Lösung ist $a^* = P(e)$.

2.3.18. Begriffe der Fehleranalyse.

(i) Der numerische Algorithmus $V: \mathfrak{G} \subset \mathfrak{R}^N \to \mathfrak{R}^M$ heißt in der Klasse $\{\mathfrak{G}, P\}$ numerisch gutartig, wenn zu jedem $e \in \mathfrak{G}$ eine Störung $de \in \mathbb{R}^N$ angegeben werden kann, so daß

$$\boldsymbol{a} = \boldsymbol{P}(\boldsymbol{e} + \boldsymbol{\delta} \boldsymbol{e}) \quad \text{und} \quad \|\boldsymbol{\delta} \boldsymbol{e}\| \leq F_{w^{\mathcal{V}}} \|\boldsymbol{e}\|$$

$$\tag{44}$$

mit einer von r und e unabhängigen Fehlerkumulationskonstanten F_w gilt.

(ii) Für $e \in \mathfrak{E}$ ist

$$\Delta \boldsymbol{a}_{\text{opt}}(\boldsymbol{e},\boldsymbol{v}) = \max \left\{ \|\boldsymbol{P}(\boldsymbol{e} + \boldsymbol{\delta}\boldsymbol{e}) - \boldsymbol{P}(\boldsymbol{e})\| : \boldsymbol{\delta}\boldsymbol{e} \in \boldsymbol{\mathsf{R}}^{N} \text{ mit } \boldsymbol{e} + \boldsymbol{\delta}\boldsymbol{e} \in \boldsymbol{\mathscr{E}} \\ \text{und } \|\boldsymbol{\delta}\boldsymbol{e}\| \leq \boldsymbol{v} \|\boldsymbol{e}\| \right\}$$
(45)

das durch die relative Darstellungsunsicherheit v der Eingangsdaten bedingte unvermeidliche Fehlerniveau der Ausgangsdaten.

(iii) Der numerische Algorithmus $V: \mathfrak{G} \subset \mathfrak{R}^N \to \mathfrak{R}^M$ heißt in der Klasse $\{\mathfrak{G}, P\}$ *numerisch stabil*, wenn der erzeugte Rundungsfehler $a^* - a = P(e) - V(e)$ des berechneten Resultates a = V(e) für jedes $e \in \mathfrak{G}$ der Ungleichung

$$\|\boldsymbol{a^*} - \boldsymbol{a}\| = \|\boldsymbol{P}(\boldsymbol{e}) - \boldsymbol{V}(\boldsymbol{e})\| \leq F_s \Delta \boldsymbol{a}_{\text{opt}}(\boldsymbol{e}, \boldsymbol{\nu}) \tag{46}$$

mit einer von ν und e unabhängigen *Fehlerkumulationskonstanten* F_s genügt. Dabei wird vorausgesetzt, daß F_w bzw. F_s minimal gewählt werden.

2.3.19. Bemerkung. (i) Numerische Gutartigkeit bedeutet: Das berechnete Resultat ist das exakte Resultat eines benachbarten Problems mit gestörten Eingangsdaten, und das relative Störungsniveau ist höchstens das F_w -fache der relativen Darstellungsunsicherheit der Eingangsdaten. Ein numerisch gutartiges Verfahren mit nicht zu großer Fehlerkumulationskonstanten nutzt die in den Eingangsdaten enthaltene Information bestmöglich aus. Numerische Gutartigkeit ist die beste Eigenschaft, die ein numerisches Verfahren in bezug auf den Rundungsfehlereinfluß haben kann.

(ii) Numerische Stabilität bedeutet: Der erzeugte Rundungsfehler ist höchstens das F_s -fache des durch die Darstellungsunsicherheit der Eingangsdaten bedingten unvermeidlichen Fehlers der Ausgangsdaten. Ein numerisch nicht stabiler Algorithmus verursacht für gewisse — i. allg. nicht für alle! — Probleme $\{e, P\}$ der Klasse $\{\mathfrak{G}, P\}$ beliebig hohe Genauigkeitsverluste. Numerische Stabilität ist daher eine Mindestforderung an ein vernünftiges numerisches Verfahren.

(iii) Die Abschätzungen (44) bzw. (46) sind zur Fehlerabschätzung für ein konkretes einzelnes Problem i. allg. nicht geeignet, weil

— die Fehlerkumulationskonstanten minimal in bezug auf die gesamte Klasse $\{\mathfrak{G}, P\}$, nicht aber bezüglich des vorliegenden Problems $\{e, P\}$ festgelegt sind. Für letzteres wird i. allg. eine geringere maximale Fehlerkumulation vorliegen als durch F_w bzw. F_s angegeben.

selbst die kleinstmöglichen Fehlerkumulationskonstanten stets das ungünstigste
 Fehlerverhalten berücksichtigen und daher i. allg. zu pessimistisch sind.

— auf Grund der Kompliziertheit der Algorithmen der linearen Algebra und der dadurch hervorgerufenen technischen Schwierigkeiten bei der Fehleranalyse i. allg. nur vereinfachte und vergröberte obere Schranken für die "wahren" minimalen Fehlerkumulationskonstanten berechnet werden können.

Sofern die Fehleranalyse sachgerecht durchgeführt wird und Vergröberungen bei der Analyse verschiedener Verfahren auf vergleichbarem Niveau liegen, erlauben die Fehlerkumulationskonstanten jedoch die qualitative und quantitative Einschätzung und damit den Vergleich verschiedener Verfahren. Gegebenenfalls wird es dabei nötig sein, auch auf empirisch, durch umfangreiche Testrechnungen gewonnene Kenntnisse zurückzugreifen, um ein realistisches Bild zu bekommen. Zur Fehlerabschätzung für ein konkretes Problem sollten a-posteriori-Abschätzungen unter Verwendung der berechneten Lösung – z. B. mittels geeigneter Residualkriterien – durchgeführt werden.

(iv) Für lokal lipschitzstetige Aufgaben folgt aus 2.1.8

$$\Delta \boldsymbol{a}_{\mathrm{opt}}(\boldsymbol{e}, \boldsymbol{\nu}) = L(\boldsymbol{e}) \boldsymbol{\nu} \|\boldsymbol{e}\|,$$

und numerische Gutartigkeit zieht numerische Stabilität nach sich, wobei $F_s = F_w$.

(v) Für gewisse Aufgabenklassen ist es sinnvoll und möglich, nur gewisse Teile der Eingangsdaten mit Störungen zu belasten. Ebenso kann numerische Stabilität auch nur in bezug auf gewisse Teile der Eingangsdaten definiert werden. Dabei ergeben sich i. allg. andere Fehlerkumulationskonstanten als bei Berücksichtigung aller Eingangsdaten.

(vi) Die Abschätzungen (44) bis (46) sind kollektive Normabschätzungen. Wenn statt der Normen Beträge verwendet werden, erhält man individuelle elementweise Abschätzungen und dementsprechend elementweise zu verstehende Gutartigkeit und Stabilität.

Zum Abschluß geben wir ein Beispiel eines numerisch instabilen Verfahrens an.

2.3.20. Beispiel. Wir betrachten für n = 2 das lineare Gleichungssystem

$$a_{11}x_1 + a_{12}x_2 = b_1,$$

 $a_{21}x_1 + a_{22}x_2 = b_2.$

In der Klasse der Aufgaben mit $a_{11} \neq 0$ und $a_{11}a_{22} - a_{12}a_{21} \neq 0$ können x_1 und x_2 wie folgt berechnet werden: Die erste Gleichung wird mit $-a_{21}/a_{11}$ multipliziert und zur zweiten addiert, womit sich

$$\left(a_{21}-rac{a_{21}}{a_{11}}\,a_{11}
ight)x_1+\left(a_{22}-rac{a_{21}}{a_{11}}\,a_{12}
ight)x_2=b_2-rac{a_{21}}{a_{11}}\,b_1$$

ergibt.

Der Koeffizient bei x_1 ist 0, so daß x_2 aus dieser Gleichung berechnet werden kann. Aus der ersten Originalgleichung ergibt sich dann x_1 gemäß

$$x_1 = (b_1 - a_{12}x_2)/a_{11}. \tag{47}$$

Man kann zeigen, daß dieses Verfahren (Gaußscher Algorithmus ohne Pivotisierung, siehe Kapitel 5) für die Berechnung von x_2 gutartig und damit stabil ist. Dagegen ist es numerisch instabil als Verfahren zur Berechnung von x_1 .

Zahlenbeispiel: Rechnung in $\Re(10, 5, 9, 9)$ mit symmetrischer Rundung,

$$\begin{array}{l} 0.0021305x_1 + 6.7034x_2 = 19.041, \\ 9.8702 \quad x_1 - 1.8132x_2 = 4.7207. \end{array}$$

Nach dem angegebenen Algorithmus erhält man $x_2 = 2.8401$, $x_1 = 1.4081$, die exakten Werte sind $x_2^* = 2.84018101..., x_1^* = 1.00003203...$ Der erzeugte Rundungsfehler $|x_1^* - x_1| = 0.408$ überschreitet das unvermeidliche Fehlerniveau $(\varDelta x_1)_{opt} = 4.9 \times 10^{-5}$ von x_1 um den Faktor 8000! Für die vorliegende Aufgabe ist dabei $(\varDelta x_1)_{opt}$ exakt berechnet worden. Werden die beiden Zeilen des Gleichungssystems vertauscht, so liefert derselbe Algorithmus $x_2 = 2.8402$, $x_1 = 1.0000$, also die gerundeten Werte der exakten Lösung! In der ersten Version wird die In stabilität durch Formel (47) hervorgerufen, bei welcher der Zähler die Gestalt fl (19.041 – 19.038) = 0.003 hat, also zu Auslöschung führt, während der erzeugte Rundungsfehler in der nachfolgenden Division durch die kleine Zahl $a_{11} = 0.0021305$ sehr verstärkt wird, vgl. die Nachbemerkung zu 2.2.13. Selbst für den gerundeten Wert $x_2 = 2.8402$ der exakten Lösung würde sich aus (47) die erste Komponente zu $x_1 = \text{fl}((19.041 - 19.038)/0.0021305) = 0.98498$ ergeben, also immer noch mit einem Fehler $|x_1^* - x_1| = 1.5 \times 10^{-2}$, der das unvermeidliche Fehlerniveau um den Faktor 300 überschreitet.

Übungsaufgaben

Ü 2.3.1. Zur Berechnung der Summe $z^* = \sum_{i=1}^n x_i \text{ von } n = 2^p$ Zahlen $x_i \in \Re$ kann die sog. paarweise Summation — auch binäre oder Kaskadensummation genannt — verwendet werden, die für $p = 3, n = 2^3 = 8$ wie folgt abläuft:

$$x_{1} \underbrace{x_{1}^{(1)} = x_{1} + x_{2}}_{x_{1}^{(2)} = x_{1}^{(1)} + x_{2}^{(1)}} \underbrace{x_{2}^{(1)} = x_{3} + x_{4}}_{x_{1}^{(1)} = x_{5} + x_{6}} \underbrace{x_{6} \\ x_{7} \underbrace{x_{7} \\ x_{1}^{(1)} = x_{1} + x_{2}^{(1)}}_{x_{1}^{(2)} = x_{3}^{(1)} + x_{4}^{(1)}} \underbrace{x_{1}^{(1)} = x_{5} + x_{6} \\ x_{2}^{(1)} = x_{1}^{(1)} + x_{2}^{(1)} \\ x_{1}^{(2)} = x_{1}^{(2)} + x_{2}^{(2)} = z^{*}}$$

Die allgemeine Vorschrift lautet: for i := 1(1)n do $x_i^{(0)} := x_i$ for j := 1(1)p do $| \max x = 2^{p-j}$ for $i := 1(1)\max x$ do $x_i^{(j)} := x_{2i-1}^{(j-1)} + x_{2i}^{(j-1)}$

$$z:=x_1^{(p)}$$

(i) Man zeige durch Induktion, daß für den berechneten Wert z

$$z = \sum_{i=1}^{n} x_i (1 + \mu_i^{(n)}) \quad \text{mit} \quad |\mu_i^{(n)}| \leq pv = (\log_2 n) v$$

gilt. Das Verfahren ist also numerisch gutartig mit $F = \log_2 n$ im Vergleich zu F = n - 1 für die übliche Summationsvorschrift 2.3.1.

(ii) Man überlege sich, daß die Kaskadensummation unter Verwendung von jeweils einem INTEGER- und einem REAL-Feld der Länge p realisiert werden kann, und gebe ein entsprechendes FORTRAN-Programm an.

(iii) Wie ist das Verfahren zu modifizieren, damit es für beliebiges *n* angewendet werden kann? Man erweitere das FORTRAN-Programm aus (ii) entsprechend.

(iv) Die Kaskadensummation ist besonders günstig für Parallelrechnung, d. h. für Rechnung mit mehreren parallel arbeitenden Prozessoren.

Ü 2.3.2. Für $b \in \mathbb{R}^m$ bzw. $A \in \mathbb{R}^{m,n}$ ($b \neq o, A \neq O$) bezeichne b_{μ} bzw. $a_{\mu\nu}$ ein betragsgrößtes Element, d. h.

$$|b_{\mu}| = \max_{i} |b_{i}|$$
 bzw. $|a_{\mu
u}| = \max_{i,j} |a_{ij}|$.

(i) Man beweise die Gültigkeit von

$$\|\boldsymbol{b}\|_{\boldsymbol{p}} \leq C_{\boldsymbol{p}} \|b_{\boldsymbol{\mu}}\|$$
 bzw. $\||\boldsymbol{A}|\|_{\boldsymbol{p}} \leq D_{\boldsymbol{p}} \|a_{\boldsymbol{\mu}\boldsymbol{\nu}}\|$

mit Konstanten

$$C_p = \begin{cases} m & \text{für } p = 1, \\ \sqrt{m} & \text{für } p = 2, \\ 1 & \text{für } p = \infty \end{cases} \quad \text{bzw.} \quad D_p = \begin{cases} m & \text{für } p = 1, \\ n & \text{für } p = \infty, \\ \sqrt{mn} & \text{für } p = 2, F. \end{cases}$$

(ii) Unter Verwendung von (i) zeige man, daß aus den Normabschätzungen

$$\frac{\|\tilde{\boldsymbol{b}} - \boldsymbol{b}\|_{p}}{\|\boldsymbol{b}\|_{p}} \leq \nu \quad \text{bzw.} \quad \frac{\|\tilde{\boldsymbol{A}} - \boldsymbol{A}\|_{p}}{\||\boldsymbol{A}\|\|_{p}} \leq \nu$$
(48)

die individuellen Abschätzungen

$$\frac{|b_i - b_i|}{|b_{\mu}|} \leq C_p \nu \text{ bzw. } \frac{|\tilde{a}_{ij} - a_{ij}|}{|a_{\mu\nu}|} \leq D_p \nu$$

folgen. Dies bedeutet, daß (48) nur für Elemente in der Größenordnung von $|b_{\mu}|$ bzw. $|a_{\mu\nu}|$ eine hohe individuelle relative Genauigkeit garantiert.

Ü 2.3.3. Man beweise die Behauptung aus Bemerkung 2.3.19(iv).

 \dot{U} 2.3.4. Bei der Lösung linearer Gleichungssysteme spielt das Residuum r = b - Ax einer Näherungslösung x eine Rolle.

(i) Man zeige, daß das gemäß

for i := 1(1)n do

 $\begin{vmatrix} y := 0 \\ \text{for } j := 1(1)n \text{ do } y := y + a_{ij} * x_j \\ r_i := b_i - y \end{cases}$

in üblicher Weise berechnete Residuum $r := \operatorname{fl}(\boldsymbol{b} - A\boldsymbol{x})$ der Gleichung

$$r_i = (-y_i + b_i) + [\varepsilon_i/(1 + \varepsilon_i)] r_i, \qquad |\varepsilon_i| \le \nu,$$
(49)

genügt mit $y_i := \operatorname{fl}((Ax)_i)$. Mit (36), (37), (40) folgere man hieraus die a-priori-Abschätzung

$$|\delta \boldsymbol{r}| \le \nu(|\boldsymbol{r}| + n |\boldsymbol{A}| |\boldsymbol{x}|) \tag{50}$$

für den erzeugten Rundungsfehler $\delta r = r - (b - Ax)$.

(ii) Man überlege sich, daß durch die erweiterte Vorschrift

for i := 1(1)n do

$$\begin{array}{l} y := d := 0 \\ \textbf{for } j := 1(1)n \ \textbf{do} \\ [g := a_{ij} * x_j, y := y + g, d := d + |g| + |y|] \\ r_i := b_i - y, d_i := v * (d + |r_i|) \end{array}$$

neben \boldsymbol{r} noch eine a-posteriori-Schranke $\boldsymbol{d} = (d_i)$ mit $|\boldsymbol{\delta r}| \leq \boldsymbol{d}$ berechnet wird. Wie erhöht sich der Aufwand gegenüber (i)?

(iii) Man gebe zu beiden Versionen numerisch äquivalente Varianten an, bei denen die Elemente von A spaltenweise abgearbeitet werden.

Bemerkungen zum Kapitel 2

B 2.1. Wenn x den exakten Wert einer reellen Größe und y eine fehlerbehaftete Näherung für x bezeichnet, wird in der klassischen Fehlerrechnung $\delta := y - x$ als absoluter und $\varepsilon := (y - x)/x$ als relativer Fehler eingeführt. In der Regel – z. B. im Abschnitt 2.2 – gehen

wir ebenso vor. An anderen Stellen — etwa im Abschnitt 2.1 — vertauschen wir die Rolle von x und y, da dann die angestrebten Aussagen leichter zu formulieren sind. Man beachte, daß bei kleinem relativem Fehler beide Definitionen in erster Ordnung denselben Wert liefern, siehe Ü 2.1.8.

B 2.2. Die in den Abschnitten 2.1 und 2.3 eingeführten Begriffe zur Charakterisierung von numerischen Problemen und Algorithmen sind — unter demselben oder ähnlichem Namen und mit derselben oder ähnlicher Bedeutung — von einer Reihe von Autoren betrachtet worden, vgl. etwa WILKINSON [63, 65, 71], BAUER [66, 74], FADDEEV/FADDEEVA [68], BABUŠKA [69, 72], KAHAN [71] und viele andere. Unsere Darstellung folgt dem Konzept von KIEŁBASIŃSKI [78], allerdings in einer für die Probleme der linearen Algebra ausreichenden Vereinfachung.

B 2.3. Die hier benutzte Notation für Computerzahlen geht auf WILKINSON [63] zurück, wo auch die Grundlagen der Rundungsfehleranalyse an vielen Beispielen erläutert worden sind. WILKINSON selbst weist auf Arbeiten von VON NEUMANN/GOLDSTINE [47] und GIVENS [54] hin, in denen erstmals eine Rückwärtsanalyse vorkommt. In einer Vielzahl von Arbeiten ist seitdem der technische Apparat zur Fehleranalyse weiterentwickelt worden. Interessante Ergebnisse zur automatischen Rundungsfehleranalyse sind bei MILLER/WRATHALL [80] zu finden, wo der Computer als analytisches Hilfsmittel eingesetzt wird.

B 2.4. Eine detaillierte Beschreibung der auf realen Computern vorhandenen Gleitpunktarithmetiken gibt STERBENZ [74], vgl. auch BROWN [81]. Zur stochastischen Rundungsfehleranalyse siehe wieder STERBENZ, aber auch VOEVODIN [69b]. Eine Einführung in die Computerarithmetik ist bei VOEVODIN [77] und in den meisten neueren Lehrbüchern der Numerischen Mathematik zu finden.

B 2.5. Die Basis und Mantissenlänge der Zahldarstellung und die Art der Rundung lassen sich bei Bedarf automatisch mit den von MALCOLM [72] und GENTLEMAN/MAROVICH [74] angegebenen FORTRAN-Programmen ermitteln.

B 2.6. Als Beispiel für die computergerechte Programmierung unter Vermeidung von Genauigkeitsverlusten sowie Über- und Unterlauf kann das von BLUE [78] angegebene Programm zur Berechnung der Euklidischen Norm dienen. Dort werden die Komponenten des Vektors nur einmal durchlaufen. Ähnlich arbeiten die von LAWSON et al. [79] und DONGARRA et al. [79] angegebenen Programme.

B 2.7. Eine Alternative zur rundungsfehlerbehafteten Gleitpunktarithmetik stellen sog. exakte Arithmetiken dar, siehe GREGORY-KRISHNAMURTY [84]. Die Realisierung solcher Arithmetiken ist jedoch aufwendig, und sie sind i. allg. nur für spezielle Aufgaben mit fehlerfreien Eingangsdaten zu empfehlen.

B 2.8. Die bei Computerrechnung entstehenden Rundungsfehler lassen sich auch mit den Methoden der Intervallmathematik erfassen, siehe etwa MOORE [79] und NICKEL [80]. Solche Methoden erlauben für ein konkretes Problem und die mittels eines konkreten Algorithmus berechneten Resultate die Angabe von strengen Fehlerschranken, allerdings mit vergleichsweise hohem Aufwand. Sie erlauben jedoch nicht den Vergleich verschiedener Algorithmen für bestimmte Aufgabenklassen. Letzteres ist dagegen mit den Methoden der Rundungsfehleranalyse möglich, die wiederum für die Berechnung von Fehlerschranken für einzelne konkrete Aufgaben nicht besonders geeignet sind. Da wir hauptsächlich am qualitativen Vergleich von Algorithmen (über ihre Fehlerkumulationskonstanten) interessiert sind, gehen wir auf die Intervallmathematik nicht ein.

B 2.9. Durch Kombination von Intervallmethoden und Rechnung in variabler, mehrfach hoher Genauigkeit — beides läßt sich sowohl software- als auch neuerdings hardwaremäßig realisieren — können Ergebnisse mit beliebig vorgebbaren Fehlerschranken garantiert werden. vgl. das in jüngster Zeit von KULISCH et al. [85] u. a. entwickelte und von IBM vertriebene System ACRITH. B 2.10. Für eine vorgegebene Problemklasse ist die Frage interessant, wieviel arithmetische Operationen mindestens zur Lösung eines beliebigen Problems der Klasse erforderlich sind und wie optimale Algorithmen beschaffen sind. Die Beantwortung dieser Fragen ist Gegenstand der sog. Kompliziertheitstheorie (engl. ,,complexity", russ. ,,теория сложности"). Da jedes Element der Eingangsdaten bei mindestens einer Operation beteiligt sein muß, ist eine untere Schranke für die Anzahl der Rechenoperationen von der Form $\sim K_1 N$, N Anzahl der skalaren Eingangsdaten. In diesem Sinne sind die im Abschnitt 2.3 angegebenen naiven Algorithmen für Summation und Skalarprodukt als optimal anzusehen. Eine nichttriviale Situation liegt bei der Matrixmultiplikation - Aufgabe 3 aus Abschnitt 2.3.D - vor. Für den Sonderfall $C = AB, A, B \in \mathbb{R}^{n,n}$, werden dort n^3 opms benötigt, während die untere Schranke von der Form $\sim K_1 n^2$ ist. Den ersten wesentlichen Beitrag zu diesem Problem hat STRASSEN [69] mit einem Verfahren geleistet, das $\sim K_2 n^{2.81}$ Operationen erfordert; inzwischen ist man etwa bei $\sim K_3 n^{2.5}$ angelangt. Man beachte jedoch, daß die Vorfaktoren K_i groß sind und daher die kleineren Exponenten erst für hohe Dimensionen n wirksam werden, wo andererseits die Voraussetzung der vollen Besetztheit von A und B meist nicht mehr realistisch ist. Unter gewissen Voraussetzungen hat MILLER [75] ferner gezeigt, daß ein numerisch stabiles Verfahren zur Matrixmultiplikation $\sim K_4 n^3$ Operationen erfordert. Aus den angegebenen Gründen haben die erwähnten "schnellen" Multiplikationsverfahren bislang keine praktische Bedeutung erlangt, und die naive Realisierung aus Abschnitt 2.3 kann bei Forderung der numerischen Stabilität als fast optimal angesehen werden.

3. Elementare Transformationsmatrizen

Die Aufgaben der linearen Algebra sind besonders einfach zu lösen, wenn die jeweiligen Koeffizientenmatrizen eine günstige Gestalt haben. Zum Beispiel ist ein Gleichungssystem mit einer Dreiecksmatrix viel leichter lösbar als ein System mit einer allgemeinen Matrix. Ganze Klassen von Algorithmen der linearen Algebra beruhen auf dieser Tatsache: Sie transformieren das zu lösende Problem durch eine endliche bzw. auch unendliche Folge von einfachen Transformationen in ein äquivalentes Problem mit einer Matrix der gewünschten günstigen Gestalt, z. B. in ein Problem mit einer Dreiecksmatrix. Dabei heißen zwei Probleme äquivalent, wenn sie dieselben Ausgangsdaten, also dieselben Lösungen besitzen. Die Transformation geschieht in der Regel durch Links- bzw. Rechtsmultiplikation mit geeigneten elementaren Transformationsmatrizen, wobei in der transformierten Matrix an den gewünschten Stellen nach und nach Nullen erzeugt werden.

Verfahren dieser Struktur werden häufig direkte Verfahren genannt; das vorliegende Buch ist vorwiegend solchen Verfahren gewidmet. Im Gegensatz dazu wird bei den sog. *iterativen Verfahren* die Datenmatrix nicht verändert, d. h., in jedem Iterationsschritt wird auf die originale Matrix der Eingangsdaten zurückgegriffen. Die Grenzen zwischen beiden Verfahrensklassen sind jedoch fließend.

Die erwähnten elementaren Transformationsmatrizen können als Werkzeuge der numerischen linearen Algebra angesehen werden, denn mit ihrer Hilfe lassen sich die unterschiedlichsten Algorithmen zur Lösung verschiedener Aufgabenklassen aufbauen.

Zur Illustration betrachten wir die Lösung des linearen Gleichungssystems

$$3x_1 - 6x_2 + 2x_3 = -3, -2x_1 + 4x_2 + x_3 = 16, x_1 - x_2 + x_3 = 4.$$
(1)

Nach dem aus der Schule bekannten Gaußschen Eliminationsverfahren kann x_1 aus der zweiten und dritten Gleichung eliminiert werden, indem die erste mit 2/3 bzw. -1/3 multipliziert und zur zweiten bzw. dritten addiert wird:

$$3x_1 - 6x_2 + 2x_3 = -3,$$

$$0 \cdot x_2 + \frac{7}{3} x_3 = 14,$$

$$x_2 + \frac{1}{3} x_3 = 5.$$
(2)

Da in der neuen zweiten Gleichung auch der Koeffizient bei x_2 verschwindet, geht das System (2) durch Vertauschung der letzten beiden Zeilen in das System

$$3x_1 - 6x_2 + 2x_3 = -3,$$

$$x_2 + \frac{1}{3}x_3 = 5,$$

$$\frac{7}{3}x_3 = 14$$
(3)

mit einer Dreiecksmatrix als Koeffizientenmatrix über. Aus (3) lassen sich die Unbekannten in der Reihenfolge $x_3 = 6$, $x_2 = 3$ und $x_1 = 1$ leicht berechnen.

Der Übergang von (1) nach (2) kann in Matrixschreibweise als Linksmultiplikation der Koeffizientenmatrix und des Vektors der rechten Seiten mit der Transformationsmatrix

$$\boldsymbol{L_1} = \begin{pmatrix} 1 & 0 & 0 \\ \frac{2}{3} & 1 & 0 \\ -\frac{1}{3} & 0 & 1 \end{pmatrix},$$

beschrieben werden, der von (2) nach (3) analog mittels

$$\boldsymbol{T}_{2,3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}.$$

Der Untersuchung dieser und weiterer Typen elementarer Transformationsmatrizen sind die folgenden Abschnitte gewidmet.

3.1. Permutationsmatrizen

Eine Permutation ist eine Abbildung (Zuordnung) $k: \{1, 2, ..., n\} \rightarrow \{k(1), k(2), ..., k(n)\}$ mit der Eigenschaft, daß die Funktionswerte k(i), i = 1, ..., n, jeden der Werte 1, ..., n genau einmal annehmen. Das *n*-Tupel $\{k(1), ..., k(n)\}$ heißt die durch kerzeugte Permutation der Zahlen $\{1, ..., n\}$; z. B. legt

$$\{k(1), k(2), k(3), k(4)\} = \{2, 4, 1, 3\}$$

für n = 4 eine Permutation k fest. Offensichtlich ist die Permutation k durch eine solche Tabelle $\{k(1), \ldots, k(n)\}$ eindeutig beschrieben.

⁷ Schwetlick, Numerische Algebra

Zu jeder Permutation k gibt es eine eindeutig festgelegte inverse Permutation l mit l(k(i)) = i (i = 1, ..., n), welche die durch k erzeugte Umordnung wieder rückgängig macht. Im Beispiel ist

$$\{l(1), l(2), l(3), l(4)\} = \{3, 1, 4, 2\}$$

die zu k inverse Permutation.

Durch eine Permutation k wird jedem Vektor $x \in \mathbb{R}^n$ nach der Vorschrift

$$\boldsymbol{x} = (x_i) \rightarrow \boldsymbol{y} = (y_i) \quad \text{mit} \quad y_i := x_{\boldsymbol{k}(i)} \qquad (i = 1, \dots, n) \tag{4}$$

ein Vektor $y \in \mathbb{R}^n$ zugeordnet, der die (durch k erzeugte) Permutation des Vektors x genannt wird. Der Übergang von x zu y gemäß (4) kann in eindeutiger Weise durch eine Matrix $P \in \mathbb{R}^{n,n}$ in der Gestalt

$$\boldsymbol{y} = \boldsymbol{P}\boldsymbol{x} \tag{5}$$

geschrieben werden; im Beispiel ist

$$\boldsymbol{P} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{pmatrix}.$$

3.1.1. Definition. Die Matrix $P = (p_{ij}) \in \mathbb{R}^{n,n}$ mit

$$p_{ij} = \begin{cases} 1 & \text{für } j = k(i), \\ 0 & \text{für } j \neq k(i) \end{cases} (i, j = 1, ..., n)$$

heißt die durch k erzeugte Permutationsmatrix.

Offensichtlich ist eine Matrix P genau dann eine Permutationsmatrix, wenn in jeder Zeile und Spalte außer Nullen genau eine 1 vorkommt. Dabei gibt k(i) den Spaltenindex der 1 in der *i*-ten Zeile an; entsprechend ist l(j) der Zeilenindex der 1 in der *j*-ten Spalte.

Die Transponierte einer Permutationsmatrix und das Produkt von Permutationsmatrizen sind wieder Permutationsmatrizen. Wegen $P^{\dagger}P = PP^{\dagger} = I$ gilt ferner

$$\boldsymbol{P}^{\intercal} = \boldsymbol{P}^{-1}, \tag{6}$$

d. h., Permutationsmatrizen sind orthogonal, und P^{\intercal} charakterisiert die inverse Permutation. Bei Anwendung von P auf eine Matrix $A \in \mathbb{R}^{n,m}$ mit den Zeilen $a^{i\intercal}$ bzw. auf eine Matrix $B \in \mathbb{R}^{m,n}$ mit den Spalten b^{j} ergibt sich

$$P\begin{pmatrix} -a^{1\mathsf{T}} - \\ \vdots \\ -a^{n\mathsf{T}} - \end{pmatrix} = \begin{pmatrix} -a^{k(1)\mathsf{T}} - \\ \vdots \\ -a^{k(n)\mathsf{T}} - \end{pmatrix} \quad \text{bzw.} \quad \begin{pmatrix} | & | \\ b^{1} \dots b^{n} \\ | & | \end{pmatrix} P = \begin{pmatrix} | & | \\ b^{l(1)} \dots b^{l(n)} \\ | & | \end{pmatrix}, \tag{7}$$

d. h., Linksmultiplikation mit P bewirkt eine Permutation der Zeilen gemäß k, und Rechtsmultiplikation eine Permutation der Spalten gemäß der inversen Permutation l. Sollen die Spalten gemäß k permutiert werden, so ist von rechts mit P^{\intercal} zu multiplizieren.

Die einfachsten Permutationen bestehen in der Vertauschung von zwei Elementen. bei Vektoren also in der Vertauschung von zwei Komponenten $i \neq j$. Die entsprechenden Matrizen haben die Gestalt

$$m{T}_{ij} = egin{pmatrix} i & j \ 1 & \vdots & \vdots \ \dots & 0 & \dots & 1 & \dots \ \vdots & \vdots & \vdots \ \dots & 1 & \dots & 0 & \dots \ \vdots & \vdots & 1 \end{pmatrix} m{i}$$
-te Zeile j-te Zeile

•3.1.2. Definition. Eine Matrix $T_{ij} \in \mathbf{R}^{n,n}$ $(i, j = 1, ..., n; i \neq j)$, die sich von der Einheitsmatrix nur in den vier Elementen

$$(T_{ij})_{ii} = (T_{ij})_{jj} = 0, \quad (T_{ij})_{ij} = (T_{ij})_{ji} = 1$$

 $(\mathbf{r}_{ij})_{ii} = (\mathbf{r}_{ij})_{jj} = 0,$ $(\mathbf{T}_{ij})_{ij} = (\mathbf{T}_{ij})_{ji} = 1$ unterscheidet, heißt Vertauschungsmatrix. Für i = j wird $\mathbf{T}_{ii} = \mathbf{T}_{jj} = \mathbf{I}$ gesetzt (triviale Vertauschung).

Vertauschungsmatrizen sind symmetrisch, d. h. $T_{ij} = T_{ij}^{\mathsf{T}}$ und folglich $T_{ij}^2 = I$, und es gilt

$$\det \left(\boldsymbol{T}_{ij} \right) = \begin{cases} 1 & \text{für } i = j, \\ -1 & \text{für } i \neq j. \end{cases}$$

$$\tag{8}$$

Wie im folgenden Beispiel gezeigt wird, läßt sich eine Permutation stets als Folge von Vertauschungen darstellen, wobei im i-ten Schritt das Element i mit einem geeigneten Element $s(i) \geq i$ vertauscht wird:

$$\begin{pmatrix} 1\\2\\3\\4 \end{pmatrix} \xrightarrow{\mathbf{P}} \begin{pmatrix} 2\\4\\1\\3 \end{pmatrix} \text{ äquivalent zu } \begin{pmatrix} 1\\2\\3\\4 \end{pmatrix} \xrightarrow{\mathbf{T}_{12}} \begin{pmatrix} 2\\1\\3\\4 \end{pmatrix} \xrightarrow{\mathbf{T}_{24}} \begin{pmatrix} 2\\4\\3\\1 \end{pmatrix} \xrightarrow{\mathbf{T}_{34}} \begin{pmatrix} 2\\4\\1\\3 \end{pmatrix},$$

also s(1) = 2, s(2) = 4, s(3) = 4 und $P = T_{34}T_{24}T_{12}$.

3.1.3. Aussage. Jede Permutationsmatrix läßt sich in eindeutiger Weise als Produkt von n-1 Vertauschungsmatrizen gemäß

$$P = T_{n-1,s(n-1)}T_{n-2,s(n-2)}\cdots T_{2,s(2)}T_{1,s(1)}$$
(9)

 $P = T_{n-1,s(n-1)}T_{n-2,s(n-2)}\cdots T_{2,s(2)}T_{1,s(1)}$ mit $i \leq s(i) \leq n \ (i = 1, ..., n - 1)$ darstellen.

Die zu (9) inverse Permutation ist dann

$$P^{-1} = P^{\mathsf{T}} = T_{1,s(1)} T_{2,s(2)} \cdots T_{n-2,s(n-2)} T_{n-1,s(n-1)}.$$
(10)

In der numerischen linearen Algebra kommen Permutationen fast ausschließlich als Folge von Vertauschungen gemäß (9) vor. In vielen Fällen ist es sinnvoll, die Informationen über diese Permutationen in Form der Tabelle $\{s(1), s(2), \ldots, s(n-1)\}$

7*

in einem Integerfeld zu speichern. Zu gegebenem x kann dann y = Px sogar auf dem Platz von x — also in situ — berechnet werden; dasselbe gilt für die inverse Permutation.

In anderen Fällen ist eine Darstellung der Permutation als Tabelle $\{k(1), \ldots, k(n)\}$ günstiger. Eine solche Tabelle entsteht mittels der s(i), wenn der nachfolgend beschriebene Algorithmus PER auf den Vektor $\boldsymbol{x} = (x_i) := (i)$ angewendet wird.

3.1.4. Algorithmus. Berechnung von $\boldsymbol{x} := \boldsymbol{P}\boldsymbol{x}$ (PER) bzw. $\boldsymbol{x} := \boldsymbol{P}^{-1}\boldsymbol{x}$ (PERINV) PER: for i := 1(1)n - 1 do PERINV: for i := n - 1(-1)1 do if s(i) > i then $\begin{vmatrix} \boldsymbol{z} := x_i \\ x_i := x_{s(i)} \\ x_{s(i)} := \boldsymbol{z} \end{vmatrix}$ $\begin{vmatrix} \boldsymbol{z} := x_i \\ x_i := x_{s(i)} \\ x_{s(i)} := \boldsymbol{z} \end{vmatrix}$

Durch den Algorithmus PER wird x mit Px, durch PERINV wird x mit $P^{-1}x$ überspeichert. Beide Algorithmen sind Beispiele für sog. nichtnumerische Algorithmen, da in ihnen nur logische und Speicheroperationen, aber keine arithmetischen Operationen ausgeführt werden.

Wir weisen nochmals darauf hin, daß Permutationen bei der Programmierung *nur* über die Tabellen $\{k(i)\}$ bzw. $\{s(i)\}$ und *niemals* über die zugehörigen Matrizen realisiert werden; letztere dienen ausschließlich der mathematischen Beschreibung.

Übungsaufgaben

Ü 3.1.1. Man überlege sich Algorithmen, die aus der Tabelle $\{k(i)\}$ die Tabellen $\{l(i)\}$ der inversen Permutation sowie die Tabelle $\{s(i)\}$ für die Darstellung (9) berechnen.

Ü 3.1.2. Es ist zu zeigen, daß jede Permutation P auch eindeutig in der Form

$$P = T_{n,r(n)}T_{n-1,r(n-1)}\cdots T_{3,r(3)}T_{2,r(2)}$$
(11)

mit $1 \leq r(i) \leq i$ (i = 2, ..., n) geschrieben werden kann.

3.2. Nichtorthogonale elementare Transformationsmatrizen

Im folgenden betrachten wir Matrizen der Gestalt

$$M_{k} = \begin{pmatrix} 1 & m_{1} & & \\ \ddots & \vdots & & \\ & \ddots & m_{k-1} & \\ & 1 & & \\ & & \ddots & \\ & & m_{k+1} & \\ & \vdots & \ddots & \\ & & & m_{n} & 1 \end{pmatrix} \quad \text{bzw.} \quad L_{k} = \begin{pmatrix} 1 & & & \\ \ddots & & & \\ & \ddots & & \\ & & & \ddots & \\ & & & l_{k+1} & \ddots & \\ & & & \vdots & \ddots & \\ & & & l_{n} & 1 \end{pmatrix}, \tag{1}$$

wie sie uns in der Einführung zum dritten Kapitel bereits begegnet sind.

3.2.1. Definition. (i) Eine Matrix $M_k \in \mathbb{R}^{n,n}$ (k = 1, ..., n) des Typs $M_k = M_k(m) = I + me^{kT}, \qquad m = (m_i) \in \mathbb{R}^n,$ mit $m_k = e^{k \mathsf{T}} m = 0$ heißt nichtorthogonale elementare Transformationsmatrix, abgekürzt: NT-Matrix. (ii) Eine NT-Matrix $L_k \in \mathbb{R}^{n,n}$ (k = 1, ..., n - 1) des Typs $oldsymbol{L}_k = oldsymbol{L}_k(oldsymbol{l}) = oldsymbol{I} + oldsymbol{l} e^{k op}, \qquad oldsymbol{l} = (oldsymbol{l}_i) \in oldsymbol{R}^n,$ die gleichzeitig eine untere Dreiecksmatrix ist, d. h., für die $l_i = e^{i\mathsf{T}} l = 0 \qquad (i = 1, \dots, k)$ gilt, heißt LNT-Matrix. Wegen $M_k(m) M_k(-m) = [I + me^{kT}] [I - me^{kT}]$ $= I + me^{k\intercal} - me^{k\intercal} - me^{k\intercal} me^{k\intercal} = I$

ist jede NT-Matrix regulär, und die Inverse

$$\boldsymbol{M}_{\boldsymbol{k}}(\boldsymbol{m})^{-1} = \boldsymbol{M}_{\boldsymbol{k}}(-\boldsymbol{m}) = \boldsymbol{I} - \boldsymbol{m}\boldsymbol{e}^{\boldsymbol{k}\mathsf{T}}$$
⁽²⁾

ist wieder eine NT-Matrix. Speziell gilt

$$\boldsymbol{L}_{\boldsymbol{k}}(\boldsymbol{l})^{-1} = \boldsymbol{L}_{\boldsymbol{k}}(-\boldsymbol{l}) = \boldsymbol{I} - \boldsymbol{l}\boldsymbol{e}^{\boldsymbol{k}\boldsymbol{\intercal}},\tag{3}$$

d. h., auch die Inverse einer LNT-Matrix ist wieder eine LNT-Matrix. Die Formeln (2) und (3) besagen: Die Inverse einer NT- bzw. LNT-Matrix entsteht durch Vorzeichenumkehr der nichttrivialen Elemente m_i bzw. l_i . Da die LNT-Matrizen in der numerischen linearen Algebra eine besondere Rolle spielen, gehen wir im folgenden ausführlich auf diese Unterklasse der NT-Matrizen ein.

3.2.2. Aussage. Das Produkt

 $L = L_1 L_2 \cdots L_{n-2} L_{n-1}$ der LNT-Matrizen

$$L_k = L_k(l^k), \qquad l^k = (0, ..., 0, l_{k+1,k}, ..., l_{nk})^T$$

ist eine untere Dreiecksmatrix mit Einsdiagonale, und für die Nichtdiagonalelemente gilt

$$(\mathbf{L})_{ij} = l_{ij}$$
 $(i > j, i, j = 1, ..., n).$

Für
$$n = 4$$
 ist also

 $\begin{pmatrix} 1 & & \\ l_{21} & 1 & \\ l_{31} & 1 & \\ l_{41} & & 1 \end{pmatrix}$
 $\begin{pmatrix} 1 & & \\ & 1 & \\ & l_{32} & 1 & \\ & & l_{43} & 1 \end{pmatrix}$
 $\begin{pmatrix} 1 & & \\ l_{21} & 1 & \\ & l_{31} & l_{32} & 1 & \\ & & l_{41} & l_{42} & l_{43} & 1 \end{pmatrix}$

Nachrechnen dieser Beziehung von rechts nach links zeigt, wie der Beweis zu 3.2.2. verläuft. Man beachte jedoch, daß das Produkt $L_{n-1}L_{n-2} \cdots L_1L_2$ zwar auch eine untere Dreiecksmatrix, aber i. allg. von L verschieden ist.

Wenn eine Matrix $B \in \mathbb{R}^{n,r}$ mit den Zeilen $b^{i\tau}$ von links mit der LNT-Matrix L_k multipliziert wird, ergibt sich die transformierte Matrix $\overline{B} = L_k B$ wie folgt:

$$\begin{pmatrix} 1 \\ \ddots \\ 1 \\ l_{k+1} & 1 \\ \vdots & \ddots \\ l_n & 1 \end{pmatrix} \begin{pmatrix} \cdots & \boldsymbol{b}^{1\mathsf{T}} & \cdots \\ \vdots \\ \cdots & \boldsymbol{b}^{k\mathsf{T}} & \cdots \\ \vdots \\ \cdots & \boldsymbol{b}^{k+1\mathsf{T}} & \cdots \\ \vdots \\ \cdots & \boldsymbol{b}^{n\mathsf{T}} & \cdots \end{pmatrix} = \begin{pmatrix} \cdots & \boldsymbol{b}^{1\mathsf{T}} & \cdots \\ \vdots \\ \cdots & \boldsymbol{b}^{k\mathsf{T}} & \cdots \\ \cdot \boldsymbol{b}^{k\mathsf{T}} & \cdots \\ \cdot \boldsymbol{b}^{k\mathsf{T}} & \cdots \\ \vdots \\ \cdots & \boldsymbol{b}^{n\mathsf{T}} & \cdots \end{pmatrix}$$
(4)

Man sieht: Die Transformation $\overline{B} = L_k B$ ändert die ersten k Zeilen von B nicht. Die Zeilen i = k + 1, ..., n von \overline{B} entstehen aus denen von B nach der Vorschrift

$$\langle \text{neue } i\text{-te Zeile} \rangle := \langle \text{alte } i\text{-te Zeile} \rangle + l_i * \langle \text{alte } k\text{-te Zeile} \rangle$$

In elementweiser Notation erhalten wir

3.2.3. Algorithmus zur LNT-Transformation. Berechnung von $\overline{B} = (\overline{b}_{ij}) := L_k B$ für $B = (b_{ij}) \in \Re^{n,r}$ und die LNT-Matrix $L_k = L_k(l), \ l = (0, ..., 0, l_{k+1}, ..., l_n)^{\intercal} \in \Re^n$:

$$\bar{b}_{ij} := \begin{cases} b_{ij} & (i = 1, ..., k), \\ b_{ij} + l_i * b_{kj} & (i = k + 1, ..., n) \end{cases} \quad (j = 1, ..., r)$$

Aufwand: (n - k) r opms

Algorithmus 3.2.3 kann auf dem Platz von B gemäß $B := L_k B$ ausgeführt werden. Ob die neuen Elemente zeilen- oder spaltenweise berechnet werden sollten, hängt von den Umständen – z. B. der Art der Speicherung – ab.

Aus (4) ist ersichtlich, daß die Transformation $\overline{B} = L_k B$ auch in der Form

geschrieben werden kann, wobei $B_1^{(k-1)}$ aus den ersten $k-1, B_2^{(n-k-1)}$ aus den letzten n-k+1 Zeilen von $B = B^{(n)}$ besteht. Mit $I^{(k-1)}$ ist die (k-1)-dimensionale Einheitsmatrix und mit $L_1^{(n-k+1)}$ die LNT-Matrix $L_1 = L_1(l^{(n-k+1)})$ der Dimension n-k+1 mit dem erzeugenden Vektor $l^{(n-k+1)} = (0, l_{k+1}, ..., l_n)^{\mathsf{T}}$ bezeichnet worden. Es genügt daher, LNT-Matrizen für den Index k = 1 zu betrachten, da der Fall $k \neq 1$ gemä β (6) auf diese Standardsituation zurückgeführt werden kann.

Die LNT-Matrix L_1 wird benutzt, um in einem vorgegebenen Vektor $a = (a_i) \in \mathbb{R}^n$ unterhalb der ersten Komponente Nullen zu erzeugen. Nach 3.2.3 ergibt sich für $\overline{a} = L_1 a$

$$ar{a}_1 = a_1, \quad ar{a}_i = a_i + l_i * a_1 \quad (i = 2, ..., n).$$

Die Bedingungen $\bar{a}_i = 0$ lassen sich daher erfüllen, wenn $a_1 \neq 0$ ist, und dann mit $l_i = -a_i/a_1$ für beliebige Werte von a_i .

3.2.4. Aussage. Für $a = (a_i) \in \mathbb{R}^n$, $a_1 \neq 0$, transformiert die durch

$$l_1 = 0, \quad l_i = a_i/a_1 \quad (i = 2, ..., n)$$
 (7)

 $\iota_1 = 0, \qquad \iota_i = u_i/u_1 \qquad (\iota = 2, ..., n)$ festgelegte LNT-Matrix $L_1 = L_1(-l), l = (l_i) \in \mathbf{R}^n$, den Vektor a in

$$\bar{a} = L_1(-l) a = a_1 e^1 = (a_1, 0, ..., 0)^{\mathsf{T}}.$$

Die Standardanwendung der LNT-Matrizen besteht in der Transformation von a in $\bar{a} = L_1 a$ gemäß 3.2.4 und der analogen Transformation einer Matrix B in $\overline{B} = L_1 B.$

3.2.5. Eliminationsschritt mittels LNT-Matrix. Aufgabe: Für $\boldsymbol{a} = (a_i) \in \Re^n$ mit $a_1 \neq 0$ und $\boldsymbol{B} = (b_{ij}) \in \Re^{n,r}$ ist $\boldsymbol{L}_1 = \boldsymbol{L}_1(-\boldsymbol{l}),$ $\boldsymbol{l} = (0, l_2, \dots, l_n)^{\mathsf{T}} \in \Re^n$ so festzulegen, daß

$$\bar{\boldsymbol{a}} = \boldsymbol{L}_1(-\boldsymbol{l}) \, \boldsymbol{a} = (\boldsymbol{a}_1, 0, \dots, 0)^{\mathsf{T}}$$

gilt. Mit dem derart bestimmten L_1 ist $\overline{B} = L_1 B$ zu berechnen. Algorithmus:

S1 (Festlegung von l): $l_i := a_i/a_1$ (i = 2, ..., n)

 S2 (Erste Zeile):
 $\bar{a}_1 := a_1, \bar{b}_{1j} := b_{1j}$ (j = 1, ..., r)

 S3 (Übrige Zeilen):
 $\bar{a}_i := 0, \bar{b}_{ij} := b_{ij} - l_i * b_{1j}$ (i = 2, ..., n, j = 1, ..., r)

Autward: (n-1) r opms + (n-1) opm

In den Anwendungen ist a in der Regel die erste Spalte einer Matrix $M \in \Re^{n,r+1}$, und **B** besteht aus den restlichen Spalten 2 bis r + 1. Der Algorithmus kann in situ durch Überspeichern von

$$\begin{pmatrix} a_1 \mid \dots \mid \mathbf{b}^{1\mathsf{T}} \mid \dots \\ a_2 \mid \dots \mid \mathbf{b}^{2\mathsf{T}} \mid \dots \\ \vdots \mid & \vdots \\ a_n \mid \dots \mid \mathbf{b}^{n\mathsf{T}} \mid \dots \end{pmatrix} \quad \text{mit} \quad \begin{pmatrix} a_1 \mid \dots \mid \mathbf{b}^{1\mathsf{T}} \mid \dots \\ l_2 \mid \dots \mid \mathbf{\bar{b}}^{2\mathsf{T}} \mid \dots \\ \vdots \mid & \vdots \\ l_n \mid \dots \mid \mathbf{\bar{b}}^{n\mathsf{T}} \mid \dots \end{pmatrix}$$

ausgeführt werden. Die erzeugten Nullen brauchen selbstverständlich nicht explizit gespeichert zu werden, und S2 wird überflüssig.

Mit den im Abschnitt 2.3 bereitgestellten Hilfsmitteln kann eine Rundungsfehleranalyse des Algorithmus 3.2.5 durchgeführt werden, siehe Ü 3.2.3. Wir geben hier nur die Ergebnisse an.

3.2.6. Rundungsfehleranalyse. Für $a \in \Re^n$ mit $a_1 \neq 0$ und $B \in \Re^{n,r}$ werde $L_1 = L_1(-l), l = (0, l_2, ..., l_n)^{\mathsf{T}}, \bar{a} \in \Re^n$ und $\bar{B} \in \Re^{n,r}$ nach 3.2.5 berechnet. Dann ist

$$\bar{\boldsymbol{a}} = (a_1, 0, \dots, 0)^{\mathsf{T}} = \boldsymbol{L}_1(\boldsymbol{a} + \boldsymbol{\delta}\boldsymbol{a}) \tag{8}$$

$$\bar{\boldsymbol{B}} = \boldsymbol{L}_1(\boldsymbol{B} + \boldsymbol{\delta}\boldsymbol{B}) \tag{9}$$

$$\delta a_1 = 0, \qquad \delta b_{1i} = 0 \tag{10}$$

$$|\delta a_i| \le \nu |a_i|, \tag{11}$$

$$\mathbf{a} = (a_1, 0, \dots, 0)^r = \mathbf{L}_1(\mathbf{a} + b\mathbf{a})$$
(6)
und

$$\mathbf{\overline{B}} = \mathbf{L}_1(\mathbf{B} + d\mathbf{B})$$
(9)
mit Störungen $d\mathbf{a}$ und $d\mathbf{B}$, für die
 $\delta a_1 = 0, \quad \delta b_{1j} = 0$
(10)
sowie
 $|\delta a_i| \leq v |a_i|,$
(11)
 $|\delta b_{ij}| \leq v(|b_{ij}| + 2 |\overline{b}_{ij}|),$
(12)
 $|\delta b_{ij}| \leq v(|b_{ij}| + 2 |\overline{b}_{ij}|),$
(12)

$$|\delta b_{ij}| \leq \nu(|b_{ij}| + 2 |l_i| |b_{1j}|) \tag{13}$$

Wir fassen jetzt die Eingangsdaten zur Matrix $M = (a \mid B)$ zusammen, analog wird $\boldsymbol{\delta M} = (\boldsymbol{\delta a} \mid \boldsymbol{\delta B})$ gebildet. Aus (11), (13) ergibt sich dann

$$|\boldsymbol{\delta}\boldsymbol{M}| \leq \boldsymbol{\nu}(|\boldsymbol{M}| + 2 |(\boldsymbol{o} | \boldsymbol{l}\boldsymbol{b}^{1\mathsf{T}})|),$$

wobei $b^{1\uparrow}$ die erste Zeile von **B** bezeichnet. Unter Beachtung von Ü 1.1.3 folgt hieraus $\|\boldsymbol{\delta}\boldsymbol{M}\| \leq \boldsymbol{\nu}(\|\boldsymbol{M}\| + 2 \|\boldsymbol{l}\| \|\boldsymbol{B}\|)$, wegen $\|\boldsymbol{B}\| \leq \|\boldsymbol{M}\|$ also

$$\|\boldsymbol{\delta}\boldsymbol{M}\| \leq (1+2\|\boldsymbol{l}\|)\|\boldsymbol{M}\|. \tag{14}$$

Kombination von (11) und (12) liefert dagegen

$$| \boldsymbol{\delta} \boldsymbol{M} | \leq arphi \left\{ | \boldsymbol{M} | + 2 \left(rac{\boldsymbol{O} \mid \boldsymbol{O}}{\boldsymbol{O} \mid |\boldsymbol{M}^{(2)}|}
ight)
ight\},$$

wobei $M^{(2)}$ aus den Zeilen 2 bis n von \overline{B} , d. h. aus den in 3.2.5 tatsächlich neu berechneten Elementen besteht. Hieraus ergibt sich

$$\|\delta M\| \le \nu(\|M\| + 2 \|M^{(2)}\|). \tag{15}$$

Als Normen können in (14), (15) die Matrixnormen mit $p = 1, \infty, F$ verwendet werden, in (14) ist für l dann die Vektornorm mit $p = 1, \infty, 2$ zu nehmen. Im Fall der Spektralnorm ist auf den rechten Seiten $||M||_2$ durch $|||M||_2$ zu ersetzen usw., vgl. Abschnitt 2.3.C.

Die aus 3.2.6 folgenden Abschätzungen (14) und (15) besagen: Der Eliminationsschritt 3.2.5 stellt einen numerisch gutartigen Algorithmus dar, wenn die Elininationskoeffizienten l_i im Sinne von $\|l\| \leq K_0$ beschränkt sind – aus (14) folgt dann $F = 1 + 2K_0$ –, oder wenn $\|M^{(2)}\| \leq K_1 \|M\|$ gilt, d. h., wenn die neu berechneten Elemente nicht zu groß werden – aus (15) liest man dann $F = 1 + 2K_1$ ab.

3.2.7. Bemerkung. (i) Wir werden später sehen, daß die Bedingung $||\mathbf{M}^{(2)}|| \leq K_1 ||\mathbf{M}||$ für gewisse Aufgabenklassen automatisch erfüllt ist, z. B. bei Anwendung von 3.2.5 im Gaußschen Algorithmus für Systeme mit diagonaldominanten oder symmetrischen definiten Matrizen, vgl. Kapitel 5. Obwohl dann die l_i durchaus groß werden können, liegt numerische Gutartigkeit vor.

(ii) In anderen Aufgabenklassen kann sowohl ||l|| als auch $||M^{(2)}||$ beliebig groß werden, so daß der Algorithmus instabil wird. In diesem Fall kann die Beschränktheit der l_i durch die im folgenden beschriebene *Stabilisierung* von 3.2.5 erzwungen werden:

S0: Festlegung des Index s, $1 \leq s \leq n$, einer betragsgrößten Komponente von a, d. h.

$$|a_s| = \max\{|a_i|: i = 1, ..., n\}.$$
(16)

Ersatz von a durch $\hat{a} = T_{1s}a$, d. h. Vertauschung der Komponenten a_1 und a_s von a. Analog Ersatz von B durch $\hat{B} = T_{1s}B$, d. h. Vertauschung der Zeilen 1 und s von B.

S1: Féstlegung von $\hat{l} = (0, \hat{l}_2, ..., \hat{l}_n)^{\mathsf{T}}$ gemäß $\hat{l}_i := \hat{a}_i / \hat{a}_1 \ (i = 2, ..., n)$. Wegen (16) ist

 $|\hat{l}_i| \leq 1$ $(i = 2, ..., n), \text{ also } \|\hat{l}\|_{\infty} \leq 1.$ (17)

S2, S3: Wie in 3.2.5 mit \hat{l}_i , \hat{a}_i und \hat{b}_{ij} statt l_i , a_i und b_{ij} .

Diese Stabilisierung ist gleichbedeutend mit dem Ersatz von $L_1(-l)$ durch $L_1(-l)$ T_{1s} .

3.2.8. Definition. Eine LNT-Matrix $L_k(\hat{l}) = l + \hat{l}e^{k\tau}, \ \hat{l} = (0, ..., 0, \hat{l}_{k+1}, ..., \hat{l}_n)^{\intercal}$, mit

 $|\hat{l}_i| \leq 1$ (i = k + 1, ..., n)

heißt stabilisierte LNT-Matrix, abgekürzt: SLNT-Matrix.

Im Fall (17) geht die Abschätzung (14) für die ∞ -Norm in

$$\|\delta M\|_{\infty} \leq 3\nu \|M\| \tag{18}$$

über. Aus dieser Ungleichung folgt: Der wie oben stabilisierte Algorithmus 3.2.5 ist in der Klasse aller Aufgaben mit $a \in \mathbb{R}^n$, $a \neq o$, $B \in \mathbb{R}^{n,\tau}$ numerisch gutartig mit F = 3 in der Maximumnorm.

Übungsaufgaben

Ü 3.2.1. Man beweise: Für $\boldsymbol{a} = (a_i) \in \mathbb{R}^n \text{ mit } a_k \neq 0$ transformiert die Matrix $\boldsymbol{M}_k = \boldsymbol{M}_k(-\boldsymbol{m})$, $\boldsymbol{m} = (m_i) \in \mathbb{R}^n \text{ mit } m_k = 0$, $m_i = a_i/a_k$ für $i \neq k$ den Vektor \boldsymbol{a} in $\bar{\boldsymbol{a}} = \boldsymbol{M}_k \boldsymbol{a} = a_k \boldsymbol{e}^k$ $= (0, ..., 0, a_k, 0, ..., 0)^{\mathsf{T}}$. Ü 3.2.2. Für $1 \leq k < i \leq s \leq n$ ist die Identität $T_{is}L_k(l) = L_k(T_{is}l) T_{is}$ zu beweisen.

Ü 3.2.3. Man leite die in 3.2.6 angegebenen Fehleraussagen ab.

Hinweis: Die Schlüsseloperation aus Schritt S3 von 3.2.5 läßt sich in der Form

$$\bar{b}_{ij} = [b_{ij} - l_i b_{1j} (1 + \varepsilon)] (1 + \vartheta) \quad \text{mit} \quad \varepsilon = \varepsilon_{ij}, \quad \vartheta = \vartheta_{ij}, \quad |\varepsilon|, \, |\vartheta| \leq \nu \tag{19}$$

schreiben, wobei $l_i = \operatorname{fl}(a_i/a_1)$ der berechnete Koeffizient ist. Hieraus folgt

$$\bar{b}_{ij} = b_{ij} + \delta b_{ij} - l_i b_{1j} \tag{20}$$

 \mathbf{mit}

$$\delta b_{ij} = b_{ij}\vartheta - l_i b_{1j} (\varepsilon + \vartheta + \varepsilon \vartheta) = b_{ij}\vartheta - l_i b_{1j} \eta, \qquad |\eta| \leq 2\nu.$$
⁽²¹⁾

Mit der Festlegung $\delta b_{1j} = 0$ ergibt sich (9) aus (20), und aus (21) kann die Abschätzung (13) gefolgert werden. Auflösen von (19) nach b_{ij} und Einsetzen in (21) führt unter Beachtung von $\bar{b}_{1j} = b_{1j}$ auf die alternative Darstellung

$$\delta b_{ij} = \delta \bar{b}_{ij} = \bar{b}_{ij} \vartheta / (1 + \vartheta) - l_i \bar{b}_{1j} \varepsilon.$$
⁽²²⁾

Schließlich erhält man durch Auflösen von (20) nach $l_i b_{1j}$ und Einsetzen in (22)

$$\delta b_{ij} = \bar{b}_{ij} \eta / [(1 + \vartheta) (1 + \varepsilon)] - b_{ij} \varepsilon / (1 + \varepsilon)$$
(23)

und damit die Abschätzung (12). Zum Nachweis von (11) beachte man

$$0 = \bar{a}_i = -l_i a_1 + a_i + \delta a_i$$
 $(i = 2, ..., n)$

und $l_i = (a_i/a_1) (1 + \varrho_i) \text{ mit } |\varrho_i| \leq \nu$.

3.3. Orthogonale elementare Transformationsmatrizen

Neben den im vorangegangenen Abschnitt behandelten nichtorthogonalen elementaren Transformationsmatrizen stellen orthogonale elementare Transformationsmatrizen (OT-Matrizen) ein weiteres universelles Hilfsmittel der linearen Algebra dar. Sie spielen insbesondere bei Algorithmen zur Lösung von Quadratmittelproblemen und symmetrischen Eigenwertproblemen eine fundamentale Rolle.

Zu den OT-Matrizen zählen speziell Spiegelungs- und Drehungsmatrizen, die beide Gegenstand dieses Abschnitts sind. Zu den OT-Matrizen gehören ferner Permutationsmatrizen und orthogonale Diagonalmatrizen. Erstere sind im Abschnitt 3.1 ausführlich diskutiert worden, letztere sind von der einfachen Gestalt $D = \text{diag}(d_i)$, $|d_i| = 1$, also $d_i = \pm 1$, und können lediglich Vorzeichenänderungen bewirken, so daß sie nicht weiter untersucht zu werden brauchen.

Zur Definition und Charakterisierung von allgemeinen orthogonalen Matrizen sei auf Abschnitt 1.1.H verwiesen.

A. Spiegelungsmatrizen

3.3.1. Definition. Eine Matrix $H \in \mathbb{R}^{n,n}$ der Gestalt

$$\boldsymbol{H} = \boldsymbol{I} - 2\boldsymbol{u}\boldsymbol{u}^{\mathsf{T}}, \quad \boldsymbol{u} \in \mathbf{R}^{n}, \quad \|\boldsymbol{u}\|_{2} = 1$$
⁽¹⁾

heißt Spiegelungsmatrix.

Nach Definition ist eine Spiegelungsmatrix symmetrisch, und wegen $u^{\mathsf{T}}u = 1$ gilt $H^2 = (I - 2uu^{\mathsf{T}}) (I - 2uu^{\mathsf{T}}) = I - 2uu^{\mathsf{T}} - 2uu^{\mathsf{T}} + 4uu^{\mathsf{T}}uu^{\mathsf{T}} = I$, also

$$\boldsymbol{H} = \boldsymbol{H}^{\mathsf{T}} = \boldsymbol{H}^{-1}. \tag{2}$$

Für $x \in \mathbb{R}^n$ läßt sich y = Hx = x - 2z mit $z = u(u^{\mathsf{T}}x)$ schreiben. Da z gerade die Projektion von x in Richtung u ist, kann die Abbildung

$$x \rightarrow y = Hx$$

geometrisch wie folgt gedeutet werden, vgl. Abb. 3.3.1: Das Bild y entsteht aus x durch Spiegelung an der Ursprungsebene mit dem Normalenvektor u. Dies erklärt die Bezeichnung Spiegelungsmatrix. Da solche Matrizen erstmalig von Householder für numerische Zwecke verwendet worden sind, werden sie auch als Householdersche Spiegelungsmatrizen oder kurz als Householder-Spiegelungen bezeichnet, vgl. B 3.2.



Abb. 3.3.1. Interpretation der Abbildung $x \rightarrow y = Hx$ als Spiegelung

Aus (1) liest man ab, daß jedes Vielfache $v = \mu u$ von $u, \mu \in \mathbf{R}, \mu \neq 0$, gemäß

$$\boldsymbol{H} = \boldsymbol{I} - \frac{\boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}}{\gamma} \quad \text{mit} \quad \gamma = \boldsymbol{v}^{\mathsf{T}}\boldsymbol{v}/2 \tag{3}$$

dieselbe Spiegelungsmatrix wie (1) definiert. Für die numerische Praxis ist die Darstellung (3) wegen der größeren Freiheit in der Wahl des Vektors v jedoch günstiger als (1).

Für $\boldsymbol{b} \in \mathbf{R}^n$ gilt

$$\bar{\boldsymbol{b}} = \boldsymbol{H}\boldsymbol{b} = \left(\boldsymbol{I} - \frac{\boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}}{\gamma}\right)\boldsymbol{b} = \boldsymbol{b} - \boldsymbol{v}\beta \quad \text{mit} \quad \beta = \boldsymbol{v}^{\mathsf{T}}\boldsymbol{b}/\gamma.$$
(4)

Wendet man (4) auf die Spalten b^{j} einer Matrix $B \in \mathbb{R}^{n,r}$ an, so ergibt sich für die spaltenweise Berechnung von $\overline{B} = HB$ die Vorschrift

 $\langle ext{neue } j ext{-te Spalte} \rangle := \langle ext{alte } j ext{-te Spalte}
angle - eta_j st \langle ext{Vektor } m{v}
angle$

mit dem für jede Spalte j zu berechnenden Faktor $\beta_j = (\boldsymbol{v}^{\mathsf{T}} \boldsymbol{b}^j) / \gamma$. Die Matrix \boldsymbol{H} braucht also nicht explizit gebildet zu werden.

3.3.2. Algorithmus zur Householder-Transformation. Für $\mathbf{B} = (\mathbf{b}^1, ..., \mathbf{b}^r) \in \Re^{n,r}$ und $\mathbf{v} \in \Re^n, \mathbf{v} \neq \mathbf{o}$, ist $\overline{\mathbf{B}} = (\overline{\mathbf{b}}^1, ..., \overline{\mathbf{b}}^r) = \mathbf{H}\mathbf{B}$ zu berechnen, wobei $\mathbf{H} = \mathbf{I} - \mathbf{v}\mathbf{v}^{\mathsf{T}}/\gamma$ mit $\gamma = \mathbf{v}^{\mathsf{T}}\mathbf{v}/2$ durch \mathbf{v} gegeben ist. S0: $\gamma := \mathbf{v}^{\mathsf{T}}\mathbf{v}/2$ S1: for j := 1(1)r do $\beta := \mathbf{v}^{\mathsf{T}}\mathbf{b}^j/\gamma$ $\overline{\mathbf{b}}^j := \mathbf{b}^j - \mathbf{v} * \beta$

Aufwand: $\sim n(2r+1)$ opms

Die Skalarprodukte sollen dabei in üblicher Weise gemäß 2.3.9 berechnet werden. Algorithmus 3.3.2 kann in situ gemäß $\overline{B} := HB$ ausgeführt werden.

3.3.3. Bemerkung. (i) In den uns interessierenden Anwendungen fällt die Größe $\gamma = \boldsymbol{v}^{\mathsf{T}} \boldsymbol{v}/2$ bei der Festlegung von \boldsymbol{v} ohnehin an, vgl. **3.3.4.** Es ist daher zweckmäßig, die Spiegelung \boldsymbol{H} durch die Daten $\{\boldsymbol{v}, \gamma\} = \{v_1, \ldots, v_n, \gamma\}$ zu repräsentieren. Bei der Bildung von $\boldsymbol{B} = \boldsymbol{B}\boldsymbol{H}$ gemäß **3.3.2** entfällt dann die Neuberechnung von γ im Schritt S0. Wir können dann außerdem \boldsymbol{v} so skalieren, daß $v_1 = \gamma$ gilt, so daß kein zusätzlicher Speicherplatz für γ gebraucht wird, siehe **3.3.5**(ii).

(ii) Algorithmus 3.3.2 kann bei Bedarf zeilenweise in situ realisiert werden, wenn Schritt S1 wie folgt abgeändert wird:

S1': for j := 1(1)r do $\beta_j := (\boldsymbol{v}^{\mathsf{T}}\boldsymbol{b}^j)/\gamma$ for i := 1(1)n do for j := 1(1)r do $\overline{b}_{ij} := b_{ij} - v_i * \beta_j$

Dabei wird ein Hilfsfeld der Länge r zum Speichern der Koeffizienten $\{\beta_1, \ldots, \beta_r\}$ benötigt. \Box

Im folgenden untersuchen wir, ob mittels einer Householder-Spiegelung ähnlich wie mit einer LNT-Matrix in einem Vektor a unterhalb der ersten Komponente Nullen erzeugt werden können, d. h., ob

$$\boldsymbol{H}\boldsymbol{a} = \boldsymbol{\bar{a}} = \boldsymbol{\varrho}\boldsymbol{e}^{1} = (\boldsymbol{\varrho}, 0, ..., 0)^{\mathsf{T}}$$
⁽⁵⁾

durch geeignete Wahl von v erzielt werden kann. Wird auf beiden Seiten von (5 die Euklidische Norm gebildet, so ergibt sich wegen deren Orthogonalinvarianz $\|\boldsymbol{a}\|_2 = \|\boldsymbol{H}\boldsymbol{a}\|_2 = \|\boldsymbol{\varrho}\boldsymbol{e}^1\| = |\boldsymbol{\varrho}|$, also

$$\varrho = \pm \|\boldsymbol{a}\|_2. \tag{6}$$

Analog zu (4) folgt weiter $Ha = a - v(v^{\intercal}a)/\gamma = \varrho e^1$, also $v(v^{\intercal}a)/\gamma = a - \varrho e^1$. Falls *a* nicht schon von der gewünschten Form $a = a_1 e^1$ ist, ergibt sich für jede der beiden Festlegungen (6) ein von *o* verschiedener Vektor $a - \varrho e^1$, so daß $v = \lambda(a - \varrho e^1)$ für jedes $\lambda \neq 0$ die gestellte Aufgabe löst. Der zugehörige γ -Wert ist

$$\begin{split} \gamma &= \lambda^2 (\boldsymbol{a} - \varrho \boldsymbol{e}^1)^{\mathsf{T}} (\boldsymbol{a} - \varrho \boldsymbol{e}^1)/2 = \lambda^2 (\|\boldsymbol{a}\|_2^2 - 2\varrho \boldsymbol{a}_1 + \varrho^2)/2 \\ &= \lambda^2 \varrho (\varrho - \boldsymbol{a}_1) = -\lambda \varrho \boldsymbol{v}_1. \end{split}$$

3.3.4. Aussage. Für $a \in \mathbb{R}^n$, $a \neq a_1 e^1$, und jedes $\lambda \neq 0$ transformiert die durch

$$\boldsymbol{v} = (v_i) = \lambda (\boldsymbol{a} - \varrho \boldsymbol{e}^1), \qquad \varrho = \pm \|\boldsymbol{a}\|_2, \tag{7}$$

festgelegte Householder-Spiegelung

$$\boldsymbol{H} = \boldsymbol{I} - \boldsymbol{v}^{\mathsf{T}} \boldsymbol{v} / \boldsymbol{\gamma}, \qquad \boldsymbol{\gamma} = \boldsymbol{v}^{\mathsf{T}} \boldsymbol{v} / 2 = -\lambda \varrho \boldsymbol{v}_{1}$$
(8)

den Vektor \boldsymbol{a} in

$$\bar{\boldsymbol{a}} = \boldsymbol{H}\boldsymbol{a} = \varrho \boldsymbol{e}^1 = (\varrho, 0, \dots, 0)^\mathsf{T}.$$

Die beiden Möglichkeiten $\rho = \pm ||\boldsymbol{a}||_2$ aus 3.3.4 sind in Abb. 3.3.2 für n = 2 und $\lambda = 1$ mit den zugehörigen Vektoren $\boldsymbol{v}^+ = \boldsymbol{a} - ||\boldsymbol{a}||_2 \boldsymbol{e}^1, \boldsymbol{v}^- = \boldsymbol{a} + ||\boldsymbol{a}||_2 \boldsymbol{e}^1$ dargestellt worden. Man beachte, daß $\boldsymbol{v}^{+\mathsf{T}}\boldsymbol{v}^- = 0$ gilt, d. h., die Normalenvektoren der beiden möglichen spiegelnden Ebenen sind orthogonal.



Abb. 3.3.2. Householder-Transformation 3.3.4 für n = 2 und $\lambda = 1$

3.3.5. Bemerkung. (i) Wird $\rho = \|\boldsymbol{a}\|_2$ gesetzt, so ergibt sich aus (7) und (8) $v_1 = \lambda(a_1 - \|\boldsymbol{a}\|_2).$

Die wörtliche numerische Auswertung von $a_1 - ||\boldsymbol{a}||_2$ ist für $a_1 \leq 0$ stabil. Für $a_1 > 0$ tritt jedoch im Fall $a_1 \approx ||\boldsymbol{a}||_2$ bei der Differenzbildung Auslöschung auf, so daß ein instabiles Verfahren entsteht. Unter Beachtung der Identität

$$a_1 - \|a\|_2 = -\left(\sum_{i=2}^n a_i^2\right) / (a_1 + \|a\|_2)$$

kann dieser negative Effekt vermieden werden, vgl. Ü 2.2.5. Insgesamt liefert die Auswertung gemäß

$$v_{1} = \begin{cases} \lambda(a_{1} - ||\boldsymbol{a}||_{2}) & \text{für } a_{1} \leq 0, \\ -\lambda \left(\sum_{i=2}^{n} a_{i}^{2}\right) / (a_{1} + ||\boldsymbol{a}||_{2}) & \text{für } a_{1} > 0, \end{cases}$$

$$v_{i} = \lambda a_{i} \qquad (i = 2, ..., n)$$
(9)

ein stabiles Verfahren zur Berechnung von v und $\gamma = -\lambda \|a\|_2 v_1$.
(ii) Eine einfachere Stabilisierung kann erreicht werden, wenn das Vorzeichen in (6) so gewählt wird, daß $\gamma = \lambda^2 \varrho(\varrho - a_1) = \lambda^2 (||\boldsymbol{a}||_2^2 - \varrho a_1)$ maximal wird. Dies ist bei der Festlegung

$$\varrho = \begin{cases}
\|a\|_2 & \text{für } a_1 \leq 0, \\
-\|a\|_2 & \text{für } a_1 > 0
\end{cases}$$
(10)

der Fall. Für den Skalierungsfaktor λ bietet sich wegen (8) der Wert $\lambda = -1/\rho$ an, denn dann ist $\gamma = v_1$, vgl. 3.3.3(i). Damit ergibt sich

$$\boldsymbol{v} = (v_i) = \boldsymbol{e}^1 - (1/\varrho) \, \boldsymbol{a}, \tag{11}$$

also

$$egin{aligned} & v_1 = \gamma = 1 - a_1/arrho = 1 + |a_1|/||a||_2, \ & v_i = -a_i/arrho & (i=2,...,n). \end{aligned}$$

Wegen $1 \leq \gamma = \|oldsymbol{v}\|_2^2/2 \leq 2$ gilt dabei

$$\sqrt[]{2} \leq \|m{v}\|_2 \leq 2$$
 ,

d. h., v ist gut skaliert, und bei der Berechnung von $\beta_j := v^{\mathsf{T}} b^j / \gamma$ in 3.3.2 ist weder mit Unter- noch mit Überlauf zu rechnen. Im folgenden soll v stets gemäß (10), (11) festgelegt werden. Dies ist für beliebiges $a \neq o$, also auch im Fall $a = a_1 e^1$, zulässig. \Box

Die typische Anwendung der Householder-Spiegelungen besteht wie die der LNT-Matrizen in der Transformation eines Vektors a in ein Vielfaches von e^1 und der gleichzeitigen Transformation einer Matrix **B**.

3.3.6. Eliminationsschritt mittels Householder-Spiegelung.

Aufgabe: Für $\boldsymbol{a} = (a_i) \in \Re^n$, $\boldsymbol{a} \neq \boldsymbol{o}$, und $\boldsymbol{B} = (\boldsymbol{b}^1, \dots, \boldsymbol{b}^r) \in \Re^{n,r}$ sind $\boldsymbol{v} = (v_i) \in \Re^n$ und $\varrho \in \Re$ so festzulegen, daß $\bar{\boldsymbol{a}} = \boldsymbol{H}\boldsymbol{a} = \varrho \boldsymbol{e}^1 = (\varrho, 0, \dots, 0)^{\mathsf{T}}$ mit $\boldsymbol{H} = \boldsymbol{I} - \boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}/\gamma$ und $\gamma = \boldsymbol{v}^{\mathsf{T}}\boldsymbol{v}/2 = v_1$ gilt. Für das so definierte \boldsymbol{H} ist $\bar{\boldsymbol{B}} = (\bar{\boldsymbol{b}}^1, \dots, \bar{\boldsymbol{b}}^r) = \boldsymbol{H}\boldsymbol{B}$ zu berechnen.

Algorithmus:

Aufwand: $\sim n[(2r+1) \text{ opms} + 1 \text{ opm}] + 1 \text{ opr, wenn } ||a||_2 \text{ mit } \sim n \text{ opms} + 1 \text{ opr berechnet wird.}$

Der Algorithmus kann in situ durch Überspeichern von a mit v und B mit \overline{B} durchgeführt werden. Das neu berechnete Element $\overline{a}_1 = \rho$ muß dann auf einem ge-

sonderten Platz aufgehoben werden, während die n-1 erzeugten Nullen $\bar{a}_2, \ldots, \bar{a}_n$ selbstverständlich nicht gespeichert werden. Bei dieser Realisierung geht

$$\begin{pmatrix} a_1 & | & | & | \\ a_2 & | & | \\ \vdots & \mathbf{b}^1 \mathbf{b}^2 \dots \mathbf{b}^r \\ a_n & | & | & | \end{pmatrix} \quad \text{in} \quad \begin{pmatrix} v_1 & | & | & | \\ v_2 & | & | \\ \vdots & | \mathbf{\bar{b}}^1 \mathbf{\bar{b}}^2 \dots \mathbf{\bar{b}}^r \\ v_n & | & | & | \end{pmatrix}$$

über. Falls a schlecht skaliert ist, sollte bei der Berechnung von $||a||_2$ Über- oder Unterlauf vermieden werden, siehe 2.3.E.

Die Rundungsfehleranalyse führt auf die folgenden Ergebnisse:

3.3.7. Rundungsfehleranalyse. Algorithmus 3.3.6 werde für $\boldsymbol{a} \in \mathbb{R}^n$, $\boldsymbol{a} \neq \boldsymbol{o}$, und $\boldsymbol{B} \in \mathbb{R}^{n,r}$ durchgeführt, und $\bar{\boldsymbol{a}} = (\varrho, 0, ..., 0)^{\mathsf{T}} \in \mathbb{R}^n$ sowie $\bar{\boldsymbol{B}} \in \mathbb{R}^{n,r}$ seien die berechneten Werte. Mit \boldsymbol{H}^* werde die Spiegelungsmatrix bezeichnet, die bei exakter Ausführung von 3.3.6 durch \boldsymbol{a} festgelegt wird. Dann gilt

$$\bar{\boldsymbol{a}} = (\varrho, 0, \dots, 0)^{\mathsf{T}} = \boldsymbol{H}^*(\boldsymbol{a} + \boldsymbol{\delta}\boldsymbol{a}) \tag{12}$$

sowie

$$\overline{B} = H^*(B + \delta B) = H^*B + \delta \overline{B}$$
(13)

mit Störungen δa , δB und $\delta \overline{B}$, die den Abschätzungen

 $|\delta a| \stackrel{<}{=} \nu[(n+2)/2] |a| \tag{14}$

und

$$\|\boldsymbol{\delta}\boldsymbol{b}^{j}\|_{2} \leq vF \|\boldsymbol{b}^{j}\|_{2}, \qquad \|\boldsymbol{\delta}\bar{\boldsymbol{b}}^{j}\|_{2} \leq vF \|\boldsymbol{\delta}\bar{\boldsymbol{b}}^{j}\|_{2}$$
(15)

 mit

$$F := F(n) := K(n)n$$
,

$$K(n) := 3.14(1 + 3/n), \text{ also } F(n) \le 4 \text{ für } n \ge 10$$
 (16)

genügen, wobei b^{j} bzw. \overline{b}^{j} die Spalten von **B** bzw. \overline{B} bezeichnen.

Aus den spaltenweisen Abschätzungen (15) ergibt sich sofort

~

$$\|\boldsymbol{\delta}\boldsymbol{B}\|_{F} \leq vF \|\boldsymbol{B}\|_{F} \quad \text{bzw.} \quad \|\boldsymbol{\delta}\boldsymbol{\bar{B}}\|_{F} \leq vF \|\boldsymbol{\bar{B}}\|_{F}$$
(17)

als kollektive Abschätzung in der Frobeniusnorm.

Die Ergebnisse von 3.3.7 besagen: Die Ausführung eines Eliminationsschrittes mittels Householder-Spiegelungen gemäß 3.3.6 ist ein numerisch gutartiger Prozeß mit F = 4n für $n \ge 10$ in der Euklidischen bzw. Frobenius-Norm, und die Störungen sind spaltenweise klein.

A. Drehungsmatrizen

Die zweite Klasse elementarer orthogonaler Matrizen besteht aus Matrizen des Typs

$$egin{aligned} G_{pq} &= egin{pmatrix} p & q \ 1 & dots & dots \ \dots & c & \dots & s & \dots \ dots & dots & dots & dots \ dots & dots & dots & dots \ dots & dots & dots & dots \ dots \ dots & dots \ dot$$

111

3.3.8. Definition. Eine Matrix $G_{pq} \in \mathbb{R}^{n,n}$ $(1 \le p < q \le n)$, die sich von der Einheitsmatrix nur in den vier Elementen

$$(G_{pq})_{pp} = (G_{pq})_{qq} = c, \qquad (G_{pq})_{pq} = -(G_{pq})_{qp} = s$$

unterscheidet und für die

$$c^2 + s^2 = 1 \tag{18}$$

gilt, heißt Drehungsmatrix. Im Fall p = q wird $G_{pp} = G_{qq} = I$ gesetzt.

'Aus (18) folgt in der Tat $GG^{\intercal} = I$, also die Orthogonalität von G. Die Transformation

$$x \rightarrow y = G_{pq}x$$

lautet komponentenweise

y q

$$y_{p} = c * x_{p} + s * x_{q},$$

$$y_{q} = -s * x_{p} + c * x_{q},$$

$$y_{i} = x_{i} \qquad (i \neq p, q; i = 1, ..., n).$$

$$e^{q} \downarrow$$

$$(19)$$



Sie verändert nur die Komponenten p und q von x und läßt sich als Drehung um den Winkel $-\varphi$ in der durch e^p und e^q aufgespannten Ebene deuten, wobei φ durch die Bedingungen

$$c = \cos \varphi$$
, $s = \sin \varphi$

festgelegt ist, siehe Abb. 3.3.3. Man nennt G_{pq} deshalb auch *ebene Drehung* oder – nach GIVENS, der diese Matrizen für die numerische lineare Algebra wiederentdeckt hat, vgl. B 3.3 – *Givens-Drehung*.

Ist $B \in \mathbf{R}^{n,r}$ eine Matrix mit den Zeilen $b^{i\tau}$, so ergeben sich die Zeilen \overline{b}_i^{τ} von $\overline{B} = G_{pq}B$ analog zu (19) gemäß

$$\langle \text{neue } p\text{-te Zeile} \rangle := c * \langle \text{alte } p\text{-te Zeile} \rangle + s * \langle \text{alte } q\text{-te Zeile} \rangle,$$

$$\langle \text{neue } q\text{-te Zeile} \rangle := -s * \langle \text{alte } p\text{-te Zeile} \rangle + c * \langle \text{alte } q\text{-te Zeile} \rangle;$$

die übrigen Zeilen von \overline{B} werden unverändert von B übernommen.

Wenn diese Vorschrift auch in situ realisiert werden soll, muß eines der neuen Elemente zwischengespeichert werden wie im folgenden Algorithmus:

3.3.9. Algorithmus zur Givens-Transformation. Bei gegebener Givens-Drehung $G_{pq} \in \Re^{n,n}$ mit p < q und den Drehungsparametern $c, s \in \Re$ ist $B = (b_{ij}) \in \Re^{n,r}$ mit $\overline{B} = G_{pq}B$ zu überspeichern.

for
$$j := 1(1)r$$
 do

$$\begin{vmatrix}
\mu := c * b_{pj} + s * b_{qj} \\
b_{qj} := -s * b_{pj} + c * b_{qj} \\
b_{pj} := \mu
\end{cases}$$
Autwand: 4r opm + 2r ops

Wie die bisher eingeführten Transformationsmatrizen werden auch Givens-Drehungen verwendet, um durch geeignete Festlegung von c und s in einem Vektor $\boldsymbol{a} = (a_i)$ an vorgegebenen Stellen Nullen zu erzeugen. Da beim Übergang von \boldsymbol{a} zu $\bar{a} = G_{pq}a$ nur die Komponenten p und q verändert werden, kann nur eines dieser Elemente zu Null gemacht werden, etwa

$$\bar{a}_q = -s * a_p + c * a_q = 0.$$
 (20)

Wegen der Orthogonalität von G_{pq} muß dann $|\bar{a}_p| = \sqrt{a_p^2 + a_q^2}$ gelten. Die Bedingung (20) läßt sich wie folgt erfüllen:

3.3.10. Aussage. Es seien p, q mit $1 \leq p < q \leq n$ und $a = (a_i) \in \mathbf{R}^n$ mit **3.3.10.** Aussage. Its set p, q mit $1 \ge p < q \ge n$ und $a = (a a_p^2 + a_q^2 > 0$ gegeben. Dann transformiert die durch $c = a_p/\varrho, \quad s = a_q/\varrho \quad \text{mit} \quad \varrho = \pm \sqrt{a_p^2 + a_q^2}$ festgelegte Givens-Drehung G_{pq} den Vektor a in $\bar{a} = (\bar{a}_i) = G_{pq}a$ mit $\bar{a}_p = \varrho, \quad \bar{a}_q = 0 \quad \text{und} \quad \bar{a}_i = a_i \quad (i \neq p, q).$

$$c=a_p/arrho, \quad s=a_q/arrho \quad ext{mit} \quad arrho=\pm \sqrt{a_p^2+a_q^2}$$

Dies läßt sich durch Einsetzen sofort überprüfen.

Für die numerische Realisierung sind die Formeln (21) wegen des möglichen Überbzw. Unterlaufs bei der Berechnung von $a_p^2 + a_q^2$ nicht zu empfehlen, vgl. Abschnitt 2.3.E. Eine zuverlässige Realisierung von 3.3.10 gemäß

$$arrho = egin{cases} a_p \, \sqrt{1 \, + \, (a_q/a_p)^2} & ext{in Fall} & |a_p| \geq |a_q| \, , \ a_q \, \sqrt{1 \, + \, (a_p/a_q)^2} & ext{in Fall} & |a_p| < |a_q| \end{cases}$$

wird durch das nachfolgende Verfahren gegeben:

3.3.11. Algorithmus. Berechnung der Drehungsparameter c, s und der Größen ϱ, ξ aus $a_p, a_q \in \mathfrak{R}$.

$$\begin{array}{l|l} \text{if } |a_p| \geq |a_q| \text{ then } \\ \left| \begin{array}{l} \text{if } a_p \neq 0 \\ \text{then } \\ \left| \begin{array}{l} \xi := a_q/a_p, \ \varrho := \text{sqrt} \ (1 + \xi^2), \ c := 1/\varrho \\ s := \xi := \xi/\varrho, \ \varrho := a_p \ast \varrho \\ \text{else } \left[c := 1, \ s := \varrho := \xi := 0 \right] \end{array} \right| \end{array}$$

Schwetlick, Numerische Algebra

else
$$\left| \begin{array}{l} \xi := a_p/a_q, \, \varrho := \mathrm{sqrt} \, (1 + \xi^2), \, s := 1/\varrho \\ c := \xi/\varrho, \, \varrho := a_q * \varrho \\ \mathrm{if} \, c \neq 0 \, \mathrm{then} \, \xi := 1/c \, \mathrm{else} \, \xi := 1 \end{array} \right|$$

Autward: 6 opm + 1 ops + 1 opr

Der Algorithmus arbeitet auch im Fall $a_p = a_q = 0$, wo c = 1, s = 0, also $G_{pq} = I$ gesetzt wird. Er stellt neben c, s, ϱ außerdem die Zahl ξ bereit, aus der c, s für spätere Anwendungen von G_{pq} einfach rekonstruiert werden kann. Es genügt daher, statt der beiden Zahlen c, s nur die eine Zahl ξ z. B. auf dem Platz von $\bar{a}_q = 0$ zu speichern, siehe B 3.3.

3.3.12. Algorithmus. Rekonstruktion von c, s aus $\xi \in \Re$.

 $\begin{array}{l|l} \text{if } |\xi| \leq 1 \text{ then } & \left| \begin{array}{l} \text{if } \xi \neq 1 \\ \text{then } c := \operatorname{sqrt} \left(1 - \xi^2 \right), s := \xi \\ \text{else } c := 0, s := 1 \\ \text{else } \left[c := 1/\xi, s := \operatorname{sqrt} \left(1 - c^2 \right) \right] \end{array} \right. \end{array}$

Aufwand:
$$2 \text{ opm} + 1 \text{ ops} + 1 \text{ opr}$$

Man beachte dabei, daß für das vom Algorithmus 3.3.11 bereitgestellte ξ genau eine der Bedingungen $|\xi| < 1/\sqrt{2}$, $|\xi| > \sqrt{2}$ oder $\xi = 1$ erfüllt ist, so daß weder bei der Berechnung von $1 - \xi^2$ noch bei der von $1 - c^2$ Auslöschung eintreten kann.

Soll ein Vektor $a = a^1 \in \mathbb{R}^n$ analog zu 3.3.4 mittels Givens-Drehungen in den Vektor $\bar{a} = (\varrho, 0, ..., 0)^{\intercal}$ transformiert werden, so kann dies schrittweise gemäß

$$a^1 \to a^2 = G_{12}a^1 \to a^3 = G_{13}a^2 \to \dots \to \bar{a} = a^n = G_{1n}a^{n-1}$$
 (22)

geschehen, wobei G_{1i} in der *i*-ten Komponente von a^{i-1} eine Null erzeugt. Dabei bleiben die vorher bereits erzeugten Nullen offensichtlich erhalten.

3.3.13. Aussage. Zu jedem Vektor $a \in \mathbb{R}^n$ gibt es eine Folge $G_{12}, G_{13}, \ldots, G_{1n}$ von Givens-Drehungen, so daß

$$G = G_{1n}G_{1,n-1} \cdots G_{12}$$
(23)
den Vektor \boldsymbol{a} in
$$\bar{\boldsymbol{a}} = G\boldsymbol{a} = \varrho \boldsymbol{e}^1 = (\varrho, 0, ..., 0)^{\mathsf{T}} \text{ mit } \varrho = \pm \|\boldsymbol{a}\|_2$$
transformiert.

Nachfolgend geben wir eine in-situ-Realisierung von 3.3.13 an. Wie in den vorangegangenen Abschnitten wird zusätzlich $\overline{B} = GB$ für gegebenes $B \in \Re^{n,r}$ berechnet, und zwar analog zu (22) gemäß

$$B^{(1)} = B, \qquad B^{(i)} = G_{1i}B^{(i-1)} \qquad (i = 2, ..., n), \qquad \overline{B} = B^{(n)}.$$
 (24)

3.3.14. Eliminationsschritt mittels einer Folge von Givens-Drehungen.

Aufgabe: Für $a = (a_i) \in \Re^n$ und $B = (b_{ij}) \in \Re^{n,r}$ sind die Drehungen $G_{1i} \in \Re^{n,n}$ (i = 2, ..., n) so festzulegen, daß $\overline{a} = G_{1n} \cdots G_{12}a = \varrho e^1$ gilt. Der Vektor a ist mit $(\varrho, \xi_2, ..., \xi_n)^{\mathsf{T}}$ zu überspeichern, wobei ξ_i die der Drehung G_{1i} gemäß 3.3.11

zugeordnete Zahl bezeichne. Die Matrix B ist durch $\overline{B} = G_{1n} \cdots G_{12}B$ zu überspeichern.

Algorithmus:

for i := 2(1)n do S0 (Festlegung von G_{1i}): Berechne $c := c(a_1, a_i), s := s(a_1, a_i), \varrho := \varrho(a_1, a_i), \xi = \xi(a_1, a_i)$ gemäß Algorithmus 3.3.11 mit $a_p = a_1$ und $a_q = a_i$ S1 (Transformation $a^{i-1} \rightarrow a^i := G_{1i}a^{i-1}$, Reservierung von $\xi = \xi_i$): $a_1 := \varrho, a_i := \xi$ S2 (Transformation $B^{(i-1)} \rightarrow B^{(i)} := G_{1i}B^{(i-1)}$): for j := 1(1)r do $\mu := c * b_{1j} + s * b_{ij}$ $b_{ij} := -s * b_{1j} + c * b_{ij}$ $b_{1j} := \mu$ Aufwand: Festlegung der Drehungsparameter und Berechnung von \bar{a} : (n-1) (6 opm + 1 ops + 1 opr) Berechnung von \bar{B} : (n-1)r (4 opm + 2 ops)

Durch 3.3.14 wird

$$\begin{pmatrix} a_1 & b_{11} \dots & b_{1r} \\ a_2 & b_{21} \dots & b_{2r} \\ \vdots & \vdots & \vdots \\ a_n & b_{n1} \dots & b_{nr} \end{pmatrix} \quad \text{mit} \quad \begin{pmatrix} \varrho & \overline{b}_{11} \dots & \overline{b}_{1r} \\ \xi_2 & \overline{b}_{21} \dots & \overline{b}_{2r} \\ \vdots & \vdots & \vdots \\ \xi_n & \overline{b}_{n1} \dots & \overline{b}_{nr} \end{pmatrix}$$

überspeichert. Man beachte dabei, daß die Produktmatrix G in 3.3.14 nicht explizit gebildet wird. Aus den Zahlen ξ_i (i = 2, ..., n) können die Drehungen G_{1i} bei Bedarf mittels Algorithmus 3.3.12 rekonstruiert werden.

3.3.15. Rundungsfehleranalyse. Für $a \in \Re^n$ und $B \in \Re^{n,r}$ werde Algorithmus 3.3.14 durchgeführt, und $\bar{a} = (\varrho, 0, ..., 0)^{\mathsf{T}} \in \Re^n$ sowie $\bar{B} \in \Re^{n,r}$ seien die berechneten Resultate. Mit $G_{1i}^* \in \mathbb{R}^{n,n}$ werde die exakte Givens-Drehung bezeichnet, die den berechneten Vektor $a^{i-1} \in \Re^n$ entsprechend 3.3.11 exakt in $G_{1i}a^{i-1}$ mit $(G_{1i}a^{i-1})_i = 0$ transformiert, und es sei $G^* = G_{1n}^* \cdots G_{12}^*$ das exakte Produkt der Drehungen G_{1i}^* . Dann gelten die Beziehungen

$$\bar{\boldsymbol{a}} = (\varrho, 0, \dots, 0)^{\mathsf{T}} = \boldsymbol{G}^*(\boldsymbol{a} + \boldsymbol{\delta}\boldsymbol{a}) \tag{25}$$

und

$$\overline{B} = G^*(B + \delta B) = G^*B + \delta \overline{B}.$$
(26)

Die Störungen δa , δB und $\delta \overline{B}$ genügen den Abschätzungen

$$\|\mathbf{d}\mathbf{a}\|_{2} \leq v(2n+2) \|\mathbf{a}\|_{2}, \tag{27}$$

$$\|\boldsymbol{\delta b^{j}}\|_{2} \leq \nu[6(n-1)] \|\boldsymbol{b^{j}}\|_{2}$$
(28)

sowie

$$\|\boldsymbol{d}\bar{\boldsymbol{b}}^{j}\|_{2} \leq \nu[6(n-1)] \|\bar{\boldsymbol{b}}^{j}\|_{2}, \qquad (29)$$

wobei b^{j} die *j*-te Spalte von **B** bezeichnet usw.

Aus (28), (29) folgen die analogen Abschätzungen

$$\|\boldsymbol{\delta}\boldsymbol{B}\|_{F} \leq \nu[6(n-1)] \|\boldsymbol{B}\|_{F} \quad \text{bzw.} \quad \|\boldsymbol{\delta}\boldsymbol{\bar{B}}\|_{F} \leq \nu[6(n-1)] \|\boldsymbol{\bar{B}}\|_{F} \tag{30}$$

in der Frobeniusnorm.

Die Analyse 3.3.15 besagt: Die Realisierung eines Eliminationsschrittes mittels einer Folge von Givens-Drehungen gemäß 3.3.14 ist ein numerisch gutartiger Algorithmus mit F = 2n + 2 bezüglich a und mit F = 6(n - 1) bezüglich der Spalten von Bin der Euklidischen Norm.

3.3.16. Bemerkungen. (i) Ein Vergleich der Aufwandskenngrößen von 3.3.6 und 3.3.14 zeigt, daß die Transformation $(a, B) \rightarrow (\bar{a}, \bar{B})$ mit $\bar{a} = \varrho e^1$ mittels n - 1 Givens-Drehungen etwa doppelt soviel Multiplikationen und *n*-mal soviel Wurzelberechnungen erfordert wie die mittels einer einzelnen Householder-Spiegelung. Im Normalfall, d. h. für voll besetztes a und B, ist daher die Elimination mittels einer Spiegelung vorzuziehen.

(ii) Wenn der steuernde Vektor a und die Matrix B schwach besetzt und ohne spezielle Struktur sind, läßt sich die Elimination mittels Givens-Drehungen wegen deren größerer Flexibilität trotz des höheren Operationszählers meist effektiver realisieren. Man beachte dabei, daß bei schwach besetzten Problemen der Organisationsaufwand für die Berücksichtigung der Besetztheitsstruktur einen erheblichen Anteil des Gesamtaufwandes ausmacht.

(iii) Wie im folgenden gezeigt wird, läßt sich der Aufwand für die Ausführung von Givens-Drehungen durch eine geschickte implizite Realisierung etwa halbieren, so daß sie von der Anzahl der Operationen her mit den Spiegelungen konkurrieren können. Bei Anwendung zur orthogonalen Dreiecksfaktorisierung sind die Schranken für die Rundungsfehler sogar deutlich günstiger als bei Verwendung von Spiegelungen, vgl. Abschnitt 10.2.

C. Implizite Givens-Drehungen.

Für $a = (a_i) \in \mathbb{R}^n$, $B = (b_{ij}) \in \mathbb{R}^{n,r}$ betrachten wir wieder die Transformation

$$a \to \bar{a} = G_{pq}a, \qquad B \to \bar{B} = G_{pq}B.$$
 (31)

Dabei soll die Givens-Drehung G_{pg} so festgelegt sein, daß

$$\vec{a}_q = -s * a_p + c * a_q = 0 \tag{32}$$

gilt, vgl. 3.3.10. Die oben bereits angekündigte implizite Realisierung von (31), (32) beruht auf einer skalierten Darstellung der zu transformierenden Größen gemäß

$$\begin{aligned} \boldsymbol{a} &= \operatorname{diag}\left(\sqrt{\varkappa_{i}}\right) \boldsymbol{d}, \qquad \boldsymbol{B} &= \operatorname{diag}\left(\sqrt{\varkappa_{i}}\right) \boldsymbol{E}, \\ \boldsymbol{\bar{a}} &= \operatorname{diag}\left(\sqrt{\varkappa_{i}}\right) \boldsymbol{\bar{d}}, \qquad \boldsymbol{\bar{B}} &= \operatorname{diag}\left(\sqrt{\varkappa_{i}}\right) \boldsymbol{\bar{E}} \end{aligned} \tag{33}$$

mit Skalierungsfaktoren $\varkappa_i, \, \overline{\varkappa}_i > 0 \, (i = 1, ..., n)$, die wir in den Vektoren

$$\mathbf{z} = (\mathbf{z}_i), \quad \mathbf{\bar{z}} = (\mathbf{\bar{z}}_i) \in \mathbf{R}^n$$

zusammenfassen. Wird G_{pq} ebenfalls in der skalierten Form

$$G_{pq} = \operatorname{diag}\left(\sqrt{\overline{\varkappa}_i}\right) F_{pq} \operatorname{diag}\left(1/\sqrt{\varkappa_i}\right)$$
 (34)

mit der Matrix

$$m{F}_{pq} = egin{pmatrix} p & q \ 1 & dots & dots \ \dots & dots & \dots \ dots & dots & dots \ dots & dots \ dots & dots \ dot$$

geschrieben, so läßt sich die Transformation (31) äquivalent als

$$\boldsymbol{d} \rightarrow \boldsymbol{\bar{d}} = \boldsymbol{F}_{pq} \boldsymbol{d}, \qquad \boldsymbol{E} \rightarrow \boldsymbol{\bar{E}} = \boldsymbol{F}_{pq} \boldsymbol{E}$$
 (36)

schreiben, und (32) geht in

$$\bar{d}_q = -\gamma * d_p + \delta * d_q = 0 \tag{36}$$

über. Mit den skalierten Größen lautet die Aufgabe wie folgt: Für gegebene Eingangsdaten $\{\mathbf{z}, \mathbf{d}, \mathbf{E}\}$ sind die Ausgangsdaten $\{\overline{\mathbf{z}}, \overline{\mathbf{d}}, \overline{\mathbf{E}}\}$ und die Parameter $\alpha, \beta, \gamma, \delta$ von \mathbf{F}_{pq} so zu bestimmen, daß die durch (34) festgelegte Matrix \mathbf{G}_{pq} eine Drehung ist und $\overline{d}_q = 0$ gilt.

Die Gleichung (34) bedeutet elementweise $\sqrt{\overline{\varkappa}_i/\varkappa_i} = 1$, also

$$\bar{\varkappa}_i = \varkappa_i \qquad (i \neq p, q; i = 1, ..., n) \tag{37}$$

sowie

$$c = \sqrt{\frac{\overline{\varkappa}_p}{\varkappa_p}} \alpha = \sqrt{\frac{\overline{\varkappa}_q}{\varkappa_q}} \delta, \quad s = \sqrt{\frac{\overline{\varkappa}_p}{\varkappa_q}} \beta = \sqrt{\frac{\overline{\varkappa}_q}{\varkappa_p}} \gamma.$$
(38)

Dabei muß

$$c^2 + s^2 = 1 \tag{39}$$

gelten, damit G_{pq} eine Drehung ist. Die noch frei wählbaren sechs Größen α , β , γ , δ und $\bar{\varkappa}_p$, $\bar{\varkappa}_q$ brauchen also nur vier Bedingungsgleichungen (36), (38), (39) zu erfüllen. Wir versuchen daher, zwei der Parameter α , β , γ , δ so vorzugeben, daß die Berechnung von $\bar{E} = F_{pq}E$ möglichst einfach wird. Da c und s i. allg. von 0 verschieden sind, scheiden Nullen als mögliche Werte wegen (38) aus. Es bietet sich daher an, zwei der Parameter als Einsen zu wählen, und zwar $\alpha = \delta = 1$ im Fall $c \neq 0$ bzw. $\beta = \gamma = 1$ im Fall $s \neq 0$. Die nichttrivialen Elemente von F_{pq} werden dann zu

$$\begin{pmatrix} \mathbf{1} & \beta \\ -\gamma & \mathbf{1} \end{pmatrix}$$
 bzw. $\begin{pmatrix} \alpha & \mathbf{1} \\ -\mathbf{1} & \delta \end{pmatrix}$,

so daß die Berechnung von $F_{pq}E$ nur noch 2r statt der sonst nötigen 4r Multiplikationen erfordert.

Wir kommen jetzt zur Festlegung der Parameter und setzen zunächst $d_p^2 + d_q^2 > 0$ voraus; wegen 3.3.10 und (33) gilt dann

$$c^{2} = \frac{a_{p}^{2}}{a_{p}^{2} + a_{q}^{2}} = \frac{\varkappa_{p} d_{p}^{2}}{\varkappa_{p} d_{p}^{2} + \varkappa_{q} d_{q}^{2}}, \quad s^{2} = \frac{a_{q}^{2}}{a_{p}^{2} + a_{q}^{2}} = \frac{\varkappa_{q} d_{q}^{2}}{\varkappa_{p} d_{p}^{2} + \varkappa_{q} d_{q}^{2}}.$$
 (40)

Fall I: $c \neq 0$, d. h. $d_p \neq 0$.

Wir setzen hier

$$\alpha = \delta = 1, \tag{41}$$

und aus (36) folgt sofort

 $\gamma = d_q/d_p. \tag{42}$

Wegen (41) liefert die erste der Gleichungen (38) die Bedingung $\bar{\varkappa}_p | \varkappa_p = \bar{\varkappa}_q | \varkappa_q$. Mit dieser und (42) folgt aus der zweiten Gleichung (38)

$$\beta = (\varkappa_q / \varkappa_p) \, \gamma = (\bar{\varkappa}_q / \bar{\varkappa}_p) \, \gamma. \tag{43}$$

Aus der ersten Gleichung (38) liest man weiter $\bar{\varkappa}_i = \varkappa_i c^2$, i = p, q, wegen

$$rac{1}{c^2} = rac{c^2 + s^2}{c^2} = 1 + rac{s^2}{c^2} = 1 + rac{arkappa_q d_q^2}{arkappa_p d_p^2} = 1 + eta \gamma$$

also

$$\bar{\varkappa}_p = \varkappa_p / \tau, \quad \bar{\varkappa}_q = \varkappa_q / \tau \quad \text{mit} \quad \tau = 1 + \beta \gamma$$
(44)

ab. Das neue Element \overline{d}_p ergibt sich schließlich zu

$$\bar{d}_p = d_p + \beta d_q = d_p + \beta \gamma d_p = d_p \tau.$$
(45)

Fall II: $s \neq 0$, d. h. $d_q \neq 0$.

Dann kann

 $\beta = \gamma = 1 \tag{46}$

gesetzt werden, und analog zum Fall I ergeben sich

$$\delta = d_p/d_q, \qquad \alpha = (\varkappa_p/\varkappa_q) \,\delta = (\bar{\varkappa}_q/\bar{\varkappa}_p) \,\delta,$$
(47)

$$\bar{\varkappa}_p = \varkappa_q / \tau, \qquad \bar{\varkappa}_q = \varkappa_p / \tau \quad \text{mit} \quad \tau = 1 + \alpha \delta$$
 (48)

sowie

$$\bar{d}_p = d_q \tau$$
. (49)

3.3.17. Aussage. Gegeben seien $d = (d_i) \in \mathbb{R}^n$ mit $d_p^2 + d_q^2 > 0$, $p \neq q$, und die Skalierungsfaktoren $z_i > 0$ (i = 1, ..., n). Dann können die Parameter $\alpha, \beta, \gamma, \delta$ von F_{pq} und die Skalierungsfaktoren $\bar{z}_i > 0$ (i = 1, ..., n) nach den Formeln (37), (41) bis (49) so festgelegt werden, daß

entweder $\alpha = \delta = 1$ oder $\beta = \gamma = 1$

gilt und die durch (34) definierte Matrix G_{pq} eine Givens-Drehung ist, die

$$oldsymbol{a} = ext{diag}\left(\sqrt{arkappa_i}
ight) oldsymbol{d} \quad ext{in} \quad oldsymbol{ar{a}} = ext{diag}\left(\sqrt{arkappa_i}
ight) oldsymbol{ar{d}} = oldsymbol{G}_{pq}oldsymbol{a}$$

transformiert, wobei $\bar{d} = F_{po}d$ mit

$$ar{d}_p=\sigma, \quad ar{d}_q=0 \quad ext{und} \quad ar{d}_i=d_i \quad (i \neq p,q)$$

$$e_{ij} = e_{ij}$$
 $(j = 1, ..., r; i \neq p, q; i = 1, ..., n)$

gilt. Für die analoge Transformation von $B = \operatorname{diag} \left(\sqrt{z_i} \right) E$ in $\overline{B} = \operatorname{diag} \left(\sqrt{\overline{z_i}} \right) \overline{E} = G_{pq} E$ mit $E = (e_{ij}) \in \mathbb{R}^{n,r}$ gilt $\overline{E} = F_{pq} E$, also $\overline{e}_{ij} = e_{ij}$ $(j = 1, ..., r; i \neq p, q; i = 1, ..., n)$ und $\overline{e}_{pj} = e_{pj} + \beta e_{qj}, \overline{e}_{qj} = -\gamma e_{pj} + e_{qj} (j = 1, ..., r)$ im Fall $\alpha = \delta = 1$ bzw. $\overline{z} = \overline{z} z_{ij} + \beta z_{ij} = \overline{z} = z_{ij} + \delta z_{ij} (i = 1, ..., r)$ im Fall $\beta = \alpha = 1$ (50)

$$k_{pj} = \alpha e_{pj} + e_{qj}, \bar{e}_{qj} = -e_{pj} + \delta e_{qj} \ (j = 1, ..., r) \text{ im Fall } \beta = \gamma = 1.$$
 (51)

Wir nennen die Transformation $\bar{d} = F_{no}d$ einschließlich der zugehörigen Transformation von \varkappa in $\bar{\varkappa}$ implizite Realisierung der durch G_{pq} gemäß 3.3.10 festgelegten Givens-Transformation und sprechen kurz von impliziten Givens-Drehungen. Gelegentlich werden diese Transformationen auch modifizierte Givens-Drehungen bzw. da bei der Realisierung im Gegensatz zu 3.3.10 keine Quadratwurzeln zu ziehen sind – quadratwurzelfreie Givens-Drehungen genannt; siehe B 3.4 für einen historischen Kommentar.

Um die Parameter von F_{pq} und die Skalierungsfaktoren in vernünftigen Größenordnungen zu halten, gehen wir wie in 3.3.11 vor und verwenden die Formeln aus $\text{Fall I für } |a_p| \geq |a_q|, \text{ d. h. } \varkappa_p d_p^2 > \varkappa_q d_q^2; \text{ die aus Fall II werden für } |a_p| < |a_q| \text{ ver-}$ wendet. Den Ausnahmefall $d_p = d_q = 0$ ordnen wir wie dort dem Fall I zu und setzen einfach $F_{pq} = G_{pq} = I$, $\bar{\varkappa} = \varkappa$. In Analogie zu 3.3.11 berechnen wir außerdem eine Größe ξ , aus der sowohl das Vorliegen von Fall I oder II wie auch die Werte der Parameter einfach rekonstruiert werden können.

3.3.18. Algorithmus. Berechnung der Parameter α , β , γ , δ von F_{pq} , der Skalierungsfaktoren $\bar{\varkappa}_p$, $\bar{\varkappa}_q$ und der Größen $\sigma = \bar{d}_p$ sowie ξ aus d_p , d_q , \varkappa_p , $\varkappa_q \in \Re$.

$$\begin{split} \mathbf{f} \, \varkappa_p d_p^2 &\geq \varkappa_q d_p^2 \, \mathbf{then} \\ \left| \begin{array}{l} \alpha := \delta := 1 \\ \mathbf{if} \, d_p \neq 0 \\ \mathbf{then} \\ \gamma := d_q | d_p, \, \xi := \varkappa_p | \varkappa_q, \, \beta := \gamma | \xi \\ \xi := \gamma / (1 + \xi) \\ \mathbf{else} \quad [\gamma := \beta := \xi := 0] \\ \tau := 1 + \gamma * \beta, \, \overline{\varkappa}_p := \varkappa_p / \tau, \, \overline{\varkappa}_q := \varkappa_q / \tau, \, \sigma := d_p * \tau \\ \beta := \gamma := 1 \\ \delta := d_p | d_q, \, \xi := \varkappa_p | \varkappa_q, \, \alpha := \delta * \xi \\ \mathbf{if} \, \alpha \neq 0 \, \mathbf{then} \, \xi := (1 + \xi) | \alpha \, \mathbf{else} \, \xi := 1 \\ \tau := 1 + \alpha * \delta, \, \overline{\varkappa}_p := \varkappa_q / \tau, \, \overline{\kappa}_q := \varkappa_p / \tau, \, \sigma := d_q * \tau \end{split}$$

Autwand: 12 opm + 2 ops (einschließlich der zur Fallunterscheidung nötigen Multiplikationen)

Man beachte dabei, daß die drei Fälle

$$\kappa_p d_p^2 \ge \kappa_q d_q^2, \qquad 0 = \kappa_p d_p^2 < \kappa_q d_q^2 \quad ext{bzw.} \quad 0 < \kappa_p d_p^2 < \kappa_q d_q^2$$

119

durch

$$|\xi| \leq 1/2, \qquad \xi=1 \quad ext{bzw.} \quad |\xi|>2$$

charakterisiert sind und daher zuverlässig aus ξ rekonstruiert werden können. Für den Fall $\varkappa_p d_p^2 \ge \varkappa_q d_q^2$, also $|\xi| \le 1/2$, verwenden wir im folgenden auch die Bezeichnung "Fall 1"; mit "Fall 2" werden die beiden restlichen Möglichkeiten bezeichnet.

Die in 3.3.18 vorgenommene Fallunterscheidung garantiert

$$|\gamma| \leq \sqrt{\varkappa_p/\varkappa_q}, |\beta| \leq \sqrt{\varkappa_q/\varkappa_p} \quad \text{im Fall 1}$$

(52)

bzw.

$$|lpha| \leq \sqrt{arkappa_p/arkappa_q}, \, |\delta| \leq \sqrt{arkappa_q/arkappa_p} \quad ext{in Fall 2,}$$

also

$$1 \leq \tau = \alpha \delta + \beta \gamma \leq 2.$$
 (53)

Die Rekonstruktion der Kenngrößen von F_{pq} ist analog zu 3.3.12 nach den folgenden Algorithmen möglich, wobei einmal auf die Anfangswerte \varkappa_i , zum anderen auf die Endwerte $\bar{\varkappa}_i$ der Skalierungsparameter zurückgegriffen werden kann, vgl. (43), (44) bzw. (47), (48).

3.3.19. Algorithmus. Rekonstruction von α , β , γ , δ und $\overline{\varkappa}_p$, $\overline{\varkappa}_q$ aus ξ , \varkappa_p , $\varkappa_q \in \Re$. if $|\xi| \leq 1$ then | if $\xi \neq 1$

$$f = \gamma \text{ find } | \mathbf{h} \zeta + 1 \\ \text{then Fall 1: } | \alpha := \delta := 1 \\ | \lambda := \varkappa_p / \varkappa_q, \gamma := \xi * (1 + \lambda), \beta := \gamma / \lambda \\ | \tau := 1 + \gamma * \beta, \overline{\varkappa}_p := \varkappa_p / \tau, \overline{\varkappa}_q := \varkappa_q / \tau \\ \text{else Fall 2: } | \beta := \gamma := 1, \alpha := \delta := 0 \\ | \overline{\varkappa}_p := \varkappa_q, \overline{\varkappa}_q := \varkappa_p \\ \text{else Fall 2: } | \beta := \gamma := 1 \\ | \lambda := \varkappa_p / \varkappa_q, \alpha := (1 + \lambda) / \xi, \delta := \alpha / \lambda \\ | \tau := 1 + \alpha * \delta, \overline{\varkappa}_p := \varkappa_q / \tau, \overline{\varkappa}_q := \varkappa_p / \tau \end{cases}$$

Aufwand: 6 opm + 2 ops

3.3.20. Algorithmus. Rekonstruktion von α , β , γ , δ und \varkappa_p , \varkappa_q aus ξ , $\bar{\varkappa}_p$, $\bar{\varkappa}_q \in \mathfrak{R}$. **if** $|\xi| \leq 1$ then $\begin{vmatrix} \text{if } \xi \neq 1 \\ \text{then Fall 1:} \end{vmatrix} | \alpha := \delta := 1$ $\lambda := \bar{\varkappa}_p / \bar{\varkappa}_q$, $\gamma := \xi * (1 + \lambda)$, $\beta := \gamma / \lambda$ $\tau := 1 + \gamma * \beta$, $\varkappa_p := \bar{\varkappa}_p * \tau$, $\varkappa_q := \bar{\varkappa}_q * \tau$ else Fall 2: $\begin{vmatrix} \beta := \gamma := 1 \\ \lambda := \bar{\varkappa}_q / \bar{\varkappa}_p, \alpha := (1 + \lambda) / \xi, \delta := \alpha / \lambda$ $\tau := 1 + \alpha * \delta, \varkappa_p := \bar{\varkappa}_q * \tau, \varkappa_q := \bar{\varkappa}_p * \tau$

Aufwand: 6 opm + 2 ops

Die impliziten Givens-Drehungen zeigen ihre Effektivität erst dann, wenn sie in einer längeren Folge auftreten. Wir demonstrieren dies für die in 3.3.14 formulierte Aufgabe und realisieren die dort ausgeführten Transformationen

$$a^{1} = a, a^{i-1} \rightarrow a^{i} = G_{1i}a^{i-1} (i = 2, ..., n), \qquad \bar{a} = a^{n},$$

$$B^{(1)} = B, B^{(i-1)} \rightarrow B^{(i)} = G_{1i}B^{(i-1)} (i = 2, ..., n), \qquad \bar{B} = B^{(n)}$$
(54)

mittels der Matrizen F_{1i} und der Skalierungsfaktoren $\varkappa^{(i)} = (\varkappa^{(i)}_l) \in \mathbb{R}^n$ in der impliziten Form

$$d^{1} = d, d^{i-1} \to d^{i} = F_{1i}d^{i-1} (i = 2, ..., n), \qquad \bar{d} = d^{n},$$

$$E^{(1)} = E, E^{(i-1)} \to E^{(i)} = F_{1i}E^{(i-1)} (i = 2, ..., n), \qquad \bar{E} = E^{(n)}.$$
(55)

Dabei hängen die skalierten Größen aus (55) mit den unskalierten Größen aus (54) über

$$\boldsymbol{a}^{i} = \operatorname{diag}\left(\sqrt{\varkappa_{l}^{(i)}}\right) \boldsymbol{d}^{i}, \qquad \boldsymbol{B}^{(i)} = \operatorname{diag}\left(\sqrt{\varkappa_{l}^{(i)}}\right) \boldsymbol{E}^{(i)}$$
(56)

und

$$\boldsymbol{G}_{1i} = \operatorname{diag}\left(\sqrt{\varkappa_{l}^{(i)}}\right) \boldsymbol{F}_{1i} \operatorname{diag}\left(1/\sqrt{\varkappa_{l}^{(i-1)}}\right)$$
(57)

zusammen, wobei $\boldsymbol{x}^{(1)} = \boldsymbol{x}$ und $\boldsymbol{x}^{(n)} = \bar{\boldsymbol{x}}$ ist.

3.3.21. Eliminationsschritt mittels einer Folge impliziter Givens-Drehungen

Aufgabe: Gegeben seien $a \in \mathbb{R}^n$, $B \in \mathbb{R}^{n,r}$ in der skalierten Form

$$\boldsymbol{a} = \operatorname{diag}\left(\sqrt{\varkappa_l}\right) \boldsymbol{d}, \qquad \boldsymbol{B} = \operatorname{diag}\left(\sqrt{\varkappa_l}\right) \boldsymbol{E}$$
 (58)

durch $\bar{d} = (d_i) \in \Re^n$, $E = (e_{ij}) \in \Re^{n,r}$ und $\varkappa = (\varkappa_l) \in \Re^n$, $\varkappa_l > 0$. Die Transformation (54) ist mit Givens-Drehungen G_{1i} in der impliziten Version (55) so durchzuführen, daß

$$ar{m{a}}=m{G}_{1m{n}}\cdotsm{G}_{12}m{a}=arrhom{e}^1=(arrho,0,...,0)^\intercal$$

gilt. Die Resultate sind mit $\bar{d} = \sigma e^1 = (\sigma, 0, ..., 0)^{\intercal}$ gemäß

$$ar{m{a}} = ext{diag}\left(\sqrt{m{x}_l}
ight)ar{m{d}}, \qquad ar{m{B}} = ext{diag}\left(\sqrt{m{x}_l}
ight)ar{m{E}}$$
(59)

skaliert durch $\overline{d} \in \mathbb{R}^n$, $\overline{E} \in \mathbb{R}^{n,r}$ und $\overline{z} \in \mathbb{R}^n$ darzustellen. Die Berechnung soll in situ durch Überspeichern von d mit $(\sigma, \xi_2, \ldots, \xi_n)^{\mathsf{T}}$, von z durch \overline{z} und von Edurch \overline{E} erfolgen, wobei ξ_i die zu F_{1i} gemäß 3.3.18 bestimmte Kenngröße bezeichnet.

Algorithmus:

for i := 2(1)n do

- S0 (Festlegung der impliziten Givens-Drehung F_{1i}): Berechne $\{\alpha, \beta, \gamma, \delta, \overline{\varkappa}_1, \overline{\varkappa}_i, \sigma, \xi\}$ aus $\{d_1, d_i, \varkappa_1, \varkappa_i\}$ gemäß Algorithmus 3.3.18 mit p = 1, q = i. S1 $(d^{i-1} \rightarrow d^i = F_{1i}d^{i-1}, \varkappa^{i-1} \rightarrow \varkappa^i$, Reservierung von $\xi = \xi_i$):
- S1 $(d^{i-1} \rightarrow d^i = F_{1i}d^{i-1}, \tilde{\kappa}^{i-1} \rightarrow \tilde{\kappa}^i$, Reservering von $\xi = \xi_i$): $d_1 := \sigma, d_i := \xi, \, \kappa_1 := \bar{\kappa}_1, \, \kappa_i := \bar{\kappa}_i$

$$S2 (E^{(i-1)} \rightarrow E^{(i)} = F_{1i}E^{(i-1)}):$$
for $j := 1(1)r$ do
$$|if |\xi| < 1$$
then Fall 1: $|\mu := e_{1j} + \beta * e_{ij}|$
 $e_{ij} := -\gamma * e_{1j} + e_{ij}|$
 $e_{1j} := \mu$
else Fall 2: $|\mu := \alpha * e_{1j} + e_{ij}|$
 $e_{ij} := -e_{1j} + \delta * e_{ij}|$
 $e_{1j} := \mu$

Durch 3.3.21 wird

$$\begin{pmatrix} d_1 & e_{11} \dots e_{1r} \\ d_2 & e_{21} \dots e_{2r} \\ \vdots & \vdots & \vdots \\ d_n & e_{n1} \dots e_{nr} \end{pmatrix} \operatorname{mit} \begin{pmatrix} \sigma & \bar{e}_{11} \dots \bar{e}_{1r} \\ \xi_2 & \bar{e}_{21} \dots \bar{e}_{2r} \\ \vdots & \vdots & \vdots \\ \xi_n & \bar{e}_{n1} \dots \bar{e}_{nr} \end{pmatrix} \operatorname{und} \begin{pmatrix} \varkappa_1 \\ \varkappa_2 \\ \vdots \\ \varkappa_n \end{pmatrix} \operatorname{mit} \begin{pmatrix} \bar{\varkappa}_1 \\ \bar{\varkappa}_2 \\ \vdots \\ \bar{\varkappa}_n \end{pmatrix}$$

überspeichert.

Die Ergebnisse der Rundungsfehleranalyse von 3.3.21 sind weitgehend analog zu denen von 3.3.14.

3.3.22. Rundungsfehleranalyse. Algorithmus 3.3.21 werde mit den Eingangsdaten $\mathbf{x}, \mathbf{d} \in \mathbb{R}^n, \mathbf{E} \in \mathbb{R}^{n,r}$ durchgeführt, und $\mathbf{\bar{x}}, \mathbf{\bar{d}} = \sigma \mathbf{e}^1 = (\sigma, 0, \dots, 0)^{\mathsf{T}}$ sowie $\mathbf{\bar{E}}$ seien die berechneten Resultate. Mit G_{1i}^* werde die exakte Givens-Drehung bezeichnet, die $(\mathbf{a}^{i-1})^*$ exakt in $G_{1i}^*(\mathbf{a}^{i-1})^*$ mit $[G_{1i}^*(\mathbf{a}^{i-1})^*]_i = 0$ transformiert, wobei $(\mathbf{a}^{i-1})^*$ das exakte Produkt diag $(\sqrt{\kappa_i^{i-1}}) \mathbf{d}^{i-1}$ ist, und $\mathbf{G}^* = \mathbf{G}_{1n}^* \dots \mathbf{G}_{12}^*$ sei das exakte Produkt der G_{1i}^* . Dann gelten für die gemäß (58) bzw. (59) definierten unskalierten Größen \mathbf{a}, \mathbf{B} bzw. $\mathbf{\bar{a}}, \mathbf{\bar{B}}$ die Aussagen von 3.3.15.

Wir bemerken dazu nur, daß sich die geringfügigen Unterschiede zwischen 3.3.14und 3.3.21 auf die Schlußabschätzungen (27) bis (29) nicht auswirken. Die Elimination mittels impliziter Givens-Drehungen ist also — in bezug auf die unskalierten Originaldaten — ein numerisch gutartiger Prozeß.

3.3.23. Bemerkungen. (i) Der entscheidende Vorteil von 3.3.21 gegenüber 3.3.14 liegt in der Halbierung der Anzahl von Multiplikationen, die für die Transformation von B bzw. E in \overline{B} bzw. \overline{E} erforderlich sind. Die Vermeidung von Quadratwurzelberechnungen ist eine weitere angenehme, aber hier weniger wesentliche Eigenschaft der impliziten Givens-Drehungen.

(ii) In den Anwendungen — siehe Abschnitt 10.2 — tritt 3.3.21 in der Regel als Teilschritt eines *n*-stufigen Verfahrens auf, so daß sich die Anfangswerte $\mathbf{z} = (\mathbf{z}_l)$ der Skalierungsfaktoren als Endwerte der vorherigen Stufe ergeben. In der ersten Stufe wird $\mathbf{z}_l = 1$ (l = 1, ..., n) gesetzt, oder die Faktoren ergeben sich in natürlicher Weise aus der praktischen Aufgabenstellung, etwa als Gewichte bei einem linearen Regressionsproblem. Für gewisse Aufgabentypen wie die Lösung linearer Quadratmittelprobleme braucht $\overline{B} = \text{diag}(\sqrt[]{\overline{z}_l}) \overline{E}$ nicht explizit bekannt zu sein; die Lösungen lassen sich allein aus \overline{E} berechnen. Bei Bedarf kann \overline{B} mit dem Aufwand von höchstens n opr + nr opm selbstverständlich gebildet werden.

(iii) Falls 3.3.21 mit
$$\varkappa_l = 1$$
 $(l = 1, ..., n)$ gestartet wird, folgt wegen (53)

$$\bar{\varkappa}_l \ge 1/2^{n-1} \qquad (l=1,...,n).$$
 (60)

Für großes n kann daher in ungünstigen Fällen bei der Berechnung der Skalierungsfaktoren Unterlauf eintreten.

(iv) Die erwähnte schnelle Verkleinerung der Skalierungsfaktoren kann durch eine zu 3.2.7(ii) analoge Pivotisierung vermieden werden. Dazu wird der Zeilenindex s, $1 \leq s \leq n$, mit

$$\varkappa_s d_s^2 = \max_{\substack{i=1,\dots,n}} \varkappa_i d_i^2 \tag{61}$$

bestimmt, und statt der Eingangsdaten $\{\boldsymbol{x}, \boldsymbol{d}, \boldsymbol{E}\}$ werden in 3.3.21 die gemäß

$$\hat{\boldsymbol{x}} = \boldsymbol{T}_{1s}\boldsymbol{x}, \quad \hat{\boldsymbol{d}} = \boldsymbol{T}_{1s}\boldsymbol{d}, \quad \hat{\boldsymbol{E}} = \boldsymbol{T}_{1s}\boldsymbol{E}$$
 (62)

permutierten Eingangsdaten verwendet, d. h., die Zeilen 1 und s werden vertauscht. Dies ist gleichbedeutend mit dem Ersatz von

$$\boldsymbol{F} = \boldsymbol{F}_{1n} \cdots \boldsymbol{F}_{12} \quad \text{durch} \quad \boldsymbol{F} \boldsymbol{T}_{1s} = \boldsymbol{F}_{1n} \cdots \boldsymbol{F}_{12} \boldsymbol{T}_{1s}. \tag{63}$$

Im Originalproblem bedeutet (61)

$$a_s^2 = \max_{i=1,\dots,n} a_i^2, \tag{64}$$

d. h., die maximale Komponente von a wird in die erste Zeile gebracht. Wegen

$$(a_1^{i-1})^2 = (\hat{a}_1)^2 + (\hat{a}_2)^2 + \dots + (\hat{a}_{i-1})^2 \qquad (i = 2, \dots, n)$$
(65)

gilt dann

d. h., für jedes i = 2, ..., n liegt stets Fall 1 vor. Die Algorithmen 3.3.18 bis 3.3.21 reduzieren sich daher auf die für Fall 1 vorgesehenen Teile. Aus (65) folgt weiter die Darstellung

$$au = au_{i} = (a_{1}^{i})^{2}/(a_{1}^{i-1})^{2} = (\hat{a}_{1}^{2} + \dots + \hat{a}_{i}^{2})/(\hat{a}_{1}^{2} + \dots + \hat{a}_{i-1}^{2})$$

für das im *i*-ten Schritt berechnete τ , aus der sich wegen $\overline{z}_1 = \varkappa_1/(\tau_2 \cdots \tau_n), \ \overline{z}_i = \varkappa_i/\tau_i$ — es liegt stets Fall 1 vor! — sofort

$$\bar{\varkappa}_1 \ge \varkappa_1/n, \quad \bar{\varkappa}_i \ge \varkappa_i/2 \quad (i=2,...,n)$$
(66)

ergibt. Auch für sehr großes n kann daher praktisch kein Unterlauf eintreten. Implizite Givens-Drehungen sollten daher möglichst in der beschriebenen pivotisierten Form angewendet werden. \Box

Wir betrachten abschließend das Problem, nach Ausführung von 3.3.21 für eine weitere Matrix $V \in \mathbb{R}^{n,m}$ die Transformation

$$\overline{V} = GV = G_{1n} \cdots G_{12}V \tag{67}$$

durchzuführen, wobei V und \overline{V} in der skalierten Form

$$V = \operatorname{diag}\left(\sqrt{\varkappa_l}\right) W, \qquad \overline{V} = \operatorname{diag}\left(\sqrt{\varkappa_l}\right) \overline{W}$$
(68)

durch \varkappa , W bzw. $\overline{\varkappa}$, \overline{W} gegeben sein sollen. Es gilt dann

$$\overline{W} = FW = F_{1n} \cdots F_{12}W, \tag{69}$$

und die Berechnung kann schrittweise gemäß

$$W^{(1)} = W, W^{(i-1)} \to W^{(i)} = F_{1i}W^{(i-1)} \ (i = 2, ..., n), \overline{W} = W^{(n)}$$
(70)

erfolgen. Da bei der Rekonstruktion der Parameter von F_{1i} gemäß 3.3.19 die bezüglich des *i*-ten Schrittes "alten" Skalierungsfaktoren \varkappa^{i-1} benötigt werden, werden wie in 3.3.21 neben (70) die Transformationen

$$\mathbf{x}^1 = \mathbf{x}, \, \mathbf{x}^{i-1} \to \mathbf{x}^i \, (i = 2, \dots, n), \, \bar{\mathbf{x}} = \mathbf{x}^n \tag{71}$$

ausgeführt. Man beachte dabei, daß bei Anwendung von 3.3.19 auf F_{1i} die dort vorkommenden "alten" Skalierungsfaktoren z die Bedeutung von z^{i-1} haben; die "neuen" Faktoren \bar{z} entsprechen den z^i .

3.3.24. Algorithmus. Für gegebenes $\mathbf{z} \in \Re^n$, $\mathbf{W} = (w_{ij}) \in \Re^{n,m}$ ist $\mathbf{z} \in \Re^n$, $\overline{\mathbf{W}} = (\overline{w}_{ij}) \in \Re^{n,m}$ unter Verwendung der Größen $\{\xi_2, \ldots, \xi_n\}$ aus Algorithmus 3.3.21 so zu berechnen, daß (69) gilt. Die Eingangsdaten $\{\mathbf{z}, \mathbf{W}\}$ sind durch $\{\overline{\mathbf{z}}, \overline{\mathbf{W}}\}$ zu überspeichern.

for
$$i := 2(1)n$$
 do

S1: Berechne $\{\alpha, \beta, \gamma, \delta, \overline{\varkappa}_1, \overline{\varkappa}_i\}$ aus $\{\xi_i, \varkappa_1, \varkappa_i\}$ gemäß Algorithmus 3.3.19 mit $\xi = \xi_i, \varkappa_p = \varkappa_1, \varkappa_q = \varkappa_i.$

 $\begin{array}{l} \mathbf{S2} \; (\mathrm{Transformation} \; \boldsymbol{\varkappa}^{i^{-1}} \rightarrow \boldsymbol{\varkappa}^{i}) : \\ \boldsymbol{\varkappa}_{1} := \bar{\boldsymbol{\varkappa}}_{1}, \; \boldsymbol{\varkappa}_{i} := \bar{\boldsymbol{\varkappa}}_{i} \\ \mathbf{S3} \; (\mathrm{Transformation} \; \boldsymbol{W}^{(i-1)} \rightarrow \boldsymbol{W}^{(i)} = \boldsymbol{F}_{1i} \boldsymbol{W}^{(i-1)}) : \\ \mathbf{for} \; j := 1(1)m \; \mathbf{do} \\ \quad \mathbf{if} \; |\boldsymbol{\xi}| < 1 \\ \mathbf{then} \; \mathrm{Fall} \; 1 : \; | \; \boldsymbol{\mu} := w_{1j} + \boldsymbol{\beta} \ast w_{ij} \\ \boldsymbol{w}_{ij} := -\boldsymbol{\gamma} \ast w_{1j} + w_{ij} \\ \boldsymbol{w}_{1j} := \boldsymbol{\mu} \\ \mathbf{else} \; \; \mathrm{Fall} \; 2 : \; | \; \boldsymbol{\mu} := \boldsymbol{\alpha} \ast w_{1j} + w_{ij} \\ \boldsymbol{w}_{ij} := -w_{1j} + \boldsymbol{\delta} \ast w_{ij} \\ \boldsymbol{w}_{ij} := \mu \\ \mathbf{Aufwand} : \; (n-1) \left[(2m+6) \; \mathrm{opm} + (2m+2) \; \mathrm{ops} \right] \end{array}$

3.3.25. Bemerkungen. (i) Im Gegensatz zu den klassischen Givens-Drehungen kann F_{1i} nicht direkt aus ξ_i reproduziert werden, da dazu entweder \varkappa^{i-1} oder \varkappa^i erforderlich ist, die intermediären Skalierungsfaktoren jedoch aus Platzgründen nicht aufgehoben werden können.

(ii) Bei praktischen Aufgaben werden die F_{1i} nur in der direkten Reihenfolge $F_{12}, F_{13}, \ldots, F_{1n}$ bzw. in der inversen Reihenfolge $F_{1n}, F_{1,n-1}, \ldots, F_{12}$ benötigt. Ein Beispiel für den ersten Fall ist 3.3.24, wo die \varkappa^i in der analogen Reihenfolge gemäß 3.3.18 aufdatiert werden; statt S3 kann in 3.3.24 natürlich eine beliebige Operation

mit F_{1i} stehen. Falls die inverse Reihenfolge benötigt wird, ist 3.3.24 wie folgt zu modifizieren:

for i := n(-1)2 do

S1: Berechne $\{\alpha, \beta, \gamma, \delta, \varkappa_1, \varkappa_i\}$ aus $\{\xi_i, \overline{\varkappa}_1, \overline{\varkappa}_i\}$ gemäß Algorithmus 3.3.20 mit $\xi = \xi_i, \, \bar{\varkappa}_p = \bar{\varkappa}_1, \, \bar{\varkappa}_q = \bar{\varkappa}_i$ S2 (Transformation $\varkappa^i \rightarrow \varkappa^{i-1}$): $\bar{\varkappa}_1 := \varkappa_1, \, \bar{\varkappa}_i := \varkappa_i$ S3: Ausführung einer beliebigen Operation mit F_{1i}

Die inverse Reihenfolge wird z. B. bei der Berechnung von $Z = G^{-1}V$ durch Auswertung von

$$\mathbf{Z} = \operatorname{diag}\left(1/\sqrt{\varkappa_{l}}\right) \mathbf{F}_{12}^{\mathsf{T}} \mathbf{F}_{13}^{\mathsf{T}} \cdots \mathbf{F}_{1n}^{\mathsf{T}} \left[\operatorname{diag}\left(\sqrt{\varkappa_{l}}\right) \operatorname{diag}\left(\sqrt{\varkappa_{l}}\right) \mathbf{W}\right]$$
(72)

benötigt; man beachte dabei die Gültigkeit von

$$\boldsymbol{G}^{-1} = \boldsymbol{G}^{\mathsf{T}} = \left[\operatorname{diag}\left(\sqrt{\boldsymbol{\varkappa}_{l}}\right)\boldsymbol{F}\operatorname{diag}\left(1/\sqrt{\boldsymbol{\varkappa}_{l}}\right)\right]^{\mathsf{T}}.$$
(73)

Übungsaufgaben

Ü 3.3.1. Man zeige, daß sich im Fall n = 2 jede Householder-Spiegelung $H = I - 2uu^{\intercal}$, $\boldsymbol{u} = (\boldsymbol{\gamma}, \sigma)^{\mathsf{T}}, \, \boldsymbol{\gamma}^2 + \sigma^2 = 1, \, \text{in der Gestalt}$

$$oldsymbol{H} = egin{pmatrix} c & s \ s & -c \end{pmatrix} \hspace{1.5cm} ext{mit} \hspace{1.5cm} c^2 + s^2 = 1 \end{array}$$

schreiben läßt. Wie hängen $\{c, s\}$ und $\{\gamma, \sigma\}$ zusammen? Geometrische Deutung! Die analog zu G_{pq} für $p \neq q$ und $n \geq 2$ gebildeten Matrizen

 $\boldsymbol{H}_{pq} = \begin{pmatrix} 1 & \vdots & \vdots \\ \dots & c & \dots & s & \dots \\ \vdots & \vdots & \vdots \\ \dots & s & \dots & -c & \dots \\ \vdots & \vdots & \vdots \\ \dots & \vdots & \vdots \\ \end{pmatrix} \begin{array}{c} p\text{-te Zeile} \\ q\text{-te Zeile} \\ \end{array}$

werden deshalb auch ebene Spiegelungen bzw. Givens-Spiegelungen genannt. Für sie gilt $\det (H_{ng}) = -1$ im Gegensatz zu $\det (G_{ng}) = 1$. Sie können sinngemäß wie Givens-Drehungen verwendet werden.

Ü 3.3.2. Gegeben seien $a, b \in \mathbf{R}^n$ mit $a \neq b$, $\|a\|_2 = \|b\|_2$. Man konstruiere eine Householder-Spiegelung H, für die Ha = b gilt.

 ${f U}$ 3.3.3. Man zeige, daß die in ${f U}$ 3.3.2 mittels einer Householder-Spiegelung gelöste Aufgabe auch durch eine Folge $G = G_{p(n)q(n)} \cdots G_{p(2)q(2)}$ von n-1 Givens-Drehungen $G_{p(i)q(i)}$ (i = 2, ..., n) gelöst werden kann, wobei jede Drehung die Komponente $b_{q(i)}$ erzeugt.

Ü 3.3.4. Man überlege sich, daß $G^{(i)} = G_{1i} \cdots G_{13}G_{12}$ (i = 2, ..., n) von der Gestalt

ist, und nütze diesen Fakt in einem Algorithmus zur Berechnung von $G = G^{(n)}$ aus. Die Drehungen G_{1i} seien dabei durch die Zahlen ξ_i repräsentiert. Der Aufwand ist $(n-1) \times (2 \text{ opm} + 1 \text{ ops} + 1 \text{ opr})$ für die Rekonstruktion der c_i , s_i und (n-1)n opm für die Produktbildung.

Ü 3.3.5. Der in 3.3.13, 3.3.14 bzw. 3.3.21 beschriebene Eliminationsschritt kann z. B. auch durch die Folge $G = G_{12}G_{23} \cdots G_{n-2,n-1}G_{n-1,n}$ von ebenen Drehungen realisiert werden. Welche Gestalt hat die explizit berechnete Matrix G in diesem Fall?

Ü 3.3.6. Man zeige, daß das neue Element b_{ij} im Schritt S 2 von 3.3.14 auch gemäß

$$b_{ij} := \begin{cases} b_{ij} - z * (\mu + b_{1j}) & \text{für } c \ge 0, \\ z * (\mu - b_{1j}) - b_{ij} & \text{für } c < 0 \end{cases}$$

mit z := s/(1 + |c|) berechnet werden kann. Der Aufwand zur Berechnung von \overline{B} ist dann $\sim (n-1) r \times 3$ opms gegenüber $(n-1) r \times (4 \text{ opm} + 2 \text{ ops})$ in 3.3.14. Dieses Vorgehen ist daher zweckmäßig, wenn eine Multiplikation wesentlich aufwendiger als eine Addition ist.

Ü 3.3.7. Man lege $\boldsymbol{v} \in \mathbf{R}^n$ analog zu (10), (11) so fest, daß $\boldsymbol{a} = (a_i) \in \mathbf{R}^n$ in $\bar{\boldsymbol{a}} = (\bar{a}_i) = (\boldsymbol{I} - \boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}/\gamma) \boldsymbol{a}, \gamma = \boldsymbol{v}^{\mathsf{T}}\boldsymbol{v}/2$, mit

$$\bar{\pmb{a}}_i = \begin{cases} 0 & \text{für } p \leq i \leq q, \\ \varrho & \text{für } i = k, \\ a_i & \text{sonst} \end{cases}$$

übergeht, wobei entweder $1 \leq k oder <math>1 \leq p \leq q < k \leq n$.

Bemerkungen zum Kapitel 3

B 3.1. Seit WILKINSON [65] gehört die systematische Verwendung elementarer Transformationsmatrizen zum Standard in der Literatur über numerische lineare Algebra. In dem zitierten Buch ist auch eine Fehleranalyse der Grundoperationen mit solchen Matrizen zu finden, der hier im wesentlichen gefolgt wird.

B 3.2. Spiegelungsmatrizen wurden erstmalig von HOUSEHOLDER [58] mit der Vorzeichenfestlegung aus 3.3.5 (ii) für ϱ zur orthogonalen Elimination verwendet. Als Skalierungsfaktor wird in der Literatur meist $\lambda = 1$ gewählt, was auf die einfache Darstellung $\boldsymbol{v} = \boldsymbol{a} - \varrho \boldsymbol{e}^1, \gamma = -\varrho v_1$ führt und vom Aufwand und der Fehleranalyse her etwas günstiger ist als die Festlegung $\lambda = -1/\varrho$. Allerdings ist \boldsymbol{v} dann unter Umständen nicht gut skaliert, so daß Über- oder Unterlauf entstehen kann. Dies wird mit der hier gebrauchten Festlegung $\lambda = -1/\varrho$, die von den LINPACK-Autoren DONGARRA et al. [79] verwendet wird, vermieden. Die stabile Realisierung mit $\varrho > 0$ gemäß 3.3.5 (i) geht auf PARLETT zurück und kann z. B. in dessen Buch [80a] nachgelesen werden.

B 3.3. Bereits JACOBI [1846!] benutzte ebene Drehungen in dem nach ihm benannten Verfahren zur Eigenwertbestimmung symmetrischer Matrizen. Für dieselbe Aufgabe sind sie später von GIVENS [54] eingesetzt worden, der sie dann auch [58] für den beschriebenen Eliminationsschritt und allgemeiner für das Problem aus Ü 3.3.2 verwendet hat. Die effektive Speicherung der Drehungsparameter mittels einer reellen Größe ist von STEWART [76b] vorgeschlagen worden; Algorithmus 3.3.12 ist eine Modifikation des dort beschriebenen Vorgehens.

B 3.4. Die implizite Realisierung von Givens-Drehungen mittels geeigneter Skalierungen geht auf GENTLEMAN [73] und HAMMERLING [74] zurück. Die effektive Speicherung der Drehungsparameter gemäß 3.3.18 scheint neu zu sein.

II. Reguläre lineare Gleichungssysteme

4. Grundlegende Fakten über reguläre Gleichungssysteme

Wir betrachten in diesem Kapitel reguläre lineare Gleichungssysteme, d. h. Systeme mit quadratischen und zudem regulären Koeffizientenmatrizen, und untersuchen,

- wie sich Störungen der Koeffizientenmatrix und der rechten Seite auf die Lösung auswirken und welchen Einfluß eine Skalierung des Systems hat,
- nach welchem Prinzip die Gleichungssysteme in Systeme einfacher Struktur mit orthogonalen oder Dreiecksmatrizen transformiert werden können,
- wie solche Systeme einfacher Struktur gelöst werden.

4.1. Störungstheorie und Skalierung

A. Störungstheorie regulärer Gleichungssysteme

Vorgelegt sei das lineare Gleichungssystem

$$Ax = b \tag{1}$$

mit der regulären Koeffizientenmatrix $A \in \mathbb{R}^{n,n}$ und der rechten Seite $b \in \mathbb{R}^n$. Wegen der vorausgesetzten Regularität von A besitzt (1) für jedes b die eindeutige Lösung

$$\boldsymbol{x} = \boldsymbol{A}^{-1}\boldsymbol{b},\tag{2}$$

vgl. Abschnitt 1.1. Wir betrachten jetzt neben (1) das gestörte System

$$(\mathbf{A} + \mathbf{d}\mathbf{A}) (\mathbf{x} + \mathbf{d}\mathbf{x}) = \mathbf{b} + \mathbf{d}\mathbf{b}$$
(3)

mit Störungen $\delta A \in \mathbb{R}^{n,n}$, $\delta b \in \mathbb{R}^n$ und fragen, wann (3) eine eindeutige Lösung $x + \delta x$ besitzt und wie gegebenenfalls $\|\delta x\|$ durch $\|\delta A\|$ und $\|\delta b\|$ abgeschätzt werden kann.

Für die eindeutige Lösbarkeit von (3) ist die Regularität von $A + \delta A$ notwendig und hinreichend.

Wegen

$$A + \delta A = A(I + A^{-1}\delta A) = A(I - P), \qquad P := -A^{-1}\delta A \tag{4}$$

ist $A + \delta A$ genau dann regulär, wenn I - P regulär ist. Es genügt daher, die Regularität von I - P zu untersuchen.

4.1.1. Aussage. Für die Matrix $P \subset \mathbb{R}^{n,n}$ gelte $\|P\| < 1.$ (5)

Dann ist I - P regulär, und die Inverse genügt der Abschätzung

$$\|(\boldsymbol{I} - \boldsymbol{P})^{-1}\| \leq 1/(1 - \|\boldsymbol{P}\|).$$
(6)

Beweis. Für jedes $\boldsymbol{x} \in \mathbf{R}^n$ gilt

$$\|(\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{x}\| \ge \|\|\boldsymbol{x}\| - \|\boldsymbol{P}\boldsymbol{x}\|\| \ge \|\boldsymbol{x}\| - \|\boldsymbol{P}\| \|\boldsymbol{x}\| = (1 - \|\boldsymbol{P}\|) \|\boldsymbol{x}\|.$$
(7)

Wegen (5) ist 1 - ||P|| > 0, so daß aus (7) für $x \neq o$ sofort $(I - P) x \neq o$, also die Regularität von I - P folgt. Wenn $x := (I - P)^{-1} y$ gesetzt wird, geht (7) in

$$\|\boldsymbol{y}\| \geq (1 - \|\boldsymbol{P}\|) \|(\boldsymbol{I} - \boldsymbol{P})^{-1} \boldsymbol{y}\|$$
 für alle $\boldsymbol{y} \in \mathbf{R}^n$

über. Nach der Definition der Matrixnorm folgt hieraus (6), vgl. (1.1.39).

Wir erinnern daran, daß wir stets eine der Vektornormen mit dem Index $p \in \{1, 2, \infty\}$ und für Matrizen die zugeordneten Matrixnormen desselben Index verwenden. Abweichungen — etwa die Verwendung der Frobeniusnorm — werden explizit angegeben.

Anwendung von 4.1.1 auf (4) führt auf das folgende Resultat:

4.1.2. Störungslemma. Es sei $A \in \mathbb{R}^{n,n}$ eine reguläre Matrix, und die Störungsmatrix $\delta A \in \mathbb{R}^{n,n}$ genüge der Bedingung

$$\varkappa := \|A^{-1}\| \|\delta A\| < 1.$$
(8)

Dann ist auch die gestörte Matrix $A + \delta A$ regulär, und es gelten

$$\|(A + \delta A)^{-1}\| \le \|A^{-1}\|/(1 - \varkappa) \tag{9}$$

sowie

$$\|(A + \delta A)^{-1} - A^{-1}\| \le [\|A^{-1}\|^2/(1 - \varkappa)] \|\delta A\|.$$
(10)

Beweis. Wegen (8) ist $||\mathbf{P}|| \leq ||\mathbf{A}^{-1}|| \, ||\delta \mathbf{A}|| < 1$. Aus 4.1.1 folgt dann sofort die Regularität von $\mathbf{I} - \mathbf{P}$, also die von $\mathbf{A} + \delta \mathbf{A}$, und die Abschätzung $||(\mathbf{A} + \delta \mathbf{A})^{-1}|| = ||(\mathbf{I} - \mathbf{P})^{-1} \mathbf{A}^{-1}|| \leq ||(\mathbf{I} - \mathbf{P})^{-1}|| \, ||\mathbf{A}^{-1}|| \leq ||\mathbf{A}^{-1}||/(1 - \varkappa)$. Die Ungleichung (10) ergibt sich schließlich aus der Identität

$$(A + \delta A)^{-1} - A^{-1} = -(A + \delta A)^{-1} \, \delta A A^{-1}.$$

Lemma 4.1.2 besagt: Eine im Sinne von $\|\mathbf{d}A\| \leq \Delta A < 1/\|A^{-1}\|$ kleine Störung läßt die reguläre Matrix A regulär, und es gilt

$$\|(A + \delta A)^{-1} - A^{-1}\| \leq L(A, \Delta A) \|\delta A\|, \qquad L(A, \Delta A) := \|A^{-1}\|^2/(1 - \|A^{-1}\| \Delta A).$$

Dies bedeutet: Die Aufgabe $A \rightarrow A^{-1}$ ist auf der Menge der regulären Matrizen lokal

lipschitzstetig; die absolute bzw. relative Konditionszahl ist

 $L(A) = \|A^{-1}\|^2 \quad bzw. \quad K(A) = \|A\| \, \|A^{-1}\|,$

vgl. Abschnitt 2.1.

Unter der Voraussetzung (8) besitzt das gestörte System (3) die eindeutige Lösung $x + \delta x$. Ausmultiplizieren von (3) liefert dann unter Beachtung von Ax = b die Beziehung

$$\delta \boldsymbol{x} = (\boldsymbol{A} + \delta \boldsymbol{A})^{-1} \left(-\delta \boldsymbol{A} \, \boldsymbol{x} + \delta \boldsymbol{b} \right). \tag{11}$$

Anwendung von 4.1.2 führt auf das folgende Resultat:

4.1.3. Satz. Gegeben seien das lineare Gleichungssystem

Ax = b, $A \in \mathbb{R}^{n,n}$ regulär,

mit der Lösung $x = A^{-1}b$ sowie das gestörte System

$$(\boldsymbol{A}+\delta\boldsymbol{A})\,(\boldsymbol{x}+\delta\boldsymbol{x})=\boldsymbol{b}+\delta\boldsymbol{b}\quad ext{mit}\quad \boldsymbol{arkappa}:=\|\boldsymbol{A}^{-1}\|\,\|\delta\boldsymbol{A}\|<1\,.$$

(i) Dann ist $A + \delta A$ regulär, und für die eindeutige Lösung $x + \delta x$ des gestörten Systems gilt

 $\delta \boldsymbol{x} = \delta \boldsymbol{x}' + O(\|\delta \boldsymbol{A}\| (\|\delta \boldsymbol{A}\| + \|\delta \boldsymbol{b}\|)),$

wobei

$$\delta \boldsymbol{x}' := \boldsymbol{A}^{-1} \{ -\delta \boldsymbol{A} \, \boldsymbol{x} + \delta \boldsymbol{b} \} \tag{12}$$

den bezüglich δA , δb linearen Teil von δx bezeichnet.

(ii) Die Störung δx genügt den Abschätzungen

$$\|\boldsymbol{\delta x}\| \leq \frac{\|\boldsymbol{A}^{-1}\|}{1-\varkappa} \left\{ \|\boldsymbol{\delta A}\| \|\boldsymbol{x}\| + \|\boldsymbol{\delta b}\| \right\}$$
(13)

und im Fall $\boldsymbol{b} \neq \boldsymbol{o}$

$$\frac{\|\boldsymbol{\delta x}\|}{\|\boldsymbol{x}\|} \leq \frac{1}{1-\varkappa} \left\{ \|A\| \|A^{-1}\| \frac{\|\boldsymbol{\delta A}\|}{\|A\|} + \frac{\|A^{-1}\| \|\boldsymbol{b}\|}{\|\boldsymbol{x}\|} \frac{\|\boldsymbol{\delta b}\|}{\|\boldsymbol{b}\|} \right\}.$$
(14)

Beweis. Aus (11) folgt $\delta x = A^{-1}(-\delta A x + \delta b) + d$ mit $d := [(A + \delta A)^{-1} - A^{-1}] \times (-\delta A x + \delta b)$, unter Beachtung von (10) also (12). Durch direkte Abschätzung von (11) ergibt sich mit (9) die Ungleichung (13), aus der (14) nach Division durch ||x|| folgt. \Box

Die Abschätzung (14) besagt: Die Lösung regulärer Gleichungssysteme ist eine lokal lipschitzstetige Aufgabenklasse, und die partiellen relativen Konditionszahlen sind

$$K_{A}(A, b) = \|A\| \, \|A^{-1}\| \tag{15}$$

sowie

$$K_{b}(A, b) = \frac{\|A^{-1}\| \|b\|}{\|x\|}.$$
(16)

Wegen $||x|| = ||A^{-1}b|| \le ||A^{-1}|| ||b||$ und $||b|| = ||Ax|| \le ||A|| ||x||$ gilt dabei

$$1 \leq K_{\boldsymbol{b}}(\boldsymbol{A}, \boldsymbol{b}) \leq K_{\boldsymbol{A}}(\boldsymbol{A}, \boldsymbol{b}).$$

9 Schwetlick, Numerische Algebra

Die Zahl

$$\mathrm{cond}\ (A) := \|A\| \, \|A^{-1}\| \ge 1$$

heißt daher schlechthin Konditionszahl von A; gegebenenfalls wird cond (A) mit dem Index der verwendeten Norm versehen. Ist $\varkappa = ||A^{-1}|| ||\delta A|| = \text{cond} (A) ||\delta A||/||A|| \ll 1$ und wird K_b durch cond (A) abgeschätzt, so geht (14) näherungsweise in

$$\frac{\|\boldsymbol{\delta x}\|}{\|\boldsymbol{x}\|} \lessapprox \operatorname{cond}\left(A\right) \left\{ \frac{\|\boldsymbol{\delta A}\|}{\|\boldsymbol{A}\|} + \frac{\|\boldsymbol{\delta b}\|}{\|\boldsymbol{b}\|} \right\}$$

über, und dies bedeutet: Relative Störungen in der Matrix und der rechten Seite können sich höchstens mit dem Verstärkungsfaktor cond (A) auf die Lösung auswirken. Demgemäß heißt eine Matrix schlecht konditioniert, wenn cond $(A) \gg 1$ z. B. cond $(A) \ge 10^4 -$ gilt; das Gleichungssystem (1) heißt dann auch schlecht konditioniert.

B. Die Qualität der Störungsabschätzungen

Wir befassen uns im folgenden mit der Güte der Abschätzungen (13) und (14) und zeigen insbesondere, daß sie sich für $\delta A \rightarrow O$ nicht verbessern lassen, d. h., daß durch (15), (16) in der Tat die relativen Konditionszahlen von $x = A^{-1}b$ gegeben sind. Dazu betrachten wir alle Systeme (3) mit den durch

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq \Delta\boldsymbol{A}, \quad \|\boldsymbol{\delta}\boldsymbol{b}\| \leq \Delta\boldsymbol{b} \tag{17}$$

charakterisierten Eingangsdaten $\{A + \delta A, b + \delta b\}$ und den im Fall $||A^{-1}|| \Delta A < 1$ eindeutigen Lösungen $x + \delta x, \delta x$ nach (11). Um den verkomplizierenden Einfluß von δA im Term $(A + \delta A)^{-1}$ auszuschalten, untersuchen wir zunächst die linearisierten Störungen $\delta x' = A^{-1} \{-\delta A x + \delta b\}$ gemäß (12) und stellen für diese den folgenden Hilfssatz bereit.

4.1.4. Lemma. Für alle dA, db gemäß (17) genügt

$$d\boldsymbol{z} := -d\boldsymbol{A}\,\boldsymbol{x} + d\boldsymbol{b} \tag{18}$$

der Abschätzung

$$\|\delta z\| \leq \Delta z := \Delta A \|x\| + \Delta b, \qquad (19)$$

und diese Abschätzung ist exakt: Wenn δA , δb die Normkugeln (17) durchlaufen, durchläuft δz die Normkugel (19). Insbesondere läßt sich jeder beliebige Vektor δz mit $||\delta z|| = \Delta z$ in der Form (18) mit $||\delta A|| = \Delta A$, $||\delta b|| = \Delta b$ darstellen.

Beweis. Aus (18) folgt sofort (19). Gegeben sei nun δz mit $||\delta z|| \leq \Delta z$, also $\delta z = \lambda(\Delta z) z$ mit $0 \leq \lambda \leq 1$ und ||z|| = 1. Dann wird δz durch $\delta b := \lambda(\Delta b) z$ und

$$-\delta A := \lambda(\varDelta A) \, oldsymbol{z} oldsymbol{y}^{\intercal}, \hspace{0.2cm} oldsymbol{y} := egin{cases} \left((\mathrm{sgn} \, (x_1), \, ..., \, \mathrm{sgn} \, (x_n)
ight)^{\intercal} & \mathrm{für} \quad p = 1 \,, \ oldsymbol{x}/||oldsymbol{x}||_2 & \mathrm{für} \quad p = 2 \,, \ (\mathrm{sgn} \, (x_j)) \, oldsymbol{e}^j & \mathrm{für} \quad p = \infty \, \mathrm{mit} \, |x_j| = ||oldsymbol{x}||_{\infty}
ight)$$

im Fall x + o bzw. $\delta A := O$ im Fall x = o reproduziert; man beachte Abschnitt 1.1.I und U 1.1.3. Die Zusatzbemerkung ergibt sich für $\lambda = 1$. \Box 4.1.5. Aussage. Für alle δA , δb gemäß (17) genügen die linearisierten Störungen

$$dm{x}':=A^{-1}dm{z}=A^{-1}\{-dA \;m{x}+dm{b}\}$$

der Abschätzung

$$\|\delta \boldsymbol{x}'\| \leq \Delta \boldsymbol{x}' := \|\boldsymbol{A}^{-1}\| \, \Delta \boldsymbol{z} = \|\boldsymbol{A}^{-1}\| \left\{ \Delta \boldsymbol{A} \|\boldsymbol{x}\| + \Delta \boldsymbol{b} \right\},\tag{20}$$

und diese Abschätzung ist scharf: Es gibt Störungen δA , δb mit $\|\delta A\| = \Delta A$, $\|\delta b\| = \Delta b$, für die in (20) das Gleichheitszeichen steht.

Beweis. Die Schranke (20) ist trivial. Zum Nachweis der Schärfe beachten wir, daß zu $A^{-1} \operatorname{ein} z$ mit ||z|| = 1 und $||A^{-1}z|| = ||A^{-1}|| \, ||z||$ existiert, vgl. Abschnitt 1.1.I. Mit $\delta z := (\Delta z) \, z$ und 4.1.4 folgt dann die Behauptung. \Box

4.1.6. Folgerung. Es gelte $||A^{-1}|| \Delta A < 1$. Dann ist die für alle δA , δb gemäß (17) gültige und zu (13) äquivalente Abschätzung

$$\|\delta x\| \le \Delta x := \Delta x'/(1 - \|A^{-1}\| \Delta A) = \|A^{-1}\| \{\Delta A \|x\| + \Delta b\}/(1 - \|A^{-1}\| \Delta A)$$
(21)

scharf für $\Delta A \rightarrow 0$: Zu ΔA , Δb mit $\Delta x > 0$ gibt es Störungen δA , δb mit $\|\delta A\| = \Delta A$, $\|\delta b\| = \Delta b$, so daß

$$\|\delta \boldsymbol{x}\|/\Delta \boldsymbol{x} \to 1 \quad \text{für} \quad \Delta \boldsymbol{A} \to 0$$
 (22)

gilt.

Dies besagt: Die Abschätzungen (13) und (14) sind für $\Delta A \rightarrow 0$ scharf bezüglich der Gesamtheit der durch (17) charakterisierten Störungen. Für eine einzelne feste Störung braucht (13) bzw. (14) natürlich nicht scharf zu sein, siehe Ü 4.1.5. Welcher Teil der Normkugel (20) durch die in bezug auf (17) zulässigen linearisierten Störungen $\delta x'$ ausgefüllt wird, läßt sich für die 2-Norm in einfacher Weise explizit angeben.

4.1.7. Satz. Wenn dA, db die Normkugeln $||dA||_2 \leq \Delta A$, $||db||_2 \leq \Delta b$ durchlaufen, durchlaufen die linearisierten Störungen dx' das Ellipsoid $\mathfrak{E}' = \mathfrak{E}'(A, b, \Delta A, \Delta b)$ mit dem Mittelpunkt O und den Halbachsen

$$\pm lpha_j oldsymbol{v}^j, \qquad lpha_j := arDelta oldsymbol{z} / \sigma_j \qquad (j=1,\,...,\,n)\,,$$

wobe
i σ_{j} und \boldsymbol{v}^{j} die Singulärwerte bzw. Singulärvektoren von
 A nach 1.2.14 bezeichnen.

Beweis. Mit der Singulärwertzerlegung $A = U\Sigma V^{\mathsf{T}}$ gilt $\delta x' = A^{-1}\delta z = V\Sigma^{-1}U^{\mathsf{T}}\delta z$. Nach 4.1.4 durchläuft δz die Normkugel $||\delta z||_2 \leq \Delta z$, und dasselbe gilt wegen der Orthogonalität von U für $\delta y := (\delta y_j) = U^{\mathsf{T}}\delta z$. Die zulässigen Störungen lassen sich dann durch

$$\delta \boldsymbol{x}' = \sum_{j=1}^{n} \boldsymbol{v}^{j} (\delta y_{j} / \sigma_{j}) \quad \text{mit} \quad \sum_{j=1}^{n} (\delta y_{j})^{2} \leq (\boldsymbol{\varDelta} \boldsymbol{z})^{2}$$
(23)

darstellen. Die Größen $\xi_j := \delta y_j / \sigma_j$ sind die Komponenten von $\delta x'$ bezüglich der orthonormierten Basis $\{v^1, \ldots, v^n\}$, mit denen sich (23) äquivalent

$$\delta \boldsymbol{x}' = \sum_{j=1}^{n} \boldsymbol{v}^{j} \boldsymbol{\xi}_{j} \quad \text{mit} \quad \sum_{j=1}^{n} \left(\frac{\boldsymbol{\xi}_{j}}{\boldsymbol{\alpha}_{j}}\right)^{2} = \sum_{j=1}^{n} \left(\frac{\delta y_{j}}{\Delta \boldsymbol{z}}\right)^{2} \leq 1$$

schreiben läßt. Dies ist gerade die Gleichung des Ellipsoids &'. 🗌

Die längste Halbachse von \mathscr{E}' hat die Länge $\alpha_n = \varDelta \boldsymbol{z}/\sigma_n = \varDelta \boldsymbol{x}'$, so daß die Normkugel $\mathscr{K}' = \mathscr{K}'(\boldsymbol{A}, \boldsymbol{b}, \varDelta \boldsymbol{A}, \varDelta \boldsymbol{b}) := \{ \boldsymbol{\delta x}' : \| \boldsymbol{\delta x}' \| \leq \varDelta \boldsymbol{x}' \}$ aus (20) in der Tat die umschriebene Kugel zu \mathscr{E}' darstellt. Das Verhältnis der längsten zur kürzesten Halbachse von \mathscr{E}' ist

$$\alpha_n / \alpha_1 = \sigma_1 / \sigma_n = \operatorname{cond}_2(A). \tag{24}$$

Je schlechter die Matrix konditioniert ist, um so mehr weicht als
o \mathcal{E}' von der Kugel \mathcal{K}' ab, siehe Abb. 4.1.1.



Abb. 4.1.1. Linearisierte Störungen für n = 2 im Fall der 2-Norm

Bei der Anwendung von (13) zur Abschätzung von $||\delta x||$ entsteht folgendes Problem: Häufig sind für die Störungen $\delta A = (\delta a_{ij}), \ \delta b = (\delta b_i)$ der Eingangsdaten individuelle Schranken mit

$$|\delta a_{ij}| \leq \Delta a_{ij}, \qquad |\delta b_i| \leq \Delta b_i \qquad (i, j = 1, ..., n)$$

$$\tag{25}$$

bekannt, vgl. Abschnitt 2.3.C. Wenn diese in der Matrix $\Delta A := (\Delta a_{ij})$ bzw. im Vektor $\Delta b := (\Delta b_i)$ zusammengefaßt werden, läßt sich (25) kurz als

$$|\mathbf{d}\mathbf{A}| \leq \Delta \mathbf{A}, \quad |\mathbf{d}\mathbf{b}| \leq \Delta \mathbf{b} \tag{26}$$

schreiben, siehe Abschnitt 1.1.J. Für die hier betrachteten Normen $p \in \{1, 2, \infty\}$ folgt aus (26)

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq \|\boldsymbol{\Delta}\boldsymbol{A}\| =: \boldsymbol{\Delta}\boldsymbol{A}, \quad \|\boldsymbol{\delta}\boldsymbol{b}\| \leq \|\boldsymbol{\Delta}\boldsymbol{b}\| =: \boldsymbol{\Delta}\boldsymbol{b}, \quad (27)$$

also (17) mit den in (27) definierten Werten von ΔA und Δb . Die $n^2 + n$ Schranken aus (25) werden dann durch die zwei Zahlen aus (27) repräsentiert. Obwohl die zugehörigen Schranken (13) in bezug auf (27) asymptotisch scharf sind, wird dies in bezug auf (26) i. allg. nicht zu erwarten sein. Es ist daher sinnvoll zu fragen, unter welchen Bedingungen die Schranke (13) das durch (26) charakterisierte individuelle Fehlerverhalten realistisch widerspiegelt. Da wir an komponentenweisen Abschätzungen für δx interessiert sind, bietet sich in (13) die Verwendung der ∞ -Norm an. Der Einfachheit halber betrachten wir wieder die linearisierten Störungen dx' und gehen analog zum Fall der Normschranken vor.

4.1.8. Lemma. Für alle δA , δb gemäß (26) genügt $\delta z = -\delta A x + \delta b$ der Abschätzung

$$|d\boldsymbol{z}| \leq \Delta \boldsymbol{z} := \Delta \boldsymbol{A} |\boldsymbol{x}| + \Delta \boldsymbol{b}, \qquad (28)$$

und diese Abschätzung ist exakt, d. h., sie beschreibt genau die Menge der bezüglich (26) zulässigen Störungen dz. Insbesondere läßt sich jeder beliebige Vektor dz mit $|dz| = \Delta z$ in der Form (18) darstellen mit Störungen dA, db, für die $|dA| = \Delta A$ und $|db| = \Delta b$ gilt.

Beweis. Wegen

 $|\delta \boldsymbol{z}| = |-\delta A |\boldsymbol{x} + \delta \boldsymbol{b}| \leq |\delta A |\boldsymbol{x}| + |\delta \boldsymbol{b}| \leq |\delta A| |\boldsymbol{x}| + |\delta \boldsymbol{b}| \leq \Delta A |\boldsymbol{x}| + \Delta \boldsymbol{b}$

gilt (28). Für den Beweis der Zusatzbemerkung setze man $\delta a_{ij} := \pm \Delta a_{ij}$, $\delta b_i := \pm \Delta b_i$ und verteile die Vorzeichen so, daß $|-\delta A x + \delta b| = |\delta A| |x| + |\delta b|$ gilt; dies ist offenbar stets möglich.

Zur Untersuchung von dx' benötigen wir das folgende Lemma.

4.1.9. Lemma. Für $B \in \mathbf{R}^{n,n}$ und δz mit $|\delta z| \leq \Delta z$ werde $\delta x' := B \, \delta z$ gebildet. Dann gilt

(i) $|\boldsymbol{\delta x'}| \leq |\boldsymbol{B}| |\boldsymbol{\delta z}| \leq \Delta x' := |\boldsymbol{B}| \Delta \boldsymbol{z},$ (29)

und diese Abschätzung ist komponentenweise scharf: Zu jedem *i* gibt es ein δz mit $|\delta z| = \Delta z$, so daß in der *i*-ten Komponente von (29) das Gleichheitszeichen steht.

(ii)
$$\|\boldsymbol{\delta}\boldsymbol{x}'\|_{\infty} \leq \|\boldsymbol{B}\|_{\infty} \|\boldsymbol{\Delta}\boldsymbol{z}\|_{\infty},$$
 (30)

und diese Abschätzung ist scharf, wenn $\Delta z_i = \zeta$ (i = 1, ..., n) gilt: Es gibt dann ein ∂z mit $|\partial z| = \Delta z$, für welches in (30) das Gleichheitszeichen steht.

Beweis. Die Herleitung von (29) ist trivial. Zum Nachweis der Schärfe setzen wir $(\delta z)_j = \pm \Delta z_j$ und verteilen die Vorzeichen so, daß $|(\boldsymbol{B}\delta z)_i| = (|\boldsymbol{B}| |\delta z|)_i$ gilt. Wegen der Monotonie und Absolutheit der ∞ -Norm folgt (30) direkt aus (29). Sei nun $\Delta z_i = \zeta$ für alle *i*. Nach Definition von $||\boldsymbol{B}||_{\infty}$ gibt es ein \boldsymbol{z} mit $|\boldsymbol{z}_j| = 1$ $(j = 1, \dots, n)$ und $||\boldsymbol{B}\boldsymbol{z}||_{\infty} = ||\boldsymbol{B}||_{\infty} ||\boldsymbol{z}||_{\infty}$, vgl. (1.1.41). Für $\delta \boldsymbol{z} := \zeta \boldsymbol{z}$ ist dann (30) scharf, und es gilt $|\delta \boldsymbol{z}| = \Delta \boldsymbol{z}$.

4.1.10. Aussage. Für alle δA , δb gemäß (26) werde

$$\delta oldsymbol{x}' = A^{-1} \delta oldsymbol{z} = A^{-1} \{ - \delta A \ oldsymbol{x} + \delta oldsymbol{b} \}$$

gebildet. Dann gilt

(i)
$$|\boldsymbol{\delta x'}| \leq |\boldsymbol{A}^{-1}| \boldsymbol{A z} \quad \text{mit} \quad \boldsymbol{A z} := \boldsymbol{A A} |\boldsymbol{x}| + \boldsymbol{A b},$$
 (31)

und diese Abschätzung ist komponentenweise scharf: Zu jedem *i* gibt es Störungen δA , δb mit $|\delta A| = \Delta A$, $|\delta b| = \Delta b$, so daß in der *i*-ten Komponente von (31) das Gleichheitszeichen steht.

(ii)
$$\|\boldsymbol{\delta x}'\|_{\infty} \leq \|\boldsymbol{A}^{-1}\|_{\infty} \|\boldsymbol{\Delta z}\|_{\infty} = \|\boldsymbol{A}^{-1}\|_{\infty} \|\boldsymbol{\Delta A}\|\boldsymbol{x}\| + \|\boldsymbol{\Delta b}\|_{\infty}.$$
 (32)

Im Fall $\Delta z_i = \zeta$ (i = 1, ..., n) ist diese Abschätzung scharf: Es gibt Störungen δA , δb mit $|\delta A| = \langle A, |\delta b| = \Delta b$, so daß in (32) das Gleichheitszeichen steht.

(iii)
$$\|\boldsymbol{\delta x'}\|_{\infty} \leq \|A^{-1}\|_{\infty} \{\|AA\|_{\infty} \|\boldsymbol{x}\|_{\infty} + \|Ab\|_{\infty} \}.$$
 (33)

Im Fall

$$|x_i| = \xi, \qquad \Delta b_i = \beta, \qquad \sum_{j=1}^n \Delta a_{ij} = \sigma \qquad (i = 1, ..., n)$$
 (34)

ist (33) scharf analog zu (ii).

Be we is. Nach 4.1.8 wird für δA , δb nach (26) der Bereich $|\delta z| \leq \Delta z$ voll ausgeschöpft. Anwendung von 4.1.9 mit $B := A^{-1}$ liefert dann die Aussagen (i) und (ii). Aus (32) folgt sofort (33). Es bleibt zu zeigen, daß (33) unter der Voraussetzung (34) scharf ist. Aus (34) folgt nun $\Delta z_i = \sum_j \Delta a_{ij} |x_j| + \Delta b_i = \sigma \xi + \beta$, d. h.. Δz erfüllt die Voraussetzungen von (ii), und (32) ist scharf. Wegen $||\Delta z||_{\infty} = \sigma \xi + \beta = ||\Delta A_{||_{\infty}} ||x_j||_{\infty} + ||\Delta b||_{\infty}$ stimmen die rechten Seiten von (32) und (33) überein, so daß auch (33) scharf ist. \Box

Im folgenden wird ein Vektor mit betragsgleichen Komponenten äquilibriert genannt. Analog heißt eine Matrix zeilenäquilibriert, falls die Zeilensummen der Beträge gleich sind.

Wenn die Störungen ∂A , ∂b den Bedingungen (26) genügen, geht (13) für die ∞ -Norm in

$$\|\boldsymbol{\delta x}\|_{\infty} \leq \frac{\|A^{-1}\|_{\infty}}{1 - \|A^{-1}\|_{\infty} \|AA\|_{\infty}} \{\|AA\|_{\infty} \|\boldsymbol{x}\|_{\infty} + \|Ab\|_{\infty}\}$$
(35)

über, vgl. 4.1.6. Aus 4.1.10 lassen sich dann hinreichende Bedingungen ablesen, unter denen die Schranke (35) realistisch ist.

4.1.11. Aussage. Es sei $||A^{-1}||_{\infty} ||A||_{\infty} < 1$, und die Komponenten x_i von x mögen sich betragsmäßig nicht zu sehr unterscheiden. Dann ist die Abschätzung (35) realistisch in bezug auf die Störungen (26), sofern A zeilenäquilibriert und A äquilibriert ist.

4.1.12. Bemerkung. Falls die Elemente von A dasselbe relative Fehlerniveau haben d. h., falls

$$|\delta a_{ij}| \leq \varepsilon |a_{ij}| =: \Delta a_{ij}$$

gilt, ist $\sum_{j} \Delta a_{ij} = \varepsilon \sum_{j} |a_{ij}|$ und ΔA folglich genau dann zeilenäquilibriert, wenn A selbst diese Eigenschaft hat. Analoge Aussagen gelten für **b** und $\Delta \mathbf{b}$.

Wir bemerken abschließend, daß (35) und allgemein die Abschätzungen aus 4.1.3 für nicht ausreichend äquilibrierte Schranken AA und Ab das bezüglich (26) wahre Fehlerniveau wesentlich überschreiten können. Zur Illustration geben wir das folgende Beispiel an, siehe auch Ü 4.1.7. 4.1.13. Beispiel. Für

$$A = \begin{pmatrix} 1 & 0.5 & 0 \\ 0.97 & 0.48 & 0 \\ 0.75 & 0 & -0.75 \end{pmatrix}, \quad \Delta A = O, \quad b = \begin{pmatrix} 3.50 \\ 3.37 \\ -6.75 \end{pmatrix}, \quad \Delta b = 10^{-3} \begin{pmatrix} 1 \\ 1 \\ 75 \end{pmatrix}$$

ist

Die letzte Abschätzung ist nach 4.1.10 scharf. Die Abschätzung (35) liefert dagegen

$$\|\delta \boldsymbol{x}\|_{\infty} \leq \|A^{-1}\|_{\infty} \|\Delta \boldsymbol{b}\|_{\infty} = 10^{-3} \times 394 \times 75$$
,

also eine um den Faktor 75 zu grobe Schranke.

Wenn die dritte Gleichung des Systems durch 75 dividiert wird, erhält man das durch

$$\hat{A} = \begin{pmatrix} 1 & 0.5 & 0 \\ 0.97 & 0.48' & 0 \\ 0.01 & 0 & -0.01 \end{pmatrix}, \quad \Delta \hat{A} = O, \quad \hat{b} = \begin{pmatrix} 3.50 \\ 3.37 \\ -0.09 \end{pmatrix}, \quad \Delta \hat{b} = 10^{-3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

und

$$\hat{A}^{-1} = \begin{pmatrix} -96 & 100 & 0\\ 194 & -200 & 0\\ -96 & 100 & -100 \end{pmatrix}$$

charakterisierte äquivalente System mit äquilibriertem 4b. Die Schranke (35) liefert hier

$$egin{aligned} \|\delta m{x}\|_{\infty} &\leq \|\hat{m{A}}^{-1}\|_{\infty} \, \|arpi m{b}\|_{\infty} = 10^{-3} imes 394 imes 1 \, . \end{aligned}$$

also eine scharfe Abschätzung des möglichen Fehlers.

C. Skalierung

Das eben diskutierte Beispiel zeigt, daß sich die für eine realistische Fehlerabschätzung günstigen Eigenschaften von ΔA bzw. Δb unter Umständen erzielen lassen, wenn Eingangsdaten und Fehlerschranken zeilenweise mit geeigneten positiven Faktoren multipliziert werden. Dieser Übergang vom Originalsystem in das äquivalente System

$$\bar{A}x = \bar{b}$$
 mit $\bar{A} := DA$, $\bar{b} := Db$ (36)

und den Fehlerschranken

$$\Delta \mathbf{A} := \mathbf{D} \Delta \mathbf{A}, \qquad \Delta \mathbf{\bar{b}} := \mathbf{D} \Delta \mathbf{b}$$
(37)

mittels der Diagonalmatrix

$$\boldsymbol{D} = \text{diag}\left(d_{i}\right), \qquad d_{i} > 0,$$

wird Skalierung – genauer: Zeilenskalierung – genannt.

Bei der Skalierung werden die Eingangsdaten $\{A, b\}$ des Originalsystems zu $\{\bar{A}, \bar{b}\}$ verändert. Ein numerischer Algorithmus V zur Gleichungsauflösung wird daher bei

Anwendung auf das skalierte System (36) i. allg. eine etwas andere Lösung $\bar{x} = V(\bar{A}, \bar{b})$ liefern als bei Anwendung auf das Originalsystem, wo sich x = V(A, b) ergibt.

Die praktisch verwendeten direkten Lösungsverfahren V sind alle numerisch gutartig in bezug auf A, siehe Kapitel 5 und 11. Es gibt dann eine Störung δA_R , so daß die berechnete Lösung x dem gestörten System

$$(A + \delta A_R) \boldsymbol{x} = \boldsymbol{b} \quad \text{mit} \quad \|\boldsymbol{\delta} A_R\| \leq \nu F \|A\|$$
(38)

genügt, wobei die Kumulationskonstante F nur vom Verfahren und der Dimension n, nicht aber von A abhängt, vgl. Abschnitt 2.3. Der erzeugte Rundungsfehler $\delta x_R := x - x^*$ gegenüber der exakten Lösung $x^* = A^{-1}b$ kann dann nach (13) bzw. (35) zu

$$\|\boldsymbol{\delta x_R}\|_{\infty} \leq \nu F \operatorname{cond}_{\infty}(\boldsymbol{A}) \, \|\boldsymbol{x^*}\|_{\infty} / [1 - \nu F \operatorname{cond}_{\infty}(\boldsymbol{A})] \tag{39}$$

abgeschätzt werden. Für das skalierte System gilt (39) mit DA anstelle von A, so daß die Frage entsteht, ob sich $\operatorname{cond}_{\infty}(DA)$ und damit die Schranke (39) durch geeignete Skalierung von A minimieren läßt. Der folgende Satz gibt eine positive Antwort.

4.1.14. Satz. Es sei

$$A = \begin{pmatrix} -a^{1\mathsf{T}} - \\ \vdots \\ -a^{n\mathsf{T}} - \end{pmatrix} \in \mathsf{R}^{n,n}$$

eine reguläre Matrix mit den Zeilen $a^{i\intercal}$, und mit $D = \text{diag}(d_i)$,

$$d_{i} := \frac{\max \sum_{j} |a_{kj}|}{\sum_{j} |a_{ij}|} = \frac{\max \|a^{k}\|_{1}}{\|a^{i}\|_{1}},$$
(40)

werde $\bar{A} := DA$ gebildet. Dann gelten

$$\|\boldsymbol{A}\|_{\infty} = \|\bar{\boldsymbol{A}}\|_{\infty}, \quad \|\boldsymbol{A}^{-1}\|_{\infty} / \max \, d_i \leq \|\bar{\boldsymbol{A}}^{-1}\|_{\infty} \leq \|\boldsymbol{A}^{-1}\|_{\infty}$$
(41)

$$\operatorname{cond}_{\infty}(A)/\operatorname{max} d_{i} \leq \operatorname{cond}_{\infty}(\bar{A}) \leq \operatorname{cond}_{\infty}(A).$$
 (42)

Be we is. Die Gleichheit von $||A||_{\infty}$ und $||\bar{A}||_{\infty}$ folgt sofort aus (40), denn \bar{A} ist zeilenäquilibriert, und die Zeilensummen von \bar{A} sind gleich der maximalen Zeilensumme von A. Wegen $||D||_{\infty} = \max d_i \operatorname{und} ||D^{-1}||_{\infty} = 1/\min d_i = 1$ folgt weiter $||\bar{A}^{-1}||_{\infty} = ||A^{-1}D^{-1}||_{\infty} \leq ||A^{-1}||_{\infty} ||D^{-1}||_{\infty} = ||A^{-1}||_{\infty} ||D^{-1}||_{\infty} \leq ||A^{-1}||_{\infty} ||D^{-1}||_{\infty} \leq ||\bar{A}^{-1}||_{\infty} ||D^{-1}||_{\infty} \leq ||\bar{A}^{-1}||_{\infty} ||D^{-1}||_{\infty} \leq ||A^{-1}||_{\infty} ||D^{-1}||_{\infty} ||D^{-1}||_{\infty} \leq ||A^{-1}||_{\infty} ||D^{-1}||_{\infty} ||D^{-1}||_{\infty$

Die Abschätzung (42) besagt: Unter allen durch Zeilenskalierung aus einer Matrix hervorgehenden Matrizen hat jede zeilenäquilibrierte die kleinste Kondition in der ∞ -Norm. Man beachte dabei, daß sich die Kondition von **D**A nicht ändert, wenn **D** durch λD , $\lambda > 0$, ersetzt wird, d. h., es kommt nur auf das Verhältnis der Skalierungsfaktoren d_i an. Die in (40) vorgenommene Normierung min $d_i = 1$ wird nur für die (absoluten) Abschätzungen (41) benötigt. In bezug auf die Schranke (39) für den erzeugten Rundungsfehler δx_R ist nach 4.1.14 also eine Zeilenäquilibrierung des Systems optimal.

In der Regel wird der Fehler δx_R durch weitere Anteile δx_D , δx_M zum Gesamtfehler $\delta x := \delta x_R + \delta x_D + \delta x_M$ ergänzt. Wenn die Eingangsdatenelemente nicht bereits Computerzahlen sind, tritt ein relativ durch v beschränkter, also sehr kleiner Darstellungsfehler δA_D , δb_D auf. Nach 2.3.15 gilt für diesen

$$\| \boldsymbol{\delta} \boldsymbol{A}_D \|_\infty \leq
u \| \boldsymbol{A} \|_\infty, \qquad \| \boldsymbol{\delta} \boldsymbol{b}_D \|_\infty \leq
u \| \boldsymbol{b} \|_\infty.$$

Der zugehörige Fehler δx_D in der Lösung ist daher durch

$$\|\boldsymbol{\delta x_{\mathcal{D}}}\|_{\infty} \leq 2\nu \operatorname{cond}_{\infty}(\boldsymbol{A}) \|\boldsymbol{x^{*}}\|_{\infty} / (1 - \nu \operatorname{cond}_{\infty}(\boldsymbol{A}))$$

$$\tag{43}$$

beschränkt, also etwa um den Faktor F/2 kleiner als der erzeugte Rundungsfehler und wird deshalb i. allg. vernachlässigt werden können.

Wenn die Eingangsdaten das Resultat realer Messungen sind oder durch komplizierte, eventuell nur näherungsweise ausgeführte Berechnungen erhalten wurden, treten weitere Fehler δA_M , δb_M auf, die deutlich über dem Niveau des Darstellungsfehlers liegen. Man beachte, daß v für reale Computer in der Größenordnung $10^{-6} \dots 10^{-14}$ liegt, während Meßfehler meist eine relative Genauigkeit von $10^{-1} \dots 10^{-3}$ aufweisen. Solche "großen" Fehler dA_M , db_M wollen wir unabhängig von ihrer konkreten Quelle Meßtehler nennen, und der durch sie hervorgerufene Fehler δx_M in der Lösung dominiert i. allg. die Anteile δx_R und δx_D deutlich. Die durch Zeilenäquilibrierung mögliche Verkleinerung von δx_R wird daher bedeutungslos, dagegen ist die realistische Abschätzung des durch die Meßfehler hervorgerufenen Fehleranteils dx_M wesentlich. Diese kann jedoch durch die Zeilenäquilibrierung verhindert werden; vgl. Beispiel 4.1.13, wo sich mit der fast zeilenäquilibrierten Matrix A eine wesentlich schlechtere Abschätzung als für die nicht äquilibrierte Matrix \hat{A} ergibt. Aus 4.1.9 und 4.1.10 wird klar, daß bei dominanten Meßfehlern nicht die Matrix A, sondern der Vektor $\Delta z = \Delta A |x| + A b$ bzw. die Schranken ΔA und Δb zu äquilibrieren sind.

4.1.15. Skalierungsregeln.

(R₁) Wenn die Eingangsdaten keine Meßfehler aufweisen, skaliere man so, daß A zeilenäquilibriert wird:

$$d_i := \max_k \sum_j |a_{kj}| / \sum_j |a_{ij}|$$
.

(R₂) Wenn die rechte Seite **b** einen Meßfehler $\delta \mathbf{b}_M$ mit $|\delta \mathbf{b}_M| \leq |\mathbf{b}|$ besitzt und A im Rahmen des Darstellungsfehlers exakt ist oder einen gegenüber $\Delta \mathbf{b}$ vernachlässigbaren Fehler aufweist, skaliere man so, daß $\Delta \mathbf{b}$ äquilibriert wird:

$$d_i := \max \Delta b_k / \Delta b_i.$$

(R₃) Wenn die Matrix A einen Meßfehler δA_M mit $|\delta A_M| \leq \Delta A$ besitzt und b im Rahmen des Darstellungsfehlers exakt ist oder einen gegenüber ΔA ver-

nachlässigbaren Fehler besitzt, skaliere man so, daß AA zeilenäquilibriert wird:

$$d_i := \max_k \sum_j \varDelta a_{kj} / \sum_j \varDelta a_{ij}$$

(R₄) Wenn die Matrix A und die rechte Seite **b** Meßfehler mit $|\delta A_M| \leq |A|$ und $|\delta b_M| \leq |A|$ aufweisen, die beide zum Fehlervektor

$$|\boldsymbol{z}| = A |\boldsymbol{x}| + A \boldsymbol{b} \tag{44}$$

beitragen, skaliere man so, daß Az äquilibriert wird:

$$d_i := \max_k \Delta z_k / \Delta z_i.$$

4.1.16. Bemerkung. (i) Die Regel (\mathbb{R}_1) ist lediglich in bezug auf die kollektiven Schranken (38), (39) optimal und berücksichtigt nicht die vom verwendeten Lösungsverfahren abhängige Struktur der individuellen äquivalenten Störungen $(\mathcal{O}A_R)_{ij}$. Sie braucht daher für ein konkretes Verfahren und eine konkrete Matrix nicht optimal zu sein; siehe auch Abschnitt 5.4. Da $\mathcal{O}A_R$ eine einzelne, feste Störung darstellt, wird die Abschätzung (39) auch dann meist zu pessimistisch sein, wenn (38) realistisch ist, d. h., wenn $||\mathcal{O}A_R|| \approx vF ||A||$ sein sollte. Aus diesem Grund ist (39) i. allg. nicht zur realen Abschätzung des erzeugten Rundungsfehlers geeignet. Dieser kann nach anderen Methoden besser abgeschätzt werden, etwa über das Residuum $r = \mathbf{b} - A\mathbf{x}$ von \mathbf{x} nach den im folgenden Teilabschnitt D angegebenen Methoden.

(ii) Durch den Fehlervektor $A\mathbf{z}$ gemäß (44) wird erfaßt, welchen Anteil die Fehler A_M und $\Delta \mathbf{b}_M$ zum Gesamtfehler liefern können. Abgesehen vom Sonderfall $\Delta A = \mathbf{0}$ — hier geht Regel (\mathbf{R}_4) in (\mathbf{R}_2) über — ist dazu i. allg. jedoch die Kenntnis von $|\mathbf{x}|$ erforderlich. Die für eine im Sinne von 4.1.10(ii) entsprechend (\mathbf{R}_4) optimale Skalierung ist daher erst a posteriori — d. h. nach Berechnung von \mathbf{x} — möglich und führt auf einen zweistufigen Prozeß:

Stufe 1: Skaliere A so, daß der bei Lösung von Ax = b erzeugte Rundungsfehler akzeptabel ist (also Skalierung nach (R_1) oder für viele praktisch auftretende Systeme keine Skalierung). Berechne x.

Stufe 2: Berechne $A\boldsymbol{z}$ und \boldsymbol{D} nach (R₄). Schätze $\|\boldsymbol{A}^{-1}\boldsymbol{D}^{-1}\|_{\infty}$ ab. Wegen $\|A\boldsymbol{z}\|_{\infty} = \|\boldsymbol{D}A\boldsymbol{z}\|_{\infty}$ folgt dann aus (32)

$$\|\boldsymbol{\delta}\boldsymbol{x}_{\boldsymbol{M}}'\|_{\infty} \leq \|\boldsymbol{A}^{-1}\boldsymbol{D}^{-1}\|_{\infty} \|\|\boldsymbol{\delta}\boldsymbol{z}\|_{\infty}$$

$$\tag{45}$$

als scharfe Abschätzung für den linearisierten Fehler $\boldsymbol{\sigma} \boldsymbol{x}'_{M}$. Division der Schranke aus (45) durch $1 - \|\boldsymbol{A}^{-1}\boldsymbol{D}^{-1}\|_{\infty}\|\boldsymbol{D} \boldsymbol{A}\boldsymbol{A}\|_{\infty}$ ergibt eine Schranke für $\|\boldsymbol{\sigma} \boldsymbol{x}_{M}\|_{\infty}$, siehe Ü 4.1.6 für $p = \infty$.

(iii) Bei Verwendung des Gaußschen Algorithmus erfordert die erste Stufe $\sim n^3/3$ opms, und die zweite kann mit $O(n^2)$ opms — also billig — realisiert werden, siehe Abschnitt 5.4 für Details.

(iv) Wenn die a-posteriori-Skalierung gemäß (ii) nicht zweckmäßig erscheint, sollte a priori nach (R_2) oder (R_3) skaliert werden. Falls Schätzungen für die Größenordnung der $|x_i|$ bekannt sind, kann (R_4) näherungsweise a priori ausgeführt werden. (v) Wenn die Störungen stochastischen Charakter haben, sollte der Fehlervektor Iz nicht gemäß (44), sondern nach der Vorschrift

$$(A\mathbf{z})_{i} := \left\{ \sum_{j=1}^{n} (\Delta a_{ij} x_{j})^{2} + \Delta b_{i}^{2} \right\}^{1/2}$$
(46)

festgelegt werden. \Box

Wir erwähnen abschließend, daß neben der Zeilenskalierung zusätzlich eine Spaltenskalierung gemäß

$$DAE(E^{-1}x) = DAE\hat{x} = Db, \qquad \hat{x} = E^{-1}x$$

$$(47)$$

mit $E = \text{diag}(e_i), e_i > 0$, durchgeführt werden kann. Dies entspricht einer Änderung des Maßstabes, in dem die Unbekannten x_j gemessen werden. Auf diese kompliziertere Skalierung soll hier jedoch nicht eingegangen werden, siehe B 4.3 für Literaturhinweise.

D. Residualkriterien

Wir kehren wieder zu den Schranken aus 4.1.3 zurück und bemerken, daß diese im Fall $\delta A = O$ durch untere Schranken ergänzt werden können: Wenn x bzw. $x + \delta x$ die Lösungen von

$$Ax = b$$
 bzw. $A(x + \delta x) = b + \delta b$

bezeichnen, gilt

$$\|\boldsymbol{\delta b}\|/\|\boldsymbol{A}\| \leq \|\boldsymbol{\delta x}\| \leq \|\boldsymbol{A}^{-1}\| \|\boldsymbol{\delta b}\|; \tag{48}$$

man beachte $\|\delta b\| = \|A\delta x\| \le \|A\| \|\delta x\|$. Wegen $\|b\|/\|A\| \le \|x\| = \|A^{-1}b\| \le \|A^{-1}\| \|b\|$ folgt hieraus im Fall $b \neq o$

$$\frac{1}{\text{cond } (A)} \frac{\|\boldsymbol{\delta b}\|}{\|\boldsymbol{b}\|} \leq \frac{\|\boldsymbol{\delta x}\|}{\|\boldsymbol{x}\|} \leq \text{cond } (A) \frac{\|\boldsymbol{\delta b}\|}{\|\boldsymbol{b}\|}.$$
(49)

Falls $\tilde{x} := x + \delta x$ gesetzt wird, stellt $-\delta b = b - A\tilde{x} =: r(\tilde{x})$ gerade das *Residuum* von \tilde{x} bezüglich des Systems Ax = b dar. In dieser Terminologie können (48), (49) wie folgt interpretiert werden.

4.1.17. Aussage. Es sei x die Lösung des regulären Gleichungssystems Ax = b. Dann gilt

$$\|\boldsymbol{b} - A\tilde{\boldsymbol{x}}\| / \|A\| \le \|\tilde{\boldsymbol{x}} - \boldsymbol{x}\| \le \|A^{-1}\| \|\boldsymbol{b} - A\tilde{\boldsymbol{x}}\|$$
(50)

für jedes \tilde{x} , und im Fall $b \neq o$ ist

$$\frac{1}{\operatorname{cond}(A)} \frac{\|\boldsymbol{b} - A\tilde{\boldsymbol{x}}\|}{\|\boldsymbol{b}\|} \le \frac{\|\tilde{\boldsymbol{x}} - \boldsymbol{x}\|}{\|\boldsymbol{x}\|} \le \operatorname{cond}(A) \frac{\|\boldsymbol{b} - A\boldsymbol{x}\|}{\|\boldsymbol{b}\|}.$$
(51)

Der relative Feh.er einer Näherung \tilde{x} für x kann also um den Faktor cond (A) nach oben oder unten von der relativen Residuumsgröße abweichen. Aus der Größe des Residuums, d. h. aus einer Einsetzprobe, kann daher i. allg. nicht auf die Größe

des Fehlers geschlossen werden, sondern es ist eine Abschätzung von $||A^{-1}||$ bzw. cond (A) erforderlich. Im Gegensatz dazu folgt jedoch aus der Kleinheit des Residuums $r(\tilde{x})$, daß \tilde{x} ein gegenüber Ax = b nur wenig gestörtes System löst und umgekehrt:

 $A \in \mathbb{R}^{n,n}, b$ $(A + \delta A) \, \tilde{x} = b + \delta b \quad \text{mit} \quad \|\delta A\| \leq AA, \quad \|\delta b\| \leq Ab$ gilt, wenn die Bedingung $\|r(\tilde{x})\| \leq \|b - A^{\tilde{\pi}^{||}}$ 4.1.18. Satz. Gegeben seien Zahlen $\Delta A, \Delta b \geq 0$ und $A \in \mathbb{R}^{n,n}, b, \tilde{x} \in \mathbb{R}^{n}$. Dann

(52)

$$\|\boldsymbol{r}(\tilde{\boldsymbol{x}})\| \leq \|\boldsymbol{b} - \boldsymbol{A}\tilde{\boldsymbol{x}}\| = \Delta \boldsymbol{b} + \Delta \boldsymbol{A} \|\tilde{\boldsymbol{x}}\|$$
(53)

erfüllt ist.

Beweis. Das System (52) wird genau dann durch \tilde{x} gelöst, wenn $r := b - A \tilde{x} = \delta A \tilde{x} - \delta b$ gilt. Mit den Bezeichnungen $\delta z := -r$ und $x := \tilde{x}$ kann Lemma 4.1.4 angewendet werden und besagt: Wenn δA und δb die Kugeln aus (52) durchlaufen, durchläuft r die Kugel $\|\mathbf{r}\| \leq \Delta \mathbf{b} + \Delta \mathbf{A} \|\hat{\mathbf{x}}\|$. Dies ist die zu beweisende Aussage.

Mit $\Delta A := \varepsilon \|A\|$, $\Delta b := \varepsilon \|b\|$ läßt sich 4.1.18 in relativen statt absoluten Größen formulieren.

4.1.19. Folgerung. Für $\varepsilon \ge 0$ und $A \in \mathbb{R}^{n,n}$, $b, \tilde{x} \in \mathbb{R}^n$ gibt es genau dann Störun-

gen δA , δb mit $(A + \delta A) \, \tilde{x} = b + \delta b$ und $\|\delta A\| \leq \varepsilon \|A\|$, $\|\delta b\| \leq \varepsilon \|b\|$, wenn $\|\Delta x\| < \varepsilon \|b\| + \|A\| \|\tilde{x}\|$ (54)

$$\|\boldsymbol{r}(\boldsymbol{\tilde{x}})\| = \|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{\tilde{x}}\| \le \varepsilon \{\|\boldsymbol{b}\| + \|\boldsymbol{A}\| \|\boldsymbol{\tilde{x}}\|\}$$
(55)
ilt.

g

Folgerung 4.1.19 besagt: Die Lösungen \tilde{x} der relativ wenig gestörten Gleichungen $(A + \delta A) \,\tilde{x} = b + \delta b$ sind durch die relative Kleinheit der entsprechenden Residuen $r(\tilde{x})$ gekennzeichnet.

4.1.20. Bemerkung. (i) Die praktische Bedeutung der Residualkriterien 4.1.18 und 4.1.19 besteht darin, daß durch numerische Berechnung von $\|\boldsymbol{r}(\boldsymbol{\tilde{x}})\|$ überprüft werden kann, ob \tilde{x} im Rahmen des Fehlerniveaus ΔA , Δb der Eingangsdaten als Lösung von Ax = b akzeptiert werden kann. Dabei müssen die bei der Berechnung auftretenden Rundungsfehler mit berücksichtigt werden. Man vergleiche dazu Abschnitt 5.4.

(ii) Die zu vorgegebenem $r(\tilde{x})$ gemäß 4.1.18 bzw. 4.1.4 konstruierten Störungen sind i. allg. nicht eindeutig festgelegt. So kann z. B. für die 2-Norm in 4.1.4 auch $-\delta A := \lambda(\Delta A) H$ gesetzt werden, wobei H eine Spiegelungsmatrix mit $Hx = ||x||_2 z$ ist. Eine solche Spiegelung läßt sich analog zu 3.3.4 konstruieren. Bei dieser Festlegung ist δA symmetrisch.

(iii) Zu 4.1.18 existiert das folgende elementweise Analogon: Gegeben seien $\Delta A \in \mathbb{R}^{n,n}, \ \Delta b \in \mathbb{R}^n \text{ mit } \Delta A \ge O, \ \Delta b \ge o \text{ und } A \in \mathbb{R}^{n,n}, \ b, \ \tilde{x} \in \mathbb{R}^n.$ Dann gibt es genau dann Störungen δA , δb mit

$$(A + \delta A) \,\tilde{\boldsymbol{x}} = \boldsymbol{b} + \delta \boldsymbol{b} \quad \text{und} \quad |\delta A| \leq \Delta A, \qquad |\delta \boldsymbol{b}| \leq \Delta \boldsymbol{b},$$
 (56)

wenn

$$|\boldsymbol{r}(\boldsymbol{\tilde{x}})| = |\boldsymbol{b} - \boldsymbol{A}\boldsymbol{\tilde{x}}| \le \Delta \boldsymbol{b} + \Delta \boldsymbol{A} |\boldsymbol{x}|$$
(57)

gilt. Im Sinne von

$$|\boldsymbol{\delta}\boldsymbol{A}| \leq \varepsilon \, |\boldsymbol{A}|, \qquad |\boldsymbol{\delta}\boldsymbol{b}| \leq \varepsilon \, |\boldsymbol{b}| \tag{58}$$

elementweise kleine Störungen sind daher gleichwertig zu einem im Sinne von

$$|\boldsymbol{r}(\tilde{\boldsymbol{x}})| = |\boldsymbol{b} - A\tilde{\boldsymbol{x}}| \leq \varepsilon \{|\boldsymbol{b}| + |\boldsymbol{A}| | \tilde{\boldsymbol{x}}|\}$$
(59)

kleinen Residuum. 🗌

Übungsaufgaben

Ü 4.1.1. Man zeige, daß unter der Voraussetzung $\|\boldsymbol{P}\| < 1$

$$(\boldsymbol{I}-\boldsymbol{P})^{-1} = \sum_{k=0}^{\infty} \boldsymbol{P}^{k} = \boldsymbol{I} + \boldsymbol{P} + \boldsymbol{P}^{2} + \dots + \boldsymbol{P}^{k} + \dots$$
(60)

gilt. Die unendliche Reihe (60) heißt Neumannsche Reihe von P.

Hinweis: Man beweise für die Teilsumme $S_l = \sum_{k=0}^{l} P^k$ die Identität

$$(I - P)^{-1} - S_l = (I - P)^{-1} P^{l+1}.$$
(61)

Ü 4.1.2. Man beweise unter Verwendung von (61) die Gültigkeit von

$$(A + \delta A)^{-1} = \sum_{k=0}^{l} (-A^{-1} \delta A)^k A^{-1} + O(||\delta A||^{l+1})$$
(62)

und gebe eine strenge Schranke für das Restglied an. Was ergibt sich für l = 1, 2, 3?

Ü 4.1.3. Man überlege sich, daß 4.1.2 und 4.1.3 auch für $\varkappa := ||A^{-1} \delta A||$ anstelle von $\varkappa = ||A^{-1}|| ||\delta A||$ gültig sind. Die Bedingung $||A^{-1} \delta A|| < 1$ ist schwächer als $||A^{-1}|| ||\delta A|| < 1$, aber schwieriger zu überprüfen.

Ü 4.1.4. Es sei $A = U\Sigma V^{\intercal}$ die Singulärwertzerlegung von A gemäß 1.2.14. Mit $c := U^{\intercal}b$ kann dann die Lösung von Ax = b in der Form $x = V\Sigma^{-1}c$ geschrieben werden. Man zeige, daß unter der Voraussetzung

$$(c_n)^2 \ge 0.01(c_1^2 + \dots + c_n^2)/n$$
 (63)

in der 2-Norm die Abschätzung

$$K_{\boldsymbol{b}}(\boldsymbol{A},\boldsymbol{b}) \leq 10 \sqrt[n]{n} \tag{64}$$

unabhängig von cond₂ (A) gilt, vgl. (16)

Ú 4.1.5. Für das Gleichungssystem
$$Ax = b$$
 mit $A = \begin{pmatrix} 1 & 1 \\ 1 & 0.99 \end{pmatrix}$, $b = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ ist $A^{-1} = \begin{pmatrix} -99 & 100 \\ 100 & -100 \end{pmatrix}$, $x = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ und cond_∞ $(A) = 400$. Für die Störungen
(i) $\delta A = 10^{-3} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix}$, $\delta b = 10^{-3} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$, (ii) $\delta A = 10^{-3} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$, $\delta b = 10^{-3} \begin{pmatrix} -1 \\ -1 \end{pmatrix}$

ist δx sowie $\|\delta x\|_{\infty}$ zu berechnen und mit den Schranken (13) zu vergleichen. Man begründe die unterschiedlichen Resultate.

Ü 4.1.6. Für eine der Normen p = 1 oder $p = \infty$ gelte $||A^{-1}||_p ||\delta A||_p < 1$. Mit den Bezeichnungen von 4.1.3 beweise man die Gültigkeit der Abschätzungen

$$|\delta \boldsymbol{x}| \leq (\boldsymbol{I} - |\boldsymbol{A}^{-1}| |\delta \boldsymbol{A}|)^{-1} |\boldsymbol{A}^{-1}| \{ |\delta \boldsymbol{A}| |\boldsymbol{x}| + |\delta \boldsymbol{b}| \}$$
(65)

und

$$\|\delta \boldsymbol{x}\|_{p} \leq \frac{\|\boldsymbol{A}^{-1}\|_{p}}{1 - \|\boldsymbol{A}^{-1}\|_{p} \|\delta \boldsymbol{A}\|_{p}} \||\delta \boldsymbol{A}| \|\boldsymbol{x}| + |\delta \boldsymbol{b}|\|_{p}.$$
(66)

Hin weis: Man gehe von $\delta x = (I + A^{-1} \delta A)^{-1} A^{-1} \{-\delta A x + \delta b\}$ aus und beachte $|(I - P)^{-1}| \le (I - |P|)^{-1}$ sowie die Monotonie und Absolutheit der verwendeten Normen.

Ü4.1.7. Man zeige am Beispiel

$$\boldsymbol{B} = \begin{pmatrix} 10 & 1 \\ 1 & 1 \end{pmatrix}, \quad \boldsymbol{\varDelta z} = \begin{pmatrix} 1 \\ 10 \end{pmatrix},$$

daß die Aussage (ii) aus Lemma 4.1.9 für nicht äquilibriertes Δz i. allg. falsch ist.

Ü 4.1.8. Man beweise die Äquivalenz aus Bemerkung 4.1.20 (iii) unter Verwendung von Lemma 4.1.8 analog zum Beweis von Satz 4.1.18.

Ú 4.1.9. Man zeige: Erfüllen die positiv definite Matrix $A \in S^{n,n}$ und die Störung $\delta A \in S^{n,n}$ die Bedingung (8), so ist $A + \delta A$ ebenfalls positiv definit.

Hinweis: Mit $A + \delta A$ ist nach 4.1.2 auch $B(\lambda) := A + \lambda \delta A$ für alle $\lambda \in [0, 1]$ regulär, d. h., die Eigenwerte von $B(\lambda)$ können für $\lambda \in [0, 1]$ das Vorzeichen nicht wechseln.

4.2. Prinzip der direkten Lösungsverfahren

Wir betrachten wieder das reguläre Gleichungssystem

 $A\boldsymbol{x} = \boldsymbol{b}.$ (1)

Das Prinzip der direkten Verfahren zur Lösung von (1) besteht darin, die Matrix A durch eine Folge regulärer elementarer Transformationsmatrizen auf obere Dreiecksform zu bringen. Wenn T das Produkt der Transformationsmatrizen bezeichnet, gilt also

$$TA = R$$
, T regulär, R obere Dreiecksmatrix. (2)

Im Kapitel 3 wurde bereits angedeutet, daß es aus Gründen der numerischen Stabilität in manchen Fällen nötig, in anderen Fällen zumindest günstig ist, die Matrix Adabei gewissen Zeilen- und gegebenenfalls auch Spaltenvertauschungen zu unterwerfen, d. h., statt (2) gilt

$$T\bar{A} = R$$
 mit $\bar{A} = P_z A P_s^{\dagger}$, R und T wie oben. (3)

Dabei sind P_z , P_s zwei Permutationsmatrizen, von denen P_z auf die Zeilen und P_s auf die Spalten von A wirkt. Offensichtlich stellt (3) eine spezielle Äquivalenztransformation von A in R dar. Wenn (3) nach \overline{A} aufgelöst wird, ergibt sich

$$\bar{A} = T^{-1}R = FR \quad \text{mit} \quad F := T^{-1}. \tag{4}$$

Praktisch wird T als Produkt von nichtorthogonalen elementaren Transformationsmatrizen, speziell als solches von LNT-Matrizen, oder als Produkt orthogonaler elementarer Transformationsmatrizen wie Householder-Spiegelungen oder Givens-Drehungen gewählt. Im ersten Fall spricht man vom *Gaußschen Algorithmus* und verwandten Verfahren, die Gegenstand der nachfolgenden Kapitel 5 und 6 sind. Als Produkt von LNT-Matrizen ist T ebenfalls eine untere Dreiecksmatrix, und dasselbe trifft für die Inverse $F = T^{-1}$ zu, die deshalb mit L bezeichnet wird. Es gilt dann

$$\bar{A} = LR, \qquad (5)$$

d. h., es liegt eine Dreiecks- oder LR-Faktorisierung vor.

Im zweiten Fall spricht man von Orthogonalisierungsverfahren, die in allgemeinerem Zusammenhang im Kapitel 10 behandelt werden. Als Produkt orthogonaler Matrizen ist T und damit auch F orthogonal, so daß F üblicherweise mit Q bezeichnet wird. Die Faktorisierung (4) ist dann von der Form

$$\boldsymbol{A} = \boldsymbol{Q}\boldsymbol{R},\tag{6}$$

sie heißt orthogonale Dreiecks- oder QR-Faktorisierung.

In der nach A aufgelösten Form lautet (3)

$$A = P_z^{\mathsf{T}} F R P_s, \tag{7}$$

so daß das Ausgangssystem zu

$$\boldsymbol{P}_{\boldsymbol{z}}^{\mathsf{T}}\boldsymbol{F}\boldsymbol{R}\boldsymbol{P}_{\boldsymbol{s}}\boldsymbol{x} = \boldsymbol{b} \tag{8}$$

äquivalent ist. Mit den Bezeichnungen

$$\bar{\boldsymbol{x}} = \boldsymbol{P}_s \boldsymbol{x}, \qquad \bar{\boldsymbol{b}} = \boldsymbol{P}_z \boldsymbol{b} \tag{9}$$

läßt sich (8) in der Form $FR\bar{x} = \bar{b}$ bzw. als

$$F\bar{c} = \bar{b} \quad \text{mit} \quad R\bar{x} = \bar{c}$$
 (10)

schreiben. Bei bekannter Faktorisierung (7) ist das Originalsystem (1) daher zu den beiden Gleichungssystemen (10) äquivalent, wobei \bar{x}, \bar{b} über (9) mit x, b verknüpft sind. Diese Äquivalenz ist der Ausgangspunkt für die später zu behandelnden direkten Lösungsverfahren und führt auf das nachfolgende algorithmische Schema.

4.2.1. Basisalgorithmus zur Lösung von Ax = b.

S1 (Faktorisierung von A): Bestimme Permutationsmatrizen P_z , P_s , eine obere Dreiecksmatrix R und eine untere Dreiecksmatrix oder orthogonale Matrix F mit

$$\bar{A} = P_z A P_s^{\mathsf{T}} = F R \tag{11}$$

S2 (Berechnung der zur rechten Seite **b** gehörenden Lösung x von Ax = b):

S2.1: Bestimme $\bar{b} := P_z b$ S2.2: Berechne \bar{c} als Lösung des linearen Gleichungssystems $F\bar{c} = \bar{b}$ S2.3: Berechne \bar{x} als Lösung des linearen Gleichungssystems $R\bar{x} = \bar{c}$ S2.4: Bestimme $x := P_s^T \bar{x}$

4.2.2. Bemerkung. (i) Unter Verwendung der Zerlegung (11) geht das Gleichungssystem

in

 $P_{a}^{\mathsf{T}}R^{\mathsf{T}}F^{\mathsf{T}}P_{a}y = f$

über und kann analog zum Schritt S2 wie folgt gelöst werden:

S3 (Berechnung der zur rechten Seite **d** gehörenden Lösung **y** von $A^{\mathsf{T}}y = d$)

S3.1: Bestimme $\bar{f} := P_{s}f$

 $A^{\mathsf{T}} y = f$

S3.2: Berechne \bar{g} als Lösung von $R^{\dagger}\bar{g} = \bar{f}$

- S3.3: Berechne \bar{y} als Lösung von $F^{\mathsf{T}}\bar{y} = \bar{y}$
- S3.4: Bestimme $\boldsymbol{y} := \boldsymbol{P}_{\boldsymbol{z}}^{\mathsf{T}} \boldsymbol{\bar{y}}$

(ii) Offensichtlich hängt Schritt S1 nur von der Matrix A ab, während die Schritte S2 bzw. S3 nur auf die in S1 berechneten Faktoren und Permutationen und auf die jeweiligen rechten Seiten zurückgreifen. Wenn mehrere Gleichungssysteme mit derselben Koeffizientenmatrix A oder A^{\dagger} zu lösen sind, braucht S1 nur einmal durchlaufen zu werden, und für jede rechte Seite genügt ein Durchlauf von S2 bzw. S3. Es ist daher zweckmäßig und heute üblich, den Faktorisierungsschritt S1 und den "Solver"-Schritt S2 bzw. S3 in getrennten Programmen zu realisieren. Man berücksichtige dabei, daß S1 für voll besetztes A i. allg. $O(n^3)$ opms, die Solver S2 bzw. S3 dagegen nur $O(n^2)$ opms erfordern.

(iii) Wenn die Faktorisierung S1 auf A^{\intercal} statt A angewendet wird, ergibt sich die Zerlegung $P_z A^{\mathsf{T}} P_s^{\mathsf{T}} = FR$, also

$$\hat{P}_z A \hat{P}_s^{\mathsf{T}} = \hat{L} \hat{F} \tag{12}$$

mit $\hat{P}_z := P_s$, $\hat{P}_s := P_z$, der unteren Dreiecksmatrix $\hat{L} := R^{\intercal}$ und der oberen Dreiecksmatrix oder orthogonalen Matrix $\hat{F} := F^{\intercal}$. Im Fall $\hat{F} = \hat{R}$ liegt wieder eine **LR**-Faktorisierung, im Fall $\hat{F} = \hat{Q}$ dagegen eine **LO**-Faktorisierung vor. Wir bemerken an dieser Stelle, da β (12) auch direkt durch zeilenweise Transformation von Aauf untere Dreiecksform durch Anwendung transponierter elementarer Transformationsmatrizen von rechts auf A berechnet werden kann. Dies ist z. B. bei linearen Quadratmittelproblemen mit Gleichheitsnebenbedingungen und allgemein in der Optimierung zweckmäßig, wo unterbestimmte lineare Gleichungssysteme zu lösen sind. \square

Auf die Realisierung der Permutationen sind wir bereits im Abschnitt 3.1 eingegangen. Die Lösung der in den Solvern S2, S3 auftretenden Gleichungssysteme mit orthogonalen oder Dreiecksmatrizen ist Gegenstand des folgenden Abschnitts.

Übungsaufgabe

Ü 4.2.1. Man beweise: Ist die reguläre Matrix A gemäß (11) faktorisiert, so sind F und R notwendig regulär, d. h., die Systeme in S2.2, S2.3 bzw. S3.2, S3.3 sind eindeutig lösbar.

4.3. Reguläre Gleichungssysteme einfacher Struktur

Unter regulären Gleichungssystemen einfacher Struktur verstehen wir Systeme mit regulären Dreiecksmatrizen oder mit orthogonalen Matrizen. Im folgenden befassen wir uns mit der Auflösung solcher Systeme und der Fehleranalyse der zugehörigen Algorithmen.

A. Gleichungssysteme mit regulären Dreiecksmatrizen

Reguläre Systeme mit Dreiecksmatrizen können in zwei zueinander symmetrischen Formen auftreten, nämlich als Systeme mit unteren Dreiecksmatrizen

$$l_{11}x_1 = b_1,$$

$$l_{21}x_1 + l_{22}x_2 = b_2, \quad l_{ii} \neq 0 \quad (i = 1, ..., n),$$

$$\vdots \qquad \vdots \qquad \vdots$$

$$l_{n1}x_1 + l_{n2}x_2 + \dots + l_{nn}x_n = b_n,$$
(1)

oder als Systeme mit oberen Dreiecksmatrizen

$$r_{11}x_{1} + r_{12}x_{2} + \dots + r_{1n}x_{n} = c_{1},$$

$$r_{22}x_{2} + \dots + r_{2n}x_{n} = c_{2}, \qquad r_{ii} \neq 0 \quad (i = 1, ..., n).$$

$$\vdots$$

$$r_{nn}x_{n} = c_{n},$$
(2)

Die Eingangsdaten sind $\{l_{ij}, b_i\}$ bzw. $\{r_{ij}, c_i\}$, als Ausgangsdaten treten die Komponenten $\{x_i\}$ des Lösungsvektors auf.

Durch die Umindizierung $i \leftrightarrow n - i + 1$, $j \leftrightarrow n - j + 1$ geht (2) in (1) über und umgekehrt, so daß wir uns im folgenden auf untere Dreieckssysteme (1) beschränken, die wir kurz als

$$Lx = b \tag{3}$$

schreiben. Als häufig auftretende Sonderfälle berücksichtigen wir dabei

- untere Einsdreiecksmatrizen: $l_{ii} = 1$ (i = 1, ..., n),
- untere Bidiagonalmatrizen: $l_{ij} = 0$ (i > j + 1),
- Diagonalmatrizen: $l_{ij} = 0$ (i > j).
- 10 Schwetlick, Numerische Algebra
Die zugehörigen Besetztheitsmuster sind für n = 4

Auf Grund der speziellen Gestalt des Systems (1) lassen sich die Unbekannten nacheinander in der Reihenfolge $x_1, x_2, ..., x_n$ gemäß

$$egin{aligned} &x_1 := b_1 / l_{11}\,, \ &x_2 := (b_2 - l_{21} * x_1) / l_{22}\,, \ &\vdots \end{aligned}$$

berechnen, was auf den folgenden Algorithmus führt:

4.3.1. Lösung eines Gleichungssystems mit regulärer unterer Dreiecksmatrix. Aufgabe: Für $\mathbf{L} = (l_{ij}) \in \Re^{n,n}$ mit $l_{ij} = 0$ (i < j) und $l_{ii} \neq 0$ sowie $\mathbf{b} = (b_i) \in \Re^n$ ist $\mathbf{x} = (x_i) \in \Re^n$ als Lösung von $\mathbf{L}\mathbf{x} = \mathbf{b}$ zu berechnen. Algorithmus: for i := 1(1)n do $x_i := \left(b_i - \sum_{j=1}^{i-1} l_{ij} * x_j\right) / l_{ii}$ Authwand $\cdot \sim n^{2}/2$ opms

Der Klammerausdruck in (4) soll in der Form

$$s := b_i$$

for $j := 1(1)i - 1$ do $s := s - l_{ij} * x_j$ (5)

ausgewertet werden. Eine Summe bzw. Laufanweisung, bei welcher die obere Grenze kleiner als die untere ist, wird dabei hier und im folgenden als leer angesehen.

Offensichtlich kann 4.3.1 in situ realisiert werden, d. h., der Vektor \boldsymbol{b} kann mit der Lösung \boldsymbol{x} überspeichert werden.

Das Zeichen "~" bedeutet bei Aufwandsangaben Gleichheit bis auf niedrigere Potenzen von *n*, d. h., $A \sim K_p n^p$ bedeutet $A = K_p n^p + K_{p-1} n^{p-1} + \cdots + K_0$ mit gewissen Konstanten K_0, \ldots, K_{p-1} .

4.3.2. Rundungsfehleranalyse. Algorithmus 4.3.1 werde für die dort spezifizierten Eingangsdaten $\{L, b\}$ durchgeführt. Dann genügt die berechnete Lösung x der Gleichung

$$(\boldsymbol{L} + \boldsymbol{\delta} \boldsymbol{L}) \, \boldsymbol{x} = \boldsymbol{b} \tag{6}$$

mit einer unteren Dreiecksmatrix $\boldsymbol{\delta L}$, die komponentenweise klein ist im Sinne von

$$|(\boldsymbol{\delta L})_{ij}| \leq vj |l_{ij}| \qquad (i, j = 1, \dots, n; \ i \geq j).$$

$$(7)$$

Be we is. Der Zähler z von (4) läßt sich bei Berechnung nach (5) in der Form $z = \text{fl}(z_0 - z_1 - \cdots - z_{i-1})$ mit $z_0 := b_i$ und $z_j := \text{fl}(l_{ij} * x_j)$ $(j = 1, \ldots, i-1)$ schreiben. Aus 2.3.3 folgt dann unter Beachtung von $z_j = l_{ij}x_j(1 + \delta_j)$ sofort $z = b_i(1 + \varepsilon_1) \cdots (1 + \varepsilon_{i-1}) - l_{i1}x_1(1 + \varepsilon_1) \cdots (1 + \varepsilon_{i-1}) (1 + \delta_1) - \cdots - l_{i,i-1}x_{i-1}(1 + \varepsilon_{i-1}) (1 + \delta_{i-1})$ mit $|\varepsilon_j|, |\delta_j| \leq \nu$. Berücksichtigung der Division durch l_{ii} führt auf

$$x_i = rac{z}{l_{ii}(1 + \mu_i)}, \quad |\mu_i| \leq v,$$

man vergleiche Ü 2.2.3. Einsetzen der Darstellung von z und Division von Zähler und Nenner durch $(1 + \varepsilon_1) \cdots (1 + \varepsilon_{i-1})$ ergibt schließlich

$$\begin{aligned} x_i &= \frac{b_i - l_{i1}x_1(1+\delta_1) - \dots - l_{i,i-1}x_{i-1} \frac{(1+\delta_{i-1})}{(1+\epsilon_1) \dots (1+\epsilon_{i-2})}}{l_{ii} \frac{(1+\mu_i)}{(1+\epsilon_1) \dots (1+\epsilon_{i-1})}} \\ &= \left(b_i - \sum_{j=1}^{i-1} l_{ij}(1+\beta_{ij}) x_j\right) \Big| (l_{ii}(1+\beta_{ii})) \end{aligned}$$

mit $|\beta_{ij}| \leq \nu j$. Wenn $(\delta L)_{ij} := l_{ij}\beta_{ij}$ gesetzt wird, ist dies gerade (6) und (7).

Man beachte, daß in (6) nur die Koeffizientenmatrix L mit den durch die Rundungsfehler hervorgerufenen Störungen belastet zu werden braucht.

Aus (7) folgt die kollektive Normabschätzung

$$\|\boldsymbol{\delta L}\|_{\infty} \leq \nu n \, \|\boldsymbol{L}\|_{\infty}.\tag{8}$$

Die Ergebnisse von 4.3.2 besagen: Algorithmus 4.3.1 ist numerisch gutartig in bezug auf L, und zwar elementweise mit F = n.

Der Fehler von x gegenüber der exakten Lösung $x^* = L^{-1}b$ kann entsprechend 4.1.14 in der Form

$$\frac{\|\boldsymbol{x} - \boldsymbol{x}^\star\|}{\|\boldsymbol{x}^\star\|} \leq \frac{\text{cond}\ (\boldsymbol{L})}{1 - (\|\boldsymbol{\delta}\boldsymbol{L}\|/\|\boldsymbol{L}\|) \text{ cond}\ (\boldsymbol{L})} \frac{\|\boldsymbol{\delta}\ \boldsymbol{L}\|}{\|\boldsymbol{L}\|}$$

abgeschätzt werden. In der ∞ -Norm gilt nach (8) $\|\boldsymbol{\delta L}\|_{\infty}/\|\boldsymbol{L}\|_{\infty} \leq \nu n$, so daß im Fall $\nu n \operatorname{cond}_{\infty}(\boldsymbol{L}) \ll 1$ mit einem relativen Fehler in der Größenordnung von

$$\operatorname{vn}\operatorname{cond}_{\infty}(L), \quad \operatorname{cond}_{\infty}(L) := \|L\|_{\infty} \|L^{-1}\|_{\infty}$$

$$\tag{9}$$

gerechnet werden muß. Praktische Erfahrungen zeigen jedoch, daß sich der tatsächliche relative Fehler fast immer in der Größenordnung

 $v\Phi(n)$

bewegt, wobei $\Phi(n)$ eine schwach mit n wachsende Funktion ist. Der gemäß (9) theoretisch mögliche Verstärkungsfaktor $\operatorname{cond}_{\infty}(L)$ tritt also fast nie auf, was auf das günstige individuelle, durch (8) nicht erfaßte Fehlerverhalten zurückzuführen ist.

Die bisherigen Überlegungen zeigen, daß Algorithmus 4.3.1 hinsichtlich des Aufwandes (proportional zur Anzahl der Eingangsdaten), des Speicherbedarfes (in-situ-Realisierung möglich) und des Fehlerverhaltens (extrem gutartig) als fast optimal anzusehen ist.

.

4.3.3. Bemerkung. (i) Wenn L speziell eine Einsdreiecks-, Bidiagonal- oder Diagonalmatrix darstellt, vereinfacht sich in 4.3.1 die Laufanweisung (4) zu

$$x_i := b_i - \sum_{j=1}^{i-1} l_{ij} * x_j$$

bzw.

 $x_i := (b_i - l_{i,i-1} * x_{i-1})/l_{ii}$

bzw.

$$x_i := b_i / l_{ii}.$$

Im Fall der Bidiagonal- bzw. Diagonalmatrix ist δL in 4.3.2 ebenfalls bidiagonal bzw. diagonal, und der Faktor j in (7) ist durch 2 bzw. 1 zu ersetzen.

(ii) Die Elemente von L werden in 4.3.1 bei Realisierung gemäß (4), (5) zeilenweise durchlaufen. Spaltenweise Abarbeitung ist ebenfalls möglich: Eine in-situ-Realisierung ist durch

for
$$j := 1(1)n$$
 do
 $\begin{vmatrix} b_j := b_j / l_{jj} \\ \text{for } i := j + 1(1)n \text{ do } b_i := b_i - l_{ij} * b_j \end{vmatrix}$
(10)

gegeben, wobei \boldsymbol{b} mit \boldsymbol{x} überspeichert wird.

(iii) Für das System (2) mit der oberen Dreiecksmatrix R lautet die zu (4), (5) analoge, zeilenweise arbeitende Vorschrift

for
$$i := n(-1)1$$
 do
 $\begin{vmatrix} s := c_i \\ \text{for } j := i + 1(1)n \text{ do } s := s - r_{ij} * x_j \\ x_i := s/r_{ii} \end{vmatrix}$

Die Unbekannten werden also in der inversen Reihenfolge $x_n, x_{n-1}, ..., x_1$ berechnet. Eine spaltenweise in-situ-Version entsprechend (10) ist

for
$$j := n(-1)1$$
 do
 $\begin{vmatrix} c_j := c_j / r_{jj} \\ \text{for } i := 1(1)j - 1$ do $c_i := c_i - r_{ij} * c_j \end{vmatrix}$
(11)

Die übrigen Aussagen können sinngemäß übernommen werden, insbesondere gilt für das berechnete \boldsymbol{x}

$$(\boldsymbol{R} + \boldsymbol{\delta}\boldsymbol{R}) \boldsymbol{x} = \boldsymbol{c} \quad \text{mit} \quad |(\boldsymbol{\delta}\boldsymbol{R})_{ij}| \leq \nu(n-j+1) |r_{ij}|.$$

B. Gleichungssysteme mit orthogonalen Matrizen

Wir betrachten jetzt ein Gleichungssystem

$$Qx = b \tag{12}$$

mit einer orthogonalen Matrix $Q \in \mathbb{R}^{n,n}$. Wegen der Orthogonalität von Q gilt

$$\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} = \boldsymbol{I} = \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}},\tag{13}$$

folglich

$$\boldsymbol{Q}^{-1} = \boldsymbol{Q}^{\mathsf{T}},\tag{14}$$

vgl. Abschnitt 1.1.H, so daß die Lösung $x^* = Q^{-1}b$ von (12) die Gestalt

$$x^* = Q^{\intercal} b$$

hat. Damit liegt bereits ein Verfahren zur Lösung von (12) vor.

4.3.4. Lösung eines Gleichungssystems mit einer orthogonalen Matrix. Aufgabe: Für $\mathbf{Q} = (q_{ij}) \in \Re^{n,n}$ und $\mathbf{b} = (b_i) \in \Re^n$ ist $\mathbf{x} = (x_i) \in \Re^n$ als Lösung von $\mathbf{Q}\mathbf{x} = \mathbf{b}$ zu bestimmen.

 $\mathcal{A} lgorithmus: \boldsymbol{x} := \boldsymbol{Q}^{\mathsf{T}} \boldsymbol{b}, \text{ d. h. } x_i := \sum_{j=1}^n q_{ji} * b_j \quad (i = 1, ..., n)$ $Autwand: \sim n^2 \text{ opms}$ (15)

Bei der Berechnung von x treten einmal bei der Ausführung der Matrix-Vektor Multiplikation (15) Rundungsfehler auf, vgl. Aufgabe 2 aus Abschnitt 2.3.D. Zum anderen wird die tätsächlich benutzte Matrix Q, deren Elemente Computerzahlen sind, i. allg. nicht orthogonal sein. Bereits die Computerdarstellung $Q = \text{rd}(Q^*)$ einer exakt orthogonalen Matrix $Q^* \in \mathbb{R}^{n,n}$ ist i. allg. nicht mehr orthogonal. In der numerischen linearen Algebra entsteht Q in der Regel als Ergebnis vorangegangener Rechnungen und erfüllt deshalb (13) und (14) erst recht nicht, so daß (15) auch bei fehlerfreier Realisierung nicht die exakte Lösung $x^* = Q^{-1}b$ liefern würde. Es ist daher nötig, die Abweichung der Matrix Q von der Orthogonalität geeignet zu berücksichtigen, etwa durch das Residuum der Gleichungen (13).

4.3.5. Rundungsfehleranalyse. Algorithmus 4.3.4 werde für die dort spezifizierten Eingangsdaten $\{Q, b\}$ durchgeführt, wobei Q die abgeschwächte Orthogonalitätsbedingung

$$\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} = \boldsymbol{I} + \boldsymbol{G} \quad \text{mit} \quad \|\boldsymbol{G}\|_2 = vM_1 < 1 \tag{16}$$

oder

$$QQ^{\intercal} = I + H \quad \text{mit} \quad \|H\|_2 = vM_2 < 1$$
(17)

erfüllen möge. Dann ist Q regulär, und das System Qx = b besitzt die eindeutige Lösung $x^* = Q^{-1}b$. Die nach 4.3.4 berechnete Lösung x genügt der Gleichung

$$Qx = b + \delta b, \tag{18}$$

wobei

$$\frac{\|\boldsymbol{\delta b}\|_{2}}{\|\boldsymbol{b}\|_{2}} \leq \nu \begin{cases} M_{1} \sqrt{\frac{1+\nu M_{1}}{1-\nu M_{1}}} + (1+\nu M_{1}) n^{3/2} & \text{im Fall (16),} \\ M_{2} + (1+\nu M_{2}) n^{3/2} & \text{im Fall (17).} \end{cases}$$
(19)

Der relative Fehler von x kann gemäß

$$\frac{\|\boldsymbol{\delta x}\|_{2}}{\|\boldsymbol{x^{*}}\|_{2}} \leq \nu \begin{cases} M_{1} + (1 + \nu M_{1}) n^{3/2} & \text{im Fall (16),} \\ M_{2} \sqrt{\frac{1 + \nu M_{2}}{1 - \nu M_{2}}} + (1 + \nu M_{2}) n^{3/2} & \text{im Fall (17)} \end{cases}$$
(20)

mit $\delta x := x - x^*$ abgeschätzt werden.

Beweis. Die Matrix Q erfülle die Bedingung (16). Dann gilt

$$\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q}\|_{2} = \|\boldsymbol{I} + \boldsymbol{G}\|_{2} \le \|\boldsymbol{I}\|_{2} + \|\boldsymbol{G}\|_{2} \le 1 + \nu M_{1}.$$

Nach Ü 1.2.9 ist $\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q}\|_{2} = \|\boldsymbol{Q}\|_{2}^{2} = \|\boldsymbol{Q}^{\mathsf{T}}\|_{2}^{2}$, so daß

$$\|\boldsymbol{Q}\|_2 \leq \sqrt{1 + \nu M_1} \tag{21}$$

folgt. Wegen $||G||_2 < 1$ und 4.1.1 ist I + G, also $Q^{\mathsf{T}}Q$ regulär mit $||(Q^{\mathsf{T}}Q)^{-1}|| \leq 1/(1 - ||G||_2)$ $\leq 1/(1 - \nu M_1)$. Dann ist auch Q regulär, siehe (1.2.28), und wie oben folgt $||(Q^{\mathsf{T}}Q)^{-1}||_2$ $= ||Q^{-1}||_2^2$, mithin

$$\|Q^{-1}\|_{2} \leq 1/\sqrt{1 - \nu M_{1}}.$$
(22)

Wir betrachten nun die durch (15) erzeugten Rundungsfehler. Nach (2.3.37) genügt der berechnete Vektor $x := Q^{\mathsf{T}} b$ der Gleichung

$$\boldsymbol{x} = (\boldsymbol{Q}^{\mathsf{T}} + \boldsymbol{\delta} \boldsymbol{Q}^{\mathsf{T}}) \boldsymbol{b} \quad \text{mit} \quad |\boldsymbol{\delta} \boldsymbol{Q}| \leq \boldsymbol{\nu} \boldsymbol{n} |\boldsymbol{Q}|.$$
⁽²³⁾

Analog zum Beweis von 2.3.15 folgt aus der letzten Abschätzung

$$\|\partial Q\|_2 \leq m \sqrt{n} \|Q\|_2 = m^{3/2} \|Q\|_2. \tag{24}$$

Multiplikation von (23) mit Q liefert

 $\boldsymbol{Q} \boldsymbol{x} = (\boldsymbol{Q} \boldsymbol{Q}^{\mathsf{T}} + \boldsymbol{Q} \boldsymbol{\delta} \boldsymbol{Q}^{\mathsf{T}}) \, \boldsymbol{b}$.

Das erste Q auf der rechten Seite kann wegen (16) durch $Q = Q^{-T} + Q^{-T}G$ ersetzt werden, womit sich

$$Qx = (I + Q^{-\intercal}GQ^{\intercal} + Q\delta Q^{\intercal}) b = b + \delta b$$

mit $\delta \boldsymbol{b} := (\boldsymbol{Q}^{-\intercal} \boldsymbol{G} \boldsymbol{Q}^{\intercal} + \boldsymbol{Q} \delta \boldsymbol{Q}^{\intercal}) \boldsymbol{b}$, wegen (21), (22) und (24) also

$$\begin{split} \| \boldsymbol{\delta b} \|_2 &\leq (\| \boldsymbol{Q}^{-1} \|_2 \, \| \boldsymbol{G} \|_2 \, \| \boldsymbol{Q} \|_2 + \| \boldsymbol{Q} \|_2 \, \| \boldsymbol{\delta Q} \|_2) \, \| \boldsymbol{b} \|_2 \\ &\leq \nu \left\{ M_1 \, \sqrt{1 + \nu M_1} / \sqrt{1 - \nu M_1} \, + \, (1 \, + \, \nu M_1) \, n^{3/2} \right\} \| \boldsymbol{b} \|_2 \end{split}$$

ergibt. Zur Abschätzung des relativen Fehlers von x setzen wir $b = Qx^*$ in (23) ein und erhalten unter Berücksichtigung von (16)

folglich

$$egin{aligned} m{x} &= (m{Q}^{\mathsf{T}}m{Q} + m{\delta}m{Q}^{\mathsf{T}}m{Q}) \, m{x}^{*} &= (m{I} + m{G} + m{\delta}m{Q}^{\mathsf{T}}m{Q}) \, m{x}^{*}, \ \|m{x} &= (m{I} - m{x}^{*}\|_{2} &\leq (\|m{G}\|_{2} + \|m{\delta}m{Q}\|_{2} \, \|m{Q}\|_{2}) \, \|m{x}^{*}\|_{2} \ &\leq
u\{M_{1} + (1 +
uM_{1}) \, n^{3/2}\} \, \|m{x}^{*}\|_{2}. \end{aligned}$$

Dabei sind (19) und (20) im Fall (16) bewiesen. Für den Fall (17) kann analog vorgegangen werden.

Satz 4.3.5 läßt sich wie folgt interpretieren: Wenn M_i nicht zu groß ist — etwa in der Größenordnung n oder $n^{3/2}$ —, also $vM_i \ll 1$ gilt und Q nur wenig von einer orthogonalen Matrix abweicht, ist Algorithmus 4.3.5 numerisch gutartig in bezug auf **b** mit $F = M_i + n^{3/2}$, und der relative Fehler von x ist durch vF beschränkt, also klein. Weicht Q dagegen stark von der Orthogonalität ab, d. h., ist M_i sehr groß, so können die Schranken (19) und (20) nicht mehr als akzeptabel angesehen werden. Der relative Fehler von x kann sehr groß werden, und Algorithmus 4.3.4 ist dann zur Lösung von Qx = b nicht mehr geeignet. 4.3.6. Bemerkung. (i) Unter der schärferen Voraussetzung

$$\nu M_i \leq 1/2 \tag{25}$$

können die Abschätzungen (19) und (20) durch die gröberen, aber einfacheren Schranken

$$\frac{\|\delta \boldsymbol{b}\|_2}{\|\boldsymbol{b}\|_2}, \frac{\|\delta \boldsymbol{x}\|_2}{\|\boldsymbol{x}\|_2} \leqq \nu \left\{ \sqrt{3} \ M_i + 1.5 n^{3/2} \right\}$$
(26)

ersetzt werden, wobei i = 1 bzw. i = 2 für (16) bzw. (17) ist.

(ii) Wie weit eine Abweichung von der Orthogonalität bei der Matrix Q toleriert werden kann, hängt auch von der Genauigkeit der rechten Seite b ab. Wenn diese einen Meßfehler δb_M mit dem Fehlerniveau

$$\|\boldsymbol{\delta} \boldsymbol{b}_M\|_2 / \|\boldsymbol{b}\|_2 = \imath L \leq 1/2$$

aufweist, kann eine Abweichung νM_i derselben Größenordnung toleriert werden, da sie die Genauigkeit von x nicht wesentlich verschlechtert.

(iii) Wenn die durch rM_i charakterisierte Qualität der Orthogonalität von Q nicht a priori bekannt ist, müßte etwa $||G||_2$ durch $||Q^{\mathsf{T}}Q - I||_F$ abgeschätzt werden, um den Vergleich von M_i mit L in (ii) durchführen zu können. Dies erfordert $O(n^3)$ Operationen. Unter der Voraussetzung (25), d. h. bei möglicherweise starker, aber nicht katastrophaler Abweichung von der Orthogonalität, gilt jedoch cond₂ (Q) $\leq \sqrt{3}$, so daß nach 4.1.17 die a-posteriori-Abschätzung

$$\frac{\|\boldsymbol{x}^* - \boldsymbol{x}\|_2}{\|\boldsymbol{x}^*\|_2} \le \sqrt{3} \frac{\|\boldsymbol{b} - \boldsymbol{Q}\boldsymbol{x}\|_2}{\|\boldsymbol{b}\|_2}$$
(27)

für das gemäß 4.3.4 berechnete x gilt. Aus dem nachträglich mit $O(n^2)$ Operationen berechneten Residuum $\|b - Qx\|_2$ läßt sich also auf die Güte von x schließen. \Box

4.3.7. Bemerkung. In vielen Anwendungen — vgl. 10.2 — ist Q nicht explizit bekannt, sondern z. B. implizit als Produkt einer Folge elementarer orthogonaler Transformationsmatrizen definiert, so daß Algorithmus 4.3.4 nicht direkt angewendet werden kann. Die Abschätzungen aus 4.3.5 und 4.3.6 bleiben jedoch gültig, wenn $x = Q^{\mathsf{T}}b$ nach irgendeiner Vorschrift berechnet wird, für die

$$\boldsymbol{x} = \boldsymbol{Q}^{\mathsf{T}}\boldsymbol{b} + \boldsymbol{\delta}_1 \boldsymbol{x}, \qquad \|\boldsymbol{\delta}_1 \boldsymbol{x}\|_2 \leq \nu K \|\boldsymbol{Q}\|_2 \|\boldsymbol{b}\|_2 \tag{28}$$

mit einer (akzeptablen) Kumulationskonstanten K gilt. In den Abschätzungen ist dann nur $n^{3/2}$ durch K zu ersetzen.

Übungsaufgaben

Ü 4.3.1. Man gebe eine zeilen- bzw. spaltenweise Realisierung von (15) zur Berechnung von $x = Q^{\mathsf{T}} b$ an.

Ü 4.3.2. Für die gemäß 4.3.4 berechnete Lösung x von Qx = b werde das Residuum r := b- Qx ebenfalls in der durch v charakterisierten Arithmetik berechnet. Man zeige, daß die berechneten Werte unter der Voraussetzung (25) der Abschätzung

$$\|m{r}\|_2/\|m{b}\|_2 \leqq
u igg \{ \sqrt{3} \ M_i + 3n^{3/2} igg\}$$

genügen.

Ü 4.3.3. Wenn die Matrix $Q = (q^1, ..., q^n)$ mit den Spalten q^j stärker von der Orthogonalität abweicht, kann nach BJÖRCK [67] das folgende Verfahren zur Lösung von Qx = b verwendet werden:

$$\begin{aligned}
\mathbf{b}^{1} &:= \mathbf{b} \\
\text{for } i &:= 1(1)n \quad \text{do} \\
& \left| \begin{array}{c} x_{i} &:= \mathbf{q}^{i\mathsf{T}}\mathbf{b}^{i} \\
\mathbf{b}^{i+1} &:= \mathbf{b}^{i} - x_{i} * \mathbf{q}^{i} \\
\end{aligned} \right.
\end{aligned} \tag{29}$$

Als Teil eines speziellen Verfahrens — nämlich des modifizierten Gram-Schmidt-Verfahrens, siehe Abschnitt 10.1 — ergibt diese Vorschrift eine wesentlich höhere Genauigkeit als das Verfahren 4.3.4, d. h. als die Vorschrift

(30)

for i := 1(1)n do $x_i := q^{i \top} b$.

Man beweise, daß (29) und (30) bei exakter Realisierung und exakt orthogonalem Q identisch sind und charakterisiere den Rechenaufwand und Speicherbedarf von (29) im Vergleich zu (30) bei in-situ-Realisierung, d. h., wenn **b** mit den **b**ⁱ überspeichert wird. Man überlege sich ferner, daß **b**ⁿ⁺¹ für beliebiges Q bei exakter Realisierung von (29) das Residuum **b** – Qxdarstellt.

Bemerkungen zum Kapitel 4

B 4.1. Der Begriff der Kondition hat eine lange Tradition in der Mathematik. Im Zusammenhang mit der Matrixinversion und der Lösung linearer Gleichungssysteme ist die Zahl cond (A) wohl erstmals von VON NEUMANN/GOLDSTINE [47] und TURING [48] verwendet worden. Die Ergebnisse der Störungstheorie sind in allen Lehrbüchern zu finden, allerdings meist mit der oberen Schranke cond (A) für $K_b(A, b)$ in (4.1.14). Die so vergröberte Abschätzung ist dann nicht mehr für jedes $\{A, b\}$ scharf für $||\delta A|| \rightarrow 0$.

B 4.2. Die Schärfe der (linearisierten) Normschranken (4.1.13), (4.1.14) hat VAN DER SLUIS [70] untersucht, wo auch gewisse spaltenweise Konditionszahlen betrachtet werden. Individuelle Abschätzungen sind bei BAUER [66] zu finden, siehe auch MILLER [75]. Ähnliche Probleme behandelt SKEEL [81]. Der Frage, welche Norm zur Abschätzung gewählt werden sollte, wenden sich KAHAN [66] und FADDEEV/FADDEEVA [70] zu. Die LINPACK-Autoren DONGARRA et al. [79] empfehlen eine Zeilen- und Spaltenskalierung, welche die Elemente von ΔA in etwa dieselbe Größenordnung bringt. Im Fall $\Delta b = o$ ist dies mit (\mathbb{R}_4) aus 4.1.12 identisch, denn dann ist $\Delta z = \Delta A |\mathbf{x}|$ für jedes \mathbf{x} äquilibriert und (30) realistisch. Die Festlegung von Dmit dem Ziel der Äquilibrierung von Δz scheint neu zu sein.

B 4.3. Das Problem der Konditionsminimierung durch geeignete Skalierung ist in einer Reihe von Arbeiten behandelt worden; wir zitieren FORSYTH/STRAUS [55], BAUER [63, 66, 69] – dort ist Satz 4.1.14 zu finden –, HEINRICH [63], BUSINGER [68], VAN DER SLUIS [69] und FENNER/LOIZOU [77]. In diesen Arbeiten wird auch auf andere als ∞ -Normen und simultane Zeilen- und Spaltenskalierung nach (4.1.17) eingegangen. Zum Beispiel wird cond (\overline{A}) minimal über $\overline{A} = DAE$, wenn \overline{A} und \overline{A}^{-1} zeilenäquilibriert sind, allerdings ist dann die Bestimmung der optimalen Skalierungsfaktoren in D und E nicht mehr so einfach wie im Fall der Zeilenskalierung (BAUER [66]).

B 4.4. Die Abschnitte 4.2 und 4.3 folgen in ihren Grundideen den Darstellungen von WILKINson [61, 63, 65].

5. Der Gaußsche Algorithmus

Der Gaußsche Algorithmus – lange vor GAUSS bekannt, siehe B 5.1 – war und ist *das* Standardverfahren zur Lösung linearer Gleichungssysteme. Seine Einfachheit erlaubt die Vermittlung und Anwendung bereits in der Schule, aber auch reale Gleichungssysteme mit Hunderten und Tausenden von Unbekannten werden heute nach Varianten dieses klassischen Verfahrens gelöst.

5.1. Die Grundform des Gaußschen Algorithmus

Zu lösen ist das lineare Gleichungssystem

$$a_{11}x_1 + \cdots + a_{1n}x_n = b_1,$$

$$\vdots \qquad \vdots \qquad \vdots$$

$$a_{n1}x_1 + \cdots + a_{nn}x_n = b_n,$$

kurz

$$A\boldsymbol{x} = \boldsymbol{b}, \tag{1}$$

mit der regulären Koeffizientenmatrix $A = (a_{ij}) \in \mathbb{R}^{n,n}$ und der rechten Seite $\mathbf{b} = (b_i) \in \mathbb{R}^n$.

A. Der Gaußsche Algorithmus als Verfahren der schrittweisen Elimination

Das Prinzip des Gaußschen Algorithmus ist die schrittweise Elimination der Unbekannten durch geeignete Kombinationen der Gleichungen und damit die Überführung des Systems (1) in ein äquivalentes System

$$\boldsymbol{R}\boldsymbol{x} = \boldsymbol{c} \tag{2}$$

von oberer Dreiecksform.

Im ersten Schritt wird x_1 aus den Gleichungen i = 2, ..., n eliminiert, indem ein Vielfaches der ersten Gleichung zur *i*-ten addiert wird, so daß dort der Koeffizient bei x_1 verschwindet. Dies ist gleichbedeutend mit dem Auflösen der ersten Gleichung nach x_1 und Einsetzen in die übrigen Gleichungen. Wenn die erste Spalte von Amit a bezeichnet wird und die restlichen zur Matrix B zusammengefaßt werden, liegt gerade der in 3.2.5 behandelte Eliminationsschritt mittels einer LNT-Matrix vor.

Für die Durchführbarkeit dieses Schrittes muß das Pivotelement a_{11} von 0 verschieden sein. Da auch für eine reguläre Matrix $a_{11} = 0$ gelten kann, wird ein nichtverschwindender Koeffizient $a_{s(1),1}$, $1 \leq s(1) \leq n$, in der ersten Spalte von A gesucht und durch Vertauschen der Zeilen 1 und s(1) des Systems in die Position (1, 1) gebracht, d. h., von dem Originalsystem $A^{(1)}x = b^1$, $A^{(1)} := A$, $b^1 := b$, wird zu dem zeilenvertauschten System $\hat{A}^{(1)}x = \hat{b}^1$ übergegangen, wobei jetzt

$$\hat{a}_{11}^{(1)} = a_{s(1),1}^{(1)} \neq 0$$

 $\hat{T} = T (\hat{I})$

gilt. Dieser Übergang ist stets möglich, denn andernfalls hätte A nur Nullen in der ersten Spalte und wäre damit singulär im Widerspruch zur Voraussetzung.

Die gesuchte Transformation ist dann nach 3.2.5 durch die LNT-Matrix

mit

$$\hat{L}_{1} := L_{1}(-\hat{l}^{1}), \qquad \hat{l}^{1} := (0, \hat{l}_{21}, ..., \hat{l}_{n1})^{\mathsf{T}}$$

$$\hat{l}_{i1} := \hat{a}_{i1}^{(1)}/\hat{a}_{11}^{(1)} \qquad (i = 2, ..., n)$$
(3)

7 \Τ

eindeutig festgelegt. Das transformierte System $A^{(2)}x = b^2$ hat die Koeffizienten

$$A^{(2)} = \hat{L}_1 \hat{A}^{(1)} = \begin{pmatrix} \frac{a_{11}^{(2)} \mid a_{12}^{(2)} \dots a_{1n}^{(2)}}{0 \mid a_{22}^{(2)} \dots a_{2n}^{(2)}} \\ \vdots \mid \vdots \mid \vdots \\ 0 \mid a_{n2}^{(2)} \dots a_{nn}^{(2)} \end{pmatrix} = : \begin{pmatrix} \frac{r_{11} \mid r_{12} \dots r_{1n}}{0 \mid M^{(2)}} \\ \end{pmatrix},$$

$$b^2 = \hat{L}_1 \hat{b}^1 = \begin{pmatrix} \frac{b_1^{(2)}}{b_2^{(2)}} \\ \vdots \\ b_n^{(2)} \end{pmatrix} = : \begin{pmatrix} \frac{c_1}{d^2} \\ \end{pmatrix},$$

wobei nur die Elemente von $M^{(2)}$ und d^2 neu berechnet zu werden brauchen:

$$a_{ij}^{(2)} := \hat{a}_{ij}^{(1)} - \hat{l}_{i1} * \hat{a}_{1j}^{(1)}, \qquad b_i^2 := \hat{b}_i^1 - \hat{l}_{i1} * \hat{b}_1^1 \qquad (i, j = 2, ..., n).$$
(4)

Wenn noch $M^{(1)} := A$, $d^1 := b$ gesetzt wird, ist also das Ausgangsproblem

$$\boldsymbol{M}^{(1)}\boldsymbol{x} = \boldsymbol{d}^{1} \tag{5}$$

äquivalent zum System $A^{(2)}x = b^2$, das in die erste Zeile

$$r_{11}x_1 + r_{12}x_2 + \dots + r_{1n}x_n = c_1 \tag{6}$$

und das Restsystem

$$M^{(2)}x^2 = d^2 \tag{7}$$

der Ordnung n-1 für $x^2 := (x_2, \ldots, x_n)^{\mathsf{T}}$ zerfällt. Da sich beim Übergang von $A^{(1)}$ zu $\hat{A}^{(1)}$ höchstens das Vorzeichen der Determinante ändert und $\hat{A}^{(1)}$ dieselbe Determinante wie $A^{(2)}$ besitzt — siehe Abschnitt 1.1.G —, gilt

$$\det (A^{(2)}) = r_{11} \det (M^{(2)}) = \pm \det (A^{(1)}), \tag{8}$$

so daß $M^{(2)}$ wie $A^{(1)}$ regulär sein muß. Bei bekannter Lösung x^2 von (7) kann die noch fehlende Komponente x_1 sofort aus (6) berechnet werden. Damit ist die Lösung des Gleichungssystems (5) der Ordnung n auf die Lösung des Gleichungssystems (7) der Ordnung n-1 zurückgeführt worden.

Im zweiten Schritt wird derselbe Eliminationsproze β mit n-1 statt n auf das reduzierte System (7) angewendet usw. Auf diese Weise entsteht eine Folge von Gleichungssystemen

$$A^{(k+1)}x = b^{k+1} \qquad (k = 1, ..., n-1)$$
(9)

 \mathbf{mit}

$$\boldsymbol{A}^{(k)} = \begin{pmatrix} a_{11}^{(k)} \dots a_{1,k-1}^{(k)} & a_{1k}^{(k)} \dots a_{1n}^{(k)} \\ \vdots & \vdots & \vdots \\ a_{k-1,k-1}^{(k)} & a_{k-1,k}^{(k)} \dots a_{k-1,n}^{(k)} \\ \hline & & \\ \boldsymbol{O} & \begin{vmatrix} a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & \vdots & \vdots \\ a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \\ \end{vmatrix} = : \begin{pmatrix} \boldsymbol{O} & \boldsymbol{M}^{(k)} \\ \hline & \boldsymbol{M}^{(k)} \end{pmatrix},$$

$$\boldsymbol{b}^{k} = \begin{pmatrix} b_{1}^{k} \\ \vdots \\ b_{k-1}^{k} \\ \vdots \\ b_{n}^{k} \end{pmatrix} = : \begin{pmatrix} \boldsymbol{C}^{k} \\ \boldsymbol{d}^{k} \end{pmatrix},$$
(10)

d. h., $A^{(k)}$ ist bereits bis zur Spalte k - 1 von oberer Dreiecksform. Im k-ten Eliminationsschritt werden dann die gewünschten Nullen in der k-ten Spalte von $A^{(k+1)}$ erzeugt, indem der beschriebene "erste Schritt" auf das Restsystem $M^{(k)}x^k = d^k$, $x^k := (x_k, ..., x_n)^{\mathsf{T}}$ der Ordnung n - k + 1 angewendet wird. Sinngemäße Übertragung der Formeln (3) und (4) liefert die folgende Darstellung des Gesamtprozesses:

5.1.1. Gaußscher Algorithmus. Initialisierung: $A^{(1)} := A$, $b^1 := b$ for k := 1(1)n - 1 do k-ter Eliminationsschritt: Pivotsuche: Finde Index s(k), $k \leq s(k) \leq n$, mit $a^{(k)}_{s(k),k} \neq 0$

Vertauschung: Festlegung von $\hat{A}^{(k)}$, \hat{b}^k aus $A^{(k)}$, b^k durch Vertauschen der Zeilen k und s(k)

Berechnung der Eliminationskoeffizienten:

$$\hat{l}_{ik} := \hat{a}_{ik}^{(k)} / \hat{a}_{kk}^{(k)}$$
 $(i = k + 1, ..., n)$

Transformation der Matrix:

$$a_{ij}^{(k+1)} := egin{cases} 0 & ext{für} \quad i=k+1,...,n, j=k, \ \hat{a}_{ij}^{(k)} - \hat{l}_{ik} st \hat{a}_{kj}^{(k)} & ext{für} \quad i,j=k+1,...,n, \ \hat{a}_{ij}^{(k)} & ext{sonst} \end{cases}$$

Transformation der rechten Seite:

$$b_{i}^{k+1} := \begin{cases} \hat{b}_{i}^{k} - \hat{l}_{ik} * \hat{b}_{k}^{k} & \text{für} \quad i = k+1, ..., n, \\ \hat{b}_{i}^{k} & \text{sonst} \end{cases}$$

Aufwand: Berechnung von $\mathbf{R} = A^{(n)}$: $\sim n^3/3$ opms; Berechnung von $\mathbf{c} = \mathbf{b}^n$: $\sim n^2/2$ opms.

Zur Erläuterung soll das folgende einfache Beispiel dienen: 5.1.2. Beispiel. Gegeben ist das Gleichungssystem Ax = b mit

$$A = A^{(1)} = \begin{pmatrix} 0 & 2 & -1 & -2 \\ 2 & -2 & 4 & -1 \\ 1 & 1 & 1 & 1 \\ -2 & 1 & -2 & 1 \end{pmatrix}, \quad b = b^{1} = \begin{pmatrix} -7 \\ 6 \\ 10 \\ -2 \end{pmatrix}.$$

Erster Schritt (k = 1): s(1) = 2, d. h. Vertauschen der Zeilen 1 und 2

$$\hat{A}^{(1)} = egin{pmatrix} 2 & -2 & 4 & -1 \ 0 & 2 & -1 & -2 \ 1 & 1 & 1 & 1 \ -2 & 1 & -2 & 1 \end{pmatrix}, \quad \hat{b}^1 = egin{pmatrix} 6 \ -7 \ 10 \ -2 \end{pmatrix}.$$

Mit $\hat{l}_{21} = 0/2 = 0$, $\hat{l}_{31} = 1/2 = 0.5$, $\hat{l}_{41} = -2/2 = -1$ folgt

$$A^{(2)} = egin{pmatrix} 2 & -2 & 4 & -1 \ \overline{0} & 2 & -1 & -2 \ 0 & 2 & -1 & 1.5 \ 0 & -1 & 2 & 0 \end{pmatrix}$$
, $b^2 = egin{pmatrix} 6 \ -7 \ 7 \ 4 \end{pmatrix}$.

Zweiter Schritt (k = 2): s(2) = 2, d. h. $\hat{A}^{(2)} := A^{(2)}$, $\hat{b}^2 := b^2$. Mit $\hat{l}_{32} = 2/2 = 1$, $\hat{l}_{42} = -1/2 = -0.5$ folgt

$$A^{(3)} = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & 0 & 3.5 \\ 0 & 0 & 1.5 & -1 \end{pmatrix}, \quad b^3 = \begin{pmatrix} 6 \\ -7 \\ 14 \\ 0.5 \end{pmatrix}.$$

Dritter Schritt (k = 3): s(3) = 4, d. h. Vertauschen der Zeilen 3 und 4

$$\hat{A}^{(3)} = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & 1.5 & -1 \\ 0 & 0 & 0 & 3.5 \end{pmatrix}, \quad \hat{b}^{3} = \begin{pmatrix} 6 \\ -7 \\ 0.5 \\ 14 \end{pmatrix}.$$

Mit $\hat{l}_{43} = 0/1.5 = 0$ ergibt sich $A^{(4)} = \hat{A}^{(3)}, b^{(4)} = \hat{b}^{(3)}$, also

$$\boldsymbol{R} = \boldsymbol{A}^{(4)} = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & 1.5 & -1 \\ 0 & 0 & 0 & 3.5 \end{pmatrix}, \quad \boldsymbol{c} = \boldsymbol{b}^4 = \begin{pmatrix} 6 \\ -7 \\ 0.5 \\ 14 \end{pmatrix}.$$

Der Lösungsvektor $x = (1, 2, 3, 4)^{\mathsf{T}}$ ergibt sich gemäß (4.3.11) aus dem Gleichungssystem $\mathbf{R}x = \mathbf{c}$ mit der oberen Dreiecksmatrix \mathbf{R} .

Das Beispiel macht deutlich, daß der Gaußsche Algorithmus in situ realisiert werden kann: A kann mit $\hat{A}^{(k)}$ und $A^{(k+1)}$, b mit \hat{b}^k und b^{k+1} überschrieben werden, und die Eliminationskoeffizienten \hat{l}_{ik} können auf dem Platz der im k-ten Schritt erzeugten Nullen, d. h. in den Positionen (i, k) von A gespeichert werden

5.1.3. Aussage. Der Gaußsche Algorithmus 5.1.1 ist in exakter Arithmetik für jede reguläre Matrix $A \in \mathbf{R}^{n,n}$ und jede rechte Seite $b \in \mathbf{R}^n$ durchführbar und liefert nach n-1 Schritten das System $A^{(n)}x = b^n$, wobei

$$A^{(n)} = \begin{pmatrix} \hat{a}_{11}^{(1)} & \dots & \hat{a}_{1n}^{(1)} \\ \hat{a}_{22}^{(2)} & \dots & \hat{a}_{2n}^{(2)} \\ \vdots \\ \hat{a}_{n-1,n-1}^{(n-1)} & \hat{a}_{n-1,n}^{(n-1)} \\ O & a_{nn}^{(n)} \end{pmatrix} = : \mathbf{R} = (r_{ij}), \quad \mathbf{b}^{n} = \begin{pmatrix} \hat{b}_{1}^{1} \\ \hat{b}_{2}^{2} \\ \vdots \\ \hat{b}_{n-1}^{n-1} \\ b_{n}^{n} \end{pmatrix} = : \mathbf{c} = (c_{i})$$

und

 $r_{kk} = a_{s(k),n}^{(k)} \neq 0$ $(k = 1, ..., n), \quad s(n) := n,$ (11)ist.

Dies folgt unmittelbar aus den Überlegungen bei der Herleitung des Verfahrens.

B. Matrixformulierung des Gaußschen Algorithmus

Unter Verwendung der in den Abschnitten 3.1 und 3.2 eingeführten elementaren Transformationsmatrizen läßt sich der Gaußsche Algorithmus auch in Matrixnotation formulieren. Der im k-ten Eliminationsschritt vorzunehmende Übergang von $A^{(k)}$ zu $\hat{A}^{(k)}$ erfolgt gemäß

$$\hat{A}^{(k)} := T_{k,s(k)} A^{(k)} \tag{12}$$

und der von $\hat{A}^{(k)}$ zu $A^{(k+1)}$ nach

$$A^{(k+1)} = \begin{pmatrix} \mathbf{R}^{(k)} \\ \mathbf{O} \mid \hat{\mathbf{L}}_{1}^{(n-k+1)} \hat{\mathbf{M}}^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{I}^{(k-1)} \mid \mathbf{O} \\ \mathbf{O} \mid \hat{\mathbf{L}}_{1}^{(n-k+1)} \end{pmatrix} \hat{\mathbf{A}}^{(k)} =: \hat{\mathbf{L}}_{k} \hat{\mathbf{A}}^{(k)}.$$
(13)

Dabei ist $\hat{M}^{(k)}$ diejenige Teilmatrix von $\hat{A}^{(k)}$, der $M^{(k)}$ in $A^{(k)}$ entspricht, d. h., die durch Zeilenvertauschung aus dieser hervorgeht, und es wurde

$$\hat{L}_{1}^{(n-k+1)} := L_{1}^{(n-k+1)}(-\hat{l}^{1,n-k+1}), \quad \hat{l}^{1,n-k+1} := (0, \, \hat{l}_{k+1,k}, \, \dots, \, \hat{l}_{nk})^{\mathsf{T}} \in \mathbf{R}^{n-k+1}$$

sowie

$$\hat{\boldsymbol{L}}_{\boldsymbol{k}} := \boldsymbol{L}_{\boldsymbol{k}}(-\hat{\boldsymbol{l}}^{\boldsymbol{k}}), \qquad \hat{\boldsymbol{l}}^{\boldsymbol{k}} := (0, ..., 0, \hat{\boldsymbol{l}}_{\boldsymbol{k}+1,\boldsymbol{k}}, ..., \hat{\boldsymbol{l}}_{\boldsymbol{n},\boldsymbol{k}})^{\mathsf{T}} \in \mathsf{R}^{\boldsymbol{n}}$$
(14)

gesetzt, vgl. (3.2.6).

Für den k-ten Eliminationsschritt ergibt sich damit

$$A^{(k+1)} := \hat{L}_k T_{k,s(k)} A^{(k)}, \qquad b^k := \hat{L}_k T_{k,s(k)} b^k \qquad (k = 1, ..., n - 1), \quad (15)$$

und durch (14), (15) ist die gewünschte Matrixdarstellung des Verfahrens gegeben.

C. Interpretation als Dreiecksfaktorisierung

Hintereinanderschaltung der Transformationen (15) führt auf

$$A^{(n)} = \hat{L}_{n-1}T_{n-1,s(n-1)}A^{(n-1)} = \hat{L}_{n-1}T_{n-1,s(n-1)}\hat{L}_{n-2}T_{n-2,s(n-2)}A^{(n-2)} = \cdots,$$

wegen $\boldsymbol{R} = A^{(n)}$ und $A = A^{(1)}$ also auf

$$\mathbf{R} = \mathbf{\hat{L}}_{n-1} \mathbf{T}_{n-1,s(n-1)} \cdots \mathbf{\hat{L}}_2 \mathbf{T}_{2,s(2)} \mathbf{\hat{L}}_1 \mathbf{T}_{1,s(1)} \mathbf{A}.$$
 (16)

Nach Ü 3.2.2 gilt wegen $1 < 2 \leq s(2) \leq n$

$$T_{2,s(2)}L_1(-\hat{l}^1) T_{1,s(1)} = L_1(-T_{2,s(2)}\hat{l}^1) T_{2,s(2)}T_{1,s(1)},$$

wegen $1 < 2 < 3 \leq s(3) \leq n$ analog

$$\begin{split} & [\boldsymbol{T}_{3,s(3)}\boldsymbol{L}_{2}(-\boldsymbol{l}^{2})] [\boldsymbol{T}_{2,s(2)}\boldsymbol{L}_{1}(-\boldsymbol{l}^{1}) \ \boldsymbol{T}_{1,s(1)}] \\ &= [\boldsymbol{L}_{2}(-\boldsymbol{T}_{3,s(3)}\boldsymbol{\hat{l}}^{2}) \ \boldsymbol{T}_{3,s(3)}] [\boldsymbol{L}_{1}(-\boldsymbol{T}_{2,s(2)}\boldsymbol{\hat{l}}^{1}) \ \boldsymbol{T}_{2,s(2)}\boldsymbol{T}_{1,s(1)}] \\ &= \boldsymbol{L}_{2}(-\boldsymbol{T}_{3,s(3)}\boldsymbol{\hat{l}}^{2}) \ \boldsymbol{L}_{1}(-\boldsymbol{T}_{3,s(3)}\boldsymbol{T}_{2,s(2)}\boldsymbol{\hat{l}}^{1}) \ \boldsymbol{T}_{3,s(3)}\boldsymbol{T}_{2,s(2)}\boldsymbol{T}_{1,s(1)}] \end{split}$$

usw., d. h., die $T_{k,s(k)}$ können nach rechts gebracht werden, wenn sie gleichzeitig auf die Argumentvektoren \hat{l}_j der L_j mit j < k angewendet werden. Damit geht (16) in

$$\boldsymbol{R} = \boldsymbol{L}_{n-1} \boldsymbol{L}_{n-2} \cdots \boldsymbol{L}_2 \boldsymbol{L}_1 \boldsymbol{P}_0 \boldsymbol{A} \tag{17}$$

über, wobei

$$\boldsymbol{L}_{k} := \boldsymbol{L}_{k}(-\boldsymbol{l}^{k}), \quad \boldsymbol{l}^{k} := (0, ..., 0, l_{k+1,k}, ..., l_{n,k})^{\mathsf{T}} := \boldsymbol{P}_{k} \boldsymbol{\hat{l}}^{k} \qquad (k = 1, ..., n)$$
(18)

und

$$\boldsymbol{P}_{k} := \begin{cases} \boldsymbol{T}_{n-1.s(n-1)} \cdots \boldsymbol{T}_{k+1.s(k+1)} & \text{für} \quad k = 0, 1, \dots, n-2, \\ \boldsymbol{I} & & \text{für} \quad k = n-1 \end{cases}$$
(19)

ist. Nach (3.2.3) ist jedes L_k regulär mit der Inversen $L_k^{-1} = L_k(l^k)$, so daß (17) mit $L_1^{-1} \cdots L_{n-1}^{-1}$ durchmultipliziert werden kann und auf

$$\boldsymbol{P}\boldsymbol{A} = \boldsymbol{L}_{1}(\boldsymbol{l}^{1})\cdots\boldsymbol{L}_{n-1}(\boldsymbol{l}^{n-1})\boldsymbol{R} = \boldsymbol{L}\boldsymbol{R}, \qquad \boldsymbol{P} := \boldsymbol{P}_{0}$$

$$(20)$$

führt. Das Produkt $L = L_1(l^1) \cdots L_{n-1}(l^{n-1})$ ist nach 3.2.2 eine untere Einsdreiecksmatrix mit den Komponenten l_{ik} von l^k als Elementen der k-ten Spalte. Damit kann (20) wie folgt interpretiert werden:

5.1.4. Satz. Für jede reguläre Matrix $A \in \mathbb{R}^{n,n}$ liefert der Gaußsche Algorithmus in exakter Arithmetik eine LR-Faktorisierung von A gemäß

mit $r_{kk} \neq 0$ (k = 1, ..., n). Dabei ist $\mathbf{R} = A^{(n)}$ die durch 5.1.1 erzeugte obere Dreiecksmatrix, \mathbf{L} die gemäß (18) bis (20) festgelegte Matrix der Eliminationskoeffizienten und $\mathbf{P} = \mathbf{T}_{n-1.s(n-1)} \cdots \mathbf{T}_{1.s(1)}$ die durch die Vertauschungen bei der Pivotsuche charakterisierte Permutation der Zeilen von \mathbf{A} .

Die Faktorisierung (21) ist durch P und A eindeutig festgelegt.

Beweis. Es bleibt die Eindeutigkeit der Faktorisierung (21) zu zeigen: Es gelte also PA= $LR = \tilde{L}\tilde{R}$ mit einer weiteren unteren Einsdreiecksmatrix \tilde{L} und einer (notwendig regulären) oberen Dreiecksmatrix \tilde{R} . Dann folgt $\tilde{L}^{-1}L = \tilde{R}R^{-1} =: D$. Als Produkt unterer Dreiecksmatrizen ist $\tilde{L}^{-1}L$ eine untere Dreiecksmatrix, und ebenso ist $\tilde{R}R^{-1}$ als Produkt oberer Dreiecksmatrizen von oberer Dreiecksform. Dann muß D diagonal sein und wegen $L = \tilde{L}D$ Diagonalelemente $d_i = 1$ besitzen. Also ist D = I und folglich $L = \tilde{L}, R = \tilde{R}$. \Box

Die Permutationsmatrix **P** ist i. allg. natürlich nicht eindeutig festgelegt, denn in 5.1.1 gibt es i. allg. mehrere als Pivot geeignete Elemente $a_{s(k),k}^{(k)}$, $k \leq s(k) \leq n$.

5.1.5. Bemerkung. Die Vorschrift (18) besagt: Der Vektor l^k entsteht aus dem im k-ten Eliminationsschritt benutzten Vektor \hat{l}^k der Eliminationskoeffizienten, indem sämtliche nachfolgenden Vertauschungen $T_{k+1,s(k+1)}, \ldots, T_{n-1,s(n-1)}$ auf letzteren angewendet werden. Bei in-situ-Realisierung von 5.1.1 und Abspeichern der \hat{l}_{ik} auf dem Platz von a_{ik} bedeutet dies, daß die vollen Zeilen von A einschließlich der bereits berechneten Eliminationskoeffizienten vertauscht werden.

5.1.6. Beispiel. Für die in 5.1.2 betrachtete Matrix ergibt sich mit $P = T_{3,s(3)}T_{2,s(2)}T_{1,s(1)}$ = $T_{34}T_{12}$ die Faktorisierung

$$\mathbf{PA} = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ -2 & 1 & -2 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & & & \\ 0 & 1 & & \\ -1 & -0.5 & 1 & \\ 0.5 & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & -2 & 4 & -1 \\ 2 & -1 & -2 \\ & 1.5 & -1 \\ & & 3.5 \end{pmatrix} = \mathbf{LR}. \square$$

Bei bekannter Dreiecksfaktorisierung (21) läßt sich die Lösung x des Systems (1) nach dem im Abschnitt 4.2 dargestellten Schema berechnen.

5.1.7. Lösung von Ax = b bei bekannter LR-Faktorisierung. Es sei PA = LR die LR-Faktorisierung der regulären Matrix $A \in \mathbb{R}^{n,n}$ gemäß 5.1.4. Dann ergibt sich die Lösung x des Gleichungssystems Ax = b für jedes $b \in \mathbb{R}^{n}$ wie folgt:

S2.1: Bestimme $\mathbf{\tilde{b}} = \mathbf{Pb}$

S2.2: Berechne c als Lösung von $Lc = \overline{b}$

S2.3: Berechne x als Lösung von Rx = c

Aufwand: $\sim n^2$ opms

5.1.8. Bemerkung. (i) Die Schritte S2.1 und S2.2 heißen Vorwärtselimination. Der dabei berechnete Vektor c ist identisch mit dem Vektor $c = b^n$, der sich nach dem Gaußschen Algorithmus 5.1.1 aus der rechten Seite b ergibt, siehe Ü 5.1.2. Der Schritt S2.3 heißt Rücksubstitution.

(ii) Zur Lösung von $A^{\mathsf{T}}y = d$ unter Verwendung von (21) siehe Ü 5.1.3.

(iii) Wegen 5.1.4 und 5.1.7 ist es zweckmäßig, bei der Lösung von (1) in zwei

Schritten vorzugehen und diese auch in zwei getrennten Programmen zu realisieren, vgl. Bemerkung 4.2.2(ii):

Schritt 1: Faktorisierung von A, d. h. Durchführung des Gaußschen Algorithmus nur für die Matrix A mit dem Ergebnis P, L, R, wobei P durch die Zahlen $\{s(1), \ldots, s(n-1)\}$ repräsentiert wird.

Schritt 2: Lösung des Systems (1) bei vorliegender Faktorisierung von A nach 5.1.7. Schritt 1 erfordert $\sim n^3/3$ opms, Schritt 2 dagegen nur $\sim n^2$ opms, so daß das beschriebene Vorgehen besonders günstig ist, wenn mehrere Gleichungssysteme mit derselben Koeffizientenmatrix zu lösen sind. Der aufwendige Faktorisierungsschritt braucht dann nur einmal durchgeführt zu werden.

(iv) Wenn $D = \text{diag}(r_{kk})$ gesetzt wird, kann (21) in der Form

$$PA = LR = LD(D^{-1}R) = LDR_1 = L_1R_1$$
(22)

geschrieben werden, wobei $\mathbf{R}_1 := \mathbf{D}^{-1}\mathbf{R}$ eine obere Einsdreiecksmatrix und $\mathbf{L}_1 := \mathbf{L}\mathbf{D}$ eine untere Dreiecksmatrix mit den Diagonalelementen r_{ii} bezeichnet. Die Faktorisierung $\mathbf{P}\mathbf{A} = \mathbf{L}_1\mathbf{R}_1$ kann auch direkt berechnet werden, indem beim Gaußschen Algorithmus die Pivotzeile zu Beginn des k-ten Schrittes durch das Pivotelement dividiert wird. Dies ist dann auch für k = n erforderlich. \Box

Aus der Faktorisierung (21) läßt sich in einfacher Weise die Determinante von A berechnen.

5.1.9. Aussage. Es sei $A \in \mathbb{R}^{n,n}$ regulär mit der LR-Faktorisierung (21) gemäß 5.1.4. Dann gilt

$$\det (A) = \det (P) \prod_{k=1}^{n} r_{kk} = (-1)^{\mu} \prod_{k=1}^{n} r_{kk}, \qquad (23)$$

wobei μ die Anzahl der nichttrivialen Vertauschungen in $P = T_{n-1,s(n-1)} \cdots T_{1,s(1)}$, d. h. die Anzahl der $k \in \{1, ..., n-1\}$ mit k < s(k) bezeichnet.

Beweis. Aus $A = P^{\mathsf{T}}LR$ folgt nach den Determinantenregeln aus 1.1.G sofort det (A)= det (P) det (L) det (R) = det $(P) \cdot 1 \cdot r_{11} \cdots r_{nn}$. Nun ist

$$\det \left(\boldsymbol{T}_{k,s(k)} \right) = \begin{cases} 1 & \text{für } k = s(k), \\ -1 & \text{für } k < s(k) \end{cases}$$

nach (3.1.8), womit sich sofort (23) ergibt.

5.1.10. Bemerkung. Der Gaußsche Algorithmus 5.1.1 gestattet also für jede Matrix $A \in \mathbb{R}^{n,n}$ die effektive Berechnung von det (A): Für reguläres A ist er bis zum (n-1)-ten Schritt durchführbar und liefert det (A) gemäß (23) als Produkt der Pivots mit eventueller Vorzeichenumkehr. Für singuläres A bricht der Gaußsche Algorithmus ab, weil in einem der Schritte k = 1, ..., n - 1 kein von 0 verschiedenes Pivot gefunden werden kann, bzw. es ist $r_{nn} = 0$. In diesem Fall ist det (A) = 0 zu setzen. Bei der numerischen Realisierung der Formel (23) kann leicht Überbzw. Unterlauf eintreten, so daß hier besondere Vorsichtsmaßnahmen erforderlich sind, siehe 5.2.5.

Übungsaufgaben

Ü 5.1.1. Man zeige: Wenn A symmetrisch ist und s(1) = 1 gewählt werden kann, ist auch $M^{(1)}$ symmetrisch.

Ü 5.1.2. Man weise nach, daß der in 5.1.7 berechnete Vektor $c = L^{-1}Pb$ identisch mit b^n aus 5.1.1 ist.

Ü 5.1.3. Man überlege sich, daß das Gleichungssystem $A^{\intercal}y = d$ unter Verwendung der LR-Faktorisierung (21) gemäß

S3.2: Berechne g als Lösung von $R^{\intercal}g = d$

S3.3: Berechne \overline{y} als Lösung von $L^{\intercal}\overline{y} = g$

S3.4: Bestimme $\boldsymbol{y} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{\bar{y}}$

gelöst werden kann; vgl. auch Bemerkung 4.2.2(i).

5.2. Pivotisierung

Die Bedingung $a_{kk}^{(k)} \neq 0$ reicht i. allg. nicht aus, um die numerische Stabilität des Gaußschen Algorithmus zu garantieren, siehe Beispiel 2.3.20. Bei der Rundungsfehleranalyse im Abschnitt 5.3 werden wir sehen, daß numerische Stabilität vorliegt, wenn die Matrizen $M^{(k)}$ gegenüber $M^{(1)} = A$ nicht beliebig groß werden. Für gewisse Matrizenklassen ist diese Bedingung auch ohne Pivotisierung — d. h. bei der Wahl des Diagonalelementes $a_{kk}^{(k)}$ als Pivot — erfüllt. Außer in diesen Spezialfällen muß jedoch durch geeignete Stabilisierungsmaßnahmen ein zu starkes Anwachsen der $M^{(k)}$ verhindert werden. Dies wird durch verschiedene Pivotisierungsstrategien erreicht.

A. Spaltenpivotisierung

Wenn im k-ten Schritt das Pivot als ein betragsgrößtes Element der ersten Spalte von $M^{(k)}$, also in der k-ten Spalte von $A^{(k)}$ gesucht wird, spricht man von Spaltenpivotisierung. Es gilt dann

$$|\hat{a}_{kk}^{(k)}| = |a_{s(k),k}^{(k)}| = \max\{|a_{ik}^{(k)}|: i = k, ..., n\}$$
(1)

und somit

$$|\hat{l}_{ik}| = |\hat{a}_{ik}^{(k)}/\hat{a}_{kk}^{(k)}| \le 1 \qquad (i = k + 1, ..., n),$$
(2)

d. h., die Transformationsmatrix $\hat{L}_k = L_k(-\hat{l}^k)$ wird eine SLNT-Matrix im Sinne von 3.2.8. Aus (2) folgt

$$|a_{ij}^{(k+1)}| = |\hat{a}_{ij}^{(k)} - \hat{l}_{ik} * \hat{a}_{kj}^{(k)}| \le |\hat{a}_{ij}^{(k)}| + |\hat{a}_{kj}^{(k)}| \qquad (i, j = k+1, ..., n),$$
(3)

mit $M^{(k)} = (m^{kk}, ..., m^{kn})$ also

$$\|m{m}^{k+1,j}\|_{\infty} \leq 2 \, \|m{m}^{k,j}\|_{\infty} \quad (j=k+1,...,n) \quad ext{und} \quad \|m{M}^{(k+1)}\|_{\infty} \leq 2 \, \|m{M}^{(k)}\|_{\infty}.$$

11 Schwetlick, Numerische Algebra

Durch Induktion über k ergeben sich daraus die folgenden Abschätzungen:

5.2.1. Aussage. Im Fall der Spaltenpivotisierung (1) gilt für k = 2, ..., n

$$\|\boldsymbol{m}^{k,j}\|_{\infty} \leq \gamma_k \|\boldsymbol{a}^j\|_{\infty}, \quad \|\boldsymbol{M}^{(k)}\|_{\infty} \leq \gamma_k \|\boldsymbol{A}\|_{\infty}$$

$$\gamma_k := 2^{k-1}. \qquad (5)$$

mit

Dabei bezeichnet a^{j} die *j*-te Spalte von A.

5.2.2. Bemerkung. (i) Es gibt Matrizen A, für die 5.2.1 scharf ist, d. h., für die sich der ungünstige Verstärkungsfaktor 2^{k-1} tatsächlich einstellt, siehe Ü 5.2.2.

(ii) Umfangreiche Beobachtungen haben gezeigt, daß bei praktisch vorkommenden Gleichungssystemen die Abschätzung (4) mit

$$\gamma_k \le \gamma \approx 10 \tag{6}$$

erfüllt ist, d. h., das exponentielle Anwachsen tritt praktisch nicht auf. 🗌

Auf Grund des im Sinne von (6) günstigen realen Verhaltens kann die Spaltenpivotisierung trotz der pessimistischen strengen Schranke (5) als i. allg. ausreichende Stabilisierung des Gaußschen Algorithmus empfohlen werden. Sie wird in fast allen derzeit bekannten Implementierungen verwendet. Wir geben deshalb eine computernahe Beschreibung der für die Lösung von Ax = b nach dem Gaußschen Algorithmus mit Spaltenpivotisierung wesentlichen Teilschritte an und gehen dabei grundsätzlich spaltenweise vor, da zweidimensionale Felder z. B. in FORTRAN spaltenweise abgespeichert werden, vgl. Kapitel 16.

5.2.3. Dreiecksfaktorisierung einer Matrix A mit Spaltenpivotisierung.

Aufgabe: Für $A = (a_{ij}) \in \Re^{n,n}$ ist nach dem Gaußschen Algorithmus mit Spaltenpivotisierung eine Dreiecksfaktorisierung PA = LR gemäß 5.1.4 zu berechnen. Die Matrix A ist mit den signifikanten Elementen von L und R zu überspeichern; die Permutationsmatrix **P** ist durch die Zahlen $\{s(1), \ldots, s(n-1)\}$ zu charakterisieren. Falls in einem Eliminationsschritt k kein $r_{kk} \neq 0$ gefunden werden kann oder wenn $r_{nn} = 0$ ist, soll abgebrochen und ie = -1 gesetzt werden, andernfalls wird ie = 0 gesetzt, und **R** ist dann regulär.

Algorithmus:

ie := 0for k := 1(1)n - 1 do Pivotsuche: z := 0for i := k(1)n do if $|a_{ik}| > z$ then $[z := |a_{ik}|, s := i]$ if z = 0 then [ie := -1, stop] s(k) := sVertauschung: if k < s then for j := 1(1)n do $[z := a_{kj}, a_{kj} := a_{sj}, a_{sj} := z]$ Berechnung der \hat{l}_{ik} : for i := k + 1(1)n do $a_{ik} := a_{ik}/a_{kk}$ Spaltenweise Berechnung von $M^{(k+1)}$: for j := k + 1(1)n do for i := k + 1(1)n do $a_{ii} := a_{ii} - a_{ik} * a_{ki}$

Test von r_{nn} : if $a_{nn} = 0$ then [ie := -1, stop]

Aufwand: $\sim n^3/3$ opms, $\sim n$ S (für $s(1), \ldots, s(n-1)$)

Bei der Realisierung als SUBROUTINE ist die stop-Anweisung durch RETURN zu ersetzen.

5.2.4. Lösung des Gleichungssystems $Ax = P^{\intercal}LRx = b$.

Autgabe: Für gegebene Dreiecksfaktorisierung PA = LR gemäß 5.2.3 mit regulärem R und gegebene rechte Seite $b = (b_i) \in \Re^n$ ist die Lösung $x \in \Re^n$ des Gleichungssystems Ax = b gemäß 5.1.7 zu berechnen und auf dem Platz von b zu speichern.

Algorithmus:

Aufwand: $\sim n^2$ opms

Vor Lösung von Ax = b mittels 5.2.4 muß abgefragt werden, ob 5.2.3 mit ie = 0 ausgeführt worden ist.

Der Vollständigkeit halber geben wir noch ein Programm zur Berechnung der Determinante aus der Dreiecksfaktorisierung an, wobei wegen der Über- bzw. Unterlaufgefahr eine nichtstandardisierte Darstellung von det (A) verwendet wird. Der Algorithmus arbeitet auch, wenn 5.2.3 mit ie = -1 abgebrochen worden ist, vgl. Bemerkung 5.1.10.

5.2.5. Berechnung von det $(\mathbf{A}) = \det (\mathbf{P}^{\mathsf{T}} \mathbf{L} \mathbf{R}).$

Aufgabe: Aus den Ausgangsdaten von 5.2.3 ist det (A) im Fall ie = 0 gemäß 5.1.9 in der Form det $(A) = d \times 16^e$ mit $1/16 \le |d| < 1$ und ganzzahligem Exponenten e zu berechnen; im Fall ie = -1, d. h. für numerisch singuläres A, ist d = 0, e := 0 zu setzen.

```
\begin{array}{l} Algorithmus:\\ e:=0\\ \text{if } ie=0\\ \text{then } & \left|\begin{array}{c} d:=1\\ \text{for } k:=1(1)n \text{ do}\\ & \left|\begin{array}{c} d:=d*a_{kk}, \text{ if } k < n \text{ then } [\text{if } k < s(k) \text{ then } d:=-d]\\ & M1: \text{ if } |d| < 1/16 \text{ then } [d:=d*16, e:=e-1, \text{ goto } M1]\\ & M2: \text{ if } |d| \geq 1 \\ & \text{then } d:=d/16, e:=e+1, \text{ goto } M2] \end{array}\right|\\ else \quad d:=0 \end{array}
```

Aufwand: $\sim n \text{ opm} + \text{Aufwand für Normalisierung von } d$

5.2.6. Bemerkung. (i) Statt der expliziten Dreiecksfaktorisierung (5.1.20) kann auch die implizite Darstellung (5.1.16) für den Lösungsprozeß verwendet werden. Die Elemente \hat{l}_{ik} von \hat{l}^k werden dann im unteren Dreieck von A gespeichert und nicht mit vertauscht, vgl. 5.1.5. In der Vertauschungslaufanweisung von 5.2.3 muß dabei $,j := 1(1)n^{\prime\prime}$ durch $,j := k(1)n^{\prime\prime}$ ersetzt werden. Der Vektor $c = b^n$ wird dann gemäß (5.1.15) berechnet, d. h., in 5.2.4 sind S2.1 und S2.2 durch

$$\begin{array}{l} \text{S2.1}+2: \text{ for } k:=1(1)n-1 \text{ do} \\ & \left| \begin{array}{c} z:=b_k, b_k:=b_{s(k)}, b_{s(k)}:=z \\ \text{ for } i:=k+1(1)n \text{ do } b_i:=b_i-a_{ik}*b_i \end{array} \right| \end{array} \right.$$

zu ersetzen. Dadurch wird der Vertauschungsaufwand in 5.2.3 etwa halbiert.

(ii) Durch Einführung eines Permutationsvektors $\{is(1), \ldots, is(n)\}$ kann die explizite physische Vertauschung in 5.2.3 und 5.2.4 umgangen werden, indem statt auf die *i*-te Zeile auf die Zeile is(i) zurückgegriffen wird. Die Anfangsbelegung ist is(i) := i $(i = 1, \ldots, n)$, und im k-ten Schritt werden is(k) und is(s) vertauscht, wenn s = s(k) der Zeilenindex des Pivotelements ist. Dies ist zweckmäßig, wenn A extern gespeichert oder wie im Fall schwach besetzter Matrizen eine physische Umordnung aufwendig ist, vgl. Abschnitt 6.4. \Box

B. Zeilenpivotisierung

Bisher wurde das Pivot in der ersten Spalte von $M^{(k)}$, d. h. in der k-ten Spalte von $A^{(k)}$ gesucht. Analog ist natürlich auch eine Pivotsuche in der ersten Zeile von $M^{(k)}$, d. h. in der k-ten Zeile von $A^{(k)}$ möglich. Wenn dabei ein betragsgrößtes Element gemäß

$$|a_{k,\hat{s}(k)}^{(k)}| = \max\{|a_{kj}^{(k)}|: j = k, ..., n\}$$
(7)

als Pivot gewählt und durch Vertauschen der Spalten k und $\hat{s}(k)$ in die Position (k, k) gebracht wird, spricht man von Zeilenpivotisierung. In der Matrixformulierung ist dann (5.1.12) durch

$$\hat{A}^{(k)} := A^{(k)} T_{k,\hat{s}(k)}$$
(8)

zu ersetzen. Analog zum Abschnitt 5.1.C ergibt sich damit eine Dreiecksfaktorisierung

$$A\hat{P}^{\mathsf{T}} = \hat{L}\hat{R} \tag{9}$$

 \mathbf{mit}

$$\hat{\boldsymbol{L}} = \boldsymbol{L}_{1}(\hat{\boldsymbol{l}}^{1}) \cdots \boldsymbol{L}_{n-1}(\hat{\boldsymbol{l}}^{n-1}), \qquad \hat{\boldsymbol{P}}^{\intercal} = \boldsymbol{T}_{1,\hat{\boldsymbol{s}}(1)} \cdots \boldsymbol{T}_{n-1,\hat{\boldsymbol{s}}(n-1)}.$$
(10)

Die Matrizen $\hat{L}_k = L_k(\hat{l}^k)$ sind i. allg. keine SLNT-Matrizen. Aus (7) folgt jedoch

$$|\hat{a}_{kj}^{(k)}/\hat{a}_{kk}^{(k)}| \leq 1 \qquad (j = k + 1, ..., n)$$
 (11)

und damit

$$|a_{ij}^{(k+1)}| = \left| \hat{a}_{ij}^{(k)} - \frac{\hat{a}_{ik}^{(k)} * \hat{a}_{kj}^{(k)}}{\hat{a}_{kk}^{(k)}} \right| \leq |\hat{a}_{ij}^{(k)}| + |\hat{a}_{ik}^{(k)}| \quad (i, j = k + 1, ..., n)$$
(12)

als Analogon zur Abschätzung (3). Die Abschätzungen aus Teilabschnitt A gelten deshalb auch für die Zeilenpivotisierung, wenn sie zeilenweise und in der 1-Norm formuliert werden.

C. Vollständige Pivotisierung

Der Spielraum für die Wahl des Pivotelementes wird maximal ausgenutzt, wenn das Pivot als ein betragsgrößtes Element der gesamten Restmatrix $M^{(k)}$ gewählt wird, d. h., wenn s(k), $\hat{s}(k)$ gemäß

$$|a_{s(k),\hat{s}(k)}^{(k)}| = \max\{|a_{ij}^{(k)}|: i, j = k, ..., n\}$$
(13)

festgelegt und das Pivotelement $a_{s(k),\hat{s}(k)}^{(k)}$ durch Vertauschen der Zeilen k, s(k) und der Spalten $k, \hat{s}(k)$ in die Position (k, k) gebracht wird. Es gilt dann

$$\hat{A}^{(k)} = T_{k,s(k)} A^{(k)} T_{k,\hat{s}(k)}, \tag{14}$$

was auf eine Dreiecksfaktorisierung des Typs

$$PA\hat{P}^{\mathsf{T}} = \vec{L}\vec{R} \tag{15}$$

führt. Wegen (13) gelten dabei sowohl die für die Spaltenpivotisierung wie auch die für die Zeilenpivotisierung gültigen Abschätzungen.

5.2.7. Bemerkung. (i) Bei vollständiger Pivotisierung gilt

$$\mu_k := \max \{ |a_{ij}^{(k)}| : i, j = k, ..., n \} \leq \varrho_k \mu_1 \qquad (k = 2, ..., n)$$

mit der Wachstumsschranke

$$\rho_k := k^{1/2} [2 \cdot 3^{1/2} \cdots k^{1/(k-1)}]^{1/2} \le 1.8 \ k^{0.25 \ln k}. \tag{16}$$

Im Gegensatz zu 5.2.1 ist diese Schranke jedoch nicht scharf, und bis jetzt sind keine reellen Matrizen bekannt geworden, für die

$$\mu_k \leq k\mu_1$$

nicht gilt.

(ii) Trotz des günstigeren Stabilitätsverhaltens wird die vollständige Pivotisierung für reguläre Matrizen praktisch nur in Ausnahmefällen verwendet, da der Aufwand im Vergleich zur Spalten- oder Zeilenpivotisierung sehr hoch ist ($\sim n^3/3$ Lese- und Vergleichsoperationen gegenüber $\sim n^2/2$) und die einfacheren Strategien i. allg. zufriedenstellende Ergebnisse liefern. \Box

D. Dreiecksfaktorisierung diagonaldominanter und symmetrischer definiter Matrizen Eine Matrix $A = (a_{ij}) \in \mathbb{R}^{n,n}$ heißt diagonaldominant — genauer: zeilendiagonaldominant —, wenn

$$|a_{ii}| > \sum_{j=1, j+i}^{n} |a_{ij}| \qquad (i = 1, ..., n)$$
(17)

gilt. Nach 1.2.10 ist eine diagonaldominante Matrix regulär, denn wegen (17) kann $\lambda = 0$ in keinem Geršgorin-Kreis von A liegen.

5.2.8. Aussage. Für eine diagonaldominante Matrix A ist die Dreiecksfaktorisierung nach dem Gaußschen Algorithmus ohne Pivotisierung durchführbar, wobei folgende Aussagen gelten:

(i) Alle Matrizen $M^{(k)}$ sind diagonal dominant,

(ii)
$$\|\boldsymbol{M}^{(k)}\|_{\infty} \leq \|\boldsymbol{M}^{(k-1)}\|_{\infty} \leq \|\boldsymbol{A}\|_{\infty}$$
 $(k = 2, ..., n).$ (18)

Beweis. Es genügt, den Übergang von $A = M^{(1)}$ zu $M^{(2)}$ zu betrachten. Wegen (17) ist $a_{11} \neq 0$ und kann als Pivot gewählt werden, womit sich

$$a_{ij}^{(2)} = a_{ij} - a_{i1} * a_{1j}/a_{11} \qquad (i, j = 2, ..., n)$$
⁽¹⁹⁾

ergibt. Nun ist (17) äquivalent zu

$$au_i := \sum_{j=1, j \neq i}^n |a_{ij}|/|a_{ii}| < 1$$
 $(i = 1, ..., n)$

Aufsummieren von (19) liefert

$$\sum_{j=2}^{n} |a_{ij}^{(2)}| \leq \sum_{j=2}^{n} |a_{ij}| + |a_{i1}| \sum_{j=2}^{n} |a_{1j}| / |a_{11}| = \sum_{j=2}^{n} |a_{ij}| + |a_{i1}| \tau_1 \leq \sum_{j=1}^{n} |a_{ij}|,$$
(20)

also $||M^{(2)}||_{\infty} \leq ||M^{(1)}||_{\infty}$. Zum Nachweis der Diagonaldominanz von $M^{(2)}$ folgern wir aus (19) mit $p := |a_{i1}|/|a_{i1}|$, $q := |a_{1i}|/|a_{11}|$, p, q < 1, zunächst $|a_{ii}^{(2)}| \geq |a_{ii}| - |a_{i1}| |a_{1i}|/|a_{11}| = |a_{ii}| (1 - pq) > 0$ und damit unter Beachtung von (20)

$$rac{\sum\limits_{j=2,j+1}^n |a_{ij}^{(2)}|}{|a_{ii}^{(2)}|} \leq rac{\sum\limits_{j=2,j+i}^n |a_{ij}|+|a_{i1}|\;(au_1-q)}{|a_{ii}|\;(1-pq)} = rac{ au_i-p+p(au_1-q)}{1-pq} = rac{ au_i-p+p(au_1-q)}{1-pq}$$
 $= au_i-prac{q(1- au_i)+(1- au_1)}{1-pq} \leq au_i < 1,$

d. h., auch $M^{(2)}$ ist diagonal dominant, und zwar mindestens so stark wie $M^{(1)}$.

Aussage 5.2.8 besagt, $da\beta$ für (zeilen-)diagonaldominante Matrizen ohne Pivotisierung gearbeitet werden kann. Analoge Aussagen gelten auch für spaltendiagonaldominante Matrizen, d. h. für Matrizen A, für die A^{\intercal} zeilendiagonaldominant ist.

Eine weitere Matrizenklasse, für deren Dreiecksfaktorisierungen keine Pivotisierung erforderlich ist, stellt die Klasse der symmetrischen definiten Matrizen dar. Wir beschränken uns auf symmetrische und positiv definite Matrizen, denn für negativ definites A ist (-A) positiv definit.

5.2.9. Aussage. Für eine symmetrische und positiv definite Matrix A ist die Dreiecksfaktorisierung nach dem Gaußschen Algorithmus ohne Pivotisierung durchführbar, wobei folgende Eigenschaften gelten:

(i) Alle Matrizen $M^{(k)}$ sind symmetrisch und positiv definit.

(ii)
$$\|M^{(k)}\|_{F} < \|M^{(k-1)}\|_{F} \le \|A\|_{F}$$
 $(k = 2, ..., n).$ (21)

Beweis. Nach Ü 1.1.6 ist $a_{11} > 0$ und kann als Pivot gewählt werden. Wir betrachten wieder den ersten Schritt und schreiben A bzw. $A^{(2)}$ in der Gestalt

$$A = A^{(1)} = \left(\frac{a_{11}}{a} \middle| \frac{a^{\mathsf{T}}}{W} \right), \quad A^{(2)} = \left(\frac{a_{11}}{o} \middle| \frac{a^{\mathsf{T}}}{M} \right) \quad \text{mit} \quad M = M^{(2)}.$$

Nach 3.2.1 ist $L_1 = L_1(-l) = l - l(e^1)^{\mathsf{T}}$ mit $l = \frac{1}{a_{11}} \begin{pmatrix} 0 \\ a \end{pmatrix}$ und folglich

$$\boldsymbol{M} = \boldsymbol{W} - \frac{\boldsymbol{a}\boldsymbol{a}^{\mathsf{T}}}{\boldsymbol{a}_{11}} = \boldsymbol{M}^{\mathsf{T}}.$$
(22)

Zum Nachweis der positiven Definitheit setzen wir $\boldsymbol{x} = \begin{pmatrix} \eta \\ \boldsymbol{y} \end{pmatrix}$, $\eta \in \mathbb{R}$, $\boldsymbol{y} \in \mathbb{R}^{n-1}$, und beachten daß ten, daß

$$\boldsymbol{x}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{x} = a_{11}\eta^2 + 2(\boldsymbol{a}^{\mathsf{T}}\boldsymbol{y}) \ \eta + \boldsymbol{y}^{\mathsf{T}}\boldsymbol{W}\boldsymbol{y} = a_{11} \left(\eta + \frac{\boldsymbol{a}^{\mathsf{T}}\boldsymbol{y}}{a_{11}}\right)^2 + \boldsymbol{y}^{\mathsf{T}}\boldsymbol{M}\boldsymbol{y} > 0$$
(23)

für alle $x \neq o$ ist. Dann muß $y^{\mathsf{T}}My > 0$ für $y \neq o$ sein, d. h., M ist positiv definit. Wir betrachten jetzt eine Spiegelung $H \in \mathbb{R}^{n-1,n-1}$ gemäß 3.3.4 mit $Ha = \pm ||a||_2 e^1$. Für $\overline{M} := HMH^{\top}$ und $\overline{W} := HWH^{\intercal}$ folgt dann aus (22) $\overline{W} = \overline{M} + (||\boldsymbol{a}||_2^2/a_{11}) e^{1}e^{1\intercal}$, also $||\overline{W}||_F^2 = ||\overline{M}_F^2||$ + $2\overline{m}_{11} ||\boldsymbol{a}||_2^2 a_{11} + ||\boldsymbol{a}||_2^4/a_{11}^2 \ge ||\overline{M}||_F^2$, denn \overline{M} ist wie M positiv definit, so daß auch $\overline{m}_{11} > 0$ gilt. Wegen der Orthogonalinvarianz der Frobeniusnorm erhält man schließlich

$$\|A\|_{F}^{2} > \|W\|_{F}^{2} = \|\overline{W}\|_{F}^{2} \ge \|\overline{M}\|_{F}^{2} = \|M\|_{F}^{2}.$$

Aus dem Beweis ergibt sich sofort die folgende Aussage:

5.2.10. Aussage. Für eine symmetrische und positiv definite Matrix A liefert der Gaußsche Algorithmus ohne Pivotisierung eine Dreiecksfaktorisierung A = LRmit $L^{\intercal} = D^{-1}R$, d. h., es gilt $A = LDL^{\intercal}$ (24) mit $D = \text{diag}(d_i), d_i = r_{ii} > 0$ (i = 1, ..., n), und einer unteren Einsdreiecks-

$$A = LDL^{\mathsf{T}} \tag{24}$$

matrix L

Auf die Berechnung von Dreiecksfaktorisierungen symmetrischer Matrizen und die dabei mögliche Halbierung von Speicherplatz und Rechenaufwand gehen wir im Abschnitt 6.1 nochmals im Detail ein.

Übungsaufgaben

 $\ddot{\mathbf{U}}$ 5.2.1. Man beweise im Fall der Spaltenpivotisierung die Gültigkeit von $\|m{M}^{(k)}\|_{\infty} \leq k\,\|m{A}\|_{\infty}$ im Fall einer oberen Hessenbergmatrix A, d. h. einer Matrix mit $a_{ii} = 0$ für i > j + 1, und $\operatorname{von} \|M^{(k)}\|_{\infty} \leq 2 \|A\|_{\infty}$ im Fall einer Tridiagonalmatrix A.

Ü 5.2.2. Man überlege sich, daß für die Matrix

$$A = \begin{pmatrix} 1 & & & | a \\ -1 & 1 & & | a \\ -1 & -1 & 1 & | a \\ \vdots & \vdots & \ddots & | \vdots \\ & & 1 & | a \\ \hline -1 & -1 \dots -1 & | a \end{pmatrix}$$

im Fall der Spaltenpivotisierung

$$\|A\|_{\infty} = |a| + n - 1, \quad \|M^{(k)}\|_{\infty} = 2^{k-1} |a| + n - k$$

gilt, so daß der Quotient $||M^{(k)}||_{\infty}/||A||_{\infty}$ der Schranke $\gamma_k = 2^{k-1}$ aus 5.2.1 für $|a| \to \infty$ beliebig nahe kommt.

Ü 5.2.3. Es sei A regulär. Man zeige: Wenn $PA^{\mathsf{T}} = LR$ die nach dem Gaußschen Algorithmus mit Spaltenpivotisierung berechnete Dreiecksfaktorisierung von A^{T} und $A\hat{P}^{\mathsf{T}} = \hat{L}\hat{R}$ die mittels Zeilenpivotisierung berechnete Faktorisierung bezeichnen, gilt

$$\hat{P} = P$$
, $\hat{L} = R^{\mathsf{T}} D^{-1}$, $\hat{R} = DL^{\mathsf{T}}$ mit $D = \text{diag}(r_{ii})$

sofern bei der Zeilenpivotisierung dieselbe Strategie wie bei der Spaltenpivotisierung verwendet wird.

 $\ddot{\mathbf{U}}$ 5.2.4. Man gebe zu 5.2.3, 5.2.4 und 5.2.5 analoge Algorithmen an, die mit impliziten Vertauschungen entsprechend Bemerkung 5.2.6 (ii) arbeiten.

Ü 5.2.5. Es ist zu zeigen, daß der Gaußsche Algorithmus genau dann ohne Pivotisierung mit regulärem R durchgeführt werden kann, wenn die Hauptabschnittsdeterminanten det (A_k) von A mit

$$A_k := \begin{pmatrix} a_{11} \dots a_{1k} \\ \vdots & \vdots \\ a_{k1} \dots a_{kk} \end{pmatrix}$$

der Bedingung

 $\det (A_k) \neq 0 \qquad (k = 1, ..., n) \tag{25}$

genügen, und daß bei Erfülltsein von (25)

$$r_{kk} = \det(A_k) / \det(A_{k-1}) \qquad (k = 1, ..., n)$$
⁽²⁶⁾

mit det $(A_0) := 1$ gilt.

Ü 5.2.6. Man spricht von *Diagonalpivotisierung*, wenn das Pivotelement als betragsmaximales Diagonalelement von $M^{(k)}$ gesucht wird, d. h., wenn $a_{s(k),s(k)}^{(k)}$ mit

$$|a_{s(k),s(k)}^{(k)}| = \max\{|a_{ii}^{(k)}|: k \le i \le n\}$$
(27)

als Pivot gewählt wird. Man zeige, daß für diagonaldominante und symmetrische definite Matrizen die Diagonalpivotisierung identisch mit der vollständigen Pivotisierung ist.

Ü 5.2.7. Man überlege sich: Die Matrix $A \in S^{n,n}$ ist positiv definit genau dann, wenn sie sich in der Form $A = LDL^{\intercal}$ mit einer unteren Einsdreiecksmatrix L und $D = \text{diag}(d_i)$, $d_i > 0$, faktorisieren läßt.

5.3. Rundungsfehleranalyse

Die bei der Pivotisierung vorgenommenen Zeilen- bzw. Spaltenvertauschungen entsprechen einer Umnumerierung der Gleichungen bzw. Unbekannten. Für die folgenden Überlegungen kann daher o. B. d. A. vorausgesetzt werden, daß diese Vertauschungen bereits vor Ausführung des Gaußschen Algorithmus vorgenommen worden sind und während der Elimination stets $s(k) = k = \hat{s}(k)$ gewählt wird. Damit werden die Vertauschungen und das Zeichen "^{*} für die vertauschten Größen überflüssig.

A. Fehleranalyse der Dreiecksfaktorisierung

Mit der obigen Vereinbarung lautet der k-te Eliminationsschritt

$$A^{(k+1)} = \begin{pmatrix} I^{(k-1)} & \mathbf{O} \\ \hline \mathbf{O} & L_1^{(n-k+1)} \end{pmatrix} \begin{pmatrix} \mathbf{R}^{(k)} \\ \hline \mathbf{O} & \mathbf{M}^{(k)} \end{pmatrix} = \begin{pmatrix} \mathbf{R}^{(k)} \\ \hline \mathbf{O} & L_1^{(n-k+1)} \mathbf{M}^{(k)} \end{pmatrix},$$
(1)

vgl. (5.1.13). Es genügt daher, die Rundungsfehler bei der Transformation

$$M^{(k)} \to \overline{M}^{(k)} := \left(rac{a_{kk}^{(k)} \mid a_{k,k+1}^{(k)} \dots a_{kn}^{(k)}}{o \mid M^{(k+1)}}
ight) = L_1^{(n-k-1)} M^{(k)}$$

zu untersuchen, wobei $L_1^{(n-k+1)} := L_1(-l^{1,n-k+1})$ ist. Dies ist gerade der in 3.2.5 betrachtete Eliminationsschritt. Der berechnete Vektor $l^{1,n-k+1} = (0, l_{k+1,k}, ..., l_{nk})^{\mathsf{T}}$ und die berechnete Matrix $\overline{M}^{(k)}$ bzw. $M^{(k+1)}$ genügen dann wegen 3.2.6 der Gleichung

$$\overline{\boldsymbol{M}}^{(k)} = \boldsymbol{L}_1^{(n-k+1)}(\boldsymbol{M}^{(k)} + \boldsymbol{\delta}\boldsymbol{M}^{(k)}), \qquad (2)$$

und für die Störungsmatrix $\delta M^{(k)}$ gilt

$$|\boldsymbol{\delta}\boldsymbol{M}^{(k)}| \leq \nu \left\{ |\boldsymbol{M}^{(k)}| + 2 \left(\frac{\boldsymbol{0} \mid \boldsymbol{o}^{\mathsf{T}}}{\boldsymbol{o} \mid |\boldsymbol{M}^{(k+1)}|} \right) \right\},\tag{3}$$

also

$$\|\delta M^{(k)}\| \leq \nu(\|M^{(k)}\| + 2 \|M^{(k+1)}\|), \qquad (4)$$

vgl. Formel (3.2.15) und deren Herleitung.

Einsetzen von (2) in (1) führt auf

$$A^{(k+1)} = L_k(A^{(k)} + \delta A^{(k)}) \tag{5}$$

mit

und

$$L_k = L_k(-l^k), \quad l^k = (0, ..., 0, l_{k+1,k}, ..., l_{nk})^{\mathsf{T}}$$

$$\delta A^{(k)} = egin{pmatrix} O & O \ O & O \end{pmatrix}, \quad \| \delta A^{(k)} \| = \| \delta M^{(k)} \|.$$

Für den berechneten Faktor R der Dreiecksfaktorisierung gilt dann

$$egin{aligned} R &= A^{(n)} = L_{n-1}(A^{(n-1)} + \delta A^{(n-1)}) \ &= L_{n-1}[L_{n-2}(A^{(n-2)} + \delta A^{(n-2)}) + \delta A^{(n-1)}] \ &= L_{n-1}L_{n-2}[A^{(n-2)} + \delta A^{(n-2)} + \delta A^{(n-1)}], \end{aligned}$$

denn wegen der speziellen Blockgestalt (6) von $dA^{(k)}$ ist

$$L_{n-2} \delta A^{(n-1)} = [I - l^{n-2} (e^{n-2})^{\mathsf{T}}] \, \delta A^{(n-1)} = \delta A^{(n-1)}$$

Sukzessive Fortsetzung dieser Überlegungen führt auf

$$R = L_{n-1}L_{n-2}\cdots L_2L_1[A + \delta A^{(1)} + \cdots + \delta A^{(n-1)}].$$

(6)

Da sich $L = L_1^{-1} \cdots L_{n-1}^{-1} = (l_{ik})$ ohne Rundungsfehler aus den Koeffizienten l_{ik} ergibt, vgl. 5.1.C, folgt unter Beachtung von (4), (6) die Darstellung

$$LR = A + \delta A$$
, $\delta A := \sum_{k=1}^{n-1} \delta A^{(k)}$

mit

$$\|\delta A\| \leq \nu \left(\|A\| + 3 \sum_{k=2}^{n} \|M^{(k)}\| \right).$$
(7)

5.3.1. Rundungsfehleranalyse der Dreiecksfaktorisierung nach Gauß. Für $A \in \mathbb{R}^{n,n}$ sei der Gaußsche Algorithmus durchführbar, und $L, R \in \mathbb{R}^{n,n}$ seien die berechneten Dreiecksfaktoren. Dann gilt

$$LR = A + \delta A \tag{8}$$

mit einer Störung dA, die der Abschätzung

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq \boldsymbol{v}\boldsymbol{F}[\boldsymbol{A}] \,\|\boldsymbol{A}\| \tag{9}$$

genügt, wobei

$$F[A] := 1 + 3 \sum_{k=2}^{n} \|M^{(k)}\| / \|A\|$$
(10)

ist.

Dies besagt: In jeder Matrizenklasse $\mathfrak{A} \subset \mathfrak{R}^{n,n}$, in welcher der Gaußsche Algorithmus durchführbar ist und die Zahlen F[A] im Sinne von

$$F[A] \leq F := F[\mathfrak{A}] := \max \{F[A] : A \in \mathfrak{A}\} < \infty$$
(11)

gleichmäßig beschränkt sind, ist die Gaußsche Dreiecksfaktorisierung numerisch gutartig mit F gemäß (11). Gutartigkeit liegt also vor, wenn die $M^{(k)}$ nicht beliebig größer als A werden können.

Mit den Wachstumsaussagen für $M^{(k)}$ aus Abschnitt 5.2 ergibt sich dann das folgende Resultat:

5.3.2. Satz. Die Gaußsche Dreiecksfaktorisierung ist durchführbar und numerisch gutartig

- (i) ohne Pivotisierung für diagonaldominante bzw. symmetrische definite Matrizen A mit $F \leq 3n$ bezüglich der ∞ bzw. F-Norm,
- (ii) mit Spalten-, Zeilen- oder vollständiger Pivotisierung für alle regulären Matrizen A mit $F \leq 3 \cdot 2^n$ bezüglich der Normen mit $p = 1, \infty, F$,

sofern die Kondition von A im Sinne von

$$\nu F \operatorname{cond} \left(A \right) < 1 \tag{12}$$

mit den angegebenen Werten von F nicht zu groß ist.

Beweis. Die Aussagen aus Abschnitt 5.2 sind nicht direkt auf A anwendbar, da sie nur in exakter Arithmetik gültig sind. Für die betrachteten Matrizenklassen ist jedoch der erste Eliminationsschritt stets durchführbar. Aus der Herleitung von 5.3.1 ist zu sehen, daß die

berechnete Matrix $M^{(2)}$ und die berechneten Koeffizienten l_{i1} das exakte Resultat eines Eliminationsschrittes für eine gestörte Matrix $\tilde{A}^{(1)} = A + \delta A^{(1)}$ sind, wobei $\delta A^{(1)}$ der Ungleichung (9) genügt. Auf Grund von (12) und des Störungslemmas 4.1.2 ist $\tilde{A}^{(1)}$ für alle solche Störungen regulär, im definiten Fall also auch definit, siehe Ü 4.1.9, so daß die Aussagen aus 5.2 auf $\tilde{A}^{(1)}$ statt A angewendet werden können und die jeweiligen Schranken für $M^{(2)}$ liefern. Analog ergeben sich $M^{(3)}$ und die l_{i2} als exaktes Resultat für eine gestörte Matrix $\tilde{A}^{(2)}$ usw. Durch Induktion folgt dann die Behauptung.

Für den diagonaldominanten Fall müßte noch gezeigt werden, daß auch die berechneten Matrizen $M^{(k)}$ diagonaldominant sind. Dies ist bei im Sinne von

$$au_i = \sum_{\substack{j=1 \ j \neq i}}^n |a_{ij}|/|a_{ii}| < \left(\frac{1-
u}{1+
u}\right)^i = 1 - 2i
u \quad (i = 1, ..., n)$$

genügend starker Diagonaldominanz der Fall, wir verzichten jedoch auf einen detaillierten Beweis. Praktisch ist diese Bedingung kaum schärfer als die eigentliche Dominanzbedingung $\tau_i < 1$ (i = 1, ..., n), vgl. den Beweis zu 5.2.8. \Box

5.3.3. Bemerkung. (i) Die Bedingung (12) garantiert, daß sämtliche Matrizen aus der Umgebung $\mathcal{M} := \{A + \delta A : \|\delta A\| \leq \nu F \|A\|$ von A regulär sind, und sichert daher die Durchführbarkeit der Dreiecksfaktorisierung mit regulärem R in der durch ν charakterisierten Gleitpunktarithmetik. Auf die Bedingung (12) kann verzichtet werden, wenn bei der Dreiecksfaktorisierung im k-ten Schritt $a_{kk}^{(k)} := \nu \|A\|$ gesetzt wird, falls mit der jeweiligen Pivotstrategie kein von 0 verschiedenes Pivot gefunden werden kann. Durch diese zusätzliche Störung wird die Gutartigkeit der Faktorisierung nicht beeinträchtigt. Allerdings können dann in der Menge \mathcal{M} auch singuläre Matrizen liegen, insbesondere kann A selbst singulär sein. Die Lösung von Ax = b ist dann eine numerisch inkorrekt gestellte Aufgabe, und die mit der obigen Modifikation erzwungene Dreiecksfaktorisierung ist i. allg. nicht zur Lösung geeignet; das berechnete x ist i. allg. von der Größenordnung $1/\nu$. In gewissen Ausnahmefällen — etwa bei der sog. inversen Iteration zur Eigenwertbestimmung, siehe Abschnitt 13.4 — ist dieser Fall jedoch sogar erwünscht.

(ii) In 5.3.2 kann der Faktor 3 bei sorgfältigerer Abschätzung durch 2 ersetzt werden.

(iii) Für praktisch auftretende Gleichungssysteme gilt 5.3.2(ii) mit der empirischen Schranke $F \leq \gamma n, \gamma \approx 10$, vgl. 5.2.2.

(iv) Im Fall der vollständigen Pivotisierung ist

 $F \leq 1.8n^2 n^{0.251nn},$

vgl. 5.2.7. Für praktisch auftretende Systeme liegt F in der Größenordnung von 1.

(v) Bei Bedarf kann μ_k bzw. $||M^{(k)}||$ im Laufe der Elimination mit berechnet werden. Ein zu schnelles Wachstum dieser Größen signalisiert einen Verlust an numerischer Gutartigkeit. \Box

B. Fehleranalyse der Gleichungsauflösung

Wir betrachten jetzt den gesamten Lösungsprozeß, bei dem sich an die Dreiecksfaktorisierung A = LR die Vorwärtselimination $c = L^{-1}b$ und die Rücksubstitution $x = R^{-1}c$ gemäß 5.1.7 anschließen. Nach 4.3.2 genügen die berechneten Größen c bzw. x den gestörten Gleichungen

$$(\boldsymbol{L} + \boldsymbol{\delta} \boldsymbol{L}) \, \boldsymbol{c} = \boldsymbol{b} \quad \text{bzw.} \quad (\boldsymbol{R} + \boldsymbol{\delta} \boldsymbol{R}) \, \boldsymbol{x} = \boldsymbol{c}$$
(13)

mit den elementweisen Schranken

$$|\delta l_{ik}| \leq vk |l_{ik}| \quad \text{bzw.} \quad |\delta r_{kj}| \leq v(n-j+1) |r_{kj}|. \tag{14}$$

Aus (13) folgt unter Berücksichtigung von (8)

$$(\boldsymbol{L} + \boldsymbol{\delta}\boldsymbol{L}) (\boldsymbol{R} + \boldsymbol{\delta}\boldsymbol{R}) \boldsymbol{x} = (\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \boldsymbol{x} = \boldsymbol{b}$$
(15)

mit

$$\delta_1 A := \delta A + \delta A', \qquad \delta A' := L \delta R + \delta L (R + \delta R), \tag{16}$$

d. h., die äquivalenten Störungen $\sigma_1 A$ für den Gesamtprozeß entstehen durch Überlagerung der Störungen σA für die Faktorisierung und der Störungen $\sigma A'$ für die Lösung der beiden Dreieckssysteme.

Wegen (14) ist

$$egin{aligned} \delta a_{ij}'| &= \left|\sum\limits_{k=1}^{\min(i,j)} [l_{ik} \delta r_{kj} + \delta l_{ik} (r_{kj} + \delta r_{kj})]
ight| \ &\leq extstyle \sum\limits_{k=1}^{\min(i,j)} |l_{ik} r_{kj}| \ (n-j+1+k) \leq extstyle (n+1) \sum\limits_{k=1}^{\min(i,j)} |l_{ik} r_{kj}| \,. \end{aligned}$$

Analog zu (3.2.19) gilt mit $|\varepsilon|, |\vartheta| \leq v$

$$l_{ik}r_{kj} = l_{ik}a_{kj}^{(k)} = a_{ij}^{(k)}/(1+\varepsilon) - a_{ij}^{(k+1)}/[(1+\varepsilon)(1+\vartheta)],$$

also $|l_{ik}r_{kj}| \leq |a_{ij}^{(k)}| + |a_{ij}^{(k-1)}|$ und daher

$$|\delta a_{ij}'| \le r(n+1) \left(|a_{ij}^{(1)}| + 2\sum\limits_{k=2}^{\min(i,j)+1} |a_{ij}^{(k)}|
ight).$$

Die letzte Ungleichung kann äquivalent in der Form

$$|\delta A'| \leq \nu(n+1) \left[|A| + 2\sum_{k=2}^{n} \left(\frac{O|O}{O||M^{(k)}|} \right) \right]$$
(17)

geschrieben werden und führt auf

$$\|\boldsymbol{\delta}\boldsymbol{A}'\| \leq \boldsymbol{\nu}(n+1) \left(\|\boldsymbol{A}\| + 2\sum_{k=2}^{n} \|\boldsymbol{M}^{(k)}\| \right) = \boldsymbol{\nu}\boldsymbol{F}'[\boldsymbol{A}] \|\boldsymbol{A}\|$$
(18)

mit

$$F'[A] := (n + 1) \left(1 + 2 \sum_{k=2}^{n} ||M^{(k)}|| / ||A|| \right).$$
(19)

Ein Vergleich mit (9), (10) zeigt, daß $F' \leq (n + 1) F$ gilt. Unter Beachtung von (16) ergibt sich daher das folgende Resultat.

5.3.4. Rundungsfehleranalyse des Gaußschen Algorithmus. Für $A \in \Re^{n,n}$ und $b \in \Re^n$ sei der Gaußsche Algorithmus durchführbar, und $x \in \Re^n$ sei die nach diesem Verfahren berechnete Lösung des Gleichungssystems Ax = b. Dann gilt

$$(A + \boldsymbol{\delta}_1 A) \, \boldsymbol{x} = \boldsymbol{b} \tag{20}$$

mit einer Störung ${oldsymbol \sigma}_1 A$, für die $\|{oldsymbol \sigma}_1 A\| \le
u F_1 [A] \, \|A\|$ m

$$\|\boldsymbol{\delta}_{1}\boldsymbol{A}\| \leq \nu F_{1}[\boldsymbol{A}] \,\|\boldsymbol{A}\| \quad \text{mit} \quad F_{1}[\boldsymbol{A}] \leq (n+2) \, F[\boldsymbol{A}]$$
(21)

und F[A] gemäß (10) gilt.

Aus der Störung $\delta_1 A$ kann auf den Fehler $\delta x = x - x^*$ geschlossen werden.

5.3.5. Folgerung. Der erzeugte Rundungsfehler $\delta x = x - x^*$, $x^* = A^{-1}b$, genügt der Abschätzung

$$\|\boldsymbol{\delta x}\| \leq \kappa \|\boldsymbol{x}\| \quad \text{mit} \quad \kappa := \nu F_1[\boldsymbol{A}] \text{ cond } (\boldsymbol{A}), \tag{22}$$

 $\| \boldsymbol{\delta x} \| \leq \varkappa \| \boldsymbol{x} \|$ und im Fall $\varkappa < 1$ gilt

$$\|\boldsymbol{\delta x}\| \leq \kappa/(1-\kappa) \,\|\boldsymbol{x^*}\|. \tag{23}$$

Beweis. Aus (20) folgt $\delta x = -A^{-1}\delta_1 A x$ und damit (22), während sich (23) aus 4.1.3 ergibt. \Box

5.3.4 besagt: In jeder Matrizenklasse $\mathfrak{A} \in \mathfrak{R}^{n,n}$, für die (11) gilt, ist die Berechnung von x nach dem Gaußschen Algorithmus numerisch gutartig mit der Fehlerkumulationskonstanten $F_1 \leq (n+2)$ F. Durch 5.3.5 wird die numerische Stabilität mit demselben F_1 ausgedrückt.

5.3.6. Bemerkung. (i) Die bei der Lösung der Dreieckssysteme auftretende spezielle Fehlerkorrelation bewirkt, daß der durch den Anteil $\delta A'$ von $\delta_1 A$ erzeugte Fehler $-A^{-1}\delta A'x$ im allgemeinen wesentlich kleiner als die Schranke $||A^{-1}|| ||\delta A'|| ||x||$ ist, und zwar um so mehr, je größer cond (A) ist. Praktisch kann daher F_1 in 5.3.5 durch \tilde{F}_1 ,

 $F_1 \approx egin{cases} F, & \mbox{falls cond} (A) \mbox{groß}, \ nF, & \mbox{falls cond} (A) \mbox{klein}, \end{cases}$

ersetzt werden. Aus den angegebenen Gründen führt auch eine höhere Genauigkeit bei der Lösung der Dreieckssysteme i. allg. nicht zu einer höheren Genauigkeit von x.

(ii) Die elementweisen Schranken (3) und (17) zeigen, daß δA und $\delta A'$ zeilenbzw. spaltenweise relativ klein sind, sofern $M^{(k)}$ zeilen- bzw. spaltenweise in vergleichbarer Größenordnung mit A liegt. Letzteres ist bei Zeilen- bzw. Spaltenpivotisierung der Fall. Man kann zeigen, daß die Zeilenpivotisierung in etwa dieselbe Wirkung wie eine im Sinne von 4.1.C optimale Zeilenskalierung hat. Die Spaltenpivotisierung wirkt dagegen etwa wie eine in bezug auf das Residuum b - Axoptimale Spaltenskalierung. \Box

Übungsaufgaben

Ü 5.3.1. Man ergänze die Beweisskizze zu 5.3.2.

m U 5.3.2. Man beweise: Wenn $d = \det{(A)}$ gemäß 5.1.9 berechnet wird, gilt für den berechneten Wert d

 $d = [\det (A + \delta A)] (1 + \varepsilon)$

mit δA gemäß 5.3.1 und $|\varepsilon| \leq \nu(n-1)$, sofern kein Unter- oder Überlauf eintritt.

5.4. Genauigkeitsabschätzung und iterative Verbesserung

Bei gegebenen Eingangsdaten $\{A, b\}$ liefert der Gaußsche Algorithmus nach erfolgreicher Ausführung die berechnete Lösung $x \in \Re^n$ und die berechneten Dreiecksfaktoren $L, R \in \Re^{n,n}$. Auf Grund der auftretenden Rundungsfehler ist x i. allg. nur ein Näherungswert für die exakte Lösung $x^* = A^{-1}b$. Daher entstehen die Fragen,

- wie der erzeugte Rundungsfehler $\delta x = x - x^*$ abgeschätzt werden kann,

- wie zu x eine Korrektur h berechnet werden kann, so daß x' = x + h eine bessere Näherung für x^* als x ist.

A. Fehlerschranken und Korrekturen

Im vorangegangenen Abschnitt wurde gezeigt, daß die berechnete Lösung x das gestörte System

$$(\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \boldsymbol{x} = \boldsymbol{b} \quad \text{mit} \quad \|\boldsymbol{\delta}_1 \boldsymbol{A}\| \leq v \boldsymbol{F}_1 \|\boldsymbol{A}\| \tag{1}$$

löst und daß

$$\|\boldsymbol{\delta x}\| \leq \boldsymbol{\varkappa} \|\boldsymbol{x}\| \quad \text{mit} \quad \boldsymbol{\varkappa} := \boldsymbol{\imath} F_1 \text{ cond } (\boldsymbol{A}) \tag{2}$$

gilt. Die Schranken in (1) und (2) sind für die praktische Fehlerabschätzung nicht besonders geeignet, da sie mit den a-priori-Werten für F_1 gemäß 5.3.2 und 5.3.4 i. allg. zu pessimistisch sind. Es bietet sich daher an, unter Verwendung des Residuums

$$\boldsymbol{r^*} = \boldsymbol{r^*}(\boldsymbol{x}) = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x} \tag{3}$$

eine a-posteriori-Abschätzung durchzuführen und dabei auf Teilabschnitt 4.1.D zurückzugreifen. Nach 4.1.18 existiert zu x eine Störung δA mit

 $(A + \delta A) \boldsymbol{x} = \boldsymbol{b} \quad \text{und} \quad \|\boldsymbol{\delta} A\| \leq \|\boldsymbol{b} - A\boldsymbol{x}\| / \|\boldsymbol{x}\|, \tag{4}$

und nach 4.1.17 gilt die Fehlerabschätzung

$$\|\delta x\| \le \|A^{-1}\| \|b - Ax\|.$$
⁽⁵⁾

Bei der Auswertung von (5) wird eine Schranke bzw. zumindest ein Schätzwert für $||A^{-1}||$ benötigt; dieses Problem wird im Teilabschnitt C behandelt. Außerdem muß das Residuum $r^* = b - Ax$ numerisch berechnet werden, was in der Form r = fl (b - Ax) wie in Ü 2.3.4 geschehen möge. Dabei wird der Rundungsfehler $\sigma r := r - r^*$ erzeugt. Wenn d eine komponentenweise Schranke für σr ist, d. h., wenn

$$|\boldsymbol{\delta r}| = |\boldsymbol{r} - \boldsymbol{r^*}| \le \boldsymbol{d} \tag{6}$$

gilt, folgt $|\mathbf{r}^*| \leq |\mathbf{r}| + |\mathbf{r}^* - \mathbf{r}| \leq |\mathbf{r}| + \mathbf{d}$, also $||\mathbf{r}^*|| \leq ||\mathbf{r}| + \mathbf{d}||$. Damit ergibt sich das nachstehende Resultat.

5.4.1. A-posteriori-Abschätzung unter Verwendung von r. Für reguläres $A \in \Re^{n,n}$ und $b, x \in \Re^n$ werde $r, d \in \Re^n$ mit

$$|\mathbf{r} := \mathrm{fl} (\mathbf{b} - A\mathbf{x}), \qquad |\mathbf{r} - (\mathbf{b} - A\mathbf{x})| \leq d$$

berechnet. Dann gilt

$$(\mathbf{A} + \mathbf{\delta}\mathbf{A})\,\mathbf{x} = \mathbf{b} \quad \text{mit} \quad \|\mathbf{\delta}\mathbf{A}\| \leq \Delta \mathbf{A} := \|\,|\mathbf{r}| + \mathbf{d}\|/\|\mathbf{x}\| \tag{7}$$

sowie

$$\|\delta x\| \le \|A^{-1}\| \, \|\, |r| + d\|. \tag{8}$$

Aufwand: $\sim n^2$ opms (für r) + $\sim 2n^2$ ops (für d bei Berechnung nach Ü 2.3.4).

5.4.2. Bemerkung. (i) Der bei der Addition $|\mathbf{r}| + \mathbf{d}$ und der Normbildung in (7), (8) auftretende Rundungsfehler kann vernachlässigt werden, da die beteiligten Vektoren nichtnegative Komponenten besitzen.

(ii) Praktisch wird häufig der Fehlerterm d vernachlässigt. Das kann zu falschen Schranken führen, wenn $|\delta r|$ und d in der Größenordnung von |r| liegen. Man beachte dabei, daß für δr die a-priori-Abschätzung

$$|\boldsymbol{\delta r}| \leq \nu(|\boldsymbol{r}| + n |\boldsymbol{A}| |\boldsymbol{x}|) \tag{9}$$

gilt, siehe Ü 2.3.4.

(iii) Falls das Residuum durch Akkumulation der Skalarprodukte in höherer Genauigkeit v_1 berechnet wird, gilt

$$|\boldsymbol{\sigma}\boldsymbol{r}| \leq \boldsymbol{v} |\boldsymbol{r}| + \boldsymbol{v}_1 n |\boldsymbol{A}| |\boldsymbol{x}|.$$
⁽¹⁰⁾

In diesem Fall ist fast immer $d \approx v |r|$, so daß die Fehlerschranke d weggelassen werden kann.

Um zu einer verbesserten Näherung x' zu kommen, gehen wir von dem Ansatz. $x^* = x + h^*$, d. h.

$$h^* = x^* - x = A^{-1}b - x = A^{-1}(b - Ax) = A^{-1}r^*$$

aus. Die zu x^* führende Korrektur h^* genügt also dem Gleichungssystem

$$Ah^* = r^* \tag{11}$$

mit derselben Koeffizientenmatrix A wie das Ausgangssystem, aber dem Residuum r^* als rechter Seite. Praktisch wird (11) unter Verwendung der bereits vorliegenden Dreiecksfaktorisierung von A mit $O(n^2)$ opms gelöst, wobei statt r^* das berechnete

Residuum r genommen wird. Die als Lösung von Ah = r tatsächlich berechnete Korrektur h weist daher den Fehler

$$h - h^* = h - A^{-1}r^* = (h - A^{-1}r) + A^{-1}\sigma r$$
(12)

auf. Der erste Term rechts entspricht dem bei der Lösung von Ah = r entstandenem Fehler und kann analog zu (1), (2) durch

$$\|\boldsymbol{h} - \boldsymbol{A}^{-1}\boldsymbol{r}\| \leq \varkappa \|\boldsymbol{h}\| \tag{13}$$

abgeschätzt werden. Der zweite Term berücksichtigt den bei der Berechnung von r hervorgerufenen Fehler und läßt sich mit d aus (6) in der Form

$$\|A^{-1}\delta r\| \le \|A^{-1}\| \, \|\delta r\| \le \|A^{-1}\| \, \|d\| \tag{14}$$

abschätzen.

Die berechnete Korrektur **h** kann nun in zweierlei Hinsicht verwendet werden: zur Berechnung der – hoffentlich! – besseren Näherung

$$\boldsymbol{x}' = \boldsymbol{x} + \boldsymbol{h} \tag{15}$$

für x^* mit der aus (12), (13), (14) folgenden Fehlerschranke

$$\|x' - x^*\| = \|h - h^*\| \le \varkappa \|h\| + \|A^{-1}\| \|d\|$$
(16)

oder durch Umstellen von (12) gemäß $x - x^* = -h + (h - A^{-1}r) + A^{-1}\sigma r$ zur Berechnung der – hoffentlich! – besseren Fehlerschranke

$$\|\boldsymbol{\delta x}\| = \|\boldsymbol{x} - \boldsymbol{x^*}\| \le (1 + \varkappa) \|\boldsymbol{h}\| + \|\boldsymbol{A}^{-1}\| \|\boldsymbol{d}\|.$$
(17)

Im Fall

$$\varkappa \ll 1 \quad \text{und} \quad \boldsymbol{d} \ll |\boldsymbol{r}|$$
 (18)

liefert (8) die Schranke

 $\|\delta x\| \leqslant \|A^{-1}\| \|r\|,$

während (17) die i. allg. bessere Schranke $||\delta x|| \leq ||h|| \approx ||A^{-1}r||$ ergibt. In ähnlicher Weise folgt, daß x' dann eine bessere Näherung für x^* als x ist.

B. Iterative Verbesserung

Es liegt nahe, die eben beschriebene Korrektur mit x' anstelle von x zu wiederholen und damit zu einer iterativen Verbesserung zu gelangen.

5.4.3. Iterative Verbesserung der Lösung von Ax = b

Aufgabe: Zu gegebenem $A \in \Re^{n,n}$, $b \in \Re^n$ und gegebener Dreiecksfaktorisierung $P_Z A P_S^{\mathsf{T}} = LR$ von A mit $L, R \in \Re^{n,n}$ und Permutationen P_S, P_Z soll die Lösung des Gleichungssystems Ax = b berechnet und iterativ verbessert werden. Algorithmus:

SO: $\boldsymbol{x}^{0} := \boldsymbol{o}$

for $k := 0(1)k_{\max}$ do

S1: Berechne $\mathbf{r}^{\mathbf{k}} := \mathrm{fl} \left(\mathbf{b} - A \mathbf{x}^{\mathbf{k}} \right)$

S2: Berechne $h = h^k$ als Lösung von $Ah = r^k$ unter Verwendung der Dreiecksfaktorisierung von A aus den Gleichungen

$$Lg = P_Z r^k, \quad R\bar{h} = g, \quad h = P_S^{\mathsf{T}}\bar{h}$$

 $Lg = P_Z r^k, \quad Rar{h} = g$ S3: Setze $x^{k+1} := x^k + h^k$

Aufwand: $\sim 2n^2$ opms pro Schritt $k \ge 1, \sim n^2$ opms für k = 0

Schritt S1 wird dabei wie in Ü2.3.4 realisiert. Im Fall der Spaltenpivotisierung erfordert S2 die Ausführung des "Solvers" 5.2.4 mit r^k statt b; die Dreiecksfaktorisierung muß dazu vorher nach 5.2.3 berechnet worden sein.

Wir bemerken noch, daß im ersten Schritt (k = 0) $\mathbf{r}^0 = \mathbf{b}$ gilt, d. h., \mathbf{x}^1 stimmt mit der nach dem Gaußschen Algorithmus berechneten ersten Näherung für x* überein.

5.4.4. Fehleranalyse. Es seien x^k , r^k , h^k $(k = 1, ..., k_{max})$ die nach Algorithmus 5.4.3 berechneten Größen, und es gelte

$$|\mathbf{r}^{\mathbf{k}} - (\mathbf{b} - A\mathbf{x}^{\mathbf{k}})| \leq \nu |\mathbf{r}^{\mathbf{k}}| + \nu \varphi |A| |\mathbf{x}^{\mathbf{k}}|$$
(19)

sowie

$$(A + \delta \hat{A}_k) h^k = r^k \quad \text{mit} \quad \|\delta \hat{A}_k\| \leq \nu F_1 \|A\|.$$

$$(20)$$

Wenn die Konstanten φ , F_1 der Bedingung

$$\bar{\varkappa} := \nu [F_1 + \varphi + 1] c \text{ ond } (A) \leq 0.2$$

$$\tag{21}$$

genügen, dann gibt es Störungen δA_k , ϑA_k und ξ^k , so daß

$$[A + (\delta A_k + \vartheta A_k)] (x^k + \xi^k) = b$$
⁽²²⁾

gilt mit

$$|\delta A_k| \leq v\varphi |A|, \qquad \|\vartheta A_k\| \leq v\vartheta_k \|A\|, \qquad |\xi^k| \leq v |x^k|, \qquad (23)$$

wobei

$$[ar{x}/(1-2ar{x})] [2arphi+(1+2ar{x}^{k-1}/
u)/ ext{cond} (A)]/(1-ar{x}^k), \quad \hat{x}=ar{x}/(1-ar{x}).$$
 (24)

 $[x/(1 - 2x)][2\psi + (1 + 2x - \psi)]$ conk Der Fehler von x^k genügt der Abschätzung

 $\|oldsymbol{x}^k - oldsymbol{x}^*\| \leq
u[1 + \dot{x}^k/
u + arphi ext{ cond } (A)/(1 - 2ar{oldsymbol{z}})] \|oldsymbol{x}^*\|.$ (25)

Der Beweis ist sehr technisch und soll deshalb übergangen werden. Als Normen können die mit $p = 1, \infty$ genommen werden; für p = 2, F sind etwas andere Konstanten erforderlich, auf deren Angabe wir verzichten.

5.4.5. Bemerkung. (i) Bedingung (19) ist bei einfachgenauer Berechnung von r^k mit $\varphi := n$, bei Akkumulation der Skalarprodukte mit höherer Genauigkeit $v_1 \ll v$ mit $\varphi := n \nu_1 / \nu$ erfüllt, siehe 5.4.2. Bedingung (20) ist mit dem von der Matrix A und der Pivotstrategie abhängigem $F_1 = F_1[A]$ aus Abschnitt 5.3 erfüllt.

(ii) Für genügend großes k folgt aus (24) wegen $\bar{\varkappa} \leq 0.2$

$$\vartheta_k \leq [2\varphi + 1/\text{cond} (A)]/3$$
,

mit (23) also für die Gesamtstörung

$$\|\delta A_k + \vartheta A_k\| \lesssim \nu [(5/3) \varphi + 1/(3 \operatorname{cond} (A))] \|A\|.$$
(26)

Es liegt also numerische Gutartigkeit vor, wobei die Kumulationskonstante nur von φ — der Genauigkeit der Residuenberechnung — abhängt, nicht aber von der durch F_1 charakterisierten Genauigkeit des Gaußschen Algorithmus als Basisverfahren. In analoger Weise folgt aus (25) für genügend großes k

$$\|\boldsymbol{x}^{k} - \boldsymbol{x}^{*}\| \leq \nu [1 + (5/3) \varphi \text{ cond } (\boldsymbol{A})] \, \|\boldsymbol{x}^{*}\|,$$
(27)

d. h., der relative Fehler ist unabhängig von F_1 beschränkt. Bei Berechnung der Residuen in höherer Genauigkeit ist φ cond (A) klein, vgl. wieder 5.2.4, so daß (27) dann praktisch die volle relative (einfache) Genauigkeit von x^k bedeutet.

(iii) Als Abbruchkriterium kann z. B. die Bedingung

 $\|m{h}^k\| > \|m{h}^{k-1}\|/2$

verwendet werden, d. h., es wird abgebrochen, wenn sich die Korrekturen nicht mehr schnell genug verkleinern. \Box

Die iterative Verbesserung kann als eine Ergänzung des Gaußschen Algorithmus angesehen werden, die bei genügend kleinem v, nicht zu schlechter Dreiecksfaktorisierung und nicht zu schlecht konditioniertem A in wenigen Korrekturschritten einen Gesamtprozeß liefert, der in einem besseren Sinne als der Gaußsche Algorithmus numerisch gutartig ist.

C. Abschätzung der Norm der Inversen

Für die Abschätzung des erzeugten Rundungsfehlers gemäß 5.4.1, aber auch für die Abschätzung des durch Meßfehler von A und b hervorgerufenen Fehlers nach Abschnitt 4.1.B bzw. C wird $||A^{-1}||$ benötigt. Da es i. allg. nur auf die Größenordnung ankommt, kann ein Schätzwert mit etwa 50% Fehler durchaus als akzeptabel angesehen werden. Dabei bietet es sich an, auf die Dreiecksfaktorisierung

$$\boldsymbol{P}_{\boldsymbol{Z}}\boldsymbol{A}\boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} = \boldsymbol{L}\boldsymbol{R} \tag{28}$$

von A zurückzugreifen.

Im Abschnitt 6.5 werden wir sehen, daß die explizite Berechnung von A^{-1} aus den Dreiecksfaktoren $\sim 2/3n^3$ opms, also das Doppelte des Aufwandes für die Dreiecksfaktorisierung selbst kostet und daher aus Aufwandsgründen ausscheidet. Wir versuchen daher, mit geringerem Aufwand einen akzeptablen Schätzwert zu bestimmen.

Als Hilfsaufgabe betrachten wir dazu die Berechnung eines Schätzwertes für $\|C^{-1}\|_1$, wenn C eine reguläre Dreiecksmatrix ist. Nach Definition der Matrixnorm gilt

$$\|C^{-1}\| = \max\{\|C^{-1}u\|/\|u\|: u \neq o\} = \max\{\|v\|/\|u\|: u \neq o, Cv = u\}.$$
 (29)

Bei gegebenem u ergibt sich dabei $v = C^{-1}u$ als Lösung des Dreieckssystems Cv = uim Fall einer unteren Dreiecksmatrix gemäß

$$v_{1} = \frac{u_{1}}{c_{11}},$$

$$v_{2} = -\frac{c_{21}}{c_{22}}v_{1} + \frac{u_{2}}{c_{22}},$$

$$v_{3} = -\frac{c_{31}}{c_{33}}v_{1} - \frac{c_{32}}{c_{33}}v_{2} + \frac{u_{3}}{c_{33}},$$

$$\vdots$$

$$v_{i} = -\frac{c_{i1}}{c_{ii}}v_{1} - \frac{c_{i2}}{c_{ii}}v_{2} - \frac{c_{i3}}{c_{ii}}v_{3} - \dots - \frac{c_{i,i-1}}{c_{ii}}v_{i-1} + \frac{u_{i}}{c_{ii}},$$

$$\vdots$$

$$v_{n} = -\frac{c_{n1}}{c_{nn}}v_{1} - \frac{c_{n2}}{c_{nn}}v_{2} - \frac{c_{n3}}{c_{nn}}v_{3} - \dots - \frac{c_{n,n-1}}{c_{nn}}v_{n-1} + \frac{u_{n}}{c_{nn}}.$$
(30)

Wir wählen jetzt $u = (u_k)$ als Vektor mit den Komponenten $u_k = \pm 1$ und verteilen die Vorzeichen so, daß $\|v\|_1 = \sum_{i=1}^n |v_i|$ möglichst groß wird. Im ersten Schritt wird $u_1 = 1$ gesetzt. Damit sind v_1 und die in der ersten Spalte nach dem Gleichheitszeichen stehenden Summanden $-(c_{i1}/c_{ii}) v_1$ von v_i (i = 2, ..., n) festgelegt. Im zweiten Schritt wird für jede der beiden Möglichkeiten $u_2 = u_2^{+,-} = \pm 1$ das zugehörige $v_2 = v_2^{+,-} = -(c_{21}/c_{22}) v_1 + u_2^{+,-}/c_{22}$ berechnet und das Vorzeichen in Abhängigkeit davon festgelegt, welche der beiden Summen

$$|v_2^+, -| + \sum_{i=3}^n \left| -\frac{c_{i1}}{c_{ii}} v_1 - \frac{c_{i2}}{c_{ii}} v_2^+ - \right|$$

die größere ist. Diese Summen werden als Näherung für $|v_2| + \sum_{i=3}^n |v_i|$ angesehen,

wobei die v_3, \ldots, v_n durch die aus den ersten beiden Spalten gebildeten Teilsummen ersetzt worden sind. Man beachte, daß v_1 bereits festgelegt ist und auf den Größenvergleich keinen Einfluß mehr hat. Vor dem k-ten Schritt sind in dieser Weise u_1, \ldots, u_{k-1} sowie v_1, \ldots, v_{k-1} fixiert. Das Vorzeichen von $u_k = \pm 1$ wird dann so gewählt, daß die aus den ersten k Spalten gebildeten Teilsummen der v_i $(i \ge k)$ betragsmäßig summiert einen möglichst großen Wert ergeben. Dieses Vorgehen führt auf den folgenden Algorithmus:

5.4.6. Heuristische Maximierung von $\|C^{-1}u\|_1/\|u\|_1$. Aufgabe: Für die reguläre Dreiecksmatrix $C = (c_{ij}) \in \Re^{n,n}$, $c_{ii} \neq 0$ (i = 1, ..., n)werden $u = (u_i) \in \Re^n$ mit $u_i = \pm 1$ und $v = (v_i) = C^{-1}u$ so bestimmt, daß $\|v\|_1 = \sum_{i=1}^n |v_i|$ möglichst groß wird.

 $\begin{array}{l} Algorithmus \ (falls \ C \ untere \ Dreiecksmatrix; \ für \ obere \ Dreiecksmatrix \ analog): \\ u_{1} := 1, v_{1} := 1/c_{11}, \ for \ i := 2(1)n \ do \ v_{i} := -c_{i1} * v_{1}/c_{ii} \\ for \ k := 2(1)n \ co \\ & \left| \begin{array}{c} z := 1/c_{kk}, u_{k} := v_{k} - z, v_{k} := v_{k} + z \\ SM := |u_{k}|, \ SP := |v_{k}| \\ for \ i := k \ + 1(1)n \ do \\ & \left| \begin{array}{c} z := c_{ik}/c_{ii}, u_{i} := v_{i} - z * u_{k}, v_{i} := v_{i} - z * v_{k} \\ SM := SM \ + |u_{i}|, \ SP := SP \ + |v_{i}| \\ if \ SM \ SP \ then \ [u_{k} := -1, \ for \ i := k(1)n \ do \ v_{i} := u_{i}] \\ & else \ u_{k} := 1 \end{array} \right| \\ Autward: \sim 2n^{2} \ ops \ + \sim 3/2n^{2} \ opm \end{array}$

Für die in 5.4.6 berechneten Vektoren u, v gilt nach (29)

$$\frac{\|\boldsymbol{v}\|_{1}}{\|\boldsymbol{u}\|_{1}} \leq \|C^{-1}\|_{1}, \tag{31}$$

und nach Konstruktion kann erwartet werden, daß der Quotient $||v||_1/||u||_1$ in der Größenordnung von $||C^{-1}||_1$ liegt.

Wir kehren jetzt wieder zum Ausgangsproblem zurück, eine Näherung ζ für $\zeta^* = ||A^{-1}||$ zu berechnen. Aus (28) und der Invarianz der Matrixnormen mit $p = 1, 2, \infty, F$ gegenüber Zeilen- und Spaltenvertauschungen folgt

$$\|A^{-1}\| = \|(LR)^{-1}\|.$$

Für die ∞ -Norm gilt dann

$$\zeta_{\infty}^{*} = \|A^{-1}\|_{\infty} = \|A^{-T}\|_{1} = \|(LR)^{-T}\|_{1} = \max\{\|z\|_{1}/\|x\|_{1} : x \neq o, (LR)^{T} z = x\}.$$
(32)

Wenn \boldsymbol{z} bei gegebenem \boldsymbol{x} in der Form

$$\boldsymbol{R}^{\mathsf{T}}\boldsymbol{y} = \boldsymbol{x}, \qquad \boldsymbol{L}^{\mathsf{T}}\boldsymbol{z} = \boldsymbol{y} \tag{33}$$

aus zwei Dreieckssystemen berechnet wird, ergibt sich

$$\frac{\|\boldsymbol{z}\|_{1}}{\|\boldsymbol{x}\|_{1}} = \frac{\|\boldsymbol{z}\|_{1}}{\|\boldsymbol{y}\|_{1}} \cdot \frac{\|\boldsymbol{y}\|_{1}}{\|\boldsymbol{x}\|_{1}} = \frac{\|(\boldsymbol{L}^{\mathsf{T}})^{-1} \boldsymbol{y}\|_{1}}{\|\boldsymbol{y}\|_{1}} \cdot \frac{\|(\boldsymbol{R}^{\mathsf{T}})^{-1} \boldsymbol{x}\|_{1}}{\|\boldsymbol{x}\|_{1}} \leq \zeta_{\infty}^{*}.$$
(34)

Da der Maximierungsalgorithmus nur auf einen der beiden Quotienten $||\boldsymbol{z}||_1/||\boldsymbol{y}||_1$ oder $||\boldsymbol{y}||_1/||\boldsymbol{x}||_1$ angewendet werden kann, muß überlegt werden, welcher der beiden Dreiecksfaktoren $\boldsymbol{L}^{\mathsf{T}}$ bzw. $\boldsymbol{R}^{\mathsf{T}}$ stärker zu $(\boldsymbol{L}\boldsymbol{R})^{-\mathsf{T}}$ beiträgt, d. h., in welchem sich die Kondition von \boldsymbol{A} vorrangig niederschlägt. Dies hängt selbstverständlich von der Art der Dreiecksfaktorisierung (28) ab.

Wir betrachten zunächst den Fall der Spalten- bzw. vollständigen Pivotisierung oder den einer spaltendiagonaldominanten Matrix A. Hier gilt

$$|l_{ij}| \leq 1 \qquad (i > j), \tag{35}$$

und es folgt

$$\| \boldsymbol{L}^{\intercal} \|_1 \leq n, \qquad \| (\boldsymbol{L}^{\intercal})^{-1} \|_1 = \| \boldsymbol{L}^{-1} \|_{\infty} \leq 2^{n-1}$$

unabhängig von der Matrix A, siehe Ü 5.4.1. Eine beliebig schlechte Kondition der Matrix A kann sich also nicht im Faktor L niederschlagen. Praktisch ist $||L^{-1}||_{\infty}$ meist viel kleiner als der maximal mögliche Wert 2^{n-1} , so daß der erste Quotient in (34) wegen

$$\frac{1}{\|\boldsymbol{L}^{\mathsf{T}}\|_{1}} \leq \frac{\|\boldsymbol{L}^{-\mathsf{T}}\boldsymbol{y}\|_{1}}{\|\boldsymbol{y}\|_{1}} \leq \|\boldsymbol{L}^{-\mathsf{T}}\|_{1}$$
(36)

nur wenig variiert. In diesem Fall sollte also der zweite Quotient $\|y\|_1/\|x\|_1$ maximiert werden.

Im Fall der Zeilen- bzw. vollständigen Pivotisierung oder dem einer zeilendiagonaldominanten Matrix A gilt

$$|r_{ij}| \le |r_{ii}| \qquad (i < j), \tag{37}$$

vgl. (5.2.7). Wenn hier

$$\boldsymbol{LR} = (\boldsymbol{LD}) (\boldsymbol{D}^{-1}\boldsymbol{R}) =: \boldsymbol{L}_{1}\boldsymbol{R}_{1} \quad \text{mit} \quad \boldsymbol{D} := \text{diag} (r_{ii})$$
(38)

gesetzt wird, ist R_1 wegen (37) eine obere Einsdreiecksmatrix mit

$$|(\boldsymbol{R}_1)_{ij}| \leq 1 \qquad (i < j),$$
(39)

d. h., $\mathbf{R}_1^{\mathsf{T}}$ hat qualitativ dieselben Eigenschaften wie \mathbf{L} im Fall der Spaltenpivotisierung, und eine mögliche schlechte Kondition schlägt sich in \mathbf{L}_1 nieder. In der zu (33), (34) analogen Darstellung mit \mathbf{L}_1 statt \mathbf{L} und \mathbf{R}_1 statt \mathbf{R} sollte daher der erste Quotient $\|\mathbf{z}\|_1 / \|\mathbf{y}\|_1 = \|\mathbf{L}_1^{-\mathsf{T}}\mathbf{y}\|_1 / \|\mathbf{y}\|_1$ maximiert werden.

Im Fall einer symmetrischen und positiv definiten Matrix A gilt $r_{ii} > 0$ und

$$LR = (LD^{1/2}) (D^{-1/2}R) = : L_2R_2 \text{ mit } D^{1/2} := \text{diag} (\sqrt{r_{ii}}), R_2 = L_2^{\mathsf{T}}, (40)$$

vgl. 5.2.10. Hier kann wie im Fall der Spaltenpivotisierung vorgegangen werden, wobei L durch L_2 und R durch $R_2 = L_2^{\mathsf{T}}$ zu ersetzen ist; siehe Ü 5.4.2. Zusammenfassend erhalten wir die folgenden Vorschriften:

5.4.7. Schätzung von $||(LR)^{-1}||_{\infty}$.

Aufgabe: Die Matrix A sei gemäß (28) durch eine untere Einsdreiecksmatrix Lund eine reguläre obere Dreiecksmatrix R gegeben. Aus den Dreiecksfaktoren L, R ist ein Schätzwert ζ_{∞} für $\zeta_{\infty}^* := ||(LR)^{-1}||_{\infty}$ zu berechnen.

Algorithmus: Bestimme x, z mit $(LR)^{\intercal} z = x$ wie folgt:

Fall 1 (Spalten- oder vollständige Pivotisierung; A spaltendiagonaldominant):

S1.1: Führe 5.4.6 mit $C := \mathbf{R}^{\intercal}$ aus, setze $\mathbf{x} := \mathbf{u}, \mathbf{y} := \mathbf{v}$

S2.1: Berechne \boldsymbol{z} aus $\boldsymbol{L}^{\mathsf{T}}\boldsymbol{z} = \boldsymbol{y}$

Fall 2 (Zeilen- oder vollständige Pivotisierung; A zeilendiagonaldominant):

S1.2: Setze $\boldsymbol{D} := \text{diag}(r_{ii}), \boldsymbol{L}_1 := \boldsymbol{L}\boldsymbol{D}, \boldsymbol{R}_1 := \boldsymbol{D}^{-1}\boldsymbol{R},$ führe 5.4.6 mit $\boldsymbol{C} := \boldsymbol{L}_1^{\mathsf{T}}$ aus, setze $\boldsymbol{y} := \boldsymbol{u}, \boldsymbol{z} := \boldsymbol{v}$

S2.2: Berechne $\boldsymbol{x} := \boldsymbol{R}_1^{\mathsf{T}} \boldsymbol{y}$
Fall 3 (A symmetrisch und positiv definit):

S1.3: Setze $D^{1/2} := \text{diag}(\sqrt[3]{r_{ii}}), L_2 := LD^{1/2}, R_2 := D^{-1/2}R = L_2^{\mathsf{T}},$ führe 5.4.6 mit $C := L_2$ aus, setze x := u, y := vS2.3: Berechne z aus $L_2^{\mathsf{T}} z = y$ S3: Setze $\zeta_{\infty} := ||z||_1/||x||_1.$ Aufwand: $\sim 5/2n^2 \text{ ops} + \sim 2n^2 \text{ opm}$ im Fall 1, in den Fällen 2 und 3 zusätzlich $\sim n^2/2$ opm zur Bildung von L_1 bzw. L_2 , falls die Dreiecksfaktorisierung nicht so modifiziert worden ist, daß sie diese Faktoren direkt liefert.

5.4.8. Bemerkung. (i) Wenn z wie in 5.4.7 bestimmt und daraus $w = (LR)^{-1} z$ gemäß

$$Lt = z, \qquad Rw = t \tag{41}$$

berechnet wird, gilt

$$\zeta'_{\infty} := \|\boldsymbol{w}\|_{\infty} / \|\boldsymbol{z}\|_{\infty} \le \|(\boldsymbol{L}\boldsymbol{R})^{-1}\|_{\infty}, \tag{42}$$

und i. allg. wird auch ζ'_{∞} in der Größenordnung von ζ^*_{∞} liegen. Die Größe

$$\hat{\zeta}_\infty := \max\left\{\zeta_\infty, \zeta'_\infty
ight\} = \max\left\{\|oldsymbol{z}\|_1/\|oldsymbol{x}\|_1, \|oldsymbol{w}\|_\infty/\|oldsymbol{z}\|_\infty
ight\} \leq \zeta^*_\infty$$

ist daher i. allg. eine bessere Schätzung für ζ_{∞}^{*} als jede der beiden einzelnen Werte ζ_{∞} bzw. ζ_{∞}' . Der Mehraufwand zur Lösung der beiden Dreieckssysteme (41) ist $\sim n^2$ opms, also klein. Umfangreiche Experimente zeigen, daß $\hat{\zeta}_{\infty}$ fast immer einen relativen Fehler von weniger als 50% aufweist, siehe B 5.6.

(ii) Obwohl z und x mit dem Ziel der Maximierung von $||z||_1/||x||_1$ als Näherung für ζ_{∞}^{*} konstruiert worden sind, stellen die Zahlen

$$\zeta_1 := \| \boldsymbol{z} \|_{\infty} / \| \boldsymbol{x} \|_{\infty}$$
 und $\zeta_1' := \| \boldsymbol{w} \|_1 / \| \boldsymbol{z} \|_1$

meist gute Schätzwerte für $\|(\boldsymbol{LR})^{-1}\|_1$ dar. Ebenso sind

 $\zeta_2 := \|\boldsymbol{z}\|_2 / \|\boldsymbol{x}\|_2$ und $\zeta'_2 := \|\boldsymbol{w}\|_2 / \|\boldsymbol{z}\|_2$

meist gute Näherungen für $||(LR)^{-1}||_2$.

(iii) Algorithmus 5.4.6 kann leicht so modifiziert werden, da β C in der faktorisierten Form C = D'C'E' mit regulären Diagonalmatrizen $D' = \text{diag}(d'_i), E'$ = diag (e'_i) akzeptiert wird; man braucht bloß c_{ij} überall durch $d'_i c'_{ij} e'_j$ zu ersetzen. Der Mehraufwand ist $\sim n^2/2$ opm, im Fall E' = I sogar nur *n* opm.

(iv) Wenn eine mittels Zeilenpivotisierung berechnete Dreiecksfaktorisierung $AP_{S}^{\intercal} = LR$ von A bekannt ist, kann die in (iii) beschriebene Modifikation zur a-posteriori-Abschätzung von $\|(\hat{D}A)^{-1}\|_{\infty}$ verwendet werden, wobei $\hat{D} = \text{diag}(\hat{d}_i)$ eine reguläre Skalierungsmatrix bezeichnet. Diese Aufgabe tritt bei der in bezug auf die Fehlerabschätzung optimalen Skalierung von A auf, siehe 4.1.15 und 4.1.16. In 5.4.7, Fall 2, ist dazu L bzw. L_1 durch $\hat{D}L$ bzw. $\hat{D}L_1$ zu ersetzen, d. h., man erhält

S1: Führe 5.4.6 in der gemäß (iii) modifizierten Form mit $C := DL^{\dagger}D = L_{1}^{\dagger}D$ durch, setze $\boldsymbol{y} := \boldsymbol{u}, \boldsymbol{z} := \boldsymbol{v}$

- S2: Berechne $\boldsymbol{x} := \boldsymbol{R}^{\mathsf{T}} \boldsymbol{D}^{-1} \boldsymbol{y} = \boldsymbol{R}_{1}^{\mathsf{T}} \boldsymbol{y}$
- S3: Setze $\zeta_{\infty,\hat{D}} := \|\boldsymbol{z}\|_1 / \|\boldsymbol{x}\|_1$

Dann ist $\zeta_{\infty,\hat{D}}$ ein Schätzwert für $\zeta_{\infty,\hat{D}}^* := \|(\hat{D}A)^{-1}\|_{\infty}$. Der Aufwand beträgt $\sim 5/2n^2$ (ops + opm).

 (\mathbf{v}) Die tatsächlich berechneten Dreiecksfaktoren genügen statt (28) der Beziehung

$$P_Z(A + \delta A) P_S^{\mathsf{T}} = LR_s$$

siehe 5.3.1, so daß ζ_{∞} in Wirklichkeit eine Schätzung für $\|(A + \delta A)^{-1}\|_{\infty}$ ist. Für genügend kleines δA — dies ist in der Regel der Fall — gilt dann nach 4.1.2

$$\|(A + \delta A)^{-1}\|_{\infty} = \|A^{-1}\|_{\infty} + O(\|\delta A\|_{\infty}),$$

so daß die Verfälschung durch σA i. allg. vernachlässigt werden kann. \Box

Übungsaufgaben

Ü 5.4.1. Es sei $L = (l_{ij}) = I - G$ eine untere Einsdreiecksmatrix. Man beweise die Gültigkeit von

(i)
$$G^n = O$$
 und $L^{-1} = I + \sum_{k=1}^{n-1} G^k$.

(ii) $||\boldsymbol{L}||_p \leq ||\boldsymbol{L}^*||_p = n$ und $||\boldsymbol{L}^{-1}||_p \leq ||(\boldsymbol{L}^*)^{-1}||_p = 2^{n-1}$ für alle L mit $|l_{ij}| \leq 1$, wobei $\boldsymbol{L}^* = (l_{ij}^*)$ die untere Einsdreiecksmatrix mit $l_{ij}^* := -1$ (i > j) und $p = 1, \infty$ ist. Was ergibt sich im Fall $|l_{ij}| \leq q$ (i > j)?

Hinweis: In (i) beachte man Ü 4.1.1, in (ii) zeige man zunächst $|L^{-1}| \leq (L^*)^{-1}$.

Ü 5.4.2. Es sei A symmetrisch und positiv definit mit der symmetrischen Dreiecksfaktorisierung $A = L_2 L_2^{T}$, L_2 wie in 5.4.7. Man zeige, daß

$$\|A\|_2 = (\|L_2\|_2)^2 = (\|L_2^\mathsf{T}\|_2)^2 \quad ext{sowie} \quad \|A^{-1}\|_2 = (\|L_2^{-1}\|_2)^2 = (\|L_2^{-\mathsf{T}}\|_2)^2$$

gilt.

Hinweis: Man beachte Ü 1.2.9.

Bemerkungen zum Kapitel 5

B 5.1. Das heute nach GAUSS benannte Eliminationsverfahren war bereits vor etwa 2000 Jahren im alten China bekannt. Die Beschreibung mittels elementarer Transformationsmatrizen geht auf TURING [48] zurück und ist seit WILKINSON [61, 63, 65] Standard in der numerischen linearen Algebra. Eine leicht verständliche Beschreibung haben FORSYTHE/ MOLER [67] gegeben. Ausführliche Darstellungen einschließlich historischer Kommentare finden sich z. B. bei FADDEEV/FADDEEVA [63], HOUSEHOLDER [64] und STEWART [73].

B 5.2. Trotz seiner einfachen Struktur stellt der Gaußsche Algorithmus einen schwer analysierbaren numerischen Prozeß dar, und die Frage nach einer in bezug auf den erzeugten Rundungsfehler optimalen Pivotstrategie ist nach wie vor ungeklärt. Die erste adäquate Rundungsfehleranalyse wurde für den symmetrischen und definiten Fall durch von NEUMANN/ GOLDSTINE [47] gegeben, allerdings nur für die damals übliche Festpunktarithmetik. Grundlegende Ergebnisse für den allgemeinen Fall — auch für Gleitpunktarithmetik — gehen auf WILKINSON [61, 63, 65] zurück. **B 5.3.** Von besonderer Bedeutung sind die von WILKINSON erhaltenen Aussagen über die numerische Gutartigkeit des Gaußschen Algorithmus in Abhängigkeit von den betrachteten Matrizenklassen und den verwendeten Pivotisierungsstrategien. Die konkreten Normschranken für die äquivalenten Störungen sind dabei mehr als eine Folgerung anzusehen, entscheidend ist die Existenz von solchen Störungsschranken.

B 5.4. Detaillierte Untersuchungen verschiedener Pivotisierungsstrategien werden bis in die jüngste Zeit durchgeführt. Aus der Vielzahl der Arbeiten seien VAN DER SLUIS [70] und SKEEL [79, 80, 81] erwähnt.

B 5.5. Die Rundungsfehleranalyse zur iterativen Verbesserung wurde erstmals von Wilkinson durchgeführt; weitere Arbeiten zu dieser Thematik sind u. a. Moler [67], JANKOWSKI/WOŹNIAKOWSKI [77] und SKEEL [80].

B 5.6. Es gibt eine Vielzahl von Arbeiten zur exakten Abschätzung von $||C^{-1}||$ im Fall einer regulären Dreiecksmatrix C, wir zitieren etwa LEMEIRE [75]. Die Schranken sind allerdings meist zu pessimistisch bzw. zu aufwendig zu berechnen. Die im Abschnitt 5.4 angegebenen heuristischen Methoden, die auf CLINE et al. [79] zurückgehen, brachten hier eine Wende, indem sie mit vertretbarem Aufwand die Berechnung akzeptabler unterer Schranken zulassen. Numerische Tests mit zufällig erzeugten Testmatrizen haben STEWART [80] und O'LEARV [80] durchgeführt. In der Originalarbeit von CLINE et al. wird ζ' gemäß 5.4.8 (i) als Schätzer verwendet, während O'LEARV ζ und $\hat{\zeta}$ empfichlt. Auf die Möglichkeit der a-posteriori-Abschätzung von $||(\hat{D}A)^{-1}||_{\infty}$ unter Verwendung einer zeilenpivotisierten LR-Faktorisierung von A gemäß 5.4.8 (iv) scheint bisher nicht hingewiesen worden zu sein. Für weitere Resultate sei auf GRIMES/LEWIS [81] und CLINE et al. [82] verwiesen. In gewissen Ausnahmefällen können die Konditionsschätzer auch unbrauchbare Werte liefern, siehe CLINE/REW [83].

6. Modifikationen des Gaußschen Algorithmus

Gegenstand dieses Kapitels sind Modifikationen des Gaußschen Algorithmus, die sich von der (n-1)-stufigen Grundform aus 5.1 bzw. 5.2 einmal durch die Ausnutzung spezieller Eigenschaften der Matrix A wie Symmetrie oder schwache Besetztheit, zum anderen durch die Art der Berechnung der Elemente der Dreiecksfaktoren L und R unterscheiden. Außerdem gehen wir auf die Berechnung inverser Matrizen ein.

6.1. Dreiecksfaktorisierungen symmetrischer Matrizen

Es sei $A \in S^{n,n}$ eine symmetrische Matrix, die durch die Elemente im unteren Dreieck gegeben sei; zur Speicherung symmetrischer Matrizen siehe Bemerkung 6.1.14 am Schluß dieses Abschnitts.

Wir suchen im folgenden Modifikationen des Gaußschen Algorithmus, bei denen sich die Symmetrie von A in einer Halbierung des Aufwands gegenüber dem nichtsymmetrischen Fall niederschlägt, d. h., die mit $\sim n^3/6$ opms auskommen und auf nur einem Dreieck von A mit $\sim n^2/2$ S realisiert werden können. Die Möglichkeiten zur Ausnutzung der Symmetrie hängen dabei wesentlich davon ab, ob A definit oder indefinit ist. Im definiten Fall genügt es, sich auf positiv definites A zu beschränken, denn für negativ definites A' ist A := -A' positiv definit.

A. Positiv definite symmetrische Matrizen

In Abschnitt 5.2 wurde gezeigt, daß der Gaußsche Algorithmus für eine positiv definite Matrix $A \in S^{n,n}$ ohne Pivotisierung durchführbar ist. Dabei gilt

$$L^{\mathsf{T}} = D^{-1}R$$
 mit $D = \text{diag}(d_k), \quad d_k := r_{kk} > 0$ $(k = 1, ..., n), (1)$

d. h., die Faktorisierung A = LR kann in der symmetrischen Form

$$A = LDL^{\mathsf{T}} \tag{2}$$

mit der Einsdreiecksmatrix L und der positiven Diagonalmatrix D geschrieben werden.

Die im k-ten Schritt

$$A^{(k+1)} := L_k A^{(k)} \qquad (k = 1, ..., n-1), \qquad A^{(1)} := A$$
(3)

gemäß der Partitionierung

$$A^{(k)} = \begin{pmatrix} a_{11}^{(k)} \dots a_{1n}^{(k)} \\ \vdots \\ a_{kk}^{(k)} \dots a_{kn}^{(k)} \\ O \vdots \vdots \\ a_{nk}^{(k)} \dots a_{nn}^{(k)} \end{pmatrix} = : \begin{pmatrix} R^{(k)} \\ O \\ M^{(k)} \end{pmatrix}_{n}^{k}$$
(4)

vorkommenden quadratischen Restmatrizen $M^{(k)}$ und $M^{(k+1)}$ sind außerdem wie $M^{(1)} = A$ symmetrisch und positiv definit, siehe 5.2.9. Von den Elementen

$$a_{ij}^{(k+1)} = a_{ji}^{(k+1)} := a_{ij}^{(k)} - l_{ik} * a_{kj}^{(k)} \qquad (i, j = k+1, ..., n)$$

$$(5)$$

der Matrix $M^{(k+1)}$ brauchen daher nur die im unteren Dreieck berechnet zu werden; die übrigen entstehen aus diesen durch Spiegelung an der Hauptdiagonalen. Der Gesamtprozeß kann dann auf dem Platz des unteren Dreiecks von A ausgeführt werden, wobei die erzeugten Nullen mit den Eliminationskoeffizienten

$$l_{ik} := a_{ik}^{(k)} / a_{kk}^{(k)} \qquad (i = k + 1, ..., n)$$
(6)

überspeichert werden. Bei spaltenweiser Berechnung der $a_{ij}^{(k+1)}$ müssen dabei die in (5) benötigten Elemente $a_{kj}^{(k)}$ (j = k + 1, ..., n) — dies sind gerade die Elemente $r_{k,k+1}, ..., r_{kn}$ der k-ten Zeile von \mathbf{R} — in einem Feld \mathbf{r} der Länge n zwischengespeichert werden. Der sonst für \mathbf{R} benötigte Speicherplatz von $\sim n^2/2$ S reduziert sich damit auf $\sim n$ S für \mathbf{r} . Eine zu 5.2.3 analoge Computerrealisierung des beschriebenen Verfahrens lautet wie folgt; mit $\mathfrak{S}^{n,n}$ bezeichnen wir dabei die Menge der symmetrischen Matrizen aus $\mathfrak{R}^{n,n}$:

6.1.1. LDL^T-F iktorisierung einer symmetrischen und positiv definiten Matrix

Aufgabe: Für die positiv definit Matrix $A \in \mathbb{S}^{n,n}$ ist die Dreiecksfaktorisierung $A = LDL^{\intercal}$ zu berechnen. Die Matrix A ist durch die Elemente des unteren Dreiecks gegeben, und die signifikanten Elemente von $L, D \in \mathbb{R}^{n,n}$ sind auf dem Platz des unteren Dreiecks von A zu speichern. Wenn $d_k = r_{kk} \leq 0$ für ein $k \in \{1, ..., n\}$

gilt, soll abgebrochen und ie = -1 gesetzt werden; andernfalls wird ie = 0 gesetzt, und es gilt $d_k > 0$ (k = 1, ..., n).

setzt, und es gilt $d_k > 0$ (k = 1, ..., n). Algorithmus: ie := 0for k := 1(1)n - 1 do $\begin{vmatrix} \text{if } a_{kk} \leq 0 \text{ then } [ie := -1, \text{stop}] \\ \text{for } i := k + 1(1)n \text{ do } [r_i := a_{ik}, a_{ik} := r_i/a_{kk}] \\ \text{for } j := k + 1(1)n \text{ do } a_{ij} := a_{ij} - a_{ik} * r_j \\ \text{if } a_{nn} \leq 0 \text{ then } [ie := -1, \text{stop}] \\ Aufwand : \sim n^3/6 \text{ opms}, \sim n \text{ S } (\text{für } r_2, ..., r_n) \end{vmatrix}$

Für n = 4 ergibt sich bei erfolgreicher Faktorisierung (ie = 0) das folgende Belegungsmuster:

$$\left. \begin{array}{c} \left. \begin{array}{c} a_{11} \\ a_{21} \\ a_{22} \\ a_{31} \\ a_{32} \\ a_{33} \\ a_{41} \\ a_{42} \\ a_{43} \\ a_{44} \end{array} \right| \rightarrow \left(\begin{array}{c} r_{12} \\ r_{23} \\ r_{34} \end{array} \right) \left(\begin{array}{c} \left. \begin{array}{c} d_1 \\ l_{21} \\ d_2 \\ l_{31} \\ l_{32} \\ l_{41} \\ l_{42} \\ l_{43} \\ d_4 \end{array} \right) \right. \right) \right.$$

6.1.2. Bemerkung. (i) In exakter Arithmetik ist 6.1.1 für eine symmetrische Matrix A genau dann mit $d_k > 0$ — d. h. ie = 0 — durchführbar, wenn A positiv definit ist, vgl. 5.2.10 und Ü 5.2.7. Das Verfahren 6.1.1 kann daher zur konstruktiven Überprüfung der positiven Definitheit einer symmetrischen Matrix A verwendet werden und liefert bei positiver Antwort gleichzeitig die dann existierende Faktorisierung $A = LDL^{T}$ mit $d_k > 0$. Analoge Aussagen gelten für Computerrechnung, sofern A nicht zu schlecht konditioniert ist, siehe 6.1.3 unten.

(ii) Bei zeilenweiser Berechnung der Elemente des unteren Dreiecks von $M^{(k-1)}$ kann auf das Hilfsfeld r verzichtet werden. Die letzten drei Zeilen der Laufanweisung für k sind dann wie folgt abzuändern:

for i := k + 1(1)n do

$$\varrho := a_{ik}, a_{ik} := \varrho/a_{kk}$$

for $j := k + 1(1)i$ do $a_{ij} := a_{ij} - \varrho * a_{jk}$

Man beachte dabei, daß (5) in der Form

$$a_{ij}^{(k+1)} = a_{ij}^{(k)} - l_{ik} * a_{jk}^{(k)} = a_{ij}^{(k)} - a_{ik}^{(k)} * l_{jk} \qquad (k+1 \le j \le i)$$

geschrieben werden kann. Statt ~ n S für r wird hier nur noch 1 S für ϱ zum Zwischenspeichern von $a_{ik}^{(k)} = a_{ki}^{(k)} = r_{ki}$ benötigt.

(iii) Die Lösung des Gleichungssystems $Ax = LDL^{\intercal}x = b$ kann für jedes $b \in \mathbb{R}^n$ analog zu 5.2.4 unter Verwendung der LDL^{\intercal} -Faktorisierung in drei Schritten durch Lösung von

S2.1: Lc = b, S2.2: Dd = c, S2.3: $L^{T}x = d$

erfolgen. Die Determinante ergibt sich analog zu 5.2.5 gemäß

$$\det (A) = \prod_{k=1}^{n} d_k. \quad \Box$$

Zur Rundungsfehleranalyse von 6.1.1 können die Ergebnisse von 5.3 herangezogen werden, vgl. 5.3.2. Allerdings ist die dort angewendete Beweistechnik nicht auf den Fall symmetrischer Matrizen ausgerichtet, und die äquivalenten Störungen δA bzw. $\delta_1 A$ sind i. allg. nicht symmetrisch, was der Aufgabenklasse nicht angepaßt ist. Mit feineren und aufwendigeren Mitteln lassen sich dagegen auch symmetrische Störungen konstruieren und bessere Schranken für die Kumulationskonstanten herleiten. Aus Platzgründen geben wir nur die Ergebnisse an.

6.1.3. Rundungsfehleranalyse. Die **LDL^T**-Faktorisierung 6.1.1 ist für die positiv definite Matrix $A \in \mathfrak{S}^{n,n}$ mit $d_k > 0$ (k = 1, ..., n) durchführbar, sofern

$$vF \|A^{-1}\|_2 \|A\|_F \leq 0.5 \quad \text{mit} \quad F := n + 2 \ln(n/2) \sim n$$
(7)

gilt. Zu den berechneten Faktoren L, D existiert eine Störung $\delta A \in S^{n,n}$ mit

$$A + \delta A = LDL^{\intercal} \quad \text{und} \quad \|\delta A\|_F \leq \nu F \, \|A\|_F. \tag{8}$$

Die zu $\boldsymbol{b} \in \Re^n$ gehörende und gemäß 6.1.2(iii) berechnete Lösung $\boldsymbol{x} \in \Re^n$ von $A\boldsymbol{x} = \boldsymbol{b}$ ist die exakte Lösung des durch die Matrix $\boldsymbol{\delta}_1 \boldsymbol{A} \in \mathbf{S}^{n,n}$ gestörten Systems

$$(\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \, \boldsymbol{x} = \boldsymbol{b}, \quad \text{wobei} \quad \|\boldsymbol{\delta}_1 \boldsymbol{A}\|_p \leq v F_1 \, \|\boldsymbol{A}\|_p \tag{9}$$

gilt mit

$$F_1 := \begin{cases} 2n^{3/2} + F \sim 2n^{3/2} & \text{für } p = F, \\ n^{3/2} + n + n^{1/2}F \sim 2n^{3/2} & \text{für } p = 2. \end{cases}$$
(10)

Die LDL^{\intercal} -Faktorisierung einer symmetrischen definiten Matrix und die Lösung von Ax = b mittels der Faktoren L und D sind also numerisch gutartige Prozesse.

Eine weitere symmetrische Faktorisierung kann erhalten werden, wenn

$$\hat{\boldsymbol{L}} := \boldsymbol{L} \boldsymbol{D}^{1/2} \quad \text{mit} \quad \boldsymbol{D}^{1/2} := \text{diag}\left(\sqrt[]{d_k}\right)$$

$$\tag{11}$$

gesetzt wird. Die Matrix \hat{L} hat dann die Gestalt

$$\hat{L} = \begin{pmatrix} \hat{l}_{11} & & \\ \hat{l}_{21} & \hat{l}_{22} & \\ \vdots & \vdots & \ddots & \\ \hat{l}_{n1} & \hat{l}_{n2} & \dots & \hat{l}_{nn} \end{pmatrix} \quad \text{mit} \quad \hat{l}_{ik} = l_{ik} \sqrt[n]{d_k} \quad (i = k, \dots, n), \tag{12}$$

und (2) geht in $A = LD^{1/2}D^{1/2}L^{\mathsf{T}} = (LD^{1/2}) (LD^{1/2})^{\mathsf{T}}$, also

$$\boldsymbol{A} = \boldsymbol{\hat{L}} \boldsymbol{\hat{L}}^{\mathsf{T}} \tag{13}$$

über. Die symmetrische Faktorisierung (13) heißt Cholesky- oder LL^r-Faktorisierung

der symmetrischen und positiv definiten Matrix A. Wegen (12) und (6) gilt

$$\hat{l}_{kk} = \sqrt[4]{a_{kk}} = \sqrt[4]{a_{kk}^{(k)}}, \quad \hat{l}_{ik} = a_{ik}^{(k)}/\hat{l}_{kk} \quad (i = k + 1, ..., n),$$
(14)

und (5) geht in die symmetrische Formel

$$a_{ij}^{(k+1)} = a_{ji}^{(k+1)} := a_{ij}^{(k)} - \hat{l}_{ik} * \hat{l}_{jk} \qquad (i, j = k+1, ..., n)$$
(15)

über. Mit (14), (15) ergibt sich dann das nachfolgende, nach CHOLESKY benannte Verfahren zur Berechnung von \hat{L} :

6.1.4. LL^{T} -Faktorisierung einer symmetrischen positiv definiten Matrix A ohne Pivotisierung.

Aufgabe: Für positiv definites $A \in \mathfrak{S}^{n,n}$ ist die Cholesky-Faktorisierung $A = \hat{L}\hat{L}^{\mathsf{T}}$ zu berechnen. Die Matrix A ist durch die Elemente des unteren Dreiecks gegeben, und die signifikanten Elemente von $\hat{L} \in \mathbb{R}^{n,n}$ sind auf dem Platz des unteren Dreiecks von A zu speichern. Wenn $a_{kk}^{(k)} \leq 0$ für ein $k \in \{1, ..., n\}$ gilt, soll abgebrochen und ie = -1 gesetzt werden, andernfalls wird ie = 0 gesetzt, und es gilt $\hat{l}_{kk} > 0$ (k = 1, ..., n).

Algorithmus:

$$\begin{split} ie &:= 0 \\ \text{for } k &:= 1(1)n \text{ do} \\ & \left| \begin{array}{l} \text{if } a_{kk} \leq 0 \text{ then } [ie := -1, \text{stop}] \\ a_{kk} &:= \text{sqrt} (a_{kk}) \\ \text{for } i &:= k + 1(1)n \text{ do } a_{ik} &:= a_{ik}/a_{kk} \\ \text{for } j &:= k + 1(1)n \text{ do} \\ & \text{for } i &:= j(1)n \text{ do } a_{ij} &:= a_{ij} - a_{ik} * a_{jk} \\ \end{split} \right. \\ Aujwand &: \sim n^3/6 \text{ opms} + n \text{ opr} \end{split}$$

Bei erfolgreicher Faktorisierung (ie = 0) ergibt sich für n = 4 das folgende Belegungsmuster für A:

$a_{21} a_{22} $	 $\hat{l}_{21} \ \hat{l}_{22}$	
$a_{31} \ a_{32} \ a_{33}$	$\hat{l}_{31} \ \hat{l}_{32} \ \hat{l}_{33}$	
$a_{41} \ a_{42} \ a_{43} \ a_{44}$	$\hat{l}_{41} \ \hat{l}_{42} \ \hat{l}_{43} \ \hat{\bar{l}}_{44}$	

6.1.5. Bemerkung. (i) Die Aussagen von Bemerkung 6.1.2 gelten sinngemäß für die Cholesky-Faktorisierung 6.1.4. Insbesondere kann die Lösung des Gleichungssystems $Ax = \hat{L}\hat{L}^{\mathsf{T}}x = b$ in zwei Schritten durch Lösung von

S2.1: $\hat{\boldsymbol{L}}\boldsymbol{c} = \boldsymbol{b}$, S2.2: $\hat{\boldsymbol{L}}^{\mathsf{T}}\boldsymbol{x} = \boldsymbol{c}$

erfolgen. Die Determinante ergibt sich zu

$$\det (A) = \left(\prod_{k=1}^{n} \hat{l}_{kk}\right)^{2}.$$
(16)

(ii) Die Rundungsfehleranalyse 6.1.3 bleibt für die Cholesky-Faktorisierung 6.1.4 und die Lösung von Ax = b nach (i) gültig, wenn F durch $F := n + 1 + 0.5 \ln n$ ersetzt wird. Die asymptotischen Beziehungen $F \sim n$ und $F_1 \sim 2n^{3/2}$ ändern sich dabei nicht.

Auf die Durchführbarkeitsbedingung (7) kann in beiden symmetrischen Faktorisierungen verzichtet werden, wenn $a_{kk} := v ||A||$ im Fall $a_{kk} \leq 0$ gesetzt wird. In der betrachteten Matrizenklasse wird dadurch die Gutartigkeit der Verfahren nicht beeinträchtigt, allerdings sind die derart berechneten Faktorisierungen i. allg. nicht mehr zur Lösung von Ax = b geeignet, vgl. 5.3.3(i).

(iii) Beide Faktorisierungen 6.1.1 und 6.1.4 können mit Diagonalpivotisierung gemäß Ü 5.2.6 realisiert werden. Zu Beginn des k-ten Schrittes wird dazu ein Index s = s(k) mit $k \leq s \leq n$ derart bestimmt, daß

$$|a_{ss}^{(k)}| = \max\{|a_{ii}^{(k)}|: k \leq i \leq n\}$$
(17)

gilt. Danach werden die Spalten k und s sowie die Zeilen k und s von $A^{(k)}$ vertauscht. Wenn die bereits berechneten l_{ij} bzw. \hat{l}_{ij} (j < k) analog vertauscht werden, ergeben sich die Faktorisierungen

$$\bar{A} = PAP^{\mathsf{T}} = LDL^{\mathsf{T}} = \hat{L}\hat{L}^{\mathsf{T}}$$
(18)

mit der Permutationsmatrix $P := T_{n-1,s(n-1)} \cdots T_{2,s(2)} T_{1,s(1)}$. Dabei gilt

$$d_k \ge l_{jk}^2 d_k + l_{j,k+1}^2 d_{k+1} + \dots + l_{j,j-1}^2 d_{j-1} + d_j$$

bzw.

$$\hat{l}_{kk}^2 \ge \hat{l}_{jk}^2 + \hat{l}_{j,k+1}^2 + \dots + \hat{l}_{jj}^2$$
 (19)

für $1 \leq k < j \leq n$, insbesondere ist

$$d_1 = (\hat{l}_{11})^2 \ge d_2 = (\hat{l}_{22})^2 \ge \dots \ge d_n = (\hat{l}_{nn})^2 > 0,$$
(20)

siehe Ü 6.1.1. Der Quotient

$$\gamma_2 := d_1/d_n = (\hat{l}_{11}/\hat{l}_{nn})^2 \le \text{cond}_2 (A) = [\text{cond}_2 (\hat{L})]^2$$
(21)

ist eine untere Schranke für $\operatorname{cond}_2(A)$, welche die Größenordnung fast immer ausreichend genau wiedergibt und einfacher zu berechnen ist als durch Schätzung von $\|\hat{L}^{-1}\|_2$ nach 5.4.7, 5.4.8, siehe Ü 6.1.2. \Box

Welche der Faktorisierungen $A = LDL^{\intercal}$ oder $A = \hat{L}\hat{L}^{\intercal}$ man verwenden sollte, hängt von den Umständen und auch etwas vom Geschmack ab. Die Cholesky-Faktorisierung läßt sich etwas einfacher programmieren und bietet mehr Freiheit in der Reihenfolge der Berechnung der $a_{ij}^{(k)}$; das Auftreten der *n* Wurzelberechnungen kann kaum als Nachteil angesehen werden. Demgegenüber ist die LDL^{\intercal} -Faktorisierung mit der verwandten Faktorisierung desselben Typs für indefinites symmetrisches A kompatibel. Für manche Anwendungen ist es außerdem zweckmäßig, die Diagonalelemente $d_k = (\hat{l}_{kk})^2$ explizit verfügbar zu haben.

B. Indefinite symmetrische Matrizen

Wir wollen im folgenden untersuchen, ob die symmetrischen Faktorisierungen (2) bzw. (13) auch für reguläres indefinites $A \in S^{n,n}$ existieren und stabil berechnet werden können. Für die Cholesky-Faktorisierung (13) läßt sich diese Frage sofort verneinen, da diese Darstellung für reguläres \hat{L} die positive Definitheit von A impliziert. Dagegen ist die LDL^{T} -Faktorisierung potentiell geeignet, auch indefinite Matrizen darzustellen, was durch das Auftreten positiver und negativer Diagonalelemente in D gekennzeichnet ist, vgl. Ü 6.1.4. Ein Beispiel ist die Matrix

$$\boldsymbol{A} = \begin{pmatrix} -1 & -1 \\ -1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} -1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} = \boldsymbol{L} \boldsymbol{D} \boldsymbol{L}^{\mathsf{T}}$$
(22)

mit den Eigenwerten $\lambda_{1,2} = \pm \sqrt{2}$. Die Zerlegung (22) kann nach 6.1.1 berechnet werden, wenn dort der Test " $a_{kk} \leq 0$ " durch " $a_{kk} = 0$ " ersetzt wird. Andererseits zeigt die Matrix

$$\boldsymbol{A} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \tag{23}$$

mit den Eigenwerten $\lambda_{1,2} = \pm 1$, daß nicht jede reguläre indefinite symmetrische Matrix eine Zerlegung $A = LDL^{\mathsf{T}}$ mit regulärem D besitzt. Algorithmus 6.1.1 bricht hier bereits im ersten Schritt wegen $d_1 = 0$ ab. Dies läßt sich auch nicht durch Diagonalpivotisierung nach 6.1.5, d. h. durch Zulassung symmetrieerhaltender Permutationen gemäß

$$\tilde{\boldsymbol{A}} = \boldsymbol{P}\boldsymbol{A}\boldsymbol{P}^{\mathsf{T}} = \boldsymbol{L}\boldsymbol{D}\boldsymbol{L}^{\mathsf{T}} \tag{24}$$

beheben, denn die Matrix (23) ist unter solchen simultanen Zeilen- und Spaltenvertauschungen invariant.

Das Versagen eines Algorithmus für ein spezielles Problem — hier die Matrix (23) — bedeutet meist, daß auch im Fall der Durchführbarkeit für benachbarte Probleme mit instabilem Verhalten zu rechnen ist. Diese Erfahrung bestätigt sich auch hier:

6.1.6. Beispiel. Es sei $0 < \varepsilon \ll 1$ und

$$oldsymbol{A} := egin{pmatrix} arepsilon & 1\ 1 & arepsilon \end{pmatrix}, \quad oldsymbol{b} := egin{pmatrix} 1\ 0 \end{pmatrix}.$$

Die Algorithmen 6.1.1 und 6.1.2(iii) mögen mit dem oben beschriebenen modifiziertem Test so durchgeführt werden, daß fl $(\alpha + \varepsilon^2 \alpha) = \alpha$ für $\alpha \neq 0$ gesetzt, aber alle übrigen Operationen exakt durchgeführt werden. Dies entspricht qualitativ einer Arithmetik mit $\varepsilon^2 < \nu/2$. Bei dieser Vereinbarung ergibt sich

$$\boldsymbol{D} = \begin{pmatrix} \epsilon & 0 \\ 0 & -1/\epsilon \end{pmatrix}, \quad \boldsymbol{L} = \begin{pmatrix} 1 & 0 \\ 1/\epsilon & 1 \end{pmatrix}, \quad \boldsymbol{x} = \begin{pmatrix} 0 \\ 1 \end{pmatrix};$$

die exakten Größen sind

$$m{D^*} = egin{pmatrix} arepsilon & 0 \ 0 & arepsilon & -1/arepsilon \end{pmatrix}, \ m{L^*} = egin{pmatrix} 1 & 0 \ 1/arepsilon & 1 \end{pmatrix}, \ m{x^*} = rac{1}{1-arepsilon^2} egin{pmatrix} -arepsilon \ 1 \end{pmatrix}.$$

Für den allein durch die Rundung fl $(\varepsilon - 1/\varepsilon) = -1/\varepsilon$ bei der Berechnung von d_2 hervorgerufenen Fehler $\delta x = x - x^*$ gilt dann $\|\delta x\|_{\infty} = \varepsilon + O(\varepsilon^3)$. Der Darstellungsfehler δA_D von A zur Maschinengenauigkeit ν kann dagegen höchstens zu einem Fehler

$$\|\mathbf{d} \boldsymbol{x}\|_{\infty} \leqq \|A^{-1}\|_{\infty} \|\mathbf{d} A_{\mathcal{D}}\|_{\infty} \|\boldsymbol{x}^{*}\|_{\infty} \gneqq \boldsymbol{v} \eqqcolon \Delta \boldsymbol{x}_{\mathsf{opt}}(A, \boldsymbol{v})$$

führen. Im Fall $\varepsilon^2 = \nu/3$ überschreitet daher der erzeugte Rundungsfehler das optimale Fehlerniveau um den Faktor $\varepsilon/\nu = 1/\sqrt{3\nu}$. Da der Quotient nicht unabhängig von ν beschränkt werden kann, liegt Instabilität vor. Die Lösung von Ax = bmittels spaltenpivotisierter **LR**-Faktorisierung würde auf die berechneten Faktoren

$$oldsymbol{L} = egin{pmatrix} 1 & 0 \ arepsilon & 1 \end{pmatrix}, oldsymbol{R} = egin{pmatrix} 1 & arepsilon \ arepsilon & 1 \end{pmatrix} ext{von} \ oldsymbol{ar{A}} = oldsymbol{P}oldsymbol{A} = egin{pmatrix} 1 & arepsilon \ arepsilon & 1 \end{pmatrix}, \ ext{also auf} \ oldsymbol{x} = egin{pmatrix} -arepsilon \ arepsilon & 1 \end{pmatrix}$$

führen. Der erzeugte Rundungsfehler ist hier $\|\delta x\|_{\infty} = \varepsilon^2 + O(\varepsilon^3) = \nu/3$ und liegt wie erwartet in der Größenordnung von $\Delta x_{opt}(A, \nu)$.

Das beschriebene instabile Verhalten ist darauf zurückzuführen, daß die Diagonalelemente der Matrizen $M^{(k)}$ im Unterschied zum positiv definiten Fall im Laufe der Faktorisierung beliebig klein im Vergleich zu den Nichtdiagonalelementen werden können. Dies führt zu einem unbeschränkten Wachstum der $M^{(k)}$, im Beispiel ist $M^{(2)} = (1/\varepsilon)$. Die zur Sicherung der Stabilität daher notwendige Zulassung von Nichtdiagonalpivots zerstört dagegen i. allg. die Symmetrie und verhindert die gewünschte Aufwandshalbierung.

Als Ausweg aus dieser Situation bietet sich an, einem Eliminationsschritt mit dem Nichtdiagonalpivot $a_{ss}^{(k)}$ $(s \neq \hat{s})$ sofort einen weiteren mit dem spiegelbildlich gelegenen Pivot $a_{ss}^{(k-1)}$ folgen zu lassen. Nach einem solchen Doppelschritt ist die Restmatrix $M^{(k+2)}$ wieder symmetrisch. Dieses Vorgehen ist im wesentlichen identisch mit der simultanen Elimination des Variablenpaares $\{x_s, x_s\}$, siehe Ü 6.1.6.

Zur Beschreibung einer solchen Blockelimination betrachten wir einen ersten Schritt, bei dem der Block (x_1, x_p) der ersten p Variablen eliminiert werden soll, und setzen dazu

$$\boldsymbol{A} = \begin{pmatrix} \boldsymbol{D} & | & \boldsymbol{C}^{\mathsf{T}} \\ \hline \boldsymbol{C} & | & \boldsymbol{B} \end{pmatrix}, \quad \boldsymbol{L}_{1} = \begin{pmatrix} \boldsymbol{I}_{p} & | & \boldsymbol{O} \\ \hline -\boldsymbol{G} & | & \boldsymbol{I}_{n-p} \end{pmatrix};$$
(25)

dabei ist $D \in S^{p,p}$, $C, G \in \mathbb{R}^{n-p,p}$, $B \in S^{n-p,n-p}$ mit $p \in \{1, 2\}$. Die Matrix G der Eliminationskoeffizienten soll so bestimmt werden, daß

$$L_1 A = \left(\frac{D}{C - GD} \begin{vmatrix} C^{\mathsf{T}} \\ B - GC^{\mathsf{T}} \end{vmatrix} \right) = \left(\frac{D}{O} \begin{vmatrix} C^{\mathsf{T}} \\ M \end{vmatrix} \right)$$
(26)

gilt. Dies ist für beliebiges C genau dann möglich, wenn der Pivotblock D regulär ist. Falls letzteres zutrifft, folgt

$$G = CD^{-1}, \quad M = B - GC^{\mathsf{T}} = B - CD^{-1}C^{\mathsf{T}},$$

insbesondere ist also M wie A symmetrisch. Um die volle Symmetrie bereits im ersten Schritt sichtbar zu machen, führen wir noch eine Rechtsmultiplikation mit L_1^{T} aus und erhalten

$$\bar{A} := L_1 A L_1^{\mathsf{T}} = \left(\frac{D \mid C^{\mathsf{T}}}{O \mid M} \right) \left(\frac{I_p \mid -G^{\mathsf{T}}}{O \mid I_{n-p}} \right) = \left(\frac{D \mid C^{\mathsf{T}} - DG^{\mathsf{T}}}{O \mid M} \right) = \left(\frac{D \mid O}{O \mid M} \right),$$
(27)

denn wegen der Symmetrie von D ist $C^{\intercal} - DG^{\intercal} = O$. Selbstverständlich braucht die Rechtsmultiplikation mit L_1^{\intercal} nicht explizit durchgeführt zu werden, da sie lediglich die Symmetrie in \overline{A} erzeugt und die signifikanten, bereits in L_1A vorkommenden Blöcke D und M nicht mehr verändert.

Die Regularität von D muß bei diesem Vorgehen gegebenenfalls durch geeignete symmetrieerhaltende Zeilen- und Spaltenvertauschungen, d. h. durch Übergang

 $A \to TAT^{\mathsf{T}}$

mit einer Permutationsmatrix T erzwungen werden: Falls ein $a_{ss} \neq 0$ existiert, kann dieses mit der Festlegung $T := T_{1s}$ in die Position (1, 1) gebracht und ein einfacher Eliminationsschritt mit p = 1 ausgeführt werden. Existiert kein nichtverschwindendes Diagonalelement, so wird a_{s1} ($2 \leq s \leq n$) in der ersten Spalte von A mit $\delta := a_{11}a_{ss} - a_{s1}^2 \neq 0$ gesucht und mittels $T := T_{2s}$ in die Position (2, 1) gebracht. Ein solches Element existiert für reguläres A stets, denn andernfalls wäre $a_{s1} = 0$ für alle s = 1, ..., n. Dann ist det (D) = $\delta \neq 0$, und es wird ein Blockeliminationsschritt mit p = 2 ausgeführt.

In beiden Fällen p = 1 bzw. p = 2 ist die Restmatrix M wie A selbst regulär und symmetrisch, so daß derselbe Eliminationsprozeß auf M statt A angewendet und A folglich sukzessive auf obere Blockdiagonalform mit Diagonalblöcken der Dimension 1 oder 2 transformiert werden kann.

Beginnend mit $A^{(1)} := A$, sei die Blockdiagonalform in den Zeilen und Spalten 1 bis k - 1 bereits erzeugt worden, d. h., $A^{(k)}$ sei von der Gestalt

$$A^{(k)} = \begin{pmatrix} \underline{D}^{(1)} & O \\ & & \\$$

Im Schritt k werden dann $p = p(k) \in \{1, 2\}$ und eine Vertauschungsmatrix T_k := $T_{t(k),s(k)}$ mit $k \leq t(k) \leq s(k) \leq n$ so festgelegt, daß der (p, p)-Pivotblock $D^{(k)}$ von

$$A^{(k)} := T_k A^{(k)} T_k \tag{29}$$

regulär ist; man beachte $T_k = T_k^{\mathsf{T}}$. Zur Vereinfachung der Schreibweise bezeichnen wir die durch die Vertauschungen (29) entstandene Matrix wieder mit $A^{(k)}$ und nicht mit $\hat{A}^{(k)}$ wie in 5.1, 5.2. Der nachfolgende, wie oben symmetrisierte Eliminationsschritt

$$A^{(k+p)} := L_k A^{(k)} L_k^{\mathsf{T}} \tag{30}$$

mit der Blockeliminationsmatrix

$$\boldsymbol{L}_{k} := \begin{pmatrix} \boldsymbol{I}_{k-1} & & \\ \hline & \boldsymbol{I}_{p} & \\ \hline & -\boldsymbol{G}_{k} & \boldsymbol{I}_{n-p-k+1} \end{pmatrix}, \quad \boldsymbol{G}_{k} = \begin{pmatrix} \boldsymbol{l}_{k+1,k} \\ \vdots \\ \vdots \\ \boldsymbol{l}_{nk} \end{pmatrix} \quad \text{für} \quad p = 1, \\
\boldsymbol{G}_{k} = \begin{pmatrix} \boldsymbol{l}_{k+2,k} & \boldsymbol{l}_{k+2,k+1} \\ \vdots & \vdots \\ \vdots & \vdots \\ \boldsymbol{l}_{nk} & \boldsymbol{l}_{n,k+1} \end{pmatrix} \quad \text{für} \quad p = 2, \quad (31)$$

erzeugt bei der Festlegung

$$G_k := C^{(k)} [D^{(k)}]^{-1}$$
(32)

die gewünschten Nullen auf dem Platz von $C^{(k)}$ und transformiert $B^{(k)}$ in die neue Restmatrix

$$M^{(k+p)} := B^{(k)} - G_k C^{(k)\top}.$$
(33)

Unter Beachtung von

$$[\mathbf{D}^{(k)}]^{-1} = 1/\delta_k, \qquad \delta_k := a_{kk}^{(k)}$$
(34)

im Fall p = 1 bzw.

$$[\mathbf{D}^{(k)}]^{-1} = \frac{1}{\delta_k} \begin{pmatrix} a_{k+1,k+1}^{(k)} & -a_{k+1,k}^{(k)} \\ -a_{k+1,k}^{(k)} & a_{kk}^{(k)} \end{pmatrix}, \quad \delta_k := a_{kk}^{(k)} a_{k+1,k+1}^{(k)} - (a_{k+1,k}^{(k)})^2$$
(35)

im Fall p = 2 lassen sich die Eliminationskoeffizienten l_{ij} (j = k bzw. j = k, k + 1)mittels (32) explizit durch die Elemente von $A^{(k)}$ darstellen.

Zusammenfassend erhalten wir den folgenden Basisalgorithmus:

6.1.7. Symmetrische Blockfaktorisierung der symmetrischen Matrix A. Initialisierung: $A^{(1)} := A, k := 1$

- k-ter Schritt: S1: Pivotwahl: Bestimme $p = p(k) \in \{1, 2\}$ und Vertauschungsmatrix T_k $:= T_{t(k), s(k)}, k \leq t(k) \leq s(k) \leq n$ derart, daß $A^{(k)} := T_k A^{(k)} T_k$ der Bedin-gung $\delta_k \neq 0$ mit δ_k gemäß (34), (35) genügt. S2: Bestimmung der Eliminationskoeffizienten und Transformation der Matrix
- $A^{(k)}:$

$$\begin{split} \text{S2.1: Falls } p &= 1, \text{ setze} \\ a_{ij}^{(k-1)} &:= a_{ji}^{(k+1)} := \begin{cases} 0 & \text{für } i = k+1, \dots, n, j = k, \\ a_{ij}^{(k)} - l_{ik} * a_{kj}^{(k)} & \text{für } i, j = k+1, \dots, n, \\ a_{ij}^{(k)} & \text{sonst} \end{cases} \\ \text{mit} \\ l_{ik} &:= a_{ik}^{(k)} / \delta_k, \quad i = k+1, \dots, n \\ \text{S2.2: Falls } p &= 2 \text{ und } k < n-1, \text{ setze} \\ a_{ij}^{(k-2)} &:= a_{ji}^{(k+2)} \\ &:= \begin{cases} 0 & \text{für } i = k+2, \dots, n; j = k, k+1, \\ a_{ij}^{(k)} - l_{ik} * a_{kj}^{(k)} - l_{i,k+1} * a_{k+1,j}^{(k)} & \text{für } i, j = k+2, \dots, n, \\ \text{sonst} \end{cases} \\ \text{mit} \\ l_{ik} &:= (a_{ik}^{(k)} * a_{k+1,k+1}^{(k)} - a_{i,k+1}^{(k)} * a_{k+1,k}^{(k)}) / \delta_k, \\ l_{i,k+1} &:= (-a_{ik}^{(k)} * a_{k+1,k}^{(k)} + a_{i,k+1}^{(k)} * a_{kk}^{(k)}) / \delta_k & \text{für } i = k+2, \dots, n. \\ \text{S3: Indexerhöhung: Setze } k &:= k+p. \text{ Falls } k < n, \text{ gehe nach S1.} \\ Aufwand: \sim n^3/6 \text{ opms} (\text{bei Ausnutzung der Symmetrie von } \mathbf{M}^{(k+p)}) \end{cases} \end{split}$$

Bei Auswahl eines (2, 2)-Pivotblockes $D^{(k)}$ wird die Stufe k + 1 übersprungen und sofort zu k + 2 übergegangen. Um zu einer einheitlichen und einfachen Bezeichnung zu gelangen, setzen wir in diesem Fall $T_{k+1} := L_{k+1} := I$, $A^{(k+1)} := A^{(k)}$. Mit dieser Vereinbarung läßt sich der Gesamtprozeß in der Form

$$\boldsymbol{D} := \boldsymbol{A}^{(n)} = \boldsymbol{L}_{n-1} \boldsymbol{T}_{n-1} \boldsymbol{L}_{n-2} \boldsymbol{T}_{n-2} \cdots \boldsymbol{L}_1 \boldsymbol{T}_1 \boldsymbol{A} \boldsymbol{T}_1 \boldsymbol{L}_1^{\mathsf{T}} \cdots \boldsymbol{T}_{n-2} \boldsymbol{L}_{n-2}^{\mathsf{T}} \boldsymbol{T}_{n-1} \boldsymbol{L}_{n-1}^{\mathsf{T}}$$
(36)

darstellen, und D ist nach Konstruktion von Blockdiagonalform

$$\boldsymbol{D} = \begin{pmatrix} \boxed{\boldsymbol{D}^{(1)}} \\ \boxed{\boldsymbol{D}^{(1+p_1)}} \\ \boxed{\boldsymbol{D}^{(l)}} \end{pmatrix}.$$
 (37)

Dabei ist l = n - 1 oder l = n in Abhängigkeit davon, ob der letzte Pivotblock die Dimension 2 oder 1 hat.

Wenn die nach dem k-ten Schritt vorkommenden Vertauschungen T_{k+1}, \ldots, T_{n-1} gemäß

$$\left(\frac{O}{G_k}\right) := T_{n-1} \cdots T_{k+1} \left(\frac{O}{G_k}\right) \quad (k = 1, \dots, n-2)$$
(38)

auf die Eliminationskoeffizienten G_k angewendet werden, geht (36) mit den so ver-

tauschten Größen in

$$\boldsymbol{D} = \boldsymbol{A}^{(\boldsymbol{n})} = \boldsymbol{L}_{\boldsymbol{n}-1} \boldsymbol{L}_{\boldsymbol{n}-2} \cdots \boldsymbol{L}_1 \boldsymbol{T}_{\boldsymbol{n}-1} \boldsymbol{T}_{\boldsymbol{n}-2} \cdots \boldsymbol{T}_1 \boldsymbol{A} \boldsymbol{T}_1 \cdots \boldsymbol{T}_{\boldsymbol{n}-2} \boldsymbol{T}_{\boldsymbol{n}-1} \boldsymbol{L}_1^{\mathsf{T}} \cdots \boldsymbol{L}_{\boldsymbol{n}-2}^{\mathsf{T}} \boldsymbol{L}_{\boldsymbol{n}-1}^{\mathsf{T}},$$

folglich in

$$LDL^{\mathsf{T}} := L_1^{-1} \cdots L_{n-1}^{-1} DL_{n-1}^{-\mathsf{T}} \cdots L_1^{-\mathsf{T}} = T_{n-1} \cdots T_1 A T_1 \cdots T_{n-1} =: PAP^{\mathsf{T}}$$
(39)

über, vgl. 5.1.C. Wie dort entsteht dabei L_k^{-1} aus L_k , indem $-G_k$ durch G_k ersetzt wird, vgl. (31).

Als Ergebnis erhalten wir das folgende Resultat.

6.1.8. Satz. Algorithmus 6.1.7 ist in exakter Arithmetik für jede reguläre Matrix $A \in S^{n,n}$ durchführbar und liefert eine Faktorisierung

$$PAP^{\mathsf{T}} = LDL^{\mathsf{T}} \tag{40}$$

 mit

$$L = \begin{pmatrix} I_{p(1)} & 0 \\ G_1 & I_{p(2)} \\ G_2 & I_{p(\ell)} \end{pmatrix} = \begin{pmatrix} I_{p(1)} & 0 \\ I_{p(2)} & 0 \\ I_{ij} & I_{p(\ell)} \end{pmatrix},$$

$$D = \begin{pmatrix} D^{(1)} & 0 \\ D^{(1+p(1))} & I_{p(\ell)} \end{pmatrix}.$$
(41)

Dabei ist L die Matrix der gemäß (38) vertauschten Eliminationskoeffizienten, D die Blockdiagonalmatrix der Pivotblöcke und $P := T_{n-1} \cdots T_1$ die durch die Vertauschungen definierte Permutationsmatrix.

6.1.9. Bemerkung. (i) Algorithmus 6.1.7 kann unter Ausnutzung der Symmetrie der $M^{(k)}$ zeilenweise auf dem Platz des unteren Dreiecks von A ausgeführt werden, wobei die l_{ij} auf dem Platz der zu 0 gemachten $a_{ij}^{(k)}$ gespeichert werden. Die Vertauschungen (38) bedeuten dann, daß die bereits berechneten l_{ij} wie die Zeilen von $A^{(k)}$ mit vertauscht werden. Für n = 5 ergibt sich etwa folgendes Belegungsschema $(p_1 = 2, p_3 = 1, p_4 = 2)$:

$$\begin{bmatrix} a_{11} \\ a_{21} & a_{22} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} & a_{44} \\ a_{51} & a_{52} & a_{53} & a_{54} & a_{55} \end{bmatrix} \rightarrow \begin{bmatrix} d_{11} \\ d_{21} \\ d_{22} \\ l_{31} \\ l_{32} \\ l_{33} \\ l_{41} \\ l_{52} \\ l_{53} \\ l_{52} \\ l_{53} \\ l_{54} \\ d_{54} \\ d_{55} \end{bmatrix} .$$

13*

Man beachte, daß die (2,2)-Pivots $D^{(1)}$, $D^{(4)}$ symmetrisch sind und nur das untere Dreieck gespeichert zu werden braucht. Da die zugehörigen Diagonalblöcke von Ldie Gestalt I_2 haben, treten keine Probleme bei der Speicherung der Außerdiagonalelemente d_{21} und d_{54} auf. Die zur Rekonstruktion von P erforderliche Information wird in einem Feld der Länge n gespeichert; für Details sei auf die Literatur verwiesen, siehe B 6.1.

(ii) Bei erfolgreicher Faktorisierung mit regulärem D kann das Gleichungssystem $Ax = P^{\mathsf{T}}LDL^{\mathsf{T}}Px = b$ in Analogie zu 4.2.1 für jedes $b \in \mathbb{R}^n$ unter Verwendung der Faktorisierung (40) wie folgt gelöst werden:

S2.1: Berechne c aus $Lc = \overline{b} := Pb$

S2.2: Berechne d aus Dd = c

S2.3: Berechne $\bar{\boldsymbol{x}}$ aus $\boldsymbol{L}^{\mathsf{T}} \bar{\boldsymbol{x}} = \boldsymbol{d}$, setze $\boldsymbol{x} := \boldsymbol{P}^{\mathsf{T}} \bar{\boldsymbol{x}}$.

Das System Dd = c zerfällt dabei in Teilsysteme der Dimension 1 oder 2 entsprechend der Blockstruktur von D. Der Aufwand ist $\sim n^2$ opms.

(iii) Die zur Darstellung (40) benötigten Vertauschungen der l_{ij} gemäß (38) können auch weggelassen werden, d. h., 6.1.7 kann in der dort angegebenen Version auf dem Platz des unteren Dreiecks von A ausgeführt werden. Dann muß direkt mit der Darstellung (36) gearbeitet werden, die sich äquivalent in der Form

$$A = \tilde{L}D\tilde{L}^{\mathsf{T}} \quad \text{mit} \quad \tilde{L} := (L_{n-1}T_{n-1}\cdots L_1T_1)^{-1} = T_1L_1^{-1}T_2L_2^{-1}\cdots T_{n-1}L_{n-1}^{-1}$$
(42)

schreiben läßt. Die Matrix \tilde{L} ist selbst i. allg. keine untere Dreiecksmatrix, allerdings lassen sich die signifikanten Elemente l_{ij} der einzelnen Faktoren L_k auf dem Platz des unteren Blockdreiecks von A speichern. Die Lösung von $Ax = \tilde{L}D\tilde{L}^{\mathsf{T}}x = b$ erfolgt dann gemäß

S2.1: Berechne \tilde{c} aus $\tilde{L}\tilde{c} = b$, d. h., berechne $\tilde{c} = \tilde{L}^{-1}b = L_{n-1}T_{n-1}\cdots L_1T_1b$ S2.2: Berechne \tilde{d} aus $D\tilde{d} = \tilde{c}$

S2.3: Berechne x aus $\tilde{L}^{\mathsf{T}}x = \tilde{d}$, d. h., berechne $x = \tilde{L}^{-\mathsf{T}}\tilde{d} = T_1L_1^{\mathsf{T}}\cdots T_{n-1}L_{n-1}^{\mathsf{T}}\tilde{d}$ vgl. 5.2.6(i). Die Anzahl der Rechenoperationen ist identisch mit der bei der Realisierung mit L gemäß (ii).

(iv) Die Determinante von A ergibt sich zu

$$\det (A) = \det (D^{(1)}) \det (D^{(1+p_1)}) \cdots \det (D^{(l)}) = \prod_{k=1}^n \delta_k;$$
(43)

formal muß dabei $\delta_{k+1} := 1$ gesetzt werden, sofern $D^{(k)}$ ein (2,2)-Pivot ist.

(v) Für singuläres A versagt 6.1.7, weil im Schritt S1 kein regulärer Pivotblock gefunden werden kann, oder der letzte Block $D^{(l)}$ ist singulär. Dies kann wie in 5.2.3 abgefragt und durch einen Fehlerindikator angezeigt werden. Bei der Determinantenberechnung ist dann wie in 5.2.5 mit det (A) := 0 abzubrechen.

Selbstverständlich ist das Verfahren 6.1.7 ohne Zusatzforderungen an die Pivotwahl für Computerrechnung ungeeignet, vgl. 6.1.6. Wie beim Gaußschen Algorithmus muß auch hier durch geeignete Wahl von p(k) und T_k ein zu starkes Anwachsen der

 $M^{(k)}$ verhindert werden. In Analogie zur vollständigen Pivotisierung beim Gaußschen Algorithmus könnte versucht werden, den Betrag von δ_k über alle zugelassenen Möglichkeiten zu maximieren. Dazu muß

$$v_k := \max \{ |a_{ii}^{(k)}|, |a_{ii}^{(k)}a_{jj}^{(k)} - (a_{ij}^{(k)})^2| : i, j = k, ..., n; i > j \}$$

bestimmt werden. Falls $v_k = |a_{ss}^{(k)}|$ ist, wird p = 1 und $T^{(k)} = T_{ks}$ gesetzt. Im Fall $v_k = |a_{ss}^{(k)} a_{ss}^{(k)} - (a_{ss}^{(k)})^2|$ wird p = 2 gesetzt und $a_{ss}^{(k)}$ in die Position (k+1, k) gebracht. Dazu sind zwei simultane Zeilen- und Spaltenvertauschungen erforderlich. Zur Berechnung aller v_k werden jedoch mindestens $\sim n^3/6$ opm benötigt. Das ist etwa derselbe Aufwand wie für die Berechnung von L und D selbst. Dieses Vorgehen scheidet daher aus Aufwandsgründen aus. Wie beim Gaußschen Algorithmus lassen sich jedoch auch hier aufwandsgünstigere teilweise Pivotisierungsstrategien angeben, die zu akzeptablem Stabilitätsverhalten führen. Bei der nachfolgend beschriebenen Strategie wird pro Schritt mindestens eine, höchstens jedoch zwei Spalten von $M^{(k)}$ abgesucht, so daß sie etwa der Spaltenpivotisierung beim Gaußschen Algorithmus entspricht.

6.1.10. Spaltenpivotisierung nach Bunch-Kautman-Parlett: Es sei α mit $0 < \alpha < 1$ fest vorgegeben, und $A^{(k)}$ bezeichne die zu Beginn des k-ten Schrittes von 6.1.7 vorliegende Matrix der Gestalt (28).

S1.1: Bestimme Index s = s(k), $k \leq s \leq n$, mit

 $\lambda_{k} := |a_{sk}^{(k)}| = \max \{|a_{ik}^{(k)}|: i = k + 1, ..., n\}$ S1.2: if $|a_{kk}^{(k)}| \ge \alpha \lambda_{k}$ S1.3: then $[p(k) := 1, T_{k} := I]$ else

 ense

 S1.4:
 Bestimme $\sigma_k := \max \{ |a_{is}^{(k)}| : i = k, ..., n : i \neq s \}$

 S1.5:
 if $|a_{kk}^{(k)}| \sigma_k \ge \alpha (\lambda_k)^2$

 S1.6:
 then $[p(k) := 1, T_k := I]$

 else
 else

 S1.7:
 if $|a_{ss}^{(k)}| \ge \alpha \cdot \sigma_k$

 S1.8:
 then $[p(k) := 1, T_k := T_{ks}]$

 else
 else $[p(k) := 2, T_k := T_{k+1,s}]$

Der Faktorisierungsalgorithmus 6.1.7 mit der Pivotstrategie 6.1.10 wird Bunch-Kaufman-Parlett-Faktorisierung - kurz: BKP-Faktorisierung - der symmetrischen Matrix A genannt, wobei für α der unten bestimmte optimale Wert $\alpha = \alpha_0$ gewählt wird. Das Wachstumsverhalten der $M^{(k)}$ für diesen Algorithmus wird durch die folgende Aussage charakterisiert.

6.1.11. Aussage. Algorithmus 6.1.7 ist für jede reguläre Matrix $A \in S^{n,n}$ mit der Pivotstrategie 6.1.10 für jedes $\alpha \in (0, 1)$ in exakter Arithmetik durchführbar. Für jeden Index $k \in \{1, ..., n-1\}$, für den der Schritt $k \rightarrow k + p_k$ mit $p_k \in \{1, 2\}$

durchgeführt wird, gilt dabei

(44)

$$\mu_{k+1} \leq \mu_k (1 + 1/\alpha) \quad \text{im Fall} \quad p_k = 1$$

$$\mu_{k+2} \leq \mu_k (1 + 2/(1 - \alpha)) \quad \text{im Fall} \quad p_k = 2,$$
(45)

wobei

$$u_{k} := \max \{ |a_{ij}^{(k)}| : i, j = k, ..., n \} \qquad (k = 1, ..., n)$$

$$(46)$$

gesetzt wurde. Die für den Übergang $k \to k+2$ gültige Schranke

$$\mu_{k+2} \leq \max \{ (1+1/\alpha)^2, \quad 1+2/(1-\alpha) \} \ \mu_k =: [\varrho(\alpha)]^2 \ \mu_k \tag{47}$$

 $\mu_{k+2} \leq \max \left\{ (1 + 1/\alpha), 1 + 1/\alpha \right\},$ wird für $\alpha = \alpha_0 := (1 + \sqrt{17})/9 = 0.6404$ minimal und hat den Wert $\varrho_0 := \varrho(\alpha_0)$ $= (1 + \sqrt{17})/2 \leq 2.57$. Insbesondere gilt bei Verwendung von $\alpha = \alpha_0$ (48)

$$\mu_k \le (2.57)^{k-1} \,\mu_1 \qquad (k = 1, \dots, n). \tag{48}$$

Beweis. Wir bemerken zunächst, daß

$$\lambda_{k} = |a_{sk}^{(k)}| = |a_{ks}^{(k)}| \le \max\{|a_{is}^{(k)}| : k \le i \le n; i \neq s\} = \sigma_{k} \le \mu_{k}$$
(49)

gilt. Der Einfachheit halber betrachten wir im folgenden nur den ersten Schritt und lassen den oberen Index k = 1 weg.

Fall 1: p = 1. Hier gilt nach S2.1 aus 6.1.7

$$a_{ij}^{(2)} = a_{ij} - a_{i1}a_{j1}/a_{11}$$
 $(i, j = 2, ..., n).$

Wenn S1.3 ausgeführt wird, folgt mit S1.2 und (49)

$$|a_{ij}^{(2)}| \leq \mu_1 + |a_{i1}| \; |a_{j1}|/|a_{11}| \leq \mu_1 + \lambda_1^2/(lpha\lambda_1) \leq \mu_1(1+1/lpha)$$

Wenn S1.6 ausgeführt wird, ergibt sich mit S1.5 und (49)

$$|a_{ij}^{(2)}| \leq \mu_1 + \lambda_1^2/(lpha\lambda_1^2/\sigma_1) = \mu_1 + \sigma_1/lpha \leq \mu_1(1+1/lpha).$$

Wird schließlich S1.8 ausgeführt, so folgt aus S1.7 - man beachte die Vertauschungen -

$$|a_{ij}^{(2)}| \leq \mu_1 + \sigma_1^2/(lpha\sigma_1) = \mu_1 + \sigma_1/lpha \leq \mu_1(1+1/lpha),$$

womit (44) bewiesen ist.

Fall 2: p = 2. Hier ist nach S2.2 aus 6.1.7

$$a_{ij}^{(3)} = a_{ij} - \left[(a_{i1}a_{22} - a_{i2}a_{21}) a_{j1} + (-a_{i1}a_{21} + a_{i2}a_{11}) a_{j2} \right] / \delta_1 \quad (i, j = 3, ..., n)$$

Wegen S1.9, S1.5 und S1.7 gilt

$$|a_{21}| = \lambda_1, \qquad |a_{11}| < \alpha \lambda_1^2 / \sigma_1, \qquad |a_{22}| < \alpha \sigma_1,$$
 (50)

folglich

$$|\delta_1| = -\delta_1 = a_{21}^2 - a_{11}a_{22} \ge \lambda_1^2 - (\alpha\lambda_1^2/\sigma_1) \ (\alpha\sigma_1) = \lambda_1^2(1 - \alpha^2) > 0, \tag{51}$$

insbesondere ist also $\delta_1 = \det(D^{(1)})$ negativ. Unter Beachtung von (50), (51) und $|a_{i1}| \leq \lambda_1$, $|a_{i2}| \leq \sigma_1 \ (i=3,...,n)$ erhalten wir

$$egin{aligned} |a_{ij}^{(3)}| &\leq \mu_1 + ig[(\lambda_1lpha\sigma_1+\sigma_1\lambda_1)\,\lambda_1 + ig(\lambda_1^2+\sigma_1(lpha\lambda_1^2/\sigma_1)ig)\,\sigma_1ig]/[\lambda_1^2(1-lpha^2)ig] \ &= \mu_1 + 2\sigma_1/(1-lpha) \leq \mu_1ig(1+2/(1-lpha)ig), \end{aligned}$$

also (45). Die Abschätzung (47) ergibt sich aus (44), (45) unter Beachtung der Tatsache, daß zwei aufeinanderfolgende Schritte mit einem (1, 1)-Pivot ein durch

$$\mu_{k+2} \leq (1+1/\alpha) \ \mu_{k+1} \leq (1+1/\alpha)^2 \ \mu_k$$

beschränktes Wachstum hervorrufen. Das Minimum von $\rho(\alpha)$ wird im Fall

$$(1 + 1/\alpha)^2 = 1 + 2/(1 - \alpha)$$

erreicht, was auf die quadratische Gleichung $4\alpha^2-\alpha-1=0$ mit der Wurzel $\alpha_0\in(0,1)$ führt. \Box

6.1.12. Rundungsfehleranalyse. Die reguläre Matrix $A \in \mathfrak{S}^{n,n}$ genüge der Bedingung

$$\kappa = \kappa F \operatorname{cond}_{\infty}(A) < 1 \quad \operatorname{mit} \quad F := 5.5n^2 \mu_{\max} / \|A\|_{\infty} \le 2.2n^2 (2.57)^n,$$
 (52)

wobei $\mu_{\max} := \max \{\mu_k : 1 \leq k \leq n\}$ ist mit μ_k gemäß (46). Dann ist die BKP-Faktorisierung durchführbar. Zu den berechneten Faktoren $L, D \in \mathbb{R}^{n,n}$ existiert eine symmetrische Störung dA mit

$$P(A + \delta A) P^{\mathsf{T}} = LDL^{\mathsf{T}} \quad \text{und} \quad \|\delta A\|_{\infty} \leq \nu F \|A\|_{\infty}.$$
(53)

Die zu $b \in \mathbb{R}^n$ gehörende und gemäß 6.1.9(ii) bzw. (iii) berechnete Lösung x von Ax = b ist die exakte Lösung des durch eine Matrix $\delta_1 A$ gestörten Systems

$$(\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \boldsymbol{x} = \boldsymbol{b}$$
, wobei $\|\boldsymbol{\delta}_1 \boldsymbol{A}\|_{\infty} \leq \boldsymbol{v} F_1 \|\boldsymbol{A}\|_{\infty}$ mit $F_1 \leq (n+1) F$ ist.
(54)

6.1.13. Bemerkung. (i) Das nach 6.1.11 eventuell mögliche exponentielle Wachstum der μ_k tritt praktisch fast nie auf; es gilt fast immer

$$\mu_k \leq \gamma \mu_1, \quad \text{also} \quad F \leq 5.5 \gamma n^2$$
(55)

mit γ in der Größenordnung von 1. Die BKP-Faktorisierung und die Lösung von $A\mathbf{x} = \mathbf{b}$ mittels der Faktoren \mathbf{L} und \mathbf{D} sind also numerisch gutartige Prozesse mit den Kumulationskonstanten F nach (52), (55) und F_1 nach (54).

(ii) Bei der Pivotstrategie 6.1.10 kann ein (2,2)-Pivotblock nur für indefinites A auftreten, denn jedem solchen Block entspricht ein Eigenwertpaar aus einem positiven und einem negativen Eigenwert, siehe Ü 6.1.5. Für positiv definites A ist daher D eine positive Diagonalmatrix, d. h., die BKP-Faktorisierung ist dann identisch mit der LDL^{T} -Faktorisierung 6.1.1, wobei allerdings eine spezielle Form der Diagonalpivotisierung verwendet wird. Man beachte, daß in 6.1.10 auch im positiv definiten Fall $T_k = T_{ks}$ mit $k < s \leq n$ gewählt werden kann und folglich $a_{kk}^{(k)}$ mit $a_{ss}^{(k)}$ vertauscht wird. \Box

6.1.14. Bemerkung. (i) Sämtliche in diesem Abschnitt beschriebenen symmetrischen Faktorisierungen können bzw. werden auf dem Platz des unteren Dreiecks von A ausgeführt. Selbstverständlich kann auch mit $\mathbf{R} := \mathbf{L}^{\mathsf{T}}$ auf dem Platz des oberen Dreiecks gearbeitet werden. Dies ist bei Implementierung in FORTRAN i. allg. günstiger, weil die Algorithmen auf eine zeilenweise Berechnung von \mathbf{L} , also eine spaltenweise Berechnung von \mathbf{R} orientiert sind und Felder in FORTRAN spaltenweise gespeichert werden.

(ii) Das jeweils nicht benötigte Dreieck von A kann zum Speichern der Originalmatrix A benutzt werden, wenn diese später etwa für die iterative Verbesserung von x analog zu 5.4 benötigt wird. Die Diagonalelemente müssen dann gesondert aufgehoben werden. Falls A nicht benötigt wird, kann das obere bzw. untere Dreieck zeilen- oder spaltenweise in einem eindimensionalen Feld der Länge n(n + 1)/2kompakt gespeichert werden. Wegen der komplizierteren Indexrechnung werden die Programme dabei unübersichtlicher, die Ausführungszeiten unterscheiden sich bei geeigneter Implementierung kaum von den mit zweidimensionalen Feldern arbeitenden Implementierungen. \Box

Übungsaufgaben

Ü 6.1.1. Man zeige, daß für das Cholesky-Verfahren 6.1.4

 $a_{jj}^{(j)} = a_{jj}^{(k)} - \hat{l}_{jk}^2 - \hat{l}_{j,k+1}^2 - \dots - \hat{l}_{j,j-1}^2 \qquad (j > k)$

gilt, und folgere hieraus (19) und (20).

Ü 6.1.2. Es sei $\hat{L} = (\hat{l}_{ij})$ eine untere Dreiecksmatrix. Man zeige

$$|\hat{l}_{ii}| \leq ||\hat{L}||_p$$
 $(i = 1, ..., n)$

und für reguläres \hat{L}

 $1/|\hat{l}_{nn}| \leq ||\hat{L}^{-1}||_n$

für jede der Normen $p \in \{1, 2, \infty\}$ und leite hieraus (21) ab. Gilt (21) auch für $p \in \{1, \infty\}$?

Ü 6.1.3. Man gebe zu 6.1.1 bzw. 6.1.4 analoge Algorithmen zur LDL^{T} - bzw. $\hat{L}\hat{L}^{T}$ -Faktorisierung mit Diagonalpivotisierung entsprechend 6.1.5 (iii) an.

Ü 6.1.4. Die symmetrische Matrix A besitze die Faktorisierung $A = LDL^{\mathsf{T}}$ mit einer unteren Einsdreiecksmatrix L und einer Diagonalmatrix $D = \text{diag}(d_i)$. Man zeige, daß die Anzahl der positiven, verschwindenden bzw. negativen Eigenwerte von A gleich der Anzahl der positiven, verschwindenden bzw. negativen Diagonalelemente von D ist, und folgere hieraus, daß A genau dann positiv definit ist, wenn D nur positive Diagonalelemente besitzt.

Hinweis: Man beachte 1.2.C und die Tatsache, daß A und D kongruent sind.

Ü 6.1.5. Es sei A symmetrisch und regulär mit der BKP-Faktorisierung $PAP^{\intercal} = LDL^{\intercal}$ gemäß 6.1.7/6.1.10. Man zeige:

(i) Wenn in D ein (2,2)-Block $D^{(k)}$ auftritt, besitzt dieser einen Eigenwert $\lambda_k > 0$ und einen Eigenwert $\lambda_{k+1} < 0$.

Hinweis: Man beachte det $(D^{(k)}) = \delta_k < 0$, vgl. (51).

(ii) Die Anzahl der positiven Eigenwerte von A ist gleich der Anzahl der positiven (1, 1)-Blöcke von D plus der Anzahl der (2, 2)-Blöcke; die Anzahl der negativen Eigenwerte von Aist gleich der Anzahl der negativen (1, 1)-Blöcke plus der Anzahl der (2, 2)-Blöcke von D.

Ü 6.1.6. Es sei A wie in (25) mit p = 2, $a_{21} \neq 0$ und $\delta = \det(D) = a_{11}a_{22} - a_{21}^2 + 0$. Man führe von A ausgehend zwei einfache Eliminationsschritte mit den Pivots (2, 1) und (1, 2) gemäß

$$egin{aligned} &A^{(2)} := m{L}_2(-m{l}^2) \, A = [m{I} - m{l}^2 m{e}^{2 op}] \, A \,, \ &A^{(3)} := m{L}_1(-m{l}^1) \, A^{(2)} = [m{I} - m{l}^1 m{e}^{1 op}] \, A^{(2)} \end{aligned}$$

mit

$$\begin{split} \boldsymbol{l}^2 &:= (l_{12}, 0, l_{32}, \dots, l_{n2})^{\mathsf{T}}, \qquad l_{i2} := a_{i2}/a_{21} \qquad (i = 1, \dots, n; \ i \neq 2), \\ \boldsymbol{l}^1 &:= (0, 0, l_{31}, \dots, l_{n1})^{\mathsf{T}}, \qquad l_{i1} := a_{i1}^{(2)}/a_{12}^{(2)} \qquad (i = 3, \dots, n) \end{split}$$

durch und zeige, daß $A^{(3)} = \overline{L}_1 A$ mit

$$A^{(3)} = \left(\frac{\overline{D} \mid \overline{C}}{O \mid M}\right), \quad \overline{L}_1 = \left(\frac{F \mid O}{-G \mid I_{n-2}}\right), \quad F = \begin{pmatrix} 1 & l_{12} \\ 0 & 1 \end{pmatrix}, \quad \overline{D} = \begin{pmatrix} 0 & a_{12}^{(2)} \\ a_{21} & a_{22} \end{pmatrix}$$

sowie

 $\bar{D} = FD, \quad \bar{C}^{\intercal} = FC^{\intercal}$

gilt. Dabei sind M und G die aus dem Blockeliminationsschritt (26) herrührenden Matrizen. Die Restmatrix M und die Eliminationskoeffizienten l_{ij} (i = 3, ..., n; j = 1, 2) sind also dieselben wie bei der Blockelimination von $\{x_1, x_2\}$, während der Diagonalblock D auf die Gestalt \overline{D} transformiert wird.

Hinweis: Man beachte $\bar{L}_1 = L_1(-l^1) L_2(-l^2) = I - l^1 e^{1T} - (l^2 - l_{12} l^1) e^{2T}$.

6.2. Direkte Dreiecksfaktorisierungen

A. Nichtsymmetrische Matrizen

Im Abschnitt 5.1 wurde gezeigt, daß der Gaußsche Algorithmus mit Pivotsuche in den Spalten zu jedem regulären $A = (a_{ij}) \in \mathbb{R}^{n,n}$ eine Faktorisierung

$$PA = LR \tag{1}$$

mit Dreiecksmatrizen

$$\boldsymbol{L} = \begin{pmatrix} 1 & & & \\ & \ddots & & \\ l_{i1} \dots l_{i,i-1} & & \\ \vdots & & \ddots & \\ l_{n1} \dots \dots l_{n,n-1} & 1 \end{pmatrix}, \quad \boldsymbol{R} = \begin{pmatrix} r_{11} \dots r_{1j} \dots r_{1n} \\ \ddots & \vdots & \vdots \\ & r_{jj} \dots r_{jn} \\ & \ddots & \vdots \\ & & r_{nn} \end{pmatrix},$$
$$r_{jj} \neq 0 \quad (j = 1, \dots, n)$$
(2)

und einer Permutationsmatrix P liefert, siehe 5.1.4. Im folgenden werden wir die Elemente $\{l_{ij}, r_{ij}\}$ von L und R direkt aus der Gleichung (1) bestimmen. Zur Vereinfachung betrachten wir zunächst den Fall P = I, d. h., wir setzen voraus, daß die Zerlegung

$$\boldsymbol{A} = \boldsymbol{L}\boldsymbol{R} \tag{3}$$

mit L, R gemäß (2) existiert. Bedingungen dafür sind in Ü 5.2.5 angegeben worden, siehe auch 5.2.D. Unter Beachtung der Dreiecksgestalt von L, R liest sich (3) elementweise als

$$a_{ij} = \sum_{p=1}^{n} l_{ip} r_{pj} = \sum_{p=1}^{\min\{i,j\}} l_{ip} r_{pj}$$
 (*i*, *j* = 1, ..., *n*),

also

$$a_{ij} = \sum_{p=1}^{i} l_{ip} r_{pj} = \sum_{p=1}^{i-1} l_{ip} r_{pj} + r_{ij} \qquad (i \leq j)$$

bzw.

 $a_{ij} = \sum_{p=1}^{j} l_{ip} r_{pj} = \sum_{p=1}^{j-1} l_{ip} r_{pj} + l_{ij} r_{jj} \quad (i > j).$

Auflösen nach r_{ij} bzw. l_{ij} liefert

$$r_{ij} = a_{ij} - \sum_{p=1}^{i-1} l_{ip} r_{pj} \qquad (i \le j)$$
(4)

bzw.
$$l_{ij} = (a_{ij} - \sum_{p=1}^{j-1} l_{ip} r_{pj})/r_{jj}$$
 $(i > j).$ (5)

Die Gleichungen (4), (5) enthalten die Produkte der gesuchten Elemente $\{l_{ij}, r_{ij}\}$, allerdings in einer ganz speziellen Form: Auf den rechten Seiten kommen nur die links von der Position (i, j) stehenden Elemente der *i*-ten Zeile von L und die oberhalb dieser Position stehenden Elemente der j-ten Spalte von R vor. Wenn die nichttrivialen Elemente von L und R gemäß $(L \setminus R)$ in einer (n, n)-Matrix zusammengefaßt werden, ergibt sich folgende Situation:



Die l_{ij} und r_{ij} können daher sukzessive aus (4), (5) in jeder Reihenfolge berechnet werden, bei der links und oberhalb der Position (i, j) bereits alle Elemente vorhanden sind. Unter den in diesem Sinne zulässigen Reihenfolgen haben die folgenden praktische Bedeutung erlangt:

		1	1		1			
		3	2 3			3	5	
2		5	4	5	2		1	
4		↓				4	6	
	 (R			(B ₂)		(B _n)		

In (R_1) werden die Zeilen von R und die Spalten von L abwechselnd berechnet, in (R_2) bzw. (R_3) wird $(L \setminus R)$ zeilen- bzw. spaltenweise aufgebaut. Als Beispiel geben wir für die Reihenfolge (R_1) die detaillierten Formeln an.

6.2.1. Direkte LR-Faktorisierung ohne Pivotisierung. Die reguläre Matrix $A \in \mathbb{R}^{n,n}$ besitze die LR-Faktorisierung (2), (3). Dann ist das Verfahren

for k := 1(1)n do

for
$$j := k(1)n$$
 do $r_{kj} := a_{kj} - \sum_{p=1}^{k-1} l_{kp} * r_{pj}$
for $i := k + 1(1)n$ do $l_{ik} := \left(a_{ik} - \sum_{p=1}^{k-1} l_{ip} * r_{pk}\right) / r_{kk}$

in exakter Arithmetik durchführbar und liefert die nichttrivialen Elemente von L und R.

In der deutschsprachigen Literatur wird das Verfahren 6.2.1 häufig nach GAUSS-BANACHIEWICZ benannt, siehe B 6.2. Es entsteht aus der stufenweisen Grundform des Gaußschen Algorithmus, indem alle zum Element l_{ij} bzw. r_{ij} beitragenden Operationen zusammengefaßt und in der Form (4) bzw. (5) geschrieben werden. Man spricht deshalb auch vom verketteten Gaußschen Algorithmus. Insbesondere werden in 6.2.1 dieselben Rechenoperationen wie in 5.1 ausgeführt, wenn die Skalarprodukte in üblicher Weise in Standardarithmetik — also nach 2.3.9 — berechnet werden.

Wegen der fehlenden Pivotisierung ist die direkte Faktorisierung nicht für beliebiges reguläres A durchführbar und numerisch gutartig, vgl. 5.2 und 5.3. Im folgenden soll gezeigt werden, daß die zur Stabilisierung führende Spaltenpivotisierung in einfacher Weise realisiert werden kann. Dazu beachten wir, daß

$$a_{kk}^{(k)} := r_{kk} = a_{kk} - \sum_{p=1}^{k-1} l_{kp} r_{pk}$$
(6)

gerade das Pivotelement des k-ten Hauptschrittes des Gaußschen Algorithmus darstellt. Wenn statt $a_{kk}^{(k)}$ das Element $a_{sk}^{(k)}$ mit $k < s \leq n$ genommen würde, müßte (6) durch

$$\hat{a}_{sk} := a_{sk} - \sum_{p=1}^{k-1} l_{sp} r_{pk} \qquad (k+1 \le s \le n)$$
(7)

ersetzt werden. Ein Vergleich mit 6.2.1 zeigt, daß (7) den Zähler von l_{sk} darstellt. also ohnehin berechnet werden muß. Die Spaltenpivotisierung läuft dann darauf hinaus, unter allen Elementen (6), (7) ein betragsgrößtes zu bestimmen und die zugehörige Zeile mit der k-ten zu vertauschen. Eine zu 5.2.3 analoge Computerrealisierung auf dem Platz von A lautet dann wie folgt.

6.2.2. Direkte LR-Faktorisierung mit Spaltenpivotisierung

Aufgabe: Für $A = (a_{ij}) \in \mathbf{R}^{n,n}$ ist durch direkte Faktorisierung mit Spaltenpivotisierung eine Darstellung PA = LR zu berechnen. Die Matrix A ist mit den signifikanten Elementen von $L, R \in \mathbb{R}^{n,n}$ zu überspeichern, die Permutationsmatrix $P = T_{n-1,s(n-1)} \cdots T_{1,s(1)}$ ist durch die Zahlen $\{s(1), \ldots, s(n-1)\}$ zu charakterisieren. Falls für ein $k \in \{1, \ldots, n\}$ kein $r_{kk} \neq 0$ gefunden werden kann, ist abzubrechen und ie = -1 zu setzen, andernfalls wird ie = 0 gesetzt, und R ist regulär. Algorithmus:

ie := 0for k := 1(1)n do for k := 1(1)n do for i := k(1)n do $|a_{ik} := a_{ik} - \sum_{p=1}^{k-1} a_{ip} * a_{pk}|$ $|if |a_{ik}| > z$ then $[z := |a_{ik}|, s := i]$ if z = 0 then [ie := -1, stop]s(k) := sVertauschung: if k < s then for j := 1(1)n do $[z := a_{kj}, a_{kj} := a_{sj}, a_{sj} := z]$ Berechnung der r_{kj} : for j := k + 1(1)n do $a_{kj} := a_{kj} - \sum_{p=1}^{k-1} a_{kp} * a_{pj}$ Berechnung der l_{ik} : for i := k + 1(1)n do $a_{ik} := a_{ik}/a_{kk}$ Aufwand: $\sim n^{3}/3$ opms, $\sim n$ S (für s(1), ..., s(n - 1))

6.2.3. Bemerkung. (i) Bei Auswertung der Skalarprodukte gemäß 2.3.9 liefert 6.2.2 dieselben Ausgangsdaten wie 5.2.3, so daß die Ergebnisse der Rundungsfehleranalyse aus 5.3 wörtlich für die nach 6.2.2 berechnete Faktorisierung gelten. Der Vorteil von 6.1.2 liegt in der etwas kompakteren Form und der Tatsache, daß die Rechenoperationen im wesentlichen in Form von Skalarprodukten anfallen. In 6.2.2 können daher zur Skalarproduktberechnung spezielle Routinen — gegebenenfalls als Assemblerprogramme und mit Akkumulation in höherer Genauigkeit — verwendet werden, was in der stufenweisen Grundform 5.2.3 nicht möglich ist.

(ii) Spaltenpivotisierung ist auch bei Verwendung der Reihenfolge (R₃) möglich. In analoger Weise kann für die Reihenfolgen (R₁) und (R₂) eine Zeilenpivotisierung durch Vertauschung der Spalten realisiert werden. Dagegen ist die vollständige Pivotisierung mit dem Prinzip der direkten Faktorisierung nicht verträglich, da dazu alle $a_{ij}^{(k)}$ (i, j = k, ..., n) bekannt sein müssen. \Box

B. Symmetrische definite Matrizen

Für eine symmetrische und positiv definite Matrix A ist der Gaußsche Algorithmus ohne Pivotisierung durchführbar, und es gilt

$$L^{\mathsf{T}} = D^{-1}R$$
 mit $D = \text{diag}(d_k), \quad d_k := r_{kk} > 0 \quad (k = 1, ..., n),$ (8)

also

$$A = LDL^{\intercal}, \tag{9}$$

vgl. 6.1. A. Wegen (8) kann in diesem Fall die Laufanweisung zur Berechnung der
 l_{ik} in 6.2.1 durch

for i := k + 1(1)n do $l_{ik} := r_{ki}/r_{kk}$

ersetzt werden, so daß sich der Rechenaufwand auf $\sim n^3/6$ opms reduziert. Allerdings werden wie im nichtsymmetrischen Fall n^2 S zum Speichern der l_{ij} und r_{ij} benötigt, bei in-situ-Realisierung also alle Plätze von A. Wünschenswert ist dagegen ein Verfahren, bei dem sich wie in 6.1.1 die Symmetrie auch im Speicherplatzbedarf niederschlägt und das mit $\sim n^2/2$ S für L und D auskommt. Der sich anbietende Gedanke, auf das explizite Auftreten der r_{kj} überhaupt zu verzichten und sie überall durch $l_{jk}d_k$ zu ersetzen, führt bei Anwendung auf 6.2.1 zu der Vorschrift

for k := 1(1)n do

$$d_k := a_{kk} - \sum_{p=1}^{k-1} l_{kp} * l_{kp} * d_p$$

for $i := k + 1(1)n$ do $l_{ik} := \left(a_{ik} - \sum_{p=1}^{k-1} l_{ip} * l_{kp} * d_p\right) \left| d_k \right|$

Wegen der hier auftretenden Dreifachprodukte werden zur Realisierung $\sim (n^3/3 \text{ opm} + n^3/6 \text{ ops})$, also fast soviel Rechenoperationen wie im nichtsymmetrischen Fall benötigt, so daß dieses Vorgehen ausscheidet und doch mit den r_{ij} gearbeitet werden muß. Zur Senkung des Speicherbedarfs bietet sich dabei eine solche Reihenfolge an, bei der im k-ten Schritt nur auf die Elemente der k-ten Spalte von **R** zurückgegriffen wird. Dies ist für die durch das Schema



charakterisierte Reihenfolge (\mathbf{R}_4) der Fall, bei der die l_{ij} zeilenweise berechnet werden. Unter Verwendung der Bezeichnung

 $r_i := r_{ik} = l_{ki}d_i$ (j = 1, ..., k - 1)

ergibt sich mit dieser Reihenfolge der folgende Algorithmus:

6.2.4. Direkte LDL^{T} -Faktorisierung. Es sei $A \in S^{n,n}$ positiv definit. Dann ist das Verfahren

for k := 1(1)n do

$$\begin{bmatrix} \text{for } j := 1(1)k - 1 \text{ do} \\ \begin{bmatrix} r_j := a_{kj} - \sum_{p=1}^{j-1} l_{jp} * r_p \\ \\ l_{kj} := r_j/d_j \end{bmatrix}$$
$$d_k := a_{kk} - \sum_{p=1}^{k-1} l_{kp} * r_p$$

in exakter Arithmetik mit $d_k > 0$ (k = 1, ..., n) durchführbar und liefert die Faktorisierung $A = LDL^{\intercal}$.

Autward: $\sim n^3/6$ opms, $\sim n$ S (für r_1, \ldots, r_{n-1})

6.2.5. Bemerkung. (i) Algorithmus 6.2.4 kann analog zu 6.1.1 auf dem Platz des unteren Dreiecks von A ausgeführt werden. Wenn auf die Möglichkeit der Berechnung von d_k über ein Skalarprodukt verzichtet wird, können die $r_j = r_{jk}$ sogar auf dem Platz von l_{kj} — bei in-situ-Realisierung also auf dem von a_{kj} — gespeichert werden. Dazu ist 6.2.4 wie folgt zu modifizieren, vgl. auch 6.1.2(ii):

for
$$k := 1(1)n$$
 do

for
$$j := 1(1)k - 1$$
 do $l_{kj} := a_{kj} - \sum_{p=1}^{j-1} l_{jp} * l_{kp}$
 $d_k := a_{kk}$
for $j := 1(1)k - 1$ do $[\varrho := l_{kj}, l_{kj} := \varrho/d_j, d_k := d_k - l_{kj} * \varrho]$

Ein zu 6.2.4 analoger Algorithmus kann für die Reihenfolge (R_3), d. h. für spaltenweise Berechnung von L und D angegeben werden, während die Diagonalpivotisierung entsprechend 6.1.5(iii) nicht mit dem Prinzip der direkten Berechnung von Lund D verträglich ist.

(ii) Bei Auswertung der Skalarprodukte nach 2.3.9 liefert 6.2.4 dieselben Ausgangsdaten wie 6.1.1, so daß die Ergebnisse der Rundungsfehleranalyse 6.1.3 wörtlich übernommen werden können. \Box

Für die Cholesky-Faktorisierung

$$\boldsymbol{A} = \boldsymbol{\hat{L}}\boldsymbol{\hat{L}}^{\mathsf{T}} \tag{10}$$

mit

$$\hat{\boldsymbol{L}} := \boldsymbol{L} \boldsymbol{D}^{1/2}, \quad \text{also} \quad \hat{l}_{ij} = l_{ij} \sqrt[j]{d_j} \qquad (i \ge j), \tag{11}$$

geht 6.2.4 in das folgende, nach CHOLESKY benannte Verfahren zur Berechnung von \hat{L} über. Die bei der LDL^{\intercal} -Faktorisierung beobachteten Speicherkonflikte treten dabei wie in der (n - 1)-stufigen Grundform 6.1.4 nicht mehr auf.

6.2.6. Direkte $\hat{L}\hat{L}^{\intercal}$ -Faktorisierung. Es sei $A \in S^{n,n}$ positiv definit. Dann ist das Verfahren

for k := 1(1)n do $\begin{vmatrix} \text{for } j := 1(1)k - 1 \text{ do } \hat{l}_{kj} := \left(a_{kj} - \sum_{p=1}^{j-1} \hat{l}_{jp} * \hat{l}_{kp}\right) | \hat{l}_{jj} \\ \hat{l}_{kk} := \operatorname{sqrt} \left(a_{kk} - \sum_{p=1}^{k-1} \hat{l}_{kp}^2\right) \end{vmatrix}$

in exakter Arithmetik mit $\hat{l}_{kk} > 0$ (k = 1, ..., n) durchführbar und liefert die signifikanten Elemente von \hat{L} in der Faktorisierung $A = \hat{L}\hat{L}^{\intercal}$.

Aufwand: $\sim n^3/6 \text{ opms} + n \text{ opr}$

Die Formel für \hat{l}_{kk} aus 6.2.6 läßt sich auch in der nach a_{kk} aufgelösten Form

$$\hat{l}_{k1}^2 + \hat{l}_{k2}^2 + \dots + \hat{l}_{kk}^2 = a_{kk}$$

schreiben. Hieraus folgt

$$|\hat{l}_{kj}| \leq \sqrt{a_{kk}} \quad (k \geq j) \quad \text{sowie} \quad \|\hat{L}\|_F = \left(\sum_{k=1}^n a_{kk}\right)^{1/2},$$
(12)

d. h., der Dreiecksfaktor \hat{L} ist auch ohne Pivotisierung a priori beschränkt, vgl. 5.2.9 und Ü 6.1.1.

6.2.7. Bemerkung. (i) Das Cholesky-Verfahren 6.2.6 kann analog zu 6.1.4 auf dem Platz des unteren Dreiecks von A ausgeführt werden. Spaltenweise Berechnung der \hat{l}_{ik} ist nach der Vorschrift

for k := 1(1)n do

$$a_{kk} := \operatorname{sqrt} \left(a_{kk} - \sum_{p=1}^{k-1} (a_{kp})^2 \right)$$

for $i := k + 1(1)n$ do $a_{ik} := \left(a_{ik} - \sum_{p=1}^{k-1} a_{ip} * a_{kp} \right) / a_{kk}$

möglich, während die Diagonalpivotisierung entsprechend 6.1.5(iii) nicht mit dem Prinzip der direkten Berechnung von \hat{L} verträglich ist.

(ii) Bei üblicher Auswertung der Skalarprodukte gemäß 2.3.9 sind die Ausgangsdaten von 6.2.6 mit denen von 6.1.4 identisch, so daß die Rundungsfehleranalyse aus 6.1.5 (ii) auch hier gültig ist.

(iii) Direkte Methoden zur Berechnung der BKP-Faktorisierung indefiniter symmetrischer Matrizen sind nicht bekannt und wegen der komplizierten Pivotstrategie wohl auch nicht möglich.

Übungsaufgaben

Ü 6.2.1. Man gebe die 6.2.1. entsprechenden Verfahren zur direkten LR-Faktorisierung für die Reihenfolgen (R_2) und (R_3) an und mache sich die Unterschiede im Zugriff auf die Elemente von A, L und R klar. Wie sind die Algorithmen zu modifizieren, wenn für (R_2) mit Zeilenpivotisierung und für (R_3) mit Spaltenpivotisierung gearbeitet werden soll?

Ü 6.2.2. Man modifiziere 6.2.4 so, daß die Elemente von L und D spaltenweise nach der Reihenfolge (R_3) berechnet werden.

6.3. Aufdatierung von Dreiecksfaktorisierungen

Bei gewissen praktischen Aufgaben — etwa in der linearen und nichtlinearen Optimierung, bei der Minimierung nichtlinearer Funktionen oder bei der Lösung nichtlinearer Gleichungssysteme — ist nicht nur ein einzelnes Gleichungssystem, sondern eine Folge

$$\boldsymbol{A}_{\boldsymbol{k}}\boldsymbol{x} = \boldsymbol{b}^{\boldsymbol{k}} \qquad (\boldsymbol{k} = 1, \dots, K) \tag{1}$$

von K Systemen zu lösen, wobei A_{k+1} aus A_k durch eine Rang-1-Modifikation

$$A_{k+1} = A_k + u^k v^{k\intercal} \tag{2}$$

entsteht. Es wäre dann wünschenswert, die Dreiecksfaktorisierung von $A_k \in \mathbb{R}^{n,n}$ nicht für jedes k mit dem Aufwand von

$$\sim K_1 n^3 \text{ opms}$$
 (3)

 $(K_1 = 1/3$ für nichtsymmetrisches A_k , $K_1 = 1/6$ für symmetrisches A_k) neu berechnen zu müssen, sondern die Faktoren von A_{k+1} aus denen von A_k mit geringerem Aufwand zu bestimmen. Im folgenden werden wir zeigen, daß dies in der Tat mit

$$\sim K_2 n^2 \text{ opms}$$
 (4)

möglich ist.

Zur Vereinfachung der Schreibweise bezeichnen wir A_k mit A und die durch die Modifikation (2) entstehende Matrix mit

$$\bar{\boldsymbol{A}} = \boldsymbol{A} + \boldsymbol{u}\boldsymbol{v}^{\mathsf{T}}; \tag{5}$$

die Faktoren werden analog gekennzeichnet. Die Matrizen A und \overline{A} seien beide regulär.

A. Nichtsymmetrische Matrizen

Wir betrachten zunächst die folgende

Hilfsaufgabe: Zu gegebenem

$$r = \begin{pmatrix} r_1 \\ r_2 \end{pmatrix}, \quad L = \begin{pmatrix} 1 & 0 \\ l_{21} & 1 \end{pmatrix} \text{ mit } |l_{21}| \leq \lambda$$

sind $F = \begin{pmatrix} \alpha & \beta \\ \gamma & \delta \end{pmatrix}$ regulär und $T = T_{11} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$ oder $T = T_{12} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ so festzu-

legen, daß

$$Lr = T(TLF^{-1})(Fr) =: T\overline{L}\overline{r}$$
(6)

mit

$$\overline{\boldsymbol{L}} := \boldsymbol{T} \boldsymbol{L} \boldsymbol{F}^{-1} = \begin{pmatrix} 1 & 0 \\ \overline{l}_{21} & 1 \end{pmatrix} \quad \text{und} \quad \overline{\boldsymbol{r}} := \boldsymbol{F} \boldsymbol{r} = \begin{pmatrix} \overline{r}_1 \\ 0 \end{pmatrix}$$
(7)

gilt. 🗌

Für die weiteren Überlegungen führen wir die Größen

$$\boldsymbol{Lr} = \begin{pmatrix} 1 & 0 \\ l_{21} & 1 \end{pmatrix} \begin{pmatrix} r_1 \\ r_2 \end{pmatrix} = \begin{pmatrix} r_1 \\ l_{21}r_1 + r_2 \end{pmatrix} =: \begin{pmatrix} q_1 \\ q_2 \end{pmatrix} =: q$$

ein und unterscheiden zwei Fälle.

Fall 1: $|q_1| \geq |q_2|$.

Hier wird $T := T_{11}$ gesetzt. Aus (7) folgt dann $\overline{L} = LF^{-1}$, also $F = \overline{L}^{-1}L$. Als Produkt unterer Einsdreiecksmatrizen muß F daher auch von Einsdreiecksform sein:

$$\boldsymbol{F} = \begin{pmatrix} 1 & 0 \\ \gamma & 1 \end{pmatrix}, \qquad \boldsymbol{F}^{-1} = \begin{pmatrix} 1 & 0 \\ -\gamma & 1 \end{pmatrix}.$$
(8)

Der noch freie Parameter γ ergibt sich aus der zweiten Bedingung (7), d. h. aus

$$\bar{\boldsymbol{r}} = \begin{pmatrix} \bar{r}_1 \\ 0 \end{pmatrix} = \boldsymbol{F}\boldsymbol{r} = \begin{pmatrix} r_1 \\ \gamma r_1 + r_2 \end{pmatrix}$$

im Fall $r_1 \neq 0$ zu

$$\gamma := -r_2/r_1$$

und kann wegen $|r_1|=|q_1|\ge |q_2|=|l_{21}r_1+r_2|\ge |r_2|-\lambda\,|r_1|$ gemäß $|\gamma|\le 1+\lambda$ abgeschätzt werden. Daraus folgt

$$l_{21} = l_{21} - \gamma = q_2/q_1$$

mit $|\tilde{l}_{21}| \leq 1$. Im Fall $r_1 = r_2 = 0$ setzen wir $\gamma = 0$, also F = I, was auf $\tilde{l}_{21} = l_{21}$ mit $|\tilde{l}_{21}| \leq \lambda$ führt.

Fall 2: $|q_1| < |q_2|$. Wir setzen jetzt $T := T_{12}$. Aus (7) folgt dann

$$m{F} = m{ar{L}}^{-1} m{T}_{12} m{L} = egin{pmatrix} 1 & 0 \ -ar{l}_{21} & 1 \end{pmatrix} egin{pmatrix} l_{21} & 1 \ 1 & 0 \end{pmatrix} = egin{pmatrix} rac{l_{21}}{-ar{l}_{21}l_{21} + 1} & rac{1}{-ar{l}_{21}} \end{pmatrix}.$$

Der noch freie Wert l_{21} ergibt sich wieder aus der zweiten Forderung (7), d. h. aus

$$\bar{\boldsymbol{r}} = \begin{pmatrix} \bar{r}_1 \\ 0 \end{pmatrix} = \boldsymbol{F}\boldsymbol{r} = \begin{pmatrix} l_{21}r_1 + r_2 \\ -\bar{l}_{21}(l_{21}r_1 + r_2) + r_1 \end{pmatrix} = \begin{pmatrix} q_2 \\ -\bar{l}_{21}q_2 + q_1 \end{pmatrix}$$

zu

 $l_{21} := q_1/q_2$

mit $|\bar{l}_{21}| \leq 1$. Damit erhalten wir

$$\boldsymbol{F} = \begin{pmatrix} \alpha & 1 \\ \gamma & \delta \end{pmatrix}, \qquad \boldsymbol{F}^{-1} = \begin{pmatrix} -\delta & 1 \\ \gamma & -\alpha \end{pmatrix}$$
(9)

 mit

$$\alpha := l_{21}, \quad \delta := -\bar{l}_{21} = -q_1/q_2, \quad \gamma := -\bar{l}_{21}l_{21} + 1 = \alpha\delta + 1 = r_2/q_2,$$

wobei $|\alpha| \leq \lambda$, $|\delta| \leq 1$ und $|\gamma| \leq 1 + \lambda$ ist.

Wenn F von links auf einen Spaltenvektor bzw. F^{-1} von rechts auf einen Zeilenvektor gemäß

$$ar{m{w}}:=m{F}m{w}$$
 bzw. $ar{m{f}}^{\intercal}:=m{f}^{\intercal}m{F}^{-1}$ mit $m{w},m{f}\in \mathbf{R}^2$

14 Schwetlick, Numerische Algebra

angewendet wird, ergibt sich im Fall 1 mit (8)

$$egin{array}{lll} \overline{w}_1 := w_1, & ar{f}_1 := f_1 - \gamma f_2 \ \overline{w}_2 := w_2 + \gamma w_1 & ar{f}_2 := f_2, \end{array}$$

d. h., die Berechnung erfordert jeweils 1 opms. Im Fall 2 würde sich mit (9) dagegen

$$egin{array}{lll} \overline{w}_1 := w_2 + lpha w_1, & egin{array}{llll} \overline{f}_1 := -\delta f_1 + \gamma f_2, & \ \overline{w}_2 := \delta w_2 + \gamma w_1 & egin{array}{lllllllll} \overline{f}_2 := f_1 - lpha f_2 & \ \end{array}$$

ergeben, d. h., der Aufwand wäre jeweils 3 opm + 2 ops. Unter Beachtung von $\gamma = \alpha \delta + 1$ läßt sich jedoch eine Multiplikation einsparen, indem nach der in exakter Arithmetik äquivalenten Vorschrift

vorgegangen wird. Dies entspricht der Darstellung von F bzw. F^{-1} gemäß

$$\boldsymbol{F} = \begin{pmatrix} 1 & 0 \\ \delta & 1 \end{pmatrix} \begin{pmatrix} \alpha & 1 \\ 1 & 0 \end{pmatrix} \quad \text{bzw.} \quad \boldsymbol{F}^{-1} = \begin{pmatrix} 0 & 1 \\ 1 & -\alpha \end{pmatrix} \begin{pmatrix} 1 & 0 \\ -\delta & 1 \end{pmatrix}, \tag{11}$$

bei der nur α und δ explizit vorkommen. Der Aufwand ist jetzt jeweils 2 opms.

In den nachfolgenden Anwendungen zur Aufdatierung von Dreiecksfaktorisierungen stellen die Matrizen L, \overline{L} der Hilfsaufgabe gewisse (2,2)-Diagonalblöcke von (n,n)-Einsdreiecksmatrizen L, \overline{L} in den Positionen $\{k, k + 1\}$ dar, und die Vektoren r, \overline{r} der Hilfsaufgabe sind die zugehörigen Zweierblöcke $(r_k, r_{k+1})^{\mathsf{T}}$ bzw. $(\overline{r}_k, \overline{r}_{k+1})^{\mathsf{T}}$ von n-Vektoren r, \overline{r} . Die Matrix F tritt demgemäß in der Gestalt

auf, und T_{12} ist durch $T_{k,k+1}$ zu ersetzen. In dieser Erweiterung erhalten wir als Zusammenfassung der obigen Überlegungen das folgende Resultat.

6.3.1. Dreiecksformerhaltender Elementareliminationsschritt. Gegeben seien eine untere Eindreiecksmatrix $L = (l_{ij}) \in \mathbb{R}^{n.n}$, ein Vektor $r = (r_i) \in \mathbb{R}^n$ und ein Index $k \in \{1, ..., n - 1\}$. Die Matrix $F_{k,k+1}$ des Typs (12) und die Vertauschungsmatrix T_{ks} mit $s = s(k) \in \{k, k + 1\}$ seien durch die Zahlen $\{k, l_{k+1,k}, r_k, r_{k+1}\}$

211

wie folgt festgelegt:

S1: $q_k := r_k, q_{k+1} := l_{k+1,k} * r_k + r_{k+1}$ S2: if $|q_{k}| \ge |q_{k+1}|$ then Fall 1: s(k):=k $lpha:=\delta:=1,\,eta:=0,\, ext{if}\,\,q_k=0 ext{ then }\gamma:=0 ext{ else }\gamma:=-r_{k+1}/r_k$ else Fall 2: $s(k):=k+1 \ lpha:=l_{k+1,k},\,eta:=1,\,\gamma:=r_{k+1}/q_{k+1},\,\delta:=-q_k/q_{k+1}$

Dann gilt

$$Lr = T_{ks}\overline{L}\overline{r}$$
 mit $\overline{L} := T_{ks}LF_{k,k+1}^{-1}$, $\overline{r} := F_{k,k+1}r$.

Die Matrix \overline{L} ist eine untere Einsdreiecksmatrix und unterscheidet sich von Lbis auf die im Fall 2 vorzunehmenden Vertauschungen der Zeilen $\{k, k+1\}$ nur in den Spalten $\{k, k+1\}$, und der Vektor \overline{r} unterscheidet sich von r nur in den Komponenten $\{k, k+1\}$. Dabei ist

$$ar{l}_{k+1,k} = egin{cases} l_{k+1,k} & ext{im Fall 1, } q_k = 0, \ q_{k+1}/q_k & ext{im Fall 1, } q_k \pm 0, \ q_k/q_{k+1} & ext{im Fall 2; } & ar{r}_k = egin{cases} q_k & ext{im Fall 1, } q_k = 0, \ q_{k+1} & ext{im Fall 2; } & ar{r}_{k+1} = 0. \end{cases}$$

Der Vollständigkeit halber sind in 6.3.1 auch die explizit nicht benötigten Parameter α , β , δ im Fall 1 bzw. β , γ im Fall 2 mit angegeben worden.

Unter der Voraussetzung

$$|l_{k+1,k}| \le \lambda \tag{13}$$

folgt aus 6.3.1

$$|\tilde{l}_{k+1,k}| \leq \max\{\lambda, 1\}, \tag{14}$$

und die signifikanten Blöcke F bzw. F^{-1} von $F_{k,k+1}$ bzw. $F_{k,k+1}^{-1}$ genügen den Abschätzungen

$$|\boldsymbol{F}| \leq \left(\frac{\max\left\{\lambda, 1\right\} \mid 1}{1+\lambda \mid 1}\right), \quad |\boldsymbol{F}^{-1}| \leq \left(\frac{1 \mid 1}{1+\lambda \mid \max\left\{\lambda, 1\right\}}\right), \quad (15)$$

insbesondere ist

$$\|\boldsymbol{F}\|_{\infty} \leq 2 + \lambda, \qquad \|\boldsymbol{F}^{-1}\|_{\infty} \leq 1 + \lambda + \max\{\lambda, 1\}.$$
(16)

Wenn die Einsdreiecksmatrix L den nach dem Gaußschen Algorithmus mit Spaltenpivotisierung berechneten Dreiecksfaktor darstellt, sind die Elemente betragsmäßig durch 1 beschränkt, speziell gilt

$$|l_{k+1,k}| \le 1 \qquad (k = 1, ..., n - 1) \tag{17}$$

für die in 6.3.1 verwendeten Subdiagonalelemente. Wegen (14) gilt dann (17) auch für die modifizierte Matrix \overline{L} , d. h., die Gültigkeit von (17) ist invariant unter den Elementareliminationsschritten 6.3.1. Die übrigen Elemente von $\{L, r\}$ können beim Übergang zu $\{\bar{L}, \bar{r}\}$ durchaus anwachsen, wegen (15), (16) jedoch nicht beliebig größer werden: Unter der Voraussetzung (17) kann das Betragsmaximum nach der Elimination höchstens um den Faktor 3 größer als vorher sein.

Algorithmus 6.3.1 wird jetzt (n-1)-mal angewendet, um sukzessive die Komponenten r_i (i = n, ..., 2) von r zu 0 zu machen und dabei die Dreiecksform eines linken L-Faktors zu erhalten.

6.3.2. Dreiecksformerhaltende Transformation von r in $\rho_1 e^1$. Gegeben seien die untere Einsdreiecksmatrix $L \in \mathbb{R}^{n,n}$ und der Vektor $r \in \mathbb{R}^n$. Die Matrizen L_1 und

untere Einsdreiecksmatrix $L \in \mathbb{R}^{n,n}$ und der Vektor $r \in \mathbb{R}^{n}$. Die der Vektor r^{1} seien rekursiv definiert durch die Vorschrift S1: Setze $L^{(n-1)} := L$, $r^{(n-1)} := r$ S2: for k := n - 1(-1)1 do $\begin{bmatrix} Lege \{F_{k,k+1}, T_{k,s(k)}\} \text{ gemäß } 6.3.1 \text{ durch} \\ \{k, (L^{(k)})_{k+1,k}, r_{k}^{(k)}, r_{k+1}^{(k)}\} \text{ fest.} \\ \text{Bestimme } L^{(k-1)} := T_{k,s(k)}L^{(k)}F_{k,k+1}^{-1}, r^{(k-1)} := F_{k,k+1}r^{(k)}$ S3: Setze $L_{1} := L^{(0)}, r^{1} := r^{(0)}$. Denn gilt

Dann gilt

$$Lr = P_1^{\mathsf{T}}(P_1LF_1^{-1}) (F_1r) = P_1^{\mathsf{T}}L_1r^1$$

mit der unteren Einsdreiecksmatrix $L_1 := P_1 L F_1^{-1}$ und dem Vektor $r^1 := F_1 r$ = $o_1 e^1$, wobei $= \varrho_1 e^1$, wobei

$$F_1 := F_{12}F_{23}\cdots F_{n-1,n}, \quad P_1 := T_{1,s(1)}T_{2,s(2)}\cdots T_{n-1,s(n-1)}$$

ist.

Im Fall n = 4 erfolgt die Transformation 6.3.2 nach folgendem Muster:

Die eingerahmten Felder geben dabei diejenigen Elemente an, die im jeweiligen

Schritt neu berechnet werden; gegebenenfalls zu vertauschende Elemente sind durch Kreise gekennzeichnet.

Es sei jetzt PA = LR eine Dreiecksfaktorisierung von A. Dann kann (5) in der Form

$$\tilde{\boldsymbol{A}} = \boldsymbol{P}^{\mathsf{T}}\boldsymbol{L}\boldsymbol{R} + \boldsymbol{u}\boldsymbol{v}^{\mathsf{T}} = \boldsymbol{P}^{\mathsf{T}}\boldsymbol{L}(\boldsymbol{R} + \boldsymbol{r}\boldsymbol{v}^{\mathsf{T}}) \quad \text{mit} \quad \boldsymbol{r} := \boldsymbol{L}^{-1}\boldsymbol{P}\boldsymbol{u} \tag{18}$$

geschrieben werden. Anwendung der Transformation 6.3.2 führt auf

$$\tilde{\boldsymbol{A}} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{P}_{1}^{\mathsf{T}} (\boldsymbol{P}_{1} \boldsymbol{L} \boldsymbol{F}_{1}^{-1}) \left(\boldsymbol{F}_{1} \boldsymbol{R} + \boldsymbol{F}_{1} \boldsymbol{r} \boldsymbol{v}^{\mathsf{T}} \right) =: \boldsymbol{P}_{2}^{\mathsf{T}} \boldsymbol{L}_{1} (\boldsymbol{R}_{1} + \boldsymbol{r}^{1} \boldsymbol{v}^{\mathsf{T}})$$
(19)

mit $P_2 := P_1 P$ und $R_1 := F_1 R$. Dabei wird R_1 analog zu r^1 rekursiv gemäß

$$\mathbf{R}^{(n-1)} := \mathbf{R}$$
, for $k := n - 1(1)1$ do $\mathbf{R}^{(k-1)} := \mathbf{F}_{k,k+1}\mathbf{R}^{(k)}$, $\mathbf{R}_1 := \mathbf{R}^{(0)}$ (20)

berechnet. Für n = 4 verläuft die Bildung von \mathbf{R}_1 nach dem Muster



Die Subdiagonalelemente von R werden also sukzessive aufgefüllt. Wir erinnern daran, daß eine Matrix der Form R_1 , d. h. eine Matrix R_1 mit

$$(\mathbf{R}_1)_{ij} = 0$$
 $(i > j + 1)$ (21)

obere Hessenbergmatrix genannt wird; in analoger Weise ist eine untere Hessenbergmatrix definiert, vgl. 1.1.C.

Mit \mathbf{R}_1 ist in (19) dann auch $\mathbf{R}_2 := \mathbf{R}_1 + \mathbf{r}^1 \mathbf{v}^{\mathsf{T}} = \mathbf{R}_1 + \varrho_1 \mathbf{e}^1$ von oberer Hessenbergform, denn lediglich zur ersten Zeile von \mathbf{R}_1 wird $\varrho_1 \mathbf{v}^{\mathsf{T}}$ addiert. Mit $\mathbf{L}_2 := \mathbf{L}_1$ läßt (19) sich somit als

$$\bar{\boldsymbol{A}} = \boldsymbol{P}_2^{\mathsf{T}} \boldsymbol{L}_2 \boldsymbol{R}_2 \tag{22}$$

schreiben. Die noch störenden Subdiagonalelemente $(\mathbf{R}_2)_{k+1,k}$ (k = 1, ..., n - 1) von \mathbf{R}_2 können dann wieder durch n - 1 Elementareliminationsschritte sukzessive zu 0 gemacht werden, wobei die k-te Spalte von \mathbf{R}_2 im k-ten Schritt die Rolle von \mathbf{r} übernimmt.

6.3.3. Dreiecksformerhaltende Transformation von \mathbf{R}_2 in $\overline{\mathbf{R}}$. Gegeben seien die untere Einsdreiecksmatrix $L_2 \in \mathbf{R}^{n,n}$ und die obere Hessenbergmatrix $\mathbf{R}_2 \in \mathbf{R}^{n,n}$. Die Matrizen $\overline{\mathbf{L}}, \overline{\mathbf{R}}$ seien rekursiv definiert durch die Vorschrift

$$\begin{array}{l} \text{S1: Setze } \boldsymbol{L}^{(1)} := \boldsymbol{L}_{2}, \, \boldsymbol{R}^{(1)} := \boldsymbol{R}_{2} \\ \text{S2: for } k := 1(1)n \, - \, 1 \, \operatorname{do} \\ \\ & \left| \begin{array}{c} \operatorname{Lege} \left\{ \tilde{\boldsymbol{F}}_{k,k+1}, \, \tilde{\boldsymbol{T}}_{k,s(k)} \right\} \, \operatorname{gemäß} \, 6.3.1 \, \operatorname{durch} \\ \left\{ k, \, (\boldsymbol{L}^{(k)})_{k+1,k}, \, (\boldsymbol{R}^{(k)})_{kk}, \, (\boldsymbol{R}^{(k)})_{k+1,k} \right\} \, \operatorname{fest.} \\ & \operatorname{Bestimme} \, \boldsymbol{L}^{(k+1)} := \tilde{\boldsymbol{T}}_{k,s(k)} \boldsymbol{L}^{(k)} \tilde{\boldsymbol{F}}_{k,k+1}^{-1}, \, \boldsymbol{R}^{(k+1)} := \tilde{\boldsymbol{F}}_{k,k+1} \boldsymbol{R}^{(k)} \end{array} \right. \end{array}$$

S3: Setze $\overline{L} := L^{(n)}, \, \overline{R} := R^{(n)}$

Dann gilt

$$\boldsymbol{L}_{2}\boldsymbol{R}_{2} = \boldsymbol{P}_{3}^{\mathsf{T}}(\boldsymbol{P}_{3}\boldsymbol{L}_{2}\boldsymbol{F}_{3}^{-1}) \left(\boldsymbol{F}_{3}\boldsymbol{R}_{2}\right) = \boldsymbol{P}_{3}^{\mathsf{T}}\boldsymbol{\overline{L}}\boldsymbol{\overline{R}}$$

$$\tag{23}$$

mit der unteren Einsdreiecksmatrix $\overline{L}:=P_3L_2F_3^{-1}$, der oberen Dreiecksmatrix $\overline{R}:=F_3R_2$ und

$$\mathbf{F}_3 := \tilde{\mathbf{F}}_{n-1,n} \cdots \tilde{\mathbf{F}}_{23} \tilde{\mathbf{F}}_{12}, \quad \mathbf{P}_3 := \tilde{\mathbf{T}}_{n-1,s(n-1)} \cdots \tilde{\mathbf{T}}_{2,s(2)} \tilde{\mathbf{T}}_{1,s(1)}.$$
(24)

Einsetzen von (23) in (22) führt auf die gesuchte Faktorisierung von \bar{A} . Zusammenfassend erhalten wir den folgenden Algorithmus:

6.3.4. Aufdatierung einer Dreiecksfaktorisierung bei Rang-1-Modifikation

Aufgabe: Bestimme eine Dreiecksfaktorisierung $\overline{P}\overline{A} = \overline{L}\overline{R}$ von

 $ar{A} = A + uv^{\intercal}$

bei gegebener Dreiecksfaktorisierung PA = LR von A mit unteren Einsdreiecksmatrizen L, \overline{L} , oberen Dreiecksmatrizen R, \overline{R} und Permutationsmatrizen P, \overline{P} .

Algorithmus:

S1: Stelle u in der Form

$$\boldsymbol{u} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L} \boldsymbol{r}, \, \boldsymbol{r} \text{ Lösung des Dreieckssystems } \boldsymbol{L} \boldsymbol{r} = \boldsymbol{P} \boldsymbol{u},$$
 (25)

dar. Schreibe \bar{A} als

$$\tilde{\boldsymbol{A}} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L} (\boldsymbol{R} + \boldsymbol{r} \boldsymbol{v}^{\mathsf{T}}) \tag{26}$$

S2: Führe n - 1 Elementareliminationsschritte gemäß 6.3.2 zur Transformation von r in $\rho_1 e^1$ aus und transformiere R gemäß (20) auf obere Hessenbergform $R_1 := F_1 R$. Schreibe (26) als

$$\bar{\boldsymbol{A}} = (\boldsymbol{P}^{\mathsf{T}}\boldsymbol{P}_{1}^{\mathsf{T}}) (\boldsymbol{P}_{1}\boldsymbol{L}\boldsymbol{F}_{1}^{-1}) (\boldsymbol{F}_{1}\boldsymbol{R} + \boldsymbol{F}_{1}\boldsymbol{r}\boldsymbol{v}^{\mathsf{T}}) = \boldsymbol{P}_{2}^{\mathsf{T}}\boldsymbol{L}_{1}(\boldsymbol{R}_{1} + \varrho_{1}\boldsymbol{e}^{1}\boldsymbol{v}^{\mathsf{T}}) = \boldsymbol{P}_{2}^{\mathsf{T}}\boldsymbol{L}_{2}\boldsymbol{R}_{2} (27)$$

mit der Permutationsmatrix $P_2 := P_1 P_2$, der unteren Einsdreiecksmatrix $L_2 := L_1 := P_1 L F_1^{-1}$ und der oberen Hessenbergmatrix $R_2 = R_1 + \varrho_1 e^1 v^{\intercal} = F_1 R + F_1 r v^{\intercal}$.

S3: Führe n-1 Elementareliminationsschritte gemäß 6.3.3 zur Transformation

von R_2 auf obere Dreiecksform \overline{R} aus. Schreibe (27) als

$$\bar{\boldsymbol{A}} = (\boldsymbol{P}_2^{\mathsf{T}} \boldsymbol{P}_3^{\mathsf{T}}) (\boldsymbol{P}_3 \boldsymbol{L}_2 \boldsymbol{F}_3^{-1}) (\boldsymbol{F}_3 \boldsymbol{R}_2) = \boldsymbol{\overline{P}}^{\mathsf{T}} \boldsymbol{\overline{L}} \boldsymbol{\overline{R}}$$
(28)

mit der Permutationsmatrix $\overline{P} := P_3 P_2$, der unteren Einsdreiecksmatrix $\overline{L} := P_3 L_2 F_3^{-1}$ und der oberen Dreiecksmatrix $\overline{R} := F_3 R_2$. Aufwand: $\sim K_2 n^2$ opms mit $5/2 \leq K_2 \leq 9/2$

6.3.5. Bemerkung. (i) Algorithmus 6.3.4 läßt sich auf dem Platz von L, R und u durchführen, indem u mit r und danach mit den Subdiagonalelementen von R_2 überspeichert wird. Die Permutationsmatrix P kann z. B. durch die Tabelle $\{k(1), \ldots, k(n)\}$ als Integervektor entsprechend 3.1 repräsentiert werden. Der Multiplikation $T_{i,s(i)}P$ entspricht dann die Vertauschung der Elemente i und s(i) der Tabelle, vgl. 5.2.6(ii). Der Aufwand ist mindestens $\sim 5/2n^2$ opms (wenn nur Fall 1 eintritt) und höchstens $\sim 9/2n^2$ opms (wenn stets Fall 2 eintritt).

(ii) Wenn bei der Faktorisierung neben Zeilenvertauschungen auch solche der Spalten gemäß

$$\boldsymbol{P}_{\boldsymbol{Z}}\boldsymbol{A}\boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} = \boldsymbol{L}\boldsymbol{R} \tag{29}$$

zugelassen werden, läßt sich für den Sonderfall $v = e^k$, d. h. für Spaltenmodifikationen

$$\bar{A} = A + ue^{k\tau} \tag{30}$$

eine LR-Faktorisierung

$$\overline{P}_{Z}\overline{A}\overline{P}_{S}^{\mathsf{T}} = \overline{L}\overline{R} \tag{31}$$

von \bar{A} wie folgt berechnen:

S1: Schreibe (30) bei gegebener Faktorisierung (29) in der Form

$$\bar{A} = P_Z^{\mathsf{T}} L R P_S + u e^{k\mathsf{T}} = P_Z^{\mathsf{T}} L [R + r e^{l\mathsf{T}}] P_S$$
(32)

 mit

 $r := L^{-1}P_z u$ und $e^l := P_s e^k$.

S2: Es bezeichne P_1 diejenige Permutationsmatrix, die

$$A = (a^1, ..., a^l, ..., a^n)$$
 in $AP_1^{\mathsf{T}} = (a^1, ..., a^{l-1}, a^{l+1}, ..., a^n, a^l)$

transformiert. Schreibe (32) als

$$\bar{\boldsymbol{A}} = \boldsymbol{P}_{\boldsymbol{Z}}^{\mathsf{T}} \boldsymbol{L} [\boldsymbol{R} \boldsymbol{P}_{1}^{\mathsf{T}} + \boldsymbol{r} \boldsymbol{e}^{l^{\mathsf{T}}} \boldsymbol{P}_{1}^{\mathsf{T}}] \boldsymbol{P}_{1} \boldsymbol{P}_{S} =: \boldsymbol{P}_{\boldsymbol{Z}}^{\mathsf{T}} \boldsymbol{L}_{2} \boldsymbol{R}_{2} \boldsymbol{P}_{2}$$
(33)

mit $L_2 := L$, der Permutationsmatrix $P_2 := P_1 P_S$ und der oberen Hessenbergmatrix

S3: Führe n - l Elementareliminationsschritte analog zu 6.3.3 zur Transformation der Subdiagonalelemente $(\mathbf{R}_2)_{k+1,k}$ (k = l, ..., n - 1) in 0 gemäß

$$\boldsymbol{L}_{2}\boldsymbol{R}_{2} = \boldsymbol{P}_{3}^{\mathsf{T}}(\boldsymbol{P}_{3}\boldsymbol{L}_{2}\boldsymbol{F}_{3}^{-1}) \left(\boldsymbol{F}_{3}\boldsymbol{R}_{2}\right) = \boldsymbol{P}_{3}^{\mathsf{T}}\boldsymbol{\overline{L}}\boldsymbol{\overline{R}}$$

durch, wobei $\overline{L} := P_3 L_2 F_3^{-1}$ eine untere Einsdreiecksmatrix, $\overline{R} := F_3 R_2$ eine obere Dreiecksmatrix und

$$\boldsymbol{F}_3 := \boldsymbol{\tilde{F}}_{n-1,n} \cdots \boldsymbol{\tilde{F}}_{l,l+1}, \qquad \boldsymbol{P}_3 := \boldsymbol{\tilde{T}}_{n-1,s(n-1)} \cdots \boldsymbol{\tilde{T}}_{l,s(l)}$$

ist. Schreibe (33) in der Form

$$\bar{\boldsymbol{A}} = \boldsymbol{P}_{\boldsymbol{Z}}^{\mathsf{T}} \boldsymbol{L}_{2} \boldsymbol{R}_{2} \boldsymbol{P}_{2} = \boldsymbol{P}_{\boldsymbol{Z}}^{\mathsf{T}} \boldsymbol{P}_{3}^{\mathsf{T}} \bar{\boldsymbol{L}} \bar{\boldsymbol{R}} \boldsymbol{P}_{2} =: \bar{\boldsymbol{P}}_{\boldsymbol{Z}}^{\mathsf{T}} \bar{\boldsymbol{L}} \bar{\boldsymbol{R}} \bar{\boldsymbol{P}}_{S}$$
(34)

mit den Permutationsmatrizen $\overline{P}_Z := P_3 P_Z$ und $\overline{P}_S := P_2 = P_1 P_S$.

Die in S2 von 6.3.4 erforderlichen n-1 Elementareliminationsschritte werden hier durch die Spaltenvertauschungen gemäß P_1 ersetzt, und in S3 reduziert sich die Anzahl der Elementareliminationsschritte von n-1 auf n-l.

(iii) In 6.3.2 und S2 von 6.3.4 genügt es, nur n-2 Schritte zur Erzeugung der Nullen in den Komponenten (r_i) (i = n, ..., 3) durchzuführen, da bereits damit die Hessenbergform von \mathbf{R}_2 erreicht wird.

(iv) Man kann zeigen, daß Algorithmus 6.3.4 im folgenden Sinne numerisch gutartig ist: Wenn $P(A + \delta A) = LR$ gilt, existiert eine Störung $\delta \overline{A}$, die den Beziehungen

$$\overline{P}(\overline{A} + \delta \overline{A}) = \overline{L}\overline{R} \quad \text{und} \quad \|\delta \overline{A}\| \leq \|\delta A\| + \nu K \operatorname{cond} \left(\overline{L}\right) \left(\|A\| + \|\overline{A}\|\right)$$

genügt mit einer nur von *n* abhängigen Konstanten K = K(n).

B. Symmetrische definite Matrizen

Wir setzen A als symmetrisch und positiv definit voraus. Damit die modifizierte Matrix \overline{A} dann auch symmetrisch ist, muß die Rang-1-Korrektur uv^{\intercal} ebenfalls symmetrisch sein, d. h., es muß $u = \alpha v$ und folglich

$$\bar{\boldsymbol{A}} = \boldsymbol{A} + \alpha \boldsymbol{v} \boldsymbol{v}^{\mathsf{T}} \tag{35}$$

mit einem $\alpha \in \mathbf{R}$ gelten. Ohne Einschränkung der Allgemeinheit kann dabei $\boldsymbol{v} \neq \boldsymbol{o}$

und $\alpha = 1$ oder $\alpha = -1$ vorausgesetzt werden, denn die Korrektur läßt sich stets in der Gestalt $\alpha \boldsymbol{v} \boldsymbol{v}^{\mathsf{T}} = \operatorname{sgn}(\alpha) \left(\sqrt{|\alpha|} \, \boldsymbol{v} \right) \left(\sqrt{|\alpha|} \, \boldsymbol{v} \right)^{\mathsf{T}} = + \tilde{\boldsymbol{v}} \tilde{\boldsymbol{v}}^{\mathsf{T}}$ schreiben.

Wir fragen als nächstes, unter welchen Voraussetzungen A positiv definit ist. Unter Verwendung der Cholesky-Faktorisierung

$$A = LL^{\dagger}$$

ergibt sich aus (35)

$$ar{A} = LL^{\intercal} + \alpha v v^{\intercal} = L(I + \alpha r r^{\intercal}) L^{\intercal}, \quad r \text{ Lösung von } Lr = v.$$
 (36)

Die Matrizen \bar{A} und $M := I + \alpha r r^{\intercal}$ sind also kongruent und besitzen folglich dieselben Definitheitseigenschaften, vgl. 1.2.C. Im Fall $\alpha = 1$ gilt

$$oldsymbol{x}^{\intercal} M oldsymbol{x} = oldsymbol{x}^{\intercal} oldsymbol{x} + (oldsymbol{x}^{\intercal} oldsymbol{r})^2 \geqq oldsymbol{x}^{\intercal} oldsymbol{x}$$

für alle $\boldsymbol{x} \in \mathbf{R}^n$, d. h., durch die Addition der positiv definiten Rang-1-Korrektur $\boldsymbol{v}\boldsymbol{v}^{\mathsf{T}}$ verschlechtert sich die positive Definitheit unabhängig von \boldsymbol{v} nicht. Im Fall $\alpha = -1$ gilt wegen $|\boldsymbol{x}^{\mathsf{T}}\boldsymbol{r}| \leq ||\boldsymbol{x}||_2 ||\boldsymbol{r}||_2$ dagegen

$$egin{aligned} \mathbf{x}^{\intercal} \mathbf{M} \mathbf{x} = \mathbf{x}^{\intercal} \mathbf{x} - (\mathbf{x}^{\intercal} \mathbf{r})^2 \geq \|\mathbf{x}\|_2^2 \left(1 - \|\mathbf{r}\|_2^2
ight), \end{aligned}$$

und für $x = \lambda r$ steht das Gleichheitszeichen. In diesem Fall verschlechtert sich also die positive Definitheit. Sie geht verloren, wenn die Bedingung

$$\|\boldsymbol{r}\|_{2}^{2} = \|\boldsymbol{L}^{-1}\boldsymbol{v}\|_{2}^{2} = \boldsymbol{v}^{\mathsf{T}}(\boldsymbol{L}\boldsymbol{L}^{\mathsf{T}})^{-1}\,\boldsymbol{v} = \boldsymbol{v}^{\mathsf{T}}\boldsymbol{A}^{-1}\boldsymbol{v} < 1 \tag{37}$$

nicht mehr erfüllt ist. Bei den Aufdatierungsalgorithmen wird deshalb auch zwischen den Fällen $\alpha = 1$ und $\alpha = -1$ unterschieden. Wegen der Symmetrie der Faktorisierung können dabei die im nichtsymmetrischen Fall verwendeten Elementareliminationsmatrizen $F_{k,k+1}$ durch Givens-Drehungen G_{ij} ersetzt werden, womit sich wegen $||G_{ij}||_2 = 1$ günstigere Stabilitätseigenschaften ergeben.

6.3.6. Aufdatierung einer Cholesky-Faktorisierung bei Addition einer Rang-1-Korrektur.

Aufgabe: Bestimme die Cholesky-Faktorisierung $\bar{A} = \bar{L}\bar{L}^{\intercal}$ von

$$\bar{A} = A + vv^{\mathsf{T}} \tag{38}$$

bei gegebener Cholesky-Faktorisierung $A = LL^{\intercal}$ der positiv definiten Matrix $A \in S^{n.n}$.

Algorithmus:

Setze $L_1 := (L \mid v)$ und schreibe (38) in der Form

$$\bar{A} = LL^{\mathsf{T}} + vv^{\mathsf{T}} = L_1 L_1^{\mathsf{T}}.$$
(39)

Bestimme *n* Givens-Drehungen $G_{j,n+1}$ (j = 1, ..., n) so, daß die Elemente $(L_1^{\mathsf{T}})_{n+1,j} = v_j$ in der (n + 1)-ten Zeile von L_1^{T} sukzessive zu 0 gemacht werden. Dann gilt

$$\boldsymbol{G}\boldsymbol{L}_{1}^{\mathsf{T}} := \boldsymbol{G}_{n,n+1} \cdots \boldsymbol{G}_{1,n+1} \boldsymbol{L}_{1}^{\mathsf{T}} =: \boldsymbol{L}_{2}^{\mathsf{T}} =: \left(\frac{\boldsymbol{L}^{\mathsf{T}}}{\boldsymbol{o}^{\mathsf{T}}}\right)$$
(40)
mit orthogonalem $G := G_{n,n+1} \cdots G_{1,n+1}$ und der oberen Dreiecksmatrix \overline{L}^{\intercal} bzw. L_2^{T} . Schreibe (39) in der Form

$$\bar{\boldsymbol{A}} = \boldsymbol{L}_{1}\boldsymbol{L}_{1}^{\mathsf{T}} = (\boldsymbol{L}_{1}\boldsymbol{G}^{\mathsf{T}}) (\boldsymbol{G}\boldsymbol{L}_{1}^{\mathsf{T}}) = \boldsymbol{L}_{2}\boldsymbol{L}_{2}^{\mathsf{T}} = (\bar{\boldsymbol{L}} \mid \boldsymbol{o}) \left(\frac{\boldsymbol{L}^{\mathsf{T}}}{\boldsymbol{o}^{\mathsf{T}}}\right) = \bar{\boldsymbol{L}}\bar{\boldsymbol{L}}^{\mathsf{T}}.$$

$$Aufwand: \sim n^{2}(2 \text{ opm} + 1 \text{ ops}) + \sim n \text{ opr}$$

$$(41)$$

6.3.7. Bemerkung. (i) In-situ-Realisierung von 6.3.6 auf dem Platz von L und v ist möglich. Die Transformation (40) verläuft dabei für n=3 nach dem folgenden Muster:

(ii) Die Drehungsparameter $\{c, s\}$ von $G_{j,n+1}$ können stets so gewählt werden, daß $(\overline{L})_{jj} > 0$ (j = 1, ..., n) gilt. Dies ist bei der Festlegung $c := l_{jj}/\varrho$, $s := v_j/\varrho$ mit $\varrho := \sqrt{l_{ij}^2 + v_i^2}$ der Fall. Bei Festlegung nach 3.3.1 kann (\overline{L})_{ij} < 0 für gewisse j gelten. In diesem Fall braucht nur das Vorzeichen der j-ten Zeile von \overline{L}^{\intercal} geändert zu werden, denn es gilt $\overline{L}\overline{L}^{\intercal} = (\overline{L}D) (D\overline{L}^{\intercal})$ für jedes $D = \text{diag} (d_i)$ mit $d_i = +1$.

Bei Subtraktion von $\boldsymbol{v}\boldsymbol{v}^{\intercal}$ ist ein etwas aufwendigerer Algorithmus erforderlich.

6.3.8. Aufdatierung einer Cholesky-Faktorisierung bei Subtraktion einer Rang-1-Korrektur.

Aufgabe: Bestimme die Cholesky-Faktorisierung $\bar{A} = \bar{L}\bar{L}^{\intercal}$ der als positiv definit vorausgesetzten Matrix

$$\bar{A} = A - vv^{\intercal} \tag{42}$$

bei gegebener Cholesky-Faktorisierung $A = LL^{\intercal}$ der positiv definiten Matrix $A \in S^{n,n}$.

Algorithmus:

S1: Berechne r als Lösung des Dreieckssystems

$$Lr = v. (43)$$

Falls $\|r\|_2 \ge 1$, stop (\overline{A} nicht positiv definit). Setze

$$\boldsymbol{r}^{1} := \left(\frac{\boldsymbol{r}}{\varrho}\right) \quad \text{mit} \quad \varrho := \sqrt{1 - \|\boldsymbol{r}\|_{2}^{2}};$$
(44)

man beachte

$$\|r^1\|_2^2 = \|r\|_2^2 + \varrho^2 = 1$$
.

Bestimme *n* Givens-Drehungen $G_{i,n+1}$ (i = n, ..., 1) so, daß die Komponenten r_i^1 von r^1 sukzessive zu 0 gemacht werden. Dann gilt $Gr^1 := G_{1,n+1} \cdots$ $G_{n,n+1}r^1 = \left(\frac{o}{\sigma}\right) = \sigma e^{n+1} \text{ mit } \sigma = \pm ||r^1|| = \pm 1. \text{ Setze } L_1 := (L \mid o) \text{ und bilde}$ $GL_1^{\mathsf{T}} = G_{1,n+1} \cdots G_{n,n+1} L_1^{\mathsf{T}} =: L_2^{\mathsf{T}} =: \left(\frac{\overline{L}^{\mathsf{T}}}{\overline{v}^{\mathsf{T}}} \right).$

Dabei ist

$$oldsymbol{v} = oldsymbol{L} oldsymbol{r} = oldsymbol{L}_1 oldsymbol{r}^1 = (oldsymbol{L}_1 G^{\intercal}) \, (Gr^1) = oldsymbol{L}_2 \sigma e^{oldsymbol{n} \, dots \, 1} = (oldsymbol{ar{L}} \mid oldsymbol{ar{v}}) \left(rac{oldsymbol{o}}{\sigma}
ight) = \sigma oldsymbol{ar{v}} \, ,$$

also $\overline{\boldsymbol{v}} = \sigma \boldsymbol{v}$ und folglich

$$\boldsymbol{L}\boldsymbol{L}^{\mathsf{T}} = \boldsymbol{L}_1\boldsymbol{L}_1^{\mathsf{T}} = (\boldsymbol{L}_1\boldsymbol{G}^{\mathsf{T}}) (\boldsymbol{G}\boldsymbol{L}_1^{\mathsf{T}}) = \boldsymbol{L}_2\boldsymbol{L}_2^{\mathsf{T}} = \overline{\boldsymbol{L}}\overline{\boldsymbol{L}}^{\mathsf{T}} + \overline{\boldsymbol{v}}\overline{\boldsymbol{v}}^{\mathsf{T}} = \overline{\boldsymbol{L}}\overline{\boldsymbol{L}}^{\mathsf{T}} + \sigma^2 \boldsymbol{v} \boldsymbol{v}^{\mathsf{T}},$$

d. h., es gilt

$$\overline{L}\overline{L}^{\mathsf{T}} = LL^{\mathsf{T}} - vv^{\mathsf{T}} = \overline{A}.$$
(45)

Autward: $\sim n^2(5/2 \text{ opm} + 3/2 \text{ ops}) + n \text{ opr}$

6.3.9. Bemerkung. (i) In-situ-Realisierung von 6.3.8 auf dem Platz von L und v ist möglich. Für n = 3 erfolgt dabei die Transformation $(\mathbf{r}^1 \mid \mathbf{L}_1^{\mathsf{T}}) \rightarrow (\sigma e^{n+1} \mid \mathbf{L}_2^{\mathsf{T}})$ nach folgendem Muster:



(ii) Die Drehungen $G_{i,n+1}$ können so gewählt werden, daß $(\overline{L})_{ii} > 0$ gilt, bzw. es läßt sich $(\bar{L}_{ii}) > 0$ (i = 1, ..., n) durch Vorzeichenwechsel in den Zeilen von \bar{L}^{\intercal} erzielen, vgl. 6.3.7(ii).

Das Stabilitätsverhalten der Aufdatierungsalgorithmen hängt davon ab, ob $\alpha = 1$ oder $\alpha = -1$ ist.

6.3.10. Bemerkung. Es gelte $A + \delta A = LL^{\dagger}$ mit einer symmetrischen und bezüglich A kleinen Störung δA . Dann gibt es im Fall $\alpha = 1$ eine symmetrische Störung $d\bar{A}$, so daß der gemäß 6.3.6 berechnete Faktor \bar{L} der Beziehung $\bar{A} + d\bar{A} = \bar{L}\bar{L}^{\intercal}$ genügt, und $d\bar{A}$ ist klein in bezug auf \bar{A} . Im Fall $\alpha = 1$ ist der Aufdatierungsalgorithmus 6.3.6 also stets numerisch gutartig. Im Fall $\alpha = -1$ trifft dies für 6.3.8 auch zu, sofern $\|\bar{A}\|$ nicht wesentlich kleiner als $\|A\|$ ist.

Etwa in der nichtlinearen Optimierung treten symmetrische und nach Konstruktion positiv definite Matrizen A, \overline{A} auf, die sich durch eine Rang-2-Modifikation unterscheiden. Zur Aufdatierung der zugehörigen Cholesky-Faktorisierungen siehe Ü 6.3.2 unten. Die BKP-Faktorisierung finiter symmetrischer Matrizen läßt sich ebenfalls aufdatieren, allerdings ist der Algorithmus relativ kompliziert, siehe B 6.3.

Übungsaufgaben

Ü 6.3.1. Es sei $M = I + \alpha r r^{\mathsf{T}}$ mit $r \in \mathbb{R}^n$, $r \neq o$, $\alpha \in \mathbb{R}$, $\alpha \neq 0$. Man zeige, daß M den einfachen Eigenwert $\lambda_1 := 1 + \alpha ||r||_2^2$ mit dem zugehörigen Eigenvektor r sowie den (n - 1)-fachen Eigenwert $\lambda_2 = \cdots = \lambda_n = 1$ besitzt, wobei der zugehörige Eigenraum aus allen zu r orthogonalen Vektoren besteht.

Ü 6.3.2. Es sei $M \in S^{n,n}$ eine Matrix mit rang (M) = 2. Man zeige:

(i) \boldsymbol{M} läßt sich in der Form

$$\boldsymbol{M} = \alpha_1 \boldsymbol{v}^1 \boldsymbol{v}^{1\mathsf{T}} + \alpha_2 \boldsymbol{v}^2 \boldsymbol{v}^{2\mathsf{T}} \quad \text{mit } \alpha_1, \alpha_2 \in \{-1, +1\} \quad \text{und} \quad \boldsymbol{v}^{1\mathsf{T}} \boldsymbol{v}^2 = 0$$
(46)

darstellen.

(ii) Für die Matrix

$$M := uv^{\mathsf{T}} + vu^{\mathsf{T}}, \quad u, v \in \mathsf{R}^n$$
 linear unabhängig, (47)

lautet die Zerlegung (46)

$$\boldsymbol{M} = \boldsymbol{v}^{1}\boldsymbol{v}^{1\mathsf{T}} - \boldsymbol{v}^{2}\boldsymbol{v}^{2\mathsf{T}} \tag{48}$$

 $_{\rm mit}$

$$\boldsymbol{v}^{1} := \varrho \left(\frac{\boldsymbol{u}}{\|\boldsymbol{u}\|_{2}} + \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|_{2}} \right), \quad \boldsymbol{v}^{2} := \varrho \left(\frac{\boldsymbol{u}}{\|\boldsymbol{u}\|_{2}} - \frac{\boldsymbol{v}}{\|\boldsymbol{v}\|_{2}} \right), \quad \varrho := \left(\frac{\|\boldsymbol{u}\| \|\boldsymbol{v}\|}{2} \right)^{1/2}.$$
(49)

(iii) Sind A und $\overline{A} := A + M$ beide symmetrisch und positiv definit, so kann die Cholesky-Faktorisierung von \overline{A} aus der von A durch zweimalige Anwendung der Rang-1-Aufdatierungsalgorithmen auf die Modifikationen

$$\boldsymbol{A} \to \boldsymbol{A}_1 := \boldsymbol{A} + \alpha_1 \boldsymbol{v}^1 \boldsymbol{v}^{1\mathsf{T}} \to \bar{\boldsymbol{A}} := \boldsymbol{A}_1 + \alpha_2 \boldsymbol{v}^2 \boldsymbol{v}^{2\mathsf{T}}$$

$$\tag{50}$$

berechnet werden. Im Fall unterschiedlichen Vorzeichens der α_i muß dabei $\alpha_1 = +1$ gewählt werden, damit die intermediäre Matrix A_1 positiv definit ist und der zugehörige Cholesky-Faktor L_1 existiert.

6.4. Dreiecksfaktorisierung schwach besetzter Matrizen

Eine Reihe von praktischen Problemen führt auf lineare Gleichungssysteme hoher Dimension mit schwach besetzten Koeffizientenmatrizen, vgl. 2.1.A1. Solche Matrizen sind durch die Eigenschaft

$$\eta(A) \ll n^2$$

charakterisiert, wobei $\eta(A)$ die Anzahl der Nichtnullelemente — im folgenden mit NNE abgekürzt — von A bezeichnet. Für realistische Anwendungen liegt n in der Größenordnung von 10² bis 10⁴ und größer, so daß es notwendig wird, die schwache Besetztheit sowohl im Hinblick auf den benötigten Speicherplatz — bei voll besetzten Matrizen $\sim K_1 n^2 S$ — als auch bezüglich des Rechenaufwandes — bei voll besetzten Matrizen $\sim K_2 n^3$ opms — auszunutzen. In vielen Fällen ist die mittlere Anzahl der NNE pro Zeile unabhängig von n konstant gleich m, so daß $\eta(A) = mn$ mit n wächst, aber der Besetztheitsgrad $\eta(A)/n^2 = m/n$ mit wachsendem n immer kleiner wird.

Die Vorgehensweise bei der Dreiecksfaktorisierung schwach besetzter Matrizen hängt wesentlich vom Besetztheitsmuster von A, d. h. von der Anordnung der NNE ab.

A. Bandmatrizen

Eine Matrix $A = (a_{ij}) \in \mathbf{R}^{n,n}$ mit

$$a_{ij} = 0$$
 für $j < i - p$ und $j > i + q$ (1)

heißt Bandmatrix des Typs $\{p, q\}$; die Zahlen p bzw. q heißen untere bzw. obere Bandbreite, und l := p + q + 1 wird Bandbreite schlechthin genannt, vgl. 1.1.C. Die durch die Diagonale sowie die benachbarten p unteren und q oberen Nebendiagonalen gebildeten Plätze werden das Band von A genannt. Ein Beispiel mit n = 8, p = 2, q = 3 stellt die Matrix

$$A = \begin{pmatrix} 11 & \frac{q}{12 & 13 & 14} & & \\ 21 & 22 & 23 & 24 & 25 & \\ 31 & 32 & 33 & 34 & 35 & 36 & \\ & 42 & 43 & 44 & 45 & 46 & 47 & \\ & & 53 & 54 & 55 & 56 & 57 & 58 \\ & & & 64 & 65 & 66 & 67 & 68 \\ & & & & & 75 & 76 & 77 & 78 \\ & & & & & & 86 & 87 & 88 \end{pmatrix}$$

dar. Die signifikanten Elemente von A können zeilenweise in einem rechteckigen Feld der Dimension (n, l) oder spaltenweise in einem Feld der Dimension (l, n)gespeichert werden; für die Matrix (1) führt dies auf

$$\begin{pmatrix} 11 & 12 & 13 & 14 \\ 21 & 22 & 23 & 24 & 25 \\ 31 & 32 & 33 & 34 & 35 & 36 \\ 42 & 43 & 44 & 45 & 46 & 47 \\ 53 & 54 & 55 & 56 & 57 & 58 \\ 64 & 65 & 66 & 67 & 68 \\ 75 & 76 & 77 & 78 \\ 86 & 87 & 88 \end{pmatrix}$$
 bzw.
$$\begin{pmatrix} 14 & 25 & 36 & 47 & 58 \\ 13 & 24 & 35 & 46 & 57 & 68 \\ 12 & 23 & 34 & 45 & 56 & 67 & 78 \\ 11 & 22 & 33 & 44 & 55 & 66 & 77 & 88 \\ 21 & 32 & 43 & 54 & 65 & 76 & 87 \\ 31 & 42 & 53 & 64 & 75 & 86 \end{pmatrix} .$$
(3)

Die p(p + 1)/2 + q(q + 1)/2 Plätze in den Ecken werden zum Speichern des Bandes nicht benötigt; für genügend großes n spielen sie in bezug auf den Gesamtspeicherplatz keine Rolle.

Wir betrachten im folgenden den ersten Schritt des Gaußschen Algorithmus, in der Terminologie von 5.1 also den Übergang von $A = A^{(1)}$ zu $A^{(2)}$, und unterscheiden dabei zwei Fälle:

Fall 1: Keine Pivotisierung, d. h., $a_{11} \neq 0$ wird als Pivot gewählt. Für die Matrix (2) ergibt sich dabei das Muster



Beim Übergang zu $A^{(2)}$ werden nur die eingerahmten p * q Elemente neu berechnet, insbesondere bleiben die Nullen außerhalb des Bandes erhalten. Das Verfahren kann daher auf dem Platz des Bandes von A ausgeführt werden. Die zu 0 gemachten Elemente werden mit den Eliminationskoeffizienten l_{ik} überspeichert.

Fall 2: Spalten pivotisierung, d. h. a_{s1} mit $|a_{s1}| = \max \{|a_{i1}|: 1 \leq i \leq 1 + p\}$ wird als Pivot gewählt, und die Zeilen 1 und s werden vertauscht.

Für die Matrix (2) ergibt sich hier im ungünstigsten Fall s = p + 1 das Muster



Beim Übergang zu $A^{(2)}$ werden die ersten 1 + p Zeilen bis zur Spalte 1 + q + paufgefüllt. Im nächsten Schritt kann derselbe Effekt wieder auftreten, wobei allerdings der entstehende Dreiecksfaktor R die obere Bandbreite q + p nicht überschreitet.

Zusammenfassend erhalten wir das folgende Resultat:

6.4.1. Aussage. Es sei $A \in \mathbb{R}^{n,n}$ eine reguläre Bandmatrix des Typs $\{p, q\}$. Dann gilt:

- (i) Wenn die Dreiecksfaktorisierung nach dem Gaußschen Algorithmus ohne Pivotisierung durchführbar ist (etwa für diagonaldominantes oder symmetrisches und positiv definites A, vgl. 5.2.D), sind die Dreiecksfaktoren L und RBandmatrizen des Typs $\{p, 0\}$ und $\{0, q\}$. Die Faktorisierung kann auf dem Platz des Bandes von A ausgeführt werden; der Aufwand ist $\sim n(p \text{ opm}$ + p * q opms). Im symmetrischen und positiv definiten Fall ist die Faktorisierung $A = LDL^{\intercal}$ auf dem unteren Teil des Bandes von A durchführbar, und der Aufwand reduziert sich auf $\sim n(p \text{ opm} + p(p + 1)/2 \text{ opms})$. Bei der Cholesky-Faktorisierung $A = LL^{\intercal}$ kommen n opr hinzu.
- (ii) Wenn die Dreiecksfaktorisierung nach dem Gaußschen Algorithmus mit Spaltenpivotisierung durchgeführt werden muß, können maximal p weitere obere Nebendiagonalen aufgefüllt werden, d. h., \mathbf{R} ist eine Bandmatrix des Typs $\{0, q + p\}$. Die Eliminationskoeffizienten \hat{l}_{ik} können auf dem Platz der zu 0 gemachten Elemente im unteren Teil des Bandes gespeichert werden. Das Verfahren ist daher auf dem Platz einer $\{p, q + p\}$ -Bandmatrix ausführbar, und der Aufwand ist maximal $\sim n(p \text{ opm} + p(q + p) \text{ opms}).$

6.4.2. Bemerkung. (i) Bei Rechnung ohne Pivotisierung kann die Dreiecksfaktorisierung sowohl nach den (n-1)-stufigen Grundformen des Gaußschen Algorithmus aus 5.1 oder 6.1 als auch nach den direkten Verfahren aus 6.2 berechnet werden. Bei Rechnung mit Spaltenpivotisierung treten wie beim Vorgehen ohne Pivotisierung für jedes $k \in \{1, ..., n-1\}$ höchstens p nichttriviale Eliminationskoeffizienten \hat{l}_{ik} $(i = k + 1, ..., \min (k + p, n))$ auf, die auf dem Platz des unteren Bandes gespeichert werden können. Vertauschung der \hat{l}_{ik} gemäß (5.1.18) analog zu den Zeilen der $A^{(j)}$ (j > k) führt auf die Dreiecksmatrix $L = (l_{ik})$ der Faktorisierung

$$PA = LR. (4)$$

Dabei ist L jedoch i. allg. keine $\{p, 0\}$ -Bandmatrix, obwohl in jeder Spalte höchstens p nichtverschwindende l_{ik} außerhalb der Diagonalen vorkommen. Die Darstellung (4) mit dem explizit gebildeten Dreiecksfaktor L ist daher für Bandmatrizen im Fall der Spaltenpivotisierung nicht zweckmäßig. Die angepaßte Darstellung – vgl. auch 5.2.6 (i) — ist die implizite Dreiecksfaktorisierung

$$\boldsymbol{A} = \boldsymbol{\hat{L}}\boldsymbol{R} \tag{5}$$

gemäß (5.1.16) mit

$$\hat{L} := T_{1,s(1)} L_1(\hat{l}^1) \cdots T_{n-1,s(n-1)} L_{n-1}(\hat{l}^{n-1})$$
(6)

und

$$\hat{l}^{k} := (0, ..., 0, \hat{l}_{k+1,k}, ..., \hat{l}_{\min(k+p,n),k}, 0, ..., 0)^{\mathsf{T}}.$$

Die Matrix \hat{L} wird durch die Zahlen $\{\hat{l}_{ik}, s(k): k = 1, ..., n - 1, i = k + 1, ..., min <math>(k + p, n)$ } repräsentiert. Für Details der Implementierung sei auf die Spezialliteratur verwiesen, siehe B 6.4. (ii) Das Gleichungssystem Ax = b kann unter Verwendung der ohne Pivotisierung berechneten Dreiecksfaktorisierung A = LR durch Lösung von Lc = b und Rx = c mit $\sim n(1 \text{ opm} + (p + q) \text{ opms})$ gelöst werden. Im Fall einer symmetrischen und definiten Matrix werden $\sim n(1 \text{ opm} + 2p \text{ opms})$ bei Verwendung der LDL^{\intercal} -Faktorisierung und $\sim n(2 \text{ opm} + 2p \text{ opms})$ bei Verwendung der Cholesky-Faktorisierung benötigt. Wenn die Faktorisierung (5) mit Spaltenpivotisierung berechnet wurde, ist c als Lösung von $\hat{L}c = b$ gemäß

$$c = L_{n-1}(-\hat{l}^{n-1}) T_{n-1,s(n-1)} \cdots L_1(-\hat{l}^1) T_{1,s(1)}b$$

zu bestimmen, man beachte (6) und 5.2.6(i). Der Aufwand ist maximal $\sim n(1 \text{ opm} + (2p + q) \text{ opms}).$

(iii) Die beschriebenen Bandalgorithmen erzeugen dieselben Ausgangsdaten wie die entsprechenden Algorithmen für vollbesetztes A; sie sind also insbesondere numerisch gutartig, wobei die Fehlerkumulationskonstanten F = F(n, p, q) kleinere Werte als im vollbesetzten Fall haben. Als Algorithmen zur Lösung von Ax = bsind sie fast optimal in dem Sinn, daß der Quotient $\varrho(n) := (\text{Anzahl der Rechen$ $operationen})/(\text{Anzahl der Eingangsdatenelemente}) unabhängig von <math>n$ durch eine Konstante beschränkt ist, vgl. B 2.10. Bei voll besetzter Matrix gilt für den Gaußschen Algorithmus dagegen $\varrho(n) \sim Kn$.

(iv) Bei Diagonal- oder vollständiger Pivotisierung wird die Bandstruktur i. allg. zerstört. Dasselbe trifft für die bei der BKP-Faktorisierung indefiniter symmetrischer Matrizen verwendete Pivotisierungsstrategie zu, so daß aus der Symmetrie einer Bandmatrix bei fehlender Definitheit i. allg. kein Vorteil gezogen werden kann.

(v) Im Band auftretende Nullen werden bei den diskutierten Bandtechniken nicht berücksichtigt. Falls jedoch eine sog. *Blockbandstruktur* mit Blöcken A_{ij} von Bandgestalt vorliegt, falls A also etwa vom Typ

ist, kann analog zum Vorgehen in 6.1.B zur *Blockelimination* übergegangen werden. Die Inversen der Pivotblöcke sollen dabei nicht explizit gebildet, sondern durch die zugehörigen Faktorisierungen repräsentiert werden. \Box

Den einfachsten Sonderfall einer Bandmatrix stellt eine Tridiagonalmatrix, d. h. eine $\{1, 1\}$ -Bandmatrix

$$\boldsymbol{A} = \begin{pmatrix} d_{1} & e_{1} & & \\ c_{2} & d_{2} & e_{2} & & \\ c_{3} & d_{3} & e_{3} & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & \\ & & c_{n-1} & d_{n-1} & e_{n-1} \\ & & & c_{n} & d_{n} \end{pmatrix}$$
(7)

dar, die zweckmäßigerweise durch die drei Vektoren $\mathbf{c} = (c_i)$, $\mathbf{d} = (d_i)$, $\mathbf{e} = (e_i) \in \mathbb{R}^n$ dargestellt wird. Der erste Schritt des Gaußschen Algorithmus führt je nach Wahl von d_1 bzw. c_2 als Pivot auf

$$A^{(2)} = \begin{pmatrix} d_1 & e_1 & & & \\ 0 & d_2 - \cdot l_{21}e_1 & e_2 & & \\ c_3 & & d_3 & e_3 & \\ & \ddots & & \ddots & \ddots & \end{pmatrix} \quad \text{bzw.} \quad A^{(2)} = \begin{pmatrix} c_2 & d_2 & & e_2 & & \\ 0 & e_1 - \hat{l}_{21}d_2 & -\hat{l}_{21}e_2 & & \\ c_3 & & d_3 & e_3 & \\ & \ddots & & \ddots & \ddots & \end{pmatrix};$$

der Eliminationskoeffizient ist

 $l_{21}:=c_2/d_1$ bzw. $\hat{l}_{21}:=d_1/c_2$.

Die entstehende (n-1)-dimensionale Restmatrix $M^{(2)}$ ist wieder tridiagonal, so daß analog zum ersten Schritt fortgefahren werden kann. Bei Rechnung ohne Pivotisierung ergibt sich das folgende Verfahren.

6.4.3. Dreiecksfaktorisierung einer Tridiagonalmatrix A ohne Pivotisierung und Lösung von Ax = b.

Aufgabe: Für die durch (7) gegebene reguläre Tridiagonalmatrix $A \in \mathbb{R}^{n,n}$ ist die Faktorisierung A = LR nach dem als durchführbar vorausgesetzte Gaußschen Algorithmus zu berechnen; die Matrix A ist mit den signifikanten Elementen von L und R zu überspeichern. Unter Verwendung von L, R ist x als Lösung von Ax = b zu berechnen und auf dem Platz von b zu speichern.

Algorithmus:

S1 (Dreiecksfaktorisierung A = LR): for k := 2(1)n do $[c_k := c_k/d_{k-1}, d_k := d_k - c_k * e_{k-1}]$ S2 (Lösung von LRx = b): S2.1 ($b := L^{-1}b$): for k := 2(1)n do $b_k := b_k - c_k * b_{k-1}$ S2.2 ($b := R^{-1}b$): $b_n := b_n/d_n$ for k := n - 1(-1)1 do $b_k := (b_k - e_k * b_{k+1})/d_k$ Aufwand: $\sim n(1 \text{ opm} + 1 \text{ opms})$ für $L, R, \sim n(1 \text{ opm} + 2 \text{ opms})$ für x

6.4.4. Bemerkung. (i) Wenn A symmetrisch und positiv definit ist, gilt $e_{k-1} = c_k$, und die Laufanweisung in S1 ist durch

for
$$k := 2(1)n$$
 do $[z := c_k, c_k := c_k/d_{k-1}, d_k := d_k - c_k * z]$

zu ersetzen. Damit ergibt sich die Faktorisierung $A = LDL^{\intercal}$.

(ii) Bei Rechnung mit Spaltenpivotisierung werden zusätzlich $\sim n S$ für die zweite Nebendiagonale von \mathbf{R} benötigt, außerdem sind die Indizes $\{s(1), \ldots, s(n-1)\}$ zu speichern. Wenn nur ein einzelnes System $A\mathbf{x} = \mathbf{b}$ zu lösen ist und S1 und S2.1 wie in 5.1.1 gemeinsam in einer Laufanweisung realisiert werden, brauchen die \hat{l}_{ik} und s(k) nicht aufgehoben zu werden. In diesem Fall kann die Lösung auf dem Platz der rechten Seite \mathbf{b} und der drei Diagonalen von \mathbf{A} berechnet werden, indem die neu entstehende zweite obere Nebendiagonale auf dem Platz von \mathbf{c} gespeichert wird. \square

15 Schwetlick, Numerische Algebra

B. Schwach besetzte Matrizen beliebiger Struktur

Algorithmen zur Dreiecksfaktorisierung schwach besetzter Matrizen beliebiger Struktur gehen davon aus, daß die NNE beliebig über die Matrix verteilt sein können; evtl. vorhandene regelmäßige Strukturen werden ignoriert.

In manchen Fällen ist es möglich, die Matrix A durch Zeilen- und Spaltenvertauschungen

$$\bar{\boldsymbol{A}} = \boldsymbol{P}_{\boldsymbol{Z}} \boldsymbol{A} \boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} \tag{8}$$

für nichtsymmetrisches A bzw.

$$\bar{\mathbf{A}} = \mathbf{P} \mathbf{A} \mathbf{P}^{\mathsf{T}} \tag{9}$$

für symmetrisches und definites A so umzuordnen, daß \overline{A} von Bandgestalt mit minimaler Bandbreite ist. Die zugehörigen Algorithmen zur Bestimmung der Permutationsmatrizen P_z , P_s bzw. P sind graphentheoretischer Natur und können im Rahmen einer Einführung nicht diskutiert werden, siehe B 6.5 für Hinweise.

Für symmetrisches definites A ist durch Übergang von Band- zu Einhüllenden-Techniken eine weitere Senkung des Speicherbedarfs und Rechenaufwandes möglich. Die *Einhüllende* der Matrix A besteht dabei aus der Diagonalen und allen Positionen (i, j), für die ein l mit $a_{il} \neq 0$ $(1 \leq l \leq j < i)$ oder $a_{lj} \neq 0$ $(1 \leq l \leq i < j)$ existiert. Im nachfolgenden Beispiel sind die zur Einhüllenden gehörenden Positionen umrandet worden.

$$\boldsymbol{A} = \begin{bmatrix} \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} & \boldsymbol{\times} \\ \boldsymbol{\times} & \boldsymbol{\times}$$

Aus der zeilenweisen Version 6.2.6 des Cholesky-Verfahrens liest man sofort das folgende Resultat ab:

6.4.5. Aussage: Die Matrix $A \in S^{n,n}$ sei positiv definit. Dann ist die Einhüllende von A gleich der Einhüllenden von $(L + L^{\intercal})$, wobei $A = LL^{\intercal}$ die Cholesky-Faktorisierung von A bezeichnet.

Dies besagt: Der Cholesky-Faktor L kann auf dem unteren Teil der Einhüllenden von A berechnet werden. Wenn die Anzahl der zur Einhüllenden von A gehörigen Plätze Profil von A — in Zeichen: p(A) — genannt wird, ist (p(A) + n)/2 gerade die Anzahl der für die Faktorisierung benötigten Speicherplätze. Es ist daher sinnvoll, durch geeignete simultane Zeilen- und Spaltenvertauschungen gemäß (9) das Profil von A zu minimieren, siehe B 6.5 für Hinweise.

Wenn Band- oder Einhüllenden-Techniken nicht angewendet werden können oder nicht zum Erfolg führen, muß zu universellen Verfahren für schwach besetzte Matrizen übergegangen werden. Grundlage ist wieder der Gaußsche Algorithmus, von dem bereits k - 1 Schritte durchgeführt sein mögen. Die zu Beginn des k-ten Schrittes vorliegende Restmatrix ist dann

$$m{M}^{(k)} = egin{pmatrix} a_{kk}^{(k)} \dots a_{ks}^{(k)} \dots a_{kn}^{(k)} \ dots & dots & dots \ dots & dots \ dots & dots \ dots & dots \ d$$

Wenn das Element $a_{\hat{s}\hat{s}}^{(k)} \neq 0$ als Pivot gewählt wird, gehen die $a_{ij}^{(k)}$ in

$$a_{ij}^{(k+1)} := a_{ij}^{(k)} - \frac{a_{i\hat{s}}^{(k)} * a_{sj}^{(k)}}{a_{s\hat{s}}^{(k)}} \quad (i, j = k, ..., n; i \neq s, j \neq \hat{s})$$
(10)

über, sofern auf die physische Vertauschung der Zeilen k, s und Spalten k, \hat{s} verzichtet wird. Ein Nullelement $a_{ij}^{(k)}$ kann also nur dann in ein Nichtnullelement $a_{ij}^{(k+1)}$ übergehen, wenn sowohl das zugehörige Pivotspaltenelement $a_{ij}^{(k)}$ als auch das Pivotzeilenelement $a_{sj}^{(k)}$ von 0 verschieden sind. Für schwach besetztes A wird das nur selten der Fall sein, so daß die meisten Nullen aus $M^{(k)}$ nach $M^{(k+1)}$ übernommen werden. Zur Illustration möge das folgende Beispiel dienen, bei dem die Positionen, in denen NNE entstehen können, eingerahmt worden sind:



Die Größe

$$\eta(\boldsymbol{M}^{(\boldsymbol{k}+1)}) - \eta(\boldsymbol{M}^{(\boldsymbol{k})}) \tag{11}$$

heißt lokale Auffüllung im k-ten Schritt (engl. "local fill in"). Wenn $\sigma_i^{(k)}$ bzw. $\hat{\sigma}_j^{(k)}$ die Anzahl der NNE in der Zeile $(a_{ik}^{(k)}, \ldots, a_{in}^{(k)})$ bzw. Spalte $(a_{kj}^{(k)}, \ldots, a_{nj}^{(k)})^{\mathsf{T}}$ von $M^{(k)}$ bezeichnet, kann die lokale Auffüllung wie folgt abgeschätzt werden:

6.4.6. Aussage. Bei der Wahl von $a_{\hat{s}\hat{s}}^{(k)} \neq 0$ $(k \leq s, \hat{s} \leq n)$ als Pivot gilt

$$\eta(M^{(k+1)}) - \eta(M^{(k)}) \leq (\sigma_s^{(k)} - 1) * (\hat{\sigma}_s^{(k)} - 1) = : \mu_{ss}^{(k)}$$

Die Zahlen $\mu_{s\hat{s}}^{(k)}$ heißen *Markowitz-Kosten* der Elemente $a_{s\hat{s}}^{(k)} \neq 0$. Man beachte, daß durch (11) die gesamte im k-ten Schritt auftretende Auffüllung erfaßt wird,

15*

denn bei der Berechnung der \hat{l}_{ik} und der Übernahme der Pivotzeile als k-te Zeile von R entstehen keine NNE.

Die bisherigen Überlegungen zeigen, daß ein universelles Verfahren zur Dreiecksfaktorisierung schwach besetzter Matrizen zwei Forderungen erfüllen muß:

- Auswahl einer Pivotreihenfolge $\{s(k), \hat{s}(k)\}$ (k = 1, ..., n - 1) bzw. äquivalent Festlegung von Permutationsmatrizen P_Z und P_S derart, daß die Dreiecksfaktorisierung von A stabil durchführbar ist und die durch

$$\eta(\boldsymbol{L}+\boldsymbol{R})-\eta(\boldsymbol{A})$$

gegebene Auffüllung während des Eliminationsprozesses möglichst klein ist.

 Auswahl einer Kompaktspeichertechnik, bei welcher nur die NNE von A und die im Laufe der Elimination durch Auffüllung entstehenden NNE einschließlich zugehöriger Indexinformationen zum Auffinden bzw. Einfügen von Elementen gespeichert werden.

Die in der ersten Forderung genannte stabile Durchführbarkeit spielt je nach Eigenschaften der Matrix A eine unterschiedliche Rolle: Für symmetrisches und positiv definites A werden nur symmetrische Zeilen- und Spaltenvertauschungen zugelassen, d. h., die Pivots müssen auf der Diagonalen von $M^{(k)}$ gewählt werden, und es gilt $s(k) = \hat{s}(k)$ bzw. $P_Z = P_S = P$, vgl. (9). Aus 5.2 und 6.1 ist bekannt, daß die Cholesky-Faktorisierung für jede solche Pivotreihenfolge ein numerisch gutartiger Prozeß ist unabhängig von den konkreten Zahlenwerten der NNE. Die Pivotreihenfolge kann also allein auf Grund des Besetztheitsmusters von A festgelegt werden. Demgemäß bestehen Implementierungen für symmetrische und definite schwach besetzte Systeme in der Regel aus drei Teilprogrammen:

- P1: Festlegung der Pivotreihenfolge auf Grund des Besetztheitsmusters von A
- P2: Numerische Berechnung des Cholesky-Faktors unter Verwendung der in P1 ermittelten Pivotreihenfolge
- P3: Lösung von Ax = b unter Verwendung des in P2 berechneten Cholesky-Faktors von A.

Der Zeitaufwand für die in P1 auszuführenden nichtnumerischen Operationen ist i. allg. wesentlich größer als der für P2. Da in der Praxis meist mehrere, unter Umständen sogar sehr viele Systeme derselben Besetztheitsstruktur, aber mit unterschiedlichen Zahlenwerten zu lösen sind, braucht der aufwendige erste Schritt P1 nur einmal ausgeführt zu werden.

Im nichtsymmetrischen Fall hängt das Stabilitätsverhalten im Unterschied zum symmetrischen und definiten Fall auch von den Zahlenwerten der a_{ij} ab, so daß die Pivotreihenfolge nicht ausschließlich auf Grund des Besetztheitsmusters festgelegt werden kann. Wenn das Pivot in bezug auf die restlichen NNE der Pivotspalte relativ klein ist, werden die l_{ik} und i. allg. auch die $a_{ij}^{(k+1)}$ sehr groß, und die Gutartigkeit bzw. Stabilität geht verloren, vgl. 5.2 und 5.3. Die Freiheit in der Wahl des Pivots $a_{ss}^{(k)}$ muß also auf Kosten einer möglicherweise größeren Auffüllung eingeschränkt werden, etwa durch die Bedingung

$$|a_{\hat{s}\hat{s}}^{(k)}| \ge \alpha \max\{|a_{\hat{s}\hat{s}}^{(k)}|: k \le i \le n\} \quad \text{mit} \quad 0 < \alpha \le 1.$$

$$(12)$$

Aus (12) folgt

$$|\hat{l}_{ik}| \leq 1/lpha$$
,

d. h., das Elementwachstum ist gemäß

 $\max \{ |a_{ij}^{(k+1)}| : k+1 \leq i, j \leq n \} \leq (1+1/\alpha) \max \{ |a_{ij}^{(k)}| : k \leq i, j \leq n \}$ (13)

beschränkt. Für $\alpha = 1$ liegt gerade der Fall der Spaltenpivotisierung in der Spalte \hat{s} vor, allerdings kann der Spaltenindex $\hat{s} \in \{k, ..., n\}$ noch frei gewählt werden. Durch Vorgabe von α kann der Nutzer einen Kompromiß zwischen den beiden gegensätzlichen Forderungen "numerische Gutartigkeit" (α nahe bei 1) und "geringe Auffüllung" (α nahe bei 0) zu erzielen versuchen.

Bei Verwendung einer einmal bestimmten Pivotreihenfolge für eine weitere Matrix desselben Besetztheitsmusters ist die Bedingung (12) nicht notwendig erfüllt. Es macht sich dann erforderlich, das Wachstum der $a_{ij}^{(k)}$ im Laufe der Elimination zu kontrollieren und bei zu starkem Anwachsen gegebenenfalls eine neue Pivotreihenfolge festzulegen. Eine weitere Möglichkeit besteht darin, den Genauigkeitsverlust bei der Faktorisierung durch einige Schritte der iterativen Verbesserung entsprechend 5.4 auszugleichen.

Wir bemerken abschließend, daß für schwach besetzte Gleichungssysteme, die bei der Diskretisierung von partiellen Differentialgleichungen etwa nach der Methode der finiten Elemente bzw. nach dem Differenzenverfahren entstehen, eine Reihe von speziellen Lösungsverfahren bekannt ist. Diese Verfahren nutzen die spezielle Gestalt der Koeffizientenmatrix zielgerichtet aus und sind daher für gewisse Problemklassen günstiger als universelle Lösungsverfahren, siehe B 6.6.

Übungsaufgaben

Ü 6.4.1. Man schreibe ein zu 6.4.3 analoges Programm für den Fall einer symmetrischen und positiv definiten Koeffizientenmatrix unter Verwendung der Faktorisierung $A = LDL^{T}$ bzw. $A = LL^{T}$.

Ü 6.4.2. Man gebe ein zu 6.4.3 analoges Programm zur Faktorisierung einer nichtsymmetrischen Tridiagonalmatrix mit Spaltenpivotisierung und zur Lösung von Ax = b an.

6.5. Berechnung und Aufdatierung der inversen Matrix

In den Anwendungen wird die Inverse A^{-1} einer regulären Matrix $A \in \mathbb{R}^{n,n}$ in den seltensten Fällen explizit benötigt, vgl. die Ausführungen am Ende dieses Abschnitts. Wenn die Berechnung von A^{-1} doch erforderlich ist, gehen wir wie bisher von der Dreiecksfaktorisierung

$$PA = LR \tag{1}$$

von A aus, die nach einer der Varianten des Gaußschen Algorithmus berechnet sein möge. Aus (1) folgt

$$\boldsymbol{A}^{-1} = \boldsymbol{R}^{-1} \boldsymbol{L}^{-1} \boldsymbol{P}. \tag{2}$$

Zur Auswertung der Darstellung (2) bestimmen wir in einem ersten Schritt die Inverse $S := R^{-1}$ aus der Matrixgleichung

$$SR = I$$
 (3)

analog zum Vorgehen beim verketteten Gaußschen Algorithmus in 6.2. Da S wie R eine obere Dreiecksmatrix ist, kann (3) elementweise in der Form

$$\sum_{p=i}^{j} s_{ip} r_{pj} = \begin{cases} 1 & \text{für} \cdot i = j, \\ 0 & \text{für} \cdot i \neq j \end{cases} \quad (i, j = 1, ..., n),$$

geschrieben werden. Für i > j ist dies trivial erfüllt, da die linksstehende Summe leer ist. Für i = j ergibt sich $s_{ij}r_{ij} = 1$, also

$$s_{jj} := 1/r_{jj}$$
 $(j = 1, ..., n).$ (4)

Für i < j folgt $\sum_{p=i}^{j-1} s_{ip} r_{pj} + s_{ij} r_{jj} = 0$, d. h. unter Beachtung von (4)

$$s_{ij} := -\left(\sum_{p=i}^{j-1} s_{ip} r_{pj}\right) s_{jj}.$$
(5)

Zur Bestimmung von s_{ij} werden die links von $\{i, j\}$ stehenden Elemente von S und die unterhalb von $\{i, j\}$ stehenden Elemente von R benötigt. Eine Berechnung auf dem Platz von R ist daher möglich, wenn S spaltenweise von links nach rechts und in den Spalten von oben nach unten bestimmt wird.

Im zweiten Schritt wird $T := SL^{-1}$ gebildet. Aus der äquivalenten Gleichung

$$TL = S \tag{6}$$

ergibt sich unter Beachtung der Dreiecksform von L

$$\sum_{p=j}^n t_{ip}l_{pj} = s_{ij},$$

also wegen $l_{ii} = 1$

$$t_{ij} := s_{ij} - \sum_{p=j+1}^{n} t_{ip} l_{pj} \qquad (i, j = 1, ..., n);$$
⁽⁷⁾

man beachte $s_{ij} = 0$ für i > j. Auf der rechten Seite treten die rechts von der Position (i, j) stehenden Elemente von T und die Eliminationskoeffizienten der *j*-ten Spalte auf. Eine spaltenweise in-situ-Realisierung ist daher möglich, wenn die l_{pj} (p = j + 1, ..., n) in einem Hilfsfeld der Länge n ausgelagert werden, damit die entsprechenden t_{pj} in die Positionen (p, j) gebracht werden können.

Im dritten Schritt entsteht $A^{-1} = TP$ durch Spaltenvertauschung aus T. Zusammenfassend erhalten wir den folgenden Inversionsalgorithmus:

6.5.1. Berechnung von A^{-1} unter Verwendung der LR-Faktorisierung von A.

Aufgabe: Für gegebene und gemäß 5.2.3 auf dem Platz von A gespeicherte Dreiecksfaktorisierung PA = LR mit $L, R \in \Re^{n,n}, R$ regulär, ist die Inverse A^{-1} zu berechnen und auf dem Platz von A zu speichern.

6.5.2. Bemerkung. (i) Wenn die $\sim n^3/3$ opms für die *LR*-Faktorisierung mit gezählt werden, erfordert die Matrixinversion nach 6.5.1 insgesamt $\sim n^3$ opms. Eine *n*-bzw. (n - 1)-stufige Version von S1 bzw. S2 ist möglich, siehe Ü 6.5.1.

(ii) Für symmetrisches A kann von den billigeren symmetrischen Faktorisierungen aus 6.1 ausgegangen werden. Im positiv definiten Fall würde sich etwa mit der Cholesky-Faktorisierung $A = LL^{\intercal}$ das folgende Vorgehen empfehlen:

S1: Berechne
$$\mathbf{S} := (\mathbf{L}^{\intercal})^{-1}$$
 mit $\sim n^3/6$ opms

S2: Bilde $A^{-1} := SS^{\mathsf{T}}$ mit $\sim n^3/6$ opms

Der Aufwand reduziert sich dabei auf die Hälfte. Für indefinites A kann unter Verwendung der BKP-Faktorisierung analog vorgegangen werden. Man beachte, daß A^{-1} symmetrisch ist, also nur ein Dreieck berechnet zu werden braucht.

(iii) Die Rundungsfehleranalyse von 6.5.1 unter Einbeziehung der *LR*-Faktorisierung zeigt, daß das Verfahren für nicht zu schlecht konditioniertes $A \in \Re^{n,n}$ durchführbar ist. Für die berechnete Inverse $X \approx A^{-1}$ gilt dabei: Zu jedem $i \in \{1, ..., n\}$ gibt es eine Störung $\delta A^{(i)}$, so daß die *i*-te Zeile $e^{i\tau}X$ von X der Beziehung

$$e^{i\mathsf{T}}X(A + \delta A^{(i)}) = e^{i\mathsf{T}} \quad \text{mit} \quad \|\delta A^{(i)}\|_2 \leq \nu F \|A\|_2 \tag{8}$$

genügt, wobei F = F(n) eine schwach mit *n* wachsende, von der Matrizenklasse und dem Faktorisierungsverfahren abhängende Kumulationskonstante ist. Das Verjahren 6.5.1 ist also numerisch gutartig in bezug auf jede Zeile von X. Dagegen gibt es i. allg. keine kleine Störung δA , so daß $X = (A + \delta A)^{-1}$ gilt, d. h., Gutartigkeit bezüglich X liegt nicht vor. Aus (8) folgt

$$\|\mathbf{X}\mathbf{A} - \mathbf{I}\|_{F} \leq vF \|\mathbf{A}\|_{2} \|\mathbf{X}\|_{F} \leq v\sqrt{n} F \|\mathbf{A}\|_{2} \|\mathbf{X}\|_{2} = v\hat{F} \operatorname{cond}_{2}(\mathbf{A})$$
(9)

mit $\hat{F} \approx \sqrt{n} F$, d. h., das Residuum XA - I ist klein, und X ist eine gute Linksinverse. Dagegen muß X keine gute Rechtsinverse sein, obwohl dies häufig wegen spezieller Fehlerkorrelationen der Fall ist: Es kann nur

$$\|AX - I\|_F = \|A(XA - I) A^{-1}\|_F \leq \|XA - I\|_F \operatorname{cond}_2(A) \leq \nu \widehat{F}[\operatorname{cond}_2(A)]^2$$
(10)

gezeigt werden. Wegen $X - A^{-1} = (XA - I) A^{-1}$ folgt schließlich aus (9)

$$\frac{\|X - A^{-1}\|_{F}}{\|A^{-1}\|_{2}} \leq \nu F \|A\|_{2} \|X\|_{F} \leq \nu \hat{F} \operatorname{cond}_{2}(A),$$
(11)

also die numerische Stabilität von 6.5.1.

Als alternative Möglichkeit soll ein nach GAUSS und JORDAN benanntes Vorgehen beschrieben werden; hinsichtlich der Namensgebung siehe B 6.6. Dabei wird die Matrix A nach der Vorschrift

$$A^{(1)} := A, \qquad A^{(k+1)} := D_k M_k A^{(k)} \qquad (k = 1, ..., n)$$
(12)

sukzessive mittels NT-Matrizen M_k und Diagonalmatrizen D_k der Gestalt

$$\boldsymbol{M}_{k} = \begin{pmatrix} 1 & -m_{1k} & & \\ & 1 & -m_{k-1,k} & & \\ & & 1 & & \\ & & -m_{k+1,k} & 1 & \\ & & & \ddots & \\ & & & -m_{nk} & & 1 \end{pmatrix} \quad \text{bzw.} \quad \boldsymbol{D}_{k} = \begin{pmatrix} 1 & & \\ \ddots & & \\ & 1 & & \\ & & d_{k} & \\ & & 1 & \\ & & \ddots & \\ & & & 1 \end{pmatrix} \quad (13)$$

äquivalent in die Einheitsmatrix I transformiert. Die Matrix $A^{(k)}$ stimmt in den ersten k - 1 Spalten bereits mit I überein, d. h., es ist

$$A^{(k)} = \begin{pmatrix} 1 & a_{1k}^{(k)} & \dots & a_{1n}^{(k)} \\ \vdots & \vdots & \vdots \\ 1 & a_{k-1,k-1}^{(k)} & \dots & a_{k-1,n}^{(k)} \\ \hline & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ \vdots & \vdots & \vdots \\ & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{pmatrix} = \begin{pmatrix} I_{k-1} \mid N^{(k)} \\ \hline O \mid M^{(k)} \end{pmatrix}.$$
(14)

Die Koeffizienten m_{ik} werden so bestimmt, daß in der k-ten Spalte von $M_kA^{(k)}$ außerhalb des Diagonalelementes Nullen entstehen. Dies führt auf

$$m_{ik} := a_{ik}^{(k)} / a_{kk}^{(k)} \quad (i = 1, ..., n, i \neq k).$$
⁽¹⁵⁾

Dabei muß das Pivotelement $a_{kk}^{(k)}$ natürlich von 0 verschieden sein, was der Einfachheit halber im folgenden vorausgesetzt werden soll. Mit der Festlegung

$$d_k := 1/a_{kk}^{(k)} \tag{16}$$

wird schließlich das Diagonalelement $a_{kk}^{(k)}$ von $M_k A^{(k)}$ zu 1 gemacht. Nach *n* Schritten ergibt sich dann $A^{(n+1)} = D_n M_n \cdots D_1 M_1 A$. Da $A^{(n+1)}$ nach Konstruktion gleich I ist, folgt

$$A^{-1} = \boldsymbol{D}_n \boldsymbol{M}_n \cdots \boldsymbol{D}_1 \boldsymbol{M}_1. \tag{17}$$

Die Darstellung (17) wird Produktform der Inversen genannt. Die signifikanten Koeffizienten m_{ik} bzw. d_k der Faktoren können auf dem Platz von a_{ik} bzw. a_{kk} gespeichert werden. Eine eventuell vorhandene Bandgestalt oder andersartige schwache Besetztheit schlägt sich in der Matrix der $\{m_{ik}, d_k\}$ nieder, während das explizit gebildete A^{-1} i. allg. voll besetzt ist.

Falls ein Ausdruck der Form $X = A^{-1}B$ zu berechnen ist, kann dies durch sukzessive Multiplikation mit den Faktoren geschehen. Speziell ergibt sich für $B = b \in \mathbb{R}^n$ gemäß

$$\boldsymbol{x} = \boldsymbol{A}^{-1}\boldsymbol{b} = \boldsymbol{D}_{\boldsymbol{n}}\boldsymbol{M}_{\boldsymbol{n}}\cdots\boldsymbol{D}_{\boldsymbol{1}}\boldsymbol{M}_{\boldsymbol{1}}\boldsymbol{b}$$
(18)

das sog. $Gau\beta$ -Jordan-Verfahren zur Lösung des linearen Gleichungssystems Ax = b. Die rechte Seite **b** wird dabei wie eine zusätzliche Spalte von A transformiert. Der Gesamtaufwand ist $\sim n^3/2$ opms, also größer als beim Gaußschen Algorithmus.

Für den Sonderfall B = I wird $X = A^{-1}$, und die zugehörige Transformation ist dann

$$X^{(1)} := I, \qquad X^{(k+1)} := D_k M_k X^{(k)} \qquad (k = 1, ..., n).$$
(19)

Die $X^{(k)}$ werden nach derselben Vorschrift wie die $A^{(k)}$ transformiert, lediglich die Startmatrizen sind verschieden, vgl. (12). Die Matrix $X^{(k)}$ unterscheidet sich nur in den ersten k - 1 Spalten von I. Im k-ten Schritt geht die k-te Spalte e^k von $X^{(k)}$ in

$$D_k M_k e^k = (-m_{1k}, ..., -m_{k-1,k}, d_k, -m_{k+1,k}, ..., -m_{nk})^{\mathsf{T}}$$

über, und die Spalten 1 bis k - 1 von $X^{(k)}$ werden genauso transformiert wie die Spalten k + 1 bis n von $A^{(k)}$. Die signifikanten Elemente von $X^{(k)}$ und $A^{(k)}$ können daher gemäß

$$oldsymbol{Z}^{(k)} := egin{pmatrix} x_{11}^{(k)} \ldots x_{1,k-1}^{(k)} & a_{1k}^{(k)} \ldots a_{1n}^{(k)} \ dots & dots & dots \ dots \$$

in einer (n, n)-Matrix $Z^{(k)}$ zusammengefaßt und simultan transformiert werden. Bei in-situ-Realisierung ergibt sich dann das folgende Verfahren der *Gauß-Jordan-Inversion*: **6.5.3.** Berechnung von A^{-1} nach dem Gauß-Jordan-Verfahren ohne Pivotisierung. Die reguläre Matrix $A \in \mathbb{R}^{n,n}$ besitze die Dreiecksfaktorisierung A = LR. Dann ist das Verfahren

for k := 1(1)n do Berechnung der k-ten Spalte: $a_{kk} := 1/a_{kk}$ for i := 1(1)n do [if $i \neq k$ then $a_{ik} := -a_{ik} * a_{kk}$] Berechnung der übrigen Spalten: for j := 1(1)n do if $j \neq k$ then $\begin{vmatrix} \text{if } j \neq k \text{ then} \\ | \text{ for } i := 1(1)n \text{ do } [\text{if } i \neq k \text{ then } a_{ij} := a_{ij} + a_{ik} * a_{kj}] \\ | a_{kj} := a_{kj} * a_{kk} \end{vmatrix}$

in exakter Arithmetik durchführbar und überspeichert A mit A^{-1} . Aufwand: $\sim n^3$ opms

6.5.4. Bemerkung. (i) Das Verfahren 6.5.3 erfordert denselben Aufwand wie 6.5.1. Sein Vorteil ist das kompakte Programm, allerdings ist das Rundungsfehlerverhalten etwas schlechter, vgl. (ii). Die nichtpivotisierte Version 6.5.3 ist i. allg. nur für diagonaldominante oder symmetrische und definite Matrizen geeignet; im letzten Fall braucht nur mit einem Dreieck gearbeitet zu werden. Für andere Matrizen muß eine der Pivotisierungsstrategien aus 5.2 angewendet werden. Man beachte, daß die Restmatrizen $M^{(k)}$ des Gauß-Jordan-Verfahrens und des Gaußschen Algorithmus bei gleicher Pivotisierung identisch sind. Im Fall der Spaltenpivotisierung gilt $|m_{ik}| \leq 1$ für i = k + 1, ..., n, aber i. allg. nicht für i = 1, ..., k - 1. Bei Spaltenpivotisierung ist eine als Stiefelsches Rotationsverfahren bezeichnete Implementierung zweckmäßig, bei der die Pivotzeile stets in die *n*-te Zeile gebracht wird und die übrigen zyklisch nach oben verschoben werden. Das Merken und Rückgängigmachen der Vertauschungen wird dabei überflüssig.

(ii) Die Rundungsfehleranalyse der $Gau\beta$ -Jordan-Inversion zeigt, daß das Verfahren für nicht zu schlecht konditioniertes A durchführbar und numerisch stabil im Sinne von

$$\frac{\|X - A^{-1}\|}{\|A^{-1}\|} \le \nu F \text{ cond } (A)$$
(20)

ist, wobei F = F(n) von der Matrizenklasse, der Pivotisierungsstrategie und der Norm abhängt. Dagegen liegt numerische Gutartigkeit i. allg. nicht vor, und beide Residuen ||AX - I||, ||XA - I|| können von der Größenordnung $vF[\text{cond}(A)]^2$ sein, d. h., die berechnete Inverse X braucht weder eine gute Rechts- noch Linksinverse zu sein. Ebenso ist die Berechnung der Lösung x von Ax = b nach dem Gauß-Jordan-Verfahren (18) ein numerisch stabiler, aber i. allg. kein numerisch gutartiger Prozeß und sollte daher durch den billigeren und gutartigen Gaußschen Algorithmus ersetzt werden. \Box

Wir diskutieren als nächstes das Verhalten von A^{-1} bei Rang-1-Änderungen von A.

6.5.5. Aufdatierung der inversen Matrix. Es sei $A \in \mathbb{R}^{n,n}$ eine reguläre Matrix. Dann ist

$$\bar{\boldsymbol{A}} := \boldsymbol{A} + \boldsymbol{u}\boldsymbol{v}^{\mathsf{T}}, \qquad \boldsymbol{u}, \boldsymbol{v} \in \boldsymbol{\mathsf{R}}^{\boldsymbol{n}}, \tag{21}$$

genau dann regulär, wenn

$$v^{\intercal}A^{-1}u \neq -1$$
 (22)

gilt. Ist die Regularitätsbedingung (22) erfüllt, so gilt

$$\bar{A}^{-1} = (A + uv^{\mathsf{T}})^{-1} = A^{-1} - \frac{A^{-1}uv^{\mathsf{T}}A^{-1}}{1 + v^{\mathsf{T}}A^{-1}u}.$$
(23)

Beweis. Falls (22) gilt, folgt (23) durch direktes Nachrechnen, indem die rechte Seite mit $A + uv^{\mathsf{T}}$ multipliziert wird. Im Fall $v^{\mathsf{T}}A^{-1}u = -1$ ist $w := A^{-1}u \neq o$ und $\bar{A}w = (1 + v^{\mathsf{T}}A^{-1}u) u = o$, also \bar{A} singulär. \Box

Die Aufdatierungsformel (23) wird nach WOODBURY und SHERMAN/MORRISON benannt, siehe B 6.7. Sie zeigt, daß eine Rang-1-Änderung von A eine Rang-1-Änderung von A^{-1} bewirkt, und \bar{A}^{-1} kann mit $\sim 3n^2$ opms aus A^{-1} berechnet werden, vgl. Ü 6.5.2.

Wir bemerken abschließend, daß inverse Matrizen in den Anwendungen meist in Ausdrücken der Form

$$\boldsymbol{X} = \boldsymbol{A}^{-1}\boldsymbol{B} \quad \text{bzw.} \quad \boldsymbol{Y} = \boldsymbol{C}\boldsymbol{A}^{-1} \tag{24}$$

mit $B, C^{\mathsf{T}} \in \mathbb{R}^{n,r}$ auftreten. Die naive Auswertung durch explizite Berechnung von A^{-1} und Multiplikation mit B bzw. C mit $\sim n^2(n + r)$ opms sollte grundsätzlich vermieden werden. Stattdessen wird zu den äquivalenten Matrixgleichungen

$$AX = B$$
 bzw. $A^{\mathsf{T}}Y^{\mathsf{T}} = C^{\mathsf{T}}$ (25)

übergegangen und $X = (x^1, ..., x^r)$ spaltenweise, $Y = (y^1, ..., y^r)^T$ zeilenweise aus den Gleichungssystemen

$$Ax^{j} = b^{j}$$
 bzw. $A^{\mathsf{T}}y^{j} = c^{j}$ $(j = 1, ..., r)$ (26)

berechnet, wobei $B = (b^1, ..., b^r)$ und $C^{\intercal} = (c^1, ..., c^r)$ ist. Dies erfordert eine Dreiecksfaktorisierung und die Behandlung *r* rechter Seiten, also $\sim n^2(n/3 + r)$ opms, und ist vom Rundungsfehlereinfluß günstiger, vgl. Ü 6.5.4.

Übungsaufgaben

Ü 6.5.1. Man überlege sich, daß S1, S2 aus 6.5.1 in den folgenden n- bzw. (n - 1)-stufigen Versionen realisiert werden können:

S1': for k := 1(1)n do

$$\begin{array}{l} a_{kk} := 1/a_{kk} \\ \text{for } i := 1(1)k - 1 \text{ do } a_{ik} := -a_{ik} * a_{kk} \\ \text{for } j := k + 1(1)n \text{ do} \\ \\ & \left| \begin{array}{c} r := a_{kj}, a_{kj} := 0 \\ \text{for } i := 1(1)k \text{ do } a_{ij} := a_{ij} + a_{ik} * r \end{array} \right| \end{array}$$

S2': for k := n - 1(-1)1 do

for
$$i := k + 1(1)n$$
 do $[l_i := a_{ik}, a_{ik} := 0]$
for $j := k + 1(1)n$ do
for $i := 1(1)n$ do $a_{ik} := a_{ik} - a_{ii} * l_i$

Hinweis: S1' entsteht aus S1, indem der Beitrag des im k-ten Hauptschritt berechneten s_{ik} zu jeder der Summen $\sum_{p=i}^{j-1} s_{ip}r_{pj}$ (k < j) aus (5) sofort auf dem Platz von s_{ij} akkumuliert wird. Überdies ist S1' gerade das Gauß-Jordan-Inversionsverfahren für den Sonderfall A = R. S2' entsteht durch Auswertung von $T = SL^{-1} = SL_{n-1}(-l^{n-1})\cdots L_1(-l^1)$.

Ü 6.5.2. Es sei $\bar{A} = A + uv^{\mathsf{T}}$ mit regulärem A, \bar{A} . Man überzeuge sich, daß der nachfolgende Algorithmus eine Umsetzung der Formel (23) mit $\sim 3n^2$ opms darstellt, bei der $B := A^{-1}$ mit $\bar{B} := \bar{A}^{-1}$ überspeichert wird:

$$w := Bu, \quad \alpha := 1 + v^{\mathsf{T}}w$$

 $u := w/\alpha, \quad w := B^{\mathsf{T}}v$
 $B := B - uw^{\mathsf{T}}.$

Ü 6.5.3. Man spezifiziere die Änderungsformel (23) für die Sonderfälle der Änderung einer Zeile, einer Spalte und eines Elementes von A.

Ü 6.5.4. Man zeige: (i) Wenn X eine im Sinne von

 $\|XA - I\| \leq \nu F \text{ cond } (A)$

gute Linksinverse ist, genügt x := Xb der Abschätzung

 $\|\boldsymbol{x} - \boldsymbol{A}^{-1}\boldsymbol{b}\| \leq \nu F \text{ cond } (\boldsymbol{A}) \|\boldsymbol{A}^{-1}\boldsymbol{b}\|,$

d. h., X liefert auch eine gute Lösung von Ax = b.

(ii) Wenn X eine im Sinne von

 $||AX - I|| \leq \nu F \text{ cond } (A)$

gute Rechtsinverse ist, genügt x := Xb dagegen der Abschätzung

 $\|\boldsymbol{x} - \boldsymbol{A}^{-1}\boldsymbol{b}\| \leq \nu F[\text{cond}(\boldsymbol{A})]^2 \|\boldsymbol{A}^{-1}\boldsymbol{b}\|,$

und der Exponent 2 kann i. allg. nicht durch 1 ersetzt werden.

Bemerkungen zum Kapitel 6

B 6.1. Die Faktorisierung $A = LL^{\intercal}$ einer symmetrisch positiv definiten Matrix und das zugehörige Verfahren zur direkten Berechnung von L wurden von dem französischen Offizier CHOLESKY eingeführt und zur Lösung von Normalgleichungen linearer Quadratmittelprobleme verwendet; zu den letzteren Begriffen siehe Kapitel 9. Nach dem Tod CHOLESKYS ist das Verfahren durch dessen Freund BENOIT 1924 veröffentlicht worden. Algorithmen zur symmetrischen Faktorisierung indefiniter symmetrischer Matrizen sind dagegen erst in den letzten 15 Jahren entwickelt worden, siehe BUNCH/PARLETT [71] für eine detaillierte Darstellung der Grundideen und historische Kommentare, vgl. auch BUNCH [71] und BUNCH/KAUFMAN [77]. Die hier verwendete Pivotisierungsstrategie geht auf BUNCH/KAUFMAN/PARLETT [76] zurück und wird auch von den LINPACK-Autoren DONGARRA et al. [79] verwendet. Instruktive Vergleiche mit alternativen Methoden — insbesondere mit der von AASEN [71] vorgeschlagenen Kongruenztransformation auf Tridiagonalform — haben BARWELL/GEORGE [76] durchgeführt.

B 6.2. Direkte Dreiecksfaktorisierungen nichtsymmetrischer Matrizen sind u. a. von BANA-CHIEWICZ, CROUT und DOOLITTLE angegeben worden. Sie unterscheiden sich in der Zuordnung der Einsdiagonalen zu L oder R und in der Berechnungsreihenfolge. Für Details und Quellenangaben muß auf die Spezialliteratur — etwa HOUSEHOLDER [64] und FADDEEV/FADDEEVA [63] — verwiesen werden.

B 6.3. Die hier beschriebenen Algorithmen zur Aufdatierung von Cholesky-Faktorisierungen gehen im Fall $\alpha = 1$ auf GOLUB [65] und im Fall $\alpha = -1$ auf SAUNDERS [72] zurück und sind auch in LINPACK aufgenommen worden. Alternative und z. T. etwas billigere Varianten haben FLETCHER/POWELL [74], BERSENEV [79], DAX [83] u. a. angegeben. SORENSEN [77] behandelt die Aufdatierung von BKP-Faktorisierungen.

Zur Aufdatierung von Dreiecksfaktorisierungen nichtsymmetrischer Matrizen wurden bisher nur Algorithmen publiziert, die entweder ohne Pivotisierung arbeiten — also nur für spezielle Matrizenklassen stabil sind — oder bei denen die Dreiecksform eines Faktors allmählich zerstört wird, siehe etwa BENNETT [65] oder BARTELS/STOER/ZENGER [71]. Eine Wende brachte der Elementareliminationsschritt 6.3.1, der erstmals von BURMEISTER [76] in einem unveröffentlichten Manuskript für den Fall der Spaltenmodifikation entsprechend 6.3.5 (ii) vorgeschlagen worden ist. Durch die Autoren wurde dieser Algorithmus auf den Fall bleibiger Rang-1-Modifikationen erweitert und auf sein Fehlerverhalten untersucht. Unabhängig und etwa gleichzeitig ist diese Idee von FLETCHER/MATTHEWS [83] gefunden worden. Eine Übersicht über Modifikationstechniken haben GILL/MURRAY [77] gegeben, vgl. auch 10.3 und die dort zitierte Literatur über die Aufdatierung von QR-Faktorisierungen.

B 6.4. Algorithmen zur Faktorisierung von Bandmatrizen sind in WILKINSON/REINSCH [71] und DONGARRA et al. [79] implementiert worden, und zwar auch für den nicht ganz einfachen Fall nichtsymmetrischer Matrizen und Rechnung mit Spaltenpivotisierung.

B 6.5. Über die Lösung schwach besetzter Gleichungssysteme existiert eine umfangreiche Literatur. Wir verweisen lediglich auf die Sammelbände von BUNCH/ROSE [76], DUFF/STE-WART [79] und DUFF [81], die Monographien von GEORGE/LIU [81] und PISSANETZKY [84] sowie den Übresichtsartikel von IKRAMOV [82], wo zahlreiche weiterführende Literaturhinweise zu finden sind. Zum erstenmal scheint die schwache Besetztheit einer Matrix bei Eliminationsverfahren von MARKOWITZ [57] in der linearen Optimierung ausgenutzt worden zu sein.

B 6.6. Der Inversionsalgorithmus 6.5.1 ist eine verkettete Variante des LINPACK-Verfahrens das in Ü 6.5.1 angegeben worden ist. Das Gauß-Jordan-Verfahren war vermutlich weder GAUSS noch dem Geodäten JORDAN bekannt, trotzdem hat sich der Name in der Literatur eingebürgert. Die Fehleranalyse geht auf PETERS/WILKINSON [75] zurück. Die Implementierung der spaltenpivotisierten Version nach dem Rotationsprinzip gemäß 6.5.4 (ii) kann bei STIEFEL [61] nachgelesen werden. Das Gauß-Jordansche Inversionsverfahren ist identisch mit dem sog. Austauschverfahren der linearen Optimierung.

B 6.7. Die Änderungsformel (23) wurde von SHERMAN/MORRISON 1953 angegeben; bereits vorher hatte WOODBURY 1950 eine analoge Formel für Rang-*m*-Änderungen $A + UV^{\intercal}$, $U, V \in \mathbb{R}^{n,m}$, aufgestellt; siehe ZIELKE [70] für eine detaillierte Diskussion solcher Änderungsformeln und historische Kommentare.

7. Zusammenfassung zum Teil II

Unter den direkten, d. h. in exakter Arithmetik in endlich vielen Schritten zum Ziel führenden Verfahren zur Lösung regulärer Gleichungssysteme nimmt der Gaußsche Algorithmus mit seinen verschiedenen Modifikationen eine herausragende Rolle ein. Zweckmäßig und heute üblich ist die Interpretation als Dreiecksfaktorisierung und die Trennung des Faktorisierungsprozesses von der Behandlung einer rechten Seite, was insbesondere bei der Lösung mehrerer Gleichungssysteme mit derselben Koeffizientenmatrix vorteilhaft ist.

Für die wesentlichen Klassen der diagonaldominanten und symmetrischen positiv definiten Matrizen ist die Dreiecksfaktorisierung ohne Pivotisierung durchführbar und numerisch gutartig. Für nicht zu diesen Klassen gehörende Matrizen muß die Gutartigkeit durch eine geeignete Pivotisierungsstrategie erzwungen werden, was allerdings zu einer Erhöhung des Organisationsaufwandes führt. Es besteht heute Übereinstimmung in der Ansicht, daß die vergleichsweise billige Spaltenpivotisierung in den meisten Fällen ausreichend ist, obwohl die das schlechteste Verhalten erfassenden theoretischen Abschätzungen recht pessimistisch sind, siehe 5.2. und 5.3.

Wenn die Pivotisierung nicht ausreicht, um eine akzeptable Dreiecksfaktorisierung zu liefern, kann i. allg. durch wenige Schritte der iterativen Verbesserung gemäß 5.4 eine akzeptable Lösung des Gleichungssystems erhalten werden. Falls die iterative Verbesserung nicht konvergiert, muß zu einer Arithmetik mit höherer Genauigkeit übergegangen werden. Eine weitere Möglichkeit stellt in diesem Fall der Ersatz der Dreiecksfaktorisierung durch orthogonale Faktorisierungen dar. Mit verdoppeltem Aufwand läßt sich die **QR**-Faktorisierung berechnen, siehe 10.2 und speziell 10.2.11 (iii). Ein Höchstmaß an Information gibt die Singulärwertzerlegung, die allerdings den vierfachen Aufwand der Dreiecksfaktorisierung erfordert, siehe 11.1.

Ein wesentlicher Vorteil des Gaußschen Algorithmus ist seine Anpassungsfähigkeit. Dies betrifft einmal die Möglichkeit, spezielle Eigenschaften der Matrix zur Aufwandsreduktion auszunutzen. Für symmetrische Matrizen kann der Aufwand durch Verwendung der Cholesky-Faktorisierung im definiten Fall oder der neueren BKP-Faktorisierung im indefiniten Fall halbiert werden, siehe 6.1. Wenn die Matrix Bandstruktur hat, überträgt sich diese auf die Dreiecksfaktoren. Falls ohne Pivotisierung gearbeitet werden kann, bleibt die Bandbreite erhalten; bei Spaltenpivotisierung hat der \mathbf{R} -Faktor i. allg. eine größere Bandbreite als der entsprechende Teil der Originalmatrix, siehe 6.4. Bei nicht zu großer Bandbreite sind daher auch Systeme sehr hoher Dimension effektiv lösbar. In Verbindung mit Kompaktspeichertechniken kann der Gaußsche Algorithmus auch zur Lösung großer Gleichungssysteme mit schwach besetzten Koeffizientenmatrizen, die keine Bandstruktur haben, verwendet werden, siehe wieder 6.4.

Zum anderen kann der Gaußsche Algorithmus speziellen Eigenschaften des Computers, der Programmiersprache und des Betriebssystems angepaßt werden. Bei spaltenweiser Speicherung von zweidimensionalen Feldern läßt sich unnötiger Seitenwechsel durch spaltenweise Orientierung der inneren Schleifen vermeiden; dasselbe Vorgehen ist auch für gewisse Vektorrechner zweckmäßig, vgl. auch Kapitel 16. Wenn Skalarprodukte billig mit höherer Genauigkeit akkumuliert werden können, empfehlen sich die direkten Faktorisierungen aus 6.2. Für Parallelcomputer stellt allerdings die QR-Faktorisierung nach GIVENS — siehe 10.2. — eine Konkurrenz für die Dreiecksfaktorisierung dar, da sie keine Pivotisierung erfordert und besser parallelisierbar ist. Dies trifft insbesondere für die sog. "systolic arrays" zu; zu diesen siehe Kapitel 16.

Schließlich lassen sich Dreiecksfaktorisierungen bei Rang-1-Änderungen der Matrix billig und numerisch stabil aufdatieren, was für gewisse Aufgabenklassen — etwa die Lösung nichtlinearer Gleichungssysteme und linearer wie nichtlinearer Optimierungsaufgaben — von Bedeutung ist. Zusammenfassend kann festgestellt werden, daß der Gaußsche Algorithmus für nicht zu schlecht konditionierte Systeme mittlerer Dimension — n in der Größenordnung von einigen Hundert — und für schwach besetzte Systeme hoher Dimension mit Bandstruktur oder unregelmäßiger Besetztheitsstruktur das bevorzugte Verfahren darstellt und kaum durch Konkurrenten gefährdet ist.

Etwas anders ist die Situation bei schwach besetzten Systemen hoher Dimension und regelmäßiger Besetztheitsstruktur, wie sie bei der Diskretisierung partieller Differentialgleichungen nach der Methode der finiten Elemente oder dem Differenzenverfahren entstehen. Neben speziellen Versionen des Cholesky-Verfahrens — siehe etwa GEORGE/LIU [81] — treten hier iterative Verfahren als ernsthafte Konkurrenten auf. Auf solche Verfahren kann im Rahmen dieser Einführung nicht eingegangen werden. Wir erwähnen lediglich die klassischen Darstellungen VARGA [62], YOUNG [71] und HAGEMAN/YOUNG [81] und von den neueren Entwicklungen die Methode der alternierenden Dreiecksfaktorisierung — siehe SAMARSKIĬ/NIKOLAEV [78] —, die Mehrgitterverfahren — siehe etwa HACKBUSCH [85] — und iterative Verfahren mit Vorkonditionierungen. Ein Beispiel für die letztgenannte Klasse sind die ICCG-Verfahren, bei denen eine sog. unvollständige Cholesky-Faktorisierung ("Incomplete Cholesky") als Vorkonditionierung für das Verfahren der konjugierten Gradienten ("Conjugate Gradients") verwendet wird. Zu den Verfahren der konjugierten Gradienten sei auf KuznEcov [83] verwiesen.

Auf sehr spezielle Gleichungssysteme konnten wir nicht eingehen, etwa auf solche mit Töplitz-Matrizen — siehe VOEVODIN/TYRTYŠNIKOV [83] — und Ljapunov-Gleichungen — siehe IKRAMOV [84]. Diese und andere Gleichungstypen werden auch von GOLUB/VAN LOAN [83] behandelt.

III. Lineare Quadratmittelprobleme

Nachdem wir im vorangegangenen Teil II Gleichungssysteme mit quadratischen und regulären Koeffizientenmatrizen betrachtet haben, lassen wir jetzt beliebige rechteckige Koeffizientenmatrizen $A \in \mathbb{R}^{m,n}$ zu. Die zugehörigen Gleichungssysteme werden dann auch rechteckig genannt: Im Fall $m \ge n$ spricht man von überbestimmten, im Fall m < n von unterbestimmten Systemen. Teil III ist der Quadratmittellösung überbestimmter linearer Gleichungssysteme gewidmet. In den einführenden Abschnitten gehen wir dabei auch auf unterbestimmte Systeme und klassische Lösungen ein.

8. Rechteckige Gleichungssysteme und Quadratmittelprobleme

8.1. Klassische Lösungen, Quadratmittellösungen, Pseudoinverse

A. Lösungen und Quadratmittellösungen rechteckiger Gleichungssysteme

Gegeben ist das rechteckige Gleichungssystem

$$A\boldsymbol{x} = \boldsymbol{b} \tag{1}$$

mit $A \in \mathbb{R}^{m,n}$, $b \in \mathbb{R}^m$. Gesucht sind alle Lösungen $x \in \mathbb{R}^n$ von (1), deren Gesamtheit die Lösungsmenge $\mathcal{L}_0 = \mathcal{L}_0(A, b) := \{x \in \mathbb{R}^n : Ax = b\}$ bildet. Zur weiteren Untersuchung von (1) ziehen wir die Singulärwertzerlegung

$$A = U\Sigma V^{\mathsf{T}} \tag{2}$$

von A mit $\Sigma = \text{diag}(\sigma_1, ..., \sigma_r, 0, ..., 0), \sigma_1 \ge \cdots \ge \sigma_r > 0, r = \text{rang}(A)$ und orthogonalen Matrizen $U \in \mathbb{R}^{m,m}, V \in \mathbb{R}^{n,n}$ heran, vgl. 1.2.14. Einsetzen in (1) führt unter Beachtung von $U^{-1} = U^{\mathsf{T}}$ auf das äquivalente System

$$\Sigma \xi = \beta \tag{3}$$

für die transformierten Größen

$$\mathbf{\tilde{\xi}} = (\xi_j) := \mathbf{V}^{\mathsf{T}} \mathbf{x}, \qquad \boldsymbol{\beta} = (\beta_i) := \mathbf{U}^{\mathsf{T}} \boldsymbol{b}, \tag{4}$$

und (3) läßt sich komponentenweise als

$$\sigma_{1}\xi_{1} = \beta_{1},$$

$$\vdots$$

$$\sigma_{r}\xi_{r} = \beta_{r},$$

$$0 = \beta_{r+1},$$

$$\vdots$$

$$0 = \beta_{m}$$
(5)

schreiben. Wegen $\boldsymbol{x} = V\boldsymbol{\xi}$ sind die ξ_i gerade die Komponenten von \boldsymbol{x} in der orthonormierten Basis des Urbildraumes \mathbb{R}^n , die durch die Spalten von V gebildet wird. Analog stellen die β_i die Komponenten von $\boldsymbol{b} = \boldsymbol{U}\boldsymbol{\beta}$ in der durch die Spalten von \boldsymbol{U} gebildeten orthonormierten Basis des Bildraumes \mathbb{R}^m dar. In den genannten Basen nimmt das Ausgangssystem die Form (5) an, d. h., die Gleichungen sind entkoppelt, und das System zerfällt in die r Gleichungen

$$\sigma_i \xi_i = \beta_i \qquad (i = 1, \dots, r) \tag{6}$$

für ξ und die m - r Bedingungen

$$0 = \beta_i \qquad (i = r + 1, ..., m)$$
(7)

an die transformierte rechte Seite β . Das System (1) ist also genau dann lösbar, wenn (7) erfüllt ist. Wegen

$$\mathscr{R}(\mathbf{\Sigma}) = \{ \boldsymbol{\beta} = (\boldsymbol{\beta}_i) \in \mathbf{R}^m : \boldsymbol{\beta}_{i+1} = \cdots = \boldsymbol{\beta}_m = 0 \}$$

ist (7) zu $\beta \in \mathcal{R}(\Sigma)$, also zu $\boldsymbol{b} = \boldsymbol{U}\boldsymbol{\beta} \in \boldsymbol{U}\mathcal{R}(\Sigma) = \mathcal{R}(\boldsymbol{A})$ äquivalent, vgl. (1.2.24). Die Lösbarkeitsbedingung (7) läßt sich daher auch als

$$\boldsymbol{b} \in \mathcal{R}(\boldsymbol{A}) \tag{8}$$

schreiben und ist uns in dieser Gestalt bereits im Abschnitt 1.1.E begegnet, denn die Auflösung von (1) bedeutet, den Vektor **b** als Linearkombination der Spalten von A mit den Koeffizienten x_i darzustellen. Wenn (8) erfüllt ist, wird (1) konsistent genannt, andernfalls heißt das System *inkonsistent*. Im Fall m = r, d. h. für eine zeilenreguläre Matrix, ist die Lösbarkeitsbedingung (8) wegen $\mathcal{R}(A) = \mathbb{R}^m$ für jede rechte Seite erfüllt.

Es sei jetzt (1) konsistent. Dann sind die ersten r Komponenten von ξ durch (6) eindeutig festgelegt, während die restlichen in (5) überhaupt nicht auftreten und deshalb beliebig gewählt werden können. Die Lösungen x des Originalsystems (1) ergeben sich wegen (4) gemäß $x = V\xi$ aus den Lösungen ξ des transformierten Systems (5), sie sind also durch

$$\boldsymbol{x} = \boldsymbol{x}^N + \boldsymbol{\hat{x}} \tag{9}$$

mit

$$\boldsymbol{x}^{N} = \boldsymbol{V}\boldsymbol{\xi}^{N}, \quad \boldsymbol{\xi}^{N} = (\xi_{1}, ..., \xi_{r}, 0, ..., 0)^{\mathsf{T}}, \quad \xi_{j} = \beta_{j}/\sigma_{j} \quad (j = 1, ..., r)$$
 (10)

und

$$\hat{\boldsymbol{x}} = \boldsymbol{V}\hat{\boldsymbol{\xi}}, \quad \hat{\boldsymbol{\xi}} = (0, ..., 0, \xi_{r+1}, ..., \xi_n)^{\mathsf{T}}, \quad \xi_j \text{ beliebig } (j = r+1, ..., n)$$
(11)

16 Schwetlick, Numerische Algebra

festgelegt. Dabei ist x^{N} eine spezielle Lösung des *inhomogenen Systems* Ax = b (nämlich diejenige kleinster Euklidischer Norm, wie wir später sehen werden). Wegen

$$\mathscr{N}(\Sigma) = \{ \boldsymbol{\xi} = (\xi_i) : \xi_1 = \dots = \xi_r = 0 \}$$

läßt sich $\hat{\xi}$ aus (11) durch $\hat{\xi} \in \mathcal{N}(\Sigma)$, \hat{x} also durch $\hat{x} = V\hat{\xi} \in V\mathcal{N}(\Sigma) = \mathcal{N}(A)$ charakterisieren. Daher stellt \hat{x} in (9), (11) gerade irgendein Element aus dem (n - r)-dimensionalen Nullraum $\mathcal{N}(A)$, d. h. irgendeine Lösung des homogenen Systems Ax = o dar. Die Lösungsmenge \mathcal{L}_0 kann folglich in der Form

$$\mathscr{L}_{m{0}}(m{A},\,m{b})=\{m{x}+m{x}^{N}+m{\hat{x}}:m{\hat{x}}\in\mathcal{N}(m{A})\}=:m{x}^{N}+\mathcal{N}(m{A})\}$$

geschrieben werden.

Wenn \mathcal{X} ein Teilraum von \mathbb{R}^n der Dimension s und $y \in \mathbb{R}^n$ ist, wird die Menge

$$\mathscr{L} = \{ oldsymbol{x} \in {\sf R}^n \colon oldsymbol{x} = oldsymbol{y} + oldsymbol{z}, \, oldsymbol{z} \in \mathscr{X} \} = oldsymbol{y} + \mathscr{X}$$

allgemein *lineare Mannigfaltigkeit* der Dimension $s = \dim(\mathcal{X})$ genannt. Mit $\mathbf{x}' \in \mathcal{L}$ gilt $\mathbf{x}'' \in \mathcal{L}$ genau dann, wenn $\mathbf{x}' - \mathbf{x}'' \in \mathcal{X}$ ist. Insbesondere kann \mathcal{L} auch in der Form $\mathcal{L} = \mathbf{y}' + \mathcal{X}$ mit beliebigem $\mathbf{y}' \in \mathcal{L}$ geschrieben werden, d. h., für die Darstellung von \mathcal{L} kann ein beliebiger Bezugspunkt $\mathbf{y}' \in \mathcal{L}$ gewählt werden.

Mit dieser Terminologie können wir sagen: Das konsistente System (1) hat die (n-r)-dimensionale lineare Mannigfaltigkeit $\mathcal{L}_0 = \mathbf{x}^* + \mathcal{N}(\mathbf{A})$ zur Lösungsmenge, wobei \mathbf{x}^* irgendeine Lösung von (1) darstellt.

In den Anwendungen sind die Eingangsdaten $\{A, b\}$ i. allg. fehlerbehaftet. Wenn (1) etwa ein streng überbestimmtes, konsistentes System ist, tritt der Teil (7) im transformierten System (3) wegen $r \leq n < m$ stets auf. Eine beliebig kleine Störung $\delta\beta_i \neq 0$ in einer der Komponenten $\beta_i = 0$, $i \in \{r + 1, ..., m\}$, von β wandelt dann das lösbare Originalsystem in ein gestörtes System mit der rechten Seite $\beta + \delta\beta$ bzw. $\mathbf{b} + \delta \mathbf{b} = \mathbf{U}(\beta + \delta\beta)$ um, das nicht mehr lösbar ist. Für solche Aufgaben-klassen ist daher die Forderung nach der exakten Erfüllung von (1) unnötig streng bzw. sogar unsachgemäß und sollte durch die Forderung ersetzt werden, die p-Norm des Residuums

$$r = r(x) = b - Ax$$

zu minimieren. Obwohl auch die Fälle p = 1 und $p = \infty$ von Interesse sind, hat der Fall p = 2, d. h. die Verwendung der Euklidischen Norm, die größte Bedeutung für die Anwendungen erlangt. Da wir im Teil III ausschließlich die Euklidische Vektornorm und meist die Spektralnorm einer Matrix verwenden, lassen wir den Index 2 zur Kennzeichnung dieser Normen der Einfachheit halber weg. Das abgeschwächte Ersatzproblem schreibt sich dann als

$$\|b - Ax\| o \operatorname{Minimum}!$$
 bzw. kurz: $Ax \cong b$, $x \in \mathbb{R}^n$

und ist äquivalent zu

$$(\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x})^{\mathsf{T}} (\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}) = \sum_{i=1}^{m} [(\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x})_i]^2 \to \operatorname{Minimum}!.$$
 (12)

Es wird lineares Quadratmittelproblem (engl. "linear least squares problem", russ. "линейная проблема наименьших кбадратоб") genannt; man spricht auch von der Methode der kleinsten (Fehler-) Quadrate bzw. vom linearen Ausgleichsproblem.

Ein Vektor $x \in \mathbf{R}^n$, für den

$$\|\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}\| \leq \|\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}'\|$$
 für alle $\boldsymbol{x}' \in \mathsf{R}^n$

gilt, heißt Lösung des Quadratmittelproblems $Ax \cong b$ bzw. Quadratmittellösung des linearen Gleichungssystems Ax = b; die Gesamtheit der Lösungen bildet die Lösungsmenge $\mathcal{I} = \mathcal{I}(A, b)$. Für konsistente Systeme gilt offensichtlich $\mathcal{I}_0(A, b) = \mathcal{I}(A, b)$, so daß eine Verallgemeinerung des klassischen Lösungsbegriffes vorliegt.

Wir betrachten jetzt einen beliebigen Vektor $x \in \mathbb{R}^n$. Wenn A wie eingangs durch die Singulärwertzerlegung ausgedrückt wird, ergibt sich mit den Transformationen (4) die Darstellung

$$\boldsymbol{r} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x} = \boldsymbol{U}\boldsymbol{U}^{\mathsf{T}}\boldsymbol{b} - \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathsf{T}}\boldsymbol{x} = \boldsymbol{U}(\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\xi}). \tag{13}$$

Unter Beachtung der Orthogonalinvarianz der Euklidischen Norm folgt daraus

$$\|\mathbf{r}\|^{2} = \|\boldsymbol{\beta} - \boldsymbol{\Sigma}\boldsymbol{\xi}\|^{2} = \sum_{i=1}^{r} (\beta_{i} - \sigma_{i}\xi_{i})^{2} + \sum_{i=r+1}^{m} \beta_{i}^{2}.$$
 (14)

Der erste Summand ist stets nichtnegativ und verschwindet genau dann, wenn

$$\sigma_i \xi_i = eta_i, \quad ext{also} \quad \xi_i = eta_i / \sigma_i \quad (i = 1, ..., r)$$

gilt. Die restlichen Komponenten haben auf das Residuum keinen Einfluß und können beliebig gewählt werden. Der Vergleich mit (9) bis (11) zeigt, daß die Quadratmittellösungen identisch sind mit den Lösungen des Gleichungssystems $Ax = \bar{b}$, zu dem die rechte Seite

$$\bar{\boldsymbol{b}} := \boldsymbol{U}\bar{\boldsymbol{\beta}}, \qquad \bar{\boldsymbol{\beta}} := (\beta_1, \dots, \beta_r, 0, \dots, 0)^{\mathsf{T}}$$
(15)

gehört.

- 8.1.1. Satz. Es sei $A \in \mathbf{R}^{m,n}$ gegeben. Dann gilt
- (i) Das Quadratmittelproblem $Ax \simeq b$ ist für jede rechte Seite $b \in \mathbb{R}^m$ lösbar.
- (ii) \boldsymbol{x} ist Lösung von $A\boldsymbol{x} \cong \boldsymbol{b}$ genau dann, wenn \boldsymbol{x} das Gleichungssystem

$$A\boldsymbol{x} = \boldsymbol{\bar{b}} \tag{16}$$

mit $\bar{\boldsymbol{b}}$ gemäß (15) löst.

(iii) Jede Lösung \boldsymbol{x} läßt sich in der Form

$$\boldsymbol{x} = \boldsymbol{x}^N + \boldsymbol{\hat{x}} \quad \text{mit} \ \boldsymbol{x}^N, \, \boldsymbol{\hat{x}} \text{ gemäß} \ (10), \ (11)$$
 (17)

schreiben.

Satz 8.1.1 besagt: Die lineare Mannigfaltigkeit

$$\mathscr{L}(\boldsymbol{A},\boldsymbol{b}) = \boldsymbol{x}^{N} + \mathscr{N}(\boldsymbol{A}) \tag{18}$$

der Dimension dim $\mathcal{N}(A) = n - r$, $r = \operatorname{rang}(A)$, stellt die Lösungsmenge von $Ax \simeq b$ dar. Die Lösung ist eindeutig, wenn r = n gilt, d. h., wenn A spaltenregulär ist.

Wir bemerken ferner, daß das optimale Residuum r = b - Ax für alle Lösungen x den Wert

$$oldsymbol{r} = oldsymbol{b} - oldsymbol{ar{b}} = oldsymbol{U}(eta - oldsymbol{ar{eta}}) = oldsymbol{U} oldsymbol{arrho}, \qquad oldsymbol{arrho} := (0, ..., 0, eta_{r+1}, ..., eta_m)^\intercal$$

besitzt, also durch A und b eindeutig festgelegt ist, vgl. (13).

Nachdem die Singulärwertzerlegung von A über die Residuumsdarstellung (14) in einfacher Weise zu Aussagen über Lösbarkeit und Lösungsmenge des Quadratmittelproblems $Ax \cong b$ geführt hat, soll im folgenden eine von der Singulärwertzerlegung unabhängige, geometrisch motivierte Charakterisierung der Quadratmittellösungen hergeleitet werden. Wir betrachten dazu den Fall $m = 3, n = \operatorname{rang}(A)$ $= r = 2, A = (a^1, a^2) \in \mathbb{R}^{3,2}, b \in \mathbb{R}^3, b \notin \mathcal{R}(A)$. Für $x = (x_1, x_2)^{\mathsf{T}} \in \mathbb{R}^2$ durchläuft $y = Ax = a^1x_1 + a^2x_2$ den Wertebereich $\mathcal{R}(A) = \operatorname{span} \{a^1, a^2\}$. Die Minimierung von ||b - Ax|| bedeutet also, den Vektor b im Sinne der Euklidischen Norm optimal durch eine Linearkombination der Spalten von A zu approximieren, siehe Abb. 8.1.1.



Abb. 8.1.1. Quadratmittellösung $ar{x}$ und minimales Residuum $m{b} - ar{m{b}}$

B. Orthogonale Projektionen, Projektoren und Komplemente

Aus Abb. 8.1.1 wird anschaulich klar, daß $\mathbf{r} = \mathbf{b} - \mathbf{\bar{b}}$, $\mathbf{\bar{b}} = A\mathbf{\bar{x}} \in \mathcal{R}(A)$, genau dann kleinste Euklidische Länge besitzt, wenn \mathbf{r} orthogonal zu allen Vektoren $\mathbf{y} \in \mathcal{R}(A)$ und $\mathbf{\bar{b}}$ in diesem Sinne die orthogonale Projektion von \mathbf{b} auf $\mathcal{R}(A)$ ist. Im folgenden soll der für das Verständnis der Quadratmittelprobleme wesentliche Begriff des orthogonalen Projektors präzise definiert und charakterisiert werden. Dabei lassen wir statt $\mathcal{R}(A)$ einen beliebigen Teilraum \mathcal{Y} zu.

8.1.2. Aussage. Es sei \mathcal{Y} ein Teilraum von \mathbb{R}^m . Dann gilt:

(i) Zu jedem $\boldsymbol{b} \in \mathbf{R}^{m}$ gibt es genau ein $\boldsymbol{\bar{b}} \in \boldsymbol{\mathcal{Y}}$ mit

$$\mathbf{y}^{\intercal}(\mathbf{b}-\bar{\mathbf{b}})=0$$
 für alle $\mathbf{y}\in\mathcal{Y},$ kurz: $\mathbf{b}-\mathbf{b}\perp\mathcal{Y}.$ (19)

Der Vektor $\bar{\boldsymbol{b}}$ heißt orthogonale Projektion von \boldsymbol{b} auf \mathcal{Y} .

(ii) $\bar{\boldsymbol{b}}$ ist orthogonale Projektion von \boldsymbol{b} auf \mathcal{Y} genau dann, wenn $\bar{\boldsymbol{b}}$ die Aufgabe

$$\|\boldsymbol{b} - \boldsymbol{y}\| \to \underset{\boldsymbol{y} \in \boldsymbol{\mathcal{Y}}}{\operatorname{Minimum}!}$$

$$\tag{20}$$

löst.

(iii) Es gibt genau eine Matrix $P \in \mathbb{R}^{m,m}$, so daß

$$\overline{\boldsymbol{b}} = \boldsymbol{P}\boldsymbol{b} \tag{21}$$

für alle $\boldsymbol{b} \in \mathbf{R}^m$ gilt. \boldsymbol{P} heißt orthogonaler Projektor bzw. kurz Projektor auf \mathcal{Y} . Wenn dim $\mathcal{Y} = r$ und $\{\boldsymbol{g}^1, \dots, \boldsymbol{g}^r\}$ eine Basis von \mathcal{Y} ist, kann \boldsymbol{P} in der Form

$$\boldsymbol{P} = \boldsymbol{G}(\boldsymbol{G}^{\mathsf{T}}\boldsymbol{G})^{-1}\,\boldsymbol{G}^{\mathsf{T}} \tag{22}$$

mit $G := (g^1, ..., g^r)$ geschrieben werden.

Be we is. Es sei $\{g^1, \ldots, g^r\}$ irgendeine Basis von \mathcal{Y} und $G := (g^1, \ldots, g^r)$. Dann ist rang (G)= rang $(G^{\mathsf{T}}G) = r$, vgl. (1.2.28), so daß $G^{\mathsf{T}}G$ regulär und P aus (22) wohldefiniert ist. Wegen $\overline{b} = Pb = G\{(G^{\mathsf{T}}G)^{-1} G^{\mathsf{T}}b\}$ ist $\overline{b} \in \mathcal{R}(G) = \mathcal{Y}$. Da jedes $y \in \mathcal{Y}$ in der Form y = Gw mit $w \in \mathbb{R}^r$ dargestellt werden kann, ergibt sich

$$\boldsymbol{y}^{\mathsf{T}}(\boldsymbol{b}-\bar{\boldsymbol{b}})=(\boldsymbol{G}\boldsymbol{w})^{\mathsf{T}}\left(\boldsymbol{b}-\boldsymbol{P}\boldsymbol{b}
ight)=\boldsymbol{w}^{\mathsf{T}}(\boldsymbol{G}^{\mathsf{T}}\boldsymbol{b}-\boldsymbol{G}^{\mathsf{T}}\boldsymbol{P}\boldsymbol{b})=\boldsymbol{o},$$

also (19). Es bleibt zu zeigen, daß $\overline{\mathbf{b}} = \mathbf{P}\mathbf{b}$ durch (19) eindeutig festgelegt ist. Wenn (19) auch für ein $\overline{\mathbf{c}} \in \mathcal{Y}$ erfüllt ist, folgt $\mathbf{y}^{\mathsf{T}}(\mathbf{b} - \overline{\mathbf{b}}) = \mathbf{y}^{\mathsf{T}}(\mathbf{b} - \overline{\mathbf{c}}) = 0$, also $\mathbf{y}^{\mathsf{T}}(\overline{\mathbf{b}} - \overline{\mathbf{c}}) = 0$ für alle $\mathbf{y} \in \mathcal{Y}$. Wegen $\overline{\mathbf{b}}, \overline{\mathbf{c}} \in \mathcal{Y}$ ist auch $\overline{\mathbf{b}} - \overline{\mathbf{c}} \in \mathcal{Y}$, so daß speziell $\mathbf{y} = \overline{\mathbf{b}} - \overline{\mathbf{c}}$ gewählt werden kann, was auf $(\overline{\mathbf{b}} - \overline{\mathbf{c}})^{\mathsf{T}} (\overline{\mathbf{b}} - \overline{\mathbf{c}}) = 0$, also $\overline{\mathbf{b}} = \overline{\mathbf{c}}$ führt. Daher ist $\overline{\mathbf{b}}$ eindeutig durch (19) festgelegt, so daß auch \mathbf{P} eindeutig durch \mathcal{Y} bestimmt ist und nicht von der speziell gewählten Basis $\{g^1, \ldots, g^r\}$ abhängt.

Es bleibt die Äquivalenz von (i) und (ii) zu zeigen. Nun gilt für beliebiges $m{y} \in m{y}$

$$\begin{split} \|\boldsymbol{b} - \boldsymbol{y}\|^2 &= \|(\boldsymbol{b} - \bar{\boldsymbol{b}}) + (\bar{\boldsymbol{b}} - \boldsymbol{y})\|^2 = \|\boldsymbol{b} - \bar{\boldsymbol{b}}\|^2 + 2(\boldsymbol{b} - \bar{\boldsymbol{b}})^{\mathsf{T}} (\bar{\boldsymbol{b}} - \boldsymbol{y}) + \|\bar{\boldsymbol{b}} - \boldsymbol{y}\|^2 \\ &= \|\boldsymbol{b} - \bar{\boldsymbol{b}}\|^2 + \|\bar{\boldsymbol{b}} - \boldsymbol{y}\|^2, \end{split}$$

denn wegen $\overline{b} - y \in \mathcal{Y}$ verschwindet das Skalarprodukt, siehe (19). Offensichtlich wird dann $\|b - y\|$ über $y \in \mathcal{Y}$ genau dann minimal, wenn $y = \overline{b}$ ist. \Box

Zur Illustration stellen wir die vorkommenden Größen für den Fall m = 2, $\mathcal{Y} = \text{span} \{g^1\}$, dim $\mathcal{Y} = r = 1$, in Abb. 8.1.2 dar. Der Projektor $\mathbf{Q} = \mathbf{I} - \mathbf{P}$ und der Teilraum \mathcal{Y}^{\perp} werden weiter unten eingeführt und können zunächst unbeachtet bleiben.



Abb. 8.1.2. Orthogonale Projektion und zugehörige Größen

Orthogonale Projektoren lassen sich in einfacher Weise charakterisieren.

8.1.3. Aussage. $P \in \mathbf{R}^{m,m}$ ist genau dann ein orthogonaler Projektor, wenn

$$\boldsymbol{P} = \boldsymbol{P}^{\mathsf{T}} \quad \text{und} \quad \boldsymbol{P}^2 = \boldsymbol{P} \tag{23}$$

gilt. Wenn (23) erfüllt ist, stellt P den Projektor auf $\mathcal{Y} = \mathcal{R}(P)$ dar.

Beweis. Es sei P der Projektor auf \mathcal{Y} . Dann kann P nach (22) dargestellt werden, und (23) ist erfüllt, wie man durch Einsetzen bestätigt. Genügt P umgekehrt den Forderungen (23), gilt für jedes $y \in \mathcal{Y} := \mathcal{R}(P)$, d. h. für jedes y = Pv, $v \in \mathbb{R}^m$,

$$\boldsymbol{y}^{\mathsf{T}}(\boldsymbol{b} - \boldsymbol{P}\boldsymbol{b}) = (\boldsymbol{P}\boldsymbol{v})^{\mathsf{T}} (\boldsymbol{b} - \boldsymbol{P}\boldsymbol{b}) = \boldsymbol{v}^{\mathsf{T}}(\boldsymbol{P} - \boldsymbol{P}^2) \boldsymbol{b} = 0,$$

d. h., $\mathbf{\tilde{b}} = \mathbf{Pb}$ ist in der Tat die Projektion von \mathbf{b} auf \mathcal{Y} .

Wenn \mathcal{Y} ein Teilraum von \mathbb{R}^m ist, wird das orthogonale Komplement \mathcal{Y}^{\perp} von \mathcal{Y} als Menge aller zu \mathcal{Y} orthogonalen Vektoren gemäß

$$\mathcal{Y}^{\perp} := \{ \boldsymbol{z} \in \mathbf{R}^{m} : \boldsymbol{z} \perp \mathcal{Y} \} = \{ \boldsymbol{z} \in \mathbf{R}^{m} : \boldsymbol{z}^{\mathsf{T}} \boldsymbol{y} = 0 \text{ für alle } \boldsymbol{y} \in \mathcal{Y} \}$$
(24)

eingeführt. Es ist sofort zu sehen, daß $\mathcal{X} := \mathcal{Y}^{\perp}$ wieder ein Teilraum von \mathbb{R}^{m} ist und durch die Bedingungen

$$\mathbf{R}^m = \mathcal{Y} + \mathcal{X} ext{ und } \mathcal{Y} \perp \mathcal{X}, ext{ d. h. } \mathbf{y}^\intercal \mathbf{z} = 0 ext{ für alle } \mathbf{y} \in \mathcal{Y}, extbf{ z} \in \mathcal{X}, ext{ (25)}$$

eindeutig festgelegt wird, insbesondere gilt $(\mathcal{Y}^{\perp})^{\perp} = \mathcal{Y}$. Die folgende Aussage zeigt, daß zwischen den Projektoren auf \mathcal{Y} und \mathcal{X} ein einfacher Zusammenhang besteht.

8.1.4. Aussage. Es sei \mathcal{Y} ein Teilraum von \mathbb{R}^m , und P bezeichne den orthogonalen Projektor auf \mathcal{Y} . Dann ist

$$\boldsymbol{Q} := \boldsymbol{I} - \boldsymbol{P} \tag{26}$$

der orthogonale Projektor auf $\mathcal{X} := \mathcal{Y}^{\perp}$.

Be we is. Q genügt (23) und ist also ein Projektor. Für $z \in \mathcal{X}$ gilt weiter $z^{\mathsf{T}}(b - Qb) = z^{\mathsf{T}}(b - (I - P)b) = z^{\mathsf{T}}Pb = 0$ für alle b wegen $Pb \in \mathcal{Y}$, d. h., Q projiziert auf \mathcal{X} . \Box

Aus 8.1.4 folgt, daß jedes $x \in \mathbf{R}^m$ gemäß

$$oldsymbol{x} = oldsymbol{y} + oldsymbol{z} = oldsymbol{P}oldsymbol{x} + oldsymbol{Q}oldsymbol{x}, \qquad oldsymbol{y} = oldsymbol{P}oldsymbol{x} \in \mathcal{Y}, \qquad oldsymbol{z} = oldsymbol{Q}oldsymbol{x} \in \mathcal{X}, \qquad (27)$$

in eindeutiger Weise als Summe zweier Vektoren aus \mathcal{Y} und \mathcal{X} dargestellt werden kann, wobei $\|Px + Qx\|^2 = \|Px\|^2 + 2(Px)^{\mathsf{T}}Qx + \|Qx\|^2 = \|Px\|^2 + 2x^{\mathsf{T}}PQx + \|Qx\|^2$, wegen $PQ = P(I - P) = P - P^2 = O$ also

$$\|\boldsymbol{x}\|^2 = \|\boldsymbol{P}\boldsymbol{x}\|^2 + \|\boldsymbol{Q}\boldsymbol{x}\|^2 \tag{28}$$

gilt. Dies ist gerade der Satz des PYTHAGORAS in dem durch x, Px und Qx aufgespannten rechtwinkligen Dreieck, siehe Abb. 8.1.2, wo **b** statt x vorkommt.

C. Weitere Charakterisierung von Quadratmittellösungen

Wir kehren jetzt wieder zum Quadratmittelproblem $Ax \cong b$ zurück. Dieses läßt sich auch in der Form (20) mit $\mathcal{Y} = \mathcal{R}(A)$ schreiben, denn jedes $y \in \mathcal{R}(A)$ ist von der Gestalt $y = Ax, x \in \mathbb{R}^n$. Nach 8.1.2(ii) ist dann x Lösung genau dann, wenn y = Ax

 $= \bar{b} = P_A b$ gilt, wobei P_A den Projektor auf $\mathcal{R}(A)$ bezeichnet. Man rechnet nach, daß $\bar{\boldsymbol{b}}$ aus (15) in der Tat die (mittels der Singulärwertzerlegung bestimmte) Projektion von **b** auf $\mathcal{R}(A)$ ist, siehe Ü 8.1.5. Nach 8.1.2(i) kann $\overline{\mathbf{b}}$ auch durch die Bedingung $\boldsymbol{b} - \boldsymbol{\bar{b}} \perp \mathcal{R}(\boldsymbol{A})$, d. h. durch $\boldsymbol{r} \perp \mathcal{R}(\boldsymbol{A})$ mit $\boldsymbol{r} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}$, charakterisiert werden. Dies ist zu $r \in \mathcal{R}(A)^{\perp}$ äquivalent. Das orthogonale Komplement zu $\mathcal{R}(A)$ läßt sich aber durch A^{\intercal} beschreiben.

8.1.5. Aussage. Für jede Matrix $A \in \mathbb{R}^{n,n}$ gilt

$$\mathscr{R}(A)^{\perp} = \mathscr{N}(A^{\intercal})$$

Beweis. Es gilt $z \in \mathcal{R}(A)^{\perp}$ genau dann, wenn $y^{\mathsf{T}}z = 0$ ist für alle $y \in \mathcal{R}(A)$, d. h. für alle $y = Ax, x \in \mathbb{R}^n$. Dies bedeutet $(Ax)^{\intercal} z = x^{\intercal} A^{\intercal} z = 0$ für alle $x \in \mathbb{R}^n$, also $A^{\intercal} z = o$ oder äquivalent $\boldsymbol{z} \in \mathcal{N}(A^{\mathsf{T}})$.

Die 'Bedingung $r \in \mathcal{R}(A)^{\perp}$ ist also zu $r \in \mathcal{N}(A^{\mathsf{T}})$, d. h. zu $A^{\mathsf{T}}r = A^{\mathsf{T}}(b - Ax) = o$ äquivalent.

Zusammenfassend erhalten wir die folgenden Charakterisierungen der Quadratmittellösungen.

8.1.6. Satz. Die folgenden Aussagen sind äquivalent:

- (i) x ist Lösung des Quadratmittelproblems $Ax \simeq b$.
- (ii) \boldsymbol{x} ist Lösung des Gleichungssystems

$$Ax = P_A b. (29)$$

$$\boldsymbol{r} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x} \perp \boldsymbol{\mathcal{R}}(\boldsymbol{A}). \tag{30}$$

(ii) $\mathbf{x} = \mathbf{P}_A \mathbf{b}$. (iii) \mathbf{x} genügt der Bedingung $\mathbf{r} = \mathbf{b} - A\mathbf{x} \perp \mathcal{R}(A)$. (iv) \mathbf{x} ist Lösung des Gleichungssystems

$$A^{\mathsf{T}}Ax = A^{\mathsf{T}}b. \tag{31}$$

Die Gleichungen (31) werden die zu $Ax \simeq b$ gehörenden Normalgleichungen genannt: sie sind wie das Gleichungssystem (29) auf Grund der speziellen Gestalt der rechten Seite stets konsistent.

D. Normallösung und Pseudoinverse

In 8.1.1 wurde festgestellt, daß die (n - r)-dimensionale lineare Mannigfaltigkeit

$$\mathscr{L}=\mathscr{L}(A,oldsymbol{b})=oldsymbol{x}^{\scriptscriptstyle N}+\mathscr{N}(A)$$

die Lösungsmenge von $Ax \simeq b$ darstellt. Im Fall r < n existieren also unendlich viele Lösungen, und es ist sinnvoll, unter diesen eine spezielle Lösung auszuzeichnen. Wenn die Norm von x als Kriterium gewählt wird, ergibt sich aus der Darstellung (17) wegen der Orthogonalität von V

$$\|\boldsymbol{x}\|^{2} = \|\boldsymbol{\xi}\|^{2} = \|\boldsymbol{\xi}^{N}\|^{2} + \|\boldsymbol{\xi}\|^{2} = \sum_{j=1}^{r} (\beta_{j}/\sigma_{j})^{2} + \sum_{j=r+1}^{n} \xi_{j}^{2}$$
(32)

für jedes $x = V_{2} \in \mathscr{X}$. Hieraus liest man ab, daß $\|x\|$ genau dann minimal wird, wenn

die frei wählbaren Parameter ξ_{r+1}, \ldots, ξ_n gleich 0 gesetzt werden, d. h., wenn $\hat{\boldsymbol{x}} = \boldsymbol{o}$ und folglich $\boldsymbol{x} = \boldsymbol{x}^N$ gilt. Dies ist offensichtlich dann und nur dann der Fall, wenn \boldsymbol{x} orthogonal zu $\mathcal{N}(A)$ ist, siehe Abb. 8.1.3.



Abb. 8.1.3. Lösungsmenge und Normallösung für n = 2, r = 1

8.1.7. Aussage. Es sei \mathcal{I} die Lösungsmenge des Quadratmittelproblems $Ax \simeq b$. Dann gilt

(i) Unter allen Lösungen gibt es genau eine mit kleinster Euklidischer Norm, nämlich x^N , d. h., die Aufgabe

$$|\boldsymbol{x}|| \to \underset{\boldsymbol{x} \in \boldsymbol{\mathscr{I}}}{\operatorname{Minimum}!}$$
 (33)

ist eindeutig durch x^N lösbar.

(ii) Für $\boldsymbol{x} \in \mathscr{L}$ gilt $\boldsymbol{x} = \boldsymbol{x}^N$ genau dann, wenn

$$\boldsymbol{x} \perp \mathcal{N}(\boldsymbol{A})$$
. (34)

Wegen (i) heißt x^N Minimum-Norm-Lösung, wegen (ii) auch Normallösung von $Ax \simeq b$. Ihrer Kürze wegen verwenden wir die letztgenannte Bezeichnung. Für spaltenreguläre Probleme ist x^N die einzige Lösung von $Ax \simeq b$.

Wir wollen jetzt untersuchen, wie x^N von b abhängt. Dazu benutzen wir wieder die Singulärwertzerlegung (2) und die Darstellung $x^N = V \xi^N$ mit ξ^N aus (10), vgl. 8.1.1. Wenn

$$\boldsymbol{\Sigma}^{+} := \begin{pmatrix} 1/\sigma_{1} & 0\\ \ddots & \\ 1/\sigma_{r} & \\ \hline & 0 & 0 \end{pmatrix} \in \mathbf{R}^{n,m} \quad \text{und} \quad \boldsymbol{A}^{+} := \boldsymbol{V}\boldsymbol{\Sigma}^{+}\boldsymbol{U}^{\mathsf{T}} \in \mathbf{R}^{n,m}$$
(35)

gesetzt wird, ergibt sich unter Beachtung von $\beta = U^{\mathsf{T}} b$

$$\boldsymbol{x}^{N} = \boldsymbol{V}\boldsymbol{\Sigma}^{+}\boldsymbol{U}^{\mathsf{T}}\boldsymbol{b} = \boldsymbol{A}^{+}\boldsymbol{b}, \qquad (36)$$

d. h., die Abbildung $\mathbf{b} \to \mathbf{x}^N$ ist linear und wird durch die Matrix A^+ beschrieben. Der folgende Satz zeigt, daß A^+ auch in völlig anderer Weise eingeführt werden kann.

8.1.8. Satz.

(i) Zu jeder Matrix $A \in \mathbb{R}^{m,n}$ gibt es genau eine Matrix $A^+ \in \mathbb{R}^{n,m}$, so daß

 $x^N = A^+ b$ für alle $b \in \mathbf{R}^m$

die Normallösung von $Ax \simeq b$ darstellt. Wenn (2) die Singulärwertzerlegung von A bezeichnet, ist A^+ durch (35) gegeben.

(ii) A^+ ist durch die Bedingungen

(P₁)
$$AA^+A = A$$
, (P₂) $A^+AA^+ = A^+$,

$$(P_3) (AA^+)^{\mathsf{T}} = AA^+, \quad (P_4) (A^+A)^{\mathsf{T}} = A^+A$$

eindeutig festgelegt.

Beweis. Die Eindeutigkeit von A⁺ ergibt sich aus der Eindeutigkeit der Normallösung. Es bleibt zu zeigen, daß A^+ äquivalent durch $(P_1), \ldots, (P_4)$ charakterisiert werden kann. Daß A^+ aus (35) diese Bedingungen erfüllt, läßt sich sofort durch Einsetzen überprüfen. Es sei umgekehrt jetzt $B \in \mathbf{R}^{n,m}$ eine Matrix, die den Bedingungen aus (ii) mit B anstelle von A^+ genügt. Dann folgt

$$\begin{split} \mathbf{B} &= \mathbf{B}\mathbf{A}\mathbf{B} = \mathbf{B}\mathbf{A}\mathbf{A}^{\mathsf{+}}\mathbf{A}\mathbf{B} = \mathbf{A}^{\mathsf{\top}}\mathbf{B}^{\mathsf{\top}}\mathbf{A}^{\mathsf{+}}\mathbf{A} = \mathbf{A}^{\mathsf{\top}}\mathbf{B}^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}}(\mathbf{A}^{\mathsf{+}})^{\mathsf{\top}}\mathbf{B} = \mathbf{A}^{\mathsf{+}}\mathbf{A}\mathbf{B} \\ & \begin{array}{c} 2 & 1 & 4 & 1 \\ \end{array} \\ & = \mathbf{A}^{\mathsf{+}}\mathbf{B}^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}\mathbf{A}^{\mathsf{+}}\mathbf{B}^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}} = \mathbf{A}^{\mathsf{+}}(\mathbf{A}^{\mathsf{+}})^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}}\mathbf{B}^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}} = \mathbf{A}^{\mathsf{+}}(\mathbf{A}^{\mathsf{+}})^{\mathsf{\top}}\mathbf{A}^{\mathsf{\top}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}\mathbf{A}^{\mathsf{+}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+}} = \mathbf{A}^{\mathsf{+}}\mathbf{A}^{\mathsf{+$$

Die Zahl i unter dem Gleichheitszeichen gibt dabei an, welche Bedingung (P.) für A^+ bzw. **B** benutzt worden ist.

Die Matrix A heißt Moore-Penrose-Inverse oder kurz Pseudoinverse von A; man spricht auch von der verallgemeinerten Inversen. Die Bedingungen $(P_1), \ldots, (P_4)$ werden Penrose-Bedingungen genannt. Zur Geschichte der verallgemeinerten Inversen siehe B 8.3.

Durch Einsetzen überzeugt man sich von der Gültigkeit der Rechenregeln

$$O^+ = O$$
, $(A^+)^{\mathsf{T}} = (A^{\mathsf{T}})^+$ und $(\lambda A)^+ = (1/\lambda) A^+$ für $\lambda \in \mathsf{R}, \ \lambda \neq 0$.
(37)

Weitere nützliche Eigenschaften der Pseudoinversen, deren Richtigkeit ebenfalls durch Einsetzen in die Penrose-Bedingungen bestätigt werden kann, sind in der folgenden Aussage zusammengefaßt.

8.1.9. Aussage. Es sei $A \in \mathbb{R}^{m,n}$ eine Matrix vom Rang r. Dann gilt

(i)

$$A^{+} = \begin{cases} (A^{\mathsf{T}}A)^{-1} A^{\mathsf{T}}, & \text{falls} \quad n = r \leq m, \\ A^{\mathsf{T}}(AA^{\mathsf{T}})^{-1}, & \text{falls} \quad m = r \leq n, \\ A^{-1}, & \text{falls} \quad m = r = n. \end{cases}$$

- $\begin{bmatrix} A^{-1}, & \text{falls } m = r = n. \\ \text{(ii) } A^+ = G^+ F^+, \text{ falls } A = FG \text{ mit } F \in \mathbb{R}^{m,r}, \ G \in \mathbb{R}^{r,n} \text{ (und folglich rang } (F) \\ = \text{rang } (G) = r). \\ \text{(iii) } A^+ = VB^+U^{\intercal}, \text{ falls } A = UBV^{\intercal} \text{ mit } B \in \mathbb{R}^{m,n} \text{ und orthogonalen Matrizen}$

 (\mathbf{n})

(iv)
$$A^+ = (\underbrace{B^+}_k | \underbrace{O}_{m-k})^{n}$$
, falls $A = (\underbrace{B}_{O})^{k}_{m-k}$.
(v) $[\operatorname{diag}(\sigma_1, \dots, \sigma_l)]^+ = \operatorname{diag}(\sigma_1^+, \dots, \sigma_l^+)$, wobei
 $\sigma^+ := \begin{cases} 1/\sigma & \operatorname{für} & \sigma \neq 0, \\ 0 & \operatorname{für} & \sigma = 0, \end{cases} \quad \sigma \in \mathbf{R}.$

Schließlich lassen sich die mit A verbundenen Projektoren explizit durch A und A^+ ausdrücken.

8.1.10. Aussage. Es sei $A \in \mathbb{R}^{m,n}$ gegeben. Dann ist

$$\boldsymbol{P}_{\boldsymbol{A}} = \boldsymbol{A}\boldsymbol{A}^{-}$$
 bzw. $\boldsymbol{P}_{\boldsymbol{A}\mathsf{T}} = \boldsymbol{A}^{+}\boldsymbol{A}$

$$P_A = AA^-$$
 bzw. $P_{A^{\mathsf{T}}} = A^+A$
der Projektor auf $\mathcal{R}(A)$ bzw. $\mathcal{R}(A^{\mathsf{T}})$, und
 $Q_A = I - AA^+$ bzw. $Q_{A^{\mathsf{T}}} = I - A^+A$
stellt den Projektor auf $\mathcal{N}(A^{\mathsf{T}})$ bzw. $\mathcal{N}(A)$ dar.

Beweis. Wir setzen $P := P_A = AA^+$. Wegen (P₃) und (P₁) gilt (23), d. h., P ist ein Projektor. Es bleibt $\mathcal{R}(\mathbf{P}) = \mathcal{R}(\mathbf{A})$ zu zeigen. Für jedes \mathbf{y} gilt nun $\mathbf{P}\mathbf{y} = \mathbf{A}\mathbf{A}^+\mathbf{y}$, also $\mathcal{R}(\mathbf{P}) \subset \mathcal{R}(\mathbf{A})$. And ererse its ist wegen (P₁) $Ax = AA^{+}Ax = PAx$, also $\mathcal{R}(A) \subset \mathcal{R}(P)$ und damit $\mathcal{R}(P)$ $= \mathcal{R}(A)$. Die übrigen Aussagen ergeben sich aus 8.1.4 und der Tatsache, daß $\mathcal{R}(A)^{\perp} = \mathcal{N}(A^{\mathsf{T}})$ sowie $\mathcal{R}(\mathbf{A}^{\mathsf{T}})^{\perp} = \mathcal{N}(\mathbf{A})$ gilt; letztere Eigenschaft folgt bei Anwendung von 8.1.5 auf \mathbf{A}^{T} .

Übungsaufgaben

U 8.1.1. Man zeige: Wenn $\{q^1, \ldots, q^r\}$ eine orthonormierte Basis von \mathcal{Y} ist, hat der orthonormale Projektor \boldsymbol{P} auf $\boldsymbol{\mathcal{Y}}$ die Gestalt

$$\boldsymbol{P} = \sum_{i=1}^{r} \boldsymbol{q}^{i} \boldsymbol{q}^{i\top}.$$
(38)

Der Projektor Q = I - P kann in der Produktform

$$Q = I - P = [I - q^1 q^{1\mathsf{T}}] [I - q^2 q^{2\mathsf{T}}] \cdots [I - q^{\mathsf{T}} q^{\mathsf{T}}]$$
(39)

geschrieben werden.

Ü 8.1.2. Man beweise, daß $P \in \mathbb{R}^{m,m}$ genau dann ein orthogonaler Projektor ist, wenn

$$\boldsymbol{P} = \boldsymbol{U}\boldsymbol{\Pi}\boldsymbol{U}^{\mathsf{T}} \tag{40}$$

mit orthogonalem $U \in \mathbb{R}^{m,m}$ und einer Diagonalmatrix $II = \text{diag}(\Pi_i), \Pi_i \in \{0, 1\}, \text{ gilt, d. h.},$ Projektoren sind symmetrische Matrizen mit den Eigenwerten 0 oder 1. Man gebe eine orthonormierte Basis von $\mathcal{R}(\mathbf{P})$ an.

Ü 8.1.3. Es sei $\mathscr{Y} \subset \mathsf{R}^m$ ein Teilraum und $\mathscr{Z} = \mathscr{Y}^{\perp}$. Man beweise, daß $P \in \mathsf{R}^{m,m}$ genau dann der orthonormale Projektor auf \mathcal{Y} ist, wenn

$$Py = y$$
 für alle $y \in \mathcal{Y}$ und $Pz = o$ für alle $z \in \mathcal{Z}$ (41)

gilt.

Hin weis: Es gilt $P = P^{\mathsf{T}}$ genau dann, wenn $u^{\mathsf{T}}(Pv) = (Pu)^{\mathsf{T}} v$ ist für alle $u, v \in \mathsf{R}^{\mathsf{m}}$.

Ü 8.1.4. Man zeige: Für jeden orthogonalen Projektor P gilt

$$\|\boldsymbol{P}\| = \begin{cases} 1, & \text{falls } \boldsymbol{P} \neq \boldsymbol{O}, \\ 0, & \text{falls } \boldsymbol{P} = \boldsymbol{O}, \end{cases}$$
(42)

sowie

$$\|\boldsymbol{P}\|_{\boldsymbol{F}} = \sqrt{r}, \qquad r = \operatorname{rang}\left(\boldsymbol{P}\right) = \dim \mathcal{R}(\boldsymbol{P}). \tag{43}$$

Ü 8.1.5. Man drücke die Projektoren P_A , Q_A , $P_{A^{\mathsf{T}}}$, $Q_{A^{\mathsf{T}}}$ aus 8.1.10 unter Verwendung der Singulärwertzerlegung durch U, V und Σ aus.

Ü 8.1.6. Man beweise

$$\|A^+\| = \|\Sigma^+\| = 1/\sigma_r.$$
(44)

Ü 8.1.7. Man zeige: Wenn $v \in \mathbf{R}^m$ als (m, 1)-Matrix aufgefaßt wird, gilt im Fall $v \neq o$

$$\boldsymbol{v}^{+} = \boldsymbol{v}^{\mathsf{T}} / \|\boldsymbol{v}\|^{2} \quad \text{sowie} \quad \boldsymbol{P}_{\boldsymbol{v}} = \boldsymbol{v} \boldsymbol{v}^{\mathsf{T}} / \|\boldsymbol{v}\|^{2}, \tag{45}$$

wobei P_v den Projektor auf $\mathcal{R}(v) = \text{span } \{v\}$ bezeichnet.

Ü 8.1.8. Wie lautet $(ab^{\intercal})^+$ für $a, b \in \mathbb{R}^m$. $a \neq o, b \neq o$?

Ü 8.1.9. Man zeige, daß die Vektoren $P_A b$, $A^- b$ und $A^{\mathsf{T}} b$ entweder alle verschwinden oder alle von 0 verschieden sind.

Ü 8.1.10. Man beweise: Unter allen Lösungen $x \in \mathcal{I}$ von $Ax \cong b$ hat

$$ar{oldsymbol{x}}:=oldsymbol{x}^{0}+A^{+}(oldsymbol{b}-Aoldsymbol{x}^{0})=(oldsymbol{I}-A^{+}A)\,oldsymbol{x}^{0}+A^{+}oldsymbol{b}$$

den kleinsten Euklidischen Abstand von x^0 . Skizze analog zu Abb. 8.1.3! Hinweis: Man führe $y := x - x^0$ als neue Variable ein.

8.2. Störungstheorie

Wir betrachten das lineare Quadratmittelproblem

$$Ax \simeq b, \tag{1}$$

dessen eindeutige Normallösung durch

$$\boldsymbol{x} = \boldsymbol{x}^N = \boldsymbol{A}^+ \boldsymbol{b} \tag{2}$$

gegeben ist. Im Fall einer spaltenregulären Matrix stellt (2) die einzige Lösung von (1) dar, andernfalls diejenige kleinster Euklidischer Norm.

Uns interessiert, wie sich Störungen δA , δb der Eingangsdaten auf die Normallösung auswirken. Das gestörte Problem

$$(A + \delta A) (x + \delta x) \cong b + \delta b \tag{3}$$

besitzt die Normallösung

$$\boldsymbol{x} + \boldsymbol{\delta}\boldsymbol{x} = (\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A})^{+} (\boldsymbol{b} + \boldsymbol{\delta}\boldsymbol{b}), \qquad (4)$$

und die Änderung δx gegenüber der Normallösung x des ungestörten Problems ist

$$\delta \boldsymbol{x} = [(\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A})^{+} - \boldsymbol{A}^{+}] \boldsymbol{b} + (\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A})^{+} \boldsymbol{\delta} \boldsymbol{b}.$$
⁽⁵⁾

Wir untersuchen daher zunächst, wie sich A^+ gegenüber Störungen von A verhält.

A. Störungstheorie der Pseudoinversen

Um zu zeigen, wie sich Störungen von A in unterschiedlicher Weise auf A^+ auswirken können, betrachten wir zwei einfache Beispiele.

8.2.1. Beispiel. Die Matrix $A \neq O$ sei durch ihre Singulärwertzerlegung

$$A = U \Sigma V^{\intercal}$$

mit $\Sigma = \text{diag}(\sigma_1, ..., \sigma_r, 0, ..., 0)$, $\sigma_1 \ge \cdots \ge \sigma_r > 0$, r = rang(A) und orthogonalem U, V gegeben. Die Störung δA werde in der Form

$$\delta A = U \delta \Sigma V^{\intercal}$$

geschrieben, so daß die gestörte Matrix $A + \delta A$ die Gestalt

$$A + \delta A = U(\Sigma + \delta \Sigma) V^{\intercal}$$

besitzt.

Für die durch

$$(\boldsymbol{\delta \Sigma})_{ij} = \begin{cases} \varepsilon & \text{für } i = j = r, \\ 0 & \text{sonst} \end{cases}$$
(6)

definierte spezielle Störung $\delta \Sigma$ ist

$$\mathbf{\Sigma} = egin{bmatrix} \sigma_1 & & & \ \sigma_{r-1} & & \ \sigma_r & & \ \hline & \mathbf{0} & & \ \end{bmatrix}, \quad \mathbf{\Sigma} + \mathbf{\delta} \mathbf{\Sigma} = egin{bmatrix} \sigma_1 & & & \ \sigma_{r-1} & & \ \sigma_r + \mathbf{\varepsilon} & \ \hline & \mathbf{0} & & \ \end{bmatrix}.$$

Im Fall

 $arkappa:=\|A^{\scriptscriptstyle +}\|\,\|{\mathfrak o} A\|=|arepsilon|/\sigma_r<1$

- man beachte $\|\delta A\| = \|\delta \Sigma\| = |\varepsilon|$ und $\|A^+\| = 1/\sigma_r$ - folgt

$$\sigma_r + \varepsilon \ge \sigma_r - |\varepsilon| = \sigma_r (1 - \varkappa) > 0, \tag{7}$$

d. h., $\sigma_r + \varepsilon$ bleibt positiv, und es gilt

rang
$$(A + \delta A) = \operatorname{rang} (\Sigma + \delta \Sigma) = \operatorname{rang} (\Sigma) = \operatorname{rang} (A)$$
.

Die Pseudoinversen sind dann durch

$$A^+ = V \Sigma^+ U^\intercal$$
 bzw. $(A + \delta A)^+ = V (\Sigma + \delta \Sigma)^+ U^\intercal$

(8)

mit

$$\Sigma^{+} = \begin{pmatrix} \frac{1}{\sigma_{1}} & & \\ & \frac{1}{\sigma_{r-1}} & 0 \\ & \frac{1}{\sigma_{r}} & \\ & 0 & 0 \end{pmatrix} \quad \text{bzw.} \quad (\Sigma + \delta \Sigma)^{+} = \begin{pmatrix} \frac{1}{\sigma_{1}} & & \\ & \frac{1}{\sigma_{r-1}} & 0 \\ & \frac{1}{\sigma_{r+\epsilon}} & 0 \\ & \frac{1}{\sigma_{r+\epsilon}} & 0 \\ & 0 & 0 \end{pmatrix}$$

gegeben, vgl. (8.1.35). Wegen (7) gilt

$$\|(\boldsymbol{\Sigma} + \boldsymbol{\delta}\boldsymbol{\Sigma})^+\| = \max\left\{\frac{1}{\sigma_1}, \frac{1}{\sigma_{r-1}}, \frac{1}{\sigma_r + \varepsilon}\right\} \leq \frac{1}{\sigma_r - |\varepsilon|} = \frac{1}{\sigma_r(1-\varepsilon)},$$

also

d. h.

$$\|(A+\delta A)^+\|=\|(\Sigma+\delta \Sigma)^+\|\leq \|A^+\|/(1-arkappa).$$

Analog folgt

$$\|(\boldsymbol{\Sigma} + \boldsymbol{\delta}\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+\| = \left|\frac{1}{\sigma_r + \varepsilon} - \frac{1}{\sigma_r}\right| = rac{|\varepsilon|}{\sigma_r(\sigma_r + \varepsilon)} \leq rac{|\varepsilon|}{\sigma_r^2(1 - \varkappa)},$$

 $\|(A + \delta A)^{+} - A^{+}\| \leq [\|A^{+}\|^{2}/(1 - \varkappa)] \|\delta A\|. \quad \Box$ (9)

Die Abschätzungen (8), (9) entsprechen wörtlich den Abschätzungen (4.1.9), (4.1.10) für die Inversen gestörter regulärer Matrizen. Wie bei der Matrixinversion liegt hier also bei der Störung (6), die den Rang von A nicht erhöht, lokale Lipschitzstetigkeit vor; insbesondere bleiben die Pseudoinversen beschränkt.

Das zweite Beispiel zeigt, daß auch bösartigere Störungen dA existieren.

8.2.2. Beispiel. Es sei A wie in Beispiel 8.2.1, und zusätzlich gelte

 $r = \operatorname{rang} (A) < l := \min (m, n),$

d. h., A sei rangdefizient. Die Störung $\delta \Sigma$ sei durch

$$(\boldsymbol{\delta}\boldsymbol{\Sigma})_{ij} = \begin{cases} \varepsilon & \text{für} \quad i = j = r+1, \\ 0 & \text{sonst} \end{cases}$$
(10)

,

.

festgelegt. Im Fall $|\varepsilon| = \|\delta \Sigma\| = \|\delta A\| \neq 0$ ergibt sich dann
also

$$\operatorname{rang} (\boldsymbol{A} + \boldsymbol{dA}) = \operatorname{rang} (\boldsymbol{\Sigma} + \boldsymbol{d\Sigma}) = r + 1 > r = \operatorname{rang} (\boldsymbol{\Sigma}) = \operatorname{rang} (\boldsymbol{A}),$$

und es folgt

$$\|(\boldsymbol{\Sigma} + \boldsymbol{\delta}\boldsymbol{\Sigma})^+\| = \max\{1/\sigma_r, 1/|\boldsymbol{\varepsilon}|\} \geq 1/|\boldsymbol{\varepsilon}|, \quad \|(\boldsymbol{\Sigma} + \boldsymbol{\delta}\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+\| = 1/|\boldsymbol{\varepsilon}|,$$

also

$$\|(A + \delta A)^+\| \ge 1/\|\delta A\|, \quad \|(A + \delta A)^+ - A^+\| \ge 1/\|\delta A\|.$$

Gegenüber der rangerhöhenden Störung (10) ist die Pseudoinverse also unstetig und unbeschränkt, und zwar um so stärker, je kleiner die Störung ist. Es liegt also ein inkorrekt gestelltes Problem mit extrem schlechtem Unstetigkeitsverhalten vor, vgl. Abschnitt 2.1.D und das Modellbeispiel 2.1.11.

Im folgenden werden wir zeigen, daß es genau die in den beiden Beispielen angegebenen Klassen von gutartigen und bösartigen Störungen gibt. Die gutartigen Störungen sind durch die Bedingung rang $(A + \delta A) \leq \text{rang}(A)$ charakterisiert und führen zu lokal lipschitzstetiger Abhängigkeit. Die bösartigen Störungen sind diejenigen mit rang $(A + \delta A) > \text{rang}(A)$. Gegenüber solchen Störungen liegt Unstetigkeit vor, und die Pseudoinversen sind nicht beschränkt.

Als Hilfsmittel benötigen wir Aussagen über das Störungsverhalten der Singulärwerte.

8.2.3. Satz. Es gelte $A, \, \delta A \in \mathbb{R}^{m,n}$, und σ_i bzw. $\sigma_i + \delta \sigma_i$ seien die gemäß

$$\sigma_1 \geqq \sigma_2 \geqq \cdots \geqq \sigma_l$$

bzw.

$$\sigma_1 + \delta \sigma_1 \geq \sigma_2 + \delta \sigma_2 \geq \cdots \geq \sigma_l + \delta \sigma_l$$

geordneten Singulärwerte von A bzw. $A + \delta A$, wobei $l := \min(m, n)$ ist. Dann gelten die Abschätzungen

$$|\delta\sigma|_i \le \|\delta A\| \tag{11}$$

sowie

$$\left\{\sum_{i=1}^{l} (\delta\sigma_i)^2\right\}^{1/2} \leq \|\boldsymbol{\delta}\boldsymbol{A}\|_F.$$
(12)

Beweis. Wir führen die Matrizen

$$\boldsymbol{B} := \begin{pmatrix} \boldsymbol{O} \mid \boldsymbol{A} \\ \boldsymbol{A}^{\mathsf{T}} \mid \boldsymbol{O} \end{pmatrix} \quad \text{und} \quad \boldsymbol{\delta} \boldsymbol{B} := \begin{pmatrix} \boldsymbol{O} \mid \boldsymbol{\delta} \boldsymbol{A} \\ \boldsymbol{\delta} \boldsymbol{A}^{\mathsf{T}} \mid \boldsymbol{O} \end{pmatrix}$$

ein. Unter Verwendung der Singulärwertzerlegung von A folgt

$$B = \left(\frac{O \mid U\Sigma V^{\mathsf{T}}}{V\Sigma^{\mathsf{T}}U^{\mathsf{T}} \mid O} \right) = SCS^{\mathsf{T}} \quad \text{mit} \quad C := \left(\frac{O \mid \Sigma}{\Sigma^{\mathsf{T}} \mid O} \right), \quad S := \left(\frac{U \mid O}{O \mid V} \right).$$

Wegen der Orthogonalität von S sind C und B ähnlich, besitzen also dieselben Eigenwerte.

Die Eigenwertgleichung $Cz = \lambda z$ läßt sich mit

$$oldsymbol{z} = egin{pmatrix} oldsymbol{y} \ oldsymbol{x} \end{pmatrix} ext{ als } egin{pmatrix} oldsymbol{\Sigma}^{\intercal}oldsymbol{y} \ oldsymbol{\Sigma}^{\intercal}oldsymbol{y} \end{pmatrix} = \lambda egin{pmatrix} oldsymbol{y} \ oldsymbol{x} \end{pmatrix}$$

schreiben, d. h., sie zerfällt in die |m - n| Gleichungen

$$0 = \lambda y_i$$
 $(i = l + 1, ..., m)$ oder $0 = \lambda x_i$ $(i = l + 1, ..., n)$

und die l zweidimensionalen Eigenwertprobleme

$$\begin{pmatrix} 0 & \sigma_i \\ \sigma_i & 0 \end{pmatrix} \begin{pmatrix} y_i \\ x_i \end{pmatrix} \lambda = \begin{pmatrix} y_i \\ x_i \end{pmatrix} \quad (i = 1, ..., l)$$

von denen jedes die Eigenwerte $\pm \sigma_i$ besitzt. **B** hat also die 2*l* Eigenwerte $\pm \sigma_i$, ergänzt durch |m - n| Nullen, und analoge Aussagen gelten mit $\pm (\sigma_i + \delta \sigma_i)$ für $\mathbf{B} + \delta \mathbf{B}$. Nun gilt $||\mathbf{B}|| = ||\mathbf{A}||$, $||\mathbf{B}||_F = \sqrt{2} ||\mathbf{A}||_F$, und analoge Beziehungen gelten für $\delta \mathbf{B}$ und $\delta \mathbf{A}$. Damit ergeben sich (11) und (12) direkt aus 13.1.B. \Box

Satz 8.2.3 besagt: Die Singulärwerte sind lipschitzstetige Funktionen der Matrix, und die Lipschitzkonstante hat den Wert 1. Es liegt also ein besonders gutartiges Stetigkeitsverhalten vor.

Wir können jetzt die angekündigte Aussage über rangerhöhende Störungen formulieren.

8.2.4. Satz. Es sei $A \in \mathbb{R}^{m,n}$ eine rangdefiziente Matrix und $\delta A \neq O$ eine Störung mit

$$\operatorname{rang}\left(A + \delta A\right) > \operatorname{rang}\left(A\right). \tag{13}$$

Dann gilt

$$\|(A + \delta A)^{+}\| \ge 1/\|\delta A\|.$$
(14)

Beweis. Es sei $r = \operatorname{rang}(A)$, und σ_i bzw. $\tilde{\sigma}_i = \sigma_i + \delta\sigma_i$ seien die wie in 8.2.3 geordneten Singulärwerte von A bzw. $A + \delta A$. Wegen (13) ist rang $(A + \delta A) \ge r + 1$, also $\tilde{\sigma}_{r+1} = \sigma_{r+1} + \delta\sigma_{r+1} = \delta\sigma_{r+1} > 0$. Nach (11) ist andererseits $|\delta\sigma_{r+1}| \le ||\delta A||$, womit sich sofort

$$\|(A + \delta A)^+\| = \max\{1/\tilde{\sigma}_i : \tilde{\sigma}_i > 0\} \ge 1/\tilde{\sigma}_{r+1} \ge 1/\|\delta A\|$$

ergibt. 🗌

Aus (14) folgt

$$\|(A + \delta A)^{+} - A^{+}\| \ge \|(A + \delta A)^{+}\| - \|A^{+}\| \ge (1/\|\delta A\|) - \|A^{+}\|, \quad (15)$$

d. h., $(A + \delta A)^+$ weicht von A^+ beliebig viel ab, wenn δA genügend klein ist.

Wir zeigen jetzt, daß Störungen, die den Rang nicht erhöhen, stets gutartig sind.

8.2.5. Satz. Es sei $A \in \mathbb{R}^{m,n}$ gegeben, und die Störung $\mathcal{O}A \in \mathbb{R}^{m,n}$ genüge den Bedingungen

$$\operatorname{rang}\left(A + \delta A\right) \leq \operatorname{rang}\left(A\right) \tag{16}$$

und

$$\varkappa := \|A^+\| \|\delta A\| < 1.$$
⁽¹⁷⁾

Dann ist

$$\operatorname{rang}\left(A + dA\right) = \operatorname{rang}\left(A\right),\tag{18}$$

und es gilt

$$|(A + \delta A)^{+}|| \leq ||A^{+}||/(1 - \varkappa).$$
⁽¹⁹⁾

Die Änderung $(A + \delta A)^+ - A^+$ ist in erster Ordnung durch

$$(A + \delta A)^{+} - A^{+} = -A^{+}\delta A A^{+} + A^{+}A^{+\intercal}\delta A^{\intercal}Q_{A} + Q_{A}^{\intercal}\delta A^{\intercal}A^{+\intercal}A^{+} + O(\|\delta A\|^{2})$$

$$(20)$$

mit

$$Q_A := I - AA^+$$
 und $Q_{A^{\mathsf{T}}} := I - A^+A$ (21)

gegeben und genügt der Abschätzung

$$\|(A + \delta A)^{+} - A^{+}\| \leq \mu [\|A^{+}\|^{2}/(1 - \varkappa)] \|\delta A\|,$$
(22)

wobei

$$\mu := egin{cases} (1+\sqrt{5})/2 = 1.62, & ext{falls} & ext{rang}\,(A) < \min\,(m,\,n), \ \sqrt{2} = 1.41, & ext{falls} & ext{rang}\,(A) = \min\,(m,\,n) < \max\,(m,\,n), \ 1, & ext{falls} & ext{rang}\,(A) = m = n. \end{cases}$$

Beweis. Es sei $r = \operatorname{rang}(A)$, und σ_i bzw. $\tilde{\sigma}_i = \sigma_i + \delta \sigma_i$ seien die geordneten Singulärwerte von A bzw. $A + \delta A$. Wegen (11) und (17) gilt

$$ilde{\sigma}_i = \sigma_i + \delta \sigma_i \geq \sigma_i - |\delta \sigma_i| \geq \sigma_r - \|\delta A\| = (1 - arkappa) / \|A^+\| > 0 \qquad (i = 1, ..., r),$$

also rang $(A + \delta A) \ge r$. Mit (16) folgt dann $\tilde{\sigma}_i = 0$ $(i = r + 1, ..., \min(m, n))$, mithin (18) und (19). Als nächstes untersuchen wir die Änderung

$$G:= ilde{A}^+-A^+, \qquad ilde{A}^+:=(A+\delta A)^+$$

Mit den Projektoren

$$P := P_A = AA^+, \quad S := P_{A^{\mathsf{T}}} = A^+A, \quad \hat{P} := \tilde{A}\tilde{A}^+, \quad \tilde{S} := \tilde{A}^+\tilde{A}$$

ergibt sich

$$\begin{split} G &= [\tilde{S} + (I - \tilde{S})] \left(\tilde{A}^+ - A^+ \right) [P + (I - P)] \\ &= \tilde{S} \tilde{A}^+ P + \tilde{S} \tilde{A}^+ (I - P) - \tilde{S} A^+ P - \tilde{S} A^+ (I - P) \\ &+ (I - \tilde{S}) \tilde{A}^+ P + (I - \tilde{S}) \tilde{A}^+ (I - P) - (I - \tilde{S}) A^+ P - (I - \tilde{S}) A^+ (I - P). \end{split}$$

Wegen $\tilde{S}\tilde{A}^+ = \tilde{A}^+, A^+P = A^+$ vereinfacht sich diese Darstellung zu

$$G = [\tilde{A}^{+}P - \tilde{S}A^{+}] + [\tilde{A}^{+}(I - P)] + [-(I - \tilde{S})A^{+}] =: G_{1} + G_{2} + G_{3}.$$
(24)

Unter Berücksichtigung von $A^{\mathsf{T}}(I-P) = O$, $(I-\tilde{S}) \tilde{A}^{\mathsf{T}} = O$ folgt schließlich

$$G_{1} = \tilde{A}^{+}AA^{+} - \tilde{A}^{+}\tilde{A}A^{+} = -\tilde{A}^{+}\delta AA^{+},$$

$$G_{2} = \tilde{A}^{+}(I - P) = \tilde{A}^{+}\tilde{P}(I - P) = \tilde{A}^{+}\tilde{A}^{+}\mathsf{T}\tilde{A}^{\intercal}(I - P)$$

$$= \tilde{A}^{+}\tilde{A}^{+}\mathsf{T}(\tilde{A}^{\intercal} - A^{\intercal}) (I - R) = \tilde{A}^{+}\tilde{A}^{+}\mathsf{T}\delta A^{\intercal}(I - P),$$

$$G_{3} = -(I - \tilde{S}) A^{+} = -(I - \tilde{S}) A^{+}AA^{+} = -(I - \tilde{S}) A^{\intercal}A^{+}A^{+}$$

$$= -(I - \tilde{S}) (A^{\intercal} - \tilde{A}^{\intercal}) A^{+}TA^{+} = (I - \tilde{S}) \delta A^{\intercal}A^{+}TA^{+}.$$
(25)

Zur Abschätzung der G_i ziehen wir (19) heran und erhalten sofort

$$\|G_1\| \le \|\tilde{A}^+\| \|A^+\| \|\delta A\| \le [\|A^+\|^2/(1-\varkappa)] \|\delta A\| =: \alpha$$

und

$$\|G_3\| \le \|A^+\|^2 \|\delta A\| \le \alpha.$$

$$(26)$$

Um den Faktor $\|\widetilde{A}^+\|^2$ bei der Abschätzung von G_2 zu vermeiden, verwenden wir Ü 8.2.1. Mit dieser folgt

$$\|\boldsymbol{G}_{2}\| = \|\tilde{\boldsymbol{A}}^{+}\tilde{\boldsymbol{P}}(\boldsymbol{I}-\boldsymbol{P})\| \leq \|\tilde{\boldsymbol{A}}^{+}\| \|\tilde{\boldsymbol{P}}(\boldsymbol{I}-\boldsymbol{P})\| = \|\tilde{\boldsymbol{A}}^{+}\| \|\boldsymbol{P}(\boldsymbol{I}-\tilde{\boldsymbol{P}})\| \leq \alpha,$$
(27)

denn wegen $\tilde{A}^{\mathsf{T}}(I - \tilde{P}) = O$ gilt

$$\begin{split} P(I - \tilde{P}) &= AA^{+}(I - \tilde{P}) = A^{+\intercal}A^{\intercal}(I - \tilde{P}) = A^{+\intercal}(A^{\intercal} - \tilde{A}^{\intercal}) (I - \tilde{P}) \\ &= -A^{+\intercal}\delta A^{\intercal}(I - \tilde{P}). \end{split}$$

Zusammenfassend ergibt sich $||G|| \leq 3\alpha$, also (22) mit $\mu = 3$. Die verbesserten Werte (23) für μ erhält man durch elementare geometrische Überlegungen unter Ausnutzung von Orthogonalitätsbeziehungen, siehe Ü 8.2.2. Aus (22) folgt speziell

$$(A + \delta A)^+ = A^+ + O(||\delta A||),$$

so da $\beta \, ilde{A}^+$ in (24), (25) durch A^+ ersetzt werden kann, um zur linearisierten Änderungsformel (20) zu gelangen.

8.2.6. Bemerkung. (i) Wie die Beispiele 8.2.1 und 8.2.2 zeigen, kann auf keine der Voraussetzungen (16) und (17) verzichtet werden, um die Beschränktheit von $(A + dA)^+$ zu sichern. Außerdem zeigen sie, daß die Abschätzungen (14) und (19) scharf sind, d. h., zu jedem ΔA mit $\|\Delta A\| = 1$ gibt es Störungen ΔA mit $\|\Delta A\|$ $= \Delta A$, so daß das Gleichheitszeichen steht.

(ii) Das Rangerhöhungsverbot (16) ist für jede beliebige Störung δA erfüllt, wenn A Vollrang hat, also spalten- oder zeilenregulär ist. Für Vollrangmatrizen ist die Bestimmung der Pseudoinversen daher ein korrekt gestelltes Problem, was übrigens bereits aus den in diesem Fall gültigen expliziten Darstellungen 8.1.9(i) für Afolgt. Für rangdefizientes A ist die Bestimmung von A^+ ein inkorrekt gestelltes Problem, sofern die Störungen nicht gemä β (16) eingeschränkt werden.

(iii) Für reguläres A geht (22) in die Abschätzung (4.1.10) des Störungslemmas über, während sich (20) wegen $Q_A = Q_{A^{\mathsf{T}}} = 0$ auf die bekannte Formel

$$(A + \delta A)^{-1} - A^{-1} = -A^{-1}\delta A A^{-1} + O(\|\delta A\|^2)$$

reduziert, siehe Ü 4.1.2.

B. Störungstheorie der Normallösung

Wir kehren jetzt zu den Quadratmittelproblemen (1), (3) zurück und untersuchen die Empfindlichkeit von x gegenüber nicht rangerhöhenden Störungen dA.

8.2.7. Satz. Gegeben seien die Quadratmittelprobleme

 $Ax \cong b$ und $(A + \delta A) (x + \delta x) \cong b + \delta b$ mit den Normallösungen

 $\boldsymbol{x} = \boldsymbol{A}^{\scriptscriptstyle +} \boldsymbol{b}$ bzw. $\boldsymbol{x} + \boldsymbol{\delta} \boldsymbol{x} = (\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A})^{\scriptscriptstyle +} (\boldsymbol{b} + \boldsymbol{\delta} \boldsymbol{b}).$

17 Schwetlick, Numerische Algebra

Die Störung δA genüge den Bedingungen rang $(A + \delta A) \leq \operatorname{rang}(A)$ sowie $\varkappa := ||A^+|| ||\delta A|| < 1$. Dann gilt (i) $\boldsymbol{\delta x} = \boldsymbol{\delta x'} + O(\|\boldsymbol{\delta A}\| (\|\boldsymbol{\delta A}\| + \|\boldsymbol{\delta b}\|)),$ wobei $\delta x' := A^+ \{ -\delta A \ x + \delta b \} + A^+ A^+ ^{\mathsf{T}} \delta A^{\mathsf{T}} r + (I - A^+ A) \ \delta A^{\mathsf{T}} A^{\mathsf{T}} r$ (28)den bezüglich δA , δb linearen Teil von δx darstellt und $\boldsymbol{r} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x} = (\boldsymbol{I} - \boldsymbol{A}\boldsymbol{A}^{+}) \boldsymbol{b}$ das Residuum von x bezeichnet. (ii) Die Störung dx genügt den Abschätzungen $\| \boldsymbol{\delta x} \| \leq rac{\| \boldsymbol{A}^{+} \|}{1- lpha} \left\{ \| \boldsymbol{\delta A} \| \left[\omega \| \boldsymbol{x} \| + \| \boldsymbol{A}^{+} \| \| \boldsymbol{r} \|
ight] + \| \boldsymbol{\delta b} \|
ight\}$ (29)und (im Fall $x \neq o$) $\frac{\|\boldsymbol{\delta x}\|}{\|\boldsymbol{x}\|} \leq \frac{1}{1-\kappa} \left\{ \left[\omega \text{ cond } (A) + \frac{\|\boldsymbol{r}\|}{\|\boldsymbol{A}\| \|\boldsymbol{x}\|} \left(\text{cond } (A) \right)^2 \right] \frac{\|\boldsymbol{\delta A}\|}{\|\boldsymbol{A}\|} \right.$ $+ \frac{\|\boldsymbol{A}^+\| \|\boldsymbol{b}\|}{\|\boldsymbol{x}\|} \cdot \frac{\|\boldsymbol{\delta}\boldsymbol{b}\|}{\|\boldsymbol{b}\|}$ (30)mit cond $(A) := ||A|| ||A^+||$ (31)und $\omega := \begin{cases} \sqrt{2}, & \text{falls rang} (A) < n, \\ 1, & \text{falls rang} (A) = n. \end{cases}$ Beweis. Mit den Bezeichnungen aus dem Beweis zu 8.2.5 und (5) gilt $\delta x = Gb + \tilde{A}^+ \delta b = G_1 b + G_2 b + G_3 b + \tilde{A}^+ \delta b$ (32)Unter Beachtung von (19), (25) und (8.1.42) folgt $\|G,b\| = \|\tilde{A}^{+} \delta A A^{+} b\| = \|\tilde{A}^{+} \delta A x\| \leq \left[\|A^{+}\|/(1-\varkappa)\right] \|\delta A\| \|x\| =: \beta,$

 $\|G_{2}b\| = \|(I - \tilde{S}) \ \delta A^{\mathsf{T}} A^{+\mathsf{T}} A^{+} b\| \leq \|A^{+}\| \|\delta A\| \|x\| \leq \beta,$

wegen der Orthogonalität von G_1b und G_2b also

$$\|G_1b + G_3b\|^2 = \|G_1b\|^2 + \|G_3b\|^2 \le 2\beta^2.$$

Wegen $(\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{b} = \boldsymbol{r} = (\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{r}$ gilt

$$G_2 \boldsymbol{b} = \tilde{A}^+ (\boldsymbol{I} - \boldsymbol{P}) \, \boldsymbol{b} = \tilde{A}^+ (\boldsymbol{I} - \boldsymbol{P}) \, \boldsymbol{r} = G_2 \boldsymbol{r}, \tag{33}$$

und mit (27) folgt

$$\|\boldsymbol{G}_{2}\boldsymbol{b}\| \leq \alpha \|\boldsymbol{r}\|, \qquad \|\bar{A}^{+}\boldsymbol{\delta}\boldsymbol{b}\| \leq [\|A^{+}\|/(1-\varkappa)] \|\boldsymbol{\delta}\boldsymbol{b}\| =: \gamma.$$

Damit erhalten wir

$$\|\delta \boldsymbol{x}\| \leq \|\boldsymbol{G_1}\boldsymbol{b} + \boldsymbol{G_3}\boldsymbol{b}\| + \|\boldsymbol{G_2}\boldsymbol{b}\| + \|\boldsymbol{\tilde{A}}^{+}\delta\boldsymbol{b}\| \leq \sqrt{2}\,\beta + \alpha\,\|\boldsymbol{r}\| + \gamma.$$

Dies ist gerade (29), und Division durch $||\boldsymbol{x}||$ liefert (30). Im Fall rang (\boldsymbol{A}) = n ist $\boldsymbol{G}_3 = \boldsymbol{O}$, so daß $\sqrt{2}$ durch 1 ersetzt werden kann. Die Linearisierungsformel (28) ergibt sich mit (33) aus (32), wenn $\tilde{\boldsymbol{A}}$ überall durch \boldsymbol{A} ersetzt wird. \Box

Satz 8.2.7 besagt: Die Bestimmung der Normallösung eines Quadratmittelproblems ist eine lokal lipschitzstetige Aufgabenklasse, sofern keine rangerhöhenden Störungen von A zugelassen werden.

Die partiellen relativen Konditionszahlen sind

$$K_{A}(A, b) = \omega \text{ cond } (A) + \frac{\|r\|}{\|A\|} \|x\|}{[\text{cond } (A)]^{2}}, \quad K_{b}(A, b) = \frac{\|A^{+}\| \|b\|}{\|x\|}, \quad (34)$$

und es gilt

$$1 \leq K_{\boldsymbol{b}}(\boldsymbol{A}, \boldsymbol{b}) \leq K_{\boldsymbol{A}}(\boldsymbol{A}, \boldsymbol{b}). \tag{35}$$

Die durch (31) eingeführte Zahl cond (A) ist für reguläres A mit der in 4.1 eingeführten Konditionszahl identisch und wird deshalb auch (*Spektral-*) Konditionszahl der rechteckigen Matrix A genannt. Für sie gilt

$$\operatorname{cond} (A) = \sigma_1 / \sigma_\tau, \quad \operatorname{cond} (A) \ge 1 \text{ für } A \neq O, \quad \operatorname{cond} (A) = 0 \text{ für } A = O.$$

(36)

8.2.8. Bemerkung. (i) Für konsistente Gleichungssysteme – zu diesen gehören diejenigen mit zeilenregulärer Matrix A, vgl. 8.1.A – ist r = o, so daß (30) in

$$\frac{\|\boldsymbol{\delta x}\|}{\|\boldsymbol{x}\|} \leq \frac{1}{1-\varkappa} \left\{ \omega \text{ cond } (\boldsymbol{A}) \ \frac{\|\boldsymbol{\delta A}\|}{\|\boldsymbol{A}\|} + \frac{\|\boldsymbol{A}^+\| \|\boldsymbol{b}\|}{\|\boldsymbol{x}\|} \ \frac{\|\boldsymbol{\delta b}\|}{\|\boldsymbol{b}\|} \right\}$$
(37)

übergeht. Bis auf den Faktor $\omega \leq \sqrt{2}$ entspricht (37) der für reguläre Systeme gültigen Abschätzung (4.1.14). Die Normallösung konsistenter rechteckiger Quadratmittelprobleme hat also dasselbe qualitative Störungsverhalten gegenüber kleinen, den Rang nicht erhöhenden Störungen wie die Lösung regulärer Gleichungssysteme. Insbesondere ist K_A — bis auf den für rang (A) < n auftretenden Faktor $\omega = \sqrt{2}$ mit cond (A) identisch und hängt nur von der Matrix A ab. Für reguläres A geht (37) bzw. (30) in (4.1.14) über.

(ii) Für inkonsistente Gleichungssysteme ist $r \neq o$, und in K_A tritt zusätzlich zu ω cond (A) der zu [cond (A)]² proportionale Term

$$\frac{\|\boldsymbol{r}\|}{\|\boldsymbol{A}\| \|\boldsymbol{x}\|} [\text{cond} (\boldsymbol{A})]^2$$
(38)

auf.

Für genügend kleines $\|r\|$ dominiert der erste Summand in (34), d. h., es gilt

$$K_A \approx \omega \text{ cond } (A),$$

und das Störungsverhalten ist ähnlich wie im konsistenten Fall. Für großes $||\mathbf{r}||$ dominiert (38), d. h., K_A wird proportional zu [cond (A)]². Im Gegensatz zur Situation bei regulären Gleichungssystemen hängt also K_A bei Quadratmittelproblemen nicht nur von A, sondern über das Residuum \mathbf{r} auch von der rechten Seite \mathbf{b} ab und kann proportional zu [cond (A)]² sein.

(iii) Die Bedingung rang $(A + \delta A) \leq \operatorname{rang} (A)$ ist für alle A erfüllt, wenn A spalten- oder zeilenregulär ist. Im ersten Fall ist x die eindeutige Quadratmittellösung von (1), im zweiten Fall ist (1) konsistent und r = o. \Box

Wir zeigen abschließend an einem Beispiel, wie sich rangerhöhende Störungen auf die Normallösung auswirken können.

8.2.9. Beispiel. Es seien A, δA wie in 8.2.2, und b, δb seien durch

$$\beta = U^{\mathsf{T}}b, \quad \delta\beta = U^{\mathsf{T}}\delta b$$

gegeben. Dann gilt für $\varepsilon \neq 0$

$$\delta x = V \delta \xi$$
 mit $\delta \xi = (\delta \xi_j), \quad \delta \xi_j = \begin{cases} \delta eta_j / \sigma_j & ext{fur } j = 1, ..., r, \ (eta_{r+1} + \delta eta_{r+1}) / arepsilon & ext{fur } j = r+1, \ 0 & ext{sonst}, \end{cases}$

. . . .

also

$$\| \boldsymbol{\delta x} \|^2 = \| \boldsymbol{\delta \xi} \|^2 = \| \boldsymbol{\Sigma}^{\scriptscriptstyle +} \boldsymbol{\delta \beta} \|^2 + [(eta_{r+1} + \delta eta_{r+1}) / arepsilon]^2.$$

Im Fall $\beta_{r+1} + \delta\beta_{r+1} = 0$ wirkt sich ε überhaupt nicht auf dx aus, obwohl die Pseudoinverse wie $1/|\varepsilon|$ wächst. Die zu $1/|\varepsilon|$ proportionale bösartige Wirkung tritt nur im Fall $\beta_{r+1} + \delta\beta_{r+1} \neq 0$ auf. \Box

Der Einfluß einer rangerhöhenden Störung δA auf δx hängt also auch von $\mathbf{b} + \delta \mathbf{b}$ ab. Die Störung wirkt nicht bösartig, wenn im diagonalisierten Problem $(\Sigma + \delta \Sigma)$ $\times (\mathbf{\xi} + \delta \mathbf{\xi}) \cong \beta + \delta \beta$ die Komponenten $(\beta + \delta \beta)_i$ für $i = r + 1, ..., l; l = \min(m, n)$, verschwinden. Wir werden später im Kapitel 11 solche rangerhöhenden, aber trotzdem gutartigen Störungen künstlich zur Regularisierung rangdefizienter Quadratmittelprobleme einführen.

C. Die Qualität der Störungsabschätzungen

Die Qualität der Abschätzung (29) bzw. (30) kann unter Verwendung der linearisierten Störung dx' aus (28) analog zu 4.1.B untersucht werden. Wie beim Beweis von 4.1.5 wird gezeigt, daß

$$A^{+} \{-\delta A \boldsymbol{x} + \delta \boldsymbol{b}\} \quad \text{durch} \quad \|A^{+}\| \{\|\delta A\| \|\boldsymbol{x}\| + \|\delta \boldsymbol{b}\|\}$$
(39)

scharf abgeschätzt wird. Ebenso wird

$$A^{+}A^{+} \sigma A^{\dagger} r \quad \text{durch} \quad \|A^{+}\|^{2} \| \sigma A \| \| r \|$$

$$\tag{40}$$

scharf abgeschätzt. Der im Fall rang (A) = n wegfallende und zu $A^+ \delta Ax$ orthogonale dritte Term

$$(\boldsymbol{I} - \boldsymbol{A}^{+}\boldsymbol{A})\,\boldsymbol{\delta}\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A}^{+}\boldsymbol{x} \tag{41}$$

von $\delta x'$ erhöht den Faktor 1 bei $\|\delta A\| \|x\|$ in (29) auf $\sqrt{2}$, spielt also praktisch keine Rolle. Sofern einer der beiden Summanden von K_A deutlich dominiert, sind die Abschätzungen von 8.2.7 also asymptotisch fast scharf. Eine genauere Untersuchung der zu den Termen (39), (40) gehörenden Fehlerellipsoide analog zu 4.1.7 zeigt, daß diese Abschätzungen den maximal möglichen Fehler stets asymptotisch qualitativ richtig wiedergeben.

D. Skalierung

Bei regulären Gleichungssystemen haben wir eine Skalierung mit einer Diagonalmatrix **D** vorgenommen, um die Kondition von **D**A zu reduzieren oder die Fehlerabschätzung bezüglich elementweiser Schranken für $\mathcal{O}A$, $\mathcal{O}b$ realistisch zu gestalten, vgl. 4.1.C. Es bleibt zu überlegen, ob für Quadratmittelprobleme ein analoges Vorgehen, d. h., der Übergang von

 $Ax \simeq b$ zu $DAx \simeq Db$ mit $D = \text{diag}(d_i), \quad d_i > 0 \quad (i = 1, ..., m),$

sinnvoll ist. Wir unterscheiden wieder zwei Fälle.

Fall 1: Ax = b konsistent, r = o. Dann sind die Lösungen des skalierten und unskalierten Quadratmittelproblems identisch und stimmen mit den Lösungen der jeweiligen Gleichungssysteme überein. Die Formulierung als Quadratmittelproblem ist also überflüssig. In diesem Fall können die Überlegungen aus 4.1.C sinngemäß übernommen werden, denn der Term (40) verschwindet in der linearisierten Störung $\delta x'$, und (41) trägt nur unwesentlich zu $\|\delta x'\|$ bei. Fall 1 tritt stets ein, wenn Azeilenregulär ist.

Fall 2: Ax = b inkonsistent, $r \neq o$. Dies ist in der Regel – d. h., für fast alle b – der Fall, wenn m > n gilt, Ax = b also streng überbestimmt ist, vgl. 8.1.A. Beim Übergang vom unskalierten zum skalierten Quadratmittelproblem ändert sich i. allg. die Lösungsmenge. Deshalb muß eine Skalierung nach den obigen numerischen Gesichtspunkten verworfen werden, sondern es ist so zu skalieren, daß die sich dabei einstellende Lösung im Hinblick auf die im Hintergrund stehende reale Aufgabe "richtig" bzw. zumindest "vernünftig" ist. Da die auf überbestimmte Gleichungen führenden Probleme in der Regel statistischer Natur sind, wird die natürliche Skalierung bzw. Wichtung des zugehörigen Quadratmittelproblems durch die statistischen Eigenschaften des Hintergrundproblems festgelegt und muß als gegeben angesehen werden, siehe Abschnitt 8.3.

D. Residualabschätzungen

Da 4.1.4 auch für rechteckige Matrizen gilt, kann 4.1.18 wörtlich übernommen werden.

8.2.10. Aussage. Es seien $\Delta A, \Delta b \ge 0$ und $A \in \mathbb{R}^{m,n}$, $b \in \mathbb{R}^{m}$, $\tilde{x} \in \mathbb{R}^{n}$ gegeben. Dann existieren genau dann Störungen $\delta A, \delta b$ mit

$$(A + \delta A) \,\tilde{\boldsymbol{x}} = \boldsymbol{b} + \delta \boldsymbol{b} \quad \text{und} \quad \|\delta A\| \leq \Delta A, \qquad \|\delta \boldsymbol{b}\| \leq \Delta \boldsymbol{b}, \tag{42}$$

wenn

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{\tilde{x}}\| \leq \Delta \boldsymbol{b} + \Delta \boldsymbol{A} \|\boldsymbol{\tilde{x}}\|$$
(43)

ist.

Falls $\hat{\mathbf{r}} = \mathbf{b} - A\hat{\mathbf{x}}$ im Sinne von (43) klein ist, stellt $\hat{\mathbf{x}}$ also eine Lösung – i. allg. jedoch nicht die Normallösung – des benachbarten Systems (42) dar. Für konsistente Systeme, bei denen r = b - Ax verschwindet, ist dies ein brauchbares Resultat. Für inkonsistente Systeme, bei denen das minimale Residuum von o verschieden ist, hat es jedoch keine praktische Bedeutung. Für solche Systeme bietet sich an, statt des inkonsistenten Systems (1) die stets konsistenten Normalgleichungen

$$g := A^{\mathsf{T}}r = A^{\mathsf{T}}(b - Ax) = A^{\mathsf{T}}b - A^{\mathsf{T}}Ax =: c - Mx$$

als Ausgangspunkt zu nehmen. Gilt

$$\|\tilde{\boldsymbol{g}}\| = \|\boldsymbol{A}^{\mathsf{T}}\tilde{\boldsymbol{r}}\| \leq \varepsilon \,\|\boldsymbol{A}\| \,\|\tilde{\boldsymbol{r}}\|,\tag{44}$$

so folgt durch Anwendung von 8.2.10 auf Mx = c mit $\Delta c := 0, \Delta M := \varepsilon ||A|| ||\tilde{x}||$ die Existenz einer Störung δM mit

$$(A^{\mathsf{T}}A + \delta M)\,\,\tilde{\boldsymbol{x}} = A^{\mathsf{T}}\boldsymbol{b} \quad \text{und} \quad \|\delta M\| \leq \Delta M,\tag{45}$$

wobei δM nach 4.1.20(ii) sogar symmetrisch gewählt werden kann. Der Vektor \tilde{x} löst also die gestörten Normalgleichungen (45), aber diese sind - leider! - nicht die Normalgleichungen eines gestörten Quadratmittelproblems, so daß bezüglich der eigentlichen Aufgabenstellung $Ax \simeq b$ keine Gutartigkeit gefolgert werden kann.

Wir gehen daher einen anderen Weg und wenden 8.2.10 auf das System

$$A^{\intercal}r = g = o$$

an. Mit $\varDelta g := 0, \, \varDelta A^{\intercal} := \varepsilon \, \|A\|$ folgt dann die Existenz einer Störung $\delta_1 A$ mit

$$(A + \boldsymbol{\delta}_1 A)^{\mathsf{T}} \, \tilde{\boldsymbol{r}} = (A + \boldsymbol{\delta}_1 A)^{\mathsf{T}} \, (\boldsymbol{b} - A \tilde{\boldsymbol{x}}) = \boldsymbol{O} \quad \text{und} \quad \|\boldsymbol{\delta}_1 A\| \leq \varepsilon \, \|A\|.$$
(46)

Das folgende Lemma zeigt, daß den gestörten Normalgleichungen (46) die Normalgleichungen eines gestörten Quadratmittelproblems zugeordnet werden können, wobei noch etwas allgemeinere Störungen zugelassen werden dürfen.

8.2.11. Lemma. Für $A \in \mathbb{R}^{m,n}$, $b \in \mathbb{R}^m$, $\tilde{x} \in \mathbb{R}^n$ und die Störungen $\delta_1 A$, $\delta_2 A$ gelte $(A + \delta_1 A)^{\mathsf{T}} [\boldsymbol{b} - (A + \delta_2 A) \, \boldsymbol{\tilde{x}}] = \boldsymbol{O}.$ Dann gibt es eine Störung $\boldsymbol{\delta}A$ mit $\|\boldsymbol{\delta}A\| \leq \{\|\boldsymbol{\delta}_1 A\|^2 + \|\boldsymbol{\delta}_2 A\|^2\}^{1/2},$ so daß $(A + \boldsymbol{\delta}A)^{\mathsf{T}} [\boldsymbol{b} - (A + \boldsymbol{\delta}A) \, \boldsymbol{\tilde{x}}] = \boldsymbol{O}$ gilt, d. h., $\boldsymbol{\tilde{x}}$ löst das Quadratmittelproblem $(A + \boldsymbol{\delta}A) \, \boldsymbol{\tilde{x}} \simeq \boldsymbol{b}.$

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq \{\|\boldsymbol{\delta}_1\boldsymbol{A}\|^2 + \|\boldsymbol{\delta}_2\boldsymbol{A}\|^2\}^{1/2},\tag{47}$$

Beweis. Im Fall $\mathbf{s} := \mathbf{b} - (\mathbf{A} + \delta_2 \mathbf{A}) \, \tilde{\mathbf{x}} = \mathbf{o}$ setze man $\delta \mathbf{A} := \delta_2 \mathbf{A}$, and ernfalls

$$\delta A := P \delta_1 A + (I - P) \delta_2 A$$
 mit $P := s s^{\intercal} / ||s||^2$.

Nach Ü 8.1.7 ist **P** der Projektor auf span {s}. Mit $D := \delta_1 A - \delta_2 A$ folgt

$$\boldsymbol{t} := \boldsymbol{b} - (\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A})\,\boldsymbol{\tilde{x}} = \boldsymbol{b} - (\boldsymbol{A} + \boldsymbol{\delta}_2\boldsymbol{A})\,\boldsymbol{\tilde{x}} - \boldsymbol{P}\boldsymbol{D}\boldsymbol{\tilde{x}} = \boldsymbol{\lambda}\boldsymbol{s},$$

wobei $\lambda := 1 - (\boldsymbol{D} \boldsymbol{\tilde{x}})^{\mathsf{T}} \boldsymbol{s} / \|\boldsymbol{s}\|^2$ ist. Mithin gilt

$$(A + \delta A)^{\mathsf{T}} \boldsymbol{t} = [A + \delta_1 A - (\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{D}]^{\mathsf{T}} \boldsymbol{t}$$
$$= \lambda [(A + \delta_1 A)^{\mathsf{T}} \boldsymbol{s} - \boldsymbol{D}^{\mathsf{T}} (\boldsymbol{I} - \boldsymbol{P}) \boldsymbol{s}] = \boldsymbol{o}. \quad \Box$$

Offensichtlich kann in (47) auch die Frobeniusnorm verwendet werden.

Bei Anwendung von 8.2.11 auf (46) ergibt sich mit $\sigma_2 A = O$ das nachfolgende gewünschte Resultat:

8.2.12. Aussage. Das zu \tilde{x} gehörende Residuum

$$\tilde{g} = A^{\mathsf{T}} \tilde{r} = A^{\mathsf{T}} (b - A \tilde{x})$$

 $\tilde{g} = A^{\intercal} \tilde{r} = A^{\intercal} (b - A \tilde{x})$ bezüglich der Normalgleichungen von $Ax \simeq b$ genüge der Abschätzung $\|\tilde{g}\| \leq \varepsilon \|A\| \|\hat{r}\|.$ Dann existiert eine Störung δA , so daß x das Quadratmittelproblem

$$\|\tilde{\boldsymbol{g}}\| \leq \varepsilon \|\boldsymbol{A}\| \|\tilde{\boldsymbol{r}}\|. \tag{48}$$

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A})\,\,\boldsymbol{\tilde{\boldsymbol{x}}} \cong \boldsymbol{b} \quad \text{mit} \quad \|\boldsymbol{\delta}\boldsymbol{A}\| \leq \varepsilon \,\|\boldsymbol{A}\| \tag{49}$$

löst.

Wir bemerken noch, daß (48) hinreichend, aber keineswegs notwendig für (49)ist. Insofern ist 8.2.12 schwächer als 8.2.10, wo (42) und (43) äquivalent sind.

Übungsaufgaben

Ü 8.2.1. Für $A, \tilde{A} \in \mathbb{R}^{m,n}$ gelte rang $(A) = \operatorname{rang} (\tilde{A})$. Man zeige

$$\| ilde{m{P}}(m{I}-m{P})\| = \|m{P}(m{I}- ilde{m{P}})\|$$
 für $m{P}:=AA^+,$ $m{ ilde{m{P}}}:=m{ ilde{A}}\,m{ ilde{A}}^+,$

Hinweis: Unter Verwendung der Singulärwertzerlegungen stelle man $P, \, ilde{P}$ gemäß

$$P = U\Pi U^{\intercal}, \qquad \widetilde{P} = \widetilde{U}\Pi\widetilde{U}^{\intercal}$$

mit $\Pi = \begin{pmatrix} I_r \mid O \\ O \mid O \end{pmatrix}$ dar und partitioniere $W := \tilde{U}^{\mathsf{T}}U$ gemäß $W = \begin{pmatrix} W_{11} \mid W_{12} \\ W_{21} \mid W_{22} \end{pmatrix}$. Man beweise $\|\tilde{P}(I_{-}, P)\| = \|W_{12}\|, \|P(I - \tilde{P})\| = \|W_{21}\|$ und folgere $\|W_{12}\| = \|W_{21}\|$ aus der Orthogonalität von W.

Ü 8.2.2. Für $\boldsymbol{y} \in \mathsf{R}^m$ und \boldsymbol{G}_i gemäß (25) mit $\|\boldsymbol{G}_i\| \leq lpha$ werde

$$z = G_1 y + G_2 y + G_3 y =: z_1 + z_2 + z_3$$

gebildet. Man zeige:

(i) $z_1 \perp z_3, z_2 \perp z_3$, folglich $||z||^2 = ||z_1 + z_2||^2 + ||z_3||^2$.

(ii) Mit $u := Py, v := (I - P)y, P := AA^+$, gilt $||u||^2 + ||v||^2 = ||y||^2$ und

 $G_1 y = G_1 u,$ $G_2 y = G_2 v,$ $G_3 y = G_3 u,$

wegen $\|\boldsymbol{z}_1 + \boldsymbol{z}_2\| \le \|\boldsymbol{G}_1 \boldsymbol{y}\| + \|\boldsymbol{G}_2 \boldsymbol{y}\| = \|\boldsymbol{G}_1 \boldsymbol{u}\| + \|\boldsymbol{G}_2 \boldsymbol{v}\| \le \alpha (\|\boldsymbol{u}\| + \|\boldsymbol{v}\|)$ also

$$m{z}_{||}^{-2} \leq lpha^2 \{ (||m{u}|| + ||m{v}||)^2 + ||m{u}||^2 \} =: lpha^2 \, ||m{y}||^2 \, N \, .$$

Mit $|\boldsymbol{u}| = ||\boldsymbol{y}|| \cos \varphi$, $||\boldsymbol{v}|| = ||\boldsymbol{y}|| \sin \varphi$ folgt

$$N=(\cosarphi+\sinarphi)^2+\cos^2arphi=(3+2\sin2arphi+\cos2arphi)/2\leqigg[ig(1+\sqrt{5}ig)/2ig]^2.$$

(iii) Im Fall r = n ist $\mathbf{z}_3 = \mathbf{0}$, also $||\mathbf{z}||^2 = ||\mathbf{z}_1 + \mathbf{z}_2||^2 \le 2\alpha^2 ||\mathbf{y}||^2$. Im Fall r = m ist $\mathbf{z}_2 = \mathbf{0}$, also $||\mathbf{z}||^2 = ||\mathbf{z}_1||^2 + ||\mathbf{z}_3||^2 \le 2\alpha^2 ||\mathbf{y}||^2$. Im Fall r = m = n ist $||\mathbf{z}|| = ||\mathbf{z}_1|| \le \alpha ||\mathbf{y}||$.

 $\ddot{\mathbf{U}}$ 8.2.3. Man beweise, daß cond (A) = 1 genau dann gilt, wenn A von der Gestalt

$$A = \sigma \overline{U} \overline{V}^{\intercal} = \sigma \sum_{i=1}^{r} \overline{u}^{i} \overline{v}^{i\intercal}$$

mit $\sigma > 0$ und spaltenorthonormalen Matrizen $\overline{U} \in \mathbb{R}^{m,r}$, $\overline{V} \in \mathbb{R}^{n,r}$, $r = \operatorname{rang}(A)$ ist.

8.3. Klassifikation und statistische Interpretation von Quadratmittelproblemen

Wir betrachten das Quadratmittelproblem

$$Ax \simeq b \quad \text{mit} \quad A \in \mathbb{R}^{m,n}, \qquad b \in \mathbb{R}^m$$
 (1)

und setzen $m \ge n$ voraus, d. h., (1) sei das zum überbestimmten und i. allg. inkonsistenten Gleichungssystem

$$Ax = b \tag{2}$$

gehörende Ersatzproblem. Vom numerischen Standpunkt lassen sich die Probleme des Typs (1) in zwei Klassen einteilen.

- 1. Spaltenreguläre Probleme: Rang $(A) = r = n \leq m$ Solche Probleme sind stets eindeutig lösbar, und die Lösung x hängt lokal lipschitzstetig von den Eingangsdaten ab. Sie sind also korrekt gestellt und können i. allg. befriedigend gelöst werden.
- 2. Rangdefiziente Probleme: Rang $(A) = r < n \leq m$ Solche Probleme haben unendlich viele Lösungen, und die Normallösung als Lösung kleinster Euklidischer Norm hängt unstetig von den Eingangsdaten ab. Sie sind daher inkorrekt gestellt und können ohne zusätzliche Information – z. B. über den Rang von A, um bösartige rangerhöhende Störungen in A auszuschließen – nicht befriedigend gelöst werden.

Die Quelle für überbestimmte inkonsistente Gleichungssysteme (2) ist meist das lineare statistische Modell

$$\boldsymbol{b} = \boldsymbol{A}\boldsymbol{x}^* + \boldsymbol{\epsilon}. \tag{3}$$

264

Dabei stellt ε einen *m*-dimensionalen zufälligen Vektor mit dem Erwartungswert

$$E(\varepsilon) = 0 \tag{4}$$

und der Kovarianzmatrix

$$\boldsymbol{E}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}^{\mathsf{T}}) = \sigma^{2}\boldsymbol{I} \tag{5}$$

dar, und für A gelte rang (A) = n < m. Der Vektor x^* ist ein nicht zufälliger, unbekannter Parametervektor, d r aus dem Zufallsvektor b geschätzt werden soll. Wenn die eindeutige Lösung

$$\boldsymbol{x} = \boldsymbol{A}^{\mathsf{+}}\boldsymbol{b} = (\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A})^{-1}\boldsymbol{A}^{\mathsf{T}}\boldsymbol{b}$$
(6)

von (1) als Schätzer für x^* verwendet wird, gilt

$$E(x) = A^{+}E(b) = A^{+}Ax^{*} = x^{*},$$
 (7)

d. h., \boldsymbol{x} ist eine erwartungstreue Schätzung, und die Kovarianzmatrix ergibt sich zu

$$E((\boldsymbol{x} - \boldsymbol{x^*}) (\boldsymbol{x} - \boldsymbol{x^*})^{\mathsf{T}}) = \sigma^2 A^+ A^{+\mathsf{T}} = \sigma^2 (A^{\mathsf{T}} A)^{-1}.$$
(8)

Bei bekanntem σ^2 sind unter Verwendung der Kovarianzmatrix $\sigma^2(A^{\mathsf{T}}A)^{-1}$ Aussagen über die Qualität der Schätzung x möglich. Ist σ^2 nicht bekannt, so kann

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^2 / (m-n) \tag{9}$$

als Schätzung für σ^2 genommen werden. Falls ε einer *m*-dimensionalen Normalverteilung mit den Parametern (4), (5) genügt, ist *x* gerade die Maximum-Likelihood-Schätzung. Wir bemerken noch, daß auch nicht erwartungstreue Schätzungen etwa solche nach dem Prinzip der "ridge regression" — praktische Bedeutung haben, vgl. Kapitel 11.3 und B 11.3.

Für eine detaillierte Beschreibung des statistischen Hintergrundes muß auf die einschlägige statistische Fachliteratur verwiesen werden, siehe B 8.5. Aus den skizzierten Problemen wird jedoch klar, daß neben der Lösung x in der Regel noch weitere numerische Informationen als zusätzliche Ausgangsdaten gewünscht werden, etwa

- (I₁) die ausgeglichenen Werte Ax von b und das minimale Residuum r = b Ax, zumindest aber dessen Norm $||r||^2$,
- (I₂) alle oder einige Elemente von $C := (A^{\intercal}A)^{-1}$ sowie Ausdrücke der Form BCB^{\intercal} bzw. $F^{\intercal}C^{-1}F$ für gegebene Matrizen B bzw. F,
- (I_3) der Rang r von A, im Fall r < n die Indizes von r Spalten von A, die $\mathcal{R}(A)$ aufspannen, die Koeffizienten von n - r Abhängigkeitsbeziehungen zwischen den Spalten von A, ggf. eine (orthogonale) Basis von $\mathcal{N}(A)$,
- (I_4) die Indizes von möglichst wenig Spalten von A, so daß ein auf diese Spalten reduziertes Modell (3) ein Residuum r mit $||r|| \leq \delta$ liefert, wobei δ eine vorgegebene bzw. im Laufe der Rechnung festgelegte Schranke darstellt,
- (I_5) ein Schätzwert für $||A^+||$ oder eine andere Größe, die Aussagen über die Empfindlichkeit der Lösung gegenüber Störungen zuläßt.

Die Aufgaben (I_3) und (I_4) stehen in Zusammenhang mit der *Modellwahl*. Ziel ist dabei, mit möglichst wenigen Spalten von A einen akzeptablen Ausgleich zu erhalten. Man beachte dabei, daß sich beim Streichen von Spalten aus einer spaltenregulären Matrix die Kondition i. allg. verkleinert, siehe Ü 8.3.3. Eine minimale Spaltenzahl ist daher nicht nur vom statistischen, sondern auch vom numerischen Standpunkt her erwünscht.

Eine Möglichkeit zur Beschaffung der in (I_4) geforderten Information besteht darin, alle Quadratmittelprobleme zu lösen, bei denen A auf k der n Spalten reduziert ist (k = 1, ..., n), und die zugehörigen Residuen zu vergleichen. Das erfordert die Lösung einer Folge von $2^n - 1$ Quadratmittelproblemen, von denen jedes durch Hinzufügen, Streichen oder Ändern einer Spalte aus dem vorhergegangenen entstanden ist. Man spricht auch von schrittweiser Regression, siehe B 8.9. Solche Probleme treten auch bei der Lösung von Quadratmittelproblemen mit linearen Ungleichungsnebenbedingungen auf; zu diesen siehe B 8.8. Da andererseits jeder Zeile von (3) ein Meßwert zugeordnet werden kann, tritt dasselbe Problem bezüglich der Zeilen von (1) beim Hinzufügen, Streichen oder Ändern von Meßwertsätzen auf. Es besteht daher ein Interesse an Algorithmen, die eine Folge von Quadratmittelproblemen, die durch Hinzufügen, Streichen oder Ändern von Spalten bzw. Zeilen auseinander hervorgehen, effektiv zu lösen gestatten, vgl. Abschnitt 2.1.B1 und Ü 2.1.3.

Wir betrachten abschließend den Fall, daß
e eine von $\sigma^2 I$ verschiedene Kovarianzmatrix

$$E(\varepsilon\varepsilon^{\mathsf{T}}) = W \tag{10}$$

besitzt, wobei W notwendig symmetrisch und positiv semidefinit ist. Wir setzen W zusätzlich als positiv definit voraus. Dann ist die Maximum-Likelihood-Schätzung für x^* die Lösung des sog. gewichteten oder verallgemeinerten Quadratmittelproblems

$$\boldsymbol{\varepsilon}^{\mathsf{T}} \boldsymbol{W}^{-1} \boldsymbol{\varepsilon} = \boldsymbol{r}^{\mathsf{T}} \boldsymbol{W}^{-1} \boldsymbol{r} = (\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x})^{\mathsf{T}} \boldsymbol{W}^{-1} (\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}) \to \underset{\boldsymbol{x} \in \mathbb{R}^n}{\operatorname{Minimum}!}$$
(11)

Die Wichtung der Residuen mit W^{-1} kann daher als natürliche Wichtung angesehen werden.

Wir wollen (11) auf ein ungewichtetes Quadratmittelproblem zurückführen und verwenden dazu die nach dem Cholesky-Verfahren mit $\sim m^3/6$ opms berechenbare symmetrische Dreiecksfaktorisierung

$$W = LL^{\mathsf{T}} \tag{12}$$

von W, wobei L eine untere Dreiecksmatrix mit positiven Diagonalelementen ist, vgl. 5.2.D und 6.1.A. Mit den Transformationen

$$\bar{\boldsymbol{b}} := \boldsymbol{L}^{-1}\boldsymbol{b}, \qquad \bar{\boldsymbol{A}} := \boldsymbol{L}^{-1}\boldsymbol{A} \tag{13}$$

und

$$\overline{\boldsymbol{r}} := \boldsymbol{L}^{-1} \boldsymbol{r} = \boldsymbol{L}^{-1} (\boldsymbol{b} - \boldsymbol{A} \boldsymbol{x}) = \overline{\boldsymbol{b}} - \overline{\boldsymbol{A}} \overline{\boldsymbol{x}}$$
(14)

geht (11) in das ungewichtete Problem

$$\boldsymbol{r}^{\mathsf{T}}\boldsymbol{W}^{-1}\boldsymbol{r} = \bar{\boldsymbol{r}}^{\mathsf{T}}\boldsymbol{I}\bar{\boldsymbol{r}} = (\bar{\boldsymbol{b}} - \bar{\boldsymbol{A}}\boldsymbol{x})^{\mathsf{T}} (\bar{\boldsymbol{b}} - \bar{\boldsymbol{A}}\boldsymbol{x}) = \|\bar{\boldsymbol{b}} - \bar{\boldsymbol{A}}\boldsymbol{x}\|^2 \to \underset{\boldsymbol{x} \in \mathbf{R}^n}{\operatorname{Minimum}!} \quad (15)$$

über. Die transformierten gewichteten Größen werden dabei direkt aus den Dreieckssystemen

$$L\bar{b}=b, \quad L\bar{A}=A$$

mit $\sim m^2 n/2$ opms berechnet. (Die Bezeichnung $\bar{\boldsymbol{b}}$ hat hier nichts mit irgendwelchen Projektionen zu tun, sondern wurde in Anlehnung an 4.1.C gewählt.)

Für die meisten praktischen Ausgleichsaufgaben ist W diagonal, d. h.

$$W = \text{diag}(w_i), \quad w_i > 0 \quad (i = 1, ..., m).$$
 (16)

Das bedeutet im Fall der Normalverteilung, daß die Komponenten ε_i des Fehlervektors unabhängig sind und die i. allg. unterschiedlichen Varianzen $E(\varepsilon_i^2) = w_i$ haben. Dies entspricht unabhängigen Messungen mit unterschiedlicher Genauigkeit. Dann ist

$$oldsymbol{L} = ext{diag}\left(\!\!\sqrt{w_i}\!
ight) \qquad (i=1,...,m)$$

ebenfalls diagonal, und mit der Skalierungsmatrix

$$\boldsymbol{D} = \text{diag}(d_i) := \boldsymbol{L}^{-1}, \qquad d_i := 1/\sqrt{w_i} \qquad (i = 1, ..., m)$$
(17)

geht (15) in

$$\|\bar{\boldsymbol{b}} - \bar{\boldsymbol{A}}\boldsymbol{x}\|^2 = \|\boldsymbol{D}(\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x})\|^2 = \|\boldsymbol{D}\boldsymbol{b} - \boldsymbol{D}\boldsymbol{A}\boldsymbol{x}\|^2 \to \underset{\boldsymbol{x} \in \mathbb{R}^n}{\operatorname{Minimum}!}$$
(18)

über. Die *i*-te Zeile des Originalproblems ist also mit $1/\sqrt{w_i}$ zu multiplizieren.

Wenn die Kovarianzmatrix W schlecht konditioniert ist, wird das gewichtete Problem (11) *steif* genannt. Mit den Eigenwerten ω_i von W läßt sich dies durch

$$\operatorname{cond} (\boldsymbol{W}) = \max \, \omega_i / \min \, \omega_i \gg 1 \tag{19}$$

charakterisieren, d. h., W hat Eigenwerte stark unterschiedlicher Größenordnung. Für diagonales W ist $w_i = \omega_i$, so daß (19) in

$$\max w_i \gg \min w_i \tag{20}$$

übergeht. Dieser Fall tritt auf, wenn gewisse Komponenten von **b** wesentlich genauer gemessen werden können als die übrigen. Man beachte, daß im Fall $d_l = 1/\sqrt[]{w_l} \to \infty$ für ein festes $l \in \{1, ..., m\}$ die Lösung von (18) in die Lösung des Quadratmittelproblems

$$\sum_{\substack{i=1\\i\neq l}}^{m} [(\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x})_i]^2 / w_i \to \text{Minimum! bei } \boldsymbol{x} \in \mathbf{R}^n, \quad (\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x})_l = 0 \quad (21)$$

mit Gleichungsnebenbedingungen übergeht.

Für steife Probleme mit nichtdiagonalem W kann die explizite Ausführung der Standardtransformation (13) nicht empfohlen werden, da sie i. allg. zu einem hohen Genauigkeitsverlust führt. Wir gehen auf solche Probleme nicht-weiter ein, siehe B 8.10 und die dort angegebene Literatur.

Übungsaufgaben

Ü 8.3.1. Man zeige:

(i) Es sei

$$W = U\Omega U^{\intercal}, \quad \Omega = \text{diag}(\omega_i) > 0, \quad U \text{ orthogonal},$$

die Eigenwertzerlegung der symmetrischen und positiv definiten Matrix W. Dann gilt

$$W = W^{1/2} W^{1/2}.$$
 (22)

wobei

$$\boldsymbol{W}^{1/2} := \boldsymbol{U}\boldsymbol{\Omega}^{1/2}\boldsymbol{U}^{\mathsf{T}}, \qquad \boldsymbol{\Omega}^{1/2} := \operatorname{diag}\left(\boldsymbol{\gamma}\boldsymbol{\omega}_{i}\right) \tag{23}$$

die eindeutig festgelegte positiv definite Quadratwurzel aus W bezeichnet.

(ii) Mit
$$\hat{D} := \Omega^{-1/2} = \text{diag} (1/\sqrt{\omega_i})$$
 und

$$\hat{\mathbf{b}} := \mathbf{U}^{\mathsf{T}} \mathbf{b}, \qquad \hat{A} := \mathbf{U}^{\mathsf{T}} \mathbf{A}, \qquad \hat{r} := \mathbf{U}^{\mathsf{T}} \mathbf{r}$$
(24)

geht (11) in das diagonal gewichtete Quadratmittelproblem

$$\boldsymbol{r}^{\mathsf{T}} \boldsymbol{W}^{-1} \boldsymbol{r} = \|\boldsymbol{W}^{-1/2} \boldsymbol{r}\|^2 = \|\hat{\boldsymbol{D}}(\hat{\boldsymbol{b}} - \hat{\boldsymbol{A}} \boldsymbol{x})\|^2 \to \text{Minimum}!$$
(25)

über.

Ü 8.3.2. Man überlege sich, daß (15) unter Verwendung von (12) in der Form

 $||\boldsymbol{z}||^2 \to \text{Minimum!} \quad \text{bei } \boldsymbol{x} \in \mathbf{R}^n, \, \boldsymbol{z} \in \mathbf{R}^m \text{ mit } \boldsymbol{b} = A\boldsymbol{x} + L\boldsymbol{z}$ $\tag{26}$

als Quadratmittelproblem der Dimension n + m mit Gleichungsnebenbedingungen geschrieben werden kann.

Ü 8.3.3. Wenn $B \in \mathbb{R}^{m,n-1}$ aus $A \in \mathbb{R}^{m,n}$, $m \ge n \ge 2$, durch Streichen einer Spalte entsteht, gilt für die geordneten Singulärwerte β_i von B und α_i von A die Einschließung

$$\alpha_1 \ge \beta_1 \ge \alpha_2 \ge \beta_2 \ge \cdots \ge \beta_{n-1} \ge \alpha_n \ge 0,$$

siehe Lawson/Hanson [74, Theorem 5.12]. Man folgere hieraus im Fall rang (A) = n

$$\operatorname{cond} (\mathbf{B}) \leq \operatorname{cond} (\mathbf{A}).$$

Bemerkungen zum Kapitel 8

B 8.1. Die Lösung des Minimierungsproblems $||\mathbf{b} - A\mathbf{x}||_p \rightarrow \text{Minimum!}$ heißt L_p -Lösung des Gleichungssystems $A\mathbf{x} = \mathbf{b}$; für p = 1 spricht man auch von einer Lösung kleinster absoluter Abweichungen, für $p = \infty$ von einer Čebyšev-Lösung. Die mit der ∞ -Norm beabsichtigte Äquilibrierung des Residuums \mathbf{r} kann meist auch billiger in der 2-Norm durch geeignete (unter Umständen adaptive) Skalierung erreicht werden. Man beachte dabei, daß die Berechnung der L_{∞} -Lösung zur Lösung eines speziellen linearen Optimierungsproblems äquivalent und daher i. allg. aufwendig ist. Die L_1 -Lösungen haben gegenüber den Quadratmittellösungen den Vorteil, nicht so empfindlich auf einzelne große Komponenten von \mathbf{r} zu reagieren, also "robuster" gegenüber "Ausreißern" zu sein. Dem steht die weitaus größere Kompliziertheit in der theoretischen Durchdringung und praktischen Berechnung gegenüber.

B 8.2. Das für identisch normalverteilte Fehler optimale Prinzip der kleinsten Fehlerquadrate wurde 1809 von GAUSS bei geodätischen und astronomischen Berechnungen verwendet und trägt deshalb auch seinen Namen. Unabhängig ist es bereits 1806 von LAGRANGE eingeführt worden, was damals zu einem erbitterten Prioritätsstreit geführt hat, siehe etwa LINNIK [62] für einen historischen Abriß.

B 8.3. Die als Pseudoinverse bezeichnete Matrix wurde unabhängig von MOORE 1920 und PENROSE 1955 über unterschiedliche Definitionen eingeführt, die sich später als gleichwertig herausstellten. Zur Geschichte der Pseudoinversen siehe ZIELKE [78], wo auch viele nützliche Literaturhinweise zu finden sind.

B 8.4. Die Darstellung der Störungsaussagen für A^+ und x^N folgt der bei LAWSON/HANSON [74]. Solche Aussagen sind für beliebiges r = Rang(A) erstmals von WEDIN [73] bewiesen worden, obwohl für gewisse Sonderfälle auch vorher schon analoge Ergebnisse bekannt waren. Eine umfangreiche Zusammenfassung gibt STEWART [77]; dort ist auch Lemma 8.2.11 in einer im wesentlichen äquivalenten Form enthalten.

B 8.5. Über den statistischen Hintergrund linearer Quadratmittelprobleme kann man sich bei LINNIK [62], DRAPER/SMITH [66], K. M. S. HUMAK [77, 83] u. a. informieren.

B 8.6. Wenn im linearen Modell auch die Matrix A durch zufällige Fehler verfälscht ist, sollte neben x auch A in die Schätzung einbezogen und ein ausgeglichener Wert B für A aus dem erweiterten Problem

(TLS) $\|\boldsymbol{b} - \boldsymbol{B}\boldsymbol{x}\|^2 + \|\boldsymbol{A} - \boldsymbol{B}\|_F^2 \rightarrow \text{Minimum!}$ bei $\boldsymbol{x} \in \mathbb{R}^n, \boldsymbol{B} \in \mathbb{R}^{m,n}$

bestimmt werden. Solche Probleme werden totale Quadratmittelprobleme genannt. Wegen des quadratischen Terms Bx lassen sich diese nicht direkt auf gewöhnliche Quadratmittelprobleme zurückführen und erfordern spezielle Lösungsverfahren, siehe GOLUB/VAN LOAN [80]. Probleme des Typs (TLS) entstehen z. B. bei der Anpassung der skalaren linearen Modellgleichung

$$b = a_1 x_1 + \cdots + a_n x_n$$

an *m* Meßwertsätze $\{b_i, (a_1)_i, \ldots, (a_n)_i\}$, wenn sowohl b_i als auch die $(a_j)_i$ fehlerbehaftet sind. Man spricht dann auch von *linearen Modellen mit Fehlern in den Variablen*.

B 8.7. Gauß-Newton-Verfahren zur Lösung nichtlinearer Quadratmittelprobleme

$$\|\boldsymbol{r}(\boldsymbol{x})\| \rightarrow \text{Minimum!} \quad \text{bei } \boldsymbol{x} \in \mathbf{R}^n$$

mit einer nichtlinearen Vektorfunktion $r: \mathbb{R}^n \to \mathbb{R}^m$ erfordern in jedem Schritt die Lösung eines linearen Quadratmittelproblems, siehe BARD [73], SCHWETLICK [79], BJÖRCK [81] und GILL/MURRAY/WRIGHT [82]. Nichtlineare Modelle mit Fehlern in den Variablen führen auf lineare Quadratmittelprobleme spezieller Struktur, die angepaßte Lösungsverfahren benötigen, siehe SCHWETLICK/TILLER [85].

B 8.8. In gewissen praktischen Aufgaben treten zusätzliche Gleichungs- oder Ungleichungsnebenbedingungen auf, die zu Quadratmittelproblemen des Typs

(ELS)
$$||\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}|| \rightarrow \text{Minimum}$$
 bei $\boldsymbol{x} \in \mathbb{R}^n, C\boldsymbol{x} = \boldsymbol{d}$

bzw.

(ILS) $\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\| \rightarrow \text{Minimum}$ bei $\boldsymbol{x} \in \mathbb{R}^n, C\boldsymbol{x} = \boldsymbol{d}, \boldsymbol{E}\boldsymbol{x} \ge \boldsymbol{f}$

führen. Während (ELS) durch geeignete Transformationen auf ein Problem ohne Nebenbedingungen zurückgeführt werden kann, siehe LAWSON/HANSON [74, Chapt. 20-22], erfordert (ILS) den Einsatz von Optimierungstechniken und ist wesentlich aufwendiger zu lösen, siehe wieder LAWSON/HANSON [74, Chapt. 23] und STOER [71], SCHITTKOWSKI/STOER [79], HASKELL/HANSON [81], LÖTSTEDT [84] für geeignete Verfahren. Wir gehen auf beide Aufgabenklassen nicht weiter ein.

B 8.9. Algorithmen zur schrittweisen Regression sind auf der Basis der Normalgleichungen von EFROYMSON [60], auf der Basis der Orthogonalisierungsverfahren von ELDÉN [72] und LAWSON/HANSON [74] beschrieben worden, siehe auch OSBORNE [76]. **B 8.10.** Steife diagonalgewichtete Quadratmittelprobleme werden bei LAWSON/HANSON [74, Chapt. 17] und BJÖRCK [78] diskutiert. Normalgleichungsverfahren sind für solche Probleme nicht geeignet, während QR-Faktorisierungen von DA mit geeigneten Zusatzmaßnahmen zu brauchbaren Resultaten führen, siehe 10.2.5.

Für nichtdiagonales, schlechtkonditioniertes W empfiehlt Björck [81] zwei Vorgehensweisen:

- Ersatz der Cholesky-Faktorisierung (8.3.12) durch (8.3.22) unter Verwendung der Eigenwertzerlegung von W und Rückführung auf das diagonalgewichtete Problem (8.3.25) gemäß Ü 8.3.1,
- Formulierung als nichtgewichtetes Problem mit Gleichungsnebenbedingungen gemäß Ü 8.3.2 und Lösung mittels angepaßter Orthogonalisierungstechniken nach PAIGE [79].

9. Normalgleichungsverfahren

Normalgleichungsverfahren beruhen auf der Charakterisierung der Lösungen des Quadratmittelproblems $Ax \simeq b$ als Lösungen der zugehörigen Normalgleichungen

$$A^{\mathsf{T}}Ax = A^{\mathsf{T}}b$$
.

9.1. Aufgabenstellung und Lösung der Normalgleichungen

Wir bemerken zunächst, daß die Normalgleichungen im Fall $A \in \mathbf{R}^{m,n}$ ein quadratisches, stets konsistentes Gleichungssystem der Ordnung *n* darstellen. Die Koeffizientenmatrix ist für jedes A symmetrisch und positiv semidefinit; sie ist positiv definit genau dann, wenn A spaltenregulär ist, vgl. 1.2.14 und den nachfolgenden Text. Wir setzen A daher im folgenden als spaltenregulär voraus. Es bietet sich dann an, die eindeutige Lösung der Normalgleichungen durch Cholesky-Faktorisierung von $A^{\mathsf{T}}A$ zu berechnen.

9.1.1. Lösung von $Ax \simeq b$ über die Normalgleichungen mittels Cholesky-Faktorisierung.

Aufgabe. Für spaltenreguläres $A \in \Re^{m,n}$ ist die Koeffizientenmatrix $M = A^{\intercal}A$ der Normalgleichungen zu berechnen und nach dem Cholesky-Verfahren gemäß $M = LL^{\intercal}$ mit einer unteren Dreiecksmatrix $L \in \Re^{n,n}$ zu faktorisieren. Zu gegebenem $b \in \Re^m$ ist die rechte Seite $c = A^{\intercal}b$ der Normalgleichungen

$$Mx = c \tag{1}$$

von $Ax \simeq b$ zu berechnen und die zugehörige Lösung $x \in \Re^n$ unter Verwendung von L zu bestimmen.

Algorithmus:

- S1 (Berechnung und Cholesky-Faktorisierung von M):
- S1.1: Berechne $M := A^{\mathsf{T}}A$
- S1.2: Bestimme die Cholesky-Faktorisierung $M = LL^{\intercal}$ von M mit einer unteren Dreiecksmatrix L mit positiven Diagonalelementen.

S2 (Berechnung von $c = A^{\mathsf{T}}b$ und Lösung von $LL^{\mathsf{T}}x = c$): S2.1: Berechne $c := A^{\mathsf{T}}b$ S2.2: Löse die Dreieckssysteme $Ld = c, L^{\mathsf{T}}x = d$ Aufwand: S1: $\sim (m + n/3) n^2/2$ opms, S2: $\sim (m + n) n$ opms

9.1.2. Bemerkung. (i) Für praktische Probleme ist m meist deutlich größer als n, so daß der Gesamtaufwand durch den Term $\sim mn^2/2$, d. h. durch die Berechnung von M bestimmt wird. Es ist zweckmäßig, ein Dreieck von M unter Verwendung von zusätzlichen $\sim n^2/2$ S gesondert zu speichern; die Originalmatrix A steht dann weiter zur Verfügung, etwa zur iterativen Verbesserung. Falls ein diagonalgewichtetes Problem $DAx \cong Db$ zu lösen ist, erfordert die Berechnung von DA und Db weitere $\sim mn$ opms, spielt also keine Rolle.

(ii) Die Normalgleichungen (1) können auch nach einem anderen Verfahren gelöst werden, etwa durch schrittweise Diagonalisierung nach dem Gauß-Jordan-Verfahren, vgl. 6.4. Dies empfehlen manche Autoren für die schrittweise Regression, siehe B 8.9.

(iii) In exakter Arithmetik gilt $M = A^{\mathsf{T}}A$, also

$$\|A\|_{2}^{2} = \|M\|_{2} \leq \sum_{i=1}^{n} (M)_{ii} = \|A\|_{F}^{2},$$

$$\|A^{+}\|_{2}^{2} = \|M^{-1}\|_{2} \leq \sum_{i=1}^{n} (M^{-1})_{ii} = \|A^{+}\|_{F}^{2}.$$
(2)

9.1.3. Rundungsfehleranalyse. Für spaltenreguläres $A \in \Re^{m,n}$ und beliebiges $b \in \Re^m$ ist Algorithmus 9.1.1 durchführbar, sofern A im Sinne von

$$\hat{\varkappa} := \nu N_1 [\text{cond } (A)]^2 \leq 0.5 \quad \text{mit} \quad N_1 := mn + F_1 \sim mn \tag{3}$$

nicht zu schlecht konditioniert ist, wobei

$$F_1 := n^{3/2} + n + n^{1/2} F \sim 2n^{3/2}, \qquad F := n + 1 + 0.5 \ln n \sim n$$

wie in 6.1.5(ii) gesetzt ist. Die berechneten Größen M, L und c, x genügen den Beziehungen

$$M = A^{\mathsf{T}}A + \delta_0 M, \qquad \|\delta_0 M\| \leq \nu m n \, \|A\|^2, \tag{4}$$

$$\boldsymbol{c} = (\boldsymbol{A} + \boldsymbol{\delta}_{\boldsymbol{0}} \boldsymbol{A})^{\mathsf{T}} \boldsymbol{b}, \qquad \|\boldsymbol{\delta}_{\boldsymbol{0}} \boldsymbol{A}\| \leq v m \sqrt{n} \|\boldsymbol{A}\|, \tag{5}$$

$$\boldsymbol{M} + \boldsymbol{\delta}_1 \boldsymbol{M} = \boldsymbol{L} \boldsymbol{L}^{\mathsf{T}}, \qquad \|\boldsymbol{\delta}_1 \boldsymbol{M}\| \leq \nu n F \|\boldsymbol{M}\| \sim \nu n^{3/2} \|\boldsymbol{M}\|, \tag{6}$$

$$(M + \sigma_2 M) x = c, \qquad \|\sigma_2 M\| \le rF_1 \|M\| \sim r2n^{3/2} \|M\|.$$
(7)

Beweis. Aus (2.3.43) folgt $|M - A^{\mathsf{T}}A| \leq \nu m |A^{\mathsf{T}}| |A|$, mit (2.3.30) also $||\delta_0 M|| \leq ||\delta_0 M||_F \leq \nu m ||A||_F^2 \leq \nu m n ||A||^2$. Analog ergibt sich (5) aus (2.3.37). Wegen (3), (4) gilt nun

$$||(A^{\mathsf{T}}A)^{-1}|| ||\boldsymbol{\delta}_{\boldsymbol{0}}\boldsymbol{M}|| \leq \nu mn [\text{cond } (A)]^{2} =: \varkappa \leq \hat{\varkappa} \leq 0.5,$$

so daß M nach Ü 4.1.9 positiv definit ist, und aus 4.1.2 folgt $||M^{-1}|| \leq ||(A^{\top}A)^{-1}||/(1-\varkappa)$ Mit $||M|| \leq (1 + \nu mn) ||A||^2 = ||A||^2$ zieht dies cond $(M) \leq [\text{cond} (A)]^2/(1-\varkappa)$ nach sich Damit erhalten wir

$$\begin{split} \nu F \|\boldsymbol{M}^{-1}\|_2 \|\boldsymbol{M}\|_F &\leq \nu \; \forall n \; F \; \text{cond} \; (\boldsymbol{M}) \leq \nu F_1 \; \text{cond} \; (\boldsymbol{M}) \leq \nu F_1 [\text{cond} \; (\boldsymbol{A})]^2 / (1-\varkappa) \\ &= (\hat{\boldsymbol{x}} - \varkappa) / (1-\varkappa) \leq (0.5 - \varkappa) / (1-\varkappa) \leq 0.5. \end{split}$$

Die Bedingung (6.1.7) ist daher für M statt A erfüllt, und (6), (7) ergeben sich dann sofort aus 6.1.3 und 6.1.5 (ii).

9.1.4. Bemerkung. (i) Die Bedingung (3) garantiert, daß die berechnete Matrix M positiv definit und die Cholesky-Faktorisierung durchführbar ist. Eine Mindestforderung an A ist dagegen, daß die Lösung x stetig von Störungen in der Größenordnung des Darstellungsfehlers abhängt. Für spaltenreguläres A führt dies nach 8.2.7 auf $||A^+|| || \delta A_D|| < 1$. Da der Darstellungsfehler δA_D nach 2.3.15 durch $||\delta A_D|| \leq v \sqrt{n} ||A||$ abgeschätzt werden kann, ergibt sich

$$\nu \sqrt{n} \operatorname{cond} (A) < 1 \tag{8}$$

als Bedingung dafür, daß das Problem $Ax \simeq b$ bezüglich des Fehlerniveaus ν korrekt gestellt ist. Das Auftreten des Faktors $[\text{cond}(A)]^2$ in (3) schränkt also die Unterklasse der vernünftig gestellten Probleme, für die 9.1.1 durchführbar ist, wesentlich ein.

(ii) Das folgende Beispiel zeigt, daß auf eine Bedingung des Typs (3) nicht verzichtet werden kann. Für

$$egin{aligned} egin{aligned} A = egin{pmatrix} 1 & 1 & 1 \ arepsilon & 0 & 0 \ 0 & arepsilon & 0 \ 0 & 0 & arepsilon \end{pmatrix} & ext{ist} \quad egin{matrix} M^{st} = A^{\intercal}A = egin{pmatrix} 1 + arepsilon^2 & 1 & 1 \ 1 & 1 + arepsilon^2 & 1 \ 1 & 1 & 1 + arepsilon^2 & 1 \ 1 & 1 & 1 + arepsilon^2 \end{pmatrix}, \end{aligned}$$

und es gilt $[\operatorname{cond} (A)]^2 \sim 3/\varepsilon^2$, denn $A^{\mathsf{T}}A$ hat die Eigenwerte $\{3 + \varepsilon^2, \varepsilon^2, \varepsilon^2\}$. Für $\varepsilon \in \mathfrak{R}$ mit $10r^2 \leq \varepsilon^2 < r/2$ ist (3) verletzt und $M = \operatorname{fl} (A^{\mathsf{T}}A)$ wegen fl $(1 + \varepsilon^2) = 1$ singulär, aber A erfüllt (8). \Box

9.1.5. Folgerung. Unter den Voraussetzungen von 9.1.3 genügt der erzeugte Rundungsfehler $\delta x = x - x^*$, $x^* = A^+ b = (A^{\mathsf{T}}A)^{-1} A^{\mathsf{T}}b$, der Abschätzung

$$\|\delta \boldsymbol{x}\| \leq \frac{\nu [\text{cond } (\boldsymbol{A})]^2}{1-\lambda} \left\{ N_1 \|\boldsymbol{x}^*\| + N_2 \frac{\|\boldsymbol{b}\|}{\|\boldsymbol{A}\|} \right\}, \quad N_2 := m \sqrt{n}.$$
(9)

Beweis. Es gilt nach 9.1.3

 $(A^{\mathsf{T}}A + \delta M) \boldsymbol{x} = A^{\mathsf{T}}\boldsymbol{b} + \delta \boldsymbol{c} \quad \text{mit} \quad \delta M := \delta_0 M + \delta_2 M, \qquad \delta \boldsymbol{c} := \delta_0 A^{\mathsf{T}}\boldsymbol{b}, \quad (10)$ wobei $\|\delta M\| \leq \nu N_1 \|A\|^2$ und $\|\delta \boldsymbol{c}\| \leq \nu m \sqrt{n} \|A\| \|b\|$ ist. Dies führt wegen $\|(A^{\mathsf{T}}A)^{-1}\| \|\delta M\| \leq \nu N_1 [\operatorname{cond} (A)]^2 = \hat{\boldsymbol{x}} < 1$ und 4.1.3 sofort auf (9). \Box

Die Abschätzung (9) ist in dem Sinne realistisch, daß eine Schranke der Form

$$\|\boldsymbol{\delta x}\| \leq \frac{\nu N[\text{cond } (\boldsymbol{A})]^2}{1-\hat{\varkappa}} \left\{ \|\boldsymbol{x}^*\| + \frac{\|\boldsymbol{b}\|}{\|\boldsymbol{A}\|} \right\} =: \frac{\Delta \boldsymbol{x}_{\text{gen}}'(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\nu})}{1-\hat{\varkappa}}$$
(11)

mit $1 \leq N \leq N_1$ sicher angenommen wird.

Wir wollen im folgenden untersuchen, ob 9.1.1 ein numerisch stabiles Verfahren ist. Die zum Fehlerniveau v gehörenden Darstellungsfehler sind nach 2.1.13 durch $\|\delta A\| \leq v \sqrt{n} \|A\|$, $\|\delta b\| \leq v \|b\|$ beschränkt. Für die zugehörigen linearisierten Störungen $\delta x'$ gilt nach 8.2.7

$$\|\boldsymbol{\delta x'}\| \leq \boldsymbol{\nu} \text{ cond } (\boldsymbol{A}) \left\{ \sqrt{n} \left[\|\boldsymbol{x^*}\| + \text{ cond } (\boldsymbol{A}) \frac{\|\boldsymbol{r}\|}{\|\boldsymbol{A}\|} \right] + \frac{\|\boldsymbol{b}\|}{\|\boldsymbol{A}\|} \right\} =: \Delta \boldsymbol{x'_{\min}}(\boldsymbol{A}, \boldsymbol{b}, \boldsymbol{\nu}),$$
(12)

und diese Schranke ist realistisch. Mit der Schranke $\Delta x'_{gen}$ aus (11) für den linearisierten erzeugten Rundungsfehler $\delta x'_{gen}$ folgt dann

$$rac{arDelta oldsymbol{x}'_{ ext{gen}}(oldsymbol{A},oldsymbol{b},oldsymbol{v})}{arDelta oldsymbol{x}'_{ ext{min}}(oldsymbol{A},oldsymbol{b},oldsymbol{v})} \geqq rac{N}{\sqrt{n}} ext{ cond } (oldsymbol{A}) igg/ igg\{1+rac{\|oldsymbol{r}\|_{ ext{min}}(oldsymbol{A},oldsymbol{b},oldsymbol{v})}{\|oldsymbol{A}\|\|\|oldsymbol{A}^{+}oldsymbol{b}\|\|+\|oldsymbol{b}\|\|} igg\}.$$

Die untere Schranke kann so groß wie $\sim 1/\sqrt{\nu}$ werden, denn in der durch (3) charakterisierten Unterklasse spaltenregulärer Quadratmittelprobleme gibt es solche, für die cond $(A) \sim 1/\sqrt{\nu}$ und trotzdem

$$\|\boldsymbol{r}\| \text{ cond } (A) \leq N_3(\|A\| \|A^+ \boldsymbol{b}\| + \|\boldsymbol{b}\|)$$

gilt, nämlich die mit genügend kleinem ||r|| oder die mit $||A^+b||$ in der Größenordnung von $||A^+|| ||b||$. Daher kann $\Delta x'_{\text{gen}}/\Delta x'_{\text{min}}$ nicht durch eine von A und ν unabhängige Zahl F nach oben beschränkt werden. Da die linearisierten Schranken für $\delta A \rightarrow O$ die korrekten Schranken beliebig genau approximieren, liegt keine Stabilität vor.

9.1.6. Aussage. Über der Klasse der spaltenregulären Quadratmittelprobleme, die der Bedingung (3) genügen, ist Algorithmus 9.1.1 nicht numerisch stabil.

9.1.7. Bemerkung. (i) Aus den angegebenen Abschätzungen folgt andererseits, daß 9.1.1 über der Unterklasse der spaltenregulären Quadratmittelprobleme, für die neben (3) noch

$$\|\boldsymbol{r}\| \ge \|\boldsymbol{b}\|/c_1 \quad \text{und} \quad \|\boldsymbol{x}^*\| = \|\boldsymbol{A}^+\boldsymbol{b}\| \le c_2 \|\boldsymbol{b}\|/\|\boldsymbol{A}\|$$
(13)

mit Konstanten $c_1, c_2 > 0$ gilt, der Quotient $\Delta x'_{gen} / \Delta x'_{min}$ durch

$$F \le c_1 (1 + c_2) N_1 \tag{14}$$

beschränkt ist. Dies bedeutet: Für stark inkonsistente Probleme mit kleinen Lösungen ist 9.1.1 numerisch stabil; es kann sogar numerische Gutartigkeit gezeigt werden. Aus (13) folgt

$$K_{\boldsymbol{A}} \ge \operatorname{cond} (\boldsymbol{A}) + [\operatorname{cond} (\boldsymbol{A})]^2 / (c_1 c_2), \tag{15}$$

d. h., die Unterklasse besteht aus denjenigen Quadratmittelproblemen, für die K_A durch [cond (A)]² bestimmt wird, vgl. 8.2.8(ii).

(ii) Wenn z und x wie in 5.4.7, Fall 3, bestimmt werden (dort ist L_2 durch L zu ersetzen), gilt

$$\zeta_2 := \|\boldsymbol{z}\|_2 / \|\boldsymbol{x}\|_2 \le \|(\boldsymbol{L}^{\mathsf{T}}\boldsymbol{L})^{-1}\| = \|\boldsymbol{M}^{-1}\| = \|\boldsymbol{A}^+\|^2.$$
(16)

Andererseits ist

$$\eta_2 := \sum_{i=1}^n (M)_{ii} = \|L\|_F^2 = \|A\|_F^2 \le n \, \|A\|_F^2, \tag{17}$$

so daß $\zeta_2 \eta_2/n$ eine untere Schranke für [cond (A)]² darstellt, welche die Größenordnung richtig wiedergibt. Falls die für die Gültigkeit von (3) notwendige und leicht überprüfbare Bedingung

$$\hat{lpha}:=
u m \zeta_2 \eta_2 < 1$$

nicht erfüllt ist, muß damit gerechnet werden, daß der nach 9.1.1 berechnete Cholesky-Faktor L wie auch die Lösung x selbst unbrauchbar sein können.

(iii) Die Nachteile von 9.1.1 lassen sich beheben, wenn die Schritte S1 und S2 bei einfachgenauen Eingangsdaten $\{A, b\}$ mit höherer Genauigkeit $v_1 \leq v^2$ realisiert werden und x danach wieder auf einfache Genauigkeit v gerundet wird. In diesem Fall kann $[\operatorname{cond}(A)]^2$ in (3) durch $\operatorname{cond}(A)$ ersetzt werden, und das berechnete x ist Lösung von $(A + dA) x \cong b + db$ mit $||dA|| \leq vF ||A||, ||db|| \leq vF ||b||,$ $F \sim mn$, d. h., es liegt numerische Gutartigkeit vor. Dazu sind nur $\sim n^2/2$ doppeltgenaue Speicherplätze erforderlich, allerdings erhöht sich der arithmetische Aufwand beträchtlich.

(iv) Für steife diagonalgewichtete Probleme $DAx \cong Db$ sind Normalgleichungsverfahren nicht geeignet. Als Modellfall betrachten wir die Gewichtsmatrix D $=\left(\frac{wI_{m1}}{|I_{m2}}\right)$ und das entsprechend partitionierte $A = \left(\frac{A_1}{A_2}\right)$ mit $w \gg 1$ und $1 \leq m_1 < n$. Die Matrix der Normalgleichungen ist dann

$$(\boldsymbol{D}\boldsymbol{A})^{\mathsf{T}} \boldsymbol{D}\boldsymbol{A} = w^2 A_1^{\mathsf{T}} A_1 + A_2^{\mathsf{T}} A_2;$$

sie unterscheidet sich für großes w nicht von der singulären Matrix $w^2A_1^{\mathsf{T}}A_1$.

Übungsaufgaben

Ü 9.1.1. Es sei $A \in \mathbb{R}^{m,n}$ spaltenregulär, und die Cholesky-Faktorisierung $A^{\mathsf{T}}A = LL^{\mathsf{T}}$ sei bekannt. Mit C werde $(A^{\mathsf{T}}A)^{-1} = L^{-\mathsf{T}}L^{-1}$ bezeichnet, vgl. 8.3. Man zeige: (i) Für $F \in \mathbb{R}^{n,l}$ kann $X := F^{\mathsf{T}}C^{-1}F = F^{\mathsf{T}}(A^{\mathsf{T}}A)F$ gemäß

$$\mathbf{X} = \mathbf{F}^{\mathsf{T}} \mathbf{L} \mathbf{L}^{\mathsf{T}} \mathbf{F} = \mathbf{\overline{F}}^{\mathsf{T}} \mathbf{\overline{F}} \quad \text{mit} \quad \mathbf{\overline{F}} := \mathbf{L}^{\mathsf{T}} \mathbf{F}$$
(18)

in $\sim (n+l) \ln/2$ opms berechnet werden, während die direkte Berechnung gemäß

$$X = F^{\mathsf{T}} A^{\mathsf{T}} A F = \hat{F}^{\mathsf{T}} \hat{F} \quad \text{mit} \quad \hat{F} := AF$$

 $\sim (2n + l) lm/2$ opms erfordert, also mindestens das (m/n)-fache des Aufwands von (18).

275

(ii) Für $B \in \mathbb{R}^{l,n}$ kann $Y := BCB^{\mathsf{T}}$ ohne explizite Bildung von $C = (A^{\mathsf{T}}A)^{-1}$ gemäß

$$Y = BL^{-\mathsf{T}}L^{-1}B = \overline{B}\overline{B}^{\mathsf{T}} \quad \text{mit} \quad \overline{B} := BL^{-\mathsf{T}}$$
(19)

berechnet werden, wobei \overline{B} aus dem Dreieckssystem

$$\overline{B}L^{\mathsf{T}} = B$$
 bzw. $L\overline{B}^{\mathsf{T}} = B^{\mathsf{T}}$

bestimmt wird. Dies erfordert insgesamt $\sim (n+l) \ln/2$ opms. Für $\boldsymbol{B} := \boldsymbol{A}$ ist

 $Y = A(A^{\mathsf{T}}A)^{-1}A = (AL^{-\mathsf{T}})(AL^{-\mathsf{T}})^{\mathsf{T}}$

der Projektor P_A auf $\mathcal{R}(A)$.

Ü 9.1.2. Man beweise: Wenn

$$A := A + uv^{\mathsf{T}}, \quad u \neq o, \quad v \neq o, \tag{20}$$

durch Rang-1-Änderung aus A entsteht, gilt

$$\bar{A}^{\mathsf{T}}\bar{A} = A^{\mathsf{T}}A - yy^{\mathsf{T}} + zz^{\mathsf{T}} \quad \text{mit} \quad y := A^{\mathsf{T}}u/||u||, \qquad z := y + ||u||v, \tag{21}$$

d, h., $\bar{A}^{\mathsf{T}}\bar{A}$ entsteht aus $A^{\mathsf{T}}A$ durch eine Rang-2-Änderung, die in der symmetrischen Form (21) als Differenz zweier positiv semidefiniter Rang-1-Terme geschrieben werden kann.

9.2. **Iterative Verbesserung**

Im vorangegangenen Abschnitt wurde gezeigt, daß das Normalgleichungsverfahren 9.1.1 zwei wesentliche Nachteile aufweist, nämlich:

- Die Klasse der durch 9.1.1 lösbaren spaltenregulären Quadratmittelprobleme wird durch die Bedingung (3) wesentlich eingeschränkt.
- Das Verfahren ist numerisch instabil.

Wir zeigen im folgenden, daß der zweite Nachteil durch iterative Verbesserung analog zu 5.4.3 behoben werden kann; der erste Nachteil bleibt dabei leider bestehen.

9.2.1. Normalgleichungsverfahren mit iterativer Verbesserung

Aufgabe: Für gegebenes spaltenreguläres $A \in \Re^{m,n}$ sei $L \in \Re^{n,n}$ der nach 9.1.1 berechnete Cholesky-Faktor von $M := A^{\mathsf{T}}A$. Unter Verwendung von L ist für gegebenes $b \in \Re^m$ die Lösung des Quadratmittelproblems $Ax \simeq b$ über die Normalgleichungen zu berechnen und iterativ zu verbessern.

S0: $x^0 = o$ for $k := 0(1)k_{\max}$ do

- S1: Berechne $r^k := \operatorname{fl}(b Ax^k)$ und $g^k := \operatorname{fl}(A^{\mathsf{T}}r^k)$
- S2: Berechne $h = h^k$ als Lösung von $Mh = g^k$ unter Verwendung der Cholesky-Faktorisierung $M = LL^{\intercal}$ von M aus den Dreieckssystemen $Lf = g^k$, $L^{\mathsf{T}}h = f$ Setze $x^{k+1} := x^k + h^k$

S3: Setze
$$x^{k+1} := x^k + h^k$$

Aufwand: $\sim (2m + n) n$ opms pro Schritt $k \geq 1$, $\sim (m + n) n$ opms für k = 0.

Für k = 0 ist $x^1 = h^0$ gerade die in 9.1.1 berechnete Lösung x.

9.2.2. Fehleranalyse. Die spaltenreguläre Matrix $A \in \Re^{m,n}$ genüge der Bedingung

$$ar{m{\kappa}} := v \overline{N}_1 [ext{cond} (A)]^2 \leq 0.23, \qquad \overline{N}_1 := (m+1) (n+1) + 2n^2.$$

$$\tag{1}$$

Dann ist 9.2.1 für jedes $\boldsymbol{b} \in \Re^m$ durchführbar, und die berechneten \boldsymbol{x}^k (k = 1, ..., k_{\max}) sind Lösungen der gestörten Quadratmittelprobleme

$$(A + \delta A_k) \, \boldsymbol{x}^k \simeq \boldsymbol{b} \,, \tag{2}$$

wobei

$$\|\boldsymbol{\delta}\boldsymbol{A}_{k}\| \leq 1.3\nu \left(\sqrt{m^{2}n + n^{3}} + \bar{\boldsymbol{x}}^{k-1} \, \overline{N}_{1} \operatorname{cond} (\boldsymbol{A})\right) \|\boldsymbol{A}\|$$
(3)
ist. Mit $\boldsymbol{x}^{*} = \boldsymbol{A}^{+}\boldsymbol{b}, \, \boldsymbol{r}^{*} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}^{*}$ gelten die Fehlerabschätzungen
$$\|\boldsymbol{x}^{k} - \boldsymbol{x}^{*}\| \leq \bar{\boldsymbol{x}}^{k} \|\boldsymbol{x}^{*}\| + 1 \, \Im \, \operatorname{cond} \, (\boldsymbol{A}) \, ||\boldsymbol{x}^{*}\| + mn^{1/2} \operatorname{cond} \, (\boldsymbol{A}) \, ||\boldsymbol{x}^{*}\| / ||\boldsymbol{A}||$$

$$\|\boldsymbol{x}^{k} - \boldsymbol{x}^{*}\| \leq \bar{\boldsymbol{x}}^{k} \|\boldsymbol{x}^{*}\| + 1.3\nu \text{ cond } (\boldsymbol{A}) \{ n^{3/2} \|\boldsymbol{x}^{*}\| + mn^{1/2} \text{ cond } (\boldsymbol{A}) \|\boldsymbol{r}^{*}\| / \|\boldsymbol{A}\| \}.$$
(4)

Der Beweis ist sehr technisch und soll übergangen werden; er benutzt Lemma 8.2.11.

9.2.3. Bemerkung. (i) Für genügend großes k folgt aus (2), (3) die numerische Gutartigkeit, aus (4) die numerische Stabilität des Verfahrens 9.2.1. Da dies für k = 1, d. h. für die nicht verbesserte Lösung gemäß 9.1.1, nicht gefolgert werden kann, vgl. 9.1.6, spielt die iterative Verbesserung für das Normalgleichungsverfahren 9.1.1 eine andere Rolle als für die Lösung regulärer Gleichungssysteme nach dem Gaußschen Algorithmus in 5.4. Für das Verfahren 9.1.1 ist sie erforderlich, um überhaupt Gutartigkeit zu garantieren, während sie beim Gaußschen Algorithmus die für das Grundverfahren ohnehin vorliegende Gutartigkeit lediglich verbessert.

(ii) Als Abbruchkriterium kann wie in 5.4.5 die Bedingung

 $\|m{h}^{k}\| > \|m{h}^{k-1}\|/2$

gewählt werden. Falls die Korrekturen nicht schnell genug kleiner werden, ist entweder die Grenzgenauigkeit erreicht, oder (1) ist nicht erfüllt, so daß die iterative Verbesserung ohnehin fragwürdig ist. Der letzte Fall kann durch Schätzung von $\bar{z} \sim \hat{x} \approx \tilde{z}$ analog zu 9.1.7(ii) diagnostiziert werden. Die Konvergenzgeschwindigkeit hängt entscheidend von $\bar{\varkappa}$ ab.

(iii) Wenn M und L mit erhöhter Genauigkeit $v_1 \leq v^2$ berechnet und L danach wieder auf einfache Genauigkeit v gerundet wird, kann \bar{x} in 9.2.2 durch

$$\vec{\varkappa}_1 := \nu \overline{N}_1 \text{ cond } (A)$$

ersetzt werden.

(iv) Wenn die Residuen r^k und g^k mit höherer Genauigkeit $v_1 \leq r^2$ berechnet und danach wieder auf einfache Genauigkeit v gerundet werden, liefert 9.2.2 für genügend großes k die volle relative Genauigkeit im Sinne von

$$\|oldsymbol{x^k} - oldsymbol{x^*}\| \leq
u N_3 \left\{\|oldsymbol{x^*}\| + \operatorname{cond}\left(A\right)\|oldsymbol{r^*}\|/\|A\|\right\}, \qquad N_3 pprox 1.$$

Ob diese Genauigkeit nötig und sinnvoll ist, hängt von der Genauigkeit der Eingangsdaten ab. \Box

Bemerkungen zum Kapitel 9

B 9.1. Die Lösung von Quadratmittelproblemen über die Normalgleichungen geht bereits auf GAUSS zurück. Die Normalgleichungen werden deshalb auch *Gaußsche Normalgleichungen* und deren Koeffizientenmatrix $A^{T}A$ *Gaußsche Transformierte* von A genannt. In der statistischen Literatur wird auch fast ausschließlich dieser klassische Weg beschrieben, obwohl er vom numerischen Standpunkt her deutliche Nachteile gegenüber den moderneren Orthogonalisierungsverfahren aufweist.

B 9.2. Die Fehleranalyse der iterativen Verbesserung gemäß 9.2.2 und speziell der Nachweis der numerischen Gutartigkeit für genügend großes k geht auf KIEŁEASIŃSKI zurück und scheint neu zu sein. Ein verwandtes Verfahren hat BJÖRCK [78] untersucht. In der letztgenannten Arbeit sind auch instruktive numerische Beispiele zu finden.

10. Orthogonalisierungsverfahren

Orthogonalisierungsverfahren zur Lösung von Quadratmittelproblemen $Ax \cong b$ benutzen die Darstellung von $A \in \mathbb{R}^{m,n}$ gemäß A = QR als Produkt einer spaltenorthonormalen Matrix $Q \in \mathbb{R}^{m,l}$ und einer oberen Dreiecks- bzw. Trapezmatrix $R \in \mathbb{R}^{l,n}$; man spricht von einer QR-Faktorisierung von A. Die verschiedenen Varianten unterscheiden sich in der Wahl der Formate — üblich sind l = n oder l = m — und in den Konstruktionsvorschriften für Q und R. Unter Verwendung einer solchen QR-Faktorisierung lassen sich die Quadratmittellösungen wie auch andere interessierende Größen, etwa P_A oder A^+ , in einfacher Weise bestimmen.

10.1. Orthogonalisierung nach Gram-Schmidt

Wir setzen $A \in \mathbb{R}^{m,n}$ als spaltenregulär voraus und suchen eine Faktorisierung

$$\boldsymbol{A} = \boldsymbol{Q}\boldsymbol{R} \tag{1}$$

mit einer spaltenorthonormalen Matrix $Q \in \mathbb{R}^{m,n}$ und einer oberen Dreiecksmatrix $R \in \mathbb{R}^{n,n}$, d. h., es gelte

$$Q^{\mathsf{T}}Q = I_n$$
 und $R = (r_{ij})$ mit $r_{ij} = 0$ $(i > j; i, j = 1, ..., n)$. (2)

Wegen der Spaltenregularität von A muß R regulär sein, also $r_{ij} \neq 0$ (j = 1, ..., n) gelten. Wenn die Spalten von A bzw. Q gemäß

$$A = (a^1, ..., a^n), \qquad Q = (q^1, ..., q^n)$$

,

mit a^{j} bzw. q^{j} bezeichnet werden, liest sich (1) spaltenweise als

$$a^{j} = \sum_{i=1}^{j} q^{i} r_{ij} = q^{1} r_{1j} + q^{2} r_{2j} + \dots + q^{j} r_{jj} \qquad (j = 1, \dots, n),$$
(3)

d. h., die Elemente der *j*-ten Spalte von \mathbf{R} sind gerade die Koeffizienten von \mathbf{a}^{j} in der orthonormalen Basis $\{\mathbf{q}^{1}, ..., \mathbf{q}^{j}\}$. Umgekehrt gilt wegen der Regularität von \mathbf{R} auch $\mathbf{Q} = A\tilde{\mathbf{R}}$ mit der oberen Dreiecksmatrix $\tilde{\mathbf{R}} = \mathbf{R}^{-1}$, also die zu (3) analoge Darstellung von \mathbf{q}^{j} als Linearkombination von $\mathbf{a}^{1}, ..., \mathbf{a}^{j}$. Hieraus folgt

span
$$\{a^1, ..., a^j\}$$
 = span $\{q^1, ..., q^j\}$ $(j = 1, ..., n);$ (4)

insbesondere ist also $\{q^1, \ldots, q^n\}$ eine orthonormale Basis von $\mathcal{R}(A)$. Andererseits folgt aus (4) die Gültigkeit der Darstellung (3). Dies bedeutet: Die Faktorisierung von A gemäß (1), (2) ist gleichbedeutend mit der Konstruktion einer orthonormalen Basis $\{q^1, \ldots, q^n\}$, für die (4) gilt. Daß eine solche Basis existiert, wird weiter unten konstruktiv bewiesen.

Unter Verwendung von (1), (2) läßt sich A^+ nach 8.1.9 in der Form

$$A^{+} = R^{+}Q^{+} = R^{-1}Q^{\mathsf{T}}$$
⁽⁵⁾

darstellen. Die eindeutige Lösung von $Ax \simeq b$ ist dann durch

$$\boldsymbol{x} = \boldsymbol{A}^{+}\boldsymbol{b} = \boldsymbol{R}^{-1}\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{b} \tag{6}$$

gegeben, und das zugehörige minimale Residuum ist

$$\boldsymbol{r} = \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x} = \boldsymbol{b} - \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{b}. \tag{7}$$

Das Problem $Ax \simeq b$ kann daher wie folgt gelöst werden.

- 10.1.1. Basisalgorithmus zur Lösung von $Ax \simeq b$ mittels Faktorisierung A = QR
- S1 (*QR*-Faktorisierung von *A*): Bestimme $Q \in \mathbb{R}^{m,n}$ mit $Q^{\mathsf{T}}Q = I_n$ und eine obere Dreiecksmatrix $R \in \mathbb{R}^{n,n}$, so daß A = QR gilt.
- S2 (Berechnung der zur rechten Seite **b** gehörenden Lösung $x = R^{-1}Q^{\mathsf{T}}b$ und des Residuums $r = b QQ^{\mathsf{T}}b$):
- S2.1: Berechne $c = Q^{\mathsf{T}}b$ und r = b Qc
- S2.2: Bestimme x als Lösung des Dreieckssystems Rx = c.

Zur Berechnung der QR-Faktorisierung gehen wir von (3) aus und nehmen an, daß bereits ein System orthonormaler Vektoren $\{q^1, ..., q^{k-1}\}$, für das (4) für j = 1, ..., k - 1 gilt, bestimmt worden ist. Wir suchen dann q^k in der Form

$$oldsymbol{q}^k = oldsymbol{p}^k/r_{kk}$$
 mit $oldsymbol{p}^k = r_{kk}oldsymbol{q}^k = oldsymbol{a}^k - \sum_{i=1}^{k-1}oldsymbol{q}^i r_{ik}$

Die Forderung, daß q^k zu $\{q^1, ..., q^{k-1}\}$ orthogonal sei, führt wegen der Orthonormalität von $\{q^1, ..., q^{k-1}\}$ auf

$$0 = r_{kk} q^{j \mathsf{T}} q^k = q^{j \mathsf{T}} a^k - \sum_{i=1}^{k-1} q^{j \mathsf{T}} q^i r_{ik} = q^{j \mathsf{T}} a^k - r_{jk} \quad (j = 1, ..., k-1),$$

so daß bis auf r_{kk} alle Koeffizienten r_{jk} festgelegt sind. Der noch fehlende Wert ergibt sich aus der Normierungsbedingung $\|\mathbf{q}^k\| = \|\mathbf{p}^k\|/|r_{kk}| = 1$ zu

$$r_{kk} = \pm \|\boldsymbol{p}^k\|.$$

Wegen der linearen Unabhängigkeit der a^i muß dabei $p^k \neq o$ gelten, denn andernfalls wäre wegen (4) $a^k \in \text{span } \{a^1, \dots, a^{k-1}\}.$

Zusammenfassend erhalten wir den folgenden, Gram-Schmidt-Orthogonalisierung genannten Algorithmus.

10.1.2. Gram-Schmidt-Orthogonalisierung. Für jede spaltenreguläre Matrix $A = (a^1, ..., a^n) \in \mathbb{R}^{m,n}$ werden durch die Vorschrift

$$p^{k} := a^{k} - \sum_{i=1}^{k-1} r_{ik} q^{i}, \qquad r_{ik} := q^{i} \tau a^{k} \qquad (k = 1, ..., n),$$

$$r_{kk} := \|p^{k}\|, \qquad q^{k} := p^{k} / r_{kk}$$
(8)

in exakter Arithmetik eine spaltenorthonormale Matrix $Q = (q^1, ..., q^n) \in \mathbb{R}^{m,n}$ und eine obere Dreiecksmatrix $\mathbf{R} = (r_{ij}) \in \mathbb{R}^{n,n}$ mit $r_{kk} > 0$ (k = 1, ..., n) erzeugt, so daß A = QR gilt.

Man sagt auch, daß $\{q^1, ..., q^n\}$ durch Gram-Schmidt-Orthogonalisierung der Spalten $\{a^1, ..., a^n\}$ von A entstanden ist.

10.1.3. Bemerkung. (i) Für spaltenreguläres A ist die QR-Faktorisierung (1), (2) bis auf das Vorzeichen der Spalten von Q und entsprechender Zeilen von R eindeutig festgelegt : Aus

$$A=QR=ar{Q}ar{R}$$

mit $Q^{\mathsf{T}}Q = \bar{Q}^{\mathsf{T}}\bar{Q} = I_n$ und oberen Dreiecksmatrizen R, \bar{R} folgt

$$\bar{\boldsymbol{Q}} = \boldsymbol{Q}\boldsymbol{D}, \quad \bar{\boldsymbol{R}} = \boldsymbol{D}\boldsymbol{R} \text{ mit } \boldsymbol{D} = ext{diag}\left(d_{i}\right), \quad d_{i} = \pm 1,$$

siehe Ü 10.1.3.

(ii) Aus (1), (2) ergibt sich

$$A^{\intercal}A = R^{\intercal}Q^{\intercal}QR = R^{\intercal}R$$
 ,

d. h., in exakter Arithmetik ist \mathbf{R}^{T} bis auf das Vorzeichen der Spalten mit dem eindeutig festgelegten Cholesky-Faktor L von $A^{\mathsf{T}}A$ identisch. Für das durch 10.1.2 definierte \mathbf{R} gilt wegen $r_{kk} > 0$ daher $\mathbf{R}^{\mathsf{T}} = L$. \Box

Das folgende Beispiel zeigt, daß der elegante und kompakte Gram-Schmidt-Prozeß in einer katastrophalen Weise numerisch instabil ist.

10.1.4. Beispiel. Es sei $\varepsilon > 0$ und

$$\boldsymbol{A} = \begin{pmatrix} 1 & 1 & 1\\ \varepsilon & 0 & 0\\ 0 & \varepsilon & 0\\ 0 & 0 & \varepsilon \end{pmatrix}$$
(9)

die Matrix aus 9.1.4(ii). Algorithmus 10.1.2 werde so ausgeführt, daß fl $(1 + \varepsilon^2) = 1$ gesetzt. aber alle übrigen Operationen exakt realisiert werden; dies entspricht qualitativ dem Fall $\varepsilon^2 < \nu/2$. Dann ergibt sich

$$Q = \begin{pmatrix} 1 & 0 & 0 \\ \varepsilon & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} & 0 \\ 0 & 0 & \frac{1}{\sqrt{2}} \end{pmatrix}, \qquad R = \begin{pmatrix} 1 & 1 & 1 \\ 0 & \varepsilon \sqrt{2} & 0 \\ 0 & 0 & \varepsilon \sqrt{2} \end{pmatrix}.$$

Wieder in exakter Arithmetik folgt

$$A = QR, \qquad Q^{\mathsf{T}}Q = \begin{pmatrix} 1 + \varepsilon^2 & -\frac{\varepsilon}{\sqrt{2}} & -\frac{\varepsilon}{\sqrt{2}} \\ -\frac{\varepsilon}{\sqrt{2}} & 1 & \frac{1}{2} \\ -\frac{\varepsilon}{\sqrt{2}} & \frac{1}{2} & 1 \end{pmatrix}$$

d. h., Q und R liefern eine exakte Faktorisierung, aber Q ist extrem nichtorthogonal. Die exakten Faktoren sind

$$Q^* = \begin{pmatrix} 1 & \frac{\varepsilon}{\sqrt{2}} & \frac{\varepsilon}{\sqrt{6}} \\ \varepsilon & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} \\ 0 & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} \\ 0 & 0 & \frac{2}{\sqrt{6}} \end{pmatrix} (1 + O(\varepsilon^2)), \quad \mathbf{R}^* = \begin{pmatrix} 1 & 1 & 1 \\ 0 & \varepsilon \sqrt{2} & \frac{\varepsilon}{\sqrt{2}} \\ 0 & 0 & \frac{\varepsilon \sqrt{6}}{2} \end{pmatrix} (1 + O(\varepsilon^2)).$$

Die Schreibweise $\mathbf{F} = \mathbf{G}(1 + O(\varepsilon^2))$ bedeutet dabei $f_{ij} = g_{ij}(1 + O(\varepsilon^2))$ für alle *i*, *j*. Für die rechten Seiten $\mathbf{b}^1 := \mathbf{e}^1$, $\mathbf{b}^2 := \mathbf{e}^3$ ergeben sich mit den berechneten Faktoren Q, \mathbf{R} die numerischen Lösungen x^1, x^2 , mit den exakten Faktoren Q^*, \mathbf{R}^* die exakten Lösungen $x^{*,1}, x^{*,2}$ zu

$$egin{aligned} & m{x}^1 = rac{1}{\epsilon} egin{pmatrix} 1 \ 0 \ 0 \end{pmatrix}, \quad m{x}^{st,1} = egin{pmatrix} 1/3 \ 1/3 \ 1/3 \end{pmatrix} (1 + O(\epsilon^2)), \ & m{x}^2 = rac{1}{\epsilon} egin{pmatrix} -1/2 \ 1/2 \ 0 \end{pmatrix}, \quad m{x}^{st,2} = rac{1}{\epsilon} egin{pmatrix} -1/3 \ 2/3 \ -1/3 \end{pmatrix} (1 + O(\epsilon^2)). \end{aligned}$$

Die allein durch die Rundung fl $(1 + \varepsilon^2) = 1$ generierten relativen Fehler sind $||\delta x^i||/||x^{*,1}|| = 1.4$, $||\delta x^2||/||x^{*,2}|| = 0.5$. Ein Darstellungsfehler von A zum Fehlerniveau ν kann im Fall $\varepsilon^2 = \nu/2$ dagegen nach 8.2.7 nur eine Störung δx^i hervorrufen, für die $||\delta x^i||/||x^{*,i}||$ von der Größenordnung $\varepsilon \approx \nu^{1/2}$ ist. \Box

Die klassische Form 10.1.2 der Gram-Schmidt-Orthogonalisierung ist also für Computerrechnung nicht geeignet.

Um zu einer numerisch brauchbaren Version zu kommen, schreiben wir (8) als

$$\boldsymbol{p}^{k} = \boldsymbol{a}^{k} - \sum_{i=1}^{k-1} \boldsymbol{q}^{i} \boldsymbol{q}^{i\mathsf{T}} \boldsymbol{a}^{k} = (\boldsymbol{I} - \boldsymbol{P}_{k-1}) \, \boldsymbol{a}^{k} \quad \text{m t} \quad \boldsymbol{P}_{k-1} := \sum_{i=1}^{k-1} \boldsymbol{q}^{i} \boldsymbol{q}^{i\mathsf{T}}. \tag{10}$$

Nach Ü 8.1.1 läßt sich der Projektor $I - P_{k-1}$ in der Produktform

$$I - P_{k-1} = (I - q^{k-1}q^{k-1}) \cdots (I - q^2q^{2\mathsf{T}}) (I - q^1q^{1\mathsf{T}})$$
(11)

darstellen. Einsetzen von (11) in (10) und rekursive Auswertung führt dann auf

$$p^{k,1} := a^k,$$

$$r_{ik} := q^{iT} p^{k,i}, \qquad p^{k,i+1} := p^{k,i} - r_{ik} q^i \qquad (i = 1, ..., k - 1),$$

$$p^k := p^{k,k}$$
(12)

als in exakter Arithmetik äquivalente Berechnungsvorschrift für p^{k} . Dabei gilt

$$q^{i^{\intercal}}p^{k,i} = q^{i^{\intercal}}(I - P_{i-1}) a^k = q^{i^{\intercal}}a^k$$
 $(i = 1, ..., k - 1),$

d. h., die in (12) auftretenden Koeffizienten r_{ik} sind mit denen aus (8) identisch, und $p^{k,i}$ ist gerade die (i-1)-te Teilsumme $a^k - r_{1k}q^1 - \cdots - r_{i-1,k}q^{i-1}$ bei der Berechnung von p^k gemäß (8). Im Unterschied zu (8) wird jedoch in (12) der auf q^i projizierte Anteil $r_{ik}q^i$ sofort von a^k abgezogen, bevor die nachfolgende Projektion gebildet wird. Wegen

$$\|\boldsymbol{a}^{k}\|^{2} = \|\boldsymbol{p}^{k,i}\|^{2} + r_{1k}^{2} + \dots + r_{i-1,k}^{2}$$
(13)

werden die $p^{k,i}$ dabei immer kürzer. In Computerarithmetik werden die r_{ik} gemäß (12) genauer berechnet werden als nach (8), denn die Schranke für den erzeugten Rundungsfehler ist im ersten Fall proportional zu $||p^{k,i}||$, im zweiten zu $||a^k||$. Es ist daher zu erwarten, daß der nach (12) berechnete Vektor p^k in geringerem Maße durch Rundungsfehler beeinflußt wird.

Derselbe Trick wird dann auch bei der Berechnung von $c = (c_i) = Q^{\intercal} b$ angewendet und führt auf

$$b^1 := b, \quad c_i := q^{i^{\intercal}} b^i, \quad b^{i-1} := b^i - c_i q^i \quad (i = 1, ..., k - 1).$$
 (14)

In exakter Arithmetik gilt

$$q^{i op} b^i = q^{i op} (I - P_{i-1}) \ b = q^{i op} b$$
 ,

und

$$oldsymbol{b}^{n+1} = (oldsymbol{I} - oldsymbol{P}_n) \,oldsymbol{b} = (oldsymbol{I} - oldsymbol{Q} oldsymbol{Q}^\intercal) \,oldsymbol{b} = oldsymbol{b} - oldsymbol{Q} oldsymbol{c} = oldsymbol{b} - oldsymbol{A} oldsymbol{x} = oldsymbol{r}$$

ist das minimale Residuum, vgl. Ü 4.3.3.

Die Realisierung von 10.1.1 gemäß (12), (14) wird *modifiziertes Gram-Schmidt-*Verfahren — kurz: MGS-Verfahren — genannt. Wir geben eine Realisierung an, bei der p^k und die $p^{k,i}$ auf dem Platz von q^k , die b^i auf dem von r und c auf dem von xgespeichert werden. **10.1.5.** Orthogonalisierung von A und Lösung von $Ax \simeq b$ nach dem modifizierten Gram-Schmidt-Verfahren

Aufgabe: Für gegebenes spaltenreguläres $A = (a^1, ..., a^n) \in \mathbb{R}^{m,n}$ ist eine QR-Faktorisierung von A mit $Q = (q^1, ..., q^n) \in \mathbb{R}^{m,n}$ und einer oberen Dreiecksmatrix $\mathbf{R} = (r_{ij}) \in \mathbb{R}^{n,n}$ nach den modifizierten Gram-Schmidt-Formeln (12) zu berechnen. Mit dieser ist $x \in \mathbb{R}^n$ und $r \in \mathbb{R}^m$ gemäß (6), (7) bei modifizierter Auswertung von $Q^{\mathsf{T}}b$ nach (14) zu bestimmen.

Algorithmus:

S1 (QR-Faktorisierung von A mittels modifizierter Gram-Schmidt-Orthogonalisierung):

for k := 1(1)n do S1.0: $q^k := a^k$ S1.1: for i := 1(1)k - 1 do $[r_{ik} := q^{iT}q^k, q^k := q^k - r_{ik} * q^i]$ S1.2: $r_{kk} := ||q^k||, q^k := q^k/r_{kk}$

S2 (Modifizierte Berechnung von $c = Q^{\dagger}b$ und $r = b - QQ^{\dagger}b$, Lösung von Rx = c):

S2.1:
$$r := b$$

for $i := 1(1)n$ do $[x_i := q^{i^{\intercal}}r, r := r - x_i * q^i]$
S2.2: for $k := n(-1)1$ do
 $\begin{vmatrix} x_k := x_k/r_{kk} \\ \text{for } i := 1(1)k - 1 \text{ do } x_i := x_i - r_{ik} * x_k \\ Aufwand: S1: \sim mn^2 \text{ opms, } S2: \sim \left(2m + \frac{n}{2}\right)n \text{ opms}$

10.1.6. Bemerkung. (i) In-situ-Realisierung durch Überspeichern von A mit Q bzw. von b mit r ist möglich. Dann sind noch $\sim n^2/2$ S für R und n S für x erforderlich.

(ii) Das nichtmodifizierte Gram-Schmidt-Verfahren ergibt sich aus 10.1.5, wenn die Anweisungen $r_{ik} := q^{iT}q^k$ bzw. $x_i := q^{iT}r$ durch $r_{ik} := q^{iT}a^k$ bzw. $x_i := q^{iT}b$ ersetzt werden. Die Programme unterscheiden sich daher fast nicht, insbesondere ist der Aufwand für MGS identisch mit dem für das klassische Verfahren. Die kleinen Unterschiede in den Computerprogrammen haben jedoch große Auswirkungen, denn im Gegensatz zum instabilen Gram-Schmidt-Prozeß ist der modifizierte numerisch gutartig. Dies ist nicht das einzige Beispiel dafür, daß Verfahren, die in exakter Arithmetik äquivalent sind, in Gegenwart von Rundungsfehlern ein völlig unterschiedliches Stabilitätsverhalten aufweisen können.

(iii) In 10.1.5 wird ${\pmb R}$ spaltenweise aufgebaut. Eine zeilenweise Berechnung ist gemäß

for
$$k := 1(1)n$$
 do $p^{k,1} := a^k$
for $k := 1(1)n$ do
 $\begin{vmatrix} r_{kk} := \|p^{k,k}\|, q^k := p^{k,k}/r_{kk} \\ \text{for } j := k + 1(1)n$ do $[r_{kj} := q^{k\top}p^{j,k}, p^{j,k+1} := p^{j,k} - r_{kj} * q^k]$

möglich, wobei alle $p^{j,k}$ selbstverständlich auf dem Platz von q^j bzw. a^j gespeichert werden. Dabei kann eine "Pivotisierung" durch Spaltenvertauschung wie folgt

durchgeführt werden: Zu Beginn des k-ten Schrittes werden die Spalten k und $\hat{s}(k)$ vertauscht, wobei $\hat{s} = \hat{s}(k)$ mit $k \leq \hat{s} \leq n$ durch

$$\| oldsymbol{p}^{\hat{s},oldsymbol{k}} \| = \max \left\{ \| oldsymbol{p}^{oldsymbol{j},oldsymbol{k}} \| \colon oldsymbol{k} \leq oldsymbol{j} \leq n
ight\}$$

festgelegt ist. Dadurch wird r_{kk} maximal unter allen Kandidaten. Die dabei benötigten Normen können wegen (13) gemäß

$$\| \boldsymbol{p}^{j,k+1} \|^2 = \| \boldsymbol{p}^{j,k} \|^2 - (r_{kj})^2$$

rekursiv aus $\|p^{j,1}\| = \|a^j\|$ bestimmt werden, wobei auf Auslöschung zu achten ist, siehe B 10.4. Der zur Normierung verwendete Wert sollte daher nochmals direkt - berechnet werden. In dieser Modifikation liefert 10.1.5 eine **QR**-Faktorisierung

$$\bar{\boldsymbol{A}} = \boldsymbol{A} \boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} = \boldsymbol{Q} \boldsymbol{R} \quad \text{mit} \quad \boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} := \boldsymbol{T}_{1,\hat{\boldsymbol{s}}(1)} \boldsymbol{T}_{2,\hat{\boldsymbol{s}}(2)} \cdots \boldsymbol{T}_{n-1,\hat{\boldsymbol{s}}(n-1)}, \tag{15}$$

für die

$$r_{kk}^2 \ge r_{kj}^2 + r_{k+1,j}^2 + \dots + r_{jj}^2 \quad (j = k+1, \dots, n; k = 1, \dots, n-1)$$
(16)

gilt, insbesondere ist

$$r_{11} \ge r_{22} \ge \cdots \ge r_{nn}$$

Wegen des Auftretens der Permutationsmatrix P_s in (15) müssen die Komponenten x_i der Lösung x im derart modifizierten Algorithmus 10.1.5 am Schluß noch gemäß $x := P_s^{\mathsf{T}} x$ permutiert werden, um die Spaltenvertauschungen wieder rückgängig zu machen.

(iv) Der in (iii) beschriebene pivotisierte MGS-Prozeß kann auch für rangdefizientes A mit rang $(A) = r \leq n \leq m$ angewendet werden. In exakter Arithmetik bricht das Verfahren nach r Schritten wegen $||\mathbf{p}^{j,r+1}|| = 0$ (j = r + 1, ..., n) ab, und $\{q^1, ..., q^r\}$ ist eine orthogonale Basis von $\mathcal{R}(A)$. Bei Computerrealisierung ist ein geeigneter Abbruchtest zu verwenden, siehe Kapitel 11. \Box

10.1.7. Rundungsfehleranalyse. Algorithmus 10.1.5 ist für spaltenreguläres $A \in \mathbb{R}^{m,n}$ und beliebiges $b \in \mathbb{R}^m$ mit $r_{kk} \neq 0$ (k = 1, ..., n) durchführbar, sofern

$$8.6\nu F \text{ cond } (A) < 1 \quad \text{mit} \quad F := 1.4mn^{3/2} \tag{17}$$

gilt. Die berechneten Faktoren Q, R genügen den Beziehungen

$$\boldsymbol{A} + \boldsymbol{\sigma}_{0}\boldsymbol{A} = \boldsymbol{Q}\boldsymbol{R} \quad \text{mit} \quad \|\boldsymbol{\sigma}_{0}\boldsymbol{A}\| \leq \boldsymbol{v}(F/m) \|\boldsymbol{A}\|$$
(18)

und

$$Q^{\mathsf{T}}Q - I = G \quad \text{mit} \quad \|G\| \leq 3\nu F \text{ cond } (A). \tag{19}$$

Die berechnete Lösung x ist exakte Lösung des Quadratmittelproblems

$$(\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \boldsymbol{x} \simeq \boldsymbol{b} + \boldsymbol{\delta}_1 \boldsymbol{b}, \qquad (20)$$

und für das berechnete Residuum r gilt

$$\boldsymbol{r} = (\boldsymbol{b} + \boldsymbol{\delta}_2 \boldsymbol{b}) - (\boldsymbol{A} + \boldsymbol{\delta}_2 \boldsymbol{A}) \boldsymbol{x}.$$
⁽²¹⁾

Die Matrizen $A + \sigma_i A$ sind spaltenregulär, und es gilt

$$\|\boldsymbol{\delta}_{i}\boldsymbol{A}\| \leq vF \|\boldsymbol{A}\|, \qquad \|\boldsymbol{\delta}_{i}\boldsymbol{b}\| \leq v(F/\gamma n) \|\boldsymbol{b}\| \qquad (i=1,2).$$

$$(22)$$

Die Störungsmatrix $\delta_0 A$ ist überdies spaltenweise klein im Sinne von

$$\|\boldsymbol{\delta}_{0}\boldsymbol{a}^{j}\| \leq \nu(1.4j) \|\boldsymbol{a}^{j}\| \qquad (j = 1, ..., n).$$
(23)

Aus Platzgründen verzichten wir auf den nicht ganz einfachen Beweis dieser Aussagen, siehe B 10.1.

Die Analyse 10.1.7 zeigt: Die Berechnung der Quadratmittellösung x und des zugehörigen Residuums r nach dem modifizierten Gram-Schmidt-Verfahren ist ein numerisch gutartiger Proze β .

10.1.8. Bemerkung. (i) Die Klasse der spaltenregulären Quadratmittelprobleme, für die 10.1.5 durchführbar ist, wird durch die Bedingung (17) charakterisiert und besteht daher praktisch aus allen spaltenregulären Problemen, die zum Fehlerniveau v vernünftig gestellt sind, vgl. 9.1.4(i). Insbesondere tritt hier der bei den Normalgleichungsverfahren vorkommende Faktor [cond (A)]² nicht auf. Bei festem Rundungsfehlerniveau ist daher das MGS-Verfahren den Normalgleichungsverfahren in bezug auf die Anwendungsbreite grundsätzlich überlegen. Dies muß allerdings mit einem fast verdoppelten Aufwand bezahlt werden.

(ii) Wenn Q nach dem klassischen Gram-Schmidt-Proze β berechnet wird, kann

$$\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} - \boldsymbol{I}\| \leq \min\left\{n^{3/2}, \nu F[\text{cond}(\boldsymbol{A})]^{n-1}\right\}$$
(24)

gezeigt werden, und diese Schranke ist realistisch. Demgegenüber ist die Schranke

$$\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} - \boldsymbol{I}\| \leq 3\boldsymbol{v}\boldsymbol{F} \operatorname{cond}\left(\boldsymbol{A}\right) < 1 \tag{25}$$

für das nach MGS berechnete Q deutlich besser, hängt aber immer noch von cond (A) ab. Das MGS-Verfahren ist daher in bezug auf den Faktor Q nicht numerisch gutartig. Aus (25) kann bei klassischer Eerechnung von $c = Q^{\mathsf{T}}b$ die Gutartigkeit bezüglich x auch nicht gefolgert werden, vgl. 4.3.5. Für das Vorliegen der Gutartigkeit bezüglich x und r erweist sich die modifizierte Berechnung dieser Größen gemäß 10.1.5 wesentlich; bei dieser wird die durch (25) nicht erfaßte günstige Korrelation der Fehler in Q und R ausgenutzt und ergibt eine von cond (A) unabhängige Kumulationskonstante, siehe B 10.1.

(iii) Ein ausreichend spaltenorthonormales Q kann mit $\sim 2m \cdot n^2$ opms, also doppeltem Aufwand, durch das Gram-Schmidt-Verfahren mit Re-Orthogonalisierung gemäß

$$\mathbf{r} \ k := 1(1)n \ do \\ \mathbf{q}^k := \mathbf{a}^k \\ \text{for } i := 1(1)k - 1 \ do \ [r_{ik} := \mathbf{q}^{i\mathsf{T}}\mathbf{q}^k, \ \mathbf{q}^k := \mathbf{q}^k - r_{ik} * \mathbf{q}^i] \\ \text{for } i := 1(1)k - 1 \ do \ [s_{ik} := \mathbf{q}^{i\mathsf{T}}\mathbf{q}^k, \ \mathbf{q}^k := \mathbf{q}^k - s_{ik} * \mathbf{q}^i, \ r_{ik} := r_{ik} + s_{ik}] \\ r_{kk} := ||\mathbf{q}^k||, \ \mathbf{q}^k := \mathbf{q}^k/r_{kk}$$

fo

berechnet werden. Für die so berechneten Matrizen Q, R gilt (18), und (19) kann zu

$$\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} - \boldsymbol{I}\| \leq \nu F_1 \tag{26}$$

mit F_1 in der Größenordnung von F verschärft werden. Die Abweichung von der Orthonormalität hängt also nicht mehr von cond (A) ab und ist stets klein. Die Berechnung von Q nach dem Gram-Schmidt-Verfahren mit Re-Orthogonalisierung ist daher ein numerisch gutartiger Prozeß. Die Aussagen (20) bis (22) gelten dann unabhängig davon, ob $c = Q^{\mathsf{T}}b$ modifiziert oder klassisch gemäß $c_i = q^{i\mathsf{T}}b$ berechnet wird, siehe wieder 4.3.5.

Übungsaufgaben

Ü 10.1.1. Man zeige, daß das zu beliebigem $x \in \mathbb{R}^n$ gehörende Residuum r = b - Ax unter Verwendung der QR-Faktorisierung (1), (2) von A in der Form

$$\boldsymbol{r} = \boldsymbol{u} + \boldsymbol{v} \tag{27}$$

 \mathbf{mit}

$$\boldsymbol{u} = \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{r} = \boldsymbol{Q}(\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{b} - \boldsymbol{R}\boldsymbol{x}) \perp \boldsymbol{v} := (\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}}) \boldsymbol{r} = (\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}}) \boldsymbol{b}$$
(28)

dargestellt werden kann und daß

$$\|\boldsymbol{r}\|^{2} = \|\boldsymbol{u}\|^{2} + \|\boldsymbol{v}\|^{2} = \|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{b} - \boldsymbol{R}\boldsymbol{x}\|^{2} + \|(\boldsymbol{I} - \boldsymbol{Q}\boldsymbol{Q}^{\mathsf{T}})\,\boldsymbol{b}\|^{2}$$
(29)

gilt. Aus (29) kann der Basisalgorithmus 10.1.1 direkt abgelesen werden.

Ü 10.1.2. Es sei A = QR eine QR-Faktorisierung der spaltenregulären Matrix $A \in \mathbb{R}^{m,n}$. Man überlege sich, daß die Normallösung y des zeilenregulären Systems

$$A^{\mathsf{T}}\boldsymbol{y} = \boldsymbol{f}, \qquad \boldsymbol{f} \in \mathsf{R}^n, \tag{30}$$

wie folgt berechnet werden kann, vgl. auch 4.2.2(i).

S3.1: Bestimme g als Lösung des Dreieckssystems $R^{\intercal}g = f$

S3.2: Berechne
$$\boldsymbol{y} := \boldsymbol{Q}\boldsymbol{g}$$
.

Ü 10.1.3. Man beweise 10.1.3 (i). Hinweis: Man betrachte $\bar{Q}^{\intercal}Q = \bar{R}R^{-1} = S$, zeige $S^{\intercal} = S^{-1}$ und benutze Ü 1.1.12.

Ü 10.1.4. Man rechne Beispiel 10.1.4 für das modifizierte Gram-Schmidt-Verfahren 10.1.5 durch und überprüfe, daß sich dann bis auf einen Faktor $1 + O(\varepsilon^2)$ die exakten Resultate ergeben und daß die Schranke (25) angenommen wird.

10.2. Orthogonalisierung nach Householder und Givens

Wir suchen für die als spaltenregulär vorausgesetzte Matrix $A \in \mathbf{R}^{m,n}$ wieder eine **OR**-Faktorisierung

$$\boldsymbol{A} = \hat{\boldsymbol{Q}}\hat{\boldsymbol{R}},\tag{1}$$

im Unterschied zum Ansatz in 10.1 jedoch mit orthogonalem $\hat{Q} \in \mathbb{R}^{m,m}$ und einer oberen Dreiecksmatrix $\hat{R} \in \mathbb{R}^{m,n}$, d. h., es gelte

$$\hat{Q}^{\mathsf{T}}\hat{Q} = I_m \quad \text{und} \quad \hat{R} = (r_{ij}) \quad \text{mit} \quad r_{ij} = 0 \quad \text{für} \quad i > j$$

 $(i = 1, ..., m; \; j = 1, ..., n).$ (2)

Wenn $\hat{\boldsymbol{O}}$ und $\hat{\boldsymbol{R}}$ gemäß

$$\hat{\boldsymbol{Q}} = (\boldsymbol{Q} \mid \tilde{\boldsymbol{Q}}), \qquad \hat{\boldsymbol{R}} = \left(\frac{\boldsymbol{R}}{\boldsymbol{O}}\right) = \left(\frac{\boldsymbol{\nabla}}{\boldsymbol{O}}\right)$$
(3)

mit $Q \in \mathbb{R}^{m,n}$, $\tilde{Q} \in \mathbb{R}^{m,m-n}$ und $R \in \mathbb{R}^{n,n}$ partitioniert werden, geht (1) in

$$A = QR$$

über. Dies bedeutet: Die Blöcke O und R sind bis auf das Vorzeichen der Spalten von Q und der Zeilen von R eindeutig festgelegt und identisch mit der in diesem Sinne eindeutigen QR-Faktorisierung von A nach (10.1.1), (10.1.2), vgl. 10.1.3. Die Spalten der Matrix \dot{Q} bilden eine orthonormale Basis des orthogonalen Komplements $\mathcal{N}(\mathbf{A}^{\mathsf{T}})$ von $\mathcal{R}(\mathbf{A}) = \mathcal{R}(\mathbf{O})$ und unterliegen keinen weiteren Beschränkungen.

Aus (1) bis (3) folgt nach 8.1.9

$$\boldsymbol{A}^{+} = \boldsymbol{\hat{R}}^{+} \boldsymbol{\hat{Q}}^{\mathsf{T}} = (\boldsymbol{R}^{-1} \mid \boldsymbol{O}) \, \boldsymbol{\hat{Q}}^{\mathsf{T}}, \tag{4}$$

so daß die eindeutige Lösung des Quadratmittelproblems $Ax \simeq b$ durch

$$\boldsymbol{x} = \boldsymbol{A}^{+}\boldsymbol{b} = (\boldsymbol{R}^{-1} \mid \boldsymbol{O}) \, \boldsymbol{\hat{Q}}^{\mathsf{T}} \boldsymbol{b} \tag{5}$$

gegeben ist. Das zugehörige minimale Residuum besitzt die Darstellung

$$\mathbf{r} = \mathbf{b} - A\mathbf{x} = \mathbf{b} - \hat{\mathbf{Q}}\hat{\mathbf{R}}\hat{\mathbf{R}}^{\dagger}\hat{\mathbf{Q}}^{\mathsf{T}}\mathbf{b} = \mathbf{b} - \hat{\mathbf{Q}}\left(\frac{\mathbf{I}_{n} \mid \mathbf{O}}{\mathbf{O} \mid \mathbf{O}}\right)\hat{\mathbf{Q}}^{\mathsf{T}}\mathbf{b} = \hat{\mathbf{Q}}\left(\frac{\mathbf{O} \mid \mathbf{O}}{\mathbf{O} \mid \mathbf{I}_{m-n}}\right)\hat{\mathbf{Q}}^{\mathsf{T}}\mathbf{b}.$$
 (6)

Das Problem $Ax \simeq b$ kann daher nach dem folgenden Basisalgorithmus gelöst werden.

- 10.2.1. Basisalgorithmus zur Lösung von $Ax \simeq b$ mittels Faktorisierung $A = \hat{Q}\hat{R}$. S1 (**QR**-Faktorisierung von A): Bestimme eine orthogonale Matrix $\hat{Q} \in \mathbb{R}^{m,m}$ und eine obere Dreiecksmatrix $\hat{R} = \left(\frac{R}{O}\right) \in \mathbb{R}^{m,n}$, so daß $A = \hat{Q}\hat{R}$ gilt.
- S2 (Berechnung der zur rechten Seite b gehörenden Lösung x und des Residuums r = b - Ax): S2.1: Berechne

$$\hat{\mathbf{c}} = \left(rac{oldsymbol{c}}{oldsymbol{ ilde{c}}}
ight) = \hat{oldsymbol{Q}}^{\intercal}oldsymbol{b} ~~ ext{und} ~~oldsymbol{r} = \hat{oldsymbol{Q}}\left(rac{oldsymbol{o}}{oldsymbol{ ilde{c}}}
ight)$$

S2.2: Bestimme x als Lösung des Dreieckssystems Rx = c.

Wir bemerken, daß die Matrix \hat{Q} für die Berechnung von x und r nicht explizit bekannt zu sein braucht. Es genügt, wenn $\hat{Q}^{\dagger}b$ bzw. $\hat{Q}f$ für gegebenes b bzw. f mit ausreichender Genauigkeit berechnet werden kann. Wegen der Orthogonalität von \hat{Q} gilt ferner $\|r\| = \|\tilde{c}\|$, so daß auf die Berechnung von r verzichtet werden kann, wenn nur nach $\|\boldsymbol{r}\|$ gefragt ist.

Um zu einer Konstruktionsvorschrift für \hat{Q} und \hat{R} zu kommen, schreiben wir (1) in der äquivalenten Form

$$\hat{Q}^{\intercal}A = \hat{R};$$

man beachte, daß \hat{Q} orthogonal ist. Das Auffinden einer Faktorisierung (1), (2) ist also gleichwertig zum Auffinden einer orthogonalen Matrix \hat{Q}^{\intercal} , die A durch Linksmultiplikation auf obere Dreiecksform \hat{R} bringt. Es bietet sich dann an, die Matrix \hat{Q}^{\intercal} als Produkt der im Abschnitt 3.3 eingeführten orthogonalen elementaren Transformationsmatrizen zu konstruieren, also die Matrix A durch eine Folge solcher Transformationen sukzessive auf Dreiecksform zu bringen.

Im folgenden soll gezeigt werden, wie A mittels Householder-Spiegelungen spaltenweise auf Dreiecksform gebracht werden kann. Wenn man vom Typ der Transformationsmatrix absieht, wird dabei wie beim Gaußschen Algorithmus vorgegangen.

Im ersten Schritt wählen wir die Householder-Spiegelung $H_1 = \overline{H}_1$ so, daß

$$A =: A^{(1)} =: M^{(1)} =: (a^1 | B^{(1)})$$

in

$$A^{(2)} = H_1 A^{(1)} = \overline{H}_1 M^{(1)} = (\overline{a}^1 | \overline{B}^{(1)}) = \left(\frac{\varrho_1}{o} \middle| \overline{B}^{(1)}\right) = \left(\frac{r_{11} | r_{12} \dots r_{1n}}{o | M^{(2)}}\right)$$

transformiert wird. Dies ist nach 3.3.6 stets möglich. Im zweiten Schritt wird die analoge Transformation mit \overline{H}_2 auf $M^{(2)}$ angewendet usw. Die nach k-1 Schritten erhaltene Matrix $A^{(k)}$ ist dann in den ersten k-1 Spalten bereits von oberer Dreiecksform

Im k-ten Schritt wird

$$M^{(k)} = egin{pmatrix} a^{(k)}_{kk} & a^{(k)}_{k,k+1} \dots a^{(k)}_{kn} \ dots & dots \ a^{(k)}_{mk} & dots \ a^{(k)}_{m,k+1} \dots a^{(k)}_{mn} \end{pmatrix} = dots (oldsymbol{a}^k \mid oldsymbol{B}^{(k)})$$

gemäß 3.3.6 durch die Householder-Spiegelung $\overline{H}_k \in \mathbb{R}^{m-k+1,m-k+1}$ in

$$\overline{H}_{k}M^{(k)} = (\overline{a}^{k} | \overline{B}^{(k)}) = \left(\frac{\varrho_{k}}{o} \middle| \overline{B}^{(k)}\right) = : \left(\frac{r_{kk} \mid r_{k,k+1} \dots r_{kn}}{o \mid M^{(k+1)}}\right)$$
(7)

transformiert. Dies ist bei der Festlegung

$$\overline{H}_{k} = I - \frac{\overline{v}^{k} \overline{v}^{k} \tau}{\gamma_{k}} \quad \text{mit} \quad \overline{v}^{k} = \begin{pmatrix} v_{kk} \\ v_{k+1,k} \\ \vdots \\ v_{mk} \end{pmatrix} = \begin{pmatrix} 1 - a_{kk}^{(k)}/\varrho_{k} \\ -a_{k+1,k}^{(k)}/\varrho_{k} \\ \vdots \\ -a_{mk}^{(k)}/\varrho_{k} \end{pmatrix} = e^{1} - \frac{1}{\varrho_{k}} a^{k},$$
(8)

. . .

$$\varrho_{k} := \begin{cases} \|\boldsymbol{a}^{k}\| & \text{für } \boldsymbol{a}_{kk}^{(k)} \leq 0, \\ -\|\boldsymbol{a}^{k}\| & \text{für } \boldsymbol{a}_{kk}^{(k)} > 0 \end{cases} \quad \text{und} \quad \gamma_{k} := v_{kk} = \|\overline{\boldsymbol{v}}^{k}\|^{2}/2 \tag{9}$$

der Fall. Mit der erweiterten Householder-Spiegelung

$$\boldsymbol{H}_{k} := \left(\frac{\boldsymbol{I}_{k-1} \mid \boldsymbol{O}}{\boldsymbol{O} \mid \boldsymbol{\overline{H}}_{k}} \right) = \boldsymbol{I} - \frac{\boldsymbol{v}^{k} \boldsymbol{v}^{k\mathsf{T}}}{\gamma_{k}}, \quad \boldsymbol{v}^{k} := \left(\frac{\boldsymbol{O}}{\boldsymbol{\overline{v}}^{k}} \right)_{k}^{k-1} = \left(\begin{array}{c} 0\\ \vdots\\ 0\\ \vdots\\ \vdots\\ v_{mk} \end{array} \right) \in \mathbf{R}^{m}$$
(10)

läßt sich der k-te Schritt durch

$$A^{(k+1)} := H_k A^{(k)} = \begin{pmatrix} R^{(k)} \\ O & \overline{H}_k M^{(k)} \end{pmatrix} = \begin{pmatrix} R^{(k)} \\ O & \overline{Q}^{(k)} \\ O & \overline{Q}^{(k)} \end{pmatrix} = \begin{pmatrix} R^{(k+1)} \\ O & M^{(k+1)} \end{pmatrix}$$

beschreiben, und $A^{(k+1)}$ ist bis zur k-ten Spalte von oberer Dreiecksform.

Die nach n Schritten erhaltene Matrix

$$\hat{R} := A^{(n+1)} = H_n A^{(n)} = H_n H_{n-1} \cdots H_1 A^{(1)} = \hat{Q}^{\mathsf{T}} A$$

 $_{mit}$

$$\hat{Q}^{\intercal} := H_n H_{n-1} \cdots H_1, \qquad \hat{Q} = H_1 H_2 \cdots H_n$$

ist dann nach Konstruktion von oberer Dreiecksform, so daß eine QR-Faktorisierung des Typs (1), (2) vorliegt. Im Fall m = n ist $A^{(n)}$ bereits eine Dreiecksmatrix, so daß nur n - 1 Schritte erforderlich sind. Der Gesamtprozeß wird Householder-Orthogonalisierung der Matrix A genannt.

10.2.2. Householder-Orthogonalisierung. Es sei $A \in \mathbb{R}^{m,n}$ eine spaltenreguläre Matrix. Dann sind die Transformationen

$$A^{(k+1)} := H_k A^{(k)}$$
 $(k = 1, ..., l),$ $A^{(1)} := A, l := \min\{m - 1, n\},$

mit Householder-Spiegelungen H_k gemäß (8) bis (10) in exakter Arithmetik durchführbar und ergeben eine QR-Faktorisierung $A = \hat{Q}\hat{R}$ mit der orthogonalen Matrix

$$\hat{\boldsymbol{Q}} := \boldsymbol{H}_1 \boldsymbol{H}_2 \cdots \boldsymbol{H}_{l-1} \boldsymbol{H}_l$$

und der oberen Dreiecksmatrix $\hat{R} = \left(\frac{R}{O}\right) = A^{(l-1)}$, wobei $\varrho_k = r_{kk} \neq 0$ (k = 1, ..., n) ist.

Beweis. Es ist noch die Durchführbarkeit zu zeigen, die wegen $\gamma_k = (|a_{kk}^{(k)}| + ||a^k||)/||a^k||$ zur Bedingung $a^k \neq o$ (k = 1, ..., n) gleichwertig ist. Es sei also erstmals $a^j = o$ für ein gewisses $j \in \{1, ..., n\}$. Dann sind die ersten j Spalten von $A^{(j)}$ linear abhängig, und wegen $H_{i-1} \cdots H_1 A = A^{(j)}$ trifft dies auch für die ersten j Spalten von A zu. Dies widerspricht der vorausgesetzten Spaltenregularität von A.

10.2.3. Bemerkung. (i) Bei Berechnung von $\overline{B}^{(k)}$ gemäß 3.3.6 erfordert die Householder-Orthogonalisierung $\sim (m - n/3) n^2$ opms. Für $m \gg n$ ist sie daher etwa so teuer wie die MGS-Orthogonalisierung. Im Fall $m \approx n$ ist der Aufwand mit dem für das Normalgleichungsverfahren aus 9.1 vergleichbar.

(ii) Mit zusätzlichen n S zum Speichern der Diagonalelemente $\varrho_k = r_{kk}$ von **R** kann die Householder-Orthogonalisierung in situ realisiert werden, indem die signifikanten Elemente v_{ik} $(i \ge k)$ der v^k auf dem Platz der zu 0 gemachten a_{ik} und die r_{ki} (j > k) auf dem Platz der nicht mehr benötigten a_{ki} gespeichert werden. Für m = 4, n = 3 sind dann Anfangs- und Endbelegung durch

$$\begin{pmatrix} \mathbf{*} & \mathbf{*} & \mathbf{*} \\ a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{pmatrix} \rightarrow \begin{pmatrix} \underline{\varrho_1} & \underline{\varrho_2} & \underline{\varrho_3} \\ v_{11} & 1 & r_{12} & r_{13} \\ v_{21} & v_{22} & r_{23} \\ v_{31} & v_{32} & v_{33} \\ v_{41} & v_{42} & v_{43} \end{pmatrix}$$

gegeben. Im Gegensatz zur MGS-Orthogonalisierung ist also für R kein zusätzlicher Speicherplatz erforderlich.

Um die tatsächlich auszuführenden Operationen deutlich zu machen, geben wir eine in-situ-Realisierung des Householder-Verfahrens entsprechend 10.2.3(ii) in detaillierter Form an.

10.2.4. Orthogonale Faktorisierung der Matrix A nach Householder.

Aufgabe: Für spaltenreguläres $A \in \mathbb{R}^{m,n}$ ist durch Householder-Orthogonalisierung $\hat{\boldsymbol{R}} = \left(\frac{\boldsymbol{R}}{\boldsymbol{O}}\right), \ \boldsymbol{R} = (r_{ij}) \in \mathbb{R}^{n,n}, \text{ und einer in der Produktform}$ $\hat{\boldsymbol{Q}} = \boldsymbol{H}_1 \boldsymbol{H}_2 \cdots \boldsymbol{H}_l, \quad l := \min \{m - 1, n\}, \quad (11)$ gegebenen Matrix $\hat{\boldsymbol{Q}}$ zu bestimmen. Die Faktoren \boldsymbol{H}_k sind in der Form

$$\boldsymbol{H}_{\boldsymbol{k}} = \boldsymbol{I} - \boldsymbol{v}^{\boldsymbol{k}} \boldsymbol{v}^{\boldsymbol{k} \mathsf{T}} / \gamma_{\boldsymbol{k}}, \quad \gamma_{\boldsymbol{k}} = \boldsymbol{v}^{\boldsymbol{k} \mathsf{T}} \boldsymbol{v}^{\boldsymbol{k}} / 2 \qquad (k = 1, ..., l), \tag{12}$$

durch $\boldsymbol{v}^{k} = (0, ..., 0, v_{kk}, ..., v_{mk})^{\mathsf{T}} \in \mathbf{R}^{m}$ festgelegt, wobei \boldsymbol{v}^{k} gemäß $\gamma_{k} = v_{kk}$ skaliert ist. Die Matrix A ist durch die signifikanten Elemente v_{ik} $(i \geq k)$ von \boldsymbol{v}^{k} bzw. r_{kj} (k < j) von \boldsymbol{R} zu überschreiben; die Diagonalelemente $\varrho_k = r_{kk}$ von \boldsymbol{R} sind gesondert zu speichern.

19 Schwetlick, Numerische Algebra
Die Householder-Spiegelung H_k kann entsprechend 3.3.13 durch eine Folge

$$G_k := G_{km} G_{k,m-1} \cdots G_{k,k+2} G_{k,k+1}$$
(13)

von m-k Givens-Drehungen G_{ki} ersetzt werden, die nacheinander die Elemente $a_{ik}^{(k)}$ (i = k + 1, ..., m) zu 0 machen. Dabei geht $A^{(k+1)} = H_k A^{(k)}$ in

$$A^{(k+1)} = G_k A^{(k)} = G_{km} \cdots G_{k,k+1} A^{(k)} \qquad (k = 1, ..., l)$$
(14)

über, und $\hat{oldsymbol{Q}}$ ist durch

$$\hat{oldsymbol{Q}}^{\intercal} = oldsymbol{G}_l \cdots oldsymbol{G}_1, \hspace{0.1 in} ext{also} \hspace{0.1 in} oldsymbol{Q} = oldsymbol{G}_1^{\intercal} \cdots oldsymbol{G}_l^{\intercal}$$

bzw. ausführlich durch

$$\hat{\boldsymbol{Q}} = \boldsymbol{G}_{12}^{\mathsf{T}} \cdots \boldsymbol{G}_{1m}^{\mathsf{T}} \boldsymbol{G}_{23}^{\mathsf{T}} \cdots \boldsymbol{G}_{2m}^{\mathsf{T}} \cdots \boldsymbol{G}_{lm}^{\mathsf{T}}$$
(15)

gegeben. Dieser Prozeß wird Givens-Orthogonalisierung von A genannt.

10.2.5. Bemerkung. (i) Bei der expliziten Realisierung von (14), d. h. bei der Realisierung der zu (7) analogen Transformation

$$\boldsymbol{M}^{(k)} = (\boldsymbol{a}^{k} \mid \boldsymbol{B}^{(k)}) \to \bar{\boldsymbol{G}}_{k} \boldsymbol{M}^{(k)} = (\bar{\boldsymbol{a}}^{k} \mid \bar{\boldsymbol{B}}^{(k)}) = \left(\frac{\varrho_{k}}{\boldsymbol{o}} \mid \bar{\boldsymbol{B}}^{(k)}\right)$$
(16)

 mit

$$G_k = \left(egin{matrix} I_{k-1} & O \ \hline O & egin{matrix} egin{matrix} O \ \hline egin{matrix} O & egin{matrix} egin{matrix} egin{matrix} O \ \hline egin{matrix} egin{matrix} egin{matrix} O \ \hline egin{matrix} egin$$

nach 3.3.14, erfordert die Givens-Orthogonalisierung $\sim (m - n/3) n^2$ (2 opm + 1 ops) und $\sim (m - n/2) n$ opr im Vergleich zu $\sim (m - n/3) n^2$ opms und n opr für die Householder-Orthogonalisierung, ist also etwa doppelt so teuer. Für n in der Größenordnung von 10 liegt der Aufwand für die Berechnung der $\sim (m - n/2) n$ Wurzeln in der Größenordnung von $(m - n/3) n^2$ opms, kann also nicht vernachlässigt werden. In-situ-Realisierung ist möglich, indem jeder Drehung G_{ki} gemäß 3.3.11 die Zahl ξ_{ik} zugeordnet und letztere auf dem Platz des zu 0 gemachten a_{ik} (i > k) gespeichert wird; die a_{kj} $(j \ge k)$ werden mit r_{kj} überschrieben. Das zu 10.2.3(ii) analoge Schema ist dann

$$\begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \\ a_{41} & a_{42} & a_{43} \end{pmatrix} \rightarrow \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ \overline{\xi_{21}} & r_{22} & r_{23} \\ \overline{\xi_{31}} & \overline{\xi_{32}} & r_{33} \\ \overline{\xi_{41}} & \overline{\xi_{42}} & \overline{\xi_{43}} \end{pmatrix}.$$

(ii) Wenn die Matrix A in der skalierten Form

$$\boldsymbol{A} = \operatorname{diag}\left(\sqrt{\varkappa_{i}}\right)\boldsymbol{D}, \quad \mathrm{d. h.} \quad a_{ij} = \sqrt{\varkappa_{i}} d_{ij} \qquad (i = 1, ..., m; \ j = 1, ..., n)$$
(17)

durch die Zahlen $\{d_{ij}, \varkappa_i\}$ dargestellt wird und die Transformationen (14) bzw. (16) in der impliziten Form gemäß 3.3.21 realisiert werden, spricht man von *impliziter Givens-Orthogonalisierung*. Die Transformationen G_{ki} sind dann in der skalierten Form durch \mathbf{F}_{ki} und zugehörige intermediäre Skalierungsvektoren $\mathbf{z}^{(k,i)}$ und $\mathbf{\bar{z}}^{(k,i)}$ gegeben. Die Matrix \mathbf{F}_{ki} kann gemäß 3.3.18 durch die dort festgelegte Zahl ξ_{ik} repräsentiert werden. Die Skalierungsvektoren $\mathbf{z}^{(k,i)}$, $\mathbf{\bar{z}}^{(k,i)}$ sind nicht zugänglich, können jedoch in direkter Reihenfolge aus den Anfangswerten $\mathbf{z} = (\varkappa_i)$ gemäß 3.3.19, in inverser Reihenfolge aus den Endwerten $\mathbf{\bar{z}} = (\mathbf{\bar{z}}_i)$ gemäß 3.3.20 unter Verwendung der ξ_{ik} rekonstruiert werden. Die Matrix $\hat{\mathbf{Q}}$ aus (15) ist in diesem Sinne implizit durch die Zahlen $\{\xi_{ik}, \varkappa_i\}$ bzw. $\{\xi_{ik}, \mathbf{\bar{x}}_i\}$ gegeben. Die Matrix $\hat{\mathbf{R}} = A^{(l+1)}$ fällt in der skalierten Form

$$egin{aligned} \hat{m{R}} = \left(rac{m{R}}{m{O}}
ight) = ext{diag} \left(\sqrt{ar{m{z}}_i}
ight) \hat{m{Z}}, & \hat{m{Z}} = \left(rac{m{Z}}{m{O}}
ight), & ext{d. h.} & r_{ij} = \sqrt{ar{m{z}}_i} z_{ij} \ (i \leq j; \; i, j = 1, ..., n) \end{aligned}$$

an und wird durch die Zahlen $\{z_{ij}, \bar{\varkappa}_i\}$ dargestellt.

Die implizite Givens-Orthogonalisierung ist also ein Prozeß, der den Eingangsdaten $\{d_{ij}, \varkappa_i\}$ die Ausgangsdaten $\{z_{ij} \ (i \leq j), \xi_{ij} \ (i > j), \overline{\varkappa}_i\}$ zuordnet.

Der Aufwand ist $\sim (m - n/3) n^2$ (1 opm + 1 ops), entspricht also dem der Householder-Orthogonalisierung; Quadratwurzeln treten nicht auf. In-situ-Realisierung analog zu 10.2.3(ii) ist möglich. Das zugehörige Belegungsschema ist

$$\begin{pmatrix} \varkappa_1 & d_{11} & d_{12} & d_{13} \\ \varkappa_2 & d_{21} & d_{22} & d_{23} \\ \varkappa_3 & d_{31} & d_{32} & d_{33} \\ \varkappa_4 & d_{41} & d_{42} & d_{43} \end{pmatrix} \rightarrow \begin{pmatrix} \overline{\varkappa}_1 & \underline{z_{11}} & z_{12} & z_{13} \\ \overline{\varkappa}_2 & \overline{\xi_{21}} & \underline{z_{22}} & z_{23} \\ \overline{\varkappa}_3 & \overline{\xi_{31}} & \overline{\xi_{32}} & \underline{z_{33}} \\ \overline{\varkappa}_4 & \overline{\xi_{41}} & \overline{\xi_{42}} & \overline{\xi_{43}} \end{pmatrix}.$$

(iii) Zur Faktorisierung von A werden üblicherweise $\varkappa_i := 1$ und D := A als Eingangsdaten für die implizite Givens-Orthogonalisierung gewählt. Ist jedoch die Matrix eines diagonalgewichteten Quadratmittelproblems $\hat{D}\hat{A}x \cong \hat{D}\hat{b}$ mit $\hat{D} = \text{diag}(\hat{d}_i), \hat{d}_i = 1/\sqrt{w_i}$, zu faktorisieren, vgl. 8.3, so sollte $D := \hat{A}$ und

$$\kappa_i := w_{\max}/w_i$$
 $(i = 1, ..., m)$ mit $w_{\max} = \max \{w_i : i = 1, ..., m\}$

gesetzt werden, denn damit ist $A = \text{diag} \left(\sqrt{\varkappa_i} \right) D = \sqrt{w_{\text{max}}} \hat{D} \hat{A}$ bis auf den skalaren Faktor $\sqrt{w_{\text{max}}}$ mit $\hat{D} \hat{A}$ identisch.

(iv) Falls die Skalierungsfaktoren im Laufe der impliziten Givens-Orthogonalisierung zu klein werden und zu Unterlauf führen könnten, vgl. 3.3.23, ist eine *Reskalierung* erforderlich. Dabei wird das zu kleine aktuelle \varkappa_i durch $\varkappa_i * \omega^2$ mit einem Faktor $\omega \gg 1$ ersetzt, und die *i*-te Zeile von $D^{(k)}$ wird durch ω dividiert. Wenn ω eine ganzzahlige Potenz der Basis ist, entstehen dabei keine Rundungsfehler.

(v) Die Reskalierung wird überflüssig, wenn eine "Spaltenpivotisierung" – d. h. Pivotsuche in den Spalten – nach der folgenden Vorschrift durchgeführt wird, vgl. 3.3.23(iv): Zu Beginn des k-ten Schrittes werden die Zeilen k und s von $\{\mathbf{z}^k, \mathbf{D}^{(k)}\}$ vertauscht, wobei s = s(k) mit $k \leq s \leq m$ durch

$$\kappa_s^{(k)}[d_{sk}^{(k)}]^2 = \max \left\{ \kappa_i^{(k)}[d_{ik}^{(k)}]^2 \colon k \le i \le m \right\}$$
(18)

festgelegt ist und $\mathbf{x}^{k} = (\mathbf{x}_{i}^{(k)})$ die Skalierungsfaktoren von $A^{(k)} = \text{diag}\left(\sqrt{\mathbf{x}_{i}^{(k)}}\right) \mathbf{D}^{(k)}$ bezeichnet. Damit ergibt sich eine orthogonale Faktorisierung

$$\boldsymbol{A} = \boldsymbol{\hat{Q}}\boldsymbol{\hat{R}} \quad \text{mit} \quad \boldsymbol{\hat{Q}} = \boldsymbol{T}_{1,s(1)}\boldsymbol{G}_{1}^{\mathsf{T}}\cdots\boldsymbol{T}_{l,s(l)}\boldsymbol{G}_{l}^{\mathsf{T}}$$
(19)

in impliziter, skalierter Form. Da wegen (18) immer Fall 1 eintritt, vereinfacht sich auch die Berechnung und Rekonstruktion der Drehungsparameter.

(vi) Die in (v) beschriebenen Zeilenvertauschungen sind auch bei der Householderoder expliziten Givens-Orthogonalisierung möglich. Dabei ist (18) durch

$$|a_{sk}^{(k)}| = \max\{|a_{ik}^{(k)}|: k \le i \le m\}$$
(20)

zu ersetzen. Mit dieser zusätzlichen Pivotisierung sind die genannten Verfahren auch für steife diagonalgewichtete Probleme geeignet. Wenn auf die Pivotisierung verzichtet wird, sollten die Zeilen zumindest nach abnehmenden Gewichten $d_i = 1/\sqrt{w_i}$ geordnet werden.

(vii) Außer den genannten Zeilenvertauschungen sind auch Spaltenvertauschungen analog zu 10.1.6(iii) möglich. Wenn die Spalten der zu Beginn des k-ten Schrittes vorliegenden Restmatrix $M^{(k)}$ gemäß

$$M^{(k)} = [m^{k,k}, m^{k,k+1}, ..., m^{k,n}]$$

bezeichnet werden, wird $\hat{s} = \hat{s}(k)$ mit $k \leq \hat{s} \leq n$ und

$$\|\boldsymbol{m}^{\boldsymbol{k},\boldsymbol{\hat{s}}}\| = \max\left\{\|\boldsymbol{m}^{\boldsymbol{k},\boldsymbol{j}}\|: \boldsymbol{k} \leq \boldsymbol{j} \leq \boldsymbol{n}\right\}$$

$$(21)$$

gesucht und eine Vertauschung der Spalten k und $\hat{s}(k)$ von $A^{(k)}$ vorgenommen. Als Ergebnis entsteht eine **QR**-Faktorisierung

$$\bar{\boldsymbol{A}} = \boldsymbol{A} \boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} = \hat{\boldsymbol{Q}} \hat{\boldsymbol{R}} \quad \text{mit} \quad \boldsymbol{P}_{\boldsymbol{S}}^{\mathsf{T}} = \boldsymbol{T}_{1,\hat{\boldsymbol{s}}(1)} \boldsymbol{T}_{2,\hat{\boldsymbol{s}}(2)} \cdots \boldsymbol{T}_{n-1,\hat{\boldsymbol{s}}(n-1)}$$
(22)

und

$$r_{kk}^2 \ge r_{kj}^2 + r_{k+1,j}^2 + \dots + r_{jj}^2 \quad (j = k+1, \dots, n; k = 1, \dots, n-1);$$
 (23)

insbesondere gilt

$$|r_{11}| \geq |r_{22}| \geq \cdots \geq |r_{nn}|.$$

Die benötigten Normen $||\boldsymbol{m}^{k,j}||$ können analog zu 10.1.6 aus $||\boldsymbol{m}^{1,j}|| = ||\boldsymbol{a}^j||$ rekursiv berechnet werden. Wie dort bricht das Verfahren in dieser pivotisierten Version im Fall rang $(\boldsymbol{A}) = r < n$ nach r Schritten wegen $||\boldsymbol{m}^{r+1,j}|| = 0$ (j = r + 1, ..., n) ab, sofern exakt gerechnet wird. Es ist daher mit geeigneten Abbruchkriterien auch für rangdefizientes \boldsymbol{A} geeignet, siehe Kapitel 11.2. \Box

Im Unterschied zur MGS-Orthogonalisierung liefert die Householder-Orthogonalisierung den orthogonalen Faktor \hat{Q} nicht explizit, sondern in der Produktdarstellung

$$\hat{\boldsymbol{Q}}^{\mathsf{T}} = \boldsymbol{H}_l \cdots \boldsymbol{H}_2 \boldsymbol{H}_1 \quad \text{bzw.} \quad \hat{\boldsymbol{Q}} = \boldsymbol{H}_1 \boldsymbol{H}_2 \cdots \boldsymbol{H}_l,$$

wobei jeder Faktor H_k gemäß (12) durch den Vektor v^k repräsentiert ist. Das ist durchaus kein Nachteil, denn zur Berechnung von x und r nach dem Basisalgorithmus 10.2.1 wie auch für weitere Anwendungen genügt es, wenn $\hat{Q}^{\mathsf{T}}b$ und $\hat{Q}f$ für gegebenes b und f einfach berechnet werden kann. Dies ist mittels der Rekursion

$$b^{1} := b, \qquad b^{k+1} := H_{k}b^{k} \qquad (k = 1, ..., l), \qquad \hat{Q}^{\mathsf{T}}b := b^{l+1},$$

$$f^{l+1} := f, \qquad f^{k} := H_{k}f^{k+1} \qquad (k = l, ..., 1), \qquad \hat{Q}f := f^{1}$$
(24)

möglich. Die zugehörigen Computerrealisierungen lauten wie folgt:

10.2.6. Berechnung von $\hat{Q}^{\dagger}b$ und $\hat{Q}f$ mittels der Produktdarstellung von \hat{Q} .

Aufgabe: Die Matrix $\hat{\boldsymbol{Q}}$ sei durch die Ausgangsdaten $\boldsymbol{v}^{k} \in \Re^{m}$ (k = 1, ..., l) von 10.2.4 gemäß (11), (12) definiert, die signifikanten Elemente v_{ik} $(i \ge k)$ seien auf dem Platz von \boldsymbol{A} gespeichert. Für $\boldsymbol{b}, \boldsymbol{f} \in \Re^{m}$ ist $\hat{\boldsymbol{c}} = \hat{\boldsymbol{Q}}^{\mathsf{T}}\boldsymbol{b}$ bzw. $\hat{\boldsymbol{g}} = \hat{\boldsymbol{Q}}\boldsymbol{f}$ zu berechnen und auf dem Platz von \boldsymbol{b} bzw. \boldsymbol{f} zu speichern.

Algorithmus:

```
\begin{split} l &:= \min \ \{m - 1, n\} \\ A1 \ (b &:= \hat{Q}^{\mathsf{T}}b): \\ \text{for } k &:= 1(1)l \text{ do} \\ & \left| \begin{array}{l} \beta &:= 0 \\ \text{for } i &:= k(1)m \text{ do } \beta &:= \beta + a_{ik} * b_i \\ \beta &:= \beta/a_{kk} \\ \text{for } i &:= k(1)m \text{ do } b_i &:= b_i - a_{ik} * \beta \\ A2 \ (f &:= \hat{Q}f): \\ \text{for } k &:= l(-1)1 \text{ do} \\ & \left| \begin{array}{l} \beta &:= 0 \\ \text{for } i &:= k(1)m \text{ do } \beta &:= \beta + a_{ik} * f_i \\ \beta &:= \beta/a_{kk} \\ \text{for } i &:= k(1)m \text{ do } f_i &:= f_i - a_{ik} * \beta \\ Autwand: \text{ jeweils } \sim (2m - n) n \text{ opms.} \end{split} \end{split}
```

Man beachte, daß die explizite Bildung von \hat{Q} auch aus Aufwandsgründen nicht zweckmäßig ist: Außer den zur Berechnung erforderlichen Rechenoperationen – zu diesen siehe Ü 10.2.2 – wären dafür m^2 S zum Speichern von \hat{Q} und jeweils $\sim m^2$ opms für die Berechnung von $\hat{Q}^{\mathsf{T}}b$ bzw. $\hat{Q}f$ erforderlich. Die implizite Produktdarstellung mit Speicherung der erzeugenden Daten v_{ik} auf dem Platz von A ist deshalb als natürliche Darstellung von \hat{Q} anzusehen. Dasselbe trifft sinngemäß für die Givens-Faktorisierung zu.

10.2.7. Rundungsfehleranalyse. Die Householder-Orthogonalisierung 10.2.4 ist für spaltenreguläres $A \in \Re^{m,n}$ mit $\varrho_k = r_{kk} \neq 0$ (k = 1, ..., n) durchführbar, sofern

$$\kappa := \nu F \text{ cond } (A) < 1 \quad \text{mit} \quad F := Kmn^{3/2}$$
(25)

gilt. Dabei ist K := K(m) = 3.14(1 + 3/m) wie in 3.3.7, insbesondere gilt

$$K(m) \leq 4$$
 für $m \geq 10$. (26)

Durch 10.2.4 ist eine exakt orthogonale Matrix $\hat{Q} \in \mathbb{R}^{m,m}$ definiert, so daß mit dem berechneten Dreiecksfaktor $\hat{R} \in \Re^{m,n}$

$$\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A} = \hat{\boldsymbol{Q}} \hat{\boldsymbol{R}} \quad \text{mit} \quad \|\boldsymbol{\delta} \boldsymbol{A}\| \leq v F \|\boldsymbol{A}\| \tag{27}$$

gilt. Die Ausgangsdaten $v^k \in \Re^m$ (k = 1, ..., l) von 10.2.4 liefern eine numerische Darstellung von \hat{Q} derart, daß die durch 10.2.6 für jedes $b, f \in \Re^m$ berechenbaren Vektoren $\hat{c} = \mathrm{fl}(\hat{Q}^{\mathsf{T}}b)$ und $\hat{g} = \mathrm{fl}(\hat{Q}f)$ den Beziehungen

$$fl(\hat{\boldsymbol{Q}}^{\mathsf{T}}\boldsymbol{b}) = \hat{\boldsymbol{Q}}^{\mathsf{T}}(\boldsymbol{b} + \boldsymbol{\delta}\boldsymbol{b}) \quad \text{mit} \quad \|\boldsymbol{\delta}\boldsymbol{b}\| \leq v(F/\sqrt{n}) \|\boldsymbol{b}\|,$$

$$fl(\hat{\boldsymbol{Q}}\boldsymbol{f}) = \hat{\boldsymbol{Q}}(\boldsymbol{f} + \boldsymbol{\delta}\boldsymbol{f}) \quad \text{mit} \quad \|\boldsymbol{\delta}\boldsymbol{f}\| \leq v(F/\sqrt{n}) \|\boldsymbol{f}\|$$
(28)

genügen.

Dabei ist die Störungsmatrix δA spaltenweise klein im Sinne von

$$\|\boldsymbol{\delta}\boldsymbol{a}^{\boldsymbol{j}}\| \leq \nu K m \boldsymbol{j} \|\boldsymbol{a}^{\boldsymbol{j}}\| \qquad (\boldsymbol{j}=1,...,n).$$
⁽²⁹⁾

Beweis. Wir setzen zusätzlich $\mu := \nu 3.14(m + 3)^2 \leq 1$ voraus, was praktisch keine Einschränkung bedeutet. Der Beweis ist auch ohne diese Voraussetzung möglich, verläuft dann aber wesentlich komplizierter.

Wir nehmen an, daß bereits k-1 Schritte des Householder-Verfahrens durchgeführt worden sind und daß

$$\boldsymbol{A}^{(k)} = \boldsymbol{\hat{Q}}_{k-1} (\boldsymbol{A} + \delta \boldsymbol{\bar{A}}^{(k)}) \tag{30}$$

$$\|\delta \bar{A}^{(k)}\|_{F} \leq \nu K m(k-1) \, \|A\|_{F} \tag{31}$$

genügt. Dann folgt

$$\|A^{(k)}\|_{F} \leq \|A\|_{F} + \|\delta\bar{A}^{(k)}\|_{F} \leq [1 + \nu Km(k-1)] \|A\|_{F}.$$
(32)

Aus (31) ergibt sich ferner

$$\|A^+\| \, \|\delta ar{A}^{(k)}\| \leq \|A^+\| \, \|\delta ar{A}^{\prime k)}\|_F \leq
u Km(k-1) \, \sqrt{n} ext{ cond } (A) < arkappa < 1$$
 .

Die Matrix $A + \delta \bar{A}^{(k)}$ ist dann nach 8.2.5 spaltenregulär, und dasselbe trifft wegen (30) auch für

$$A^{(k)} = \begin{pmatrix} \mathbf{R}^{(k)} \\ \mathbf{O} \mid \mathbf{M}^{(k)} \end{pmatrix}$$
 mit $\mathbf{M}^{(k)} = (\mathbf{a}^k \mid \mathbf{B}^{(k)})$

zu. Insbesondere ist $a^k \neq o$, so daß der k-te Schritt mit $\varrho_k = r_{kk} \neq 0$ ausgeführt werden kann. Es sei jetzt H_k^* diejenige Spiegelungsmatrix, die der Spalte a^k nach den Formeln (8) bis (10) in exakter Arithmetik zugeordnet wird. Nach 3.3.7 genügt dann die berechnete Matrix $A^{(k+1)}$ der Beziehung

$$A^{(k+1)} = H_k^* (A^{(k)} + \delta A^{(k)})$$
(33)

 $_{mit}$

$$\|\delta A^{(k)}\|_F \leq
u \ 3.14(m-k+4) \ \|A^{(k)}\|_F \leq
u Km\lambda \ \|A\|_F;$$

die letzte Ungleichung folgt aus (32) mit

$$\lambda := rac{m-k+4}{m+3} \left[1 +
u Km(k-1)
ight] = 1 - rac{k-1}{m+3} \left(1 - \mu \; rac{m-k+4}{m+3}
ight) \leq 1.$$

Es gilt also

$$\|\delta A^{(k)}\|_F \leq \nu K m \|A\|_F.$$
(34)

Einsetzen von (30) in (33) liefert

$$\mathbf{1}^{(k+1)} = H_k^*[\hat{Q}_{k-1}(A + \delta \bar{A}^{(k)}) + \delta A^{(k)}] = \hat{Q}_k(A + \delta \bar{A}^{(k+1)})$$

 $_{\rm mit}$

$$\hat{oldsymbol{D}}_k = oldsymbol{H}_k^{oldsymbol{\star}} \hat{oldsymbol{Q}}_{k-1}, \qquad \delta oldsymbol{ar{A}}^{(k+1)} = \delta oldsymbol{ar{A}}^{(k)} + \hat{oldsymbol{Q}}_{k-1}^{\mathsf{T}} \delta oldsymbol{A}^{(k)},$$

d. h., (30) gilt auch für den Index k + 1. Aus (31) und (34) folgt außerdem

 $\|\delta \bar{A}^{(k+1)}\|_F \leq \|\delta \bar{A}^{(k)}\|_F + \|\delta A^{(k)}\|_F \leq \nu Kmk \|A\|_F,$

also (31) für k + 1 statt k. Damit sind (30), (31) für jedes k = 1, ..., l + 1 gültig, insbesondere ist

$$\hat{\boldsymbol{R}} = A^{(l+1)} = \hat{\boldsymbol{Q}}_l(A + \delta \bar{A}^{(l+1)}) =: \hat{\boldsymbol{Q}}^{\mathsf{T}}(A + \delta A)$$

 \mathbf{mit}

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq \|\boldsymbol{\delta}\boldsymbol{A}\|_{F} = \|\boldsymbol{\delta}\bar{\boldsymbol{A}}^{(l+1)}\|_{F} \leq \nu Kml \, \|\boldsymbol{A}\|_{F} \leq \nu F \, \|\boldsymbol{A}\|.$$

Der Beweis von (28) verläuft analog.

Der Leser lasse sich nicht dadurch irritieren, daß die in 10.2.7 erwähnte und im Beweis definierte exakt orthogonale Matrix $\hat{Q} = H_1^* \cdots H_l^*$ nicht zugänglich ist. Wesentlich ist:

- Es existiert eine exakt orthogonale Matrix $\hat{\mathbf{Q}}$, die mit dem berechneten Dreiecksfaktor $\hat{\mathbf{R}}$ die exakte \mathbf{QR} -Faktorisierung der wenig gestörten Matrix $\mathbf{A} + \mathbf{\delta A}$ liefert.
- Die nicht zugängliche exakt orthogonale Matrix \hat{Q} wird numerisch durch die berechneten Vektoren v^k (k = 1, ..., l) im folgenden Sinne gut repräsentiert: Die mittels der v^k nach 10.2.6 berechneten Vektoren $fl(\hat{Q}^{\mathsf{T}}b)$ bzw. $fl(\hat{Q}f)$ sind das exakte Produkt der Matrizen \hat{Q}^{T} bzw. \hat{Q} mit den wenig gestörten Argumenten $b + \delta b$ bzw. $f + \delta f$.

Die Householder-Orthogonalisierung von A und die Berechnung von $\hat{Q}^{\mathsf{T}}b$ bzw. $\hat{Q}f$ mittels der $\mathbf{v}^{\mathsf{*}}$ sind also numerisch gutartige Prozesse. 10.2.8. Bemerkung. (i) Für die in 10.2.5 beschriebene explizite oder implizite Givens-Orthogonalisierung und die zu 10.2.6 analoge Berechnung von $\hat{Q}^{\dagger}b$ bzw. $\hat{Q}f$ bleibt 10.2.7 gültig, wobei F durch

$$F := 6(m+n) n^{1/2}$$
(35)

zu ersetzen ist. Zum Nachweis von (35) ist eine sorgfältige Abschätzungstechnik nötig; eine Übertragung des Beweises von 10.2.7 unter Verwendung von 3.3.15 bzw. 3.3.22 würde auf $F = 6 mn^{3/2}$ führen; siehe B 10.3 für Literaturhinweise. Die explizite oder implizite Givens-Orthogonalisierung ist also ein numerisch gutartiger Prozeß mit der wesentlich günstigeren Fehlerkumulationskonstanten (35) als bei der Householder-Orthogonalisierung. Trotzdem wird die Householder-Faktorisierung derzeit der nach GIVENS meist vorgezogen, da sie für vollbesetztes A einerseits billiger als die explizite Givens-Faktorisierung, andererseits einfacher zu verstehen und zu realisieren ist als die implizite Givens-Faktorisierung. Unter Beachtung aller Gesichtspunkte wie Aufwand, Rundungsfehlerverhalten und Flexibilität in der Anwendung besitzt die implizite Givens-Orthogonalisierung jedoch deutliche Vorteile gegenüber den übrigen Orthogonalisierungsverfahren. Zugehörige hochwertige Software liegt allerdings in allgemein zugänglicher Form noch nicht vor.

(ii) Bei Bedarf kann \hat{Q} spaltenweise gemäß $\tilde{q}^{j} := \operatorname{fl}(\hat{Q}e^{j})$ (j = 1, ..., m) nach 10.2.6 berechnet werden. Wegen der speziellen Form von $f = e^{j}$ vereinfacht sich dabei die Berechnung etwas, siehe U 10.2.2 für Hinweise zur Realisierung und Aufwandsangaben. Für die derart berechnete Matrix $\tilde{Q} = (\tilde{q}^{1}, ..., \tilde{q}^{m})$ gilt

$$\|\tilde{\boldsymbol{Q}}^{\mathsf{T}}\tilde{\boldsymbol{Q}} - \boldsymbol{I}\|_{F} \leq 2\nu K(m) \ m^{3/2}n, \tag{36}$$

d. h., sie ist unabhängig von cond (A) ausreichend orthogonal.

Wir kehren jetzt wieder zur Berechnung der Quadratmittellösung x und des Residuums r nach dem Basisalgorithmus 10.2.1 zurück. Eine Computerrealisierung von S1 ist bereits durch 10.2.4 gegeben. Im folgenden wird eine analoge Realisierung von S2 unter Verwendung der Auswertungsvorschriften 10.2.6 für Ausdrücke der Form $\hat{Q}^{\dagger}b$ bzw. $\hat{Q}f$ ohne explizite Berechnung von \hat{Q} angegeben.

10.2.9. Lösung des Quadratmittelproblems $\hat{Q}\hat{R}x \simeq b$.

Aufgabe: Unter Verwendung der gemäß 10.2.4 berechneten und auf dem Platz von A gespeicherten Householder-Faktorisierung $A = \hat{Q}\hat{R}$ von $A \in \Re^{m,n}$ mit $r_{kk} = \varrho_k \neq 0$ (k = 1, ..., n) ist die zur rechten Seite $b \in \Re^m$ gehörende Lösung $x \in \Re^n$ von $Ax \simeq b$ und das zugehörige Residuum $r = b - Ax \in \Re^m$ sowie dessen Norm $nr = ||r|| \in \Re$ zu berechnen. Dabei ist b mit r zu überspeichern. Algorithmus:

S2.1.1: Berechne $\boldsymbol{b} := \operatorname{fl}(\hat{\boldsymbol{Q}}^{\intercal}\boldsymbol{b})$ nach Algorithmus A1 aus 10.2.6

S2.1.2: Partitioniere **b** gemäß

$$\boldsymbol{b} = \begin{pmatrix} \boldsymbol{c} \\ - \\ \tilde{\boldsymbol{c}} \end{pmatrix}_{\substack{n+1, \\ \vdots \\ m}}^{1}, \quad \text{setze} \quad \boldsymbol{x} := \boldsymbol{c}, \quad \boldsymbol{b} := \begin{pmatrix} \boldsymbol{o} \\ \tilde{\boldsymbol{c}} \end{pmatrix}$$

Berechne $n\boldsymbol{r} := \|\tilde{\boldsymbol{c}}\|$

S2.1.3: Berechne $\boldsymbol{b} := \operatorname{fl}(\hat{\boldsymbol{Q}}\boldsymbol{b})$ nach Algorithmus A2 aus 10.2.6.

S2.2: Berechne $\boldsymbol{x} := \operatorname{fl}(\boldsymbol{R}^{-1}\boldsymbol{x})$ analog zu S2.2 von 10.1.5

Aufwand: $\sim (2mn - n^2/2)$ opms für x und nr, $\sim (2mn - n^2)$ opms für r.

10.2.10. Rundungsfehleranalyse. Die spaltenreguläre Matrix $A \in \Re^{m,n}$ genüge der Bedingung

$$\varkappa = \nu F \operatorname{cond} (A) < 1 \quad \operatorname{mit} \quad F := K m n^{3/2}, \tag{37}$$

K := K(m) wie in 10.2.7. Dann kann die Householder-Faktorisierung $A = \hat{Q}\hat{R}$ gemäß 10.2.4 mit $r_{kk} \neq 0$ (k = 1, ..., n) berechnet werden, und Algorithmus 10.2.9 ist mit den so bestimmten Faktoren für jedes $\boldsymbol{b} \in \Re^m$ durchführbar. Die berechnete Lösung \boldsymbol{x} ist exakte Lösung des spaltenregulären Quadratmittelproblems

$$(A + \boldsymbol{\delta}_1 A) \boldsymbol{x} \simeq \boldsymbol{b} + \boldsymbol{\delta} \boldsymbol{b} \quad \text{mit} \quad \|\boldsymbol{\delta}_1 A\| \leq \nu F_1 \|A\|, \qquad F_1 := F + n^{3/2} \sim F,$$
(38)

und für das berechnete Residuum r gilt

$$\boldsymbol{r} + \boldsymbol{\delta}\boldsymbol{r} = (\boldsymbol{b} + \boldsymbol{\delta}\boldsymbol{b}) - (\boldsymbol{A} + \boldsymbol{\delta}_1\boldsymbol{A})\boldsymbol{x}$$
(39)

mit

$$\|\boldsymbol{\delta b}\| \leq \boldsymbol{v}(F/\sqrt{n}) \|\boldsymbol{b}\|, \qquad \|\boldsymbol{\delta r}\| \leq \boldsymbol{v}(F/\sqrt{n}) \|\boldsymbol{r}\|.$$
(40)

Beweis. Aus 10.2.7 folgt die Durchführbarkeit der Householder-Orthogonalisierung und die Darstellung $A + \delta A = \hat{Q}\hat{R}$ sowie

$$\mathrm{fl}(\hat{Q}^{\mathsf{T}}b) = \hat{c} = \left(\frac{c}{\tilde{c}}\right) = \hat{Q}^{\mathsf{T}}(b + \delta b), \quad r = \mathrm{fl}\left(\hat{Q}\left(\frac{o}{\tilde{c}}\right)\right) = \hat{Q}\left\{\left(\frac{o}{\tilde{c}}\right) + \delta c\right\}$$

mit geeigneten Störungen δA , δb , δc . Nach 4.3.2 löst das berechnete x das gestörte Dreieckssystem ($R + \delta R$) x = c, wobei

$$|\delta oldsymbol{R}|| \leq ||\delta oldsymbol{R}||_F \leq
u n^{3/2} \, ||oldsymbol{R}||_F$$

Wenn

$$\delta_1 A := \delta A + \hat{Q}\left(rac{\delta R}{O}
ight), \;\; \delta r := - \hat{Q} \delta C$$

gesetzt wird, ergibt sich daher

$$m{A}+\delta_1m{A}=\hat{m{Q}}\left\{\hat{m{R}}+\left(rac{\deltam{R}}{m{O}}
ight)
ight\}=\hat{m{Q}}\left(rac{m{R}+\deltam{R}}{m{O}}
ight), \ \ m{r}+\deltam{r}=\hat{m{Q}}\left(rac{m{o}}{m{ ilde c}}
ight),$$

also

$$(\boldsymbol{A} + \boldsymbol{\delta}_{1}\boldsymbol{A})^{\mathsf{T}}(\boldsymbol{r} + \boldsymbol{\delta}\boldsymbol{r}) = \left(\frac{\boldsymbol{R} + \boldsymbol{\delta}\boldsymbol{R}}{\boldsymbol{O}}\right)^{\mathsf{T}}\left(\frac{\boldsymbol{o}}{\tilde{\boldsymbol{c}}}\right) = \boldsymbol{o}.$$
(41)

Überdies ist $A + \delta_1 A$ wie $R + \delta R$ spaltenregulär, und es gilt

$$r + \delta r = \hat{Q}\hat{c} - \hat{Q}\left(rac{c}{o}
ight) = b + \delta b - \hat{Q}\left(rac{R+\delta R}{O}
ight)x = (b + \delta b) - (A + \delta_1 A)x,$$

also (39), und (41) stellt gerade die Normalgleichungen des gestörten Problems (38) dar. Die Abschätzungen (40) ergeben sich wegen $||\delta r|| = ||\delta c||$ direkt aus (28). Die Schranke für $\delta_1 A$ erhält man aus (27) unter Beachtung von ||A|| = ||R|| gemäß

$$\|\delta_1 A\| \le \|\delta A\| + \|\delta R\| \le vF \|A\| + vn^{3/2} \|R\| = v(F + n^{3/2}) \|A\|.$$

10.2.11. Bemerkung. (i) Die Rundungsfehleranalyse 10.2.10 gilt mit günstigeren Kumulationskonstanten

$$F := 6(m+n) n^{1/2}, \qquad F_1 := F + n^{3/2} = (6m+7n) n^{1/2}$$
 (42)

auch für den Fall, daß die Faktorisierung $A = \hat{Q}\hat{R}$ mittels expliziter oder impliziter Givens-Drehungen berechnet und $\hat{Q}^{\mathsf{T}}b$ bzw. $\hat{Q}f$ in zu 10.2.6 analoger Weise bestimmt wird, vgl. 10.2.8(i). Die Berechnung von x und r gemä β 10.2.9 ist daher sowohl bei Verwendung der Householder- als auch der expliziten oder impliziten Givens-Faktorisierung ein numerisch gutartiger Proze β . Wegen der deutlich kleineren Kumulationskonstanten ist die Givens-Faktorisierung dabei favorisiert.

(ii) Die für die Durchführbarkeit der Householder- bzw. Givens-Faktorisierung mit $r_{kk} \neq 0$ hinreichende Bedingung (25), (37) ist mit der analogen Bedingung (10.1.17) für das MGS-Verfahren vergleichbar und schränkt die Klasse der zum Fehlerniveau v vernünftig gestellten spaltenregulären Quadratmittelprobleme nicht wesentlich ein. Auch im Fall $r_{kk} = 0$ kann die Durchführbarkeit mit nichtverschwindenden Diagonalelementen für jedes $A \neq O$ erzwungen werden, indem $r_{kk} := v ||A||$ gesetzt wird, vgl. 5.3.3(i) für die analoge Situation bei der *LR*-Faktorisierung. Allerdings ist die so modifizierte *QR*-Faktorisierung wie dort nur für spezielle Zwecke verwendbar und insbesondere nicht direkt zur Lösung rangdefizienter Quadratmittelprobleme geeignet.

(iii) Im Fall m = n liefern alle bisher behandelten Orthogonalisierungsverfahren für ausreichend reguläres A eine Faktorisierung A = QR mit orthogonalem Q und einer regulären Dreiecksmatrix R. Für die berechneten bzw. in Produktform implizit definierten Faktoren gilt dabei

$$A + \delta A = QR \quad \text{mit} \quad \|\delta A\| \leq \nu F \|A\|, \tag{43}$$

und die gemäß 10.1.5 bzw. 10.2.9 berechnete Lösung \boldsymbol{x} löst das gestörte reguläre System

$$(\boldsymbol{A} + \boldsymbol{\delta}_1 \boldsymbol{A}) \boldsymbol{x} = \boldsymbol{b} \quad \text{mit} \quad \|\boldsymbol{\delta}_1 \boldsymbol{A}\| \leq v F_1 \|\boldsymbol{A}\|.$$
(44)

Die Kumulationskonstanten F, F_1 und die Anzahl der zur Berechnung der QR-Faktorisierung erforderlichen Rechenoperationen sind in der nachfolgenden Tabelle zusammengestellt. Für die Householder-Faktorisierung wurde die für $m \ge 10$ gültige Schranke $K(m) \le 4$ verwendet. Zum Vergleich sind außerdem die entsprechenden Werte für die Gauß-Faktorisierung $P_Z A = LR$ mit Spaltenpivotisierung aus 5.3 mit aufgenommen worden.

MGS $1.4n^{3/2}$ $1.4n^{3/2}$ n^3 opms Householder $4n^{5/2}$ $4n^{5/2}$ $\sim n^3$ opms Givens explizit $12n^{3/2}$ $13n^{3/2}$ $\sim 2n^3/3$ opms Givens implicit $12n^{3/2}$ $13n^{3/2}$ $\sim 2n^3/3(2 \text{ opm} + 1 \text{ ops}) + n^2/2$	Verfahren	F_1	Aufwand für Faktorisierung
Gauß $ \begin{array}{c} 12n^{3/2} \\ C(n) n^{1/2} \\ C(n) n^{3/2} \\ \end{array} \begin{array}{c} 2n^{3/3} \\ \sim n^{3/3} \\ \sim n^{3/3} \\ 0 \\ \text{ms} \end{array} \rangle$	MGS Householder Givens explizit Givens implizit Gauß	$\begin{array}{c c c c c c c c c c c c c c c c c c c $	$\sim n^3$ opms $\sim 2n^3/3$ opms $\sim 2n^3/3(2 \text{ opm} + 1 \text{ ops}) + n^2/2 \text{ op}$ $\sim 2n^3/3 \text{ opms}$ $\sim n^3/3 \text{ opms}$

Bei der Gauß-Faktorisierung gilt $C(n) \leq 2 \cdot 2^n$ für die volle Klasse der genügend regulären Systeme. Die praktisch vorkommenden Gleichungssysteme gehören zur Unterklasse, für die $C(n) \approx n$ gesetzt werden kann, vgl. 5.3. Der Faktor $n^{1/2}$ wird durch den Übergang zur Spektralnorm hervorgerufen.

Wir wenden uns abschließend der iterativen Verbesserung von Quadratmittellösungen zu und beachten dazu, daß die exakte Lösung $x^* = A^+ b$ und das zugehörige exakte Residuum $r^* = b - Ax^*$ die eindeutige Lösung des regulären Gleichungssystems

$$\left(\frac{I \mid A}{A^{\mathsf{T}} \mid O}\right) \left(\frac{r}{x}\right) = \left(\frac{b}{o}\right) \tag{45}$$

sind, siehe Ü 10.2.4. Wenn $\{r, x\}$ Näherungen für $\{r^*, x^*\}$ sind, stellt

$$\left(\frac{e}{f}\right) := \begin{pmatrix} b - Ax - r \\ -A^{\mathsf{T}}r \end{pmatrix} \tag{46}$$

das zugehörige Residuum bezüglich (45) dar. Analog zu 5.4 kann dann der Ansatz

$$\left(\frac{\boldsymbol{r}^{*}}{\boldsymbol{x}^{*}}\right) = \left(\frac{\boldsymbol{r}}{\boldsymbol{x}}\right) + \left(\frac{\boldsymbol{s}}{\boldsymbol{h}}\right) \tag{47}$$

gemacht werden. Die Korrektur $\left(\frac{s}{h}\right)$ genügt in exakter Arithmetik dem System

$$\left(\frac{I \mid A}{A^{\mathsf{T}} \mid O}\right) \left(\frac{\mathbf{s}}{\mathbf{h}}\right) = \left(\frac{e}{f}\right),\tag{48}$$

das unter Verwendung der Faktorisierung $A = \hat{Q}\hat{R}$, $\hat{R} = \left(\frac{R}{O}\right)$, wie folgt gelöst werden kann:

S1: Berechne
$$\hat{Q}^{\mathsf{T}} e =: \left(\begin{matrix} u \\ - \\ v \end{matrix} \right)_{\substack{n+1 \\ \vdots \\ m}}^{1}$$

- S2: Bestimme \boldsymbol{w} aus $\boldsymbol{R}^{\mathsf{T}}\boldsymbol{w} = \boldsymbol{f}$
- S3: Bestimme **h** aus $\mathbf{Rh} = \mathbf{u} \mathbf{w}$

S4: Berechne
$$s := \hat{Q}\left(\frac{w}{v}\right).$$

Dann stellt

$$\left(\frac{\overline{r}}{\overline{x}}\right) := \left(\frac{r}{x}\right) + \left(\frac{s}{h}\right) \tag{50}$$

bei Computerrechnung eine i. allg. verbesserte Näherung für $\left(\frac{r^{\star}}{x^{*}}\right)$ dar.

10.2.12. Bemerkung. (i) Der Übergang von $\{x, r\}$ zu $\{\bar{x}, \bar{r}\}$ gemäß (50) kann als einzelner Schritt eines zu 5.4.3 analogen Verfahrens zur iterativen Verbesserung

299

(49)

von Quadratmittellösungen angesehen werden. Man beachte dabei, daß r und \overline{r} hier Näherungen für $r^* = b - Ax^*$ bezeichnen und nicht die zu x und \overline{x} gehörenden Residuen.

(ii) Wenn eine Faktorisierung A = QR des Typs (10.1.1), (10.1.2) bekannt ist, läßt sich eine zu (49) analoge Vorschrift zur Lösung von (48) angeben, siehe Ü 10.2.4.

(iii) Wegen der Gutartigkeit der Orthogonalisierungsverfahren aus 10.1 und 10.2 und der günstigen Werte der Fehlerkumulationskonstanten ist eine iterative Verbesserung i. allg. zwecklos, sofern die Residuen $\{e, f\}$ nur in einfacher Genauigkeit vberechnet werden. Wenn die Eingangsdaten $\{A, b\}$ jedoch als exakt betrachtet werden können und die Residuen in höherer Genauigkeit $v_1, v_1 \leq v^2$, berechnet werden, kann durch die iterative Verbesserung die volle erreichbare Genauigkeit im Sinne von

$$\begin{pmatrix} \|\boldsymbol{\delta r}\| \\ \|\boldsymbol{\delta x}\| \end{pmatrix} \leq \left[\boldsymbol{\nu} + \boldsymbol{v_1} F \text{ cond } (A) \right] \begin{pmatrix} \|\boldsymbol{r^*}\| + \|\boldsymbol{x^*}\| / \|A^+\| \\ \|\boldsymbol{x^*}\| + \|A^+\| \|\boldsymbol{r^*}\| \end{pmatrix}$$

erhalten werden. Da diese Situation in den Anwendungen selten vorkommt, gehen wir auf die iterative Verbesserung nicht im Detail ein. Man beachte auch, daß bei den Normalgleichungsverfahren durch die iterative Verbesserung die Gutartigkeit überhaupt erst erreicht wird, also dort eine andere Situation vorliegt als bei den a-priorigutartigen Orthogonalisierungsverfahren. \Box

Übungsaufgaben

Ü 10.2.1. Man schreibe das Residuum r = b - Ax für beliebiges $x \in \mathbb{R}^n$ unter Verwendung der Faktorisierung $A = \hat{Q}\hat{R}$ in der Form

$$\boldsymbol{r} = \hat{\boldsymbol{Q}}(\hat{\boldsymbol{c}} - \hat{\boldsymbol{R}}\boldsymbol{x}), \qquad \hat{\boldsymbol{c}} := \hat{\boldsymbol{Q}}^{\mathsf{T}}\boldsymbol{b} = \left(\frac{\boldsymbol{c}}{\tilde{\boldsymbol{c}}}\right), \qquad \hat{\boldsymbol{R}} = \left(\frac{\boldsymbol{R}}{\boldsymbol{O}}\right)$$
(51)

und folgere hieraus

$$\|\mathbf{r}\|^{2} = \|\hat{\mathbf{c}} - \hat{\mathbf{R}}\mathbf{x}\|^{2} = \|\mathbf{c} - \mathbf{R}\mathbf{x}\|^{2} + \|\tilde{\mathbf{c}}\|^{2}.$$
(52)

Aus (52) läßt sich der Basisalgorithmus 10.2.1 direkt ablesen.

Ü 10.2.2. Es sei $\mathbf{\hat{Q}} = \mathbf{H}_1 \cdots \mathbf{H}_l$ mit $l = \min \{m - 1, n\}$ gemäß 10.2.4 in Produktform durch die erzeugenden Vektoren v^k (k = 1, ..., l) gegeben.

(i) Man zeige die Gültigkeit von

$$\dot{\mathbf{Q}} \boldsymbol{e}^{j} = \boldsymbol{H}_{1} \cdots \boldsymbol{H}_{l} \boldsymbol{e}^{j} = \boldsymbol{H}_{1} \cdots \boldsymbol{H}_{\tilde{l}(j)} \boldsymbol{e}^{j} \qquad (j = 1, ..., m)$$

$$\tilde{\boldsymbol{j}}_{(j)} = (j \quad \text{für} \quad j = 1, ..., l - 1,$$
(53)

 \mathbf{mit}

$$\tilde{l}(j) := \begin{cases} j & \text{für } j = 1, ..., l-1, \\ l & \text{für } j \ge l. \end{cases}$$

(ii) Man nutze (53) zur spaltenweisen Berechnung von $\hat{Q} = (q^1, ..., q^m)$ gemäß $q^j = \hat{Q}e^j$ in zu 10.2.6 analoger Weise aus und überlege sich, daß dazu $\sim 2(m^2 - mn + n^2/3) n$ opms erforderlich sind. Wenn nur die ersten n Spalten $\{q^1, ..., q^n\}$ benötigt werden, reduziert sich der Aufwand auf $\sim (m - n/3) n^2$ opms.

Ü 10.2.3. Man überlege sich, daß die Normallösung des zeilenregulären Systems

$$A^{\mathsf{T}}\boldsymbol{y} = \boldsymbol{f} \tag{54}$$

unter Verwendung der Faktorisierung $A = \hat{Q}\hat{R}$ mit in Produktform gegebenem \hat{Q} wie folgt berechnet werden kann, vgl. auch 4.2.2 (i) und Ü 10.1.2:

S3.1: Bestimme \boldsymbol{g} als Lösung von $\boldsymbol{R}^{\mathsf{T}}\boldsymbol{g} = \boldsymbol{f}$, setze $\boldsymbol{\hat{g}} := \left(\frac{\boldsymbol{g}}{\boldsymbol{o}}\right)$. S3.2: Berechne $\boldsymbol{y} := \operatorname{fl}(\boldsymbol{\hat{Q}}\boldsymbol{\hat{g}})$ nach 10.2.6.

 $Ü 10.2.4. \ \mathrm{Für} \ A \in \mathbf{R}^{m,n} \ \mathrm{mit} \ m \geq n \ \mathrm{sei} \ M(A) := \left(\frac{I \ | A}{A^{\intercal} \ | \ O} \right).$

(i) Man zeige, daß M(A) genau dann regulär ist, wenn rang (A) = n gilt und daß im Fall der Spaltenregularität von A

$$M(A)^{-1} = \left(rac{I - AA^+ \mid A^{+\intercal}}{A^+ \mid -A^+A^{+\intercal}}
ight)$$

gilt.

(ii) Man überprüfe, daß die Lösung von (48) durch die Vorschriften (49) gegeben ist.

(iii) Man zeige, daß das System (48) unter Verwendung einer Faktorisierung A = QR des Typs (10.1.1), (10.1.2) wie folgt gelöst werden kann:

S1: Berechne $\boldsymbol{u} := \boldsymbol{Q}^{\mathsf{T}} \boldsymbol{e}$

.

S2: Bestimme w aus $R^{\mathsf{T}}w = f$

S3: Bestimme h aus Rh = u - w

S4: Berechne s := e - Q(u - w).

10.3. Aufdatierung von QR-Faktorisierungen

Bei praktischen Aufgaben ist häufig nicht nur ein einzelnes Quadratmittelproblem zu lösen, sondern eine Folge

$$A_k x \simeq b^k \qquad (k = 1, \dots, K) \tag{1}$$

von K Problemen, die im folgenden Sinne benachbart sind: A_{k+1} entsteht aus A_k durch

- eine Rang-1-Modifikation

- Hinzufügen oder Streichen einer Spalte

- Hinzufügen oder Streichen einer Zeile,

vgl. 2.1.B1 und 9.3. Es ist dann günstig, die aufwendige QR-Faktorisierung nur für die Startmatrix A_1 zu berechnen und die Faktoren von A_{k+1} aus denen von A_k durch Aufdatierung mit geringem Aufwand zu bestimmen.

Zur Vermeidung von mehrfachen Indizes bezeichnen wir im folgenden die Matrix A_k mit A und die modifizierte Matrix A_{k+1} mit \overline{A} . Es gelte $A \in \mathbb{R}^{m,n}$ mit $m \ge n$, und A wie \overline{A} seien spaltenregulär.

Je nach Art der Faktorisierung unterscheiden wir zwei Fälle:

Fall 1: Gegeben ist die Faktorisierung

$$A = QR$$

mit $Q \in \mathbb{R}^{m,n}$, $Q^{\mathsf{T}}Q = I_n$ und der oberen Dreiecksmatrix $R \in \mathbb{R}^{n,n}$.

Fall 2: Gegeben ist die Faktorisierung

$$A = \hat{Q}\hat{R}$$

mit orthogonalem $\hat{Q} \in \mathbb{R}^{m,m}$ und der oberen Dreiecksmatrix $\hat{R} = \left(\frac{R}{O}\right) \in \mathbb{R}^{m,n}$.

Fall 1 entspricht der in 10.1, Fall 2 der in 10.2 behandelten Situation.

Gesucht ist die analoge Faktorisierung

 $ar{A}=ar{Q}ar{R}$ bzw. $ar{A}=ar{Q}ar{R}$

der modifizierten Matrix \bar{A} . Für die Berechnung der Faktoren \bar{Q} , \bar{R} bzw. \hat{Q} , \hat{R} nach den in 10.1 und 10.2 beschriebenen Methoden sind

$$\sim (K_1 m n^2 + K_2 n^3) \text{ opms} \tag{2}$$

erforderlich. Wir werden unten sehen, daß die Aufdatierung der Faktoren Q, Rbzw. \hat{Q}, \hat{R} von A dagegen mit einem Aufwand von

$$\sim (K_4 m n + K_5 n^2) \text{ opms}$$
 (3)

möglich, also billiger als die Neuberechnung ist. Falls bei Faktorisierungen des Typs $A = \hat{Q}\hat{R}$ die Matrix \hat{Q} explizit berechnet wird, kommt in (2) bzw. (3) ein Term $\sim K_3 m^2 n$ bzw. $\sim K_6 m^2$ hinzu. Allerdings wird \hat{Q} in den meisten Anwendungen überhaupt nicht explizit benötigt, siehe 10.3.6 und 10.3.7 unten.

Die im folgenden für den Fall 1 angegebenen Aufdatierungsalgorithmen beruhen sämtlich auf der Zerlegung eines Vektors $d \in \mathbb{R}^m$ in der Form

$$\boldsymbol{d} = \boldsymbol{P}\boldsymbol{d} + (\boldsymbol{I} - \boldsymbol{P})\,\boldsymbol{d} = \boldsymbol{Q}\boldsymbol{r} + \boldsymbol{p} \tag{4}$$

mit $P := QQ^{\intercal}$ als Projektor auf $\mathcal{R}(Q) = \mathcal{R}(A)$ und

$$\boldsymbol{r} := \boldsymbol{Q}^{\mathsf{T}} \boldsymbol{d}, \qquad \boldsymbol{p} := \boldsymbol{d} - \boldsymbol{Q} \boldsymbol{r}. \tag{5}$$

Dabei gilt

$$Qr \in \mathcal{R}(Q), \quad p \perp \mathcal{R}(Q)$$

und

$$\|d\|^2 = \|Qr\|^2 + \|p\|^2 = \|r\|^2 + \|p\|^2$$
,

vgl. 8.1.B. Die analoge Darstellung für den Fall 2 lautet

$$\boldsymbol{d} = \hat{\boldsymbol{Q}}\hat{\boldsymbol{r}} \quad ext{mit} \quad \hat{\boldsymbol{r}} := \hat{\boldsymbol{Q}}^{\intercal}\boldsymbol{d}, \qquad \|\hat{\boldsymbol{r}}\| = \|\boldsymbol{d}\|; \tag{6}$$

wegen $\mathcal{R}(\hat{Q}) = \mathbf{R}^m$ tritt der komplementäre Vektor p nicht auf.

10.3.1. Aufdatierung bei Rang-1-Modifikation

Aufgabe: Bestimme Faktorisierung von

$$\boldsymbol{A} = \boldsymbol{A} + \boldsymbol{u}\boldsymbol{v}^{\mathsf{T}}, \qquad \boldsymbol{u} \in \mathsf{R}^{m}, \quad \boldsymbol{v} \in \mathsf{R}^{n}. \tag{7}$$

Algorithmus 1: Gegeben ist die Faktorisierung A = QR.

S1: Zerlege *u* gemäß

$$\boldsymbol{u} = \boldsymbol{Q}\boldsymbol{r} + \boldsymbol{p} \quad \text{mit} \quad \boldsymbol{r} := \boldsymbol{Q}^{\mathsf{T}}\boldsymbol{u}, \qquad \boldsymbol{p} := \boldsymbol{u} - \boldsymbol{Q}\boldsymbol{r} \perp \boldsymbol{\mathcal{R}}(\boldsymbol{Q}).$$
 (8)

Berechne $\varrho := \|p\|$. Falls $\varrho > 0$, setze $q := p/\varrho$. Andernfalls bestimme beliebigen Vektor $q \perp \mathcal{R}(Q)$ mit $\|q\| = 1$. Schreibe (7) mit $p = \varrho q$ in der Form

$$\bar{\boldsymbol{A}} = \boldsymbol{Q}\boldsymbol{R} + (\boldsymbol{Q}\boldsymbol{r} + \varrho\boldsymbol{q})\,\boldsymbol{v}^{\mathsf{T}} = (\boldsymbol{Q} \mid \boldsymbol{q})\left(\left(\frac{\boldsymbol{R}}{\boldsymbol{o}^{\mathsf{T}}}\right) + \left(\frac{\boldsymbol{r}}{\varrho}\right)\boldsymbol{v}^{\mathsf{T}}\right) =: \boldsymbol{Q}_{1}(\boldsymbol{R}_{0} + \boldsymbol{r}^{0}\boldsymbol{v}^{\mathsf{T}})$$
(9)

mit der spaltenorthonormalen Matrix $Q_1 := (Q \mid q)$ und $R_0 := \left(\frac{R}{o^{\intercal}}\right), r^0$ $:= \left(\frac{r}{o}\right).$

S2: Bestimme *n* Givens-Drehungen $G_{n,n+1}, \ldots, G_{12}$ so, daß die Komponenten $(r^0)_i$ $(i = n + 1, \ldots, 2)$ von r^0 sukzessive zu 0 gemacht werden, vgl. 3.3.14 und U 3.3.5. Dann gilt

$$Gr^0 := G_{12} \cdots G_{n,n+1} r^0 = \varrho_1 e^1, \qquad \varrho_1 = \pm \|r^0\|,$$

und

$$GR_0 = G_{12} \cdots G_{n,n+1}R_0 =: R_1$$

ist von oberer Hessenbergform. Damit geht (9) in

$$\bar{\boldsymbol{A}} = (\boldsymbol{Q}_1 \boldsymbol{G}^{\mathsf{T}}) \left(\boldsymbol{G} \boldsymbol{R}_0 + \boldsymbol{G} \boldsymbol{r}^0 \boldsymbol{v}^{\mathsf{T}} \right) = (\boldsymbol{Q}_1 \boldsymbol{G}^{\mathsf{T}}) \left(\boldsymbol{R}_1 + \varrho_1 \boldsymbol{e}^1 \boldsymbol{v}^{\mathsf{T}} \right) =: \boldsymbol{Q}_2 \boldsymbol{R}_2 \tag{10}$$

mit spaltenorthonormalem $Q_2 := Q_1 G^{\intercal}$ und der oberen Hessenbergmatrix $R_2 := R_1 + \varrho_1 e^1 v^{\intercal}$ über.

S3: Bestimme *n* Givens-Drehungen $\tilde{G}_{12}, ..., \tilde{G}_{n,n+1}$ so, daß die Subdiagonalelemente $(\mathbf{R}_2)_{i+1,i}$ (j = 1, ..., n) von \mathbf{R}_2 sukzessive zu 0 gemacht werden. Dann ist

$$ilde{G} oldsymbol{R}_2 := oldsymbol{ ilde{G}}_{n,n+1} \cdots oldsymbol{ ilde{G}}_{12} oldsymbol{R}_2 =: oldsymbol{R}_3 =: \left(oldsymbol{rac{oldsymbol{R}}{oldsymbol{o}}}
ight)$$

von oberer Dreiecksform, und $Q_3 := Q_2 \tilde{G}^{\intercal} =: (\bar{Q} \mid \bar{q})$ ist spaltenorthonormal. Schreibe (10) in der Form

$$ar{A}=Q_2R_2=(Q_2 ilde{G}^{\intercal})\,(ilde{G}R_2)=Q_3R_3=(ar{Q}\midar{q})\left(rac{ar{R}}{o^{\intercal}}
ight)=ar{Q}ar{R}.$$

Man beachte dabei, daß im Ausnahmefall $\varrho = 0$ in S2 und S3 jeweils $G_{n,n+1} = \tilde{G}_{n,n+1}$ = I gilt. Der Vektor q trägt daher zu \bar{Q} und \bar{R} nichts bei und braucht deshalb überhaupt nicht explizit bestimmt zu werden.

Algorithmus 2: Gegeben ist die Faktorisierung $A = \hat{Q}\hat{R}$

S1: Stelle u in der Form

$$\boldsymbol{u} = \hat{\boldsymbol{Q}}\hat{\boldsymbol{r}}, \qquad \hat{\boldsymbol{r}} := \hat{\boldsymbol{Q}}^{\mathsf{T}}\boldsymbol{u}$$
 (11)

dar und schreibe (7) als

$$\bar{\boldsymbol{A}} = \hat{\boldsymbol{Q}}\hat{\boldsymbol{R}} + \hat{\boldsymbol{Q}}\hat{\boldsymbol{r}}\boldsymbol{v}^{\mathsf{T}} = \hat{\boldsymbol{Q}}(\hat{\boldsymbol{R}} + \hat{\boldsymbol{r}}\boldsymbol{v}^{\mathsf{T}}). \tag{12}$$

S2: Bestimme m - 1 Givens-Drehungen $G_{m-1,m}, \ldots, G_{12}$ so, daß

$$G\hat{\boldsymbol{r}} := \boldsymbol{G}_{12}\cdots \boldsymbol{G}_{m-1,m}\hat{\boldsymbol{r}} = \varrho_1 \boldsymbol{e}^1, \qquad \varrho_1 = \pm \|\hat{\boldsymbol{r}}\| = \pm \|\boldsymbol{u}\|,$$

gilt. Dann ist

$$G\hat{R} = G_{12} \cdots G_{m-1,m} \left(rac{R}{O}
ight) = G_{12} \cdots G_{\hat{l},\hat{l}+1} \left(rac{R}{O}
ight) =: \hat{R}_{1}$$

von oberer Hessenbergform, wobei $\hat{l} := \min \{m - 1, n\}$. Damit geht (12) in

$$\bar{\boldsymbol{A}} = (\hat{\boldsymbol{Q}}\boldsymbol{G}^{\mathsf{T}}) \left(\boldsymbol{G}\hat{\boldsymbol{R}} + \boldsymbol{G}\hat{\boldsymbol{r}}\boldsymbol{v}^{\mathsf{T}}\right) = (\hat{\boldsymbol{Q}}\boldsymbol{G}^{\mathsf{T}}) \left(\hat{\boldsymbol{R}}_{1} + \varrho_{1}\boldsymbol{e}^{1}\boldsymbol{v}^{\mathsf{T}}\right) =: \hat{\boldsymbol{Q}}_{2}\hat{\boldsymbol{R}}_{2}$$
(13)

mit orthogonalem $\hat{Q}_2 := \hat{Q}G^{\intercal}$ und der oberen Hessenbergmatrix $\hat{R}_2 := \hat{R}_1 + \varrho_1 e^1 v^{\intercal}$ über.

S3: Bestimme \hat{l} Givens-Drehungen $\tilde{G}_{12}, \dots, \tilde{G}_{\hat{l},\hat{l}+1}$ so, daß \hat{R}_2 sukzessive auf obere Dreiecksform

$$ilde{G} \hat{R}_2 := ilde{G}_{\hat{l},\hat{l}+1} \cdots ilde{G}_{12} \hat{R}_2 =: \dot{ar{R}} =: \left(rac{R}{O}
ight)$$

transformiert wird. Schreibe (13) mit der orthogonalen Matrix $\hat{ar{Q}} := \hat{m{Q}}_2 ilde{m{G}}^\intercal$ als

$$ar{A}=(\hat{oldsymbol{Q}}_2 ilde{G}^{\intercal})\,(ilde{G}\hat{oldsymbol{R}}_2)=ar{oldsymbol{Q}}ar{oldsymbol{R}}\,.$$

Im Fall m = 5, n = 3 verläuft die Transformation aus S2 von Algorithmus 2 nach folgendem Muster:

Es entsteht also tatsächlich eine obere Hessenbergmatrix. Die eingerahmten Felder geben dabei diejenigen Elemente an, die sich im jeweiligen Teilschritt verändern.

10.3.2. Aufdatierung bei Hinzufügen einer Spalte Aufgabe: Bestimme Faktorisierung von

$$\bar{\boldsymbol{A}} = (\boldsymbol{A} \mid \boldsymbol{a}), \qquad \boldsymbol{a} \in \mathbf{R}^{\boldsymbol{m}}. \tag{14}$$

Algorithmus 1: Gegeben ist die Faktorisierung A = QR. Zerlege a gemäß

$$a = Qr + p$$
 mit $r := Q^{\mathsf{T}}a$, $p = a - Qr$. (15)

Setze $\varrho := \| \boldsymbol{p} \|, \, \boldsymbol{q} := \boldsymbol{p} / \varrho$ und schreibe (14) in der Form

$$ar{m{A}} = (m{Q}m{R} \mid m{Q}m{r} + arrhom{q}) = (m{Q} \mid m{q}) \left(egin{matrix} m{R} \mid m{r} \ m{o^{\intercal} \mid arrho} \end{pmatrix} =: ar{m{Q}}m{m{R}}$$

mit der spaltenorthonormalen Matrix $\overline{Q} := (Q \mid q)$ und der oberen Dreiecksmatrix $\overline{R} = \left(\frac{R \mid r}{o^{\intercal} \mid \varrho}\right)$. Wegen rang $(\overline{A}) = n + 1$ ist dabei $a \notin \mathcal{R}(A)$, also $p \neq o$.

Algorithmus 2: Gegeben ist die Faktorisierung $A = \hat{Q}\hat{R}$. S1: Stelle *a* in der Form

$$\boldsymbol{a} = \hat{\boldsymbol{Q}}\hat{\boldsymbol{r}} \quad \text{mit} \quad \hat{\boldsymbol{r}} := \hat{\boldsymbol{Q}}^{\mathsf{T}}\boldsymbol{a} =: \left(\frac{\boldsymbol{r}}{\tilde{\boldsymbol{r}}}\right)_{\substack{n=1\\ \vdots\\m}}^{\frac{1}{n}}$$
(16)

dar. Schreibe (14) als

$$\bar{\boldsymbol{A}} = (\hat{\boldsymbol{Q}}\hat{\boldsymbol{R}} \mid \hat{\boldsymbol{Q}}\hat{\boldsymbol{r}}) = \hat{\boldsymbol{Q}}(\hat{\boldsymbol{R}} \mid \hat{\boldsymbol{r}}) =: \hat{\boldsymbol{Q}}\hat{\boldsymbol{R}}_{1}, \qquad \hat{\boldsymbol{R}}_{1} := (\hat{\boldsymbol{R}} \mid \hat{\boldsymbol{r}}).$$
(17)

S2: Bestimme orthogonales $\tilde{H} \in \mathbb{R}^{m-n,m-n}$ als Householder Spiegelung oder als Produkt von m - n - 1 Givens-Drehungen so, daß

$$ilde{H} ilde{r}=arrho e^1, \qquad arrho=\pm \| ilde{r}\|,$$

gilt. Setze

$$H := \left(\frac{I_n \mid O}{O \mid \tilde{H}}\right)$$

und schreibe (17) in der Form

$$ar{A} = \hat{Q}\hat{R}_1 = (\hat{Q}H^{\intercal})(H\hat{R}_1) =: \overline{\hat{Q}}\overline{\hat{R}}$$

mit orthogonalem $\hat{oldsymbol{Q}} := \hat{oldsymbol{Q}} H^\intercal$ und der oberen Dreiecksmatrix

$$\hat{m{ar{R}}} := H \hat{m{R}}^1 = \left(rac{I_n \mid O}{O \mid H}
ight) \left(rac{m{R} \mid m{r}}{O \mid m{ ilde{r}}}
ight) = \left(rac{m{R} \mid m{r}}{O \mid arrho e^1}
ight).$$

Offensichtlich stellen die Algorithmen 10.3.2 gerade den (n + 1)-ten Schritt des Gram-Schmidt- bzw. Householder- oder Givens-Verfahrens dar.

10.3.3. Aufdatierung bei Streichen einer Spalte Aufgabe: Bestimme Faktorisierung von \overline{A} , wobei

$$\bar{A} = (\bar{a}^1, ..., \bar{a}^{n-1}) := (a^1, ..., a^{l-1}, a^{l+1}, ..., a^n), \quad n > 1,$$
 (18)

20 Schwetlick, Numerische Algebra

durch Streichen der *l*-ten Spalte a^l aus $A = (a^1, ..., a^l, ..., a^n)$ entsteht $(1 \le l \le n)$. Algorithmus 1: Gegeben ist die Faktorisierung A = QR. Bilde

$$\boldsymbol{R}_{0} := \begin{pmatrix} r_{11} \cdots r_{1,l-1} & r_{1,l+1} & \cdots & r_{1,n} \\ \vdots & \vdots & & & \\ & r_{l-1,l-1} & r_{l-1,l+1} & \cdots & r_{l-1,n} \\ & & r_{l,l+1} & \cdots & r_{l,n} \\ \boldsymbol{O} & & r_{l+1,l+1} & \cdots & r_{l+1,n} \\ & & & \vdots & \\ & & & & r_{n,n} \end{pmatrix} \in \boldsymbol{R}^{m,n-1}$$
(19)

aus $R \in \mathbb{R}^{n,n}$ durch Streichen der *l*-ten Spalte. Bestimme n - l Givens-Drehungen $G_{l,l+1}, \ldots, G_{n-1,n}$ so, daß die Subdiagonalelemente $(R_0)_{j+1,j}$ $(j = l, \ldots, n - 1)$ von R_0 sukzessive zu 0 gemacht werden. Dann ist

$$GR_0 := G_{n-1,n} \cdots G_{l,l+1}R_0 =: R_1 =: \left(\frac{\overline{R}}{o^{\mathsf{T}}}\right)$$

mit einer oberen Dreiecksmatrix $\overline{R} \in \mathbb{R}^{n-1,n-1}$. Bilde $Q_1 := QG^{\intercal} = : (\overline{Q} \mid \overline{q})$ und schreibe (18) in der Form

$$ar{A} = oldsymbol{Q} oldsymbol{R}_0 = (oldsymbol{Q} oldsymbol{G}^{\intercal}) (oldsymbol{G} oldsymbol{R}_0) = oldsymbol{Q}_1 oldsymbol{R}_1 = (oldsymbol{ar{Q}} \mid oldsymbol{ar{q}}) igg(oldsymbol{ar{R}} oldsymbol{ar{Q}} oldsymbol{B} oldsymbol{A}) = oldsymbol{ar{Q}} oldsymbol{ar{R}}$$

Algorithmus 2: Gegeben ist die Faktorisierung $A = \hat{Q}\hat{R}$. Bilde $\hat{R}_0 = \left(\frac{R_0}{O}\right) \in \mathbb{R}^{m,n-1}$

analog zu Algorithmus 1 aus \hat{R} durch Streichen der *l*-ten Spalte. Bestimme n - l Givens-Drehungen $G_{l,l+1}, \ldots, G_{n-1,n}$ so, daß

$$G\hat{R}_0 := G_{n-1,n} \cdots G_{l,l+1}\hat{R}_0 =: \hat{\overline{R}} =: \left(\frac{R}{O}\right) \in \mathbb{R}^{m,n-1}$$

von oberer Dreiecksform ist. Bilde $\hat{\bar{Q}} := \hat{Q}G^{\intercal}$ und schreibe (18) in der Form

$$ar{A}=\hat{Q}\hat{R}_{0}=(\hat{Q}G^{\intercal})~(G\hat{R}_{0})=\ddot{ar{Q}}\ddot{ar{R}}$$
 .

10.3.4. Aufdatierung bei Hinzufügen einer Zeile

Aufgabe: Bestimme Faktorisierung von

$$ar{A} = \left(rac{a^{ au}}{A}
ight), \quad a \in \mathbb{R}^n.$$
 (20)

Algorithmus 1: Gegeben ist die Faktorisierung A = QR. Schreibe (20) in der Form

$$\bar{A} = \left(\frac{a^{\mathsf{T}}}{QR}\right) = \left(\frac{1 | o^{\mathsf{T}}}{o | Q}\right) \left(\frac{a^{\mathsf{T}}}{R}\right) =: Q_1 R_1$$
(21)

mit der spaltenorthonormalen Matrix Q_1 und der oberen Hessenbergmatrix R_1 .

Bestimme *n* Givens-Drehungen $G_{12}, \ldots, G_{n,n+1}$ so, daß

$$GR_1 := G_{n,n+1} \cdots G_{12}R_1 =: \left(\frac{\overline{R}}{o^{\mathsf{T}}}\right)$$

von oberer Dreiecksform ist. Schreibe (21) mit der spaltenorthonormalen Matrix

$$Q_1G^\intercal =: (ar{Q} \mid ar{q})$$

in der Form

$$ar{A} = (oldsymbol{Q}_1 G^{\intercal}) (GR_1) = (ar{oldsymbol{Q}} \mid ar{oldsymbol{q}}) \left(rac{oldsymbol{R}}{oldsymbol{o}^{\intercal}}
ight) = ar{oldsymbol{Q}} ar{oldsymbol{R}}$$

Algorithmus 2: Gegeben ist die Faktorisierung $A = \hat{Q}\hat{R}$. Schreibe (20) in der Form

$$\boldsymbol{A} = \left(\frac{\boldsymbol{a}^{\mathsf{T}}}{\hat{\boldsymbol{Q}}\hat{\boldsymbol{R}}}\right) = \left(\frac{1 | \boldsymbol{o}^{\mathsf{T}}}{\boldsymbol{o} | \hat{\boldsymbol{Q}}}\right) \left(\frac{\boldsymbol{a}^{\mathsf{T}}}{\hat{\boldsymbol{R}}}\right) =: \hat{\boldsymbol{Q}}_{1}\hat{\boldsymbol{R}}_{1}$$
(22)

mit orthogonalem \hat{Q}_1 und der oberen Hessenbergmatrix \hat{R}_1 . Transformiere \hat{R}_1 mittels *n* Givens-Drehungen $G_{12}, \ldots, G_{n,n+1}$ auf obere Dreiecksform

$$G\hat{R}_1 := G_{n,n+1} \cdots G_{12}\hat{R}_1 =: \hat{\overline{R}} =: \left(\frac{R}{O}\right)$$

Schreibe (22) mit der orthogonalen Matrix $\hat{oldsymbol{Q}} := \hat{oldsymbol{Q}}_1 G^\intercal$ in der Form

$$ar{A}=(\hat{m{Q}}_1G^{\intercal})\,(G\hat{m{R}}_1)=ar{ar{m{Q}}}ar{m{R}}^{-}.$$

10.3.5. Aufdatierung bei Streichen einer Zeile

Aufgabe: Bestimme Faktorisierung der Matrix \overline{A} , die aus A durch Streichen der *l*-ten Zeile $(1 \leq l \leq m)$ entsteht, d. h.

$$\bar{A} = \begin{pmatrix} -a^{1\mathsf{T}} & -\\ \vdots \\ -a^{l-1\mathsf{T}} -\\ -a^{l+1\mathsf{T}} -\\ \vdots \\ -a^{m\mathsf{T}} & - \end{pmatrix}, \quad \text{wobei} \quad A = \begin{pmatrix} -a^{1\mathsf{T}} & -\\ \vdots \\ -a^{l\mathsf{T}} & -\\ -\vdots & -\\ -a^{m\mathsf{T}} - \end{pmatrix}, \quad m > n.$$
(23)

Algorithmus 1: Gegeben ist die Faktorisierung A = QRS0: Bestimme Permutationsmatrix P_Z so, daß

$$\tilde{A} := P_Z A = \begin{pmatrix} -a^{\mathsf{T}} & -\\ -a^{\mathsf{I}\mathsf{T}} & -\\ -\vdots & -\\ -a^{\mathsf{I}-\mathsf{I}\mathsf{T}} & -\\ -a^{\mathsf{I}-\mathsf{I}} & -\\ -\vdots & -\\ -a^{\mathsf{m}\mathsf{T}} & - \end{pmatrix} = \begin{pmatrix} a^{\mathsf{T}} \\ \bar{A} \end{pmatrix} \quad \text{mit} \quad a := a^{\mathsf{I}}$$
(24)

20*

gilt. Dann ist

$$\tilde{A} = P_Z A = P_Z Q R = \tilde{Q} R \quad \text{mit} \quad \tilde{Q} := P_Z Q.$$
 (25)

S1: Betrachte die Rang-1-Modifikation

$$A_1 := \left(\frac{\boldsymbol{o}^{\mathsf{T}}}{\boldsymbol{A}}\right) = \tilde{\boldsymbol{A}} - \boldsymbol{e}^1 \boldsymbol{a}^{\mathsf{T}}.$$
(26)

Zerlege e^1 gemäß

$$e^1 = \tilde{Q}r + p, \quad r := \tilde{Q}^{\mathsf{T}}e^1, \quad p := e^1 - \tilde{Q}r.$$
 (27)

Dabei ist r^{\intercal} die erste Zeile von \tilde{Q} , braucht also nicht gesondert berechnet zu werden. Setze

$$\varrho := \|\boldsymbol{p}\| = \sqrt{1 - \|\boldsymbol{\tilde{Q}}^{\mathsf{T}} \boldsymbol{e}^{\mathsf{I}}\|^2} = \sqrt{1 - \|\boldsymbol{r}\|^2}.$$
(28)

Falls $\varrho > 0$, setze $\boldsymbol{q} := \boldsymbol{p}/\varrho$. Andernfalls bestimme $\tilde{\boldsymbol{p}} \neq \boldsymbol{o}, \ \tilde{\boldsymbol{p}} \perp \mathcal{R}(\tilde{\boldsymbol{Q}})$ wie folgt: Wähle Index $s \in \{2, ..., m\}$ mit $\|\tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^s\| < 1$. Ein solcher Index existiert, denn wegen $\|\tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^i\| \leq 1 \ (i = 1, ..., m)$ wäre sonst $\|\tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^i\| = 1 \ (i = 2, ..., m)$. Da nach (28) im Fall $\varrho = 0$ auch $\|\tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^1\| = 1$ gilt, folgte $n = \|\tilde{\boldsymbol{Q}}\|_F^2 = \sum_{i=1}^m \|\tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^i\|^2 = m$ im Widerspruch zu m > n. Setze

$$ilde{m{r}}:= ilde{m{Q}}^{\intercal}m{e}^{s}, \quad ilde{m{p}}:=m{e}^{s}- ilde{m{Q}} ilde{m{r}}, \quad ilde{arepsilon}:=\| ilde{m{p}}\|, \quad m{q}:= ilde{m{p}}/ ilde{arepsilon}, \tag{29}$$

wobei $\tilde{\varrho}^2 = 1 - \| \tilde{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{e}^s \|^2 > 0$ ist nach Wahl von s. In beiden Fällen ist

$$oldsymbol{p}=oldsymbol{e}^1- ilde{oldsymbol{Q}}oldsymbol{r}=arepsilonoldsymbol{q}\,,\qquad \|oldsymbol{q}\|=1\,,\qquad oldsymbol{q}\perp\mathcal{R}(oldsymbol{Q})\,.$$

Schreibe (26) in der Form

$$A_{1} = \tilde{A} - e^{1}a^{\mathsf{T}} = \tilde{Q}R - (\tilde{Q}r + \varrho q) a^{\mathsf{T}} = (\tilde{Q} \mid q) \left(\left(\frac{R}{o^{\mathsf{T}}} \right) - \left(\frac{r}{\varrho} \right) a^{\mathsf{T}} \right)$$
$$=: Q_{1}(R_{0} - r^{0}a^{\mathsf{T}})$$
(30)

mit der spaltenorthonormalen Matrix ${oldsymbol Q}_1 := ({oldsymbol ilde Q} \mid {oldsymbol q})$ und

$$\boldsymbol{R}_0 := \left(rac{\boldsymbol{R}}{\boldsymbol{o}^{\mathsf{T}}}
ight), \ \ \boldsymbol{r}^0 := \left(rac{\boldsymbol{r}}{arrho}
ight).$$

Dabei gilt

$$\boldsymbol{Q}_{1}^{\mathsf{T}}\boldsymbol{e}^{\mathsf{I}} = \left(\frac{\boldsymbol{\tilde{Q}}^{\mathsf{T}}}{\boldsymbol{q}^{\mathsf{T}}}\right)\boldsymbol{e}^{\mathsf{I}} = \left(\frac{\boldsymbol{\tilde{Q}}^{\mathsf{T}}\boldsymbol{e}^{\mathsf{I}}}{\boldsymbol{q}^{\mathsf{T}}\boldsymbol{e}^{\mathsf{I}}}\right) = \left(\frac{\boldsymbol{r}}{\varrho}\right) = \boldsymbol{r}^{\mathsf{0}},\tag{31}$$

denn im Fall $\rho > 0$ ist nach (28), (29)

$$q_1 := oldsymbol{e}^{1 op} oldsymbol{q} = oldsymbol{e}^{1 op} oldsymbol{p} = (1 - \| oldsymbol{ ilde{oldsymbol{Q}}}^{ op} oldsymbol{e}^{1} \|^2) / arrho = arrho$$
 ,

und im Fall $\rho = 0$ gilt wegen der Spaltenorthonormalität von Q_1

$$1 \ge \| ilde{m{Q}}_1^{\mathsf{T}} e^1\|^2 = \| ilde{m{Q}}^{\mathsf{T}} e^1\|^2 + q_1^2 = 1 + q_1^2,$$

also $q_1 = 0 = \varrho$.

S2: Wähle G wie in 10.3.1 als Folge von Givens-Drehungen so, daß

$$Gr^{0} := G_{12} \cdots G_{n,n+1} r^{0} = \varrho_{1} e^{1}, \qquad \varrho_{1} = \pm ||r^{0}|| = \pm 1$$
 (32)

gilt. Dann ist

 $\boldsymbol{G}\boldsymbol{R}_0 = \boldsymbol{G}_{12} \cdots \boldsymbol{G}_{n,n+1} \boldsymbol{R}_0 =: \boldsymbol{R}_1$

von oberer Hessenbergform. Schreibe (30) in der Form

$$A_1 = (Q_1 G^{\mathsf{T}}) (GR_0 - Gr^0 a^{\mathsf{T}}) = (Q_1 G^{\mathsf{T}}) (R_1 - \varrho_1 e^1 a^{\mathsf{T}}) =: Q_2 R_2$$
(33)

mit der oberen Hessenbergmatrix

$$oldsymbol{R}_2 := oldsymbol{R}_1 - arrho_1 e^1 a^\intercal = : \left(rac{oldsymbol{ar{r}}^\intercal}{oldsymbol{ar{R}}}
ight)$$

und der spaltenorthonormalen Matrix

$$oldsymbol{Q}_2 := oldsymbol{Q}_1 G^\intercal = egin{pmatrix} \pm 1 & oldsymbol{o}^\intercal \ \overline{oldsymbol{q}} & oldsymbol{ar{q}} \end{pmatrix} = egin{pmatrix} \pm 1 & oldsymbol{o}^\intercal \ \overline{oldsymbol{o}} \end{pmatrix}.$$

Man beachte dabei, daß aus (31), (32) die Beziehung $Q_2^{\mathsf{T}} e^1 = G Q_1^{\mathsf{T}} e^1 = G r^0$ $= \varrho_1 e^1$ folgt und daß $\|Q_2 e^1\|^2 = 1 + \|\overline{q}\|^2 = 1$, also $\overline{q} = \overline{o}$ gilt, denn Q_2 hat wie Q_1 normierte Spalten. Damit geht (33) in

$$egin{aligned} \left(egin{aligned} m{\sigma}^{\mathsf{T}} \ ar{m{A}} \end{matrix}
ight) &= A_1 = m{Q}_2 R_2 = \left(egin{aligned} \pm 1 & m{\sigma}^{\mathsf{T}} \ m{o} & ar{m{Q}} \end{matrix}
ight) \left(egin{aligned} ar{m{r}} \ ar{m{R}} \end{matrix}
ight) &= \left(egin{aligned} \pm ar{m{r}} \ ar{m{Q}} ar{m{R}} \end{matrix}
ight), \ ext{also } ar{m{r}} &= m{o} ext{ und} \ ar{m{A}} &= ar{m{Q}} ar{m{R}} \end{aligned}$$

über.

Algorithmus 2: Gegeben ist die Faktorisierung $A = \hat{Q}\hat{R}$.

S0: Bestimme P_z wie im Schritt S0 von Algorithmus 1. Dann gilt

$$\tilde{A} = P_Z A = P_Z \hat{Q} \hat{R} = \tilde{Q} \hat{R} \quad \text{mit} \quad \tilde{Q} := P_Z \hat{Q}.$$

S1: Schreibe (26) in der Form

$$A_1 = \tilde{A} - e^1 a^{\mathsf{T}} = \tilde{Q}(\hat{R} - \hat{r}a^{\mathsf{T}}) \quad \text{mit} \quad \hat{r} := \tilde{Q}^{\mathsf{T}} e^1$$
(34)

S2: Bestimme m - 1 Givens-Drehungen so, daß

$$G\hat{r} = G_{12} \cdots G_{m-1,m}\hat{r} = \varrho_1 e^1, \qquad \varrho_1 = \pm \|\hat{r}\| = \pm \|e^1\| = \pm 1$$

gilt. Dann ist

$$G\hat{R} := G_{12} \cdots G_{m-1,m} \hat{R} = G_{12} \cdots G_{\hat{l},\hat{l}+1} \hat{R} =: \hat{R}_1, \qquad \hat{l} := \min\{m-1, n\},$$

von oberer Hessenbergform. Schreibe (34) als

$$A_{1} = (\tilde{\boldsymbol{Q}}\boldsymbol{G}^{\mathsf{T}}) \left(\boldsymbol{G}\hat{\boldsymbol{R}} - \boldsymbol{G}\hat{\boldsymbol{r}}\boldsymbol{a}^{\mathsf{T}}\right) = (\tilde{\boldsymbol{Q}}\boldsymbol{G}^{\mathsf{T}}) \left(\hat{\boldsymbol{R}}_{1} - \varrho_{1}\boldsymbol{e}^{1}\boldsymbol{a}^{\mathsf{T}}\right) = : \hat{\boldsymbol{Q}}_{2}\hat{\boldsymbol{R}}_{2}$$
(35)

mit der oberen Hessenbergmatrix

$$\hat{m{R}}_2 := \hat{m{R}}_1 - arrho_1 m{e}^1 m{a}^{\mathsf{T}} = \left(rac{m{ar{r}}^{\mathsf{T}}}{\hat{m{R}}}
ight)$$

und der orthogonalen Matrix

$$\hat{Q}_2 := ilde{Q}G^\intercal = \left(rac{\pm 1 \left| oldsymbol{o}^\intercal
ight)}{ar{q} \left| ar{\hat{Q}}
ight)} = \left(rac{\pm 1 \left| oldsymbol{o}^\intercal
ight)}{oldsymbol{o} \left| ar{\hat{Q}}
ight)};$$

man beachte $\hat{Q}_2^{\mathsf{T}} e^1 = \varrho_1 e^1$ und die Spaltennormiertheit von \hat{Q}_2 . Damit geht (35) in

$$\begin{pmatrix} \boldsymbol{o}^{\mathsf{T}} \\ \boldsymbol{\bar{A}} \end{pmatrix} = \boldsymbol{A}_{1} = \begin{pmatrix} \pm 1 & \boldsymbol{o}^{\mathsf{T}} \\ \boldsymbol{o} & \boldsymbol{\dot{\bar{Q}}} \end{pmatrix} \begin{pmatrix} \boldsymbol{\bar{r}}^{\mathsf{T}} \\ \boldsymbol{\bar{R}} \end{pmatrix} = \begin{pmatrix} \pm \boldsymbol{\bar{r}}^{\mathsf{T}} \\ \boldsymbol{\bar{Q}} \boldsymbol{\bar{R}} \end{pmatrix},$$

also $\boldsymbol{\bar{r}} = \boldsymbol{o}$ und
 $\boldsymbol{\bar{A}} = \boldsymbol{\dot{\bar{Q}}} \boldsymbol{\bar{R}}$ (36)

Die Algorithmen aus 10.3.5 sind auf den Fall $u = -e^1$ zugeschnittene Varianten von 10.3.1, wobei im Algorithmus 1 der für p = o benötigte Vektor $q \perp \mathcal{R}(Q)$ jedoch explizit berechnet wird.

10.3.6. Bemerkung. Die Anwendung der beschriebenen Aufdatierungsalgorithmen hängt von der Art der Faktorisierung und der Aufgabenstellung ab:

(i) Für Faktorisierungen $A_k = Q_k R_k$ nach Fall 1 sollten beide Faktoren Q_k und R_k explizit mitgeführt werden. Aus ihnen lassen sich alle weiteren benötigten Größen wie Quadratmittellösungen x^k , Residuen $b^k - Ax^k$ u. a. einfach berechnen. Der Aufwand für die Aufdatierung ist vom Typ (3), also klein.

(ii) Die explizite Mitführung beider Faktoren Q_k , R_k bzw. \hat{Q}_k , \hat{R}_k empfiehlt sich ebenfalls, wenn eine Folge regulärer linearer Gleichungssysteme $A_k x^k = b^k$ mit

$$A_{k+1} = A_k + \boldsymbol{u}^k(\boldsymbol{v}^k)^\mathsf{T} \in \mathsf{R}^{n,n} \tag{37}$$

zu lösen ist und die rechten Seiten nicht a priori bekannt sind. Man beachte dabei, daß die Faktorisierungen gemäß Fall 1 und Fall 2 für m = n bis auf das Vorzeichen der Spalten von Q und Zeilen von R identisch sind. Probleme dieses Typs treten in der Optimierung und bei der Lösung nichtlinearer Gleichungssysteme auf.

(iii) Für Faktorisierungen $A_k = \hat{Q}_k \hat{R}_k$ gemäß Fall 2 ist die explizite Mitführung von \hat{Q}_k aus Aufwandsgründen nicht zu empfehlen, sofern $m \gg n$ ist. Ein Kompromiß für nicht zu großes K wäre die Berechnung von \hat{Q}_1 in Produktform nach dem Householder- oder Givens-Verfahren und das Abspeichern der zugehörigen Daten v_{ik} bzw. ξ_{ik} auf dem Platz von A_1 nach 10.2.4. Die bei der Aufdatierung benötigten Givens-Drehungen G_{ij} werden durch die zugehörigen ξ_{ij} repräsentiert; für diese muß dann zusätzlicher Speicherplatz bereitgestellt werden. Der Mehraufwand an Speicherplatz zum Abspeichern der Produktform von \hat{Q}_k und an Rechenoperationen zur Berechnung von $\hat{Q}_k^{\mathsf{T}} b$ bzw. $\hat{Q}_k f$ in zu 10.2.6 analoger Form ist bei dieser Version proportional zu K.

(iv) Wenn die vorkommenden rechten Seiten b^j , die Vektoren u^j der Rang-1-Modifikationen und die hinzuzufügenden Spalten a^j a priori bekannt sind, können die Dreiecksfaktoren $\hat{R}_k = \hat{Q}_k^{\mathsf{T}} A_k$, die transformierten Größen $\hat{Q}_k^{\mathsf{T}}(b^j, u^j, a^j)$ und damit die Quadratmittellösungen x^k und die Residuumsnormen $nr_k := ||b^k - A_k x^k||$ ohne explizite Kenntnis des Orthogonalfaktors \hat{Q}_k berechnet werden. Dasselbe trifft für das Streichen von Spalten und Hinzufügen von Zeilen zu. Das Streichen der *l*-ten Zeile nach dem in 10.3.5 angegebenen Algorithmus 2 erfordert die Kenntnis von $\tilde{Q}^{\mathsf{T}}e^1 = \hat{Q}^{\mathsf{T}}e^l$, und dieser Vektor kann ebenfalls wie b^j , u^j oder a^j aufdatiert werden. Es gibt jedoch andere Algorithmen, die Zeilenstreichungen ohne Rückgriff auf $\hat{Q}^{\mathsf{T}}e^l$ ermöglichen. Für Details und Literaturhinweise sei auf das nachfolgende Beispiel 10.3.7 und B 10.6 verwiesen.

(v) Die Realisierung der Givens-Drehungen kann sowohl in expliziter wie auch in impliziter Form erfolgen. Im letzten Fall sind die Skalierungsfaktoren mitzuführen, siehe B 10.3.

10.3.7. Beispiel. Gegeben sind die Quadratmittelprobleme

$$A_k x^k \simeq b \qquad (k = 1, 2, 3) \tag{38}$$

 $_{\rm mit}$

$$A_1 := (a^1, a^2, a^3), \qquad A_2 := (a^1, a^2, a^3, a^4), \qquad A_3 := A_2 + uv^{\mathsf{T}}.$$
 (39)

Sie werden durch die Eingangsdaten

$$\boldsymbol{B} := (\boldsymbol{a}^1, \boldsymbol{a}^2, \boldsymbol{a}^3, \boldsymbol{a}^4, \boldsymbol{u}, \boldsymbol{b}) = (\boldsymbol{A}_1 \mid \boldsymbol{a}^4 \mid \boldsymbol{u} \mid \boldsymbol{b}) \in \boldsymbol{\mathsf{R}}^{\boldsymbol{m}, \boldsymbol{6}} \quad \text{und} \quad \boldsymbol{v} \in \boldsymbol{\mathsf{R}}^4 \tag{40}$$

dargestellt. Wenn $\hat{Q} \in \mathbf{R}^{m,m}$ irgendeine orthogonale Matrix bezeichnet, sind die durch die transformierten Eingangsdaten

$$\hat{oldsymbol{Q}}^{\intercal}oldsymbol{B} = (\hat{oldsymbol{Q}}^{\intercal}oldsymbol{A}_1 \mid \hat{oldsymbol{Q}}^{\intercal}oldsymbol{a}^4 \mid \hat{oldsymbol{Q}}^{\intercal}oldsymbol{u} \mid \hat{oldsymbol{Q}}^{\intercal}oldsymbol{b})$$

festgelegten Probleme zu den durch **B** charakterisierten äquivalent: Sie haben dieselben Lösungen x^k und dieselben Residuumsnormen $nr_k = \|\hat{Q}^{\intercal}(b - A_k x^k)\|$ $= \|b - A_k x^k\|$. Die Probleme (38), (39) können daher wie folgt gelöst werden:

S1 (k = 1): Bestimme \hat{Q}_1^{T} als Produkt von Householder-Spiegelungen oder Givens-Drehungen gemäß 10.2.4 so, daß

$$m{B}_1 := \hat{m{Q}}_1^{\mathsf{T}} m{B} =: (\hat{m{R}}_1 \mid \hat{m{r}}^{4,1} \mid \hat{m{r}}^1 \mid \hat{m{c}}^1) = \hat{m{Q}}_1^{\mathsf{T}} (m{A}_1 \mid m{a}^4 \mid m{u} \mid m{b})$$

mit einer oberen Dreiecksmatrix $\hat{R}_1 = \hat{Q}_1^{\mathsf{T}} A_1$ gilt. Berechne x^1 und nr_1 aus \hat{R}_1 und $\hat{c}^1 = \hat{Q}_1^{\mathsf{T}} b$ gemäß 10.2.9.

S2 (k = 2): Bestimme **H** entsprechend 10.3.2 (Algorithmus 2, S2) mit

$$\hat{r} := \hat{r}^{4,1} = \hat{Q}_1^{\mathsf{T}} a^4$$

Transformiere \boldsymbol{B}_1 in

$$B_2 := HB_1 =: (\hat{R}_2 \mid \hat{r}^2 \mid \hat{c}^2) = H\hat{Q}_1^{\mathsf{T}}B = \hat{Q}_2^{\mathsf{T}} ((A_1 \mid a^4) \mid u \mid b) = \hat{Q}_2^{\mathsf{T}} (A_2 \mid u \mid b)$$

mit $\hat{Q}_2 := \hat{Q}_1 H^{\mathsf{T}}$ und der oberen Dreiecksmatrix $\hat{R}_2 := (H\hat{R}_1 | H\hat{r}^{4,1}) = \hat{Q}_2^{\mathsf{T}} A_2$. Berechne x^2 und nr_2 aus \hat{R}_2 und $\hat{c}^2 = \hat{Q}_2^{\mathsf{T}} b$.

S3 (k = 3): Bestimme G entsprechend 10.3.1 (Algorithmus 2, S2) mit

$$\hat{\boldsymbol{r}} := \hat{\boldsymbol{r}}^2 = \hat{\boldsymbol{Q}}_2^{\mathsf{T}} \boldsymbol{u}$$

Transformiere B_2 in

$$\boldsymbol{B}_{21} := \boldsymbol{G}\boldsymbol{B}_2 =: (\hat{\boldsymbol{R}}_{21} \mid \varrho_1 \boldsymbol{e}^1 \mid \hat{\boldsymbol{c}}^{2,1}) = \boldsymbol{G}\hat{\boldsymbol{Q}}_2^{\mathsf{T}}\boldsymbol{B} = \hat{\boldsymbol{Q}}_{21}^{\mathsf{T}}(\boldsymbol{A}_2 \mid \boldsymbol{u} \mid \boldsymbol{b})$$

mit der oberen Hessenbergmatrix $\hat{R}_{21} := G\hat{R}_2$ und $\hat{Q}_{21} := \hat{Q}_2 G^\intercal$.

Führe die Rang-1-Modifikation $A_3 = A_2 + uv^{\intercal}$ mit den durch \hat{Q}_{21}^{\intercal} transformierten Daten B_{21} gemäß

$$\tilde{B}_{22} := (\hat{R}_{21} + \varrho_1 e^1 v^\top | \varrho_1 e^1 | \hat{c}^{2,1}) = \hat{Q}_{21}^\top (A_2 + uv^\top | u | b) = \hat{Q}_{21}^\top (A_3 | u | b)$$

durch, wobei $\hat{R}_{22} := \hat{R}_{21} + \varrho_1 e^1 v^{\intercal}$ von oberer Hessenbergform ist, und streiche die nicht mehr benötigte fünfte Spalte von \tilde{B}_{22} , d. h., setze

$$\hat{B}_{22} := (\hat{R}_{22} \mid \hat{c}^{2,1}) = \hat{Q}_{21}^{\mathsf{T}}(A_3 \mid b).$$

Bestimme \tilde{G} entsprechend 10.3.1 (Algorithmus 2, S3) mit $\hat{R}_2 := \hat{R}_{22}$. Transformiere B_{22} in

$$oldsymbol{B}_3 := oldsymbol{ ilde{G}} oldsymbol{B}_{22} = (oldsymbol{\hat{R}}_3 \mid oldsymbol{\hat{c}}^3) = oldsymbol{ ilde{G}} oldsymbol{\hat{Q}}_{21}^{ op} oldsymbol{B} = oldsymbol{\hat{Q}}_3^{ op} (A_3 \mid oldsymbol{b})$$

mit $\hat{Q}_3 := \hat{Q}_{21}\tilde{G}^{\mathsf{T}} = \hat{Q}_2 G^{\mathsf{T}}\tilde{G}^{\mathsf{T}}$ und der oberen Dreiecksmatrix $\hat{R}_3 = \hat{Q}_3^{\mathsf{T}}A_3$, wobei \bar{B} aus B durch Streichen der fünften Spalte entsteht.

Berechne x^3 und nr_3 aus \hat{R}_3 und $\hat{c}^3 = \hat{Q}_3^{\mathsf{T}} \hat{b}$.

Die Grundidee des Vorgehens in 10.3.7 ist die folgende: Statt die Orthogonalfaktoren \hat{Q}_k nach der Vorschrift

$$\hat{\boldsymbol{Q}}_{k+1} := \hat{\boldsymbol{Q}}_k \boldsymbol{G}_k^\mathsf{T} \tag{41}$$

mit den im k-ten Schritt vorkommenden orthogonalen Transformationen G_k^{T} explizit aufzudatieren und damit die transformierten Eingangsdaten

$$\boldsymbol{B}_{k+1} := \hat{\boldsymbol{Q}}_{k+1}^{\mathsf{T}} \boldsymbol{B} \tag{42}$$

zu berechnen, werden diese gemäß $B_{k+1} = (\hat{Q}_k G_k^{\mathsf{T}})^{\mathsf{T}} B = G_k \hat{Q}_k^{\mathsf{T}} B$, d. h.

$$\boldsymbol{B}_{k+1} := \boldsymbol{G}_k \boldsymbol{B}_k \tag{43}$$

aufdatiert. Die Matrix G_k ist dabei das Produkt sämtlicher auf \hat{R}_k angewandter Drehungen bzw. Spiegelungen. In dieser Weise lassen sich auch Spaltenstreichungen, Zeilenhinzufügungen und — falls e^{l_j} mit in B aufgenommen wird — das Streichen von Zeilen l_j realisieren. In allen Fällen wird weder \hat{Q}_k^{T} noch G_k explizit benötigt; die Transformationen werden sukzessive auf das aktuelle B_k angewendet und brauchen auch nicht gespeichert zu werden.

10.3.8. Bemerkung. (i) Für die Realisierung der Aufdatierungsalgorithmen ist wesentlich, daß die in exakter Arithmetik gültigen Orthogonalitätsbeziehungen auch bei Computerrechnung in ausreichendem Maße gewährleistet sind. Die Startfaktoren Q_1 bzw. \hat{Q}_1 müssen daher ausreichend orthonormal sein. Im Fall 1 sollte Q_1 daher entweder nach dem MGS-Verfahren mit Re-Orthogonalisierung oder durch Householder- bzw. Givens-Orthogonalisierung und spaltenweise Berechnung von $Q = (q^1, ..., q^n)$ bestimmt werden; im Fall 2 empfiehlt sich die Householder- oder Givens-Orthogonalisierung.

In den zum Fall 1 gehörenden Aufdatierungsalgorithmen müssen die dort zu berechnenden orthogonalen Zerlegungen (4), (5) ein p liefern, das genügend orthogonal zu $\mathcal{R}(Q)$ ist. Wenn p klein im Verhältnis zu d ist, tritt bei der Berechnung von p = d - Qr Auslöschung auf, die die Orthogonalität wesentlich beeinträchtigen kann. In diesem Fall sollte re-orthogonalisiert werden, etwa nach der Vorschrift

$$m{r}:=m{Q}^{\intercal}m{d},\,m{p}:=m{d}-m{Q}m{r}$$

 \mathbf{f} $\|m{p}\| \leq 0.1$ $\|m{d}\|$ then $[m{s}:=m{Q}^{\intercal}m{p},\,m{p}:=m{p}-m{Q}m{s},\,m{r}:=m{r}+m{s}]$

Dabei ist es günstig, wenn die Realisierung in modifizierter Form analog zu 10.1.5 erfolgt. Der Faktor 0.1 hat dabei nur exemplarischen Charakter. Er bewirkt, daß die aufwandsverdoppelnde Re-Orthogonalisierung nur dann ausgeführt wird, wenn die angenäherte Orthogonalität um mehr als den Faktor 10 gestört werden könnte.

In 10.3.5 (Algorithmus 1, S1) sollte die Bedingung $\|\tilde{Q}^{\intercal}e^{i}\| < 1$ als erfüllt angesehen werden, wenn die berechnete Norm der *i*-ten Zeile von \tilde{Q} der Bedingung

$$\|\tilde{Q}^{\mathsf{T}} e^{i}\| \leq 1 - m r$$

genügt, andernfalls wird $\|\tilde{Q}^{\intercal}e^{i}\|$ als 1 angesehen. In Entsprechung zur obigen Regel sollte p bzw. \tilde{p} re-orthogonalisiert werden, wenn

$$\sqrt{0.99} \leq \|\tilde{\boldsymbol{Q}}^{\intercal} \boldsymbol{e}^{i}\| \leq 1 - mv$$
 für $i = 1$ bzw. $i = s$

gilt.

(ii) Man kann zeigen, daß die Aufdatierungsalgorithmen für das Hinzufügen von Spalten oder Zeilen in folgendem Sinne numerisch gutartig sind: Wenn $A + \delta A = QR$ mit kleinem δA gilt, ist auch $\bar{A} + \delta \bar{A} = \bar{Q}\bar{R}$ mit kleinem $\delta \bar{A}$. Wenn δA spaltenweise klein in bezug auf A ist, trifft dies auch für $\delta \bar{A}$ zu. Gutartigkeit liegt auch für Rang-1-Modifikationen und das Streichen von Spalten oder Zeilen vor, sofern die modifizierte Matrix \bar{A} in der Norm nicht signifikant kleiner als A ist. Wenn $\|\bar{A}\|$ wesentlich kleiner als $\|A\|$ ist, wird $\|\delta \bar{A}\|$ im allgemeinen nicht klein in bezug auf $\|\bar{A}\|$ sein, so daß die berechneten Faktoren \bar{Q}, \bar{R} eine weit von A entfernte Matrix reproduzieren. In diesem Fall empfiehlt sich eine Neuberechnung der Faktorisierung aus \bar{A} . Aus Platzgründen soll auf die präzise Formulierung und den Beweis dieser Aussagen verzichtet werden. \Box

Übungsaufgaben

U 10.3.1. Man überlege sich: Wenn $ar{A}$ aus A durch Modifikation der l-ten Spalte gemäß

$$\bar{A} = A + ue^{l\tau} \tag{44}$$

entsteht und P_S diejenige Permutationsmatrix bezeichnet, die $A = (a^1, ..., a^l, ..., a^n)$ in $\tilde{A} = AP_S^{\mathsf{T}} = (a^1, ..., a^{l-1}, a^{l+1}, ..., a^n, a^l)$ transformiert, kann eine QR-Faktorisierung von AP_S^{T} wie folgt berechnet werden:

- S1: Wie 10.3.1 (Algorithmus 1, S1)
- S2: Schreibe (9) in der Form

$$\bar{A}P_{S}^{\mathsf{T}} = Q_{1}(R_{0} + r^{0}e^{l\mathsf{T}})P_{S}^{\mathsf{T}} = Q_{1}(R_{0}P_{S}^{\mathsf{T}} + r^{0}e^{n\mathsf{T}}) =: Q_{2}R_{2}$$

$$\tag{45}$$

mit $Q_2 := Q_1$ und der oberen Hessenbergmatrix $R_2 := R_0 P_S^{\mathsf{T}} + r^0 e^{n\mathsf{T}}$, die eine zu (19) analoge Gestalt hat.

S3: Bestimme n - l + 1 Givens-Drehungen $\tilde{G}_{l,l+1}, ..., \tilde{G}_{n,n+1}$ so, daß die Subdiagonalelemente $(\mathbf{R}_2)_{j+1,j}$ (j = l, ..., n) sukzessive zu 0 gemacht werden. Dann ist

$$ilde{G} R_2 := ilde{G}_{n,n+1} \cdots ilde{G}_{l,l+1} R_2 =: R_3 =: \left(rac{ extbf{R}}{oldsymbol{o}^{\intercal}}
ight)$$

von oberer Dreiecksform, und $Q_3 := Q_2 \tilde{G}^{\intercal} =: (\bar{Q} \mid \bar{q})$ ist spaltenorthonormal. Schreibe (45) in der Form

$$ar{A} m{P}_S^\intercal = m{Q}_2 m{R}_2 = (m{Q}_2 m{ ilde{G}}^\intercal) (m{ ilde{G}} m{R}_2) = m{Q}_3 m{R}_3 = (m{ar{Q}} \mid m{ar{q}}) \left(rac{m{R}}{m{o}^\intercal}
ight) = m{ar{Q}}m{ar{R}}$$

In analoger Weise kann im Fall 2 vorgegangen werden. Die in 10.3.1(S2) erforderlichen Givens-Drehungen werden dabei auf Kosten der durch P_S bewirkten Spaltenvertauschungen eingespart.

Ü 10.3.2. Man überlege sich, daß im Schritt S2 von 10.3.1 und 10.3.5 die letzte Drehung G_{12} eingespart werden kann. Wie lauten dann die Formeln für R_2 bzw. \hat{R}_2 ?

Ü 10.3.3. Gegeben sei das Quadratmittelproblem $Ax \cong b$. Wenn $\overline{A}\overline{x} = \overline{b}$ durch Streichen der *l*-ten Zeile $(a^{l_{\mathsf{T}}}, b_l)$ von $(A \mid b)$ entsteht und $\overline{\overline{A}}\overline{\overline{x}} = \overline{\overline{b}}$ durch Hinzufügen des i-fachen der *l*-ten Zeile gemäß

$$(\bar{\bar{A}} \mid \bar{\bar{b}}) = \left(\frac{\mathrm{i} a^{l \intercal} \mid \mathrm{i} b_l}{A \mid b} \right)$$
, i imaginäre Einheit mit i² = -1

gebildet wird, gilt für jedes $\boldsymbol{x} \in \mathbf{R}^n$

$$\|\bar{\boldsymbol{b}}-\bar{A}\boldsymbol{x}\|^2=\|\bar{\boldsymbol{b}}-\bar{A}\boldsymbol{x}\|^2.$$

In bezug auf die Residuumnorm hat das Hinzufügen des i-fachen einer Zeile also denselben Effekt wie das Streichen der Zeile. Man kann zeigen, daß die Aufdatierung der QR-Faktorisierung von A bei Übergang zu \overline{A} gemäß 10.3.4 in reeller Arithmetik ausgeführt werden kann.

Ü 10.3.4. Man ermittle die Koeffizienten K_4 und K_5 in (3) für die angegebenen Aufdatierungsalgorithmen unter der Voraussetzung, daß im Fall 1 stets mit Re-Orthogonalisierung gearbeitet wird und im Fall 2 die Matrix \hat{Q}_k nicht berechnet wird.

Bemerkungen zum Kapitel 10

B 10.1. Das Gram-Schmidt-Verfahren gehört zum klassischen Bestand der Analysis und wird seit Jahrzehnten als konstruktives Hilfsmittel zur Beschaffung orthogonaler Basen verwendet. Hinweise auf die unter Umständen extreme Instabilität sind bei RICE [66] zu finden, wo auch die bei gleichem Aufwand wesentlich günstigere modifizierte Version beschrieben ist. Der Nachweis der numerischen Stabilität des MGS-Verfahrens (und nicht der numerischen Gutartigkeit, wie in LAWSON/HANSON [74, S. 130/132] irrtümlich berichtet wird) geht auf

BJÖRCK [67a] zurück und nutzt die spezielle Fehlerkorrelation in Q und R aus. Unter Benutzung des Stewart-Lemmas 8.2.11 ist die numerische Gutartigkeit von KIELBASIŃSKI [unveröffentlichtes Manuskript] gezeigt worden.

B 10.2. Die Verwendung von Householder-Spiegelungen zur Erzeugung von Nullen geht auf HOUSEHOLDER [58] zurück. Ihre Anwendung auf rechteckige Matrizen und zur Lösung von Quadratmittelproblemen wurde von GOLUB [65] und BUSINGER/GOLUB [65] beschrieben. Die Fehleranalyse der **QR**-Faktorisierung mittels Householder-Transformationen wurde von WIL-KINSON [65] durchgeführt.

B 10.3. Im Fall m = n ist die orthogonale Dreiecksfaktorisierung mittels Givens-Drehungen von GIVENS [58] beschrieben worden. Wegen des doppelt so hohen Aufwands hat sie – abgesehen von Sonderfällen wie z. B. schwach besetztem A oder Modifikationstechniken gemäß Abschnitt 10.3 – zunächst nur wenig praktische Bedeutung erlangt. Das änderte sich mit dem Bekanntwerden der impliziten Realisierungsmöglichkeit, siehe B 3.4. Systematische Darstellungen der Nutzung von Givens-Drehungen sind bei GENTLEMAN [73, 75], WILKINSON [77] und VOEVODIN [77] zu finden. Dort können auch die günstigen Kumulationskonstanten aus 10.2.8 nachgelesen werden, siehe auch RATH [82].

B 10.4. Zur Aufdatierung der Spaltennormen bei der Realisierung der Orthogonalisierungsverfahren mit Spaltenvertauschungen gemäß 10.1.6 (iii) bzw. 10.2.5 (vii) sei auf BJÖRCK [67] bzw. auf die LINPACK-Dokumentation von DONGARRA et al. [79] verwiesen.

B 10.5. Grundlegende Ergebnisse zur iterativen Verbesserung von Quadratmittellösungen gehen auf Björck [67b, 68, 78] und Björck/Golub [67] zurück. Dabei ist wesentlich, daß sowohl x als auch r verbessert werden. Die sich auf den ersten Blick anbietende, zu 5.4 analoge Verbesserung allein von x nach der Vorschrift

Bestimme **h** als Lösung von $Ah \cong b - Ax$, setze $\bar{x} := x + h$

ist nur im Fall kleiner Residuen brauchbar, siehe GOLUB [65] und GOLUB/WILKINSON [66].

B 10.6. Über Aufdatierungsalgorithmen von *QR*-Faktorisierungen existiert eine umfangreiche Literatur. Wir zitieren die Arbeiten von LAWSON/HANSON [74], GILL/GOLUB/MURRAY/SAUN-DERS [74], GILL/MURRAY/SAUNDERS [75], DANIEL/GRAGG/KAUFMAN/STEWART [76] und PAIGE [80], wo neben den hier angegebenen Verfahren auch zahlreiche andere gefunden werden können. Für statistische Anwendungen sei auf GOLUB [69], ELDÉN [72], GOLUB/STYAN [73] und LAWSON/HANSON [74] verwiesen. Speziell wird erstmals bei GOLUB [69] das Streichen einer Zeile durch das in bezug auf die Quadratmittellösung äquivalente Hinzufügen des i-fachen der zu streichenden Zeile realisiert (i imaginäre Einheit). Rundungsfehleranalysen für Aufdatierungsalgorithmen sind z. B. bei DANIEL et al. [76] und PAIGE [80] zu finden. Allerdings ist das Problem des Rundungsfehlereinflusses bei Aufdatierungsalgorithmen noch nicht völlig geklärt, und es gibt bis jetzt auch keine umfassende Darstellung des Gesamtkomplexes.

11. Rangdefiziente Quadratmittelprobleme

Im Fall einer rangdefizienten Matrix A hat das Quadratmittelproblem $Ax \simeq b$ unendlich viele Lösungen, deren Gesamtheit eine lineare Mannigfaltigkeit bildet. Die Lösungsmenge kann durch die Normallösung $x = A^+b$ — die Lösung kleinster Euklidischer Norm — und eine orthogonale Basis des Nullraumes von A charakterisiert werden, vgl. 8.1.

Im Abschnitt 8.2 haben wir gesehen, daß die Normallösung eine unstetige und unbeschränkte Funktion der Eingangsdalen $\{A, b\}$ ist, d. h., die Berechnung der

315

Normallösung ist ein inkorrekt gestelltes Problem. Ursache für dieses irreguläre Verhalten sind rangerhöhende Störungen der Matrix A. Zur Überwindung der Inkorrektheit bieten sich daher zwei Möglichkeiten an:

- Ausschließung der bösartigen Störungen bei Beibehaltung der Zielfunktion ||Ax b||, d. h. Einschränkung des Definitionsgebietes der Aufgabe
- Änderung der Zielfunktion bei Beibehaltung des Definitionsgebietes derart, daß die für die ursprüngliche Aufgabe bösartigen Störungen für die modifizierte Aufgabe gutartig sind.

In beiden Fällen wird dabei zu einem modifizierten, korrekt gestelltem Problem übergegangen. Man spricht von *Regularisierung*, und die zugehörigen Algorithmen werden *Regularisierungsverfahren* genannt.

In diesem Kapitel werden wir beide genannten Möglichkeiten zur Regularisierung ausnutzen; die erste führt auf die *diskreten*, die zweite auf die *kontinuierlichen Regularisierungsverfahren*. In beiden Verfahrensklassen hängt der Regularisierungsprozeß und damit die regularisierte Lösung von einem Parameter ab. Die Wahl dieses Parameters kann nicht allein nach numerischen Gesichtspunkten erfolgen, sondern erfordert zusätzliche Informationen über erwünschte Eigenschaften der Lösung. Weil sich solche Eigenschaften nicht immer mathematisch präzise formulieren lassen und der Einfluß des Parameters auf die Lösung a priori nicht eingeschätzt werden kann, ist die Ausführung der Regularisierungsverfahren im Dialog mit dem sachkundigen Anwender oder Bearbeiter einer völlig automatisierten Ausführung i. allg. vorzuziehen.

Da die numerisch berechnete Singulärwertzerlegung — im folgenden mit SVD (engl. "Singular Value Decomposition") abgekürzt — die zuverlässigste Grundlage für die Lösung rangdefizienter Quadratmittelprobleme darstellt und für gewisse diskrete Regularisierungsverfahren explizit benötigt wird, gehen wir zunächst auf die Berechnung dieser Zerlegung ein.

11.1. Numerische Berechnung der Singulärwertzerlegung

Für eine gegebene Matrix $A \in \mathbf{R}^{m,n}$, $m \ge n$, erfolgt die Berechnung der Singulärwertzerlegung

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{\Sigma}\boldsymbol{V}^{\mathsf{T}},\tag{1}$$

 $U \in \mathbf{R}^{m,m}$ und $V \in \mathbf{R}^{n,n}$ orthogonal, $\Sigma = \text{diag}(\sigma_1, \ldots, \sigma_n) \in \mathbf{R}^{m,n}, \sigma_1 \ge \cdots \ge \sigma_n \ge 0$, in zwei Stufen:

In der ersten Stufe wird A durch endlich viele orthogonale Äquivalenztransformationen auf obere Bidiagonalform gebracht:

11.1.1. Bidiagonulisierung einer Matrix mittels Householder-Spiegelungen. Die Matrix $A \in \mathbf{R}^{m,n}$ werde nach der Vorschrift

$$A^{(1)} := A, \qquad A^{(k+1)} := H_k A^{(k)} \hat{H}_k \qquad (k = 1, ..., n)$$
(2)

mit Householder-Spiegelungen $H_k \in \mathbb{R}^{m,m}$, $\hat{H}_k \in \mathbb{R}^{n,n}$ transformiert. Die Spiege-

lung H_k sei so festgelegt, daß die Elemente in der k-ten Spalte von $\bar{A}^{(k)} := H_k A^{(k)}$ unterhalb des Diagonalelementes zu 0 gemacht werden, und \hat{H}_k werde so gewählt, daß die Elemente in der k-ten Zeile von $A^{(k+1)} := \bar{A}^{(k)} \hat{H}_k$ rechts vom Nebendiagonalelement annulliert werden. Dann gilt in exakter Arithmetik

$$U_1^{\mathsf{T}}AV_1 = A^{(n+1)} = \left(\frac{B}{O}\right) \tag{3}$$

mit einer Bidiagonalmatrix

$$\boldsymbol{B} = \begin{pmatrix} \alpha_1 & \gamma_2 & & \\ & \alpha_2 & \gamma_3 & \\ & \ddots & \ddots & \\ & & \ddots & \gamma_n \\ & & & \ddots & \gamma_n \\ & & & & \alpha_n \end{pmatrix} \in \mathbf{R}^{n,n}$$
(4)

und den orthogonalen Transformationsmatrizen

$$U_1^{\mathsf{T}} := H_n \cdots H_2 H_1 \in \mathbb{R}^{m,m}, \qquad H_n = I_m \quad \text{im Fall } m = n,$$

$$V_1 := \hat{H}_1 \hat{H}_2 \cdots \hat{H}_n \in \mathbb{R}^{n,n}, \qquad \hat{H}_{n-1} = \hat{H}_n = I_n.$$
(5)

Aufwand:

 $\sim 2(m-n/3)\,n^2\,{\rm opms}\,+\sim 2n$ opr und $\sim 2n\,{\rm S}$ für ${\pmb B}$ ohne explizite Berechnung von ${\pmb U}_1,{\pmb V}_1$

 $\sim 2n^3/3$ opms und $\sim n^2$ S für explizite Berechnung von V_1

 $\sim 2(m^2 - mn + n^2/3) n$ opms und m^2 S für explizite Berechnung von U_1

Für m = 5, n = 4 verlaufen die ersten Schritte der Transformation nach dem Muster

$$\rightarrow \left(\begin{matrix} \times & \times & 0 & 0 \\ 0 & \times & \times & \boxed{0} \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ 0 & 0 & \times & \times \\ \overline{A}^{(2)} \hat{H}_2 = A^{(3)} \end{matrix} \right) \rightarrow \cdots .$$

Die im jeweiligen Teilschritt neu entstehenden Nullen sind dabei eingerahmt worden.

11.1.2. Bemerkung. (i) Wenn $a_{ik}^{(k)} = 0$ (i = k, ..., m) gilt, ist die Spiegelung H_k nicht eindeutig bestimmt. In diesem Fall soll $H_k := I_m$ gesetzt werden, und im Fall der Unbestimmtheit von \hat{H}_k soll analog $\hat{H}_k := I_n$ gesetzt werden.

(ii) Der Bidiagonalisierungsprozeß kann analog zur **QR**-Faktorisierung in situ unter Verwendung von zusätzlichen $\sim 2n$ S realisiert werden, vgl. 10.2.3. Die Matrizen U_1, V_1 sind dann in der Produktform (5) durch die zu H_k bzw. \hat{H}_k gehörenden Vektoren v^k bzw. \hat{v}^k gegeben. Die Berechnung von U_1f , $U_1^{\mathsf{T}}f$ bzw. V_1g , $V_1^{\mathsf{T}}g$ für $f \in \mathbb{R}^m, g \in \mathbb{R}^n$ ist analog zu 10.2.6 mit $\sim (2m - n) n$ opms bzw. $\sim n^2$ opms möglich.

(iii) Die Aufwandsangaben für die explizite Berechnung von U_1 bzw. V_1 beziehen sich auf die Auswertung gemäß $U_1 := H_1(\cdots(H_{n-1}(H_n))\cdots)$ bzw. $V_1 := \hat{H}_1(\cdots(\hat{H}_{n-2} \times (\hat{H}_{n-1}))\cdots)$, vgl. Ü 10.2.2. Im Fall $m \gg n$ sollte die Berechnung von U_1 wegen des Terms $2m^2n$ tunlichst vermieden werden, während die $\sim 2n^3/3$ opms für V_1 kaum ins Gewicht fallen. Die Teilmatrix der ersten n Spalten von U_1 kann nach der Vorschrift

$$H_1\left(\cdots \left(H_{n-1}\left(H_n\left(rac{I_n}{O}
ight)
ight)
ight)\cdots
ight) ext{ mit}\sim (m-n/3)\ n^2 ext{ opms ermittelt werden}.$$

(iv) Für genügend großes m/n wird der Bidiagonalisierungsaufwand reduziert, wenn A zunächst auf obere Dreiecksform gebracht und die Dreiecksmatrix danach bidiagonalisiert wird, siehe Ü 11.1.1.

(v) Die Spiegelungen H_k bzw. \hat{H}_k können durch jeweils m - k bzw. n - k - 1 explizite oder implizite Givens-Drehungen ersetzt werden.

In der zweiten Stufe des Prozesses zur Berechnung der SVD wird die Bidiagonalmatrix **B** orthogonal äquivalent auf Diagonalform $\tilde{\boldsymbol{\Sigma}} = \boldsymbol{U}_2^{\mathsf{T}} \boldsymbol{B} \boldsymbol{V}_2$ transformiert, d. h., es wird die SVD von **B** bestimmt. Dazu bietet sich an, den Beweis zu 1.2.14 mit **B** statt **A** konstruktiv nachzuvollziehen, also die Zahlen $\lambda_i := \sigma_i^2$ als Eigenwerte der symmetrischen Tridiagonalmatrix

$$\boldsymbol{T} := \boldsymbol{B}^{\mathsf{T}} \boldsymbol{B} = \begin{pmatrix} \alpha_1^2 & \alpha_1 \gamma_2 \\ \alpha_1 \gamma_2 & \alpha_2^2 + \gamma_2^2 & \alpha_2 \gamma_3 \\ \vdots & \vdots & \ddots & \vdots \\ \ddots & \vdots & \ddots & \ddots & \vdots \\ \ddots & \vdots & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \alpha_{n-1} \gamma_n & \alpha_n^2 + \gamma_n^2 \end{pmatrix}$$
(6)

zu bestimmen und U_2 , V_2 mittels der zugehörigen Eigenvektoren festzulegen. Als einen für diese Eigenwertaufgabe besonders geeigneten Algorithmus werden wir in 13.7 den sogenannten **QR**-Algorithmus mit Verschiebungen kennenlernen. Bei Berechnung der λ_i nach diesem oder einem anderen Verfahren ist mindestens mit einem Fehler

$$|\delta\lambda_i| \leq \nu \|\boldsymbol{T}\|_2 = \nu \|\boldsymbol{B}\|_2^2 = \nu \|\boldsymbol{A}\|_2^2 = \nu \sigma_1^2$$
(7)

in der Größenordnung des optimalen Fehlerniveaus von λ_i bezüglich T zu rechnen,

siehe 13.1.2. Für $\sigma_i := \sqrt{\lambda_i}$ ergibt sich damit mindestens ein Fehler

$$|\delta\sigma_i| = \left|\sqrt{\lambda_i + \delta\lambda_i} - \sqrt{\lambda_i}\right| = \frac{|\delta\lambda_i|}{\sqrt{\lambda_i + \delta\lambda_i} + \sqrt{\lambda_i}} \approx \frac{|\delta\lambda_i|}{2\sigma_i} \le \nu \frac{\sigma_1^2}{2\sigma_i}.$$
 (8)

Andererseits ist das optimale Fehlerniveau von σ_i bezüglich **B** nach 8.2.3 durch

$$|\delta\sigma_i| \leq \nu \|\boldsymbol{B}\|_2 = \nu \|\boldsymbol{A}\|_2 = \nu\sigma_1 \tag{9}$$

charakterisiert, also um den Faktor $0.5\sigma_1/\sigma_i$ kleiner als beim Umweg über T. Die kleineren Singulärwerte werden daher nicht mit ausreichender Genauigkeit berechnet; bezüglich dieser Werte liegt keine numerische Stabilität vor. Ein ähnlicher negativer Effekt ist uns bereits beim Übergang vom Quadratmittelproblem $Ax \simeq b$ zu den theoretisch äquivalenten Normalgleichungen $A^{\mathsf{T}}Ax = A^{\mathsf{T}}b$ begegnet, vgl. 9.1.

Im folgenden soll versucht werden, den ungünstigen Einfluß des expliziten Übergangs zu $B^{T}B$ zu vermeiden und B direkt durch eine theoretisch unendliche Folge

$$\boldsymbol{B}^{(1)} := \boldsymbol{B}, \qquad \boldsymbol{B}^{(k+1)} := \boldsymbol{G}_k \boldsymbol{B}^{(k)} \hat{\boldsymbol{G}}_k^{\mathsf{T}} \qquad (k = 1, 2, \ldots)$$
(10)

von orthogonalen Äquivalenztransformationen auf Diagonalform \tilde{B} zu transformieren. Dabei werden G_k , \hat{G}_k als Produkte

$$G_k := G_{n-1,n} \cdots G_{23} G_{12}, \qquad \hat{G}_k^{\mathsf{T}} := \hat{G}_{12}^{\mathsf{T}} \hat{G}_{23}^{\mathsf{T}} \cdots \hat{G}_{n-1,n}^{\mathsf{T}}$$
(11)

von Givens-Drehungen so festgelegt, daß die Matrizen $B^{(k)}$ wie B von Bidiagonalform

$$\boldsymbol{B}^{(k)} = \begin{pmatrix} \alpha_1^{(k)} & \gamma_2^{(k)} & & \\ & \alpha_2^{(k)} & & \\ & \ddots & \ddots & \\ & \ddots & \ddots & \\ & & \ddots & \gamma_n^{(k)} \\ & & & \alpha_n^{(k)} \end{pmatrix}$$
(12)

sind und $T^{(k)} := B^{(k) \top} B^{(k)}$ in exakter Arithmetik gerade die Matrizen wären, die der erwähnte **QR**-Algorithmus aus $T = T^{(1)}$ erzeugt. Man beachte, daß die $T^{(k)}$ dieselbe Struktur wie **T** haben. Der Einfachheit halber lassen wir bei den Drehungen und den Elementen von $B^{(k)}$ bzw. $T^{(k)}$ den oberen Iterationsindex k wie in (11) weg. Wir setzen ferner voraus, daß alle α_i und γ_i von 0 verschieden sind. Dies ist keine Einschränkung, denn andernfalls zerfällt **B** —gegebenenfalls nach einer Zusatztransformation — in Teilmatrizen niedrigerer Dimension, siehe Ü 11.1.3.

Die genannten Forderungen — Bidiagonalität von $B^{(k)}$ und Äquivalenz mit dem *QR*-Algorithmus — lassen sich wie folgt erfüllen: Zunächst wird wie in 13.7 der Verschiebungsparameter $\varrho = \varrho_k$ als derjenige Eigenwert der (2,2)-Teilmatrix

$$\begin{pmatrix} \alpha_{n-1}^2 + \gamma_{n-1}^2 & \alpha_{n-1}\gamma_n \\ \alpha_{n-1}\gamma_n & \alpha_n^2 + \gamma_n^2 \end{pmatrix}$$

in der rechten unteren Ecke von $T^{(k)}$ gewählt, der dem Diagonalelement $\alpha_n^2 + \gamma_n^2$ am nächsten liegt. Danach wird die erste Drehung \hat{G}_{12} so bestimmt, daß das zweite Element $\alpha_1\gamma_2$ der ersten Spalte $(\alpha_1^2 - \varrho, \alpha_1\gamma_2, 0, ..., 0)^{\mathsf{T}}$ von $\mathbf{T}^{(k)} - \varrho \mathbf{I} = \mathbf{B}^{\mathsf{T}}\mathbf{B} - \varrho \mathbf{I}$ bei der Transformation in $\hat{\mathbf{G}}_{12}(\mathbf{B}^{\mathsf{T}}\mathbf{B} - \varrho \mathbf{I})$ zu 0 gemacht wird. Mit dem so festgelegten $\hat{\mathbf{G}}_{12}$ wird $\mathbf{B}^{(k)}\hat{\mathbf{G}}_{12}^{\mathsf{T}}$ gebildet, wobei allerdings die Bidiagonalform in der Position {2, 1} zerstört wird. Alle weiteren Drehungen werden dann in der durch

$$\boldsymbol{B}^{(k+1)} := \boldsymbol{G}_{n-1,n} \Big\{ \cdots \, \boldsymbol{G}_{23} \{ [\boldsymbol{G}_{12}(\boldsymbol{B}^{(k)} \hat{\boldsymbol{G}}_{12}^{\mathsf{T}})] \, \hat{\boldsymbol{G}}_{23}^{\mathsf{T}} \} \cdots \hat{\boldsymbol{G}}_{n-1,n}^{\mathsf{T}} \Big\}$$
(13)

festgelegten Reihenfolge lediglich dazu benutzt, die Bidiagonalform sukzessive wieder herzustellen: Durch Multiplikation von links mit G_{12} wird in $\{2, 1\}$ wieder eine Null erzeugt, allerdings entsteht dabei in $\{1, 3\}$ ein neues Nichtnullelement. Dieses wird durch Rechtsmultiplikation mit \hat{G}_{23}^{T} annulliert, wobei in $\{3, 2\}$ ein von 0 verschiedenes Element entsteht usw. Die Reihenfolge, in der Nichtnullelemente entstehen und im nächsten Teilschritt wieder annulliert werden, ist für n = 5 durch das folgende Schema gegeben:



Wegen der schnellen, praktisch meist kubischen oder besseren Konvergenz des **QR**-Algorithmus konvergieren die Produkte $\alpha_l^{(k)}\gamma_{l+1}^{(k)}$ – das wären die Nebendiagonalelemente von $T^{(k)}$ – sehr schnell gegen 0. Daher muß für nicht zu großes k eines der Elemente $\alpha_l^{(k)}$ bzw. $\gamma_{l+1}^{(k)}$ im Sinne von

$$|\alpha_l^{(k)}| \leq \nu \|\boldsymbol{B}\| \quad \text{bzw.} \quad |\gamma_{l+1}^{(k)}| \leq \nu \|\boldsymbol{B}\| \tag{14}$$

genügend klein sein und kann im Rahmen der Computergenauigkeit als 0 angesehen werden. Dies entspricht einer Störung ∂B der Matrix B mit $||\partial B|| \leq v ||B||$, beeinträchtigt also die Genauigkeit praktisch nicht. Beim Nullsetzen von $\gamma_{l+1}^{(k)}$ zerfällt die Matrix $B^{(k)}$ direkt, beim Nullsetzen von $\alpha_l^{(k)}$ nach einer Zusatztransformation in Blöcke niedrigerer Dimension, siehe wieder Ü 11.1.3. Das Verfahren wird dann mit jedem der Blöcke einzeln und damit billiger fortgesetzt. Meist tritt (14) für l = n - 1ein, siehe 13.7. Falls $|\gamma_n^{(k)}|$ klein ist und 0 gesetzt wird, zerfällt $B^{(k)}$ in den (1, 1)-Block $\alpha_n^{(k)}$, der als bereits feststehendes Diagonalelement angesehen wird, und die durch Streichen der *n*-ten Zeile und Spalte entstehende (n - 1)-dimensionale Restmatrix; man spricht von *Deflation*. Wenn $\alpha_{n-1}^{(k)}$ gleich 0 gesetzt werden kann, ist eine analoge Deflation nach Ausführung der in Ü 11.1.3 beschriebenen Zusatztransformationen sogar zweimal möglich, wobei einmal eine 0 abgespalten wird.

Alle diese Maßnahmen haben zur Folge, daß sich der theoretisch unendliche Diagonalisierungsprozeß (10) praktisch wie ein finiter Prozeß verhält: Nach k Schritten kann

$$\tilde{\boldsymbol{B}} := \boldsymbol{B}^{(k+1)} = \operatorname{diag}\left(\alpha_{\boldsymbol{i}}^{(k+1)}\right) = \boldsymbol{G}_{\boldsymbol{k}} \cdots \boldsymbol{G}_{\boldsymbol{1}} \boldsymbol{B} \hat{\boldsymbol{G}}_{\boldsymbol{1}}^{\mathsf{T}} \cdots \hat{\boldsymbol{G}}_{\boldsymbol{k}}^{\mathsf{T}}$$
(15)

im Rahmen der Computergenauigkeit als Diagonalmatrix angesehen werden. Der

Abbruchindex k ist dabei meist nicht viel größer als n, d. h., pro Diagonalelement ist im Mittel nur eine Schleife (10) erforderlich. In den seltensten Fällen überschreitet k den Wert 2n.

Da die Diagonalelemente von \mathbf{B} i. allg. weder nichtnegativ noch der Größe nach geordnet sind, setzen wir abschließend

$$\tilde{\boldsymbol{\Sigma}} = \operatorname{diag}(\sigma_1, \dots, \sigma_n) := \boldsymbol{P}(\boldsymbol{D}\tilde{\boldsymbol{B}}) \boldsymbol{P}^{\mathsf{T}}$$
(16)

 \mathbf{mit}

$$oldsymbol{D} = ext{diag}\left(d_{i}
ight), \quad d_{i} := \left\{egin{array}{cc} 1 & ext{für} & lpha_{i}^{(k+1)} \geqq 0 \ -1 & ext{für} & lpha_{i}^{(k+1)} < 0 \end{array}
ight.$$

und einer Permutationsmatrix **P**, die $\sigma_1 \ge \cdots \ge \sigma_n$ sichert.

Die gesamte Transformation von B in $\tilde{\Sigma}$ läßt sich dann durch

$$\tilde{\Sigma} = PDG_k \cdots G_1 B\hat{G}_1^{\mathsf{T}} \cdots \hat{G}_k^{\mathsf{T}} P^{\mathsf{T}} = U_2^{\mathsf{T}} BV_2$$
(17)

mit

$$U_2^{\mathsf{T}} := \boldsymbol{P} \boldsymbol{D} \boldsymbol{G}_k \cdots \boldsymbol{G}_1, \qquad \boldsymbol{V}_2^{\mathsf{T}} := \boldsymbol{P} \hat{\boldsymbol{G}}_k \cdots \hat{\boldsymbol{G}}_1 \tag{18}$$

beschreiben. Aus (3) und (17) ergibt sich schließlich

$$A = U_1 \left(\frac{B}{O}\right) V_1^{\mathsf{T}} = U_1 \left(\frac{U_2 \tilde{\Sigma} V_2^{\mathsf{T}}}{O}\right) V_1^{\mathsf{T}} = U_1 \left(\frac{U_2 | O}{O | I}\right) \left(\frac{\tilde{\Sigma}}{O}\right) V_2^{\mathsf{T}} V_1^{\mathsf{T}}$$
(19)

als die gesuchte SVD von A. Zusammenfassend erhalten wir das folgende Resultat.

11.1.3. Berechnung der SVD. Für $A \in \mathbb{R}^{m,n}$ werde der nachfolgende zweistufige Algorithmus ausgeführt:

S1: Transformiere A nach dem Verfahren 11.1.1 gemäß

$$U_1^{\mathsf{T}}AV_1 = \left(\frac{B}{O}\right)$$

auf Bidiagonalform.

S2: Transformiere **B** nach den Formeln (10) bis (18) gemäß

$$U_2^{\mathsf{T}}BV_2 = \tilde{\Sigma}$$

auf Diagonalform $\tilde{\Sigma} = \text{diag}(\sigma_j), \sigma_1 \ge \cdots \ge \sigma_n \ge 0.$ Dann ist die Singulärwet tzerlegung $A = U \tilde{\Sigma} V^{\intercal}$ durch die Faktoren

$$\boldsymbol{U} := \boldsymbol{U}_1 \left(\frac{\boldsymbol{U}_2 \mid \boldsymbol{O}}{\boldsymbol{O} \mid \boldsymbol{I}_{m-n}} \right), \quad \boldsymbol{\Sigma} := \left(\frac{\boldsymbol{\tilde{\Sigma}}}{\boldsymbol{O}} \right), \quad \boldsymbol{V} := \boldsymbol{V}_1 \boldsymbol{V}_2$$
(20)

gegeben.

Aufwand: $\sim [2(m - n/3) n^2 \text{ opms} + 2n^2 \text{ opr}]$ und $\sim 2n \text{ S}$ für Σ ohne Berechnung von U und V; zur Darstellung von U, V siehe Bemerkung 11.1.6.

11.1.4. Bemerkung. (i) Die Berechnung von Σ ist mit zusätzlichen $\sim 2n$ S in situ möglich, wobei U_1, V_1 ohne Mehraufwand in der Produktform (5) anfallen, siehe

11.1.2(ii). Wir setzen generell voraus, daß in der Stufe S2 höchstens 2n QR-Schritte erforderlich sind, und zwar zwei pro Singulärwert. Wegen der Dimensionsreduzierung mittels Deflation sind dann insgesamt $\sim 2[2(n-1)+2(n-2)+\cdots+2] \sim 2n^2$ Givens-Drehungen auszuführen; de facto sind es meist weniger. Die Stufe S2 erfordert daher einen Aufwand von höchstens $\sim 2n^2$ opr $+ O(n^2)$ übrigen Operationen, die gegenüber den $\sim 2(m - n/3) n^2$ opms der Bidiagonalisierungsstufe S1 vernachlässigt werden können, d. h., die Berechnung der Singulärwerte aus der Bidiagonalform kostet praktisch nichts.

(ii) Die Kleinheitsbedingungen (14) werden meist durch raffiniertere Bedingungen wie

$$|\alpha_l^{(k)}| \le \nu(|\gamma_l^{(k)}| + |\gamma_{l+1}^{(k)}|) \quad \text{bzw.} \quad |\gamma_{l+1}^{(k)}| \le \nu(|\alpha_l^{(k)}| + |\alpha_{l+1}^{(k)}|) \tag{21}$$

ersetzt, die bei günstig strukturiertem B die Berechnung der kleinen Singulärwerte mit höherer Genauigkeit erlauben als bei Verwendung von (14).

Die Art der Darstellung der Matrizen U, V hängt von der beabsichtigten Nutzung der SVD ab. Wir erinnern daran, daß die Normallösung des Quadratmittelproblems $Ax \simeq b$ mittels der SVD zumindest in der Theorie wie folgt berechnet werden kann, vgl. 8.1.D:

11.1.5. Basisalgorithmus zur Lösung von $Ax \simeq b$ mittels SVD

S1 (Faktorisierung der Matrix A und Rangbestimmung) S1.1: Berechne Singulärwertzerlegung $A = U\Sigma V^{\intercal}$ S1.2: Bestimme $r = \operatorname{rang}(A)$ aus der Bedingung

$$\sigma_1 \ge \cdots \ge \sigma_r > 0 = \sigma_{r+1} = \cdots = \sigma_n \tag{22}$$

S2 (Berechnung der zur rechten Seite $m{b}$ gehörenden Normallösung $m{x}=A^+m{b}$ $= V \Sigma^+ U^{\mathsf{T}} b$, des Residuums $r = b - Ax = U(I - \Sigma \Sigma^+) U^{\mathsf{T}} b$ sowie dessen

 $= V \Sigma^{+} U^{+} u^{+}$

11.1.6. Bemerkung. (i) Zur Berechnung der Normallösung x werden die ersten r Spalten von V benötigt, und die Spalten r + 1 bis n bilden eine orthonormale Basis des Nullraumes $\mathcal{N}(A)$, der gemäß $\mathcal{L} = \mathcal{L}(A, b) = x + \mathcal{N}(A)$ in die Lösungsmenge \mathcal{L} von $Ax \simeq b$ eingeht. Es ist daher zweckmäßig, die Matrix V explizit zu berechnen: Zunächst wird V_1 wie in 11.1.2 gebildet und danach durch Rechtsmultiplikation mit den Drehungen $\hat{G}_{i,i+1}^{\mathsf{T}}$ aufdatiert. Dies erfordert insgesamt $\sim n^3(2/3 \text{ opms} + 4 \text{ opm})$ +2 ops), vgl. 11.1.4, und die Berechnung von V5 kostet dann $\sim n^2$ opms. Die alternative implizite Produktdarstellung

$$\boldsymbol{V} = \hat{\boldsymbol{H}}_1 \cdots \hat{\boldsymbol{H}}_{n-2} \boldsymbol{V}_2 \tag{23}$$

mit explizit berechnetem V_2 erfordert $\sim n^3(4 \text{ opm} + 2 \text{ ops})$, ist also etwas billiger. Sie ist jedoch unhandlicher, und die Berechnung von $V\xi = \hat{H}_1(\cdots \hat{H}_{n-2}(V_2\xi) \cdots)$ entsprechend 10.2.6 kostet $\sim 2n^2$ opms.

(ii) Die zu (i) analoge explizite Berechnung von U benötigt $\sim [2(m^2 - mn + n^2/3)n$ opms $+ mn^2(4 \text{ opm} + 2 \text{ ops})]$ und m^2 S, scheidet also für $m \gg n$ aus Aufwandsgründen aus. Da zur Festlegung von x jedoch nur die Teilmatrix $U_{1,r} = (u^1, ..., u^r)$ der ersten r Spalten von U gebraucht wird, genügt es, $U_{1,n}$ gemäß

$$\boldsymbol{U}_{1,n} := \boldsymbol{H}_1\left(\cdots\left(\boldsymbol{H}_{n-1}\left(\boldsymbol{H}_n\left(\frac{\boldsymbol{U}_2}{\boldsymbol{O}}\right)\right)\right)\cdots\right)$$
(24)

mit $\sim [(2m-n) n^2 \text{ opms} + n^3(4 \text{ opm} + 2 \text{ ops})]$ und mn S explizit zu berechnen. Das Residuum kann dann gemäß $r := b - U_{1,r}(U_{1,r}^{\mathsf{T}}b)$ bestimmt werden, vgl. (10.1.7).

(iii) Billiger als die explizite Berechnung von $U_{1,n}$ ist die Verwendung der impliziten Produktform

$$\boldsymbol{U} = \boldsymbol{H}_{1}\boldsymbol{H}_{2}\cdots\boldsymbol{H}_{n}\left(\frac{\boldsymbol{U}_{2}\mid\boldsymbol{O}}{\boldsymbol{O}\mid\boldsymbol{I}_{m-n}}\right)$$
(25)

von U mit explizit berechnetem U_2 ; der Aufwand ist $\sim n^3(4 \text{ opm} + 2 \text{ ops}) + n^2 \text{ S}$. Analog zu 10.2.6 kann dann $U^{\intercal}b$ bzw. Uq mit $\sim 2mn$ opms berechnet werden.

(iv) Wenn die rechte Seite **b** a priori bekannt ist und **U** später nicht mehr benötigt wird, kann $\beta = U^{\intercal}b$ parallel zu $\Sigma = U^{\intercal}AV$ bestimmt werden, indem alle von links auf **A** angewendeten Transformationen sofort auch auf **b** angewendet werden. Damit wird die explizite Berechnung von **U**, $U_{1,n}$ oder U_2 vermieden, und der Aufwand zur Bildung von $U^{\intercal}b$ reduziert sich auf $\sim [(2m - n)n \text{ opms} + n^2(4 \text{ opm} + 2 \text{ ops})]$. Das Residuum muß dann gemäß r := b - Ax aus den Originaldaten berechnet werden.

(v) Für die verschiedenen Darstellungen von U sind in der nachfolgenden Tabelle die Aufwandszahlen in opms für die Faktorisierung $A = U\Sigma V^{\intercal}$ und die Berechnung von x und r für eine rechte Seite b zusammengestellt worden. Dabei ist der durch die Givens-Drehungen verursachte Term (4 opm + 2 ops) generell durch 3 opms ersetzt worden, vgl. Ü 3.3.6. Zum Vergleich wurden die entsprechenden Angaben für die **QR**-Faktorisierung nach HOUSEHOLDER gemäß 10.2.4 und 10.2.9 mit angeführt.

Variante	Aufwand für			
	$\{\boldsymbol{U}, \boldsymbol{\Sigma}, \boldsymbol{V}\}$ bzw. $\{\hat{\boldsymbol{Q}}, \hat{\boldsymbol{R}}\}$	<i>x</i>	r	
V, U explizit nach (i), (ii)	$4m^2n + 3mn^2 + 11n^3/3$	$m^2 + n^2$	m^2	
$V, U_{1,n}$ explizit nach (i), (ii)	$4mn^2 + 5n^2$	$mn + n^2$	mn	
V explizit, U Produktform	$2mn^2 + 6n^3$	$2mn + n^2$	2mn	
nach (iii)			1	
V explicit, $U^{T}b$ direkt nach (iv)	$2mn^2+3n^3$	$2mn + 3n^2$	mn	
QR -Faktorisierung $A = \hat{Q}\hat{R}$	$mn^2 - n^3/3$	$2mn - n^2/2$	$2mn - n^2$	

Der Faktorisierungsaufwand für die zu empfehlenden Varianten (iii), (iv) ist im Fall $m \gg n$ also etwa doppelt so groß wie für die Householder-Faktorisierung, vgl. jedoch Ü 11.1.1.

Wir geben abschließend ohne Beweis die Ergebnisse der Rundungsfehleranalyse an.

11.1.7. Rundungsfehleranalyse. Für $A \in \mathbb{R}^{m,n}$ werde Algorithmus 11.1.3 ausgeführt, und es sei k die Anzahl der in S2 benötigten QR-Schritte; i. allg. ist $k \approx n$. Dann gibt es exakt orthogonale Matrizen $U \in \mathbb{R}^{m,m}$, $V \in \mathbb{R}^{n,n}$, so daß die berechnete Diagonalmatrix $\Sigma \in \mathbb{R}^{m,n}$ der Beziehung

$$A + \delta A = U \Sigma V^{\intercal} \quad ext{mit} \quad \| \delta A \|_F \leq
u F \| A \|_F, \qquad F \sim 8mn + 12kn, \quad (26)$$

genügt. Die exakten Orthogonalfaktoren U, V werden durch die gemäß 11.1.6 explizit berechneten Matrizen \tilde{U}, \tilde{V} im folgenden Sinne gut dargestellt: Für alle $\boldsymbol{\xi} \in \mathbb{R}^n, \ \boldsymbol{b} \in \mathbb{R}^m$ gilt für die mittels \tilde{U}, \tilde{V} berechneten Vektoren $\boldsymbol{x} := \operatorname{fl}(\tilde{V}\boldsymbol{\xi}),$ $\boldsymbol{\beta} := \operatorname{fl}(\tilde{U}^{\mathsf{T}}\boldsymbol{b})$

$$\boldsymbol{\beta} = \boldsymbol{U}^{\mathsf{T}}(\boldsymbol{b} + \boldsymbol{\delta}\boldsymbol{b}) \text{ mit } \|\boldsymbol{\delta}\boldsymbol{b}\| \leq vm(4n + 6k) \|\boldsymbol{b}\|$$

und

 $\boldsymbol{x} = \boldsymbol{V}(\boldsymbol{\xi} + \boldsymbol{\delta}\boldsymbol{\xi}) \quad ext{mit} \quad \|\boldsymbol{\delta}\boldsymbol{\xi}\| \leq \nu n(4n + 6k) \|\boldsymbol{\xi}\|.$ (27)

Dasselbe trifft zu, wenn β , x mittels der impliziten Produktdarstellungen (23), (25) berechnet oder wenn U, V durch U^{\intercal} , V^{\intercal} ersetzt wird.

Dies besagt: Die numerische Berechnung der Diagonalmatrix Σ und die Darstellungen der Faktoren U, V in der in 11.1.6 beschriebenen Form sind numerisch gutartige Prozesse mit akzeptablen Kumulationskonstanten.

Übungsaufgaben

Ü 11.1.1. Man zeige: Im Fall m>5n/3 ist die Bidiagonalisierung von $A\in {\sf R}^{m,n}$ nach der folgenden Modifikation von 11.1.1

S1.1: Transformiere A gemäß Algorithmus 10.2.4 nach der Vorschrift

$$Q^{\mathsf{T}}A = \left(\frac{R}{O}\right) \tag{28}$$

mittels Householder-Spiegelungen $Q^{\uparrow} := H_n \cdots H_1$ auf obere Dreiecksform. S1.2: Transformiere $R \in \mathbb{R}^{n,n}$ gemäß 11.1.1 nach der Vorschrift

$$\overline{U}_{1}^{\mathsf{T}} R V_{1} = B \tag{29}$$

mittels Householder-Spiegelungen $\overline{U}_1^{\intercal} := \overline{H}_{n-1} \cdots \overline{H}_2, V_1 := \hat{H}_1 \cdots \hat{H}_{n-2}$ auf Bidiagonal-form **B**

billiger als die direkte Anwendung von 11.1.1 auf A. Dabei gilt

$$U_1^{\mathsf{T}} A V_1 = \left(\frac{B}{O}\right) \quad \text{mit} \quad U_1 := Q \left(\frac{\overline{U}_1 | O}{O | I_{m-n}}\right), \tag{30}$$

und der Aufwand ist $\sim [(m + n) n^2 \text{ opms} + 3n \text{ opr}].$

Ü 11.1.2. Man zeige: Wenn **B** die gemä β (3), (4) zu **A** gehörende Bidiagonalmatrix bezeichnet, gilt

$$||A||_2^2 \leq [\max \{|\alpha_i| + |\gamma_i|: i = 1, ..., n\}] [\max \{|\alpha_i| + |\gamma_{i+1}|: i = 1, ..., n\}],$$

wobei $\gamma_1 = \gamma_n = 0$ zu setzen ist.

/

Hinweis: Man beachte $||\mathbf{A}||_2 = ||\mathbf{B}||_2$ und Ü 1.2.9.

Ü 11.1.3. Es sei B eine Bidiagonalmatrix gemäß (4). Man überlege sich:

(i) Ist $\gamma_{l+1} = 0$ für ein $l \in \{1, ..., n-1\}$, so zerfällt **B** gemäß

$$\boldsymbol{B} = \begin{pmatrix} \alpha_1 & \gamma_2 & & \\ \vdots & \vdots & \vdots & \\ & \ddots & \gamma_{l-1} & 0 \\ & & \alpha_l & \\ \hline & & & \alpha_l & \\ & & & \ddots & \vdots \\ & & & \ddots & \ddots & \\ & & & \ddots & \gamma_n \\ & & & & & \alpha_n \end{pmatrix}$$

in zwei Bidiagonalblöcke der Dimension l bzw. n - l. (ii) Ist $\alpha_l = 0$ für ein $l \in \{1, ..., n-1\}$, so existieren Drehungen $G_{l,l+1}, ..., G_{ln}$, so daß

$$\boldsymbol{B} = \begin{pmatrix} \alpha_{1} & \gamma_{2} & & & \\ \vdots & \vdots & & & \\ \alpha_{l-1} & \gamma_{l} & & \\ 0 & \gamma_{l+1} & & \\ \hline & 0 & \gamma_{l+2} & & \\ \vdots & \vdots & \ddots & \\ & & \ddots & \gamma_{n} \\ & & & & \alpha_{n} \end{pmatrix} \quad \text{in } \boldsymbol{\overline{B}} = \begin{pmatrix} \alpha_{1} & \gamma_{2} & & & \\ \vdots & \ddots & & & \\ \alpha_{l+1} & \gamma_{l} & & \\ 0 & 0 & & \\ \hline & 0 & 0 & & \\ \hline & 0 & 0 & & \\ \hline & & & \alpha_{l+1} & \overline{\gamma}_{l+2} \\ \vdots & & \vdots & \ddots \\ & & & \ddots & \overline{\gamma}_{n} \\ & & & & \overline{\alpha}_{n} \end{pmatrix}$$
(31)

übergeht mit $\overline{B} = G_{ln} \cdots G_{l,l+1}B$.

(iii) Ist $\alpha_n = 0$, so existieren Drehungen $\hat{G}_{n-1,n}, \ldots, \hat{G}_{1n}$, so daß

übergeht mit $\overline{B} = B\hat{G}_{n-1,n}^{\mathsf{T}} \cdots \hat{G}_{1n}^{\mathsf{T}}$.

11.2. **Diskrete Regularisierung**

A. Diskrete Regularisierung mittels SVD

Wir betrachten das Quadratmittelproblem $Ax \simeq b$ mit $A \in \Re^{m,n}$, $b \in \Re^m$. Die Computermatrix A repräsentiere die rangdefiziente exakte Matrix $A^* \in \mathbb{R}^{m,n}$, die den unbekannten Rang $r = \operatorname{rang}(A^*) < n$ und die unbekannten Singulärwerte

$$\sigma_1^* \ge \sigma_2^* \ge \cdots \ge \sigma_r^* > 0 = \sigma_{r+1}^* = \cdots = \sigma_n^* \tag{1}$$
besitzen möge. Dabei wird zugelassen, daß sich A nicht nur durch den Darstellungsfehler δA_D , sondern auch durch einen Meßfehler δA_M von A^* unterscheidet: Es gelte

$$A = A^* + \delta A_D + \delta A_M \quad \text{mit} \quad \|\delta A_D\|_F \le r \|A^*\|_F =: \Delta A_D,$$
$$\|\delta A_M\|_F \le \Delta A_M. \tag{2}$$

Bei der Bestimmung der Normallösung nach dem Basisalgorithmus 11.1.5 wird zunächst die SVD von A gemäß 11.1.3 berechnet. Dann gilt für die berechneten Faktoren $\{U, \Sigma, V\}$ nach 11.1.7

$$A + \delta A_R = U \Sigma V^{\intercal} \quad \text{mit} \quad \|\delta A_R\|_F \leq vF \|A\|_F =: \Delta A_R, \quad F \sim 8mn + 12kn.$$
(3)

Unter Beachtung von (2) ergibt sich daraus

$$A^* + \delta A = U\Sigma V^{\intercal} \tag{4}$$

 \mathbf{m} t

$$\boldsymbol{\delta A} := \boldsymbol{\delta A}_{D} + \boldsymbol{\delta A}_{M} + \boldsymbol{\delta A}_{R}, \quad \boldsymbol{\varepsilon} := \|\boldsymbol{\delta A}\|_{F} \leq \Delta \boldsymbol{A}_{D} + \Delta \boldsymbol{A}_{M} + \Delta \boldsymbol{A}_{R} =: \Delta \boldsymbol{A}.$$
(5)

Die berechnete SVD ist also die exakte Zerlegung einer um δA gestörten Matrix, und das Störungsniveau ε liegt in der Größenordnung des durch Darstellungs- und Meßfehler hervorgerufenen Niveaus. In diesem Sinne ist die berechnete SVD akzeptabel. Insbesondere folgt aus (4), (5) nach 8.2.3

$$|\sigma_i - \sigma_i^*| \le \varepsilon \qquad (i = 1, ..., n), \tag{6}$$

d. h., die exakten Singulärwerte σ_i^* werden ausreichend gut approximiert. Trotzdem ist die berechnete SVD nicht direkt zur Bestimmung von $(A^*)^+$ und $x^* = (A^*)^+ b$ geeignet: Auf Grund der Rundungs- und Meßfehler werden fast immer alle berechneten Singulärwerte

$$\sigma_1 \geqq \sigma_2 \geqq \cdots \geqq \sigma_r \geqq \sigma_{r+1} \geqq \cdots \geqq \sigma_n \tag{7}$$

sogar positiv sein, so daß S 1.2 von 11.1.5 den Rang n und damit einen zu hohen Wert liefert. Für i > r folgt nun aus (1) sofort $\delta \sigma_i = \sigma_i - \sigma_i^* = \sigma_i$, mit (6) also

$$\varepsilon \ge \sigma_{r+1} \ge \cdots \ge \sigma_n. \tag{8}$$

Die berechneten σ_i (i > r) sind demnach kleine positive Zahlen, die wegen des zu groß festgelegten Ranges mit $1/\sigma_i$ in die Pseudoinverse $(A^* + \delta A)^+ = V\Sigma^+ U^{\intercal}$ eingehen, also eine bösartige, rangerhöhende Störung von Σ^* darstellen. Die Bösartigkeit wirkt sich auch auf die Normallösung $\boldsymbol{x} = (A^* + \delta A)^+ \boldsymbol{b} = V\Sigma^+ U^{\intercal}\boldsymbol{b}$ aus, sofern in $\boldsymbol{\beta} = U^{\intercal}\boldsymbol{b}$ nicht alle Komponenten $\beta_{r+1}, \ldots, \beta_n$ verschwinden, was wegen der Rundungsfehler unwahrscheinlich ist; siehe 8.2.B.

Der beschriebene unerwünschte Effekt kann vermieden werden, indem 꾿 durch

Nullsetzen der "kleinen" berechneten σ_i regularisiert wird:

Statt
$$\Sigma = \begin{pmatrix} \sigma_1 & & & \\ \vdots & & & \\ & \sigma_p & & \\ & & \sigma_{p+1} & \\ & & \vdots & \\ & & & \sigma_n \end{pmatrix}$$
 wird $\Sigma_a := \begin{pmatrix} \sigma_1 & & & \\ \vdots & & & \\ & \sigma_p & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & \\ & & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \\$

verwendet, wobei der sog. Pseudorang $p = p(\alpha)$ in Abhängigkeit vom Regularisierungsparameter $\alpha \ge 0$ durch die Bedingung

$$\sigma_1 \ge \cdots \ge \sigma_p > \alpha \ge \sigma_{p+1} \ge \cdots \ge \sigma_n \tag{10}$$

festgelegt wird. Die derart regularisierende Singulärwertzerlegung $U\Sigma_a V^{\intercal}$ ist die exakte SVD der regularisierten Matrix

$$A_{\alpha} := U\Sigma_{\alpha}V^{\intercal} = A^{*} + \delta A_{\alpha}, \qquad (11)$$

durch die eine regularisierte Pseudoinverse $A_{\alpha}^{+} := V \Sigma_{\alpha}^{+} U^{\top}$ und eine regularisierte Normallösung $x_{\alpha} := A_{\alpha}^{+} b$ definiert werden. Mit dieser Modifikation geht 11.1.5 in das folgende regularisierte Verfahren über:

11.2.1. Lösung von $Ax \simeq b$ mittels diskreter Regularisierung der SVD.

S1 (Festlegung der regularisierten SVD):

- S1.1: Bestimme Singulärwertzerlegung $A = U \Sigma V^{\intercal}$ nach 11.1.3.
- S1.2: Wähle Regularisierungsniveau $\alpha \ge 0$, bestimme Pseudorang $p = p(\alpha)$ gemäß (10), definiere A_{α} nach (9), (11).
- S2 (Berechnung der zu b gehörenden regularisierten Normallösung $x_{\alpha} = A_{\alpha}^{\dagger} b$ = $V \Sigma_{\alpha}^{\dagger} U^{\dagger} b$, des Residuums $r_{\alpha} = b - A x_{\alpha} = U(I - \Sigma \Sigma_{\alpha}^{\dagger}) U^{\dagger} b$ und dessen Norm $n r_{\alpha} = ||r_{\alpha}||$):
- S2.1: Berechne $\beta := U^{\mathsf{T}} b$
- S2.2: Setze $\boldsymbol{\varrho}_{\alpha} := (\boldsymbol{I} \boldsymbol{\Sigma}\boldsymbol{\Sigma}_{\alpha}^{+}) \boldsymbol{\beta} = (0, ..., 0, \beta_{p+1}, ..., \beta_{m})^{\mathsf{T}},$ berechne $n\boldsymbol{r}_{\alpha} := \|\boldsymbol{\varrho}_{\alpha}\|$ und $\boldsymbol{r}_{\alpha} := U\boldsymbol{\varrho}_{\alpha}$
- S2.3: Setze $\boldsymbol{\xi}_{\mathfrak{a}} := \boldsymbol{\Sigma}_{\mathfrak{a}}^{+} \boldsymbol{\beta} = (\beta_{1}/\sigma_{1}, ..., \beta_{p}/\sigma_{p}, 0, ..., 0)^{\mathsf{T}}$, berechne $\boldsymbol{x}_{\mathfrak{a}} := V \boldsymbol{\xi}_{\mathfrak{a}}$

Zur Festlegung des Regularisierungsparameters α beachten wir, daß die regularisierte Matrix A_{α} durch die Störung

$$\delta A_{\alpha} := A_{\alpha} - A^{*} = \delta A - U(\Sigma - \Sigma_{\alpha}) V^{\intercal}, \quad \|\delta A_{\alpha}\| \leq \varepsilon + \sigma_{p(\alpha)+1} \leq \varepsilon + \alpha$$
(12)

aus A^* entsteht, siehe (10), (11) und (4), (5). Nun sind die regularisierten Größen A^+_{α}, x_{α} nach 8.2.5 bzw. 8.2.7 lokal lipschitzstetige Funktionen von δA_{α} , wenn die Bedingungen

$$p(\alpha) = \operatorname{rang} \left(\Sigma_{\alpha} \right) = \operatorname{rang} \left(A_{\alpha} \right) = \operatorname{rang} \left(A^* + \delta A_{\alpha} \right) \le \operatorname{rang} \left(A^* \right) = r \quad (13)$$

und

$$\|(A^*)^+\| \| \boldsymbol{\delta} \boldsymbol{A}_{\mathfrak{a}} \| < 1 \tag{14}$$

erfüllt sind; in (13) steht dann sogar das Gleichheitszeichen.

Die Ungleichung (13) ist wegen (8) im Fall

$$\varepsilon \le \alpha$$
 (15)

sicher gültig. Durch die Festlegung (15) werden also rangerhöhende und damit bösartige Störungen δA_{a} ausgeschlossen. Aus (12) folgt andererseits $\|(A^{*})^{+}\| \| \delta A_{a}\|$ $\leq (\varepsilon + \alpha)/\sigma_r^*$ so daß (14) im Fall

$$\alpha < \sigma_r^* - \varepsilon \tag{16}$$

erfüllt ist. Die beiden gegensätzlich wirkenden Bedingungen (15) und (16) lassen sich genau dann gleichzeitig erfüllen, wenn $\varepsilon \leq \alpha < \sigma_r^* - \varepsilon$, also

$$2\varepsilon < \sigma_r^* = 1/\|(A^*)^+\| \tag{17}$$

gilt; der kleinste zulässige a-Wert ist dann

$$\alpha := \alpha_{\text{opt}} := \varepsilon. \tag{18}$$

Für die mit diesem α berechneten regularisierten Größen gelten folgende Aussagen:

11.2.2. Fehleranalyse. Die Matrix $A^* \in \mathbb{R}^{m,n}$ mit rang $(A^*) = r$ und den Singulärwerten (1) werde durch die Computermatrix $A \in \Re^{m,n}$ im Sinne von (4), (5) dargestellt. Für A und $b \in \Re^m$ werde Algorithmus 11.2.1 mit $\alpha := \varepsilon$ durchgeführt. Dann genügen die in S1 berechneten Faktoren der regularisierten SVD der Beziehung

$$A^* + \delta A_{\varepsilon} = A_{\varepsilon} = U \Sigma_{\varepsilon} V^{\intercal} \quad \text{mit} \quad \|\delta A_{\varepsilon}\| \leq 2\varepsilon, \tag{19}$$

und es gilt

$$\operatorname{rang} \left(A_{\varepsilon} \right) = \operatorname{rang} \left(A^{\ast} + \delta A_{\varepsilon} \right) \leq \operatorname{rang} \left(A^{\ast} \right). \tag{20}$$

Zu den in S2 berechneten regularisierten Größen x_{ϵ} , r_{ϵ} existieren Störungen derart, daß

$$\hat{\boldsymbol{x}}_{\varepsilon} := \boldsymbol{x}_{\varepsilon} + \boldsymbol{\delta} \boldsymbol{x}_{\varepsilon} \quad \text{mit} \quad \|\boldsymbol{\delta} \boldsymbol{x}_{\varepsilon}\| \leq \nu n (4n + 6k) \|\boldsymbol{x}_{\varepsilon}\| \tag{21}$$

die Normallösung des gestörten Quadratmittelproblems

$$(A^* + \delta A_{\epsilon}) \hat{\boldsymbol{x}}_{\epsilon} \simeq \boldsymbol{b} + \delta \boldsymbol{b}_{\epsilon} \quad \text{mit} \quad \|\boldsymbol{\delta} \boldsymbol{b}_{\epsilon}\| \leq \nu m (4n + 6k) \|\boldsymbol{b}\|$$
(22)

ist, und es gilt

$$\begin{split} \hat{\boldsymbol{r}}_{\varepsilon} &:= \boldsymbol{r}_{\varepsilon} + \boldsymbol{\delta}\boldsymbol{r}_{\varepsilon} = (\boldsymbol{b} + \boldsymbol{\delta}\boldsymbol{b}_{\varepsilon}) - (\boldsymbol{A}^{*} + \boldsymbol{\delta}\boldsymbol{A}_{\varepsilon}) \, \hat{\boldsymbol{x}}_{\varepsilon} \\ & \text{nit} \quad \|\boldsymbol{\delta}\boldsymbol{r}_{\varepsilon}\| \leq \boldsymbol{v}\boldsymbol{m}(4n + 6k) \, \|\boldsymbol{r}_{\varepsilon}\|. \end{split}$$
(23)
Wenn zusätzlich die Bedingung (17) erfüllt ist, gilt
$$\text{rang} \, (\boldsymbol{A}^{*} + \boldsymbol{\delta}\boldsymbol{A}_{\varepsilon}) = \text{rang} \, (\boldsymbol{A}^{*}), \qquad (24) \\ \text{und} \, \boldsymbol{A}_{\varepsilon}^{-}, \, \hat{\boldsymbol{x}}_{\varepsilon} \, \text{und} \, \hat{\boldsymbol{r}}_{\varepsilon} \, \text{sind lokal lipschitzstetige Funktionen der Störungen} \, \boldsymbol{\delta}\boldsymbol{A}_{\varepsilon}, \, \boldsymbol{\delta}\boldsymbol{b}_{\varepsilon}. \end{split}$$

Beweis. Die zentralen Aussagen über A_{ε}^{-} und die Rangbeziehungen wurden bereits hergeleitet; die restlichen ergeben sich aus den Formeln für x_{ε} und r_{ε} unter Beachtung von (11.1.27). \Box

11.2.2 besagt: Die regularisierte, durch $\{U, \Sigma_{\epsilon}, V\}$ gegebene SVD ist die exakte SVD der regularisierten Matrix A_{ϵ} , die sich von A^* durch eine Störung δA_{ϵ} in der Größenordnung des Fehlerniveaus ϵ unterscheidet. Sie stellt daher eine akzeptable SVD von A^* dar. Die regularisierte Normallösung x_{ϵ} und das zugehörige Residuum r_{ϵ} sind die leicht gestörte Normallösung und das leicht gestörte Residuum eines benachbarten Quadratmittelproblems. Die Störung δA_{ϵ} ist dabei nach Konstruktion nicht rangerhöhend, und bei Gültigkeit von (17) liegt sogar lokal lipschitzstetige Abhängigkeit vor.

11.2.3. Bemerkung. (i) Wenn $b \in \Re^m$ der durch Darstellungs- und Meßfehler verfälschte Repräsentant der unbekannten rechten Seite $b^* \in \mathbb{R}^m$ ist, braucht nur

$$oldsymbol{b} = oldsymbol{b}^{st} + \delta oldsymbol{b}, \hspace{0.2cm} \delta oldsymbol{b}_{D} + \delta oldsymbol{b}_{M}, \hspace{0.2cm} \| \delta oldsymbol{b}_{D} \| \leq
u \, \| oldsymbol{b}_{M} \|, \hspace{0.2cm} \| \delta oldsymbol{b}_{M} \| \leq \Delta oldsymbol{b}_{M}$$

gesetzt zu werden.

(ii) Bei Erfülltsein von (17), d. h. im Fall $\varkappa := 2\varepsilon ||(A^*)^+|| < 1$, folgt aus (19) und 8.2.5

$$\frac{\|(\boldsymbol{A}^{*})^{\scriptscriptstyle +} - \boldsymbol{A}_{\varepsilon}^{\scriptscriptstyle +}\|}{\|(\boldsymbol{A}^{*})^{\scriptscriptstyle +}\|} \leq \varepsilon \, \frac{\left(1 + \sqrt{5}\right) \, \|(\boldsymbol{A}^{*})^{\scriptscriptstyle +}\|}{1 - \varkappa},\tag{25}$$

und für $x^* := (A^*)^+ b$ sowie die exakt aus $\{U, \Sigma_{\epsilon}, V\}$ berechnete regularisierte Lösung $A_{\epsilon}^+ b$ gilt nach 8.2.7

$$\frac{\|\boldsymbol{x}^{*} - \boldsymbol{A}_{\varepsilon}^{+}\boldsymbol{b}\|}{\|\boldsymbol{x}^{*}\|} \leq \varepsilon \frac{2 \|(\boldsymbol{A}^{*})^{+}\|}{1 - \varkappa} \left[\sqrt{2} + \frac{\|(\boldsymbol{A}^{*})^{+}\| \|\boldsymbol{v}^{*}\|}{\|\boldsymbol{x}^{*}\|} \right]$$
(26)

mit $r^* = b - A^*x^*$. Unter Verwendung von (21), (22) kann $x^* - x_{\varepsilon}$ in analoger Weise abgeschätzt werden. Zur Schranke aus (26) kommen dabei noch Terme proportional zu ν hinzu.

(iii) Es scheint, als ob das Verfahren 11.2.1 mit $\alpha = \varepsilon$ rangdefiziente Probleme stets befriedigend löst. Das ist leider nicht immer der Fall: Zum einen ist das Fehlerniveau ε meist nicht genau bekannt. Bestenfalls kennt man Schranken die eine

$$\varepsilon_{\min} \leq \varepsilon \leq \varepsilon_{\max}$$
,

Unsicherheit der Rangbestimmung in den Grenzen

$$p_{\min} := p(\varepsilon_{\max}) \leq p(\varepsilon) \leq p(\varepsilon_{\min}) =: p_{\max}$$

nach sich ziehen. Zum anderen ist das gutartige, lokal lipschitzstetige Verhalten der regularisierten Größen an die Bedingung (17) gebunden, d. h., es liegt nur vor, wenn das Fehlerniveau deutlich kleiner als σ_r^* ist. Da r nicht bekannt ist, läßt sich diese Bedingung praktisch nicht überprüfen. Wenn (17) nicht erfüllt ist, brauchen sich (13) und (14) für kein $\alpha > 0$ gleichzeitig erfüllen zu lassen, d. h., *jeder* Regularisierungsversuch liefert entweder eine rangerhöhende oder eine zu große Störung. Eine solche Konfliktsituation kann nur durch Verwendung einer genügend hohen Genauigkeit vermieden werden, sofern keine Meßfehler auftreten. Im Fall $\Delta A_M = 0$ ist nämlich $\varepsilon \leq \Delta A_D + \Delta A_R \leq \nu(1+F) ||A^*||_F$, so daß (17) für genügend kleines ν erfüllt ist. Bei Abwesenheit von Meßfehlern sollte dabei $\varepsilon_{\min} \approx \nu ||A|| = \nu \sigma_1$ als untere Schranke für den ε -Bereich gewählt werden. Bei signifikanten Meßfehlern braucht (17) auch bei noch so hoher Computergenauigkeit nicht erfüllt zu sein.

Die obigen Ausführungen legen es nahe, den Regularisierungsparameter $\alpha \ge 0$ als variabel anzusehen und alle zugehörigen regularisierten Lösungen x_a als mögliche Kandidaten für die gesuchte Lösung zu betrachten. Da $p(\alpha)$ nur die Werte $\{0, 1, \ldots, n\}$ annehmen kann, gibt es i. allg. n + 1 verschiedene x_a , deshalb auch der Name *diskrete Regularisierung*. Um den Einfluß von α auf die regularisierten Größen x_a , r_a zu untersuchen, gehen wir von (4) aus und nehmen der Einfachheit halber an, daß $\{x_\alpha, r_a\}$ fehlerfrei aus $\{b, U, \Sigma, V\}$ berechnet werden. Die in Gleitpunktarithmetik berechneten Größen unterscheiden sich von diesen nur unwesentlich. Mit den Bezeichnungen von 11.2.1 gilt dann

$$\|\boldsymbol{x}_{\alpha}\|^{2} = \|\boldsymbol{\xi}_{\alpha}\|^{2} = \sum_{j=1}^{p(\alpha)} (\beta_{j}/\sigma_{j})^{2} \leq \sum_{j=1}^{\operatorname{rang}(\boldsymbol{A})} (\beta_{j}/\sigma_{j})^{2} = \|\boldsymbol{x}^{N}\|^{2}$$
(27)

bzw.

$$\|\boldsymbol{r}_{a}\|^{2} = \|\boldsymbol{\varrho}_{a}\|^{2} = \sum_{i=p(a)+1}^{m} (\beta_{i})^{2} \ge \sum_{i=\mathrm{rang}(\hat{\boldsymbol{A}})+1}^{m} (\beta_{i})^{2} = \|\boldsymbol{r}^{N}\|^{2}$$
(28)



mit $\hat{A} := A^* + \delta A$, $x^N := x_0 := A^+ b$, $r^N := r_0 := b - Ax^N$. Die Normen $||x_a||$ und $||r_a||$ sind also stückweise konstante Funktionen, wobei $||x_a||$ monoton fällt, $||r_a||$ dagegen monoton wächst, siehe Abb. 11.2.1 und 11.2.2.

Um ein realistischeres Bild zu erhalten, müßte man sich die α -, $||\boldsymbol{x}_{\alpha}||$ - und $||\boldsymbol{r}_{\alpha}||$ -Achse logarithmisch geteilt denken. Im Fall $\Delta A_M = 0$, $\nu \approx 10^{-6}$ ist eine typische Verteilung der berechneten Singulärwerte durch die Tabelle

i	1	2		r	r+1	r+2		n
σ_i	15.783	13.021		0.00352	$0.75 imes10^{-4}$	$0.29 imes 10^{-4}$		$0.83 imes 10^{-5}$

gegeben, wobe
i $\sigma_{r+1}, \ldots, \sigma_n$ in der Größenordnung von $\nu ||\mathcal{A}|| \approx \nu \sigma_1 = 1.5 \times 10^{-5}$ liegen und σ_r deutlich größer ist, als
op = rleicht erkannt werden kann. Mit fallende
m α , d. h. mit wachsendem p, fällt
 $||\boldsymbol{r}_a||$, während $||\boldsymbol{x}_a||$ i. allg. nicht zu stark wächst. Wen
n α in die Größenordnung von ε kommt, wächst
 $||\boldsymbol{x}_a||$ i. allg. sehr stark an, meist um mehrere Zehnerpotenzen, während $||\boldsymbol{r}_a||$ sich entsprechend der Zahlenwerte der
 β_i weiter — aber i. allg. nicht drastisch — verkleinert, siehe (27), (28). An Hand des Diagramms 11.2.2 kann der Bearbeiter dann denjenigen Pseudorang $p(\alpha)$ wählen, dessen zugehöriges
 \boldsymbol{x}_a und \boldsymbol{r}_a seinen Vorstellungen über die erwartete Lösung am besten entspricht. Geeignete Auswahlkriterien wären etwa

$$\|\boldsymbol{x}_{a}\| \to \text{Minimum!} \quad \text{bei} \quad \|\boldsymbol{r}_{a}\| \leq \tau$$
(29)

oder

$$\|\boldsymbol{r}_{\boldsymbol{\alpha}}\| \to \text{Minimum!} \quad \text{bei} \quad \|\boldsymbol{x}_{\boldsymbol{\alpha}}\| \leq \omega$$
(30)

mit vorgegebenen, aus der realen Aufgabenstellung abzuleitenden Schranken $\tau, \omega > 0$. In (29) muß dabei $\tau \ge ||\mathbf{r}^N||$ gefordert werden, damit überhaupt zulässige Lösungen existieren, vgl. (28). Zur Lösung von (30) wird p — von 0 beginnend — solange erhöht, daß die Nebenbedingung gerade noch erfüllt ist. Für (29) wird p solange erhöht, bis die Nebenbedingung erstmals erfüllt ist. Selbstverständlich erfolgt die Steuerung wie beschrieben direkt mit p und nicht mit α .

B. Diskrete Regularisierung mittels anderer Faktorisierungen

Um die relativ aufwendige Berechnung der SVD zu vermeiden, können auch andere Faktorisierungen zur diskreten Regularisierung herangezogen werden. Als Beispiel betrachten wir im folgenden die QR-Faktorisierung mittels Householder-Spiegelungen, vgl. 10.2. Wegen der erwarteten Rangdefizienz kann dabei nicht spaltenweise wie im Vollrangfall vorgegangen werden, sondern es sind Spaltenvertauschungen gemäß 10.2.5(vii) erforderlich. Wenn das Verfahren mit $A \in \Re^{m,n}$ gestartet wird, erhalten wir nach dem k-ten Schritt

$$A^{(k+1)} = \left(\begin{array}{c} R^{(k+1)} \\ O \\ R^{(k+1)} \\ k + 1 \end{array} \right) = H_k \cdots H_1 A T_{1,\hat{s}(k)} \cdots T_{k,\hat{s}(k)}$$

mit der Restmatrix

 $M^{(k+1)} = (m^{k+1,k+1}, m^{k+1,k+2}, ..., m^{k+1,n}).$

Im Fall $k = r = \operatorname{rang} (A^*)$ müßte bei fehlerfreiem $A = A^*$ und in exakter Arithmetik $M^{(k+1)} = O$ gelten. Bei Anwesenheit von Rundungs- und Meßfehlern wird jedoch auch in diesem Fall fast immer $M^{(k+1)} \neq O$ sein, so daß der nächste Schritt ausgeführt werden kann. Allerdings ist zu erwarten, daß $M^{(k+1)}$ klein ist. Dann wird auch das nächste Diagonalelement klein sein, denn wegen der Spaltenvertauschungsstrategie gilt

$$|r_{k+1,k+1}| = \max\{\|\boldsymbol{m}^{k+1,j}\|: j = k+1, ..., n\},$$
(31)

vgl. (10.2.21). Sollte doch $M^{(k+1)} = O$ sein, so wird einfach $H_j := I_m$, $T_{j,\hat{s}(j)} := I_n$ (j = k + 1, ..., n) gesetzt. Mit dieser Festlegung erhalten wir auch im rangdefizienten Fall stets eine Faktorisierung von A, und nach 10.2.5 und 10.2.7 gilt

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}_{R}) \boldsymbol{P}^{\mathsf{T}} = \boldsymbol{Q}\boldsymbol{R} \quad \text{mit} \quad \|\boldsymbol{\delta}\boldsymbol{A}_{R}\| \leq \nu F \|\boldsymbol{A}\| =: \Delta \boldsymbol{A}_{R}, \qquad F \leq 4mn^{3/2}.$$
(32)

Dabei beschreibt $P^{\intercal} := T_{1,\hat{s}(1)} \cdots T_{n-1,\hat{s}(n-1)}$ die Spaltenvertauschungen, $Q := H_1 \cdots H_n$ ist eine exakt orthogonale Matrix und R die berechnete obere Dreiecksmatrix mit

$$|r_{11}| \geq |r_{22}| \geq \cdots \geq |r_{nn}|.$$

Der in 10.2.7 geforderte Vollrang und die Bedingung (10.2.25) werden dort nur gebraucht, um die Durchführbarkeit mit $r_{kk} \neq 0$ zu sichern, was bei der obigen Modifikation überflüssig wird. Wenn zusätzlich wieder Darstellungs- und Meßfehler analog zu (2) zugelassen werden, ergibt sich schließlich

$$(A^* + \delta A) P^{\mathsf{T}} = QR \quad \text{mit} \quad \varepsilon := \|\delta A\| \leq \Delta A_D + \Delta A_M + \Delta A_R =: \Delta A.$$
(33)

Da fast immer $r_{kk} \neq 0$ (k = 1, ..., n), also rang $(\mathbf{R}) = n$ sein wird, ist $\mathbf{\delta}\mathbf{A}$ i. allg. eine rangerhöhende bösartige Störung. Zum Zweck der Regularisierung wird dann analog zu (10) der Pseudorang $p = p(\alpha)$ durch

$$|r_{11}| \ge \cdots \ge |r_{pp}| > \alpha \ge |r_{p+1,p+1}| \ge \cdots \ge |r_{nn}|$$
(34)

festgelegt. Wegen (31) ist dies gleichwertig mit der Definition von p als kleinstem Index aus $\{0, 1, ..., n\}$, für den

$$\|\boldsymbol{m}^{\mathbf{k}+1,j}\| \leq \alpha \qquad (j=p+1,...,n)$$
 (35)

gilt, d. h., eine im Sinne von (35) spaltenweise kleine Restmatrix $M^{(k+1)}$ wird als vernachlässigbar angesehen. Die Regularisierung besteht dann im Ersatz von

$$\boldsymbol{R} = \left(\frac{\boldsymbol{R}_{11} \mid \boldsymbol{R}_{12}}{O \mid \boldsymbol{R}_{22}}\right)_{m+1}^{p} \quad \text{durch} \quad \boldsymbol{R}_{a} := \left(\frac{\boldsymbol{R}_{11} \mid \boldsymbol{R}_{12}}{O \mid O}\right) = \left(\frac{\boldsymbol{R}_{1}}{O}\right)_{m}^{p} \quad (36)$$

 mit

$$\boldsymbol{R}_{1} := (\boldsymbol{R}_{11} \mid \boldsymbol{R}_{12}) := \begin{pmatrix} r_{11} \dots r_{1p} & r_{1,p+1} \dots r_{1n} \\ \vdots & \vdots \\ 0 & r_{pp} & r_{p,p+1} \dots r_{pn} \end{pmatrix}.$$
(37)

Die derart regularisierte QR-Faktorisierung ist die exakte Faktorisierung der regularisierten Matrix

$$A_{\alpha} := QR_{\alpha}P = A^* + \delta A_{\alpha}, \qquad (38)$$

wobei wegen (33)

$$\delta A_{\alpha} := A_{\alpha} - A^* = \delta A - Q(R - R_{\alpha}) P, \quad \|\delta A_{\alpha}\| \leq \varepsilon + \sqrt{n - p} \alpha \quad (39)$$

gilt. Man beachte, daß aus (35)

$$\|\boldsymbol{R} - \boldsymbol{R}_{a}\|_{F} = \|\boldsymbol{R}_{22}\|_{F} = \|\boldsymbol{M}^{(p+1)}\|_{F} \leq \sqrt{n-p} \ x$$

folgt.

Aus der Faktorisierung (38) ergibt sich nach 8.1.9

$$A^{+}_{\alpha} = P^{\mathsf{T}} R^{+}_{\alpha} Q^{\mathsf{T}} = P^{\mathsf{T}} (R^{+}_{1} \mid O) Q^{\mathsf{T}}.$$

$$\tag{40}$$

Zur Darstellung von \mathbf{R}_1^+ transformieren wir \mathbf{R}_1 gemäß

$$\boldsymbol{R}_{1}\boldsymbol{\bar{Q}}^{\mathsf{T}} = (\boldsymbol{R}_{11} \mid \boldsymbol{R}_{12}) \, \boldsymbol{\bar{H}}_{p} \cdots \boldsymbol{\bar{H}}_{2} \boldsymbol{\bar{H}}_{1} = (\boldsymbol{\bar{R}}_{11} \mid \boldsymbol{O}) =: \boldsymbol{\bar{R}}_{1}$$
(41)

mittels p Householder-Spiegelungen $\overline{H}_k = I - \overline{v}^k (\overline{v}^k)^{\intercal} / \overline{\gamma}_k$ in \overline{R}_1 , wobei \overline{R}_{11} eine obere Dreiecksmatrix sein soll. Dies ist möglich, wenn \overline{v}^k in der Form

 $\overline{v}^{k} = (0, ..., 0, \overline{v}_{kk}, 0, ..., 0, \overline{v}_{p+1,k}, ..., \overline{v}_{nk})^{\mathsf{T}}$ (k = p, ..., 1)

angesetzt wird und damit Nullen in den Positionen (k, p + 1) bis (k, n) von R_1 erzeugt werden, vgl. Ü 3.3.7. Aus (41) folgt dann wieder nach 8.1.9

$$\boldsymbol{R}_{1}^{+} = (\boldsymbol{\bar{R}}_{1}\boldsymbol{\bar{Q}})^{+} = \boldsymbol{\bar{Q}}^{\mathsf{T}}\boldsymbol{\bar{R}}_{1}^{-} = \boldsymbol{\bar{Q}}^{\mathsf{T}} \left(\frac{\boldsymbol{\bar{R}}_{11}^{-1}}{\boldsymbol{O}} \right), \tag{42}$$

mit (40) also

$$A_{\alpha}^{\perp} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{\bar{Q}}^{\mathsf{T}} \left(\frac{\boldsymbol{\bar{R}}_{11}^{-1} \mid \boldsymbol{O}}{\boldsymbol{O} \mid \boldsymbol{O}} \right) \boldsymbol{Q}^{\mathsf{T}}.$$

$$\tag{43}$$

Die Produktdarstellung (43) ist die Grundlage für das folgende Verfahren:

- 11.2.4. Lösung von $Ax \simeq b$ mittels diskreter Regularisierung der QR-Faktorisierung· S1 (Festlegung der regularisierten QR-Faktorisierung):
- S1.1: Bestimme Faktorisierung $AP^{\intercal} = QR$ gemäß 10.2.5(vii) mit der oben beschriebenen Modifikation im Fall $M^{(k+1)} = O$
- S1.2: Wähle Regularisierungsniveau $\alpha \ge 0$, bestimme Pseudorang $p = p(\alpha)$ gemäß (34) und lege R_{α} nach (36) fest.

- S1.3: Transformiere $(\mathbf{R}_{11} \mid \mathbf{R}_{12})$ gemäß (41) in $(\mathbf{\overline{R}}_{11} \mid \mathbf{O}) = \mathbf{R}_1 \mathbf{\overline{Q}}^{\mathsf{T}}$.
- S2 (Berechnung der regularisierten Normallösung $x_{\alpha} = A_{\alpha}^{+}b$, des Residuums $r_{\alpha} = b - Ax_{\alpha}$ und dessen Norm $nr_{\alpha} = ||r_{\alpha}||$: S2.1: Berechne

$$oldsymbol{c} = \left(rac{oldsymbol{c}^1}{oldsymbol{c}^2}
ight) := oldsymbol{Q}^{\intercal}oldsymbol{b}\,, \qquad oldsymbol{c}^1 \in \, oldsymbol{\mathsf{R}}^p\,, \quad oldsymbol{c}^2 \in \,oldsymbol{\mathsf{R}}^{m-p}$$

S2.2: Bestimme z^1 aus dem Dreieckssystem $\overline{R}_{11}z^1 = c^1$, setze $z := \left(\frac{z^1}{o}\right)$ mit $o \in \mathbf{R}^{n-p}$

S2.3: Berechne
$$\boldsymbol{y} := \left(rac{\boldsymbol{y}^1}{\boldsymbol{y}^2}
ight) := ar{\boldsymbol{Q}}^\intercal \boldsymbol{z} = ar{\boldsymbol{H}}_p \cdots ar{\boldsymbol{H}}_1 \boldsymbol{z}, ext{ setze } \boldsymbol{x}_a := \boldsymbol{P}^\intercal \boldsymbol{y}$$

S2.4: Berechne $d^2 := c^2 - R_{22}y^2$, setze $nr_a := ||d^2||$ und $d := \left(\frac{o}{d^2}\right) \in \mathbb{R}^m$, berechne $r_a := Od$

Der Ausdruck für r_{a} ergibt sich dabei aus $r_{a} = b - Ax_{a} = Q(c - Ry)$ unter Beachtung von

$$m{R}m{y} = igg(rac{m{R}_{11}m{y}^1 + m{R}_{12}m{y}^2}{m{R}_{22}m{y}^2} igg) = igg(rac{m{c}^1}{m{R}_{22}m{y}^2} igg).$$

Wir fragen jetzt wieder, wann die für das gutartige Verhalten von δA_{α} zu fordernden Bedingungen

$$p(\alpha) = \operatorname{rang} \left(\boldsymbol{R}_{\alpha} \right) = \operatorname{rang} \left(\boldsymbol{A}_{\alpha} \right) = \operatorname{rang} \left(\boldsymbol{A}^{\boldsymbol{*}} + \boldsymbol{\delta} \boldsymbol{A}_{\alpha} \right) \leq \operatorname{rang} \left(\boldsymbol{A}^{\boldsymbol{*}} \right) = r$$
 (44)

und

$$\|(A^*)^+\| \| \delta A_a \| < 1 \tag{45}$$

erfüllt sind. Zur Überprüfung von (44) betrachten wir die ersten p Spalten der Gleichung (33). Wenn die entsprechenden Teilmatrizen mit dem unteren Index 0 gekennzeichnet werden, ergibt sich

$$(A^*P^{\mathsf{T}})_0 = (QR)_0 - (\delta A \cdot P^{\mathsf{T}})_0 = QR_0 - (\delta A \cdot P^{\mathsf{T}})_0.$$

Da \mathbf{R}_0 nach Konstruktion Vollrang p hat, folgt im Fall

$$\|(\boldsymbol{Q}\boldsymbol{R}_0)^+\|\|(\boldsymbol{\delta}\boldsymbol{A}\cdot\boldsymbol{P}^{\mathsf{T}})_0\|<1\tag{46}$$

aus 8.2.5 $p = \operatorname{rang} (\boldsymbol{QR}_0) = \operatorname{rang} ((\boldsymbol{A^*P^\intercal})_0)$, wegen $\operatorname{rang} ((\boldsymbol{A^*P^\intercal})_0) \leq \operatorname{rang} (\boldsymbol{A^*P^\intercal})$ = rang (A*) also (44). Nun gelten

$$\|(QR_0)^+\| = \|R_0^+\| = \left\|\left(\frac{R_{11}}{O}\right)^+\right\| = \|(R_{11}^{-1} \mid O)\| = \|R_{11}^{-1}\|$$

und $\|(\mathbf{\delta} A \cdot \mathbf{P}^{\mathsf{T}})_{\mathbf{0}}\| \leq \|\mathbf{\delta} A \cdot \mathbf{P}^{\mathsf{T}}\| = \|\mathbf{\delta} A\| = \varepsilon$, so daß (46) und damit (44) sicher erfüllt sind, wenn $\boldsymbol{R}_{11} = \boldsymbol{R}_{11}(\alpha)$ der Bedingung

$$\|\boldsymbol{R}_{11}(\boldsymbol{\alpha})^{-1}\|\,\varepsilon < 1 \tag{47}$$

genügt. Zur Abschätzung von $||\mathbf{R}_{11}^{-1}||$ beachten wir den folgenden Satz: Für jede reguläre und im Sinne von $||\mathbf{R}_{11}e^{j}||^2 = \sum_{i=1}^{j} r_{ij}^2 = 1$ (j = 1, ..., p) spaltennormierte obere Dreiecksmatrix $\mathbf{R}_{11} \in \mathbf{R}^{p,p}$, deren Elemente der Ungleichung (10.2.23) genügen und die daher durch **QR**-Faktorisierung mit Spaltenvertauschungen aus einer spaltennormierten Matrix \mathbf{A} entstanden sein kann, gilt

$$|r_{pp}| \ge \sigma_p(\boldsymbol{R}_{11}) = 1/\|\boldsymbol{R}_{11}^{-1}\| \ge |r_{pp}|/2^{p-1},$$
(48)

siehe B 11.2. Die Abschätzung nach unten ist dabei leider fast scharf, siehe Ü 11.2.2. Da (48) auch für nicht spaltennormiertes R_{11} gilt, folgt mit (34)

$$\|oldsymbol{R}_{11}^{-1}\|\,arepsilon\leq 2^{p-1}arepsilon/|r_{pp}|<2^{p-1}arepsilon/lpha$$
 ,

und die für die Gültigkeit von (44) hinreichende Bedingung (47) ist erfüllt, wenn

$$2^{p-1}\varepsilon \le \alpha \tag{49}$$

gilt. Wahrend bei Verwendung der SVD bereits durch die Festlegung $\varepsilon \leq \alpha$ eine Rangerhöhung ausgeschlossen wird, tritt hier zusätzlich der Faktor 2^{p-1} auf, der bei hohem Pseudorang $p = p(\alpha)$ ein gegenüber ε wesentlich vergrößertes Regularisierungsniveau bedingt. Daß ein solches ungünstiges Verhalten tatsächlich vorliegen kann, zeigt das folgende Beispiel.

11.2.5. Beispiel. Es sei $\mathbf{R}_{11} := \mathbf{R}(n, c, s)$ die Matrix aus Ü 11.2.2 für n := 15 und $c := 1 - s^2$, $s := 0.01^{1/14}$. Bei dieser Festlegung ergibt sich

$$1 = |r_{11}| \ge \cdots \ge |r_{15,15}| = 0.01$$
 .

Für jedes $\alpha < 0.01$ – also für relativ große α – ist $p(\alpha) = 15$, d. h., R_{11} wird als regulär angesehen. Nach Ü 11.2.2 gilt jedoch

$$\sigma_{15}(\boldsymbol{R_{11}}) \leq |r_{nn}|/K(n,c) = 0.01/815 = 1.2 \times 10^{-5},$$

so daß R_{11} bereits im Rahmen des Fehlerniveaus $\varepsilon \approx 10^{-5}$ singulär ist. \Box

Wir wenden uns jetzt der zweiten Bedingung (45) zu. Sie ist wegen (39) sicher erfüllt, falls $(\varepsilon + \sqrt{n-p} \alpha)/\sigma_r^* < 1$, also die (16) entsprechende Ungleichung

$$\alpha < (\sigma_r^* - \varepsilon) / \sqrt{n - p} \tag{50}$$

gilt. Die beiden für das gutartige Verhalten von A hinreichenden Bedingungen (49), (50) sind daher genau dann gleichzeitig erfüllbar, wenn

$$\left(1 + \sqrt{n-r} \ 2^{r-1}\right)\varepsilon \leq \sigma_r^* = 1/||(\boldsymbol{A^*})^-||$$
(51)

gilt.

11.2.6. Bemerkung. (i) Die Bedingung (51) ist für großes r wesentlich einschränkender als die bei Verwendung der SVD auftretende analoge Bedingung (17). Wenn sie nicht erfüllt ist, kann es spaltennormierte Matrizen A geben, so daß die hinreichenden Bedingungen (44) und (45) durch kein α befriedigt werden können, d. h., jeder Regularisierungsversuch würde entweder eine rangerhöhende oder eine

zu große Störung einführen. Es liegt in diesem Fall eine unlösbare Konfliktsituation in der Wahl von α vor. Die numerische Rangbestimmung und Regularisierung auf der Grundlage der **QR**-Faktorisierung mit Spaltenvertauschungen ist daher mit einer wesentlich höheren Unsicherheit verbunden als die mittels SVD.

(ii) Selbstverständlich braucht die beschriebene Konfliktsituation nicht immer aufzutreten. Für viele praktisch vorkommende Matrizen werden die Bedingungen (44), (45) für $\alpha \approx \varepsilon$ erfüllt sein. In diesem Fall gelten für die berechneten Größen A_a und x_a zu 11.2.2 analoge Aussagen, jedoch ist r_a i. allg. nicht das leicht gestörte Residuum eines leicht gestörten Nachbarproblems. Überdies gilt dann

$$\frac{1}{\sqrt{n-p+1}} \frac{|r_{11}|}{|r_{pp}|} \leq \frac{\max\{|\bar{r}_{jj}|: j = 1, ..., p\}}{\min\{|\bar{r}_{jj}|: j = 1, ..., p\}} \leq \operatorname{cond}(\boldsymbol{R}_{1}) = \operatorname{cond}(\boldsymbol{A}_{i}) \approx \operatorname{cond}(\boldsymbol{A^{*}}),$$
(52)

man vgl. 6.1.5(iii) und beachte die Bildung von \overline{R}_1 . Die unteren Schranken können als Konditionsschätzer verwendet werden.

(iii) Wenn alle zu $p \in \{0, ..., n\}$ gehörenden regularisierten Lösungen gesucht sind, sollten diese in der Reihenfolge p := n(-1)0 berechnet werden. Zur Realisierung von (41) sind dann Givens-Drehungen günstiger, weil beim Übergang von p zu p-1 nur die ersten p-1 Elemente der p-ten Spalte von $\mathbf{R}_1 = \mathbf{R}_1^{(p)}$ durch Rechtsmultiplikation mit $\mathbf{G}_{p-1,p}^{\mathsf{T}}, ..., \mathbf{G}_{1p}^{\mathsf{T}}$ annulliert zu werden brauchen.

(iv) Eine zu 11.2.4 analoge Regularisierung ist auf der Grundlage der MGS-Faktorisierung mit Spaltenvertauschungen gemäß 10.1.6 (iii) möglich.

(v) Da die
$$(n, n)$$
-Dreiecksmatrix $\tilde{R} = \left(\frac{R_{11} \mid R_{12}}{O \mid R_{22}}\right)$ im Faktor $R = \left(\frac{\tilde{R}}{O}\right)$ von

 $AP^{\mathsf{T}} = QR$ in exakter Arithmetik bis auf das Vorzeichen der Spalten mit dem Faktor L der mit Diagonalpivotisierung gemäß 6.1.5(iii) berechneten Cholesky-Faktorisierung $P(A^{\mathsf{T}}A) P^{\mathsf{T}} = LL^{\mathsf{T}}$ übereinstimmt, kann 11.2.4 äquivalent als diskrete Regularisierung der Normalgleichungen formuliert werden, siehe Ü 11.2.3. Die Ungleichung (10.2.23) geht dann in (6.1.19) über. Bei Computerrechnung treten allerdings die durch die explizite Bildung von $M = A^{\mathsf{T}}A$ hervorgerufenen Probleme auf, siehe Kapitel 9.

(vi) Wenn $r = \operatorname{rang} (A^*)$ bekannt ist, braucht für die Gutartigkeit von δA_{α} nur noch (45) bzw. (50) gefordert zu werden. Für $\alpha := \varepsilon$ ist letztere Forderung kaum einschränkender als (17), so daß die Bedenken gegen eine Verwendung der QR-Faktorisierung in diesem Fall entfallen. \Box

Übungsaufgaben

Ü 11.2.1. Die Matrix $A \in \mathbb{R}^{m,n}$ besitze die exakte Singulärwertzerlegung $A = U\Sigma V^{\intercal}$. Die Zahl p = p(x, A) sei gemäß (10) festgelegt, und $A_{\alpha} = A_{\alpha}(x, A)$ sei die durch (9), (11) definierte regularisierte Matrix. Man zeige, daß

$$p(\alpha, A) = \min \{ \operatorname{rang} (A + \delta A) : ||\delta A|| \leq \alpha \}$$

337

gilt und $\boldsymbol{B} = \boldsymbol{A}_{a}$ die Lösung von

$$\min \{ ||A - B||_F : \operatorname{rang} (B) \leq p(\alpha, A) \}$$

ist.

Ü 11.2.2. Für c, s > 0 werde die Matrix

betrachtet. Man zeige:

(i) Es gilt $\boldsymbol{u} = \boldsymbol{R}\boldsymbol{v}$ für $\boldsymbol{u} := s^{n-1}\boldsymbol{e}^n$ und

$$m{v}:=(c(1+c)^{m{n}-2},\,c(1+c)^{m{n}-3},\,...,\,c(1+c),\,c,\,1)^{\intercal}.$$

folglich

$$\sigma_n(\mathbf{R}) = 1/||\mathbf{R}^{-1}|| \le ||\mathbf{u}||/||\mathbf{v}|| =: |r_{nn}|/K(n, c)$$

 $_{\rm mit}$

$$[K(n, c)]^{2} = 1 + \frac{c^{2}[(1 + c)^{2n-2} - 1]}{(1 + c)^{2} - 1}.$$

(ii) Im Fall c := s := 1 sind alle Diagonalelemente gleich 1, aber es gilt

 $\sigma_n(\boldsymbol{R})/|r_{nn}| \leq 1/K(n, 1) \leq \sqrt{3}/2^{n-1}.$

Die Elemente von R erfüllen jedoch die Bedingung (10.2.23) nicht, d. h., R kann nicht durch QR-Faktorisierung mit Spaltenvertauschungen entstanden sein.

(iii) Bei der Festlegung 0 < c < 1, $s := \sqrt{1 - c^2}$ ist (10.2.23) dagegen erfüllt, und **R** ist spaltennormiert. Außerdem gilt

$$\sigma_n(\mathbf{R})/|r_{nn}| \leq 1/K(n, c) \xrightarrow[c \to 1-0]{} 1/K(n, 1) \leq \sqrt{3}/2^{n-1}.$$

Ü 11.2.3. Es sei $PMP^{\intercal} = LL^{\intercal}$ die in exakter Arithmetik nach 6.1.5 (iii) berechnete Cholesky-Faktorisierung von $M := A^{\intercal}A$; im Fall max $\{a_{jj}^{(k+1)}: k+1 \leq j \leq n\} = 0$ werde dabei L_{k+1} $= \cdots L_{n-1} = I$ gesetzt (bei Computerrechnung ist der Test max $\{a_{jj}^{(k+1)}: k+1 \leq j \leq n\} \leq 0$ zu verwenden und zusätzlich $M^{(k+1)} := O$ zu setzen). Man zeige:

(i) Wenn $p = p(\alpha)$ aus

$$l_{11} \geq \cdots \geq l_{pp} > \alpha \geq l_{p+1,p+1} \geq \cdots \geq l_{nn}$$

bestimmt und

$$\boldsymbol{L} = \begin{pmatrix} \boldsymbol{L}_{11} \mid \boldsymbol{O} \\ \boldsymbol{L}_{21} \mid \boldsymbol{L}_{22} \end{pmatrix} \quad \text{zu} \quad \boldsymbol{L}_{\mathfrak{a}} := \begin{pmatrix} \boldsymbol{L}_{11} \mid \boldsymbol{O} \\ \boldsymbol{L}_{21} \mid \boldsymbol{O} \end{pmatrix} \quad \text{mit} \quad \boldsymbol{L}_{11} \in \mathsf{R}^{p,p}$$

22 Schwetlick, Numerische Algebra

regularisiert wird, gilt

$$(\boldsymbol{A}_{\boldsymbol{\alpha}}\boldsymbol{P}^{\mathsf{T}})^{\mathsf{T}} (\boldsymbol{A}_{\boldsymbol{\alpha}}\boldsymbol{P}^{\mathsf{T}}) = \boldsymbol{P}\boldsymbol{A}_{\boldsymbol{\alpha}}^{\mathsf{T}}\boldsymbol{A}_{\boldsymbol{\alpha}}\boldsymbol{P}^{\mathsf{T}} = \boldsymbol{L}_{\boldsymbol{\alpha}}\boldsymbol{L}_{\boldsymbol{\alpha}}^{\mathsf{T}},$$

wobei $A_{\alpha}P^{\intercal} = (A_1(\alpha) \mid A_2(\alpha))$ aus $AP^{\intercal} = (A_1 \mid A_2)$ nach der Vorschrift $A_1(\alpha) := A_1, A_2(\alpha)$ $:= A_1 L_{11}^{-\intercal} L_{21}^{\intercal}$ gebildet wird.

(ii) Die Normalgleichungen

$$(A_{\mathfrak{a}}P^{\mathsf{T}})^{\mathsf{T}} (A_{\mathfrak{a}}P^{\mathsf{T}}) y = (A_{\mathfrak{a}}P^{\mathsf{T}})^{\mathsf{T}} b, \qquad y := \left(\frac{y^{1}}{y^{2}}\right) := Px_{\mathfrak{a}}$$

des regularisierten Quadratmittelproblems $A_{\alpha}x_{\alpha}\cong b$ reduzieren sich auf

$$oldsymbol{L}_{11}oldsymbol{L}_{11}^{\mathsf{T}}oldsymbol{y}^1 + oldsymbol{L}_{11}oldsymbol{L}_{21}^{\mathsf{T}}oldsymbol{y}^2 = A_1^{\mathsf{T}}oldsymbol{b} =: oldsymbol{c}^1$$

und haben die Lösungsmenge

$$\boldsymbol{L} = \boldsymbol{L}(\boldsymbol{A}_{a}, \boldsymbol{b}) = \{ \boldsymbol{x}_{a} = \boldsymbol{P}^{\mathsf{T}} \boldsymbol{y} \colon \boldsymbol{y}^{1} = \boldsymbol{L}_{11}^{-\mathsf{T}} \boldsymbol{L}_{11}^{-1} \boldsymbol{c}^{1} - \boldsymbol{L}_{11}^{-\mathsf{T}} \boldsymbol{L}_{21}^{\mathsf{T}} \boldsymbol{y}^{2}, \boldsymbol{y}^{2} \in \mathsf{R}^{n-p} \text{ beliebig} \}$$

(iii) Die Normallösung ist durch dasjenige y^2 gegeben, das

$$\|\boldsymbol{y}\|^{2} = \|\boldsymbol{y}^{1}\|^{2} + \|\boldsymbol{y}^{2}\|^{2} = \left\| \left[\frac{\boldsymbol{L}_{11}^{-\mathsf{T}}\boldsymbol{L}_{11}^{-1}\mathbf{c}^{1}}{\boldsymbol{o}} \right] - \left[\frac{\boldsymbol{L}_{11}^{-\mathsf{T}}\boldsymbol{L}_{21}^{\mathsf{T}}}{\boldsymbol{I}} \right] \boldsymbol{y}^{2} \right\|^{2}$$
(53)

minimiert und das sich als eindeutige Lösung der zu (53) gehörenden Normalgleichungen

$$(\boldsymbol{L}_{21}\boldsymbol{L}_{11}^{-1}) (\boldsymbol{L}_{21}\boldsymbol{L}_{11}^{-1})^{\mathsf{T}} + \boldsymbol{I} \} \boldsymbol{y}^{2} = \boldsymbol{L}_{21}\boldsymbol{L}_{11}^{-1}\boldsymbol{L}_{11}^{-\mathsf{T}}\boldsymbol{L}_{11}^{-\mathsf{T}}\boldsymbol{L}_{11}^{-\mathsf{T}}\boldsymbol{c}^{\mathsf{T}}$$

ergibt.

ł

11.3. Kontinuierliche Regularisierung

Wie bisher sei $A^* \in \mathbb{R}^{m,n}$ eine rangdefiziente Matrix mit $r = \text{rang}(A^*) < n$, die im Sinne von (11.2.2) durch $A \in \mathfrak{N}^{m,n}$ repräsentiert wird. Die kontinuierliche Regularisierung als eine zweite Regularisierungstechnik besteht dann darin, die Originalaufgabe

 $\|\boldsymbol{b} - A\boldsymbol{x}\|^2 \rightarrow \text{Minimum!},$

 $\boldsymbol{b} \in \Re^{\boldsymbol{m}}$, durch die regularisierte Aufgabe

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^{2} + \alpha^{2} \|\boldsymbol{x}\|^{2} = \left\| \left(\frac{\boldsymbol{b}}{\boldsymbol{o}} \right) - \left(\frac{\boldsymbol{A}}{\alpha \boldsymbol{I}} \right) \boldsymbol{x} \right\|^{2} \to \text{Minimum!}$$
(1)

zu ersetzen, d. h., statt $Ax \simeq b$ wird das Problem

$$A_{\mathfrak{a}}x \simeq b_{\mathfrak{a}} \quad \text{mit} \quad A_{\mathfrak{a}} := \left(\frac{A}{\alpha I}\right) \in \mathbb{R}^{m+n,n}, \qquad b_{\mathfrak{a}} := \left(\frac{b}{o}\right) \in \mathbb{R}^{m+n}$$
 (2)

gelöst. Die zugehörigen Normalgleichungen $A_a^{\mathsf{T}} A_a x = A_a^{\mathsf{T}} b_a$ haben die Form

$$\boldsymbol{M}_{\boldsymbol{\alpha}}\boldsymbol{x} := (\boldsymbol{A}^{\mathsf{T}}\boldsymbol{A} + \boldsymbol{\alpha}^{2}\boldsymbol{I})\,\boldsymbol{x} = \boldsymbol{A}^{\mathsf{T}}\boldsymbol{b}\,. \tag{3}$$

Aus der Gültigkeit von

$$oldsymbol{x}^\intercal M_{lpha}oldsymbol{x} = \|oldsymbol{A}_{lpha}oldsymbol{x}\|^2 = \|oldsymbol{A}oldsymbol{x}\|^2 + lpha^2 \|oldsymbol{x}\|^2 \ge lpha^2 \|oldsymbol{x}\|^2$$
 für alle $oldsymbol{x} \in oldsymbol{R}^n$

folgt

$$\sqrt{\lambda_i(M_{\alpha})} = \sigma_i(A_{\alpha}) \ge \alpha \qquad (i = 1, ..., n).$$
(4)

Insbesondere ist A_{α} im Fall $\alpha > 0$ spaltenregulär, so daß A_{α}^{+} durch

$$A_{\alpha}^{+} = (A_{\alpha}^{\mathsf{T}} A_{\alpha})^{-1} A_{\alpha}^{\mathsf{T}} = M_{\alpha}^{-1} (A^{\mathsf{T}} \mid \alpha I) =: (\mathbf{S}_{\alpha} \mid T_{\alpha})$$
⁽⁵⁾

gegeben ist und die äquivalenten Probleme (1) bis (3) die eindeutige Lösung

$$\boldsymbol{x}_{\alpha} = \boldsymbol{A}_{\alpha}^{+} \boldsymbol{b}_{\alpha} = \boldsymbol{S}_{\alpha} \boldsymbol{b} + \boldsymbol{T}_{\alpha} \boldsymbol{o} = (\boldsymbol{A}^{\mathsf{T}} \boldsymbol{A} + \alpha^{2} \boldsymbol{I})^{-1} \boldsymbol{A}^{\mathsf{T}} \boldsymbol{b}$$
(6)

besitzen.

A. Kontinuierliche Regularisierung mittels SVD

Die Berechnung und Analyse der regularisierten Größen S_{α} und x_{α} sowie des zugehörigen Residuums $r_{\alpha} = b - Ax_{\alpha}$ ist am leichtesten unter Verwendung der SVD möglich. Der Einfachheit halber ignorieren wir dabei zunächst die bei Computerrechnung auftretenden Rundungsfehler und betrachten nur den Einfluß der Regularisierung. Es sei dazu

$$A = U\Sigma V^{\intercal}, \qquad \Sigma = \text{diag} (\sigma_i), \qquad \sigma_1 \ge \cdots \ge \sigma_p > 0 = \sigma_{p+1} = \cdots = \sigma_n$$
(7)

die exakte SVD einer beliebigen Matrix vom Rang p, die gemäß $A = A^* + \delta A$ durch eine Störung δA aus A^* entsteht und nicht notwendig die oben eingeführte Computerdarstellung von A^* zu sein braucht. Später werden wir insbesondere $\delta A = o$ setzen, so daß (7) die exakte SVD von A^* ist, und x_a wird dann die exakte, zu A^* gehörende regularisierte Lösung x_a^* .

Einsetzen von (7) in (5) liefert für die regularisierte Pseudoinverse

$$S_{\alpha} = V[(\Sigma^{\mathsf{T}}\Sigma + \alpha^{2}I)^{-1}\Sigma^{\mathsf{T}}] U^{\mathsf{T}} =: V\Omega_{\alpha}U^{\mathsf{T}}$$
(8)

mit

$$\boldsymbol{\Omega}_{\alpha} = \operatorname{diag}\left(\frac{\sigma_{1}}{\sigma_{1}^{2}+\alpha^{2}}, \dots, \frac{\sigma_{p}}{\sigma_{p}^{2}+\alpha^{2}}, 0, \dots, 0\right) \in \mathbb{R}^{n,m}.$$
(9)

Unter Verwendung dieser Darstellung ergibt sich der folgende Algorithmus zur Berechnung von x_{α} , vgl. 11.2.1:

- 11.3.1. Lösung von $Ax \simeq b$ mittels kontinuierlicher Regularisierung der SVD.
- S1 (Festlegung der regularisierten Pseudoinversen S_{α}):
- S1.1: Berechne Singulärwertzerlegung $A = U \Sigma V^{\intercal}$
- S1.2: Wähle Regularisierungsniveau $\alpha > 0$, definiere S_{α} , Ω_{α} gemäß (8), (9).
- S2 (Berechnung der zu \boldsymbol{b} gehörenden regularisierten Lösung \boldsymbol{x}_{a} , des Residuums $\boldsymbol{r}_{a} = \boldsymbol{b} \boldsymbol{A}\boldsymbol{x}_{a}$ und der Normen $n\boldsymbol{x}_{a} = \|\boldsymbol{x}_{a}\|$ sowie $n\boldsymbol{r}_{a} = \|\boldsymbol{r}_{a}\|$):
- S2.1: Berechne $\beta := U^{\intercal}b$

S2.2: Bestimme
$$\boldsymbol{\varrho}_{\alpha} := (\boldsymbol{I} - \boldsymbol{\Sigma}\boldsymbol{\Omega}_{\alpha}) \boldsymbol{\beta} = \left(\frac{\alpha^{2}\beta_{1}}{\sigma_{1}^{2} + \alpha^{2}}, ..., \frac{\alpha^{2}\beta_{p}}{\sigma_{p}^{2} + \alpha^{2}}, \beta_{p+1}, ..., \beta_{m}\right)^{\prime}$$
,
berechne $n\boldsymbol{r}_{\alpha} := \|\boldsymbol{\varrho}_{\alpha}\|$ und $\boldsymbol{r}_{\alpha} := \boldsymbol{U}\boldsymbol{\varrho}_{\alpha}$

S2.3: Bestimme
$$\boldsymbol{\xi}_{\alpha} := \boldsymbol{\Omega}_{\alpha} \boldsymbol{\beta} = \left(\frac{\sigma_1 \beta_1}{\sigma_1^2 + \alpha^2}, ..., \frac{\sigma_p \beta_p}{\sigma_p^2 + \alpha^2}, 0, ..., 0 \right)^{\mathsf{T}}$$
,
berechne $n \boldsymbol{x}_{\alpha} := \|\boldsymbol{\xi}_{\alpha}\|$ und $\boldsymbol{x}_{\alpha} := V \boldsymbol{\xi}_{\alpha}$

Aus (7) bis (9) folgt

$$A^{+} - S_{a} = V(\Sigma^{+} - \Omega_{a}) U^{\intercal} = V \operatorname{diag} \left(\frac{\alpha^{2}}{\sigma_{1}(\sigma_{1}^{2} + \alpha^{2})}, ..., \frac{\alpha^{2}}{\sigma_{p}(\sigma_{p}^{2} + \alpha^{2})}, 0, ..., 0 \right) U^{\intercal},$$

wegen $||A^+|| = 1/\sigma_p$ also

$$\|A^{+} - S_{\alpha}\| = \|\Sigma^{+} - \Omega_{\alpha}\| = \frac{\alpha^{2}}{\sigma_{p}(\sigma_{p}^{2} + \alpha^{2})} \leq \frac{\alpha^{2}}{\sigma_{p}^{3}} = \alpha^{2} \|A^{+}\|^{3}.$$
(10)

Die Abschätzung von x_{a} gegen die Normallösung $x^{N} = A^{+}b = V\Sigma^{+}\beta$ führt auf

$$\|\boldsymbol{x}^N - \boldsymbol{x}_{a}\|^2 = \|(\boldsymbol{\Sigma}^+ - \boldsymbol{\Omega}_{a}) \ \boldsymbol{\beta}\|^2 = \sum_{j=1}^p \left[\frac{\alpha^2 \beta_j}{\sigma_j (\sigma_j^2 + \alpha^2)}
ight]^2 \leq \left[\frac{\alpha^2}{\sigma_p^2}
ight]^2 \sum_{j=1}^p \left(\frac{\beta_j}{\sigma_j}
ight)^2,$$

d. h.

$$|\boldsymbol{x}^{N} - \boldsymbol{x}_{a}|| \leq \alpha^{2} \, \|\boldsymbol{A}^{+}\|^{2} \, \|\boldsymbol{x}^{N}\|.$$
(11)

Für den Fehler von r_{a} gegenüber $r^{N} := b - Ax^{N}$ ergibt sich analog

$$\|m{r}^N-m{r}_{\mathfrak{a}}\|^2=\|m{\Sigma}(m{\Sigma}^+-m{\Omega}_{\mathfrak{a}})\,m{eta}\|^2=\sum\limits_{j=1}^p\left[rac{lpha^2m{eta}_j}{\sigma_j^2+lpha^2}
ight]^2\leq\left[rac{lpha^2}{\sigma_p}
ight]^2\sum\limits_{j=1}^p\left(rac{m{eta}_j}{\sigma_j}
ight)^2,$$

mithin

$$\|\boldsymbol{r}^{N}-\boldsymbol{r}_{\boldsymbol{\alpha}}\| \leq \alpha^{2} \|\boldsymbol{A}^{+}\| \|\boldsymbol{x}^{N}\|.$$
(12)

Aus den Abschätzungen (10) bis (12) folgt für $\alpha \rightarrow 0$ die Gültigkeit von

$$S_{\mathfrak{a}} o A^+, \qquad x_{\mathfrak{a}} o x^N \quad ext{und} \quad r_{\mathfrak{a}} o r^N,$$

d. h., die exakt berechneten regularisierten Größen approximieren Pseudoinverse, Normallösung und deren Residuum beliebig genau. Insbesondere können S_{α} , x_{α} und r_{α} auch für $\alpha = 0$ durch die entsprechenden Grenzwerte definiert werden. Man beachte ferner, daß der durch die Störung $\left(\frac{O}{\alpha I}\right)$ von $\left(\frac{A}{O}\right)$ bedingte Block

$$T_{\alpha} = M_{\alpha}^{-1} \alpha I = V \operatorname{diag} \left(\frac{\alpha}{\sigma_1^2 + \alpha^2}, ..., \frac{\alpha}{\sigma_p^2 + \alpha^2}, \frac{1}{\alpha}, ..., \frac{1}{\alpha} \right) V^{\mathsf{T}}$$

von A_x^- im Fall p < n für $\alpha \to 0$ unbeschränkt wächst, sich also bösartig auswirkt. Diese Bösartigkeit überträgt sich jedoch nicht auf x_{α} , da der entsprechende Block von b_{α} nach Konstruktion o ist, vgl. (6).

Für die Normen $\|\boldsymbol{x}_{a}\| = \|\boldsymbol{\xi}_{a}\|$ bzw. $\|\boldsymbol{r}_{a}\| = \|\boldsymbol{\varrho}_{a}\|$ gilt

$$\|\boldsymbol{x}_{a}\|^{2} = \sum_{j=1}^{p} \left[\frac{\sigma_{j}\beta_{j}}{\sigma_{j}^{2} + \alpha^{2}}\right]^{2} \leq \sum_{j=1}^{p} \left(\frac{\beta_{j}}{\sigma_{j}}\right)^{2} = \|\boldsymbol{x}^{N}\|^{2}$$
(13)

bzw.

$$\|\boldsymbol{r}_{\alpha}\|^{2} = \sum_{i=1}^{p} \left[\frac{\alpha^{2} \beta_{i}}{\sigma_{i}^{2} + \alpha^{2}} \right]^{2} + \sum_{i=p+1}^{m} (\beta_{i})^{2} \ge \sum_{i=p+1}^{m} (\beta_{i})^{2} = \|\boldsymbol{r}^{N}\|^{2},$$
(14)

und ihr Verhalten als Funktion von α wird durch die folgende Aussage charakterisiert:

11.3.2. Aussage. Die Normen $||\boldsymbol{x}_{\alpha}||$ bzw. $||\boldsymbol{r}_{\alpha}||$ sind für $0 \leq \alpha < \infty$ monoton fallende bzw. wachsende Funktionen von α mit

$$\lim_{lpha
ightarrow+0} \|oldsymbol{x}_{lpha}\| = \|oldsymbol{x}^{N}\| \geqq \|oldsymbol{x}_{lpha}\| \geqq 0 = \lim_{lpha
ightarrow\infty} \|oldsymbol{x}_{lpha}\|$$

bzw.

$$\lim_{\alpha \to +0} \|\boldsymbol{r}_{\alpha}\| = \|\boldsymbol{r}^{N}\| \leq \|\boldsymbol{r}_{\alpha}\| \leq \|\boldsymbol{b}\| = \lim_{\alpha \to \infty} \|\boldsymbol{r}_{\alpha}\|.$$

Im Fall $x^N = A^+ b \neq o$ liegt strenge Monotonie vor, andernfalls sind beide Normen konstant.

Beweis. Differentiation von (13) bzw. (14) liefert

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \|\boldsymbol{x}_{\alpha}\|^{2} = -4\alpha \sum_{j=1}^{p} \frac{\sigma_{j}^{2} \beta_{j}^{2}}{(\sigma_{j}^{2} + \alpha^{2})^{3}} \leq 0$$
(15)

bzw.

$$\frac{\mathrm{d}}{\mathrm{d}\alpha} \|\boldsymbol{r}_{\alpha}\|^{2} = 4\alpha^{3} \sum_{j=1}^{p} \frac{\sigma_{j}^{2}\beta_{j}^{2}}{(\sigma_{j}^{2} + \alpha^{2})^{3}} \ge 0, \qquad (16)$$

und im Fall $x^{N} \neq o$, d. h. $(\beta_{1}, ..., \beta_{p}) \neq (0, ..., 0)$, gilt strenge Ungleichheit. \Box

In Abb. 11.3.1 ist ein mögliches Bild der Funktionen $||x_{\alpha}||, ||r_{\alpha}||$ skizziert worden.



Abb. 11.3.1. $||\boldsymbol{x}_{\alpha}||$ und $||\boldsymbol{r}_{\alpha}||$ als Funktionen von α

Wie im Fall der diskreten Regularisierung bietet es sich an, alle $\{x_{\alpha}: 0 \leq \alpha\}$ als mögliche Kandidaten für eine regularisierte Lösung anzusehen und unter diesen einen auf Grund der konkreten Problemstellung besonders geeigneten auszuwählen. Da die Kandidaten eine einparametrige stetige Schar bilden — deshalb auch der Name kontinuierliche Regularisierung —, wird man sich praktisch auf endlich viele Repräsentanten beschränken müssen, etwa auf die Berechnung von x_{α} für $\alpha = 10^{l_{p}} ||A||$ = $10^{l_{p}}\sigma_{1}$, l = 0, 1, ..., so daß $10^{l_{p}} \leq 1$ ist. Eine andere Möglichkeit besteht wieder in der Festlegung von α analog zu (11.2.29) bzw. (11.2.30). Der folgende Satz zeigt, daß die im Sinne von (11.2.29) bzw. (11.2.30) optimalen x_{α} die Nebenbedingungen mit Gleichheit erfüllen und gewisse zusätzliche Ext remaleigenschaften besitzen.

11.3.3. Satz. Es sei $x^N \neq o$. Dann gilt:

(i) Zu jedem $\omega \min 0 < \omega \leq ||x^N||$ gibt es gen au ein $\alpha = \alpha_1(\omega) \geq 0$ mit $||x_x|| = \omega$, und x_{α} ist eindeutige Lösung der Aufgabe

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\| \to \text{Minimum!} \quad \text{bei} \quad \|\boldsymbol{x}\| \leq \omega.$$
 (17)

- (ii) Zu jedem τ mit $||\mathbf{r}^{N}|| \leq \tau < ||\mathbf{b}||$ gibt es genau ein $\alpha = \alpha_{2}(\tau) \geq 0$ mit $||\mathbf{r}_{\alpha}|| = \tau$, und \mathbf{x}_{α} ist eindeutige Lösung von
 - $\|\boldsymbol{x}\| \to \text{Minimum!}$ bei $\|\boldsymbol{b} \boldsymbol{A}\boldsymbol{x}\| \leq \tau.$ (18)

Beweis. Wegen 11.3.2 sind $||\boldsymbol{x}_{\alpha}||$ und $||\boldsymbol{r}_{\alpha}||$ streng monotone Funktionen mit den Wertebereichen (0, $||\boldsymbol{x}^{N}||$] und $[||\boldsymbol{r}^{N}||, ||\boldsymbol{b}||)$, so daß α in beiden Fällen eindeutig festgelegt ist.

Zu (i): Es sei $||\boldsymbol{x}_{\alpha}|| = \omega$. Da \boldsymbol{x}_{α} die Aufgabe (1) löst, gilt $||\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}||^2 + \alpha^2 ||\boldsymbol{x}||^2 \ge ||\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}_{\alpha}||^2 + \alpha^2 ||\boldsymbol{x}_{\alpha}||^2$, also

$$\|oldsymbol{b}-oldsymbol{A}oldsymbol{x}_{lpha}\|^2=\|oldsymbol{b}-oldsymbol{A}oldsymbol{x}_{lpha}\|^2\geqlpha^2[\|oldsymbol{x}_{lpha}\|^2-\|oldsymbol{x}\|^2]$$

für alle x, d. h., x_{α} löst (17). Im Fall $\alpha > 0$ ist x_{α} die einzige Lösung von (1), so daß für $x \pm x_{\alpha}$ strenge Ungleichheit gilt und x_{α} auch die einzige Lösung von (17) ist. Im Fall $\alpha = 0$ ist $x_{\alpha} = x^{N}$. Alle weiteren Lösungen müssen denselben Zielfunktionswert wie x^{N} haben, also Lösungen von $Ax \cong b$ sein. Wegen der Eindeutigkeit der Normallösung haben jedoch alle solche $x \pm x^{N}$ eine größere Norm als $||x^{N}|| = \omega$, sind also nicht zulässig.

Zu (ii): Es gelte $||\boldsymbol{r}_{\alpha}|| = \tau$. Nach (i) gilt dann $||\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}|| > ||\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}_{\alpha}|| = ||\boldsymbol{r}_{\alpha}||$ für alle $\boldsymbol{x} \neq \boldsymbol{x}_{\alpha}$ mit $||\boldsymbol{x}|| \leq ||\boldsymbol{x}_{\alpha}||$, d. h., es gibt kein zulässiges \boldsymbol{x} , das einen kleineren Zielfunktionswert als \boldsymbol{x} . liefert. \Box

11.3.4. Bemerkung. (i) Zur Bestimmung von α kann ein skalares Nullstellenverfahren — etwa das Newton-Verfahren — auf die Gleichungen

 $f(\alpha) := \|\boldsymbol{x}_{\alpha}\| - \omega = 0$ bzw. $g(\alpha) = \|\boldsymbol{r}_{\alpha}\| - \tau = 0$

oder i. allg. besser auf die äquivalenten Gleichungen

$$\hat{f}(\alpha) := \frac{1}{\|\boldsymbol{x}_{\alpha}\|} - \frac{1}{\omega} = 0$$
 bzw. $\hat{g}(\alpha) = \frac{1}{\|\boldsymbol{r}_{\alpha}\|} - \frac{1}{\tau} = 0$

angewendet werden. Die Ableitungen lassen sich analog zu (15), (16) leicht berechnen.

(ii) Im Fall $||x^{N}|| < \omega$ ist x^{N} die normkleinste Lösung von (17). Im Fall $\tau < ||r^{N}||$ hat (18) keine Lösung, und für $||b|| \leq \tau$ ist x = o die einzige Lösung. \Box

Wir kehren jetzt wieder zur ursprünglichen Fragestellung zurück, die unbekannte, zu A^* gehörende Normallösung $x^* = (A^*)^+ b$ durch die zur gestörten Matrix $A = A^* + \delta A$ gehörende regularisierte Lösung $x_{\alpha} = S_{\alpha}b$ zu approximieren. Um den Regularisierungsparameter α möglichst gut festlegen zu können, benötigen wir eine Abschätzung von $||x^* - x_{\alpha}||$ als Funktion von α . Zur Ableitung einer solchen Schranke fügen wir

$$\boldsymbol{x}^{*}_{a} = \boldsymbol{S}^{*}_{a}\boldsymbol{b} = [(\boldsymbol{A}^{*})^{\mathsf{T}}\boldsymbol{A}^{*} + \alpha^{2}\boldsymbol{I}]^{-1} (\boldsymbol{A}^{*})^{\mathsf{T}}\boldsymbol{b}$$

als Hilfspunkt ein und schätzen gemäß

$$\|oldsymbol{x^*}-oldsymbol{x_{lpha}}\| \leq \|oldsymbol{x^*}-oldsymbol{x_{lpha}}^*\| + \|oldsymbol{x_{lpha}}^*-oldsymbol{x_{lpha}}\|$$

ab. Eine Schranke für den ersten Term ergibt sich aus (11), wenn dort $A = A^*$ gesetzt wird. Zur Abschätzung von $x^*_{\alpha} - x_{\alpha} = (S^*_{\alpha} - S_{\alpha}) b$ stellen wir $S^*_{\alpha} - S_{\alpha}$ für $\alpha > 0$ in der Form

$$egin{aligned} M_{lpha}(oldsymbol{S}^{st}_{lpha}-oldsymbol{S}^{st}_{lpha}-A^{\intercal}&=(A^{\intercal}A+lpha^{2}oldsymbol{I})\,oldsymbol{S}^{st}_{lpha}-A^{\intercal}\ &=-A^{\intercal}(oldsymbol{I}-Aoldsymbol{S}^{st}_{lpha})+lpha^{2}oldsymbol{S}^{st}_{lpha}\ &=-A^{\intercal}(oldsymbol{I}-Aoldsymbol{S}^{st}_{lpha})+(A^{st})^{\intercal}\,(oldsymbol{I}-A^{st}oldsymbol{S}^{st}_{lpha})\ &=-A^{\intercal}(A^{st}-A)\,oldsymbol{S}^{st}_{st}+(A^{st}-A)^{\intercal}\,(oldsymbol{I}-A^{st}oldsymbol{S}^{st}_{st}), \end{aligned}$$

also

$$S_{\alpha}^{*} - S_{\alpha} = -S_{\alpha}(A^{*} - A) S_{\alpha}^{*} + M_{\alpha}^{-1}(A^{*} - A)^{\mathsf{T}} (I - A^{*}S_{\alpha}^{*})$$
(19)

dar. Anwendung auf **b** liefert mit $r^*_{\alpha} := b - A^* x^*_{\alpha}$

$$\boldsymbol{x}_{\alpha}^{*} - \boldsymbol{x}_{\alpha} = -\boldsymbol{S}_{\alpha}(\boldsymbol{A}^{*} - \boldsymbol{A}) \, \boldsymbol{x}_{\alpha}^{*} + \boldsymbol{M}_{\alpha}^{-1}(\boldsymbol{A}^{*} - \boldsymbol{A})^{\mathsf{T}} \, \boldsymbol{r}_{\alpha}^{*}. \tag{20}$$

Es gelte nun

$$\mathbf{A} = \mathbf{A}^* + \mathbf{\delta} \mathbf{A}_D + \mathbf{\delta} \mathbf{A}_M, \qquad \|\mathbf{A}^* - \mathbf{A}\| = \|\mathbf{\delta} \mathbf{A}_D + \mathbf{\delta} \mathbf{A}_M\| = \varepsilon, \qquad (21)$$

vgl. (11.2.2). Dann ergibt sich aus (20)

$$egin{aligned} \|m{x}^{*}-m{x}_{lpha}\|&\leq \|m{x}^{*}-m{x}_{lpha}^{*}\|+\|m{x}_{lpha}^{*}-m{x}_{lpha}\| \ &\leq \|m{x}^{*}-m{x}_{lpha}^{*}\|+arepsilon[\|m{S}_{lpha}\|\,\|m{x}_{lpha}^{*}\|+\|m{M}_{lpha}^{-1}\|\,\|m{r}_{lpha}^{*}\|], \end{aligned}$$

wegen $||\mathbf{r}_{\alpha}^{*}|| \leq ||\mathbf{r}^{*}|| + ||\mathbf{r}_{\alpha}^{*} - \mathbf{r}^{*}||$ und der auch für die aus A^{*} berechneten Größen gültigen Ungleichungen (11) bis (13) also

$$\|\boldsymbol{x}^{*} - \boldsymbol{x}_{\alpha}\| \leq \{\alpha^{2} \| (\boldsymbol{A}^{*})^{+} \|^{2} + \varepsilon [\|\boldsymbol{S}_{\alpha}\| + \|\boldsymbol{M}_{\alpha}^{-1}\| \left(\|\boldsymbol{r}^{*}\| / \|\boldsymbol{x}^{*}\| + \alpha^{2} \| (\boldsymbol{A}^{*})^{+}\| \right)] \} \|\boldsymbol{x}^{*}\|.$$
(22)

Aus (4) folgt nun $\|M_{\alpha}^{-1}\| \leq 1/\alpha^2$. Zur Abschätzung von S_{α} beachten wir (8), (9) und erhalten

$$\|\mathbf{S}_{\mathbf{a}}\| = \|\mathbf{\Omega}_{\mathbf{a}}\| = \max\left\{\frac{\sigma}{\sigma^2 + \alpha^2} : \sigma = \sigma_i\right\} \leq \frac{1}{2\alpha} =: \gamma(\alpha), \tag{23}$$

denn die Funktion $\varphi_{\alpha}(\sigma) := \sigma/(\sigma^2 + \alpha^2)$ hat für $0 \leq \sigma$ genau ein Maximum an der Stelle $\sigma = \alpha$, und der maximale Funktionswert ist $\varphi_{\alpha}(\alpha) = 1/(2\alpha)$. Unter Verwendung

der Abkürzungen $\mu := \|(A^*)^+\|, \varrho := \|r^*\|/\|x^*\|$ kann dann in (22) weiter gemäß

$$\|\boldsymbol{x}^{*} - \boldsymbol{x}_{a}\| \leq \left\{ \alpha^{2} \mu^{2} + \varepsilon \left[\frac{\varrho}{\alpha^{2}} + \frac{1}{2\alpha} + \mu \right] \right\} \|\boldsymbol{x}^{*}\| = : \eta_{\varepsilon}(\alpha) \|\boldsymbol{x}^{*}\|$$
(24)

abgeschätzt werden. Die Schrankenfunktion $\eta_{\epsilon}(\alpha)$ hat das in Abb. 11.3.2 skizzierte Bild.



Wünschenswert ist daher die Festlegung von $\alpha = \alpha_{opt}^*$, so daß $\eta_{\epsilon}(\alpha)$ minimal wird. Wir unterscheiden dazu zwei Fälle:

Fall 1: $\varrho \neq 0$, d. h. $r^* = b - A^*x^* \neq o$ (inkonsistenter Fall) Um die auf eine Gleichung vierten Grades führende exakte Minimierung zu umgehen, vernachlässigen wir den gegenüber ϱ/α^2 für $\alpha \to 0$ kleinen Term $1/(2\alpha)$ und bestimmen $\alpha = \alpha_{out}$ als Minimumstelle von $\alpha^2 \mu^2 + \varepsilon [\varrho/\alpha^2 + \mu]$, was auf

$$\alpha_{\text{opt}} = \left(\frac{\epsilon \varrho}{\mu^2}\right)^{1/4} \quad \text{mit} \quad \eta_{\epsilon}(\alpha_{\text{opt}}) = \epsilon^{1/2} (2\mu \varrho^{1/2}) \left(1 + \frac{\epsilon^{1/4}}{4\varrho^{3/4} \mu^{1/2}}\right) \tag{25}$$

führt. Durch elementare Überlegungen überzeugt man sich davon, daß die Ordnung $O(\varepsilon^{1/2})$ auch bei der Wahl von $\alpha = \alpha_{opt}^*$ als exakter Minimumstelle von η_{ε} nicht verbessert werden kann.

Fall 2: $\rho = 0$, d. h. $r^* = b - A^*x^* = o$ (konsistenter Fall) Hier ergibt sich α als Minimumstelle von $\alpha^2 \mu^2 + \varepsilon [1/(2\alpha) + \mu]$ zu

$$\alpha_{\text{opt}} = \alpha_{\text{opt}}^* = \left(\frac{\varepsilon}{4\mu^2}\right)^{1/3} \quad \text{mit} \quad \eta_{\epsilon}(\alpha_{\text{opt}}) = \varepsilon^{2/3} \left(\frac{3\mu^{2/3}}{4^{2/3}}\right) \left(1 + \frac{\varepsilon^{1/3} 4^{2/3} \mu^{1/3}}{3}\right), \quad (26)$$

d. h., es kann eine Fehlerordnung $O(\varepsilon^{2/3})$ garantiert werden.

Bei feinerer Abschätzung von S_{α} als durch $\gamma(\alpha)$ gemäß (23) läßt sich im Fall 2 sogar die Fehlerordnung $O(\varepsilon)$ erreichen, sofern das Fehlerniveau ε im Sinne von

$$2\varepsilon < \sigma_r^* = 1/||(A^*)^+|| = 1/\mu \tag{27}$$

genügend klein ist. Wir erinnern daran, daß (27) gerade die für eine korrekte α -Wahl bei diskreter Regularisierung mittels SVD geforderte Bedingung (11.2.17) ist. Aus

(27) folgt nämlich wegen $|\sigma_i^* - \sigma_i| \leq \varepsilon$, also $\sigma_i^* - \varepsilon \leq \sigma_i \leq \sigma_i^* + \varepsilon$ (i = 1, ..., n) die Ungleichung

$$0 \leq \sigma_n \leq \cdots \leq \sigma_{r+1} \leq \varepsilon < \sigma_r^* - \varepsilon \leq \sigma_r \leq \cdots \leq \sigma_1.$$
⁽²⁸⁾

Die Singulärwerte von A liegen also entweder links von ε oder rechts von $\sigma_r^* - \varepsilon$, d. h., zwischen ε und $\sigma_r^* - \varepsilon$ tritt eine sog. Spektrumslücke auf. Die Schranke (23) kann dann gemäß

$$\|\mathbf{S}_{\mathbf{a}}\| \leq \max\left\{\frac{\sigma}{\sigma^{2} + \alpha^{2}} : 0 \leq \sigma \leq \varepsilon \text{ und } \sigma_{\mathbf{r}}^{*} - \varepsilon \leq \sigma\right\} =: \gamma_{\varepsilon}(\alpha) \leq \gamma(\alpha) = \frac{1}{2\alpha}$$
(29)

mit



verfeinert werden, siehe Abb. 11.3.3. Man kann elementar überprüfen, daß $\alpha^2 \mu^2 + \epsilon [\gamma_{\epsilon}(\alpha) + \mu]$ für dasjenige α minimal wird, das im Intervall $\epsilon < \alpha < \sigma_r^* - \epsilon =: \overline{\sigma}$ liegt und zu gleichen Randwerten $\varphi_{\alpha}(\varepsilon)$ und $\varphi_{\alpha}(\overline{\sigma})$ führt. Dies ergibt

$$\alpha_{\text{opt}} = (\epsilon \overline{\sigma})^{1/2} \quad \text{und} \quad \gamma_{\epsilon}(\alpha_{\text{opt}}) = 1/(\overline{\sigma} + \epsilon) = 1/\sigma_{r}^{*} = \mu,$$
(31)

also

$$\begin{aligned} \|\boldsymbol{x}^{*} - \boldsymbol{x}_{a_{opt}}\| &\leq \{\boldsymbol{x}_{opt}^{2} \mu^{2} + \boldsymbol{\varepsilon}[\boldsymbol{\gamma}_{\boldsymbol{\varepsilon}}(\boldsymbol{\alpha}_{opt}) + \mu]\} \|\boldsymbol{x}^{*}\| = \{\boldsymbol{\varepsilon}(\bar{\boldsymbol{\sigma}}\mu) \ \mu + \boldsymbol{\varepsilon}[\mu + \mu]\} \|\boldsymbol{x}^{*}\| \\ &\leq 3\boldsymbol{\varepsilon}\mu \ \|\boldsymbol{x}^{*}\|, \end{aligned}$$
(32)

man beachte $\bar{\sigma}\mu = (\sigma_r^* - \varepsilon)/\sigma_r^* \leq 1$.

Zusammenfassend erhalten wir das folgende Resultat.

11.3.5. Aussage. Unter der Voraussetzung $||A^* - A|| \leq \varepsilon$ kann durch geeignete Wahl von $\alpha = \alpha_{opt} > 0$ für den Fehler von $x_{\alpha} := S_{\alpha}b$ gegenüber $x^* := (A^*)^+ b$ die Schranke

$$\|\boldsymbol{x^*} - \boldsymbol{x}_a\| / \| \boldsymbol{x^*} \| \leq K \varepsilon^{\theta} (1 + N \varepsilon^{\theta/2})$$
 (33)

mit

$$heta := \left\{ egin{array}{ll} 1/2 & {
m für} & {m r}^{m *} \ne {m o}\,, \ 2/3 & {
m für} & {m r}^{m *} = {m o}\,, \ 1 & {
m für} & {m r}^{m *} = {m o}\,\,{
m und}\,\,2arepsilon\,\,\|({m A}^{m *})^+\| < 1 \end{array}
ight.$$

garantiert werden, wobei K > 0 und $N \ge 0$ die durch (25), (26) bzw. (32) gegebenen, nur von A^* und **b** abhängenden Zahlenwerte haben.

11.3.6. Bemerkung. (i) Im konsistenten Fall wird bei Gültigkeit von (27) mit $\alpha = \alpha_{opt}$ nach (32) die Genauigkeit

$$\|oldsymbol{x^*}-oldsymbol{x_{lpha}}\|/\|oldsymbol{x^*}\|\leq 3arepsilon\,\|(A^*)^+\|$$

garantiert. Das ist praktisch die volle erreichbare Genauigkeit, denn nach 8.2.7 muß selbst bei Ranggleichheit von A^* und A mit einem Fehler

$$\|\boldsymbol{x^*} - \boldsymbol{x}\| / \|\boldsymbol{x}\| \leq \varepsilon \sqrt{2} \|(\boldsymbol{A^*})^+\| / (1 - \varepsilon \|(\boldsymbol{A^*})^+\|) \leq 2 \sqrt{2} \varepsilon \|(\boldsymbol{A^*})^+\|$$

gerechnet werden.

(ii) Im inkonsistenten Fall wird für beliebiges ε eine Genauigkeit $O(\varepsilon^{1/2})$ garantiert. Bei diskreter Regularisierung muß dagegen ε im Sinne von (27) klein sein, allerdings ist die Fehlerordnung dann $O(\varepsilon)$, vgl. 11.2.3(ii). Grob gesprochen läßt sich also für $r^* \neq o$ mittels kontinuierlicher Regularisierung höchstens etwa die Hälfte der gültigen Ziffern erreichen, die mittels diskreter Regularisierung erzielt werden können.

(iii) Bei Computerrechnung ist A in (21) durch $A = U\Sigma V^{\intercal} - \delta A_R$ zu ersetzen, vgl. (11.2.3). Das Fehlerniveau ε kann sich dabei maximal um $\|\delta A_R\| \leq \nu F \|A\|$ vergrößern, siehe (11.2.5). Außerdem sind die bei der Berechnung von x_a bzw. r_a aus $\{U, \Sigma, V, b\}$ auftretenden Rundungsfehler mit zu berücksichtigen. Wir verzichten auf eine detaillierte, zu 11.2.2 analoge Angabe der Resultate einer solchen Fehleranalyse und bemerken lediglich, daß sie sich qualitativ nicht wesentlich von denen aus 11.3.5 unterscheiden. Bei Bedarf kann auch ein Fehler $b = b^* + \delta b$ der rechten Seiten mit berücksichtigt werden.

(iv) Falls α nicht nach dem bereits beschriebenen heuristischen Vorgehen durch Beobachtung von $||\boldsymbol{x}_{a}||$ und $||\boldsymbol{r}_{a}||$ für verschiedene α -Werte bzw. durch Vorgabe von $\omega = ||\boldsymbol{x}_{a}||$ oder $\tau = ||\boldsymbol{r}_{a}||$ bestimmt wird, können die angegebenen Formeln für α_{opt} zur Festlegung von α verwendet werden. Die in den Formeln auftretenden σ_{i}^{*} sind dann durch die berechneten σ_{i} zu ersetzen, wobei $r = p(\varepsilon)$ durch (11.2.10) festgelegt wird; unter der Bedingung (27) ist diese Festlegung korrekt. Für andere Prinzipien zur Festlegung von α sei auf B 11.3 verwiesen. \Box

Als Mittel zur Regularisierung rangdefizienter oder schlecht konditionierter Quadratmittelprobleme ist daher die diskrete Regularisierung i. allg. vorzuziehen. Es gibt jedoch Aufgaben, die aus anderen Gründen die Lösung kontinuierlicher Regularisierungsaufgaben für einen oder mehrere α -Werte erfordern, vgl. B 11.3. In diesem Fall ist es vom Aufwand her zweckmäßig, mit der Bidiagonalform statt mit der SVD zu arbeiten.

B. Kontinuierliche Regularisierung mittels der Bidiagonalform

Es sei

$$A = U_1 \left(rac{B}{O}
ight) V_1^{\mathsf{T}}, \quad B \in \mathsf{R}^{n,n} ext{ bidiagonal}, aga{35}$$

die nach 11.1.1 berechnete Bidiagonalform von A. Mit den Bezeichnungen $\beta := \left(\frac{\beta}{\bar{\beta}}\right)$:= $U_1^{\mathsf{T}} b$ und $\xi := V_1^{\mathsf{T}} x$ läßt sich dann (1) in der Form

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^{2} + \alpha^{2} \|\boldsymbol{x}\|^{2} = \left\| \left(\frac{\tilde{\beta}}{\bar{\beta}} \right) - \left(\frac{\boldsymbol{B}}{\boldsymbol{O}} \right) \boldsymbol{\xi} \right\|^{2} + \alpha^{2} \|\boldsymbol{\xi}\|^{2}$$
$$= \|\tilde{\boldsymbol{\beta}} - \boldsymbol{B}\boldsymbol{\xi}\|^{2} + \alpha^{2} \|\boldsymbol{\xi}\|^{2} + \|\bar{\boldsymbol{\beta}}\|^{2}$$
(36)

schreiben, so daß sich die Berechnung von $x_{lpha} = V_1 \xi_{lpha}$ auf die von ξ_{lpha} als Lösung von

$$\|\tilde{\boldsymbol{\beta}} - \boldsymbol{B}\boldsymbol{\xi}\|^{2} + \alpha^{2} \|\boldsymbol{\xi}\|^{2} = \left\| \left(\frac{\tilde{\boldsymbol{\beta}}}{\boldsymbol{o}} \right) - \left(\frac{\boldsymbol{B}}{\alpha \boldsymbol{I}} \right) \boldsymbol{\xi} \right\|^{2} \to \text{Minimum!}$$
(37)

reduziert. Zur Lösung von (37) beachten wir, daß der Block αI von links nach rechts durch Multiplikation mit $\sim 2n$ Givens-Drehungen bei Erhalt der Bidiagonalform von **B** annulliert werden kann; für n = 4 ist der erste Hauptschritt durch das folgende Muster gegeben:



Wie üblich sind dabei die neu zu berechnenden Elemente eingerahmt worden. Mit dieser Technik ergibt sich das nachstehend beschriebene Verfahren.

11.3.7. Lösung von $Ax \simeq b$ mittels kontinuierlicher Regularisierung der Bidiagonalform.

- S1 (Festlegung der regularisierten Bidiagonalform): S1.1: Berechne Bidiagonalform (35) von A nach 11.1.1 S1.2: Wähle Regularisierungsniveau $\alpha > 0$

- S2 (Berechnung der zu **b** gehörenden Größen $x_a, r_a = b Ax_a$ und deren Normen $nx_a = ||x_a||, nr_a = ||r_a||$):
 - ${}_{m{eta}} := \left(rac{ ilde{m{eta}}}{ar{m{eta}}}
 ight) := m{U}_1^{\mathsf{T}}m{b}$

S2.1: Berechne

S2.2: Bestimme
$$\bar{Q}^{\intercal} := G_n \cdots G_2 G_1$$
 mit
 $G_i := G_{n+i, n+i+1} G_{i, n+i}$ $(i = 1, ..., n - 1),$ $G_n := G_{n, 2n}$
als Produkt von $2n - 1$ Givens-Drehungen so, daß

$$\bar{Q}^{\mathsf{T}}\left(\frac{B}{\alpha I}\right) = \left(\frac{\bar{B}}{O}\right)$$

mit einer Bidiagonalmatrix \bar{B} gilt.

 $\begin{array}{l} \text{S2.3: Berechne}\left(\frac{\tilde{\gamma}}{\hat{\gamma}}\right) := \bar{\boldsymbol{Q}}^{\intercal}\left(\frac{\tilde{\beta}}{o}\right), \text{ bestimme } \boldsymbol{\xi}_{\mathfrak{a}} \text{ als Lösung des Bidiagonalsystems} \\ \overline{\boldsymbol{B}}\boldsymbol{\xi}_{\mathfrak{a}} = \tilde{\gamma}, \text{ setze } n\boldsymbol{x}_{\mathfrak{a}} := \|\boldsymbol{\xi}_{\mathfrak{a}}\| \text{ und } \boldsymbol{x}_{\mathfrak{a}} := \boldsymbol{V}_{1}\boldsymbol{\xi}_{\mathfrak{a}} \\ \text{S2.4: Berechne}\left(\frac{\tilde{\varrho}}{\hat{\varrho}}\right) := \bar{\boldsymbol{\mathcal{Q}}}\left(\frac{o}{\hat{\gamma}}\right), \text{ setze } \varrho_{\mathfrak{a}} := \left(\frac{\tilde{\varrho}}{\tilde{\beta}}\right), n\boldsymbol{r}_{\mathfrak{a}} := \|\varrho_{\mathfrak{a}}\|, \boldsymbol{r}_{\mathfrak{a}} := \boldsymbol{U}_{1}\varrho_{\mathfrak{a}} \end{array}$

11.3.8. Bemerkungen. (i) Die Matrizen U_1, V_1 können in Produktform auf dem Platz von A gespeichert werden, vgl. 11.1. Das Problem der effektiven Berechnung und Darstellung der in der SVD vorkommenden Faktoren U, V wird daher umgangen. Für $m \gg n$ ist der Aufwand zur Berechnung von $\{U_1, B, V_1\}$ etwa $\sim mn^2$, entspricht also dem einer QR-Faktorisierung von A. Die Berechnung von x_a bzw. r_a aus $\{U_1, B, V_1, \beta = U_1^T b\}$ erfordert $\sim K_1 n^2$ bzw. $\sim K_2 mn$ Operationen, die Berechnung allein von $||x_a||$ und $||r_a||$ sogar nur $\sim K_3 n$ Operationen.

(ii) Statt der Bidiagonalform könnte auch die **QR**-Faktorisierung von A mit Spaltenvertauschungen verwendet werden, siehe U 11.3.3. Der Faktorisierungsaufwand ist für $m \gg n$ etwa derselbe wie in 11.3.7, allerdings erfordert die Berechnung von x_{α} bzw. r_{α} aus $\{Q, R, P, c = Q^{\mathsf{T}}b\}$ dann $\sim K_4 n^3$ bzw. $\sim K_5 m n^2$ Operationen, so daß diese Version teurer als 11.3.7 ist.

(iii) Bei Berechnung von x_{α} aus den regularisierten Normalgleichungen (3) ist für jedes α eine erneute Cholesky-Faktorisierung von M_{α} mit $\sim n^3/6$ opms erforderlich.

(iv) In verschiedenen Anwendungen ist es zweckmäßig, statt (1) das gewichtete regularisierte Problem

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^{2} + \alpha^{2} \|\boldsymbol{D}\boldsymbol{x}\|^{2} = \left\| \left(\frac{\boldsymbol{b}}{\boldsymbol{o}}\right) - \left(\frac{\boldsymbol{A}}{\alpha \boldsymbol{D}}\right)\boldsymbol{x} \right\|^{2} \to \text{Minimum}!$$
(38)

zu verwenden, wobei $D \in \mathbb{R}^{n,n}$ eine reguläre Gewichtsmatrix bezeichnet. Für den praktisch wichtigen Sonderfall $D = \text{diag}(d_i), d_i > 0$, lassen sich 11.3.1 und 11.3.7 fast wörtlich übertragen; für allgemeinere Wichtung siehe Ü 11.3.4. \Box

Übungsaufgaben

Ü 11.3.1. Man zeige, daß für jedes $m{x}=m{x}_{lpha}+m{h}$ die Identität

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^{2} - \|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}_{\alpha}\|^{2} = \alpha^{2}[-2\boldsymbol{h}^{\mathsf{T}}\boldsymbol{x}_{\alpha} - \|\boldsymbol{h}\|^{2}] + \alpha^{2}\|\boldsymbol{h}\|^{2} + \|\boldsymbol{A}\boldsymbol{h}\|^{2}$$
(39)

gilt, und folgere hieraus, daß x_{α} mit $||x_{\alpha}|| = \omega$ die einzige Lösung von (17) ist.

Hinweis: Es gilt $\|\boldsymbol{x}\| \leq \|\boldsymbol{x}_{a}\|$ genau dann, wenn $2\boldsymbol{h}^{\mathsf{T}}\boldsymbol{x}_{a} + \|\boldsymbol{h}\|^{2} \leq 0$ ist.

Ü 11.3.2. Man zeige, daß sich aus (19) die Abschätzung

$$\|(\boldsymbol{A}^{*})^{+} - \boldsymbol{S}_{a}\| \leq \left\{ \alpha^{2} \mu^{2} + \varepsilon \left[\frac{\hat{\varrho}}{\alpha^{2}} + \frac{1}{2\alpha} + \mu \right] \right\} \|(\boldsymbol{A}^{*})^{+}\|,$$

$$\hat{\varrho} := \frac{\|\boldsymbol{I} - (\boldsymbol{A}^{*}) (\boldsymbol{A}^{*})^{+}\|}{\|(\boldsymbol{A}^{*})^{+}\|} \leq \frac{1}{\mu}$$
(40)

ergibt, und folgere hieraus die Existenz eines $\alpha = \alpha_{out}$ mit

$$||(A^*)^+ - S_{\alpha}||/||(A^*)^+|| \leq K \varepsilon^{1/2} (1 + N \varepsilon^{1/4}).$$
(41)

Man vergleiche die Schranke (41) mit der bei diskreter Regularisierung unter der Voraussetzung (27) erreichbaren Genauigkeit, siehe 11.2.3 (ii).

Ü 11.3.3. Man überlege sich, daß x_{a} , r_{a} unter Verwendung der QR-Faktorisierung mit Spaltenvertauschungen wie folgt berechnet werden können:

S 1: Berechne Faktorisierung $AP^{\intercal} = QR = Q\left(\frac{\tilde{R}}{O}\right)$ analog zu 11.2.B:

S2.1: Berechne

$$oldsymbol{c} := \left(rac{oldsymbol{ ilde{c}}}{oldsymbol{ar{c}}}
ight) := oldsymbol{Q}^{\intercal}oldsymbol{b}$$
 .

Dann gilt $\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^2 + \alpha^2 \|\boldsymbol{x}\|^2 = \|\boldsymbol{\tilde{c}} - \boldsymbol{\tilde{R}}\boldsymbol{y}\|^2 + \alpha^2 \|\boldsymbol{y}\|^2 + \|\boldsymbol{\bar{c}}\|^2 \text{ mit } \boldsymbol{y} := \boldsymbol{P}\boldsymbol{x}.$ S2.2: Bestimme $\bar{\boldsymbol{Q}}^{\intercal} := \boldsymbol{G}_1 \cdots \boldsymbol{G}_{n-1} \boldsymbol{G}_n$ mit

$$G_i := G_{n,n+i} \cdots G_{i+1,n+i} G_{i,n+i}$$
 $(i = n, ..., 1)$

als Produkt von n(n + 1)/2 Givens-Drehungen so, daß

$$\bar{Q}^{\mathsf{T}}\left(\frac{\bar{R}}{\alpha I}\right) = \left(\frac{\bar{R}}{O}\right)$$

mit einer oberen Dreiecksmatrix \bar{R} gilt. Berechne $\left(\frac{\tilde{f}}{\hat{f}}\right) := \bar{Q}^{\intercal}\left(\frac{\tilde{c}}{o}\right)$. Dann gilt

$$\|oldsymbol{ ilde{c}} - oldsymbol{ ilde{h}}oldsymbol{y}\|^2 + lpha^2 \|oldsymbol{y}\|^2 = \|oldsymbol{ ilde{f}} - oldsymbol{ ilde{R}}oldsymbol{y}\|^2 + \|oldsymbol{ ilde{f}}\|^2$$

S2.3: Berechne y als Lösung des Dreieckssystems $\overline{R}y = \tilde{f}$, setze $x_{\alpha} := P^{\mathsf{T}}y$. S2.4: Berechne

$$\left(rac{ ilde{m{a}}}{ ilde{m{d}}}
ight):=ar{m{Q}}\left(rac{m{o}}{ ilde{m{f}}}
ight), \ \ ext{setze} \ m{d}:=\left(rac{ ilde{m{d}}}{m{m{c}}}
ight).$$

Dann gilt $d = c - Ry = Q^{\mathsf{T}}(b - Ax_{\alpha})$. Bestimme $nr_{\alpha} := ||d||$ und $r_{\alpha} := Qd$.

Ü 11.3.4. Es sei $W \in \mathbb{R}^{n,n}$ eine symmetrische und positiv definite Gewichtsmatrix. Man führe das gewichtete regularisierte Problem

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{x}\|^2 + \alpha^2 \boldsymbol{x}^\mathsf{T} \boldsymbol{W}^{-1} \boldsymbol{x} \to \text{Minimum!}$$
(42)

unter Verwendung der Cholesky-Faktorisierung $W = LL^{\intercal}$ von W auf das ungewichtete Problem

$$\|\boldsymbol{b} - \boldsymbol{A}\boldsymbol{L}\boldsymbol{y}\|^2 + \alpha^2 \, \|\boldsymbol{y}\|^2 \to \text{Minimum!} \tag{43}$$

mit $\boldsymbol{x} = \boldsymbol{L} \boldsymbol{y}$ zurück.

Bemerkungen zum Kapitel 11

B 11.1. Das Bidiagonalisierungsverfahren 11.1.1 geht auf GOLUB/KAHAN [65], das darauf aufbauende Verfahren 11.1.3 zur Berechnung der SVD auf GOLUB/REINSCH [71] zurück; vgl. auch GOLUB/KAHAN [68]. Zur Implementierung sei auf DONGARRA et al. [79] verwiesen. Die für $m \gg n$ nur halb so teure Version aus Ü 11.1.1 wird bei LAWSON/HANSON [74] erwähnt; ein FORTRAN-Programm hat CHAN [82] angegeben.

B 11.2. Die hier "diskrete Regularisierung der SVD" genannte Regularisierungstechnik wird bei GOLUB/REINSCH [70] beschrieben. Algorithmus 11.2.4 zur diskreten Regularisierung der QR-Faktorisierung ist LAWSON/HANSON [74] entlehnt, wo auch detaillierte Hinweise zur Implementierung zu finden sind. Eine inkorrekte Formel für die Berechnung des Residuums — Formel (14.7) in LAWSON/HANSON — wurde dabei korrigiert. Die Abschätzung (11.2.48) geht auf FADDEEV et al. [68] zurück, ein Beweis kann bei LAWSON/HANSON nachgelesen werden. Die Beispiele aus Ü 11.2.2 stammen von KAHAN [66], siehe wieder LAWSON/HANSON.

B 11.3. Die Idee der kontinuierlichen Regularisierung geht u. a. auf LEVENBERG [44] und TYHONOV [63] zurück; ausführliche Darstellungen und weiterführende Literaturhinweise geben Lawson/Hanson [74], TYHONOV/ARSENIN [79], TYHONOV et al. [83] und HOFMANN [86]. Bei der Lösung nichtlinearer Quadratmittelprobleme mittels sog. regularisierter Gauß-Newton-Verfahren — häufig auch nach MARQUARDT benannt — wird entweder α oder ω gemäß 11.3.3 (i) vorgegeben. Im letzten Fall spricht man vom sog. "trust region approach". Die Steuerung von α durch $||\mathbf{r}_{\alpha}||$ über Defektprinzipien ist von Morozov [66, 73, 74/84] u. a. untersucht worden. Andere, auf dem Prinzip der "cross validation" beruhende Varianten sind bei GOLUB/HEATH/WAHBA [79] und FRIEDRICH/HOFMANN/TAUTENHAHN [79] zu finden. In der statistischen Literatur ist die kontinuierliche Regularisierung unter dem Namen "ridge regression" bekannt, und α heißt "ridge parameter". Die Identität (11.3.19), die daraus folgende Abschätzung und die Abschätzung (11.3.33) mit $\theta = 1/2$ bzw. $\theta = 2/3$ im Fall $r^* = o$ bzw. $r^* \pm o$ gehen auf VOEVODIN [69, 79] zurück, vgl. auch VOEVODIN/KUZNECOV [84, § 26]. Die unter der Voraussetzung (27) gültige feinere Abschätzung mit $\gamma_{\epsilon}(\alpha)$ nach (30) und die daraus für $r^* = o$ folgende Abschätzung (33) mit $\theta = 1$ hat KIEŁBASIŃSKI [76, unveröffentlichtes Manuskript] angegeben. Auf die Nützlichkeit der Bidiagonalform zur Berechnung von x_a ist von VOEVODIN [69] hingewiesen worden, allerdings wird dort das bidiagonale Problem durch Übergang zu den tridiagonalen Normalgleichungen gelöst.

12. Zusammenfassung zum Teil III

Das Vorgehen bei der Lösung linearer Quadratmittelprobleme $Ax \cong b$ hängt wesentlich davon ab, ob die Matrix $A \in \mathbb{R}^{m,n}$ Vollrang n hat oder ob A möglicherweise rangdefizient ist, also $r = \operatorname{rang}(A) < n$ gilt.

Im gutartigen Vollrangfall mit nicht zu großem $||A^+|| = 1/\sigma_n(A)$ sind Orthogonalisierungsverfahren die heute bevorzugte Methode, wobei die Householder-Faktorisierung am weitesten verbreitet und am besten durch Software hoher Qualität verfügbar gemacht ist. Vom Aufwand und dem Stabilitätsverhalten vergleichbar ist die modifizierte Gram-Schmidt-Orthogonalisierung. Sie sollte verwendet werden, wenn eine orthonormale Basis von $\mathcal{R}(A)$ explizit benötigt wird; andernfalls hat die Householder-Faktorisierung leichte Vorteile. Aufwandsmäßig liegt auch die **QR**-Faktorisierung mittels impliziter Givens-Drehungen in derselben Größenordnung. Die Fehlerkumulationskonstanten sind jedoch deutlich kleiner als bei Verwendung von Householder-Spiegelungen. Implizite wie auch die doppelt so teuren expliziten Givens-Drehungen erlauben außerdem auf Grund ihrer flexibleren Einsatzmöglichkeiten die Ausnutzung spezieller Strukturen der Matrix A, was etwa bei der Aufdatierung von **QR**- und **LL**^T-Faktorisierungen, der Diagonalisierung von Bidiagonalmatrizen und der **QR**-Faktorisierung schwachbesetzter Matrizen zum Tragen kommt. Leider ist zur Zeit keine Software auf der Grundlage impliziter Givens-Drehungen allgemein verfügbar, ein Zustand, der sich möglicherweise bis zum Erscheinen dieses Buches geändert hat.

Die im Fall $m \gg n$ nur halb so aufwendigen Normalgleichungsverfahren haben gegenüber den Orthogonalisierungsverfahren den Nachteil, daß die Klasse der damit behandelbaren Probleme bei gleicher Computergenauigkeit wesentlich kleiner ist und daß speziell im Fall von fast konsistenten Systemen mit kleinem Residuum r = b - Ax ein größerer Genauigkeitsverlust eintritt. Der zweite Nachteil läßt sich durch iterative Verbesserung oder durch Aufstellung und Lösung der Normalgleichungen in höherer Genauigkeit beheben. Für nicht zu schlecht konditionierte Aufgaben mit nicht zu kleinen Residuen sind die Normalgleichungsverfahren daher eine Alternative zu den Orthogonalisierungsverfahren; wegen ihres geringen Aufwands sind sie besonders für Klein- und Kleinstrechner attraktiv.

Grundsätzlich anders ist die Situation bei rangdefizienten Problemen mit nicht zu großem $||A^+|| = 1/\sigma_r(A)$, die sich von Vollrangproblemen mit einer Singulärwert-verteilung gemäß

$$\sigma_1 \geq \cdots \geq \sigma_r \gg \sigma_{r+1} \geq \cdots \geq \sigma_n > 0$$

vom numerischen Standpunkt nicht unterscheiden. Die Inkorrektheit dieser Aufgaben — ϵ -Störungen können sich wie $1/\epsilon$ in der Normallösung und der Pseudoinversen auswirken – macht die Anwendung von Regularirierungstechniken erforderlich. Wesentlicher Bestandteil der diskreten Regularisierungsverfahren ist eine numerische Rangbestimmung, die zuverlässig nur mittels der numerisch berechneten Singulärwertzerlegung möglich ist, und zwar nur dann, wenn das Fehlerniveau ε bekannt ist und deutlich unter $\sigma_t(A)$ liegt. In ähnlicher Weise lassen sich OR-Faktorisierungen mit Spaltenvertauschungen bzw. LL^T-Faktorisierungen der Normalgleichungsmatrix mit Zeilen- und Spaltenvertauschungen zur Rangbestinimung und diskreten Regularisierung verwenden, allerdings erhöht sich dabei das Risiko einer falschen Rangfestlegung u. U. wesentlich. Eine alternative Regularisierungsmöglichkeit bilden die kontinuierlichen Regularisierungsverfahren. In beiden Verfahrensklassen ist vielfach ein interaktives Vorgehen angebracht, bei dem der Wert des Regularisierungsparameters α (bzw. direkt der Pseudorang $p = p(\alpha)$) in Abhängigkeit von Kenngrößen wie $||x_{\alpha}||$ und $||r_{\alpha}||$ der regularisierten Lösung vom Bearbeiter festgelegt wird.

Das in der Statistik häufig erforderliche sequentielle Ändern des Problems durch Ändern/Streichen/Hinzufügen von Spalten/Zeilen läßt sich mittels geeigneter Aufdatierungstechniken ökonomisch realisieren, und zwar sowohl bei den verschiedenen QR-Faktorisierungen als auch bei der LL^{T} -Faktorisierung der Matrix der Normalgleichungen.

Probleme mit extrem unterschiedlichen Zeilenskalierungen — sog. steife Probleme — erfordern besondere Techniken bzw. Zusatzmaßnahmen, da eine Umskalierung wegen des festgelegten statistischen Hintergrundes i. allg. nicht möglich ist. Spaltenskalierungen verändern das Residuum nicht und sind zulässig. Falls nur Darstellungsfehler vorliegen, sollten die Normen der Spalten in derselben Größenordnung liegen.

Auf Quadratmittelprobleme hoher Dimension mit schwach besetzten Koeffizientenmatrizen konnte im Rahmen dieser Einführung nicht eingegangen werden. Bewährt haben sich hier **QR**-Faktorisierungen mittels Givens-Drehungen, aber auch iterative Verfahren für die Normalgleichungen. Wir verweisen auf die Übersichtsartikel von BJÖRCK [76], DUFF/REID [76], GEORGE/HEATH/PLEMMONS [81], GEORGE/ HEATH/NG [83] und HEATH [83] und die darin zitierte Literatur.

Auf die Erzeugung von Testmatrizen mit vorgegebenen Eigenschaften wie Kondition und Rang geht ZIELKE [85] in einer umfangreichen Arbeit ein.

IV. Eigenwertprobleme

13. Das spezielle symmetrische Eigenwertproblem

Es sei $S^{n,n} := \{A \in \mathbb{R}^{n,n} : A = A^{\mathsf{T}}\}$ die Menge der symmetrischen (n, n)-Matrizen. Das spezielle symmetrische Eigenwertproblem lautet dann:

Für eine gegebene Matrix $A \in S^{n,n}$ sind die Eigenwerte $\lambda_j = \lambda_j[A]$, d. h. die Zahlen λ_j , für die

 $Au^{j} = \lambda_{i}u^{j}$ mit $u^{j} \neq o$

gilt, und gegebenenfalls die zugehörigen Eigenvektoren u^{j} zu berechnen.

Je nach Herkunft der Aufgabe treten dabei spezifische Forderungen auf wie

- Berechnung des vollen Eigensystems, d. h. aller Eigenwerte und Eigenvektoren
- nur Berechnung aller Eigenwerte
- nur Berechnung gewisser Eigenwerte und gegebenenfalls zugehöriger Eigenvektoren, etwa Berechnung der l kleinsten oder größten Eigenwerte oder Berechnung der Eigenwerte aus einem vorgegebenen Intervall [a, b],

vgl. auch 2.1.C1. Das spezielle Eigenwertproblem tritt in den Anwendungen so oft auf, daß es zweckmäßig ist, für jede der o.g. spezifischen Aufgaben — und manche andere — spezielle Algorithmen zu entwickeln. Derzeit genutzte Programmpakete enthalten deshalb eine Menge von Basisprozeduren, die sich in flexibler Weise zur Behandlung verschiedener spezifischer Aufgaben einsetzen lassen.

Eine wesentliche Erfahrung des numerischen Rechnens, die sich auch theoretisch begründen läßt, soll schon hier erwähnt werden: Die Eigenwerte $\{\lambda_j\}$ von $A = (a_{ij}) \in \mathbf{S}^{n,n}$ sollen i. allg. niemals als Nullstellen des charakteristischen Polynoms $p(\lambda) = \det(A - \lambda I)$

 $=\sum_{k=0}^{n} c_k \lambda^k = 0 \text{ unter Verwendung der numerisch berechneten oder auch exakt bekannten Koeffizienten } \{c_k\} bestimmt werden. Der Berechnungsweg$

 $\{a_{ii}\} \rightarrow \{c_k\} \rightarrow \{\lambda_l\}$

ist i. allg. instabil, häufig sogar in katastrophaler Weise. Der Weg über das charakteristische Polynom scheidet daher i. allg. aus, und zwar auch für n = 2, siehe Ü 13.1.1. Wir bemerken abschließend, daß das spezielle symmetrische Eigenwertproblem wesentlich günstigere Eigenschaften im Hinblick auf Anzahl der benötigten arithmetischen Operationen, Speicherplatzanforderungen und erreichbare Genauigkeit aufweist als das spezielle Eigenwertproblem nichtsymmetrischer Matrizen.

13.1. Grundlegende Eigenschaften, Störungstheorie, Residualkriterien

A. Allgemeine Eigenschaften

Die Eigenschaften des durch die Matrix $A \in \mathbb{R}^{n,n}$ bzw. $A \in \mathbb{S}^{n,n}$ definierten speziellen Eigenwertproblems sind bereits im Abschnitt 1.2 zusammengestellt worden. Wir erinnern daran, daß sich der Fall einer symmetrischen Matrix A von dem einer nichtsymmetrischen im wesentlichen dadurch unterscheidet, daß

- alle Eigenwerte reell sind
- eine orthonormale Basis von Eigenvektoren existiert
- die Anzahl der positiven, negativen und verschwindenden Eigenwerte invariant unter Kongruenztransformationen ist.

In diesem Abschnitt seien die entsprechend ihrer Vielfachheit gezählten Eigenwerte $\{\lambda_i\}$ von A gemäß

$$\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_{n-1} \leq \lambda_n \tag{1}$$

geordnet, und $\{u^1, \ldots, u^n\}$ sei ein zugehöriges orthonormales System von Eigenvektoren, d. h., es gelte

$$\boldsymbol{A}\boldsymbol{u}^{j} = \lambda_{j}\boldsymbol{u}^{j} \quad \text{mit} \quad \boldsymbol{u}^{i\mathsf{T}}\boldsymbol{u}^{j} = \begin{cases} 1 & \text{für} \quad i = j, \\ 0 & \text{für} \quad i \neq j \end{cases} \quad (i, j = 1, ..., n).$$
(2)

Wenn

$$\boldsymbol{U} := (\boldsymbol{u}^1, \dots, \boldsymbol{u}^n) \quad \text{und} \quad \boldsymbol{\Lambda} := \text{diag} (\lambda_1, \dots, \lambda_n) \tag{3}$$

gesetzt wird, läßt sich (2) in der Gestalt

$$AU = UA$$
 mit $U^{\intercal}U = I$

schreiben, d. h., U ist orthogonal, und es gilt

$$U^{-1}AU = U^{\dagger}AU = \Lambda.$$
⁽⁴⁾

Die Matrix A wird also durch Ähnlichkeitstransformation mit der orthogonalen Matrix U auf Diagonalform Λ gebracht, und diese Transformation ist gleichzeitig eine Kongruenztransformation, vgl. 1.2.4 und 1.2.13. Wenn A mittels einer orthogonalen Matrix $Q \in \mathbb{R}^{n,n}$ gemäß

$$\bar{A} = Q^{\mathsf{T}} A Q \tag{5}$$

ähnlich in $\bar{A} \in S^{n,n}$ transformiert wird, ändern sich die Eigenwerte nicht, d. h., \bar{A} hat auch die Eigenwerte $\{\lambda_i\}$, und die Spalten der orthogonalen Matrix

$$ar{m{U}}=(ar{m{u}}^1,...,ar{m{u}}^n)=m{Q}^\intercalm{U}$$

bilden ein orthonormales System zugehöriger Eigenvektoren. Ist umgekehrt ein Eigenpaar $\{\lambda, \overline{u}\}$ von \overline{A} bekannt, so erhält man durch die Rücktransformation

$$\boldsymbol{u} = \boldsymbol{Q}\boldsymbol{\bar{u}} \tag{6}$$

ein Eigenpaar $\{\lambda, u\}$ von A, siehe 1.2.1.

Wir werden speziell gewählte Ähnlichkeitstransformationen (5) benutzen, um das durch A definierte Ausgangsproblem in das Eigenwertproblem einer Matrix \tilde{A} mit einfacherer Gestalt zu transformieren.

Eine weitere, praktisch wie theoretisch nützliche Transformation stellt die sog. Spektralverschiebung oder kurz Verschiebung (engl. "shift", russ. "сдвиг")

$$\bar{A} = A - \mu I \tag{7}$$

um $\mu \in \mathbf{R}$ dar. Offensichtlich ist $\{\lambda, u\}$ Eigenpaar von A genau dann, wenn $\{\overline{\lambda}, u\}$ mit

$$\bar{\lambda} = \lambda - \mu$$
(8)

Eigenpaar von \bar{A} ist. Durch (7) werden also die Eigenvektoren nicht verändert, während die Eigenwerte um das durch den reellen Verschiebungsparameter μ charakterisierte Stück auf der reellen Achse verschoben werden.

B. Störungstheorie

Wir betrachten jetzt neben der Matrix $A \in S^{n,n}$ mit den Eigenpaaren $\{\lambda_j, u^j\}$ gemäß (2) noch die gestörte Matrix

$$A + \delta A \quad \text{mit} \quad \delta A \in S^{n,n},$$

und es sei $\{\mu, u\}$ ein beliebiges Eigenpaar von $A + \delta A$ mit ||u|| = 1, d. h., es gelte

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}) \boldsymbol{u} = \boldsymbol{\mu}\boldsymbol{u}, \quad \|\boldsymbol{u}\| = 1.$$
(9)

Als Norm verwenden wir stets die Euklidische Vektornorm bzw. die Spektralnorm, sofern nichts Gegenteiliges explizit gesagt wird.

Wenn u gemäß

$$\boldsymbol{u} = \boldsymbol{U}\boldsymbol{z} = \sum_{j=1}^{n} z_{j} \boldsymbol{u}^{j} \quad \text{mit} \quad \|\boldsymbol{z}\| = 1$$
(10)

als Linearkombination der Eigenvektoren von A ausgedrückt wird, geht (9) nach Multiplikation mit U^{\intercal} von links in

$$U^{\intercal}(A + \delta A) Uz = \mu z \quad \text{mit} \quad ||z|| = 1$$

über, woraus wegen (4) und der Orthogonalität von U

$$\boldsymbol{y} := (\boldsymbol{\mu} \boldsymbol{I} - \boldsymbol{\Lambda}) \, \boldsymbol{z} = \boldsymbol{U}^{\mathsf{T}} \boldsymbol{\delta} \boldsymbol{A} \boldsymbol{U} \boldsymbol{z}, \tag{11}$$

also

$$|\boldsymbol{y}|| \leq \|\boldsymbol{\delta}A\| \|\boldsymbol{z}\| = \|\boldsymbol{\delta}A\| \tag{12}$$

folgt. Nun gilt

$$\|\boldsymbol{y}\|^{2} = \sum_{j=1}^{n} (\mu - \lambda_{j})^{2} z_{j}^{2} \ge \left\{ \min_{j} |\mu - \lambda_{j}| \right\}^{2} \sum_{j=1}^{n} z_{j}^{2} = \left\{ \min_{j} |\mu - \lambda_{j}| \right\}^{2}$$

so daß sich aus (12) sofort der folgende Satz ergibt:

13.1.1. Bauer-Fike-Theorem. Es gelte A, $\delta A \in S^{n,n}$, es seien $\{\lambda_1, \ldots, \lambda_n\}$ die Eigenwerte von A, und μ sei irgendein Eigenwert von $A + \delta A$. Dann gilt

min $|\mu - \lambda_j| \leq ||\mathbf{d}A||.$

Eine genauere Analyse zeigt, daß auch die folgende schärfere Aussage gültig ist:

13.1.2. Satz. Es gelte $A, \delta A \in S^{n,n}$. Es seien $\{\lambda_i\}$ die gemäß

 $\lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_n$ geordneten Eigenwerte von A, und $\{\mu_j\}$ seien die analog gemäß $\mu_1 \leq \mu_2 \leq \cdots \leq \mu_n$ geordneten Eigenwerte von $A + \delta A$. Dann gilt

 $|\mu_j - \lambda_j| \leq ||\boldsymbol{\delta} \boldsymbol{A}|| \qquad (j = 1, ..., n).$

Unter Verwendung der Frobeniusnorm kann schließlich die folgende simultane Abschätzung bewiesen werden:

13.1.3. Wielandt-Hoffman-Theorem. Unter den Voraussetzungen von 13.1.2 gilt

$$\sqrt{\sum\limits_{j=1}^n (\mu_j - \lambda_j)^2} \leq \| {\pmb \sigma} A \|_F$$

Die letzten beiden Sätze besagen: Die Eigenwerte symmetrischer Matrizen sind lipschitzstetige Funktionen der Matrix, und die Lipschitzkonstante hat in den angegebenen Normen den Wert 1. Das symmetrische Eigenwertproblem ist in diesem Sinne also extrem gut konditioniert.

Um zu einer besseren Einsicht in die individuelle Empfindlichkeit der einzelnen Eigenwerte und Eigenvektoren zu gelangen, betrachten wir wieder ein Eigenpaar $\{\mu, u\}$ von $A + \delta A$, das der Beziehung (9) genügt. Aus (10) folgt dann

$$\boldsymbol{z} = (z_i) = \boldsymbol{U}^{\mathsf{T}}\boldsymbol{u} \quad \text{mit} \quad \boldsymbol{z}_i = \boldsymbol{u}^{i\mathsf{T}}\boldsymbol{u} = \cos\left(\boldsymbol{\boldsymbol{\zeta}}\left(\boldsymbol{u}^i, \boldsymbol{u}\right)\right). \tag{13}$$

Mit $\boldsymbol{y} = (y_i)$ gemäß (11) ergibt sich ferner

$$y_i = (\mu - \lambda_i) z_i = \boldsymbol{u}^{i\mathsf{T}} \boldsymbol{\delta} \boldsymbol{A} \boldsymbol{u} \qquad (i = 1, ..., n).$$
(14)

Es sei jetzt j ein Index mit $z_i \neq 0$, d. h., u sei nicht orthogonal zu u^j . Aus (14) folgt dann

$$|\mu - \hat{\lambda}_j| = |\boldsymbol{u}^{j\intercal} \boldsymbol{\delta} \boldsymbol{A} \boldsymbol{u}| / |\boldsymbol{z}_j| \leq ||\boldsymbol{\delta} \boldsymbol{A} \boldsymbol{u}|| / |\boldsymbol{z}_j|.$$

Im Fall $\mu \neq \lambda_i$ läßt sich abschätzen, wie **u** und **u**ⁱ von der Orthogonalität abweichen, denn (13) und (14) ziehen dann

$$\left|\cos\left(\sphericalangle\left(\boldsymbol{u}^{i},\boldsymbol{u}
ight)
ight)
ight| = |\boldsymbol{z}_{i}| = |\boldsymbol{u}^{i\mathsf{T}}\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}|/|\mu - \lambda_{i}| \leq \|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\|/|\mu - \lambda_{i}|$$

nach sich. Ist für festes j sogar $\gamma_j = \min \{ |\mu - \lambda_i| : i \neq j \} > 0$, so ergibt sich mit (10), (14) und (11)

$$[\sin (\langle (u^j, u) \rangle]^2 = 1 - z_j^2 = \sum_{i+j} z_i^2 = \sum_{i+j} \frac{|y_i|^2}{|\mu - \lambda_i|^2} \le \frac{\|y\|^2}{\gamma_j^2} = \frac{\|\delta A u\|^2}{\gamma_j^2}$$

als Abschätzung für den Winkel zwischen u und u^{j} .

Zusammenfassend erhalten wir das folgende Ergebnis:

13.1.4. Aussage. Es gelte A, $\delta A \in S^{n,n}$, und es sei $\{\lambda_i, u^i\}$ ein orthonormales Eigensystem von A. Dann gilt für jedes Eigenpaar $\{\mu, u\}$ von $A + \delta A$ mit ||u|| = 1(i) im Fall $\boldsymbol{u}^{j\intercal}\boldsymbol{u} \neq 0$

$$|\mu - \lambda_{j}| \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\|}{|\boldsymbol{u}^{j\top}\boldsymbol{u}|} \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\|}{|\boldsymbol{u}^{j\top}\boldsymbol{u}|},$$
(ii) im Fall $\mu \neq \lambda_{i}$

$$|\cos\left(\langle (\boldsymbol{u}^{i}, \boldsymbol{u})\right)| \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\|}{|\boldsymbol{u}^{j\top}\boldsymbol{u}|} \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\|}{|\boldsymbol{u}^{j\top}\boldsymbol{u}|},$$
(15)

$$\left|\cos\left(\not\prec\left(\boldsymbol{u}^{i},\boldsymbol{u}\right)\right)\right| \leq \frac{\left\|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\right\|}{\left|\boldsymbol{\mu}-\lambda_{i}\right|} \leq \frac{\left\|\boldsymbol{\delta}\boldsymbol{A}\right\|}{\left|\boldsymbol{\mu}-\lambda_{i}\right|},$$
(iii) im Fall $\gamma_{j} := \min\left\{\left|\boldsymbol{\mu}-\lambda_{i}\right|: i \neq j\right\} > 0$

$$\sin\left(\not\prec\left(\boldsymbol{u}^{j},\boldsymbol{u}\right)\right) \leq \frac{\left\|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\right\|}{\left|\boldsymbol{u}-\lambda_{i}\right|} \leq \frac{\left\|\boldsymbol{\delta}\boldsymbol{A}\right\|}{\left|\boldsymbol{u}-\boldsymbol{u}\right|}.$$
(17)

$$\sin\left(\langle (\boldsymbol{u}^{j},\boldsymbol{u})\right) \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\boldsymbol{u}\|}{\gamma_{j}} \leq \frac{\|\boldsymbol{\delta}\boldsymbol{A}\|}{\gamma_{j}}.$$
(17)

13.1.5. Bemerkung. (i) Mit dem i. allg. nicht normierten Vektor $v := u/u^{j} u = u/z_{j}$ läßt sich (15) äquivalent als

$$|\mu - \lambda_j| \leq \| \mathbf{\delta} A \mathbf{v} \| \leq \| \mathbf{\delta} A \| \| \mathbf{v} \|$$

schreiben. Dabei gilt im Fall $z_i > 0$

$$\boldsymbol{v} = \boldsymbol{u}^{j} + \boldsymbol{\delta}\boldsymbol{u}^{j} \quad \text{mit} \quad \boldsymbol{u}^{j} \perp \boldsymbol{\delta}\boldsymbol{u}^{j}, \tag{18}$$

also

$$\|\boldsymbol{v}\|^2 = 1 + \|\boldsymbol{\delta}\boldsymbol{u}^j\|^2, \qquad \|\boldsymbol{\delta}\boldsymbol{u}^j\| = \|\boldsymbol{v}\|\sin\left(\boldsymbol{\boldsymbol{\triangleleft}}\left(\boldsymbol{u}^j,\boldsymbol{u}\right)\right), \tag{19}$$

vgl. Abb. 13.1.1; im Fall $z_j < 0$ ersetze man u^j durch $-u^j$.

(ii) Es seien jetzt $\{\mu_i\}$ die gemäß 13.1.2 geordneten Eigenwerte von $A + \delta A$, und $\{v^{j}\}$ sei ein zugehöriges orthonormales System von Eigenvektoren. Wenn $\{\mu, u\}$



Abb. 13.1.1. Eigenvektoren u^{j} , u und v

 $= \{\mu_j, v^j\}$ gesetzt wird, besagt Aussage (ii): Falls μ_j genügend von λ_i getrennt und δA genügend klein ist, sind die zugehörigen Eigenvektoren v^j von $A + \delta A$ und u^i von A fast orthogonal. Da nach 13.1.2

$$|\mu - \lambda_i| = |\mu_j - \lambda_i| \ge |\lambda_j - \lambda_i| - |\lambda_j - \mu_j| \ge |\lambda_j - \lambda_i| - \|\delta A\|$$
(20)

gilt, ist $\mu = \mu_i$ genügend von λ_i getrennt, wenn λ_i genügend von λ_i getrennt und $\boldsymbol{\sigma} \boldsymbol{A}$ genügend klein ist.

Dagegen besagt Aussage (iii): Falls μ_j genügend von allen λ_i $(i \neq j)$ getrennt und σA genügend klein ist, weichen die Richtungen von v^j und u^j nur wenig voneinander ab. Da aus (20)

$$\gamma_i \ge \min\left\{ |\lambda_i - \lambda_i| : i \neq j \right\} - \|\boldsymbol{\delta}\boldsymbol{A}\| \tag{21}$$

folgt, ist letzteres sicher der Fall, wenn λ_j von allen übrigen Eigenwerten λ_i genügend getrennt und σA genügend klein ist, also eine genügend große *Spektrumslücke* um λ_j vorliegt.

(iii) Wenn gewisse Eigenwerte λ_i $(i \neq j)$ im Sinne von $|\lambda_i - \lambda_j| = O(||\boldsymbol{\delta}\boldsymbol{A}||)$ pathologisch dicht zu λ_j benachbart sind, gilt $|\mu_j - \lambda_i| = O(||\boldsymbol{\delta}\boldsymbol{A}||)$ und somit $\gamma_j = O(||\boldsymbol{\delta}\boldsymbol{A}||)$, so daß (16) und (17) i. allg. keine brauchbaren Schranken liefern. Daß dies ohne Vorliegen einer genügend großen Spektrumslücke auch nicht erwartet werden kann, zeigt das Beispiel aus Ü 2.1.9. Ist ein solcher "*Haufen"* benachbarter Eigenwerte jedoch ausreichend von den restlichen Eigenwerten getrennt, so lassen sich analoge Schranken für den geeignet zu definierenden Winkel zwischen den Räumen, die durch die entsprechenden Eigenvektoren \boldsymbol{u}^j bzw. \boldsymbol{v}^j aufgespannt werden, angeben.

(iv) In gewissen wichtigen Fällen, auf die wir im folgenden eingehen werden, gilt $\| \boldsymbol{\delta} A \boldsymbol{u} \| \ll \| \boldsymbol{\delta} A \| \| \boldsymbol{u} \| = \| \boldsymbol{\delta} A \|$. Die Schranken mit $\| \boldsymbol{\delta} A \boldsymbol{u} \|$ sind dann wesentlich besser als diejenigen, die durch die Vergröberung $\| \boldsymbol{\delta} A \boldsymbol{u} \| \le \| \boldsymbol{\delta} A \|$ aus diesen entstehen. Wegen der Symmetrie von $\boldsymbol{\delta} A$ und der Gültigkeit von $\| \boldsymbol{u}^{i \dagger} \boldsymbol{\delta} A \boldsymbol{u} \| = \| (\boldsymbol{\delta} A \boldsymbol{u}^{i})^{\dagger} \boldsymbol{u} \| \le \| \boldsymbol{\delta} A \boldsymbol{u} \|$ kann überdies der Term $\| \boldsymbol{\delta} A \boldsymbol{u} \|$ in (15) und (16) durch $\| \boldsymbol{\delta} A \boldsymbol{u}^{j} \|$ bzw. $\| \boldsymbol{\delta} A \boldsymbol{u}^{i} \|$ ersetzt werden.

Zur Anwendung der Störungsabschätzungen betrachten wir eine Matrix $A^* \in S^{n,n}$, die durch ihre Computerdarstellung $A \in \Re^{n,n}$, $A = A^{\intercal}$, repräsentiert wird. Dann gilt

$$\|\boldsymbol{\delta}\boldsymbol{A}_{\boldsymbol{D}}\| = \|\boldsymbol{A}^* - \boldsymbol{A}\| \leq \boldsymbol{v} \, \mathbf{n} \, \|\boldsymbol{A}^*\|,$$

siehe 2.3.15, so daß aus 13.1.2 die Abschätzung

$$|\lambda_j^* - \lambda_j| \leq \nu \sqrt{n} ||A^*|| = \nu \sqrt{n} \max\{|\lambda_i^*|: i = 1, ..., n\}$$

$$(22)$$

folgt. Durch (22) ist das optimale Fehlerniveau von λ_j gegeben.

Es seien jetzt $\{\mu_j\}$ die nach einem numerisch gutartigen Verfahren aus A berechneten Eigenwertnäherungen, d. h., die $\{\mu_j\}$ sind die exakten Eigenwerte einer gestörten Matrix $A + \delta A, \delta A \in S^{n,n}$, wobei

$$\|\boldsymbol{\delta}\boldsymbol{A}\| \leq vF \|\boldsymbol{A}\| \tag{23}$$

gilt. Dann liefert 13.1.2 für $\delta \lambda_i := \mu_i - \lambda_i$ die zu (22) analoge Schranke

$$|\delta\lambda_j| \leq ||\boldsymbol{\delta}A|| \leq vF \max\{|\lambda_i|: i=1,...,n\},$$

und für den relativen Fehler gilt

$$\frac{|\delta\lambda_j|}{|\lambda_j|} \leq vF \frac{\max\left\{|\lambda_i|: i = 1, \dots, n\right\}}{|\lambda_j|}.$$
(24)

Die Schranke (24) garantiert nur für die betragsgroßen Eigenwerte λ_j einen kleinen relativen Fehler, während die betragskleinen μ_j durchaus stark von den zugehörigen exakten λ_j abweichen können. Die Forderung, auch die betragskleinen Eigenwerte von A mit hoher relativer Genauigkeit zu berechnen, ist deshalb i. allg. unrealistisch und nur schwer bzw. überhaupt nicht zu erfüllen. Wenn jedoch A eine spezielle Struktur hat und die Störung δA dieser Struktur angepaßt ist, liefern die individuellen Abschätzungen aus 13.1.4 auch kleine Schranken für den relativen Fehler der betragskleinen Eigenwerte. Beispiele solcher Matrizen sind die sog. gestuften (engl. "graded") Matrizen, bei denen die Beträge der Elemente von links oben nach rechts unten abnehmen. Zur Illustration betrachten wir das folgende Zahlenbeispiel.

13.1.6. Beispiel. Gegeben sei die gestufte Matrix

$$A = \begin{pmatrix} 10\,000 & 100 & 1 \\ 100 & 100 & 1 \\ 1 & 1 & 1 \end{pmatrix} \quad \begin{array}{c} \lambda_1 = 0.989\,9 \\ \text{mit} & \lambda_2 = 99.0 \\ \|A\| = \lambda_3 = 10\,001 \\ \end{array} \quad \text{und} \quad u^1 = \begin{pmatrix} 0.000\,001 \\ -0.010\,100 \\ 0.999\,950 \\ \end{array} \right).$$

Für $v := 10^{-6}$, F := 10 und eine im Sinne von (23) kleine Störung ist $||\delta A|| \leq 10^{-1}$, und (24) führt auf

$$|\delta\lambda_1/\lambda_1| \leq 10^{-6} \times 10 \times 10^4 = 10^{-1},\tag{25}$$

d. h., der kleinste Eigenwert λ_1 kann bis zu 10% gestört werden. Wenn jedoch die Störung δA in zu A analoger Weise gestuft ist, d. h., wenn

$$|\delta A| \le \nu F |A| = 10^{-5} |A| \tag{26}$$

gilt, ist wieder $\|\delta A\| \leq 10^{-5} \|A\| \approx 10^{-1}$, aber

$$|\delta A u^1| \leq 10^{-5} |A| \, |u^1| pprox 10^{-5} egin{pmatrix} 2 \ 2 \ 1 \end{pmatrix}, \qquad ||\delta A u^1|| \lessapprox 3 imes 10^{-5}$$

Aus (17) ergibt sich mit $\gamma_1 \approx 98$ zunächst $s := \sin (\langle (u^1, u) \rangle \leq 0.1/98 \approx 1 \times 10^{-3}$, wobei u

den zu $\mu = \lambda_1 + \delta \lambda_1$ gehörenden normierten Eigenvektor von $A + \delta A$ bezeichnet. Dies führt auf

$$|\boldsymbol{u}^{1\mathsf{T}}\boldsymbol{u}| = \left|\cos\left(\measuredangle\left(\boldsymbol{u}^{1},\boldsymbol{u}\right)\right)\right| = \sqrt{1-s^{2}} \gtrsim 1-5 \times 10^{-7}$$

und weiter mit (15) unter Beachtung von 13.1.5(iv) auf

$$|\delta\lambda_1/\lambda_1| \leq ||\delta A u^1||/|\lambda_1 u^{1\mathsf{T}} u| \leq 3 \times 10^{-5}, \tag{27}$$

also eine wesentlich bessere Schranke als (25). Man beachte, daß hier

$$\|\boldsymbol{\delta} \boldsymbol{A} \boldsymbol{u}^{1}\| pprox 3 imes 10^{-5} \ll 10^{-1} pprox \|\boldsymbol{\delta} \boldsymbol{A}\| \, \|\boldsymbol{u}^{1}\| = \|\boldsymbol{\delta} \boldsymbol{A}\|$$

gilt. Ebenso kann die oben benutzte Schranke für s verbessert werden: Aus (19) folgt $\| \delta u^{1} \|$ $=\sqrt{1+\|\delta u^1\|^2}$ s. also

$$\|\delta u^1\| = s / \sqrt{1 - s^2} \lessapprox 1 imes 10^{-3} / (1 - 5 imes 10^{-7}) pprox 1 imes 10^{-3}$$

und mit (18) – vgl. auch Abb. 13.1.1 –

$$egin{aligned} & \|\delta Aoldsymbol{u}\| = \|oldsymbol{u}^{1} oldsymbol{T}oldsymbol{u}\| \|\delta Aoldsymbol{u}\| & \leq \|\delta Aoldsymbol{u}^{1}\| + \|\delta A\| \|\delta u^{1}\| \lessapprox 3 imes 10^{-5} + 1 imes 10^{-4} \ & = 1.3 imes 10^{-4}. \end{aligned}$$

Nach (17) ergibt sich daraus die verbesserte Schranke

$$s \leq \|\boldsymbol{\delta} \boldsymbol{A} \boldsymbol{u}\|/\gamma_1 \lesssim 1.3 \times 10^{-6}$$

Durch wiederholte Anwendung dieses Prinzips erhält man schließlich $s \leq 3.1 \times 10^{-7}$ und $\|\delta u^1\| \lesssim 1.3 imes 10^{-7}$, d. h., u und u^1 weichen nur sehr wenig voneinander ab.

Das Beispiel zeigt, daß für die Berechnung der betragskleinen Eigenwerte gestufter Matrizen solche numerisch gutartigen Algorithmen verwendet werden sollten, bei denen die den erzeugten Rundungsfehlern äquivalente Störung σA in zu A analoger Weise gestuft ist.

C. Residualkriterien

Wie bei linearen Gleichungssystemen und Quadratmittelproblemen können auch bei Eigenwertaufgaben geeignete Residuen zur Kontrolle und gegebenenfalls Verbesserung der Genauigkeit verwendet werden. Wir beginnen mit der folgenden Aussage:

13.1.7. Aussage. Es sei $A \in S^{n,n}$, $\mu \in \mathbb{R}$, $u \in \mathbb{R}^n$ mit ||u|| = 1. Dann sind die folgen-(i) Es existiert ein $\delta A \in S^{n,n}$ mit $(A + \delta A) u = \mu u$ und $\| e^{(ij)} = F_{n,n} u^{(ij)}$ den Aussagen äquivalent:

$$(A + \delta A) u = \mu u \quad \text{und} \quad \|\delta A\| \leq \varepsilon.$$
(28)

(ii) Es gilt

$$\|\boldsymbol{r}\| \leq \varepsilon \quad \text{mit} \quad \boldsymbol{r} := \boldsymbol{A}\boldsymbol{u} - \mu\boldsymbol{u}. \tag{29}$$

Der Vektor $\mathbf{r} = (\mathbf{A} - \mu \mathbf{I}) \mathbf{u}$ heißt *Residuum* des Paares $\{\mu, \mathbf{u}\}$ in bezug auf das Eigenwertproblem der Matrix A.

Beweis. Aus (28) folgt $\boldsymbol{r} = -\delta \boldsymbol{A}\boldsymbol{u}$, also $\|\boldsymbol{r}\| \leq \|\delta \boldsymbol{A}\| \leq \varepsilon$. Gelte umgekehrt (29), d. h., es sei $\eta := \|\boldsymbol{r}\| \leq \varepsilon$. Dann gibt es eine Matrix $\boldsymbol{H} \in \mathbf{S}^{n,n}$ mit $\|\boldsymbol{H}\| = 1$ und $\boldsymbol{r} = \eta \boldsymbol{H}\boldsymbol{u}$, etwa $\boldsymbol{H} := \boldsymbol{I}$ im Fall $\eta = 0$ oder im Fall $\eta > 0$ die Spiegelungsmatrix, die \boldsymbol{u} in \boldsymbol{r}/η abbildet, siehe U 3.3.2. Dies zieht $\boldsymbol{A}\boldsymbol{u} - \mu\boldsymbol{u} = \eta \boldsymbol{H}\boldsymbol{u}$, also (28) mit $\delta \boldsymbol{A} := -\eta \boldsymbol{H}$, $\|\delta \boldsymbol{A}\| = \eta \leq \varepsilon$ nach sich. \Box

Aus 13.1.7 folgt insbesondere, $da\beta \langle \mu, u \rangle$ genau dann ein Eigenpaar der im Sinne von

 $\| \mathbf{\delta} A \| \leq \mathbf{v} F \| A \|$

zu A benachbarten Matrix $A + \delta A$ ist, wenn

 $\|\boldsymbol{r}\| = \|\boldsymbol{A}\boldsymbol{u} - \boldsymbol{\mu}\boldsymbol{u}\| \leq \boldsymbol{\nu}F \|\boldsymbol{A}\|$

gilt, d. h., wenn das relative Residuum $\|\mathbf{r}\|/\|A\|$ durch vF beschränkt ist.

Da verschiedene Algorithmen nur Näherungen für Eigenwerte, andere nur Näherungen für Eigenvektoren liefern, entstehen sofort die Fragen:

- Wie soll μ bei gegebener Eigenvektornäherung u gewählt werden, damit $\{\mu, u\}$ ein möglichst gutes Eigenpaar darstellt?
- Wie soll \boldsymbol{u} bei gegebener Eigenwertnäherung μ gewählt werden, damit { μ , \boldsymbol{u} } ein möglichst gutes Eigenpaar ist?

Entsprechend 13.1.7 soll dabei die Güte von $\{\mu, u\}$ durch die Norm des Residuums r charakterisiert werden.

Wir untersuchen im folgenden das erste Problem.

13.1.8. Aussage. Es seien $A \in S^{n,n}$ und $u \in \mathbb{R}^n$ mit ||u|| = 1 gegeben. Dann gilt:

(i) Das Quadratmittelproblem

$$\|A\boldsymbol{u} - \boldsymbol{\mu}\boldsymbol{u}\| \to \underset{\boldsymbol{\mu} \in \mathbf{R}}{\operatorname{Minimum}}$$
(30)

hat die eindeutige Lösung

$$\mu = \boldsymbol{u}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{u} =: \varrho. \tag{31}$$

Das zugehörige kürzeste Residuum $r := Au - \varrho u$ ist zu u orthogonal, d. h., es gilt

$$\boldsymbol{u}^{\mathsf{T}}\boldsymbol{r}=\boldsymbol{0}. \tag{32}$$

- (ii) Das Paar $\{\varrho, u\}$ ist Eigenpaar einer Matrix $A + \delta A$, wobei $\delta A \in S^{n,n}$ ist mit $\|\delta A\| \leq \|r\|$.
- (iii) Es gibt ein Eigenpaar $\{\lambda_i, u^j\}$ von A mit

$$|\rho - \lambda_i| = \min\{|\rho - \lambda_i| : i = 1, ..., n\} \leq ||\mathbf{r}||.$$
(33)

Im Fall $\gamma_i := \min \{ |\varrho - \lambda_i| : i \neq j \} > 0$ gilt

$$\sin\left(\langle \langle (\boldsymbol{u}, \boldsymbol{u}^{j}) \rangle \leq \|\boldsymbol{r}\|/\gamma_{j}$$
(34)

und

$$|\varrho - \lambda_j| \le \|\mathbf{r}\|^2 / \gamma_j. \tag{35}$$
Beweis. Die Normalgleichungen zu (30) lauten $u^{\mathsf{T}}u\mu = u^{\mathsf{T}}Au$, woraus sofort (31) und (32) folgen, vgl. 8.1.6. Die Aussage (ii) ergibt sich aus 13.1.7 und impliziert (33) nach 13.1.1, während sich (34) aus 13.1.4 ergibt. Es bleibt der Nachweis von (35). Wir setzen dazu $\eta_i := \lambda_i - \varrho$ (i = 1, ..., n) und numerieren die Eigenwerte $\{\lambda_i\}$ von A so, daß

$$|\eta_1| = |\lambda_1 - \varrho| \le |\eta_2| = |\lambda_2 - \varrho| \le \dots \le |\eta_n| = |\lambda_n - \varrho|$$
(36)

gilt. Dann ist

$$\nu_1 = |\eta_2|, \tag{37}$$

und (35) liest sich als

 $|\eta_1\eta_2| \le \|\boldsymbol{r}\|^2. \tag{38}$

Mit der Eigenwertzerlegung (4) und $\boldsymbol{z} = \boldsymbol{U}^{\mathsf{T}}\boldsymbol{u}$ ergibt sich nun

$$\boldsymbol{r} = (\boldsymbol{A} - \boldsymbol{\varrho} \boldsymbol{I}) \, \boldsymbol{u} = \boldsymbol{U} (\boldsymbol{\Lambda} - \boldsymbol{\varrho} \boldsymbol{I}) \, \boldsymbol{U}^{\mathsf{T}} \boldsymbol{u} = \boldsymbol{U} (\boldsymbol{\Lambda} - \boldsymbol{\varrho} \boldsymbol{I}) \, \boldsymbol{z}, \tag{39}$$

wegen der Orthogonalität von U also

$$\|\boldsymbol{r}\|^{2} = \|(\boldsymbol{\Lambda} - \boldsymbol{\varrho}\boldsymbol{I})\,\boldsymbol{z}\|^{2} = \sum_{i=1}^{n} \eta_{i}^{2} z_{i}^{2}.$$
(40)

Die Orthogonalitätsrelation (32) geht in

$$0 = \boldsymbol{u}^{\mathsf{T}} \boldsymbol{r} = \boldsymbol{z}^{\mathsf{T}} (\boldsymbol{.} \mathbf{1} - \varrho \boldsymbol{I}) \, \boldsymbol{z} = \sum_{i=1}^{n} \eta_i z_i^2$$
(41)

über, und außerdem gilt

$$1 = \boldsymbol{u}^{\mathsf{T}}\boldsymbol{u} = \boldsymbol{z}^{\mathsf{T}}\boldsymbol{z} = \sum_{i=1}^{n} z_{i}^{2}.$$
(42)

Wenn (41) mit $\sigma(|\eta_2| - |\eta_1|)$, $\sigma := \text{sgn}(\eta_1)$, und (42) mit $-|\eta_1\eta_2|$ multipliziert und zu (40) addiert wird, folgt

$$\|\mathbf{r}\|^{2} - |\eta_{1}\eta_{2}| = \sum_{i=1}^{n} \{\eta_{i}^{2} + \eta_{i}\sigma(|\eta_{2}| - |\eta_{1}|) - |\eta_{1}\eta_{2}|\} z_{i}^{2}$$
$$= \sum_{i=1}^{n} (\eta_{i} - \eta_{1}) (\eta_{i} + \sigma |\eta_{2}|) z_{i}^{2} \ge 0.$$
(43)

also die zu beweisende Ungleichung (38). Man beachte, daß die beiden Klammerausdrücke in (43) wegen (36) für jedes i dasselbe Vorzeichen haben, das Produkt also nicht negativ ist. \Box

Die dem Vektor u mit ||u|| = 1 gemäß (31) zugeordnete optimale Eigenwert-Näherung $\varrho = \varrho(u) = u^{\mathsf{T}}Au$ wird *Rayleigh-Quotient* von u bezüglich A genannt; für jeden nicht notwendig normierten Vektor $v \neq o$ ist er durch

$$\varrho(\boldsymbol{v}) := \frac{\boldsymbol{v}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{v}}{\boldsymbol{v}^{\mathsf{T}} \boldsymbol{v}} \tag{44}$$

definiert. Man beachte, daß $\varrho(\lambda v) = \varrho(v)$ für alle $\lambda \neq 0$, $v \neq o$ gilt, also $\varrho(v)$ nur von der Richtung, nicht aber von Orientierung und Länge des Argumentvektors v abhängt.

Wir bemerken an dieser Stelle, daß die Schranken (34), (35) wie die Schranke (17) aus 13.1.4 meist nur theoretische Bedeutung haben, da realistische untere Schranken für γ_i nur selten bekannt sind. Das folgende Zahlenbeispiel zeigt, daß sich solche Schranken beschaffen lassen, wenn λ_i genügend gut von den übrigen Eigenwerten λ_i getrennt ist und die λ_i genügend gut durch gewisse Rayleigh-Quotienten ϱ_i approximiert werden.

13.1.9. Beispiel. Betrachtet werde die Matrix

$$\boldsymbol{A} = \begin{pmatrix} 0.00 & 0.08 & 0.01 \\ 0.08 & 1.00 & 0.02 \\ 0.01 & 0.02 & 1.20 \end{pmatrix}.$$

Da die Diagonalelemente bis auf a_{11} dominant sind, bietet es sich an, die Einheitsvektoren e^1, e^2, e^3 als Näherungen u für die gesuchten Eigenvektoren u^1, u^2, u^3 zu wählen. Wenn ϱ_j den zu $u = e^j$ gehörenden Rayleigh-Quotienten und $\eta_j = ||r^j||$ die Norm des zugehörigen kürzesten Residuums bezeichnet, ergeben sich die Zahlenwerte

Aus (33) folgt dann

$$|\varrho_1 - \lambda_1| \le 0.0806, \quad |\varrho_2 - \lambda_2| \le 0.0825, \quad |\varrho_3 - \lambda_3| \le 0.0224$$
 (45)

bzw.

$$-0.0806 \leq \lambda_1 \leq 0.0806, \qquad 0.9175 \leq \lambda_2 \leq 1.0825, \qquad 1.1776 \leq \lambda_3 \leq 1.2224; (46)$$

da die Einschließungsintervalle (46) disjunkt sind, muß in jedem genau ein Eigenwert λ_i liegen. Aus (46) ergeben sich die folgenden unteren Schranken für $\gamma_j = \min \{|\varrho_j - \lambda_i|: i \neq j\}$:

$$\gamma_1 \ge 0.9175, \quad \gamma_2 \ge 0.1776, \quad \gamma_3 \ge 0.1175.$$
 (47)

Die Abschätzung (35) liefert dann $|\varrho_j - \lambda_j| \leq \eta_j^2/\gamma_j$, mit den angegebenen unteren Schranken (47) also

$$|arrho_1 - \lambda_1| \leq 0.0071\,, \qquad |arrho_2 - \lambda_2| \leq 0.038\,, \qquad |arrho_3 - \lambda_3| \leq 0.0043$$

und damit deutlich bessere Schranken als (45).

Wir wenden uns abschließend der zweiten der oben gestellten Fragen zu, nämlich wie u mit ||u|| = 1 zu wählen ist, damit $\{\mu, u\}$ bei gegebener Eigenwertnäherung μ ein möglichst gutes Eigenpaar von A im Sinne eines minimalen Residuums wird.

13.1.10. Aussage. Die Matrix $A \in S^{n,n}$ besitze die Eigenwerte $\{\lambda_i\}$, und $\{u^i\}$ sei eine orthonormale Basis zugehöriger Eigenvektoren. Es gelte $\mu \in \mathbb{R}$, und $j \in \{1, ..., n\}$ sei durch

$$\min\{|\lambda_i - \mu| : i = 1, ..., n\} = |\lambda_j - \mu|$$
(48)

definiert. Dann hat die Aufgabe

$$\|A\boldsymbol{u} - \boldsymbol{\mu}\boldsymbol{u}\| \to \text{Minimum! bei } \|\boldsymbol{u}\| = 1$$
 (49)

die Lösung $u = u^{j}$, und für das minimale Residuum gilt

$$\|A\boldsymbol{u}^{j} - \boldsymbol{\mu}\boldsymbol{u}^{j}\| = |\boldsymbol{\lambda}_{j} - \boldsymbol{\mu}|.$$
⁽⁵⁰⁾

Beweis. Mit (3) und der Transformation u = Uz läßt sich (49) äquivalent in der Form

$$\|Az - \mu z\|^2 = \sum_{i=1}^n (\lambda_i - \mu)^2 z_i^2 \to \text{Minimum! bei } \|z\|^2 = \sum_{i=1}^n z_i^2 = 1$$
(51)

schreiben, aus der die angegebene Lösung sofort abgelesen werden kann. 🗌

Die Bestimmung eines optimalen $u = u^i$ läuft also darauf hinaus, einen normierten Eigenvektor zu demjenigen Eigenwert λ_i zu finden, der am dichtesten an μ liegt. Im Abschnitt 13.4 werden wir zeigen, daß diese Aufgabe näherungsweise sehr effektiv mittels der sog. inversen Iteration gelöst werden kann.

Übungsaufgaben

Ü 13.1.1. Es wird die Matrix $A = \begin{pmatrix} \alpha & \beta \\ \beta & \gamma \end{pmatrix} \in S^{2,2}$ und deren charakteristisches Polynom

$$\det \left(\boldsymbol{A} - \lambda \boldsymbol{I} \right) = \lambda^2 + c_1 \lambda + c_0, \qquad c_1 := -(\alpha + \gamma), \qquad c_0 := \alpha \gamma - \beta^2$$

betrachtet.

(i) Man berechne $\{c_1, c_0\}$ in dezimaler Gleitpunktarithmetik mit acht Mantissenstellen für $\alpha := \gamma := 1$, $\beta := 7 \times 10^{-5}$ und bestimme $\{\lambda_1, \lambda_2\}$ als Nullstellen des so errechneten charakteristischen Polynoms. Man vergleiche die Ergebnisse mit den exakten Eigenwerten

$$\lambda_{\max}^* = 1.00007, \quad \lambda_{\min}^* = 0.99993.$$

(ii) Man berechne λ_1 , λ_2 für dieselben Daten und in derselben Arithmetik nach dem folgenden Algorithmus

$$V: \xi := \alpha + \gamma, \quad \delta := \alpha * \gamma - \beta^2, \quad \omega := \sqrt{(\alpha - \gamma)^2 + 4 * \beta^2}$$

if $\xi = 0$ then $[\lambda_{\max} := \omega/2; \lambda_{\min} := -\lambda_{\max}]$
else $[\lambda_{\max} := \operatorname{sgn}(\xi) (|\xi| + \omega)/2, \lambda_{\min} := \delta/\lambda_{\max}]$

Man überprüfe, daß (V) in exakter Arithmetik die Eigenwerte von A liefert, wobei $|\lambda_{\max}| \ge |\lambda_{\min}|$ ist.

(iii) Man zeige, daß die nach (ii) berechneten Größen einen gemäß

$$\begin{split} |\lambda_{\max} - \lambda_{\max}^*| &\leq \nu(5 \; |\lambda_{\max}^*|) = 5\nu \, ||A||_2, \\ |\lambda_{\min} - \lambda_{\min}^*| &\leq \nu(7 \; |\lambda_{\min}^*| + |\lambda_{\max}^*|) \leq 8\nu \, ||A||_2 \end{split}$$

beschränkten Fehler haben.

Ü 13.1.2. Es sei $A = U + iV \in \mathbb{C}^{n,n}$, $U, V \in \mathbb{R}^{n,n}$, eine *Hermitesche Matrix*, d. h., es gelte $A^H = (A)^{\mathsf{T}} = A$, vgl. 1.2. Der Matrix $A \in \mathbb{C}^{n,n}$ werde die Matrix

$$\boldsymbol{B} := \begin{pmatrix} \boldsymbol{U} \mid -\boldsymbol{V} \\ \boldsymbol{V} \mid \boldsymbol{U} \end{pmatrix} \in \mathbb{R}^{2n, 2n}$$
(52)

zugeordnet. Man zeige:

(i) U ist symmetrisch, V ist schiefsymmetrisch, B ist symmetrisch.

(ii) Für $u, v \in \mathbb{R}^n$ werde $z := \left(\frac{u}{v}\right)$ und $z' := \left(\frac{-v}{u}\right)$ definiert, und es sei $\lambda \in \mathbb{R}$. Dann ist $\{\lambda, z\}$ genau dann Eigenpaar von B, wenn (λ, z') ein Eigenpaar ist.

(iii) Für $\lambda \in \mathbf{R}$ und $u, v \in \mathbf{R}^*$ ist $\{\lambda, u + iv\}$ genau dann Eigenpaar von A, wenn $\{\lambda, z\}$ und $\{\lambda, z'\}$ Eigenpaare von B sind.

Das Eigenwertproblem der komplexen Hermiteschen Matrix A kann also auf das Eigenwertproblem der reellen symmetrischen Matrix B doppelter Dimension zurückgeführt werden. Vom Aufwand her ist jedoch die direkte Lösung des Eigenwertproblems für A in komplexer Arithmetik i. allg. nur halb so teuer wie bei Übergang zu B, da Speicherplatz und Operationszahlen wie $\sim K_1 n^2$ bzw. $\sim K_2 n^3$ von n abgängen.

Ü 13.1.3. Man zeige: Die Aufgabe (49) ist äquivalent zu

 $\|\delta A\| \to \text{Minimum!}$

bei

$$\delta \boldsymbol{A} \in \mathbf{S}^{n,n}, \quad \boldsymbol{u} \in \mathbf{R}^n \quad \text{mit} \quad (\boldsymbol{A} + \delta \boldsymbol{A}) \, \boldsymbol{u} = \mu \boldsymbol{u} \quad \text{und} \quad \|\boldsymbol{u}\| = 1, \tag{53}$$

und die minimalen Zielfunktionswerte stimmen überein.

13.2. Das Jacobi-Verfahren

Das Jacobi-Verfahren ist das älteste numerische Verfahren zur Lösung des symmetrischen Eigenwertproblems und zugleich das erste "moderne" im Sinne der Eignung für Computerrechnung. Obwohl es inzwischen effektivere Algorithmen gibt, findet dieses Verfahren wegen seiner Eleganz, Einfachheit und numerischen Gutartigkeit auch heute noch seine Anwendungsgebiete.

A. Das Basisverfahren

Die Grundidee des Verfahrens besteht darin, die Matrix $A \in S^{n,n}$ durch eine Folge orthogonaler Ähnlichkeitstransformationen

$$A_1 := A, \qquad A_{k+1} := G_k^{\mathsf{T}} A_k G_k \qquad (k = 1, 2, ...)$$
 (1)

mittels geeignet gewählter Drehungsmatrizen $G_k := G_{pq}$ sukzessive auf Diagonalform zu bringen. Aus (1) folgt

$$A_{k+1} = V_{k+1}^{\mathsf{T}} A V_{k+1} \tag{2}$$

mit den orthogonalen Transformationsmatrizen

$$V_1 := I, \qquad V_{k+1} := G_1 \cdots G_{k-1} G_k = V_k G_k,$$
(3)

insbesondere sind alle A_k wie A symmetrisch und haben dieselben Eigenwerte $\{\lambda_i\}$.

Die im k-ten Schritt verwendete Givens-Drehung hängt von den Indizes p, q und den Drehungsparametern c, s mit $c^2 + s^2 = 1$ ab, siehe 3.3.B. Diese Größen sollen so festgelegt werden, daß $A_{k+1} = (a_{ij}^{(k+1)})$ einer Diagonalmatrix möglichst nahe kommt. Als Maß für die Abweichung bietet sich in natürlicher Weise die Zahl

$$\omega_{\mathbf{k}} := \omega(\mathbf{A}_{\mathbf{k}}) := \sum_{i < j} [a_{ij}^{(k)}]^2 \tag{4}$$

an. Wenn

$$M_k := \text{diag}(a_{11}^{(k)}, \dots, a_{nn}^{(k)}) \text{ und } R_k := A_k - M_k$$
 (5)

gesetzt wird, ergibt sich aus (2) und (4)

$$\boldsymbol{M}_{k} = \boldsymbol{A}_{k} - \boldsymbol{R}_{k} = \boldsymbol{V}_{k}^{\mathsf{T}} (\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A}_{k}) \boldsymbol{V}_{k}$$

$$\tag{6}$$

 $_{\rm mit}$

$$\delta A_k := -V_k R_k V_k^{\mathsf{T}} \in \mathsf{S}^{n,n}, \qquad \|\delta A_k\|_F = \|R_k\|_F = \sqrt{2\omega_k}.$$
⁽⁷⁾

Die Diagonalelemente von A_k sind also die exakten Eigenwerte der gestörten Matrix $A + \delta A_k$, und die Spalten von V_k sind die zugehörigen Eigenvektoren. Nach 13.1.3 gilt dann

$$|a_{ii}^{(k)} - \lambda_{l(i)}| \leq \sqrt{\sum_{i=1}^{n} (a_{ii}^{(k)} - \lambda_{l(i)})^2} \leq \|\delta A_k\|_F = \sqrt{2\omega_k}.$$
(8)

Für jedes k ist dabei $\{l(1), \ldots, l(n)\} = \{l(1, k), \ldots, l(n, k)\}$ eine Permutation derart, daß die $\{\lambda_{l(i)}\}$ in derselben Weise wie die $\{a_{ii}^{(k)}\}$ geordnet sind. Die $\{a_{ii}^{(k)}\}$ approximieren also die entsprechend geordneten Eigenwerte von A mit einem durch $\sqrt{2\omega_k}$ beschränkten absoluten Fehler.

Mit dem Kriterium (4) lautet die Aufgabe dann, ein $G_k = G_{pq}(c, s)$ so zu wählen, daß

$$\omega_k - \omega_{k+1} =: \Phi_k(p, q, c, s) \to \text{Maximum!}$$
(9)

bei

$$1 \leq p < q \leq n$$
 und $c^2 + s^2 = 1$

gilt. Zur Vereinfachung der Schreibweise lassen wir wie üblich den Index k weg und kennzeichnen die zum Index k + 1 gehörenden Größen durch einen Querstrich, also $A = A_k$, $\bar{A} = A_{k+1}$ usw. Bei der Berechnung von $\bar{A} = G_{pq}^{\mathsf{T}} A G_{pq}$ aus A ändern sich nur die Elemente der Zeilen und Spalten p und q von A, und zwar nach den Vorschriften

$$\begin{split} \bar{a}_{pp} &:= c^2 a_{pp} - 2cs a_{pq} + s^2 a_{qq}, \\ \bar{a}_{qq} &:= s^2 a_{pp} + 2cs a_{pq} + c^2 a_{qq}, \\ \bar{a}_{pq} &:= \bar{a}_{qp} = cs (a_{pp} - a_{qq}) + (c^2 - s^2) a_{pq} \end{split}$$
(10)

für die Elemente in Kreuzungspunkten der Zeilen und Spalten p, q und nach

$$\begin{split} \bar{a}_{ip} &:= \bar{a}_{pi} := ca_{ip} - sa_{iq}, \\ \bar{a}_{iq} &:= \bar{a}_{qi} := sa_{ip} + ca_{iq} \end{split}$$

$$(i \neq p, q) \tag{11}$$

für die übrigen Elemente. Wegen der Orthogonalität von G_{pq} folgt aus (11)

$$(\bar{a}_{ip})^2 + (\bar{a}_{iq})^2 = (a_{ip})^2 + (a_{iq})^2 \qquad (i \neq p, q).$$

Die Quadratsumme der Nichtdiagonalelemente wird also nur durch die Elemente in den Positionen $\{p, q\}$ und $\{q, p\}$ verändert, so daß die Änderung Φ durch

$$\Phi(p, q, c, s) = \omega_k - \omega_{k+1} = \omega(A) - \omega(\bar{A}) = (a_{pq})^2 - (\tilde{a}_{pq})^2$$
(12)

gegeben ist. Bei festem $\{p, q\}$ ist die Reduktion von ω_k maximal, wenn $\{c, s\}$ so gewählt wird, daß

$$\bar{a}_{pq} = 0 \tag{13}$$

gilt. Dies führt auf

$$\omega(\mathbf{A}) = \omega(\mathbf{A}) - (a_{pq})^2. \tag{14}$$

Eine Drehung G_{pq} , die (13) realisiert, heißt *Jacobi-Drehung* in der durch $\{p, q\}$ festgelegten Ebene. Die zugehörigen Parameter $\{c, s\}$ müssen dazu so gewählt werden, daß

$$cs(a_{pp} - a_{qq}) + (c^2 - s^2) a_{pq} = 0$$
⁽¹⁵⁾

gilt. Wir setzen im folgenden $a_{pq} \neq 0$ voraus, denn andernfalls läßt sich $\omega(A)$ nach (12) überhaupt nicht reduzieren. Dann kann in (15) nicht c = 0 gelten, so daß

$$t := s/c \tag{16}$$

als neue Größe eingeführt werden kann. Damit geht (15) in die quadratische Gleichung

$$t^2 + 2\delta t - 1 = 0 \tag{17}$$

 $_{\rm mit}$

$$\delta := \frac{a_{qq} - a_{pp}}{2a_{pq}} \tag{18}$$

über. Die betragskleinste Lösung von (17) kann nach der Vorschrift

$$t := \begin{cases} 1/(|\delta| + \sqrt{1 + \delta^2}) & \text{für } \delta \ge 0, \\ -1/(|\delta| + \sqrt{1 + \delta^2}) & \text{für } \delta < 0 \end{cases}$$
(19)

numerisch stabil berechnet werden, siehe Ü 13.2.3. Aus (16) folgt dann

$$c := 1/\sqrt{1+t^2}, \quad s := tc,$$
 (20)

wenn c positiv gewählt wird. Unter Verwendung von t, s und

$$\tau := s/(1+c) \tag{21}$$

lassen sich die Transformationsformeln (10), (11) in der in exakter Arithmetik äquivalenten, für Computerrechnung jedoch wesentlich günstigeren Korrekturform

$$\begin{split} \bar{a}_{pp} &:= a_{pp} - t a_{pq}, \\ \bar{a}_{qq} &:= a_{qq} + t a_{pq}, \\ \bar{a}_{pq} &:= \bar{a}_{qp} := 0 \end{split}$$

bzw.

$$egin{aligned} ar{a}_{ip} &:= ar{a}_{pi} := a_{ip} - s(a_{iq} + au a_{ip}), \ ar{a}_{iq} &:= ar{a}_{qi} := a_{iq} + au(a_{ip} + ar{a}_{ip}) \end{aligned}$$
 $(i = p, q) \tag{23}$

schreiben. Gegenüber der sonst üblichen Darstellung

$$\bar{a}_{iq} = \bar{a}_{qi} = a_{iq} + s(a_{ip} - \tau a_{iq})$$

führt die zweite Gleichung (23) auf die Einsparung einer Multiplikation, vgl. Ü 3.3.6.

Wenn die Wahl des Pivotelementes $a_{pq} \neq 0$ noch offengelassen wird, ergibt sich das folgende Basisverfahren.

13.2.1. Jacobi-Basisverfahren

- S0 (Initialisierung): Setze $A_1 := A$, k := 1, berechne $\omega_1 := \omega(A)$, wähle $\varepsilon > 0$ S1 (Abbruchtest): if $2\omega_k \leq \varepsilon^2$ then stop S2 (Pivotwahl): Wähle p = p(k), q = q(k) mit $1 \leq p < q \leq n$ und $a_{pq}^{(k)} \neq 0$. S3 (Bestimmung der Drehungsparameter): Bestimme $t = t_k$, $c = c_k$, $s = s_k$, $\tau = \tau_k$ nach (18) bis (21)
- S4 (Berechnung von A_{k+1}): Berechne $A_{k+1} := \tilde{A}$ gemäß (22), (23) aus $A := A_k$ S5 (Berechnung von ω_{k+1}): Setze $\omega_{k+1} := \omega_k [a_{pq}^{(k)}]^2$ S6 (Vorbereitung des neuen Schrittes): Setze k := k + 1, goto S1

- Aufwand: [$\sim n(4 \text{ ops} + 3 \text{ opm}) + 2 \text{ opr}$] pro Schritt

13.2.2. Bemerkung. (i) Das Verfahren kann in situ – und zwar auf dem Platz nur eines Dreiecks einschließlich der Diagonalen - realisiert werden.

(ii) Wenn der Abbruchtest S1 nach k Schritten erfüllt ist, gilt wegen (8)

$$|\lambda_{l(i)} - a_{ii}^{(k+1)}| \le \varepsilon.$$
⁽²⁴⁾

Im allgemeinen sollte $\varepsilon \ge \nu \|A\|_F$ gewählt werden, da eine höhere Genauigkeit ohnehin nicht erreichbar ist, vgl. 13.1.B. Mit feineren Abbruchtests läßt sich unter Umständen für die betragskleinen Eigenwerte eine höhere relative Genauigkeit erzielen.

(iii) Wenn ω_k in die Größenordnung von $k\nu\omega_1$ kommt, sollte ω_{k+1} nicht durch Aufdatierung von ω_k , sondern direkt aus A_{k+1} berechnet werden. Es genügt, dies nur einmal für ein geeignetes $k = k_0$ durchzuführen.

(iv) Im Konvergenzfall $\omega_k \to 0$ gilt $|a_{ij}^{(k)}| \to 0$ $(i \neq j)$. Nach (18) wird dann $|\delta| = |\delta_k|$ häufig groß sein. Wegen

$$|\tau| \le \frac{|s|}{1 + \sqrt{2}/2} < |s| \le |t| \le \frac{1}{1 + |\delta|}$$
(25)

sind dann τ_k , s_k und t_k betragsklein, so daß die Korrekturformeln (22), (23) wesentliche Genauigkeitsvorteile gegenüber den naiven Realisierungen (10), (11) bringen, vgl. U 13.2.4. Der bei der Berechnung der Diagonalelemente gemäß (22) erzeugte Rundungsfehler kann weiter verkleinert werden, wenn die auftretenden Korrekturen $\pm ta_{pq}$ nicht sofort zu a_{pp} bzw. a_{qq} geschlagen, sondern etwa über N = n(n-1)/2aufeinanderfolgende Schritte in einem Hilfsfeld akkumuliert und erst dann insgesamt addiert werden.

(v) Wenn das vollständige Eigenwertproblem zu lösen ist, d. h., wenn neben den $\{\lambda_i\}$ auch die zugehörigen Eigenvektoren $\{u^i\}$ berechnet werden sollen, müssen die

369

Transformationsmatrizen V_k gemäß (3) mit berechnet werden, vgl. (6). Dies ist auf dem Platz einer (n, n)-Matrix in zu (23) analoger Weise möglich und kostet zusätzliche $\sim n(3 \text{ ops} + 3 \text{ opm})$ pro Schritt, vgl. Ü 3.3.6. \Box

B. Pivotisierungsstrategien, Konvergenz, Fehleranalyse

Zur Realisierung des Basisverfahrens 13.2.1 muß der Teilschritt S 2 noch spezifiziert werden, d. h., es muß festgelegt werden, welche Pivotindizes p = p(k), q = q(k) mit $a_{pq}^{(k)} \neq 0$ im k-ten Schritt gewählt werden sollen. Aus (14) ist abzulesen, daß die bereits von JACOBI verwendete Maximalpivotwahl

$$\{p, q\} \text{ so, da} \ |a_{pq}^{(k)}| = \max\{|a_{ij}^{(k)}| \colon 1 \le i < j \le n\}$$
(26)

die größte Reduktion von ω_k bringt, d. h., das durch (26) festgelegte Paar $\{p, q\}$ und die zugehörigen Jacobi-Parameter $\{c, s\}$ lösen (9). Es gilt dann

$$\omega_{k} = \omega(A_{k}) \leq N(a_{pq}^{(k)})^{2}, \qquad N := n(n-1)/2,$$
(27)

also

$$\omega_{k+1} = \omega_k - (a_{pq}^{(k)})^2 \le (1 - 1/N) \, \omega_k = \varkappa \omega_k, \qquad \varkappa := 1 - 1/N < 1.$$
(28)

Das durch (26) festgelegte Verfahren wird klassisches Jacobi-Verfahren genannt. Die maximale Reduktion von ω_k muß bei diesem Verfahren jedoch mit dem Aufwand von $N \sim n^2/2$ Lese- und Vergleichsoperationen pro Schritt zur Suche von $a_{pq}^{(k)}$ bezahlt werden. Es entsteht daher die Frage, ob sich dieser Pivotisierungsaufwand reduzieren läßt.

Als billigste Variante bietet sich an, das Pivot zyklisch alle Positionen $\{i, j\}, i < j$, durchlaufen zu lassen, etwa spaltenweise in der durch

$$\{1, 2\}, \{1, 3\}, \{2, 3\}, \dots, \{j - 2, j - 1\}, \{1, j\}, \{2, j\}, \dots, \{j - 1, j\}, \{1, j + 1\}, \dots, \{n - 1, n\}$$
(29)

gegebenen Reihenfolge. Die zu den N Pivotpositionen (29) gehörenden Jacobi-Schritte werden ein Zyklus (engl. "sweep") genannt, und das zugehörige Verfahren heißt zyklisches Jacobi-Verfahren. Selbstverständlich wird dabei die k-te Drehung übergangen, wenn $a_{pq}^{(k)} = 0$ ist. Wir bemerken dabei ausdrücklich, daß nach einem Zyklus keineswegs alle Nichtdiagonalelemente annulliert sind, denn die erzeugten Nullen werden in den nachfolgenden Schritten i. allg. zerstört.

Als Nachteil des zyklischen Jacobi-Verfahrens muß angesehen werden, daß ein Schritt mit einem Pivot, für das $0 < (a_{pq}^{(k)})^2 \ll \omega_k$ gilt, trotz des Aufwands von O(n)Operationen keine bemerkbare Reduktion von ω_k bringt. Durch Einführung eines sog. Schwellwertes (engl. "threshold") für die zulässige Pivotgröße kann dieser Effektivitätsverlust, der vor allem in den ersten Zyklen auftritt, vermieden werden. Bei einem solchen Schwellwert-Jacobi-Verfahren werden die Nichtdiagonalelemente wie beim zyklischen Verfahren etwa in der Reihenfolge (29) durchlaufen, aber die Drehung wird nur ausgeführt, wenn der Pivotbetrag nicht unter dem Schwellwert liegt. Falls das quadratische Mittel $\sqrt{\omega_k/N}$ der Nichtdiagonalelemente als Schwellwert gewählt wird, führt dies auf die Festlegung

$$\{p, q\} \text{ erstes Folgepaar von } \{p(k-1), q(k-1)\} \text{ im Zyklus (29)}$$

mit $(a_{pg}^{(k)})^2 \ge \omega_k / N.$ (30)

Die Abschätzungen (27), (28) sind dann ebenfalls erfüllt. Für Hinweise auf andere Schwellwertstrategien siehe B 13.2.

Die folgende Aussage zeigt, daß (28) die Konvergenz von A_k gegen eine Diagonalmatrix M impliziert.

13.2.3. Aussage. Das Verfahren 13.2.1 werde in exakter Arithmetik mit $\varepsilon = 0$ und einer Pivotwahl durchgeführt, für die (28) mit einem $\varkappa < 1$ erfüllt ist, also etwa als klassisches oder Schwellwert-Jacobi-Verfahren. Dann gilt

$$\lim A_k = M = \operatorname{diag}\left(\lambda_{m(1)}, \dots, \lambda_{m(n)}\right) \tag{31}$$

mit einer geeigneten Permutation $\{m(i)\}$, und die Konvergenzgeschwindigkeit ist durch die Ungleichungen

$$\begin{aligned} &i \neq j, \ k \ge 1: \ |a_{ij}^{(k)}| \le \|\boldsymbol{R}_k\|_F \\ &i = j, \ k \ge k_0: \ |a_{ii}^{(k)} - \lambda_{m(i)}| \le \|\boldsymbol{M}_k - \boldsymbol{M}\|_F \end{aligned} \right\} \le \sqrt{2\omega_k} \le (\sqrt{\varkappa})^{k-1} \sqrt{2\omega_1} \quad (32) \end{aligned}$$

charakterisiert, wobei \mathbf{R}_k und \mathbf{M}_k wie in (5) und k_0 ein genügend hoher Index ist.

Beweis. Für $i \neq j$ folgt (32) sofort aus (28). Für i = j setzen wir $\mu_k := a_{ii}^{(k)}$. Aus (22), (25) und (32), Fall $i \neq j$, ergibt sich dann mit $\varrho := \sqrt{\varkappa}$

$$|\mu_{k+1} - \mu_k| \leq |t_k| |a_{p(k)q(k)}^{(k)}| \leq e^{k-1} \sqrt[\gamma]{2\omega_1},$$

also für jedes r = 0, 1, ...

$$\begin{aligned} |\mu_{k+r+1} - \mu_k| &\leq |\mu_{k+r+1} - \mu_{k+r}| + \dots + |\mu_{k+1} - \mu_k| \leq [\varrho^{k+r-1} + \dots + \varrho^{k-1}] \sqrt[n]{2\omega_1} \\ &\leq \varrho^{k-1} [1 + \varrho + \varrho^2 + \dots] \sqrt[n]{2\omega_1} = [\varrho^{k-1}/(1 - \varrho)] \sqrt[n]{2\omega_1}. \end{aligned}$$

Die Folge $\{a_{ii}^{(i)}\} = \{\mu_k\}$ ist somit eine Cauchy-Folge und daher konvergent gegen ein $\mu =: \mu_i$, d. h., es gilt (31) mit $M = \text{diag}(\mu_i)$. Die μ_i müssen dann die Eigenwerte von A sein, so daß $\mu_i = \lambda_{m(i)}$ mit einer geeigneten Permutation $\{m(i)\}$ gilt. Für genügend großes k kann daher die in (8) vorkommende und vom Index k abhängende Permutation $\{l_k(i)\}$ identisch zu $\{m(i)\}$ gewählt werden, so daß (32) auch für i = j gilt. Man beachte, daß aus (8) allein nicht die Konvergenz der Diagonalelemente folgt, da diese beim Übergang von k zu k + 1 z. B. ihre Plätze wechseln können, vgl. Ü 13.2.2. \Box

13.2.4. Bemerkung. (i) Ungleichung (32) besagt, daß die Elemente von A_k mindestens linear mit dem Konvergenzfaktor $\sqrt[]{\varkappa} < 1$ konvergieren. Dies ist eine theoretisch nützliche, praktisch aber schwache Aussage, da \varkappa nahe bei 1 liegt. Soll etwa aus (32) abgeschätzt werden, nach wieviel vollen Zyklen r bzw. Schritten k = Nr das Abbruchkriterium $2\omega_{k+1} \leq \varepsilon^2 := \varepsilon_{rel}^2 ||A||_F^2$ mit $\varepsilon_{rel} \geq r$ erfüllt ist, so ergibt sich wegen $2\omega_{k+1} \leq 2\omega_1 \varkappa^k \leq ||A||_F^2 \varkappa^k$ die Bedingung $\varkappa^k \leq \varepsilon_{rel}^2$, also $k = Nr \leq (-2 \ln \varepsilon_{rel})/(-\ln \varkappa)$ und wegen $-\ln \varkappa \geq 1 - \varkappa = 1/N$ schließlich

$$r = k/N \ge 2\ln\left(1/\varepsilon_{\rm rel}\right). \tag{33}$$

Im Fall $\varepsilon_{\rm rel} := 10^{-6}$ bzw. 10^{-12} liefert dies die Schranke $r_{\rm max} = 28$ bzw. 55 für die Anzahl der benötigten Zyklen. Umfangreiche numerische Experimente zeigen jedoch, daß auch für sehr kleines $\varepsilon_{\rm rel}$ und ν der Abbruchtest fast immer nach vier und praktisch nie später als nach zehn Zyklen erfüllt ist, also

$$r = k/N \le 4 \dots 10 \tag{34}$$

unabhängig von ε gilt.

(ii) Das in (i) beschriebene reale Konvergenzverhalten ist ein typisches Kennzeichen sog. überlinearer Konvergenz, die in Form der asymptotisch quadratischen Konvergenz beim Jacobi-Verfahren in der Tat vorliegt: Falls $\eta_k := \max \{|a_{ij}^{(k)}|: i < j\}$ genügend klein, also k genügend groß ist, gilt

$$\eta_{k+N} \le C(\eta_k)^2 \tag{35}$$

mit einer Konstanten C > 0. Der Beweis dieser Aussage ist kompliziert und kann hier nicht gegeben werden, siehe B 13.2; die wesentlichen Ideen sind in Ü 13.2.5 zu finden. \Box

Bei Computerrechnung und Abbruch nach k = Nr Schritten wird der durch (32) beschriebene Abbruchfehler durch die Rundungsfehler überlagert. Eine detaillierte Fehleranalyse zeigt, daß zu der berechneten Matrix A_{k+1} eine Störung $\delta \bar{A}_{k+1} \in S^{n.n}$ mit

$$A_{k+1} = V_{k+1}^{\mathsf{T}}(A + \delta \bar{A}_{k+1}) V_{k+1}, \quad \|\delta \bar{A}_{k+1}\|_F \leq \nu F \|A\|_F, \quad F \leq 18n^{3/2}r$$
(36)

existiert, wobei $V_{k+1} := G_1^*G_2^* \cdots G_k^*$ exakt orthogonal und G_k^* die der berechneten Matrix A_k zugeordnete exakte Jacobi-Drehung ist. Auf Grund der speziellen Struktur des Verfahrens und bei sorgfältiger Implementierung gemäß 13.2.1 und 13.2.2 (iv) ist jedoch die äquivalente Störung wesentlich kleiner als durch (36) ausgedrückt wird. Praktisch gilt fast immer

$$F \leq k/n \sim (n/2) r$$
, und meistens ist sogar $F \leq 10$, (37)

siehe B 13.2. Wenn das berechnete A_{k+1} analog zu (5), (6) dargestellt wird, ergibt sich aus (36)

$$M_{k+1} = A_{k+1} - R_{k+1} = V_{k+1}^{\mathsf{T}} (A + \delta \bar{A}_{k+1} - V_{k+1} R_{k+1} V_{k+1}^{\mathsf{T}}) V_{k+1}$$

=: $V_{k+1}^{\mathsf{T}} (A + \delta \hat{A}_{k+1}) V_{k+1}$ (38)

mit

$$\|\delta \hat{A}_{k+1}\|_{F} \leq \|\delta \bar{A}_{k+1}\|_{F} + \|R_{k+1}\|_{F} \leq \nu F \|A\|_{F} + \sqrt{2\omega_{k+1}} \leq (\nu F + \varepsilon_{rel}) \|A\|_{F}.$$
(39)

Dies führt auf das folgende Ergebnis, bei dem sowohl Abbruchfehler als auch Rundungsfehler erfaßt sind:

13.2.5. Fehleranalyse. Das Jacobi-Verfahren werde in Computerarithmetik mit $\varepsilon = \varepsilon_{\text{rel}} \|A\|_F \ge \nu \|A\|_F$ und einer Pivotwahl, für die (28) gilt, durchgeführt. Dann bricht das Verfahren nach k = Nr Schritten ab, wobei $r \le 4...10$ ist, und die

berechneten Diagonale
lemente $\mu_i=a_{ii}^{(k+1)}$ von A_{k+1} sind die exakten Eigenwerte
einer Matrix

$$A + \delta A \quad \text{mit} \quad \delta A \in \mathbf{S}^{n,n}, \qquad \|\delta A\|_F \leq (\nu F + \varepsilon_{\text{rel}}) \|A\|_F \tag{40}$$

und F nach (36) bzw. (37).

13.2.6. Bemerkung. (i) Aus (40) folgt nach 13.1.3 die zu (8) analoge Abschätzung

$$|a_{ii}^{(k+1)} - \lambda_{l(i)}| \leq \sqrt{\sum_{i=1}^{n} (a_{ii}^{(k+1)} - \lambda_{l(i)})^2} \leq \|\delta A\|_F \leq (\nu F + \varepsilon_{\text{rel}}) \|A\|_F$$
(41)

mit einer geeigneten Permutation $\{l(i)\}$, die der Ordnung der $\{a_{ii}^{(k+1)}\}$ entspricht.

(ii) Für $\varepsilon_{rel} := v$ besagt (40), da β das Jacobi-Verfahren zur Berechnung der Eigenwerte von A ein numerisch gutartiger Proze β ist, und (41) drückt die numerische Stabilität aus.

(iii) Aus (38) folgt, daß die Spalten v^j von $V := V_{k+1}$ die zu den $\mu_j = a_{jj}^{(k-1)}$ gehörenden orthonormalen Eigenvektoren von $A + \delta A$ sind. Wenn $\tilde{V} := \tilde{V}_{k+1}$ die gemäß (3) berechnete Näherung für V bezeichnet, gilt

$$\|\tilde{\boldsymbol{V}} - \boldsymbol{V}\|_F \leq v F_1 \quad \text{mit} \quad F_1 \leq 6n^2, \tag{42}$$

d. h., V wird durch das berechnete \tilde{V} ausreichend genau approximiert. Insbesondere ist \tilde{V} wegen

$$\|\tilde{V}^{\mathsf{T}}\tilde{V} - I\|_{2} \leq \|\tilde{V} - V\|_{2} (1 + 2 \|\tilde{V} - V\|_{2}) \leq 2\nu F_{1}$$
(43)

ausreichend orthogonal.

(iv) Falls der Eigenwert $\lambda_{l(i)} \approx a_{ii}^{(k)}$ von A einfach und genügend von den übrigen Eigenwerten getrennt ist, stellt die Spalte \tilde{v}^i von \tilde{V} nach 13.1.4 eine ausreichend gute Näherung für den exakten Eigenvektor u^i von A dar. Im Fall mehrfacher Eigenwerte oder Eigenwerthaufen, die von den übrigen Eigenwerten genügend getrennt sind, approximiert der von den entsprechenden Spalten \tilde{v}^i aufgespannte Raum den entsprechenden Eigenraum ausreichend genau.

(v) Bei Abbruch nach vier Zyklen — also $4N \sim 2n^2$ Schritten — kostet die Berechnung der Eigenwerte allein $\sim 6n^3$ opms, sofern ein Schritt mit $\sim 3n$ opms gezählt wird, was nach 13.2.1 berechtigt ist. Wenn auch die Eigenvektoren gesucht sind, muß V_k nach (3) aufdatiert werden. Bei Realisierung nach Ü 3.3.6 kostet dies ebenfalls $\sim 3n$ opms pro Schritt, also insgesamt nochmals $\sim 6n^3$ opms. Zur Berechnung des vollständigen Eigensystems werden daher $\sim 12n^3$ opms benötigt.

(vi) In Abschnitt 13.7 werden wir den **QR**-Algorithmus kennenlernen, der die Eigenwerte mit einem Neuntel, das Eigensystem mit einem Drittel des hier benötigten Aufwandes bei vergleichbarer Qualität zu berechnen erlaubt. Aus diesem Grund ist die praktische Anwendung des Jacobi-Verfahrens auf Spezialfälle beschränkt, etwa die Berechnung der Eigenwerte mit geringer Genauigkeit – als Beispiel $\varepsilon_{\rm rel} \approx 10^{-1}$ – oder die Berechnung der Eigenwerte mit voller Genauigkeit, wenn die Nichtdiagonalelemente klein sind – z. B. im Sinne von max $\{|a_{ij}|: i \neq j\} \leq \sqrt{\nu} ||A||_F$. In beiden Fällen besteht die Chance, mit ein bis zwei Zyklen zum Ziel zu kommen. Auch für kleines $n - \text{etwa } n \leq 10 - \text{kann sich das Jacobi-Verfahren gegenüber dem <math>QR$ -Algorithmus als konkurrenzfähig erweisen.

(vii) Wegen seiner Einfachheit und leichten Parallelisierbarkeit wird sich die Bedeutung des Jacobi-Verfahrens für nichtklassische Computerarchitekturen — etwa Parallelcomputer — möglicherweise erhöhen. \Box

Übungsaufgaben

Ü 13.2.1. Man zeige: Wenn der Drehwinkel φ der Jacobi-Drehung G_{pq} durch $c = \cos \varphi$, $s = \sin \varphi$ festgelegt wird, gilt $t = \tan \varphi$, $\delta = \cot (2\varphi)$, $\tau = \tan (\varphi/2)$.

Ú 13.2.2. Man überprüfe, daß die zur betragsgrößten Lösung $\hat{t} = -1/t$ von (17) gehörenden Parameter $\hat{c} := -s$, $\hat{s} := c$ eine weitere Jacobi-Drehung $G_{pq}(\hat{c}, \hat{s})$ liefern. Dabei gilt für $|\delta| \to \infty$

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$
, aber $\begin{pmatrix} \hat{c} & \hat{s} \\ -\hat{s} & \hat{c} \end{pmatrix} \rightarrow \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$.

Ü 13.2.3. Man zeige, daß die aus den Eingangsdaten $\{a_{pp}, a_{qq}, a_{pq}\}$ numerisch gemäß (18), (19), (20) berechneten Werte $\tilde{t}, \tilde{c}, \tilde{s}$ den Ungleichungen

$$|t - \tilde{t}| \leq 6 \nu |t|, \quad |c - \tilde{c}| \leq 9 \nu |c|, \quad |s - \tilde{s}| \leq 9 \nu |s|$$

$$(44)$$

genügen. Wie sollte die Berechnungsvorschrift (19) modifiziert werden, damit für großes $|\delta|$ kein Überlauf eintritt?

Ü 13.2.4. Es sei $\boldsymbol{a} := (a_{ip}, a_{iq})^{\mathsf{T}}, \, \bar{\boldsymbol{a}} := (\bar{a}_{ip}, \bar{a}_{iq})^{\mathsf{T}}$. Man zeige, daß im Fall $|s| \ll 1$ der bei der Berechnung von $\bar{\boldsymbol{a}}$ nach (11) erzeugte Rundungsfehler $\delta \bar{\boldsymbol{a}}$ durch

$$\|\delta oldsymbol{ar{a}}\| \leq 12
u \|oldsymbol{ar{a}}\|$$

beschränkt ist, während bei Verwendung von (23) eine Schranke in der Größenordnung

$$\|\delta ar{a}\| \leq
u \|ar{a}\|$$

zu erwarten ist. In ähnlicher Weise wirkt sich der Ersatz von (10) durch (22) aus.

Ü 13.2.5. Es sei min $\{|\lambda_i - \lambda_j| : i \neq j\} =: 2\gamma > 0$, und k sei so groß, daß

 $\eta_k := \max \{ |a_{ij}^{(k)}| : i < j \} \ll \gamma \le \min \{ |a_{ii}^{(1)} - a_{jj}^{(1)}| : i \neq j \}$

für alle $l \ge k$ gilt. Man überlege sich:

(i) Wegen (25), (22), (23) ändert die Jacobi-Drehung G_{pq} im k-ten Schritt außer a_{pq} alle übrigen Elemente um höchstens η_k^2/γ , so daß $\eta_{k+1} \leq \eta_k + (\eta_k)^2/\gamma \approx \eta_k$ gilt.

(ii) Wenn ab Schritt k ein voller Zyklus (29) von N Jacobi-Schritten durchlaufen wird, kann ein in diesem Zyklus annulliertes Element \bar{a}_{pq} durch höchstens n-2 nachfolgende Jacobi-Drehungen modifiziert werden, so daß

$$\eta_{k+N} \leq n(\eta_k)^2/\gamma$$

gilt.

13.3. Vektor- und Teilraumiteration

Gegenstand dieses Abschnittes sind iterative Verfahren zur Approximation derjenigen Eigenvektoren der Matrix $A \in S^{n,n}$, die zu sog. dominanten, d. h. betragsmäßig größten und von den übrigen getrennten Eigenwerten gehören. Bei der Realisierung dieser Verfahren fallen in natürlicher Weise auch Approximationen für die entsprechenden Eigenwerte an. Die Algorithmen selbst haben keinen allzu großen direkten Anwendungsbereich. Sie sind jedoch Grundlage für leistungsfähige Weiterentwicklungen wie die inverse Iteration und der **QR**-Algorithmus, aber auch für die meisten modernen Verfahren zur Lösung des Eigenwertproblems schwach besetzter Matrizen hoher Dimension.

Im folgenden seien die Eigenwerte $\{\lambda_i\}$ von A nach absteigenden Beträgen gemäß

$$|\lambda_n| \leq |\lambda_{n-1}| \leq \cdots \leq |\lambda_2| \leq |\lambda_1| \tag{1}$$

geordnet, und $\{u^j\}$ sei ein orthonormales System zugehöriger Eigenvektoren. Die Eigenwerte $\{\lambda_1, \ldots, \lambda_p\}$ heißen dann *dominant*, wenn

$$|\lambda_{p+1}| < |\lambda_p| \tag{2}$$

gilt. Den p dominanten Eigenwerten entspricht der durch die zugehörigen Eigenvektoren aufgespannte Teilraum

$$\mathscr{S}_p = \mathscr{S}_p(A) := \operatorname{span} \left\{ u^1, \dots, u^p \right\},\tag{3}$$

der *p*-ter dominanter Teilraum von A genannt wird. Offensichtlich ist $\{u^1, ..., u^p\}$ eine orthonormale Basis von \mathscr{S}_p , und es gilt dim $\mathscr{S}_p = p$. Aus $Au^j = \lambda_j u^j \in \mathscr{S}_p$ (j = 1, ..., p) folgt außerdem

$$\boldsymbol{A}\boldsymbol{\mathscr{S}}_{p} = \{\boldsymbol{y} = \boldsymbol{A}\boldsymbol{x} \colon \boldsymbol{x} \in \boldsymbol{\mathscr{S}}_{p}\} \subset \boldsymbol{\mathscr{S}}_{p}, \tag{4}$$

d. h., \mathscr{S}_p ist invariant unter der durch A vermittelten linearen Abbildung. Ein Teilraum \mathscr{S}_p , für den (4) gilt, heißt *invarianter Teilraum* von A. Wir werden später sehen, daß umgekehrt jedem invarianten Teilraum von A der Dimension p genau p Eigenwerte und Eigenvektoren von A zugeordnet werden können. Wegen der Trennung der dominanten Eigenwerte von den restlichen wird man daher erwarten können, aus einer genügend guten Approximation von \mathscr{S}_p auch genügend gute Approximationen für die Eigenpaare $\{\lambda_i, u^j\}$ (j = 1, ..., p) beschaffen zu können.

Um zu den gesuchten Approximationen für \mathcal{S}_p zu kommen, beachten wir, daß die Eigenwertzerlegung

$$A = U \Lambda U^{\mathsf{T}} = [\lambda_1 u^1 u^{1\mathsf{T}} + \dots + \lambda_p u^p u^{p\mathsf{T}}] + [\lambda_{p+1} u^{p+1} u^{p+1\mathsf{T}} + \dots + \lambda_n u^n u^{n\mathsf{T}}]$$
(5)

von A die analoge Zerlegung

$$A^{k} = U \Lambda^{k} U^{\mathsf{T}} = [\lambda_{1}^{k} u^{1} u^{1\mathsf{T}} + \dots + \lambda_{p}^{k} u^{p} u^{p\mathsf{T}}] + [\lambda_{p+1}^{k} u^{p+1} u^{p+1\mathsf{T}} + \dots + \lambda_{n}^{k} u^{n} u^{n\mathsf{T}}]$$
(6)

für die k-te Potenz A^k von A nach sich zieht. Durch die Potenzbildung wird die Dominanz der ersten p Terme mit wachsendem k immer größer, so daß A^k für genügend großes k praktisch durch die Summanden in der ersten Klammer repräsentiert wird. Wenn A^k auf einen beliebigen Vektor $v \in \mathbb{R}^n$ angewendet wird, ergibt sich daher

$$A^{\mathbf{k}}\boldsymbol{v} = \sum_{j=1}^{n} \lambda_{j}^{k} \boldsymbol{u}^{j}[\boldsymbol{u}^{j\top}\boldsymbol{v}] \approx \left[\lambda_{1}^{k} \boldsymbol{u}^{1}[\boldsymbol{u}^{1\top}\boldsymbol{v}] + \dots + \lambda_{p}^{k} \boldsymbol{u}^{p}[\boldsymbol{u}^{p\top}\boldsymbol{v}]\right],$$
(7)

sofern die Koeffizienten $u^{j\top}v$ (j = 1, ..., p) nicht sämtlich verschwinden. Für geeignet gewähltes v und genügend großes k ist also $A^{k}v$, "fast" eine Linearkombination der ersten p Eigenvektoren, d. h., $A^{k}v$ liegt "fast" in \mathscr{F}_{p} . Wenn (7) für p geeignet gewählte und linear unabhängige Vektoren $\{v\}$ gebildet wird, werden die zugehörigen Vektoren $\{A^{k}v\}$ einen Teilraum aufspannen, der \mathscr{F}_{p} für genügend großes k ausreichend gut approximiert. Die Berechnung von $A^{k}v$ erfolgt selbstverständlich rekursiv gemäß $A^{k+1}v := A(A^{k}v)$. Da $||A^{k}v||$ im Fall $|\lambda_{1}| > 1$ unbeschränkt wächst, im Fall $|\lambda_{1}| < 1$ dagegen gegen 0 geht, muß bei der praktischen Berechnung eine Normierung vorgenommen werden, um unnötigen Über- bzw. Unterlauf zu vermeiden. Weil für die Teilraumapproximation nur die Richtungen wesentlich sind, hat das auf die Approximationsgüte keinen Einfluß.

Wir beginnen mit der genauen Untersuchung des einfachsten Falles p = 1.

A. Vektoriteration

Hier genügt es, einen Vektor v zu iterieren. Das zugehörige Verfahren lautet wie folgt:

13.3.1. Basisverfahren der Vektoriteration S0 (Initialisierung): Wähle v^1 mit $||v^1|| = 1$, setze k := 1S1 (Iteration): Berechne $w^{k+1} := Av^k$ S2 (Normierung): Setze $v^{k+1} := w^{k+1}/||w^{k+1}||$ S3: Setze k := k + 1, goto S1 Aufwand pro Schritt: eine Auswertung von $Ax + \sim n(2 \text{ opm} + 1 \text{ ops}) + 1$ opr

13.3.2. Bemerkung. (i) Im Unterschied zu anderen Eigenwertverfahren wird die Matrix A in 13.3.1 nicht verändert und braucht auch explizit nicht verfügbar zu sein. Es genügt, daß eine Prozedur vorhanden ist, die zu gegebenem x das Bild y = Ax berechnet und im übrigen als "black box" behandelt werden kann. Dies ist besonders bei schwach besetzten Matrizen hoher Dimension, wie sie etwa bei der Diskretisierung von Eigenwertaufgaben bei Differentialgleichungen entstehen, vorteilhaft.

(ii) Statt der 2-Norm kann auch die ∞ -Norm zur Normierung verwendet werden, etwa gemäß

S2': Bestimme s = s(k + 1) mit $|(w^{k+1})_s| = \max \{|(w^{k+1})_i|: 1 \le i \le n\},$ setze $\omega_{k+1} := (w^{k+1})_s$ und $v^{k+1} := w^{k+1}/\omega_{k+1}.$

Bei dieser Festlegung gilt $||v^{k+1}||_{\infty} = (v^{k+1})_s = 1$. Für die Konvergenzanalyse und die Eigenwertapproximation ist jedoch die 2-Norm günstiger.

(iii) Wegen der Gültigkeit der Darstellung

 $\boldsymbol{v}^{k+1} = \boldsymbol{A}^{k} \boldsymbol{v}^{1} / \tau_{k+1}, \qquad \tau_{k+1} := \| \boldsymbol{w}^{2} \| \cdot \| \boldsymbol{w}^{3} \| \cdots \| \boldsymbol{w}^{k+1} \|$ (8)

heißt Algorithmus 13.3.1 auch *Potenzmethode* (engl. "power method", russ. "степенной метод"). In der deutschsprachigen Literatur ist auch der Name "von Mises-Iteration" gebräuchlich, siehe B 13.3.

Das Konvergenzverhalten der Vektoriteration wird durch den folgenden Satz beschrieben:

13.3.3. Satz. Es sei λ_1 der dominierende Eigenwert von A, und u^1 bezeichne den zugehörigen normierten Eigenvektor. Der Startvektor v^1 genüge der Bedingung

$$\sigma := \boldsymbol{u}^{1\mathsf{T}} \boldsymbol{v}^1 > 0. \tag{9}$$

Dann gilt für die gemäß 13.3.1 in exakter Arithmetik erzeugte Folge $\{v^k\}$

$$0 \leq \tan \varphi_{k+1} \leq \varkappa \tan \varphi_k \leq \varkappa^k \tan \varphi_1 = \varkappa^k \sqrt{1 - \sigma^2} / \sigma$$
 (10)

mit

$$\varkappa := |\lambda_2/\lambda_1| < 1 \tag{11}$$

und

$$\varphi_k := \not\triangleleft (\boldsymbol{u}^1, \vartheta_k \boldsymbol{v}^k) \in [0, \pi/2), \qquad \vartheta_k := (\operatorname{sgn} \lambda_1)^{k+1} \in \{+1, -1\}.$$
(12)

Beweis. Wir führen die zueinander komplementären und bezüglich A invarianten Teilräume $\mathscr{I}_1 := \operatorname{span} \{u^1\}$ und $\mathscr{I}_2 := \operatorname{span} \{u^2, ..., u^n\} = \mathscr{I}_1^\perp$ ein, vgl. 8.1.B. Die Zerlegung von $\vartheta_1 v^1$ in die Komponenten bezüglich \mathscr{I}_1 und \mathscr{I}_2 kann dann in der Form

$$\boldsymbol{v}^{1} = \boldsymbol{\vartheta}_{1} \boldsymbol{v}^{1} = \boldsymbol{u}^{1} \cos \varphi_{1} + \boldsymbol{z}^{1} \sin \varphi_{1} \quad \text{mit} \quad \boldsymbol{z}^{1} \in \boldsymbol{\mathcal{X}}_{2}, \qquad \|\boldsymbol{z}^{1}\| = 1,$$
(13)

geschrieben werden. Wegen (9) gilt dabei $\sigma = \cos \not\lt (\boldsymbol{u}^1, \boldsymbol{v}^1) = \cos \varphi_1 > 0$, also $0 \leq \varphi_1 < \pi/2$ und $\tan \varphi_1 = \sqrt{1 - \sigma^2}/\sigma \geq 0$. Linksmultiplikation von (13) mit sgn $(\lambda_1) \boldsymbol{A}$ führt auf

$$\operatorname{sgn}(\lambda_1) A \vartheta_1 v^1 = \vartheta_2 A v^1 = \vartheta_2 w^2 = \operatorname{sgn}(\lambda_1) \left\{ \lambda_1 u^1 \cos \varphi_1 + A z^1 \sin \varphi_1 \right\}$$

mit $Am{z}^1\inm{\mathcal{J}}_2$, also $\|m{w}^2\|\geq|\lambda_1|\cosarphi_1>0$ und somit auf die Darstellung

$$\boldsymbol{\vartheta}_2 \boldsymbol{v}^2 = \boldsymbol{\vartheta}_2 \boldsymbol{w}^2 / ||\boldsymbol{w}^2|| = \boldsymbol{u}^1 \cos \varphi_2 + \boldsymbol{z}^2 \sin \varphi_2 \quad \text{mit} \quad \boldsymbol{z}^2 \in \boldsymbol{\mathcal{Z}}_2, \quad ||\boldsymbol{z}^2|| = 1,$$
 (14)

wobei

$$\cos \varphi_2 := \frac{|\lambda_1| \cos \varphi_1}{||\boldsymbol{w}^2||}, \qquad \sin \varphi_2 := \frac{||\boldsymbol{A}\boldsymbol{z}^1|| \sin \varphi_1}{||\boldsymbol{w}^2||}, \qquad \boldsymbol{z}^2 := \frac{\operatorname{sgn}\left(\lambda_1\right) \boldsymbol{A}\boldsymbol{z}^1}{||\boldsymbol{A}\boldsymbol{z}^1||}$$

im Fall $Az^1 \neq o$; für $Az^1 = o$ setzen wir einfach $z^2 := u^2$. Aus (14) folgt

$$0 \leq arphi_2 < \pi/2 \quad ext{und} \quad 0 \leq an arphi_2 = [||Am{z}^1||/|\lambda_1|] an arphi_1.$$

Nach (13) gilt nun

$$\boldsymbol{z}^1 = \sum_{j=2}^n c_j \boldsymbol{u}^j \quad ext{mit} \quad \sum_{j=2}^n c_j^2 = 1,$$

woraus sich

$$A\boldsymbol{z}^1 = \sum_{j=2}^n \lambda_j c_j \boldsymbol{u}^j$$
, also $||A\boldsymbol{z}^1||^2 = \sum_{j=2}^n \lambda_j^2 c_j^2 \leq \lambda_2^2 \sum_{j=2}^n c_j^2 = \lambda_2^2$

ergibt. Daher ist $||Az^1||/|\lambda_1| \leq |\lambda_2/\lambda_1|$, d. h., (10) gilt für k = 1. Die Gültigkeit für alle $k \geq 1$ folgt in analoger Weise durch Induktion.

13.3.4. Bemerkung. (i) Die Charakterisierung der "Entfernung" zweier Richtungen u, v durch die Größe $|\tan(\not\prec(u, v))|$ scheint adäquat zu sein: Orthogonale Richtungen sind unendlich entfernt, während Vielfache eines Vektors die Entfernung 0, also dieselbe Richtung haben. Wir weisen darauf hin, daß der Vorzeichenfaktor ϑ_k im Verfahren 13.3.1 selbst überhaupt nicht vorkommt. Er wurde in 13.3.3 nur eingeführt, um $\vartheta_k v^k$ wie u^1 zu orientieren. Dies garantiert $0 \leq \varphi_k \leq \pi/2$ und erlaubt die einfache Formulierung von Konvergenzaussagen, siehe (ii) unten. Auf die Güte der Approximation der Eigenrichtung span $\{u^1\}$ durch v^k hat ϑ_k keinen Einfluß, denn offensichtlich ist v^k eine genau so gute Approximation wie $-v^k$. Ebenso ist die Forderung $\sigma > 0$ keine Einschränkung; im Fall $\sigma < 0$ braucht nur u^1 durch $-u^1$ ersetzt zu werden. Benötigt wird eigentlich nur $\sigma \neq 0$.

(ii) Wegen

$$\|\vartheta_k \boldsymbol{v}^k - \boldsymbol{u}^1\| = \sqrt{2(1 - \cos\varphi_k)} = 2\sin\left(\varphi_k/2\right) \le \tan\varphi_k \tag{15}$$

folgt aus (10) die Konvergenz der geeignet orientierten Vektoren $\{\pm v^k\}$ gegen u^1 .

(iii) Wenn $\vartheta_k v^k$ analog zu (13) in der Form

$$\vartheta_k \boldsymbol{v}^k = \boldsymbol{u}^1 \cos \varphi_k + \boldsymbol{z}^k \sin \varphi_k \quad \text{mit} \quad \boldsymbol{z}^k \in \mathcal{X}_2, \qquad \|\boldsymbol{z}^k\| = 1, \tag{16}$$

dargestellt wird, hat z^k für genügend großes k im Fall $|\lambda_2| > |\lambda_3|$ die Richtung von u^2 , vgl. auch (7). Es gilt dann tan $\varphi_{k+1} \approx \varkappa \tan \varphi_k$, d. h., die Schranke (10) ist realistisch; schnellere Konvergenz ist nicht zu erwarten. \Box

Die zur Eigenvektornäherung v^k gehörende beste Eigenwert
approximation ist nach 13.1.8 der Rayleigh-Quotient

$$\varrho_k := \varrho(\boldsymbol{v}^k) = \boldsymbol{v}^{k\mathsf{T}} \boldsymbol{A} \boldsymbol{v}^k = \boldsymbol{v}^{k\mathsf{T}} \boldsymbol{w}^{k+1}. \tag{17}$$

Aus (16) folgt nun

$$\boldsymbol{A}(\vartheta_k \boldsymbol{v}^k) = \boldsymbol{u}^1 \lambda_1 \cos \varphi_k + \boldsymbol{A} \boldsymbol{z}^k \sin \varphi_k \quad \text{mit} \quad \boldsymbol{A} \boldsymbol{z}^k \in \mathcal{X}_2 \perp \mathcal{S}_1 ,$$
(18)

also

$$arphi_k = arrho(artheta_k oldsymbol{v}^k) = (oldsymbol{u}^1 \cos arphi_k + oldsymbol{z}^k \sin arphi_k)^{\mathsf{T}} (oldsymbol{u}^1 \lambda_1 \cos arphi_k + oldsymbol{A} oldsymbol{z}^k \sin arphi_k)$$

= $\lambda_1 \cos^2 arphi_k + oldsymbol{z}^{k\mathsf{T}} oldsymbol{A} oldsymbol{z}^k \sin^2 arphi_k = \lambda_1 - [\lambda_1 - oldsymbol{z}^{k\mathsf{T}} oldsymbol{A} oldsymbol{z}^k] \sin^2 arphi_k.$

Wegen

$$|\lambda_1 - \boldsymbol{z^{k^{\intercal}}} A \boldsymbol{z^{k}}| \leq |\lambda_1| + \|\boldsymbol{z^{k}}\| \cdot \|A \boldsymbol{z^{k}}\| \leq |\lambda_1| + |\lambda_2| \leq 2 |\lambda_1| = 2 \|A\|$$

ergibt sich daraus die Abschätzung

$$|\varrho_k - \lambda_1| \leq 2 \|A\| \sin^2 \varphi_k \leq 2 \|A\| \tan^2 \varphi_k.$$
⁽¹⁹⁾

Die Rayleigh-Quotienten konvergieren also linear mit dem Konvergenzfaktor \varkappa^2 und damit schneller als die Eigenvektornäherungen, wo der Fehler durch tan φ_k beschränkt ist und \varkappa als Konvergenzfaktor auftritt, siehe (15). Aus der für genügend großes k gültigen Beziehung

$$\boldsymbol{w}^{k+1} = A\boldsymbol{v}^k \approx \lambda_1 \boldsymbol{v}^k \tag{20}$$

lassen sich auch einfacher berechenbare Eigenwertnäherungen ableiten. Durch Normbildung folgt $||\boldsymbol{v}^{k+1}|| \approx ||\lambda_1 \boldsymbol{v}^k|| = |\lambda_1|$, so daß

$$\gamma_{k} := \operatorname{sgn}\left(\lambda_{1}\right) \left\|\boldsymbol{w}^{k+1}\right\| \tag{21}$$

eine Näherung für λ_1 darstellt. Der Fehler ist durch

$$|\gamma_{k} - \lambda_{1}| \leq 2 \|A\| \sin^{2}(\varphi_{k}/2) \leq (\|A\|/2) \tan^{2}\varphi_{k}$$
(22)

beschränkt, siehe Ü 13.3.2, also vergleichbar mit dem von ϱ_k . Da $||w^{k+1}||$ ohnehin berechnet werden muß, braucht in (21) nur noch das Vorzeichen sgn (λ_1) bestimmt zu werden. Zu diesem Zweck lesen wir (20) komponentenweise als $(w^{k+1})_j \approx \lambda_1 (v^k)_j$. Die Quotienten

$$\omega_{\mathbf{k},j} := (\mathbf{w}^{\mathbf{k}+1})_j / (\mathbf{v}^{\mathbf{k}})_j, \qquad (\mathbf{v}^{\mathbf{k}})_j \neq 0, \tag{23}$$

stellen also ebenfalls Näherungen für λ_1 dar, und für genügend großes k ist

$$\operatorname{sgn}\left(\omega_{\boldsymbol{k},j}\right) = \operatorname{sgn}\left(\lambda_{1}\right). \tag{24}$$

Zu empfehlen ist die Festlegung von j = s(k) als Index der betragsgrößten Komponente $(\boldsymbol{v}^k)_{s(k)}$ von \boldsymbol{v}^k . Zur Approximation von λ_1 selbst sind die Quotienten $\omega_{k,j}$ nicht so gut geeignet, denn die im Fall $|(\boldsymbol{u}^1)_j| > \tan \varphi_k$ gültige Fehlerschranke

$$|\omega_{k,j} - \lambda_1| \leq 2 \|\boldsymbol{A}\| \tan \varphi_k / [|(\boldsymbol{u}^1)_j| - \tan \varphi_k]$$
⁽²⁵⁾

enthält tan φ_k statt tan² φ_k und liefert daher nur Konvergenz mit dem Faktor \varkappa , siehe Ü 13.3.3.

Bei der Computerrealisierung von 13.3.1 wird der Abbruchfehler durch die erzeugten Rundungsfehler überlagert, so daß 13.3.3 nur in der folgenden modifizierten Form gilt:

13.3.5. Fehleranalyse. Die Vektoriteration 13.3.1 werde in Computerarithmetik so durchgeführt, daß die berechneten Vektoren $\{w^{k+1}\}$ der Beziehung

$$\boldsymbol{w}^{k+1} = (\boldsymbol{A} + \boldsymbol{\delta} \boldsymbol{A}_k) \, \boldsymbol{v}^k \quad \text{mit} \quad \|\boldsymbol{\delta} \boldsymbol{A}_k\| \leq \boldsymbol{v} F \, \|\boldsymbol{A}\| \tag{26}$$

genügen, und mit \varkappa aus (11) und $\varepsilon := \nu(F + n/2)$ seien die Bedingungen

$$\kappa + 10\varepsilon < 1$$
 sowie $|\boldsymbol{u}^{\mathbf{1}\mathsf{T}}\boldsymbol{v}^{\mathbf{1}}| > \tilde{\varepsilon} := \varepsilon/(1 - \kappa - \varepsilon)$ (27)

erfüllt. Dann gilt für die den berechneten Vektoren $\{\bm{v^k}\}$ gemäß (12) zugeordneten Winkel $\{\varphi_k\}$ die Abschätzung

$$0 \leq \tan \varphi_k \leq \tau_k \quad \text{mit} \quad \tau_k \geq \tau_{k+1} \quad \text{und} \quad \lim_{k \to \infty} \tau_k = \bar{\varepsilon}.$$
(28)

Wenn k_0 ein im Sinne von tan $\varphi_{k_0}<1$ genügend großer Index ist, gilt zudem für $k \ge k_0$

$$\tan \varphi_{k+1} \leq \hat{\varkappa} \tan \varphi_k + \hat{\varepsilon} \leq \hat{\varkappa}^{k-k_0+1} \tan \varphi_{k_0} + \hat{\varepsilon}/(1-\hat{\varkappa})$$

mit $\hat{\varkappa} := (\varkappa + \varepsilon)/(1-\sqrt{2}\varepsilon) = \varkappa, \quad \hat{\varepsilon} := \varepsilon/(1-\sqrt{2}\varepsilon) = \varepsilon.$ (29)

Für alle $k \ge 1$ ist außerdem

$$\|\vartheta_k v^k - u^1\| \leq \tan \varphi_k + \nu n/2, \tag{30}$$

und die gemäß

$$arrho_k := oldsymbol{v}^{k^+ 1} \quad ext{bzw.} \quad \gamma_k := (ext{sgn } \lambda_1) \sqrt{oldsymbol{w}^{k+1} oldsymbol{w}^{k+1}}$$

berechneten Eigenwertnäherungen genügen den Abschätzungen

$$|\varrho_k - \lambda_1| \leq \|A\| \left[2\sin^2 \varphi_k + \nu(F+2n)\right] \tag{31}$$

bzw.

$$|\gamma_k - \lambda_1| \leq ||A|| |\eta_k/(1+\sqrt{1-\eta_k}) \quad \text{mit} \quad \eta_k := \sin^2 \varphi_k + 2\nu(F+n).$$
 (32)

Der Beweis verläuft ähnlich wie der zu 13.3.3 und (15), (19), (22), wobei (26) und die Ergebnisse aus 2.3.B beachtet werden müssen.

13.3.6. Bemerkung. (i) Die erste der Bedingungen (27) fordert, daß $\varkappa = |\lambda_2/\lambda_1|$ nicht zu nahe bei 1 liegt. Die zweite soll ausschließen, daß der Startvektor v^1 fast orthogonal zu u^1 ist, denn dann kann die in 13.3.5 beschriebene Konvergenz nicht gefolgert werden. Praktisch tritt jedoch für jeden Startvektor v^1 Konvergenz auf: Auch für exakt orthogonale u^1, v^1 erzeugen die Rundungsfehler fast immer nichtverschwindende Komponenten von v^k (k > 1) in Richtung von u^1 , die sich in den nachfolgenden Schritten weiter verstärken.

(ii) Für voll besetztes A ist (26) bei "normaler" Berechnung von $w^{k+1} = \operatorname{fl}(Av^k)$ mit $F = n^{3/2}$ erfüllt, vgl. 2.3.D und 2.3.15. In den meisten Anwendungen ist A jedoch schwach besetzt, so daß sich wesentlich kleinere Werte für F ergeben.

(iii) Aus (26) folgt mit der Abkürzung $\hat{v}^{k} := \vartheta_{k} v^{k}$

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}_{k})\,\hat{\boldsymbol{v}}^{k} = \gamma_{k}\hat{\boldsymbol{v}}^{k+1} = \gamma_{k}\hat{\boldsymbol{v}}^{k} + \gamma_{k}(\hat{\boldsymbol{v}}^{k+1} - \hat{\boldsymbol{v}}^{k}), \qquad (33)$$

mithin

$$\begin{split} \|\boldsymbol{A}\hat{\boldsymbol{v}}^{k} - \gamma_{k}\hat{\boldsymbol{v}}^{k}\| &= \|\boldsymbol{A}\boldsymbol{v}^{k} - \gamma_{k}\boldsymbol{v}^{k}\| \leq \|\boldsymbol{\delta}\boldsymbol{A}_{k}\| + \|\hat{\boldsymbol{v}}^{k+1} - \hat{\boldsymbol{v}}^{k}\| \, |\gamma_{k}| \\ &\leq \{\boldsymbol{v}F + \|\hat{\boldsymbol{v}}^{k+1} - \hat{\boldsymbol{v}}^{k}\|\} \, \|\boldsymbol{A}\|. \end{split}$$

Nach 13.1.7 existiert dann eine Störung $\delta \hat{A}_k \in S^{n,n}$ mit

$$\|\boldsymbol{\delta}\hat{\boldsymbol{A}}_{k}\| \leq \varepsilon_{k} \, \|\boldsymbol{A}\|, \qquad \varepsilon_{k} := \nu F + \|\hat{\boldsymbol{v}}^{k+1} - \hat{\boldsymbol{v}}^{k}\|, \tag{34}$$

so daß $\{\gamma_k, v^k\}$ Eigenpaar von $A + \delta \hat{A}_k$ ist. Dies liefert eine nützliche a-posteriori-Abschätzung, die z. B. als Abbruchkriterium verwendet werden kann. Mit (29), (30) läßt sich ε_k für $k \ge k_0$ weiter gemäß

$$\varepsilon_{k} \leq \nu F + \|\hat{\boldsymbol{v}}^{k+1} - \boldsymbol{u}^{1}\| + \|\hat{\boldsymbol{v}}^{k} - \boldsymbol{u}^{1}\| \leq \nu [F + (2F + n)/(1 - \hat{\boldsymbol{z}})] + 2\hat{\boldsymbol{z}}^{k-k_{0}} \tan \varphi_{k_{0}}$$

$$(35)$$

abschätzen. Sofern $(2F + n)/(1 - \hat{z})$ akzeptabel, also \hat{z} nicht zu nahe bei 1 ist, liegt daher für genügend großes k numerische Gutartigkeit vor. Ähnliche Aussagen gelten für $\{\varrho_k, v^k\}$, denn nach 13.1.8 ist dann $||Av^k - \varrho_k v^k|| \leq ||Av^k - \gamma_k v^k||$.

(iv) Eine ausreichende Genauigkeit läßt sich mit vertretbarem Aufwand nur erreichen, wenn \hat{x} genügend klein ist. Diese Forderung schränkt die direkte Anwendung von 13.3.1 auf ganz spezielle Aufgaben ein. Man beachte, daß durch eine Spektralverschiebung $\tilde{A} = A - \mu I$ höchstens die an den Enden des Spektrums gelegenen Eigenwerte zu dominanten gemacht werden können, vgl. auch Ü 13.3.1.

(v) Die richtig orientierten Näherungen $\{\hat{v}^k\}$ ergeben sich nach der Vorschrift

$$\hat{\boldsymbol{v}}^{k+1} := A \hat{\boldsymbol{v}}^{k}, \quad \gamma_{k} := (\operatorname{sgn} \lambda_{1}) \| \hat{\boldsymbol{v}}^{k+1} \|, \quad \hat{\boldsymbol{v}}^{k+1} := \hat{\boldsymbol{v}}^{k+1} / \gamma_{k} \\
(k = 1, 2, \ldots).$$
(36)

wobei sgn λ_1 etwa durch (24) mit j = j(k) als Index der betragsgrößten Komponente von \hat{v}^k festgelegt wird. \Box

Wir beschließen diesen Teilabschnitt mit einigen Kommentaren über Modifikationen und Erweiterungen der Vektoriteration.

13.3.7. Bemerkung. (i) Wenn λ_1 ein doppelter dominanter Eigenwert ist, d. h., wenn $\lambda_1 = \lambda_2$ und $|\lambda_1| > |\lambda_3|$ gilt, erzeugt 13.3.1 eine Folge $\{\vartheta_k v^k\}$, die mit dem Konvergenzfaktor $\varkappa = |\lambda_3/\lambda_1|$ gegen einen normierten Eigenvektor $\overline{\boldsymbol{u}}^1 \in \mathscr{S}_2 = \text{span} \{\boldsymbol{u}^1, \boldsymbol{u}^2\}$ konvergiert. Wird 13.3.1 mit einem zweiten Startwert $\tilde{\boldsymbol{v}}^1 \neq \boldsymbol{v}^1$ gestartet, so konvergieren die zugehörigen Iterierten $\{\vartheta_k \tilde{\boldsymbol{v}}^k\}$ fast immer gegen einen weiteren normierten Eigenvektor $\boldsymbol{u}^2 \in \mathscr{S}_2$, der von $\boldsymbol{\bar{u}}^1$ linear unabhängig ist. Ein nachfolgender Orthogonalisierungsschritt

$$\boldsymbol{q}^2 := \tilde{\boldsymbol{u}}^2 - [\bar{\boldsymbol{u}}^{1\mathsf{T}} \tilde{\boldsymbol{u}}^2] \, \bar{\boldsymbol{u}}^1, \qquad \bar{\boldsymbol{u}}^2 := \boldsymbol{q}^2 / \|\boldsymbol{q}^2\| \tag{37}$$

liefert dann einen Eigenvektor $\overline{\boldsymbol{u}}^2$ derart, daß { $\overline{\boldsymbol{u}}^1$, $\overline{\boldsymbol{u}}^2$ } eine orthonormale Basis von \mathscr{S}_2 bildet. Um Genauigkeitsverluste bei der numerischen Realisierung von (37) zu vermeiden, empfiehlt es sich, solche Orthogonalisierungsschritte schon in gewissen Abständen mit den Iterierten $\tilde{\boldsymbol{v}}^k$ durchzuführen, damit $\tilde{\boldsymbol{u}}^2$ hinreichend linear unabhängig von $\overline{\boldsymbol{u}}^1$ ist. In analoger Weise kann im Fall eines *p*-fachen dominanten Eigenwertes λ_1 vorgegangen werden.

(ii) Es gelte $|\lambda_1| > |\lambda_2| > |\lambda_3|$, und das Eigenpaar $\{\lambda_1, u^1\}$ sei bereits bestimmt worden. Wir betrachten jetzt die Matrix

$$\bar{\boldsymbol{A}} := \boldsymbol{A} - \lambda_1 \boldsymbol{u}^1 \boldsymbol{u}^{1\mathsf{T}} = 0 \boldsymbol{u}^1 \boldsymbol{u}^{1\mathsf{T}} + \sum_{j=2}^n \lambda_j \boldsymbol{u}^j \boldsymbol{u}^{j\mathsf{T}}, \qquad (38)$$

die die Eigenwerte $\{0, \lambda_2, ..., \lambda_n\}$ und die zugehörigen Eigenvektoren $\{u^1, u^2, ..., u^n\}$ hat. Durch die Transformation von λ_1 in 0 ist also λ_2 zum dominanten Eigenwert von \bar{A} geworden, und das Paar $\{\lambda_2, u^2\}$ kann durch Vektoriteration mit \bar{A} bestimmt werden. Dieses Deflation genannte Vorgehen ist gemäß $\overline{A} := \overline{A} - \lambda_2 u^2 u^{2\intercal}$ usw. fortsetzbar, sofern die Eigenwerte betragsmäßig getrennt sind. Bei der numerischen Realisierung dieses Prozesses müssen die verwendeten Näherungen allerdings eine hohe Genauigkeit aufweisen.

(iii) Der durch die explizite Deflation (38) erzielte Effekt kann auch implizit erreicht werden: Wenn 13.3.1 in exakter Arithmetik mit einem Startvektor v^1 , für den $\sigma = u^{1} v^{1} = 0$ gilt, durchgeführt wird, ist $v^{k} \in \mathcal{X}_{2} = \text{span} \{u^{2}, \dots, u^{n}\}$, vgl. (7). Überdies sind die v^k mit denjenigen Iterierten identisch, die sich aus v^1 mit \bar{A} statt A ergeben, d. h., v^{k} approximiert u^{2} . Bei Computerrechnung muß die Bedingung $v^k \in \mathcal{X}_2$, d. h. $v^k \perp u^1$, durch in gewissen Abständen vorgenommene Orthogonalisierungsschritte

$$q^{k} := v^{k} - [u^{1} v^{k}] u^{1}, \qquad v^{k} := q^{k} / \|q^{k}\|, \qquad (39)$$

erzwungen werden. Andernfalls würden die durch Rundungsfehler erzeugten Komponenten von v^k in Richtung u^1 schnell anwachsen und die Konvergenz gegen u^2 verhindern, vgl. 13.3.6(i). Dieser Prozeß läßt sich fortsetzen, um mittels Orthogonalisierung bezüglich span $\{u^1, u^2\}$ eine Vektorfolge in $\mathcal{X}_3 = \text{span} \{u^3, ..., u^n\}$ zu konstruieren usw. \square

B. Teilraumiteration

Bei der Vektoriteration wurde eine Folge $\{v^k\}$ normierter Vektoren erzeugt derart, daß span $\{v^k\}$ den dominanten Teilraum $\mathscr{S}_1 = \text{span }\{u^1\}$ approximiert. Falls die zu p dominanten Eigenwerten $\{\lambda_1, \ldots, \lambda_p\}$ gehörenden Eigenvektoren $\{u^1, \ldots, u^p\}$ gesucht sind, könnten unter der Zusatzvoraussetzung $|\lambda_{p+1}| < |\lambda_p| < \cdots < |\lambda_2| < |\lambda_1|$ die eben beschriebenen Deflationstechniken angewendet werden. Es ist i. allg. jedoch wesentlich günstiger, direkt Approximationen \mathcal{Y}_k für den p-ten dominanten Teilraum \mathscr{I}_p zu berechnen. Dazu bietet es sich an, einen Startteilraum \mathscr{Y}_1 der Dimension p mit der Matrix A zu iterieren.

13.3.8. Basisverfahren der Teilraumiteration.

- S0 (Initialisierung): Wähle Teilraum $\mathcal{Y}_1 \subset \mathbb{R}^n$ mit dim $(\mathcal{Y}_1) = p$, setze k := 1. S1 (Iteration): Definiere $\mathcal{Y}_{k+1} := A \mathcal{Y}_k = \{Ay : y \in \mathcal{Y}_k\}$
- S2: Setze k := k + 1, goto S1.

Für die Konvergenzanalyse benötigen wir ein Maß für die Entfernung zweier Teilräume. Es sei dazu $\mathscr{I} \subset \mathbb{R}^n$ ein Teilraum und $y \in \mathbb{R}^n$, $y \neq o$. Dann definieren wir den Winkel $\not\prec$ ($\boldsymbol{y}, \boldsymbol{\mathscr{S}}$) zwischen \boldsymbol{y} und $\boldsymbol{\mathscr{S}}$ durch

$$\measuredangle (\boldsymbol{y}, \boldsymbol{\mathscr{S}}) := \min \left\{ \measuredangle (\boldsymbol{y}, \boldsymbol{s}) : \boldsymbol{s} \in \boldsymbol{\mathscr{S}}, \, \boldsymbol{s} \neq \boldsymbol{o} \right\}.$$
(40)

Das Minimum in (40) existiert stets, und es gilt $0 \leq \langle (\boldsymbol{y}, \boldsymbol{\mathcal{S}}) \rangle \leq \pi/2$ sowie

$$\sin\left(\langle \langle (\boldsymbol{y}, \mathscr{S}) \rangle = \sin\left(\langle \langle (\boldsymbol{y}, P_{\mathscr{S}} \boldsymbol{y}) \rangle = \min\left\{ \|\boldsymbol{y} - \boldsymbol{s}\| \colon \boldsymbol{s} \in \mathscr{S} \} / \|\boldsymbol{y}\|, \quad (41)$$

wobei $P_{\mathscr{F}}$ den orthogonalen Projektor auf \mathscr{F} bezeichnet, vgl. 8.1.2 und insbesondere Abb. 8.1.1. Damit (41) auch im Fall $P_{\mathscr{F}} y = 0$ sinnvoll bleibt, setzen wir $\not\triangleleft (y, o)$ $:= \pi/2$ für jedes y.

Wenn $\mathcal{Y} \subset \mathbb{R}^n$ ein Teilraum ist, wird der maximale Winkel $\sphericalangle (\mathcal{Y}, \mathcal{S})$ – kurz: Winkel – zwischen \mathcal{Y} und \mathcal{S} durch

$$\langle (\mathcal{Y}, \mathcal{S}) := \max \{ \langle (\mathbf{y}, \mathcal{S}) : \mathbf{y} \in \mathcal{Y}, \mathbf{y} \neq \mathbf{o} \}$$

$$(42)$$

definiert, der minimale Winkel $\measuredangle_{\min}(\mathcal{Y}, \mathcal{S})$ dagegen durch

$$\boldsymbol{\mathbf{x}}_{\min}\left(\boldsymbol{\mathcal{Y}},\boldsymbol{\mathscr{S}}\right) := \min\left\{\boldsymbol{\mathbf{x}}\left(\boldsymbol{\mathcal{Y}},\boldsymbol{\mathscr{S}}\right): \boldsymbol{\mathcal{Y}}\in\boldsymbol{\mathcal{Y}}, \boldsymbol{\mathcal{Y}}=\boldsymbol{o}\right\}.$$
(43)

Es ist sofort einzusehen, daß für beliebige Teilräume $\mathcal{Y}, \mathcal{S} \subset \mathbb{R}^n$

$$\measuredangle_{\min}(\mathcal{Y},\mathcal{S}) = \measuredangle_{\min}(\mathcal{S},\mathcal{Y})$$
(44)

gilt, aber i. allg. ist

$$\langle \langle (\mathcal{Y}, \mathscr{S}) \neq \langle (\mathscr{S}, \mathcal{Y}) \rangle;$$

man betrachte eine Gerade und eine Ebene im Raum **R**³. Unter der zusätzlichen Voraussetzung dim $(\mathcal{Y}) = \dim (\mathcal{S})$ gilt jedoch

$$\langle (\mathcal{Y}, \mathcal{S}) = \langle (\mathcal{S}, \mathcal{Y}) = \langle (\mathcal{Y}^{\perp}, \mathcal{S}^{\perp})$$

$$(45)$$

sowie

$$\sin\left(\sphericalangle\left(\mathcal{Y},\mathcal{S}\right)\right) = \cos\left(\sphericalangle_{\min}\left(\mathcal{Y},\mathcal{S}^{\perp}\right)\right). \tag{46}$$

Zum Beweis von (45), (46) verweisen wir auf die Spezialliteratur, siehe B 13.4 für Hinweise.

Die Größe $\tan (\not\prec (\mathcal{Y}, \mathscr{S}))$ kann im Fall dim $(\mathcal{Y}) = \dim (\mathscr{S})$ zur Charakterisierung der Entfernung zwischen \mathcal{Y} und \mathscr{S} verwendet werden. Dabei gilt $\tan (\not\prec (\mathcal{Y}, \mathscr{S})) = 0$ genau dann, wenn $\mathcal{Y} = \mathscr{S}$ ist, und die Entfernung ist unendlich, wenn einer der Teilräume eine Richtung enthält, die orthogonal zum anderen Teilraum ist.

Man beachte, daß der Winkel zwischen zwei Teilräumen stets im Intervall $[0, \pi/2]$ liegt, während der Winkel zwischen zwei Vektoren Werte im gesamten Intervall $[0, \pi]$ annehmen kann, denn die Vektoren können entgegengesetzt orientiert sein. Letzteres spielt bei der Definition (40) keine Rolle, denn mit **s** liegt auch $-\mathbf{s}$ in \mathcal{S} , so daß mit $\not\prec (\mathbf{y}, \mathbf{s})$ auch der Winkel $\not\prec (\mathbf{y}, -\mathbf{s}) = \pi - \not\prec (\mathbf{y}, \mathbf{s})$ in der Menge liegt, über der das Minimum gebildet wird. Insbesondere gilt für zwei Vektoren \mathbf{u}, \mathbf{v}

 \measuredangle (span { \boldsymbol{u} }, span { \boldsymbol{v} }) = \measuredangle ($\boldsymbol{u}, \boldsymbol{v}$)

nur im Fall $u^{\mathsf{T}}v \geq 0$, also $0 \leq \langle (u, v) \leq \pi/2$; im Fall $u^{\mathsf{T}}v \leq 0$ ist

 $\measuredangle (\operatorname{span} \{\boldsymbol{u}\}, \operatorname{span} \{\boldsymbol{v}\}) = \pi - \measuredangle (\boldsymbol{u}, \boldsymbol{v}) = \measuredangle (-\boldsymbol{u}, \boldsymbol{v}) = \measuredangle (\boldsymbol{u}, -\boldsymbol{v}).$

Mit den eingeführten Größen läßt sich der nachfolgende Konvergenzsatz formulieren:

13.3.9. Satz. Es seien $\{\lambda_1, \ldots, \lambda_p\}$ dominante Eigenwerte von A, und \mathscr{S}_p sei der zugehörige p-te dominante Teilraum. Der Startteilraum $\mathscr{Y}_1 \subset \mathbb{R}^n$ genüge den

Bedingungen

$$\dim (\mathcal{Y}_1) = p \quad \text{und} \quad \sigma := \cos \left(\sphericalangle (\mathcal{Y}_1, \mathcal{S}_p) \right) > 0. \tag{47}$$

Dann gilt für die durch 13.3.8 definierte Folge $\{\mathcal{Y}_k\}$ von Teilräumen dim $(\mathcal{Y}_k) = p$ und

$$\tan \varphi_{k+1} \leq \varkappa \tan \varphi_k \leq \varkappa^k \tan \varphi_1 = \varkappa^k \sqrt{1 - \sigma^2} / \sigma$$
(48)

 mit

$$\varkappa := |\lambda_{p+1}/\lambda_p| < 1 \quad \text{und} \quad \varphi_k := \langle (\mathcal{Y}_k, \mathcal{S}_p). \tag{49}$$

Der Beweis verläuft analog zu dem von 13.3.3, wobei \mathcal{Y}_k in Komponenten bezüglich der komplementären Teilräume $\mathcal{S}_p = \operatorname{span} \{u^1, \ldots, u^p\}$ und $\mathcal{Z}_{p+1} := \operatorname{span} \{u^{p+1}, \ldots, u^n\} = \mathcal{S}_p^{\perp}$ zu zerlegen ist. \Box

Zur Implementierung von 13.3.8 liegt es nahe, den Teilraum \mathcal{Y}_k gemäß

$$\mathcal{Y}_{k} = \mathcal{R}(V_{k}), \qquad V_{k} := (\boldsymbol{v}^{k,1}, \dots, \boldsymbol{v}^{k,p}) \in \mathbf{R}^{n,p}$$

$$\tag{50}$$

durch die Spalten einer geeigneten Matrix V_k aufzuspannen. Dann spannen die Spalten von

$$W_{k+1} := (w^{k+1,1}, \dots, w^{k+1,p}) := AV_k = (Av^{k,1}, \dots, Av^{k,p})$$
(51)

den Teilraum $\mathcal{Y}_{k+1} = A \mathcal{Y}_k$ auf. Falls diese Spalten zur Vermeidung von Überbzw. Unterlauf gemäß

$$V_{k+1} := W_{k+1} \operatorname{diag}\left(1/\|w^{k+1,1}\|, \dots, 1/\|w^{k+1,p}\|\right)$$
(52)

normiert werden, ändert sich der Wertebereich nicht, d. h., es gilt $\mathcal{Y}_{k+1} = \mathcal{R}(V_{k+1})$ = $\mathcal{R}(W_{k+1})$. Wenn die Startbasis normiert ist, sind die Vektoren $\{v^{k,j}\}$ dann identisch mit denjenigen, die sich nach der Vektoriteration 13.3.1 aus $v^{1,j}$ ergeben. Im Fall $|\lambda_1| > |\lambda_2|$ approximieren daher alle diese Vektoren für genügend großes k die Richtung u^1 oder $-u^1$ beliebig genau, so daß die Basis $\{v^{k,1}, \ldots, v^{k,p}\}$ mit wachsendem k immer schiefwinkliger wird und vom numerischen Standpunkt immer schlechter zur Repräsentation von \mathcal{Y}_k geeignet ist. Aus diesem Grunde ist es zweckmäßig, mit einer orthonormalen Basis von \mathcal{Y}_{k+1} zu arbeiten, also V_{k+1} als orthonormale Matrix festzulegen. Dabei empfiehlt es sich aus später klar werdenden Gründen, mit mehr als p Vektoren zu arbeiten, also Teilräume der Dimension $m \geq p$ zu betrachten.

13.3.10. Teilraumiteration mittels orthonormaler Basen.

- S0 (Initialisierung): Wähle m mit $p \leq m \leq n$ und spaltenorthonormale Startmatrix $V_1 \in \mathbb{R}^{n,m}$, setze k := 1.
- S1 (Iteration): Berechne $W_{k+1} := AV_k$
- S2 (Orthonormalisierung): Bestimme spaltenorthonormale Matrix $Q_{k+1} \in \mathbb{R}^{n,m}$ und reguläre obere Dreiecksmatrix $R_{k+1} \in \mathbb{R}^{m,m}$ mit

$$W_{k+1} = Q_{k+1}R_{k+1}$$

S3: Setze $V_{k+1} := Q_{k+1}$ S4: Setze k := k + 1, goto S1 Aufwand pro Schritt: m Auswertungen von Ax (für S1), $\sim K_1 nm^2$ opms (für S2)

13.3.11. Bemerkung. (i) Wenn die Voraussetzungen von 13.3.9 mit m statt p erfüllt sind, gilt 13.3.9 auch für die in exakter Arithmetik ausgeführte Teilraumiteration 13.3.10. Insbesondere ist dann dim $(\mathcal{Y}_{k+1}) = \operatorname{rang}(W_{k+1}) = m$, so daß S2 in der Tat mit spaltenorthonormalem Q_{k+1} und regulärem R_{k+1} durchgeführt werden kann. Für eine ausreichend gute Approximation von \mathcal{Y}_{k+1} durch die Spalten von \mathcal{Q}_{k+1} brauchen dabei keine allzu hohen Orthogonalitätsforderungen gestellt zu werden, so daß sich das modifizierte Gram-Schmidt-Verfahren 10.1.5 zur Realisierung anbietet. Bei dieser Version ist in den Aufwandsangaben $K_1 = 1$ zu setzen. Bei höheren Orthogonalitätsforderungen muß i. allg. mit Re-Orthogonalisierung gearbeitet werden, wobei sich der Aufwand verdoppelt, vgl. 10.1.8(iii).

(ii) Wenn für festes p mit $1 \leq p \leq m$ die Partitionierungen $V_k = (V'_k | V''_k)$ mit $V'_k \in \mathbb{R}^{n,p}, \quad V''_k \in \mathbb{R}^{n,m-p}$ und analog $Q_k = (Q'_k | Q''_k), \quad W_k = (W'_k | W''_k)$ sowie $m{R}_k = \left(egin{array}{c|c|c|} m{R}''_k & m{R}''_k \ m{O} & m{R}'''_k \end{array}
ight)$ vorgenommen werden, gilt $m{W}'_{k+1} = A V'_k$ und $m{W}'_{k+1} = m{Q}'_{k+1} m{R}'_{k+1},$ d. h., die Spalten von V_k' bilden eine orthonormale Basis desjenigen Teilraums ${\mathscr Y}_k'$

der Dimension p, der durch Teilraumiteration aus \mathcal{Y}'_1 entsteht. Algorithmus 13.3.10 liefert also gleichzeitig auch Approximationen für die invarianten Teilräume \mathcal{S}_n mit $p \leq m$, sofern die Eigenwerte $\{\lambda_1, \ldots, \lambda_p\}$ dominant sind.

(iii) Für die Computerrealisierung von 13.3.10 lassen sich zu 13.3.5 analoge Aussagen beweisen. Aus Aufwands- und Genauigkeitsgründen ist auch hier wesentlich, daß $\varkappa = |\lambda_{p+1}/\lambda_p|$ deutlich kleiner als 1 ist, vgl. 13.3.6.

C. Der Rayleigh-Ritz-Algorithmus

Im folgenden wenden wir uns dem Problem zu, wie den durch Q_k repräsentierten Teilraumapproximationen \mathcal{Y}_k möglichst gute Näherungen $\{\mu_i, v^j\}$ für die gesuchten Eigenpaare von A zugeordnet werden können. Wenn in Analogie zu 13.1.8 die Quadratsummen der Residuumsnormen als Kriterium für die Güte der Approximationen verwendet werden, ergibt sich das folgende Resultat.

13.3.12. Satz. Es sei $\mathcal{Y} = \mathcal{R}(Q)$ ein *m*-dimensionaler Teilraum, der durch die Spalten der spaltenorthonormalen Matrix $Q = (q^1, ..., q^m) \in \mathbb{R}^{n,m}$ $(1 \le m \le n)$ aufgespannt wird. Aus Q und $A \in S^{n,n}$ werde die Matrix

$$\boldsymbol{P} := (\varrho_{ij}) := \boldsymbol{Q}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{Q} \in \mathbf{S}^{m,m}$$
(53)

gebildet, und

$$\boldsymbol{P} = \boldsymbol{X}\boldsymbol{M}\boldsymbol{X}^{\mathsf{T}} \tag{54}$$

sei die zugehörige Eigenwertzerlegung mit $M = \text{diag}(\mu_1, ..., \mu_m)$ und ortho-

385

gonalem $X \in \mathbb{R}^{m,m}$. Dann löst $\{M, V\}$ mit

$$\boldsymbol{V} := (\boldsymbol{v}^1, \dots, \boldsymbol{v}^m) := \boldsymbol{Q}\boldsymbol{X} \tag{55}$$

die Aufgabe

$$\|AV - VM\|_F^2 = \sum_{j=1}^m \|Av^j - \mu_j v^j\|^2 \to \text{Minimum}!$$
(56)

bei

$$M = ext{diag} (\mu_j) \in \mathbb{R}^{m,m}, \quad V \in \mathbb{R}^{n,m} ext{ spaltenorthonormal mit } \mathcal{R}(V) = \mathcal{R}(Q).$$

(57)

Be we is. Wegen $\mathcal{R}(V) = \mathcal{R}(Q)$ ist V = QX mit regulärem $X \in \mathbb{R}^{m,m}$, und wegen der geforderten Spaltenorthonormalität muß $I = V^{\mathsf{T}}V = X^{\mathsf{T}}Q^{\mathsf{T}}QX = X^{\mathsf{T}}X$ gelten, also X orthogonal sein. Für jedes solche V gilt dann für die Residualmatrix R = AV - VM

$$\|oldsymbol{R}\|_F = \|oldsymbol{R}oldsymbol{X}^{\intercal}\|_F = \|oldsymbol{A}oldsymbol{Q} - oldsymbol{Q}(oldsymbol{X}oldsymbol{M}oldsymbol{X}^{\intercal})\|_F,$$

so daß sich die Aufgabe mit der Bezeichnung $B := (b^1, ..., b^m) := AQ$ auf

$$|\boldsymbol{R}||_{F}^{2} = ||\boldsymbol{B} - \boldsymbol{Q}\boldsymbol{P}||_{F}^{2} = \sum_{j=1}^{m} ||\boldsymbol{b}^{j} - \boldsymbol{Q}\boldsymbol{\varrho}^{j}||^{2} \to \text{Minimum!}$$
(58)

bei

$$P := (q^1, \ldots, q^m) := XMX^{\mathsf{T}}, \quad M \text{ diagonal}, \quad X \text{ orthogonal},$$

reduziert. Wenn die Nebenbedingungen zunächst ignoriert werden, hat (58) die Lösungen $Q^{j} := (Q^{\mathsf{T}}Q)^{-1} Q^{\mathsf{T}} b^{j} = Q^{\mathsf{T}} b^{j}$, also $P = Q^{\mathsf{T}} B = Q^{\mathsf{T}} A Q$. Da die Matrix P symmetrisch ist, kann sie mittels ihrer Eigenwertzerlegung (54) in der geforderten Form dargestellt werden.

Die in 13.3.12 beschriebene Zuordnung $\{A, Q\} \rightarrow \{M, V\}$ wird Rayleigh-Ritz-Algorithmus — kurz: RR-Algorithmus — genannt. Die Größen $\{\mu_i\}$ und $\{v^i\}$ heißen Ritzsche Eigenwerte bzw. Eigenvektoren bezüglich des Teilraumes $\mathcal{Y} = \mathcal{R}(Q)$. Die Matrix $\mathbf{P} = (\varrho_{ij})$ mit den Elementen $\varrho_{ij} = (q^i)^{\mathsf{T}} A q^j$ stellt eine Verallgemeinerung des Rayleigh-Quotienten dar und wird Projektion von A auf $\mathcal{Y} = \mathcal{R}(Q)$ genannt, vgl. Ü 13.3.5.

13.3.13. Rayleigh-Ritz-Algorithmus. Es sei $A \in S^{n,n}$, und $Q \in \mathbb{R}^{n,m}$ sei spaltenorthonormal.

- S1: Bilde $B := AQ \in \mathbb{R}^{n,m}$, berechne $P := Q^{\mathsf{T}}B \in \mathbb{S}^{m,m}$
- S2: Berechne Eigenwertzerlegung $P = XMX^{\intercal}$ mit diagonalem M und orthogonalem X
- S3: Berechne V := QX

Aufwand: m Auswertungen von Ax + \sim $nm^2/2$ opms (für S1),

 $\sim K_2 m^3$ opms (für S2, \breve{K}_2 abhängig vom verwendeten Verfahren),

 $\sim nm^2$ opms (für S3)

13.3.14. Bemerkung. (i) Falls \mathcal{Y} ein invarianter Teilraum von A ist, also $\mathcal{R}(AQ) \subset \mathcal{R}(Q)$ gilt, muß AQ = QZ mit einer Matrix $Z \in \mathbb{R}^{m,m}$ sein. Dann folgt $P = Q^{\mathsf{T}}AQ$

 $= Q^{\intercal}QZ = Z$, mithin AQ = QP und durch Rechtsmultiplikation mit X schließlich

$$AV = VM$$
, d. h. $Av^{j} = \mu_{j}v^{j}$ $(j = 1, ..., m)$. (59)

Im Fall eines invarianten Teilraumes verschwindet also die Residualmatrix R= AV - VM, und die Ritzschen Eigenpaare $\{\mu_i, v^i\}$ von A bezüglich \mathcal{Y} sind identisch mit gewissen *m* Eigenpaaren $\{\lambda_{l(i)}, \boldsymbol{u}^{l(j)}\}$ von *A*, wobei $\mathcal{Y} = \text{span} \{\boldsymbol{u}^{l(1)}, \dots, \boldsymbol{u}^{l(m)}\}$ ist.

(ii) Wenn die Spalten von Q bereits gute Näherungen für die gesuchten Eigenvektoren darstellen, weicht P nur wenig von einer Diagonalmatrix ab. Zur Berechnung der Eigenwertzerlegung in S2 bietet sich dann das Jacobi-Verfahren an.

Falls $\mathbf{R} \neq \mathbf{O}$, also \mathcal{Y} nicht invariant ist, lassen sich unter Verwendung von \mathbf{R} die folgenden Aussagen über die Güte der Ritzschen Näherungen $\{\mu_i, v^i\}$ machen:

13.3.15. Satz. Es sei $Q \in \mathbb{R}^{n,m}$ spaltenorthonormal, und $M = \text{diag}(\mu_i)$ sowie $V = (v^1, ..., v^m)$ seien die in exakter Arithmetik berechneten Ritzschen Eigen- $V = (v^{i}, ..., v^{m}) \text{ selen die in exakter Artenmetik berechneten Intzschen Engen werte bzw. Eigenvektoren mit der zugehörigen Residualmatrix <math>\boldsymbol{R} := (\boldsymbol{r}^{1}, ..., \boldsymbol{r}^{m})$:= AV - VM. Dann gilt: (i) Es gibt Störungen $\delta A_{j} \in \mathbf{S}^{n,n}$ mit $(A + \delta A_{j}) v^{j} = \mu_{j} v^{j}$ und $\|\delta A_{j}\| \leq \|\boldsymbol{r}^{j}\|$ (j = 1, ..., m). (60) (ii) Es gibt Eigenwerte $\{\lambda_{l(1)}, ..., \lambda_{l(m)}\}$ von \boldsymbol{A} mit

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}_j) \, \boldsymbol{v}^j = \mu_j \boldsymbol{v}^j \quad \text{und} \quad \|\boldsymbol{\delta}\boldsymbol{A}_j\| \leq \|\boldsymbol{r}^j\| \qquad (j = 1, ..., m). \tag{60}$$

$$|\mu_j - \lambda_{l(j)}| \leq \|\boldsymbol{R}\| \quad \text{und} \quad \sqrt{\sum_{j=1}^n (\mu_j - \lambda_{l(j)})^2} \leq \sqrt{2} \|\boldsymbol{R}\|_F.$$
(61)

Beweis. Aussage (i) folgt direkt aus 13.1.7, zum Nachweis von (ii) sei auf die Spezialliteratur verwiesen, siehe B 13.4.

13.3.16. Bemerkung. (i) Aus 13.3.15(i) folgt nach 13.1.2, daß im Intervall $[\mu_i - ||r^i||,$ $\mu_i + \|\mathbf{r}^i\|$ mindestens ein Eigenwert λ_i von A liegt. Wenn alle diese Intervalle disjunkt sind, ist dies eine schärfere Aussage als (ii). Falls μ_i den Eigenwert λ_l genügend gut approximiert und λ_l von den übrigen Eigenwerten ausreichend getrennt ist, kann auch der Winkel zwischen $\pm v^{j}$ und dem entsprechenden Eigenvektor u^{l} abgeschätzt werden. Bei mehrfachen Eigenwerten und Eigenwerthaufen gelten analoge Abschätzungen für die zugehörigen Teilräume, auf die hier jedoch nicht eingegangen werden kann, vgl. 13.1.5.

(ii) Bei Computerrealisierung des RR-Algorithmus ist wesentlich, daß die numerisch berechnete Matrix Q die im Sinne von

$$\|\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{Q} - \boldsymbol{I}_m\| \leq \nu F_0, \quad F_0 \text{ akzeptabel}, \tag{62}$$

höchstmögliche Qualität der Orthogonalität hat, vgl. 10.1.8. Falls (62) erfüllt ist und die Eigenwertzerlegung in S2 von 13.3.13 nach einem numerisch gutartigen Verfahren vorgenommen wird, gelten die Aussagen von 13.3.15 auch für die numerisch berechneten Ritzschen Näherungen. Dabei ist in den Schranken ein Term der Größenordnung $vF_1 ||A||$ mit akzeptablem F_1 zu addieren. \Box

Wir gehen abschließend auf mögliche Kombinationen der Teilraumiteration mit dem RR-Algorithmus ein.

13.3.17. Bemerkung. (i) Die einfachste Möglichkeit zur Kopplung von Teilraumiteration und RR-Algorithmus besteht darin, die nach Schritt S2 von 13.3.10 vorliegende spaltenorthonormale Matrix $Q_{k+1} =: Q$ als Ausgangspunkt für den RR-Algorithmus zu wählen. Dann sollte natürlich nicht Q_{k+1} , sondern die denselben Teilraum $\mathcal{Y}_{k+1} = \mathcal{R}(Q_{k+1}) = \mathcal{R}(V)$ repräsentierende Matrix $V =: V_{k+1}$ der Ritzschen Eigenvektoren als neue Basismatrix gewählt werden. Bei Konvergenz der Ritzschen Eigenvektoren gegen Eigenvektoren von A ist dann P fast diagonal, vgl. 13.3.14 (ii). In dieser Version ist S3 von 13.3.10 also durch

S3': Setze $Q := Q_{k+1}$, berechne Ritzsche Näherungen $\{M, V\}$ nach 13.3.13, setze $M_{k+1} := M, V_{k+1} := V$

zu ersetzen. Pro Schritt sind dann 2m Auswertungen von Ax erforderlich, nämlich m in S1 von 13.3.10 und m in S1 von 13.3.13.

(ii) Unter Verwendung der im folgenden Schritt ohnehin zu berechnenden Matrix AV_{k+1} läßt sich die Residualmatrix $\mathbf{R} = \mathbf{R}_{k+1} = AV_{k+1} - V_{k+1}M_{k+1}$ bilden und damit gemäß 13.3.15 eine Einschätzung der Genauigkeit der Ritzschen Näherungen $\{M_{k+1}, V_{k+1}\}$ vornehmen. Dies kann als Grundlage für ein Abbruchkriterium zur Teilraumiteration genommen werden.

(iii) Für die in (i) beschriebene Variante der sogenannten *simultanen Iteration* gilt unter der Voraussetzung

$$0 \leq \lambda_n \leq \cdots \leq \lambda_{m+1} < \lambda_m < \cdots < \lambda_2 < \lambda_1 \tag{63}$$

die Abschätzung

$$\sin\left(\langle \langle (\boldsymbol{u}^{j}, \boldsymbol{v}^{k+1, j}) \rangle = O(|\lambda_{m+1}/\lambda_{j}|^{k}) \qquad (j = 1, ..., m),$$
(64)

sofern die Ritzschen Eigenwerte in Analogie zu (63), also gemäß $\mu_m \leq \cdots \leq \mu_2 \leq \mu_1$ numeriert werden. Die Konvergenz ist also um so besser, je größer m + 1 - j ist. Zur Berechnung der ersten p Eigenwerte ist es daher zweckmäßig, die Spaltenzahl metwas größer als p zu wählen. Trotz des höheren Aufwandes pro Schritt ist der Gesamtaufwand wegen der schnelleren Konvergenz dann meist niedriger als im Fall p = m, sofern λ_{m+1} merklich kleiner als λ_{p+1} ist.

(iv) Die Aussagen von (iii) lassen sich auch auf den indefiniten Fall und nicht notwendig getrennte dominante Eigenwerte $\{\lambda_1, \ldots, \lambda_m\}$ übertragen. Die Formulierung und erst recht die Beweise sind allerdings wesentlich komplizierter, so daß hier darauf verzichtet wird. Bei mehrfachen Eigenwerten bzw. Eigenwerthaufen ist natürlich nicht mit individueller Konvergenz der Ritzschen Eigenvektoren zu rechnen, allerdings approximieren die zugehörigen Teilräume die entsprechenden invarianten Eigenvektorräume von A, vgl. 13.1.B.

(v) Mit nur m Auswertungen von Ax pro Schritt kommt die folgende, in der

Prozedur ritzit von RUTISHAUSER [71] verwendete Modifikation aus, bei der S3 von 13.3.10 durch

S3.1': Bilde $\overline{P} := R_{k+1}R_{k+1}^{\mathsf{T}}$

S3.2': Berechne Eigenwertzerlegung $\overline{P} = \overline{X}\overline{M}^2\overline{X}^{\intercal}$ von \overline{P} mit diagonalem $\overline{M}^2 = \text{diag}(\overline{\mu}_1^2, \dots, \overline{\mu}_m^2)$ und orthogonalem \overline{X}

S3.3': Setze $V_{k+1} := Q_{k+1} \overline{X}$

ersetzt wird. Zur Analyse dieser zu empfehlenden Version muß auf die Literatur verwiesen werden, siehe B 13.4.

(vi) Aus Aufwandsgründen ist es zweckmäßig, den RR-Algorithmus nur in gewissen Abständen mit der Teilraumiteration zu koppeln.

(vii) Für alternative Möglichkeiten zur Festlegung der approximierenden Teilräume siehe B 13.5.

Übungsaufgaben

Ü 13.3.1. Es seien $\lambda_n \leq \lambda_{n-1} \leq \cdots \leq \lambda_2 < \lambda_1$ die Eigenwerte von A. Dann hat die durch die Verschiebung $\bar{A} := A - \mu I$ entstehende Matrix \bar{A} die Eigenwerte $\bar{\lambda}_i = \lambda_i - \mu$. Man zeige, daß die Größe

$$arkappa = arkappa(\mu; \lambda_1, ..., \lambda_n) := \max \{ |ar{\lambda}_i| : i \neq 1 \} / |ar{\lambda}_1|$$

für $\mu_{opt} := (\lambda_2 + \lambda_n)/2$ minimal wird, und der minimale Wert ist

$$\varkappa_{\text{opt}} = \varkappa(\mu_{\text{opt}}; \lambda_1, \dots, \lambda_n) = \frac{1-\tau}{1+\tau} \quad \text{mit} \quad 0 < \tau := \frac{\lambda_1 - \lambda_2}{\lambda_1 - \lambda_n} \leq 1.$$

Wie kann λ_n unter der Voraussetzung $\lambda_n < \lambda_{n-1}$ in einen dominanten Eigenwert von \overline{A} transformiert werden?

Ü 13.3.2. Man zeige unter Verwendung von (18), daß

$$arphi_k^2 = \|oldsymbol{w}^{k+1}\|^2 = \lambda_1^2 - [\lambda_1^2 - \|oldsymbol{A}oldsymbol{z}^k\|^2]\sin^2arphi_k$$

gilt und folgere hieraus (22).

Ü 13.3.3. Man stelle $\omega_{k,i}$ unter Verwendung von (16) und (18) in der Form

$$\omega_{k,j} = \frac{(\boldsymbol{\vartheta}_k \boldsymbol{w}^{k+1})_j}{(\boldsymbol{\vartheta}_k \boldsymbol{v}^k)_j} = \lambda_1 + \frac{[(\boldsymbol{A}\boldsymbol{z}^k)_j - \lambda_1(\boldsymbol{z}^k)_j]\sin\varphi_k}{(\boldsymbol{u}^1)_j\cos\varphi_k + (\boldsymbol{z}^k)_j\sin\varphi_k}$$

dar und leite hieraus die Abschätzung (25) her.

Ü 13.3.4. Es seien $P \in \mathsf{R}^{n,p}, Q \in \mathsf{R}^{n,q}$ spaltenorthonormal. Man zeige, daß

 $\|\boldsymbol{P}^{\intercal}\boldsymbol{Q}\| \leq 1$

gilt und der minimale Winkel zwischen den Teilräumen $\mathcal{Y} := \mathcal{R}(\boldsymbol{P}), \, \mathcal{X} := \mathcal{R}(\boldsymbol{Q})$ durch

$$\cos\left(\measuredangle_{\min}\left(\mathscr{Y}, \mathscr{Z} \right) \right) = \| \boldsymbol{P}^{\mathsf{T}} \boldsymbol{Q} \|$$

gegeben ist.

Hinweis: Man beachte Ü 1.2.10.

Ü 13.3.5. Es sei $\mathcal{Y} = \mathcal{R}(\mathcal{Q})$ der durch die Spalten der spaltenorthonormalen Matrix $\mathcal{Q} \in \mathbb{R}^{n,p}$ aufgespannte Teilraum. Dann kann die durch \mathcal{A} vermittelte lineare Abbildung $\boldsymbol{x} \in \mathbb{R}^n$ $\rightarrow \mathcal{A}\boldsymbol{x} \in \mathbb{R}^n$ wie folgt bezüglich \mathcal{Y} approximiert werden:

$$oldsymbol{x} \in \mathbf{R}^{oldsymbol{n}} o oldsymbol{y} := oldsymbol{P} y oldsymbol{x} \in \mathcal{Y} o oldsymbol{A} oldsymbol{y} \in \mathbf{R}^{oldsymbol{n}} o oldsymbol{z} := oldsymbol{P} y oldsymbol{A} oldsymbol{y} \in \mathcal{Y}$$
 .

Dabei bezeichnet $P_{\mathcal{Y}} := QQ^{\mathsf{T}}$ den Projektor auf \mathcal{Y} . Man zeige: Wenn

 $y = Pyx = QQ^{\mathsf{T}}x = Q\hat{y}, \quad \hat{y} := Q^{\mathsf{T}}x \text{ und } z := Q\hat{z}$

gesetzt wird, gilt

$$\hat{\boldsymbol{z}} = \boldsymbol{P}\hat{\boldsymbol{y}} \quad \text{mit} \quad \boldsymbol{P} := \boldsymbol{Q}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{Q}.$$

Die Matrix **P** heißt Projektion von A auf $\mathcal{Y} = \mathcal{R}(\mathbf{Q})$.

13.4. Inverse Iteration

Die Vektor- und Teilraumiteration aus 13.3 hat den Nachteil, daß sie nur Approximationen für Eigenvektoren bzw. Eigenräume liefert, die zu dominanten Eigenwerten gehören; außerdem ist die Konvergenz i. allg. langsam. Meistens sind jedoch andere Eigenvektoren gesucht, etwa der zu einem vorgegebenen Eigenwert λ_i gehörende Eigenvektor \boldsymbol{u}^i , wobei eine Näherung μ_i für λ_i bekannt ist, vgl. 13.1.10, oder die zu p betragskleinsten Eigenwerten gehörenden Eigenvektoren. Solche Eigenwerte lassen sich i. allg. auch durch eine Spektralverschiebung nicht zu dominanten machen. Eine Möglichkeit, auch die beschriebenen Aufgaben sehr effektiv mittels Vektoriteration zu lösen, beruht auf den folgenden beiden Fakten:

- Wenn $A = U \Lambda U^{\intercal} = U$ diag $(\lambda_j) U^{\intercal}$ die Eigenwertzerlegung von A bezeichnet und A regulär ist, gilt $A^{-1} = U \Lambda^{-1} U^{\intercal} = U$ diag $(1/\lambda_j) U^{\intercal}$, d. h., die Inverse besitzt die reziproken Eigenwerte, während die Eigenvektoren unverändert bleiben. Die betragskleinsten Eigenwerte von A gehen also in die betragsgrößten von A^{-1} über.
- Durch eine Spektralverschiebung

$$\bar{A} := A_{\mu} := A - \mu I \tag{1}$$

kann jeder Eigenwert λ_j von A zum betragskleinsten Eigenwert $\overline{\lambda}_j = \lambda_j [A_\mu] = \lambda_j - \mu$ von A_μ gemacht werden, sofern $|\lambda_j - \mu|$ genügend klein ist, vgl. Abb. 13.4.1.



Abb. 13.4.1. Spektralverschiebung $\bar{\lambda}_j = \lambda_j - \mu$

Im folgenden seien die Eigenwerte λ_j von A gemäß

$$|\lambda_1 - \mu| \le |\lambda_2 - \mu| \le \dots \le |\lambda_n - \mu| \tag{2}$$

numeriert, und es gelte

$$0 < |\lambda_1 - \mu|, \tag{3}$$

d. h., A_{μ} sei regulär. Dann existiert $\hat{A} := (A - \mu I)^{-1}$, und es ist

$$\hat{\boldsymbol{A}} = (\boldsymbol{A} - \boldsymbol{\mu} \boldsymbol{I})^{-1} = \boldsymbol{U}(\boldsymbol{\Lambda} - \boldsymbol{\mu} \boldsymbol{I})^{-1} \boldsymbol{U}^{\mathsf{T}} = \boldsymbol{U} \operatorname{diag} \left(1/(\lambda_{j} - \boldsymbol{\mu}) \right) \boldsymbol{U}^{\mathsf{T}}.$$
 (4)

Die Matrix \hat{A} besitzt also die Eigenpaare $\{\hat{\lambda}_j, u^j\}, \ \hat{\lambda}_j := 1/(\lambda_j - \mu)$. Insbesondere ist $\hat{\lambda}_1 = 1/(\lambda_1 - \mu)$ wegen (2) betragsgrößter Eigenwert von \hat{A} . Der zugehörige Eigenvektor u^1 kann dann im Fall der Dominanz von $\hat{\lambda}_1$, d. h. im Fall

$$|\lambda_1 - \mu| < |\lambda_2 - \mu|, \tag{5}$$

durch Vektoriteration mit der Matrix \hat{A} gefunden werden; man spricht von *inverser* Iteration.

13.4.1. Basisverfahren der inversen Iteration.

- S0 (Initialisierung): Wähle μ mit $\mu \neq \lambda_j$ (j = 1, ..., n) und v^1 mit $||v^1|| = 1$, setze k := 1S1 (Iteration): Berechne $w^{k+1} := \hat{A}v^k = (A - \mu I)^{-1}v^k$ S2 (Normierung): Setze $v^{k+1} := w^{k+1}/||w^{k+1}||$ S3: Setze k := k + 1, goto S1

13.4.2. Bemerkung. (i) Bei der Realisierung von S1 wird i. allg. nicht die explizit berechnete Inverse A_u^{-1} verwendet, sondern w^{k+1} wird als Lösung des Gleichungssystems

$$A_{\mu}\boldsymbol{w}^{k+1} = (\boldsymbol{A} - \boldsymbol{\mu}\boldsymbol{I}) \, \boldsymbol{w}^{k+1} = \boldsymbol{v}^k \tag{6}$$

bestimmt. Dazu bietet sich eine Dreiecksfaktorisierung von A_{μ} an, etwa die Spaltenpivotisierte Gauß-Faktorisierung $P_{\mu}A_{\mu} = L_{\mu}R_{\mu}$ aus 5.2.A oder die symmetrische BKP-Faktorisierung $P_{\mu}A_{\mu}P_{\mu}^{\intercal} = L_{\mu}D_{\mu}L_{\mu}^{\intercal}$ aus 6.1.B. Letztere nutzt zwar die Symmetrie von A_{μ} voll aus, kann aber aus einer eventuellen Bandgestalt von A keine Vorteile ziehen. In beiden Fällen gilt für die berechnete Lösung w^{k+1} von (6)

$$(A_{\mu} + \delta A_k) \boldsymbol{w}^{k+1} = \boldsymbol{v}^k \quad \text{mit} \quad \|\delta A_k\| \leq \nu F \|A_{\mu}\|, \tag{7}$$

wobei F die Kumulationskonstante des Lösungsverfahrens bezeichnet. Man beachte, daß A_{μ} i. allg. nicht definit ist, so daß die einfache Cholesky-Faktorisierung nicht verwendet werden kann. Da die Koeffizientenmatrix für alle k dieselbe ist, braucht die Faktorisierung nur einmal berechnet zu werden.

(ii) In wichtigen Anwendungen der inversen Iteration ist A tridiagonal, siehe 13.6. In diesem Fall sollte die spaltenpivotisierte Gauß-Faktorisierung nach 6.4.4 verwendet werden, deren Berechnung dann wie die Lösung von (6) nur O(n) opms kostet. Dasselbe gilt für Bandmatrizen.

(iii) Bei Realisierung in exakter Arithmetik wird das Konvergenzverhalten der Folge $\{v^k\}$ durch 13.3.3 bei Anwendung auf A_{μ}^{-1} statt A beschrieben: Falls (2), (3), (5) und (13.3.9) gelten, konvergiert $\{\vartheta_k v^k\}$ gegen u^1 , und der Konvergenzfaktor ist

$$\hat{\boldsymbol{x}} := |\boldsymbol{\lambda}_1 - \boldsymbol{\mu}| / |\boldsymbol{\lambda}_2 - \boldsymbol{\mu}|. \tag{8}$$

Die Konvergenz ist also um so schneller, je besser die Eigenwertnäherung μ ist.

(iv) Bei Computerrealisierung wird das ideale Verhalten der Folge $\{v^k\}$ durch die Rundungsfehler gestört. Wegen der erforderlichen Gleichungsauflösung ist die Fehleranalyse etwas aufwendiger als bei der einfachen Vektoriteration. Wir betrachten zunächst den Fall, daß A_{μ} im Sinne von

$$\nu F \operatorname{cond} (A_{\mu}) = \nu F |\lambda_n - \mu| / |\lambda_1 - \mu| =: \overline{\varkappa} < 1$$
(9)

mit F aus (7) ausreichend regulär ist, d. h., daß μ keine zu gute Approximation für λ_1 darstellt. Dann folgt aus 4.1.2 die Regularität von $A_{\mu} + \delta A_k$, und man erhält

$$\boldsymbol{w}^{k+1} = (A_{\mu} + \boldsymbol{\delta}A_{k})^{-1} \, \boldsymbol{v}^{k} = (\hat{A} + \boldsymbol{\delta}\hat{A}_{k}) \, \boldsymbol{v}^{k}$$

mit $\boldsymbol{\delta}\hat{A}_{k} := (A_{\mu} + \boldsymbol{\delta}A_{k})^{-1} - A_{\mu}^{-1}$ (10)

und

$$\|\delta \hat{A}_{k}\| \leq [\|\hat{A}\|^{2}/(1-\bar{z})] \|\delta A_{k}\| \leq \nu \hat{F} \|\hat{A}\|, \qquad \hat{F} := [\text{cond } (A_{\mu})/(1-\bar{z})] F.$$
(11)

Die Voraussetzungen von 13.3.5 sind also mit $\{\hat{A}, \hat{F}\}$ statt $\{A, F\}$ erfüllt, sofern

$$\hat{k} + 10\nu(\hat{F} + n/2) = |\lambda_1 - \mu|/|\lambda_2 - \mu| + 10\nu(\hat{F} + n/2) < 1$$
 (12)

gilt.

(v) Die in (iv) beschriebene Situation liegt in der Regel vor, wenn die inverse Iteration mit $\mu = 0$ – also ohne Verschiebungen – zur Bestimmung des betragskleinsten Eigenwertes verwendet wird. Falls die p betragskleinsten Eigenwerte und zugehörigen Eigenvektoren gesucht sind, sollte natürlich die Teilraumiteration mit $\hat{A} := A_0^{-1} := A^{-1}$ statt A benutzt werden. Die Dominanzforderung (13.3.2) geht dann in

$$0 < |\lambda_1| \le \dots \le |\lambda_p| < |\lambda_{p+1}| \tag{13}$$

über, und die Aussagen aus 13.3.B,C sind sinngemäß gültig.

Falls $|\lambda_1 - \mu|$ sehr klein ist, kann nicht mehr wie in (iv) geschlossen werden, da (9) i. allg. nicht erfüllt ist. Wir gehen auf diesen wichtigen Spezialfall im folgenden genauer ein. Es sei dazu μ eine im Sinne von

$$|\mu - \lambda| \le \nu F_0 \, \|A\| =: \varepsilon_0 \tag{14}$$

sehr gute Näherung für irgendeinen Eigenwert λ von A. Wir numerieren die Eigenwerte wieder gemäß (2), d. h., $\lambda = \lambda_1$ sei der μ am nahesten gelegene Eigenwert. Gesucht ist ein Vektor v mit ||v|| = 1, so daß $\{\mu, v\}$ eine akzeptable Eigenpaarapproximation darstellt. Wenn die Norm des Residuums $\mathbf{r} = (A - \mu I) v$ als Kriterium für die Güte von v verwendet wird, ist nach 13.1.10 jeder zu λ_1 gehörende Eigenvektor u^1 optimal, so daß sich die inverse Iteration zur Berechnung von v anbietet. Dabei ist es zweckmäßig, die folgende spezielle Version zu verwenden.

13.4.3. Berechnung einer akzeptablen Eigenvektorapproximation v bei gegebener Eigenwertnäherung μ mittels inverser Iteration.

Aufgabe: Zur gegebenen Eigenwertapproximation μ mit (14) ist v mit $\|v\| = 1$ so zu bestimmen, daß $\{\mu, v\}$ ein akzeptables Eigenpaar von A ist.

Algorithmus:

S0 (Initialisierung):

S0.1: Berechne Faktorisierung $A_{\mu} = P_{\mu}^{\mathsf{T}} L_{\mu} R_{\mu}$ gemäß 5.2.3 mit regulärem R_{μ}

S0.2: Berechne w^1 als Lösung von $R_u w^1 = e := (1, ..., 1)^T$

- S0.3: Setze $v^1 := w^1 / ||w^1||$
- S0.4: Setze k := 1 und $\varepsilon := \sqrt{n} (\varepsilon_0 + \varepsilon_1) = \nu \sqrt{n} (F_0 + F_1) ||A||$ mit F_0 aus (14) und $F_1 := (2 + \nu F_0) F = 2F$, $\varepsilon_1 := \nu F_1 ||A||$, F aus (7)

S1 (Iteration): Bestimme \boldsymbol{w}^{k+1} als Lösung von $\boldsymbol{P}_{\mu}^{\mathsf{T}} \boldsymbol{L}_{\mu} \boldsymbol{R}_{\mu} \boldsymbol{w}^{k+1} = \boldsymbol{v}^{k}$ S2 (Normierung): Setze $\omega_{k+1} := \|\boldsymbol{w}^{k+1}\|$ und $\boldsymbol{v}^{k+1} := \boldsymbol{w}^{k+1}/\omega_{k+1}$

S3 (Abbruchtest): if $\varepsilon * \omega_{k+1} < 1$ then [k := k + 1, goto S1]

Aufwand: $\sim n^3/3$ opms (für S 0) + $\sim kn^2$ opms (für k Iterationsschritte S1, S2, falls A voll besetzt ist)

13.4.4. Aussage. Für die Computerrealisierung von 13.4.3 gilt:

(i) Falls das Verfahren nach k Schritten wegen $\varepsilon * \omega_{k+1} \ge 1$ mit $v := v^{k+1}$ abbricht, gibt es eine Störung $\delta A \in \mathbf{S}^{n,n}$, so daß $\{\mu, v\}$ der Beziehung

$$(\boldsymbol{A} + \boldsymbol{\delta}\boldsymbol{A}) \boldsymbol{v} = \mu \boldsymbol{v} \quad \text{mit} \quad \|\boldsymbol{\delta}\boldsymbol{A}\| \leq \nu \left\{ \sqrt{n} F_{0} + \left(\sqrt{n} + 1\right) F_{1} \right\} \|\boldsymbol{A}\| \sim \varepsilon \|\boldsymbol{A}\| \quad (15)$$

genügt, d. h., $\{\mu, v\}$ ist ein akzeptables Eigenpaar von A.

(ii) Das Abbruchkriterium $\varepsilon * \omega_{k+1} \ge 1$ ist fast immer nach wenigen Schritten $k = 1, 2, 3, \dots$ erfüllt.

Beweis. Aus (14) folgt $|\mu| \leq |\lambda_1| + \varepsilon_0 \leq (1 + \nu F_0) ||A||$, so daß in (7) weiter gemäß

$$\|\delta A_k\| \le \nu F \|A_{\mu}\| \le \nu F[\|A\| + |\mu|] \le \nu F(2 + \nu F_0) \|A\| = \nu F_1 \|A\| = \varepsilon_1$$
(16)

abgeschätzt werden kann.

Zu (i): Aus (7) und der Iterationsvorschrift folgt

$$\boldsymbol{r} = (\boldsymbol{A} - \mu \boldsymbol{I}) \, \boldsymbol{v} = (\boldsymbol{A} - \mu \boldsymbol{I}) \, \boldsymbol{v}^{k+1} = \boldsymbol{v}^k / \omega_{k+1} - \delta \boldsymbol{A}_k \boldsymbol{v}^{k+1}$$

Wegen $\varepsilon * \omega_{k+1} \geq 1$ und (16) zieht dies

$$\|m{r}\| \leq 1/\omega_{k+1} + arepsilon_1 \leq arepsilon + arepsilon_1 =
uigg\{ \sqrt{n} \, \, F_0 + \left(\sqrt{n} \, + \, 1
ight) \, F_1 igg\} \, \|A\|$$

und wegen 13.1.7 schließlich die Behauptung nach sich.

Zu (ii): Wir skizzieren den Beweis nur. Nach (7) und (16) gilt

$$(\tilde{A} - \mu I) \boldsymbol{w}^2 = \boldsymbol{v}^1 \quad \text{mit} \quad \tilde{A} := A + \delta A_1, \qquad \|\delta A_1\| \leq \varepsilon_1 \tag{17}$$

für das berechnete w^2 . Wir nehmen an, daß δA_1 und damit \tilde{A} symmetrisch ist, siehe dazu 13.4.5 (ii) unten. Es sei $\{\hat{\lambda}_i, \tilde{u}^j\}$ ein orthonormales Eigensystem von \tilde{A} . Wegen (17) und 13.1.2 kann dieses System so indiziert werden, daß $|\lambda_1 - \tilde{\lambda_1}| \leq ||\delta A_1|| \leq \epsilon_1$ gilt, woraus unter Beachtung von (14)

$$|\tilde{\lambda}_1 - \mu| \le |\tilde{\lambda}_1 - \lambda_1| + |\lambda_1 - \mu| \le \varepsilon_1 + \varepsilon_0 = \varepsilon/\sqrt{n}$$
(18)

folgt. Mit $\tilde{U} := (\tilde{u}^1, ..., \tilde{u}^n), \boldsymbol{z} := \tilde{U}^{\mathsf{T}} \boldsymbol{v}^1, \boldsymbol{y} := \tilde{U}^{\mathsf{T}} \boldsymbol{w}^2$ läßt sich nun (17) in der Form $(\tilde{\lambda}_i - \mu) y_i = z_i$ (i = 1, ..., n) schreiben. Im Fall $z_i \neq 0$ muß daher $\tilde{\lambda}_i - \mu \neq 0$ und $y_i = z_i/(\tilde{\lambda}_i - \mu)$ gelten. Daher ist

$$\omega_2^2 = ||m{v}^2||^2 = ||m{y}||^2 = \sum_{i=1}^n y_i^2 \ge \sum_{i:\, z_i \neq 0} z_i^2 / (\tilde{\lambda}_i - \mu)^2 \ge n \left\{ \sum_{i \in I} z_i^2 \right\} / \varepsilon^2,$$

wobei *I* die Menge derjenigen Indizes bezeichnet, für die $z_i \neq 0$ und $|\tilde{\lambda}_i - \mu| \leq \epsilon/\sqrt{n}$ gilt. Wegen (18) gehört mindestens i = 1 zu dieser Indexmenge, sofern $z_1 = \tilde{u}^{1\top}v^1 \neq 0$ ist, was fast immer erwartet werden kann. Wenn v^1 zufällig gewählt worden wäre, ergäbe sich unter einfachen Voraussetzungen an die Verteilung der Erwartungswert $E(|z_1|) = 1/\sqrt{n}$. In den meisten Fällen wird daher

$$|z_1| \ge 1/\sqrt{n} \tag{19}$$

gelten. Aus (19) folgt aber $\varepsilon * \omega_2 \ge \sqrt{n} |z_1| \ge 1$, und die Iteration bricht mit $v = v^2$ nach dem ersten Schritt ab. Der Startvektor v^1 wird aber nicht zufällig gewählt, sondern $v^1 = w^1/||w^1||$ ergibt sich aus w^1 , und w^1 ist Lösung von

$$(\boldsymbol{A} - \mu \boldsymbol{I}) \, \boldsymbol{w}^{\scriptscriptstyle 1} = \boldsymbol{v}^{\scriptscriptstyle 0} := \boldsymbol{P}_{\mu}^{\mathsf{T}} \boldsymbol{L}_{\mu} \boldsymbol{e}$$
 .

Dies entspricht einem zusätzlichen Schritt mit dem speziellen Startvektor v^0 und erhöht die Chance für einen Abbruch mit $v = v^2$ wesentlich. Falls der Abbruchtest in S3 für k = 1 noch nicht erfüllt sein sollte, zeigt eine genauere — allerdings recht komplizierte — Analyse, daß ω_{k+1} in den weiteren Schritten sehr schnell wächst, so daß im ungünstigsten Fall nach 2, 3 Schritten ein akzeptables Paar $\{\mu, v\}$ vorliegt. \Box

13.4.5. Bemerkung. (i) Es ist cond $(A_{\mu}) = |\lambda_n - \mu|/|\lambda_1 - \mu| \ge [|\lambda_n - \mu|/(||A|||F_0)]/\nu$ wegen (14), d. h., für akzeptables F_0 ist cond (A_{μ}) in der Größenordnung $1/\nu$. Die für die Durchführbarkeit der Gaußschen Dreiecksfaktorisierung mit regulärem \mathbf{R}_{μ} hinreichende Bedingung (9) ist daher i. allg. nicht erfüllt, vgl. 5.3.2. Praktisch wird \mathbf{R}_{μ} natürlich fast immer regulär sein, allerdings werden gewisse Pivots $(\mathbf{R}_{\mu})_{jj}$ betragsmäßig sehr klein werden. Sollte ausnahmsweise sogar $(\mathbf{R}_{\mu})_{jj} = 0$ gelten, so wird $(\mathbf{R}_{\mu})_{ji} := \nu ||\mathbf{A}||$ gesetzt, vgl. 5.3.3 (i). Dabei bleibt (7) gültig.

(ii) Die im Beweis von 13.4.4 vorausgesetzte Symmetrie der Störung $\mathcal{O}A_k$ ist keine Einschränkung: Ist (17) mit einer nichtsymmetrischen Störung erfüllt, so folgt nach 4.1.18 für das Residuum die Abschätzung $\|\mathbf{r}^{k+1}\| = \|\mathbf{v}^k - A_{\mu}\mathbf{w}^{k+1}\| \leq \varepsilon_1$. Nach 4.1.20(ii) kann dann $\mathcal{O}A_k$ durch eine äquivalente symmetrische Störung $\mathcal{O}A_k$, für die (17) ebenfalls gilt, ersetzt werden.

(iii) Praktisch wird man ε in der Größenordnung von $\nu \sqrt{n} ||A||$ wählen und dabei ||A|| z. B. durch $||A||_{\infty}$ nach oben abschätzen.

(iv) Wenn mehrere im Sinne von (14) akzeptable Eigenwertnäherungen μ, μ', μ'', \ldots vorliegen, kann 13.4.3 für jede von diesen durchgeführt werden und liefert Eigenvektornäherungen v, v', v'', \ldots Die Eigenpaarapproximationen $\{\mu, v\}, \{\mu', v'\}, \{\mu'', v''\}, \ldots$ sind dann alle exakte Eigenpaare gestörter Matrizen $A + \delta A, A + \delta A', \ldots$ mit $\|\delta A\|, \|\delta A'\|, \ldots \leq \varepsilon \|A\| =: \varepsilon$, also akzeptabel. Da die Störungen $\delta A, \delta A', \ldots$ jedoch i. allg. verschieden sind, brauchen die berechneten Eigenvektoren auch im Fall verschiedener Eigenwerte nicht exakt orthogonal zu sein; nach 13.1.B ist mit einer Abweichung von

$$|\boldsymbol{v}^{\mathsf{T}}\boldsymbol{v}'| \lesssim \|\boldsymbol{\delta}\boldsymbol{A} - \boldsymbol{\delta}\boldsymbol{A}'\|/|\boldsymbol{\mu} - \boldsymbol{\mu}'| \le 2\hat{\boldsymbol{\epsilon}}/|\boldsymbol{\mu} - \boldsymbol{\mu}'| \quad \text{usw.}$$
(20)

zu rechnen. Insbesondere bei dicht benachbarten und speziell bei mehrfachen Eigenwerten als deren Grenzfall ist daher eine Re-Orthogonalisierung nach dem MGS-Verfahren aus 10.1 gemäß

$$oldsymbol{v}' := oldsymbol{v}' - (oldsymbol{v}^{\mathsf{T}}oldsymbol{v}')oldsymbol{v}, \qquad oldsymbol{v}' := oldsymbol{v}'' - (oldsymbol{v}^{\mathsf{T}}oldsymbol{v}')oldsymbol{v}, \qquad oldsymbol{v}'' := oldsymbol{v}'' - ((oldsymbol{v}')^{\mathsf{T}}oldsymbol{v}'')oldsymbol{v}, \qquad oldsymbol{v}'' := oldsymbol{v}'' / ||oldsymbol{v}''|$$

erforderlich. Die nach der Orthogonalisierung vorliegenden Paare $\{\mu, v\}, \{\mu', v'\}, \dots$ sind in fast demselben Sinne akzeptabel wie vorher. \Box

Normalerweise deutet ein Wert von cond (A_{μ}) in der Größenordnung $1/\nu$ darauf hin, daß die berechnete Lösung w von $A_{\mu}w = v$ einen sehr großen Fehler haben kann und daher praktisch wertlos ist, vgl. 4.1.A. Im Fall der inversen Iteration 13.4.3 trifft der erste Teil dieser Behauptung auch zu, allerdings zeigt der Beweis zu 13.4.4, daß sich dieser i. allg. große Fehler praktisch nicht auf die Güte der berechneten Richtung w/||w|| auswirkt, sofern ein numerisch gutartiges Verfahren verwendet wird und w im Sinne von ||w|| = O(1/(v ||A||)) genügend groß ist. Das berechnete wist dann i. allg. eine sehr schlechte Näherung für die exakte Lösung $A_{\mu}^{-1}v$ des Gleichungssystems, aber die Richtung w/||w|| approximiert die von u^1 gut, falls λ_1 ausreichend isoliert ist. Dies ist auch durch die folgenden, mehr geometrischen Überlegungen einzusehen: Für das berechnete w gilt wegen der Gutartigkeit $(A_{\mu} + \delta A) w$ v = v mit $\| \mathbf{d} A \| \leq v F \| A_{\mu} \| =: \epsilon$ und akzeptablem F. Wir betrachten den Fehler $\delta w := w - w^*$ von w bezüglich der exakten Lösung $w^* = A_u^{-1}v$. Nach 4.1.7 liegt der linearisierte Fehler $\sigma w'$ in einem Ellipsoid \mathcal{E}' um den Nullpunkt und mit den Halbachsen $\alpha_i v^j = (\varDelta z / \sigma_i) v^j$, wobei $\varDelta z := \varepsilon \| w^* \|$ gilt und σ_j und v^j die Singulärwerte bzw. Singulärvektoren von A_{μ} sind. Wegen der Symmetrie ist $\sigma_i = |\hat{\lambda}_i - \mu|$ und $v^j = \operatorname{sgn} (\lambda_j - \mu) u^j$. Unter der zusätzlichen Voraussetzung $|\lambda_1 - \mu| \ll |\lambda_2 - \mu|$ ist $\alpha_1 \gg \alpha_2 \ge \cdots \ge \alpha_n$, d. h., das Ellipsoid entartet dann in ein nadelförmiges Gebilde, dessen Längsachse die Richtung des zu λ_1 gehörenden Eigenvektors u^1 hat. Die bezüglich δA linearisierte Lösung $w' := w^* + \delta w'$ liegt dann in dem verschobenen Ellipsoid $w^* + \mathcal{E}'$, siehe Abb. 13.4.2. Der Einfachheit halber lassen wir im folgenden das Zeichen "/" weg.

Man sieht, daß dw fast immer die Richtung von u^1 haben wird. Bei gleichsinniger Orientierung von w^* und dw wirkt sich selbst ein sehr großer Fehler günstig auf die Approximation von u^1 aus, bei ungleichsinniger verschlechtert er diese nur wenig,



Abb. 13.4.2. Fehlerellipsoid bei inverser Iteration

sofern ||w|| genügend groß bleibt. Nur im praktisch unwahrscheinlichen Sonderfall $\delta w \approx -w^*$ — also für kleines ||w|| — weicht die durch w festgelegte Richtung wesentlich von u^1 ab. Falls ein Lösungsverfahren verwendet wird, das nur numerisch stabil mit derselben Konstanten F ist, müßte das Ellipsoid durch die umschriebene Kugel ersetzt werden. Es könnten dann auch große Fehler orthogonal zu u^1 auftreten, die zu einem völlig unbrauchbaren w führen.

Wenn λ_1 ein *p*-facher Eigenwert ist, entartet \mathscr{E}' etwa im Fall p = 2 in ein scheibenförmiges Gebilde, das sich nur wenig von einem Kreis mit dem großen Radius α_1 in der durch $\{u^1, u^2\}$ aufgespannten Ebene — dem Eigenraum zu λ_1 — unterscheidet. Auch dann approximiert w fast immer einen in dieser Ebene liegenden Eigenvektor zu λ_1 gut. Dasselbe ist der Fall, wenn λ_1 zu einem aus p benachbarten Eigenwerten bestehenden Eigenwerthaufen gehört, der von den restlichen genügend getrennt ist.

Falls μ keine gute Eigenwertapproximation ist, liegt es nahe, während der inversen Iteration 13.4.1 die Näherung μ sukzessive mit zu verbessern. Da bei festem $\boldsymbol{v} = \boldsymbol{v}^k$ der Rayleigh-Quotient $\varrho(\boldsymbol{v}^k)$ nach 13.1.8 die im Sinne einer minimalen Residuumsnorm optimale Eigenwertnäherung ist, bietet sich dazu die Festlegung $\mu = \mu_k = \varrho(\boldsymbol{v}^k)$ an. Dies liefert die *Rayleigh-Quotienten-Iteration* — kurz: *RQ-Iteration* — genannte Modifikation der inversen Iteration.

13.4.6. Rayleigh-Quotienten-Iteration.

- S0 (Initialisierung): Wähle v^1 mit $||v^1|| = 1$ und $\varepsilon > 0$ in der Größenordnung $v \sqrt{n} ||A||$, setze k := 1
- S1 (Iteration)
- S1.1 (Berechnung der Verschiebung ϱ_k): Bestimme $\varrho_k := \varrho(\boldsymbol{v}^k) = \boldsymbol{v}^{k\intercal} A \boldsymbol{v}^k$
- S1.2 (Inverse Iteration mit $\mu = \varrho_k$): Berechne w^{k+1} aus $(A \varrho_k I) w^{k+1} = v^k$

S2 (Normierung): Berechne $\omega_{k+1} := \| \boldsymbol{w}^{k+1} \|$, setze $\boldsymbol{v}^{k+1} := \boldsymbol{w}^{k+1} / \omega_{k+1}$

S3 (Abbruchtest): if $\varepsilon * \omega_{k+1} < 1$ then [k := k + 1, goto S1]

Aufwand pro Schritt: $\sim n^2$ opms (für S1.1) + $\sim K_1 n^3$ opms (für S1.2), falls A voll besetzt ist.

Im Unterschied zu 13.4.1 bzw. 13.4.3 sind bei der RQ-Iteration die Koeffizientenmatrizen der in jedem Schritt zu lösenden Gleichungssysteme verschieden, so daß die Dreiecksfaktorisierungen jedesmal neu berechnet werden müssen. In den Anwendungen ist A jedoch meist von Bandgestalt und insbesondere tridiagonal, so daß dies nur $\sim Kn$ opms kostet, vgl. 6.4.A.

13.4.7. Aussage. Die RQ-Iteration 13.4.6 werde mit $\varepsilon = 0$ in exakter Arithmetik durchgeführt. Dann gilt:

(i) Die Normen der Residuen $r^k := (A - \varrho_k I) v^k$ fallen monoton, d. h., es ist

$$\|\boldsymbol{r}^{k+1}\| \leq \|\boldsymbol{r}^k\|,\tag{21}$$

und die Folge $\{\varrho_k, \vartheta_k v^k\}$ mit $\vartheta_k \in \{+1, -1\}$ so, daß $(\vartheta_k v^k)^{\mathsf{T}} u^j \ge 0$ gilt, konvergiert fast immer gegen ein Eigenpaar $\{\lambda_j, u^j\}$ von A.

(ii) Falls $\{\varrho_k\}$ gegen den Eigenwert λ_j konvergiert, ist die Konvergenz kubisch im Sinne der Gültigkeit von

$$\tan \varphi_{k+1} \leq \eta_1 (\tan \varphi_k)^3 \quad \text{mit} \quad \eta_1 \approx 2 \|A\|/\gamma \tag{22}$$

für $k \ge k_0, \, k_0$ hinreichend groß, wobei $\varphi_k := \measuredangle (u^j, \vartheta_k v^k)$ und

$$\gamma := \min \left\{ |\lambda_i - \lambda_i| \colon \lambda_i \neq \lambda_i \right\}$$
(23)

ist. Es gilt zudem

$$|\lambda_j - \varrho_k| \le 2 \, \|A\| \tan^2 \varphi_k \tag{24}$$

sowie

$$\|\boldsymbol{r}^{\boldsymbol{k}}\| \leq \eta_2 \tan \varphi_{\boldsymbol{k}} \quad \text{mit} \quad \eta_2 \approx 2 \, \|\boldsymbol{A}\| \,. \tag{25}$$

Beweis. Wir beginnen mit (ii): Es gelte also lim $\varrho_k = \lambda_j$, und o. B. d. A. sei $\lambda_j = \lambda_1$. Falls φ_k wie oben definiert ist, liefert (13.3.19) gerade die Abschätzung (24). Zum Beweis von (22) beachten wir, daß \boldsymbol{v}^{k+1} aus \boldsymbol{v}^k durch einen Schritt der Vektoriteration mit $\hat{\boldsymbol{A}} := (\boldsymbol{A} - \varrho_k \boldsymbol{I})^{-1}$ entsteht. Nach 13.3.3 gilt daher

$$\tan \varphi_{k+1} \le \varkappa \tan \varphi_k \tag{26}$$

mit

$$\varkappa := \frac{\max\{|\hat{\lambda}_{i}|: i \neq 1\}}{|\hat{\lambda}_{1}|} = \frac{|\lambda_{1} - \varrho_{k}|}{\min\{|\lambda_{i} - \varrho_{k}|: i \neq 1\}}.$$
(27)

Es sei zunächst λ_1 einfach, also $\gamma = \min \{|\lambda_1 - \lambda_i| : i \neq 1\} > 0$, und k_0 sei so groß, daß $|\lambda_1 - \varrho_k| < \gamma$ für $k \ge k_0$ gilt. Dann folgt aus (26), (27) und $|\lambda_i - \varrho_k| \ge |\lambda_1 - \lambda_i| - |\lambda_1 - \varrho_k| \ge \gamma - |\lambda_1 - \varrho_k| > 0$

$$an arphi_{k+1} \leq rac{2 \, \|oldsymbol{A}\| \, (an arphi_k)^3}{\gamma - 2 \, \|oldsymbol{A}\| \, an^2 arphi_k},$$

also (22). Zum Nachweis von (25) stellen wir \mathbf{r}^k in der Form

$$\boldsymbol{\vartheta}_k \boldsymbol{r}^k = (\boldsymbol{A} - \varrho_k \boldsymbol{I}) \, \hat{\boldsymbol{v}}^k = \boldsymbol{A}(\hat{\boldsymbol{v}}^k - \boldsymbol{u}^1) + \lambda_1 (\boldsymbol{u}^1 - \hat{\boldsymbol{v}}^k) + (\lambda_1 - \varrho_k) \, \hat{\boldsymbol{v}}^k$$

mit $\hat{\boldsymbol{v}}^k := \boldsymbol{\vartheta}_k \boldsymbol{v}^k$ dar. Mit (13.3.15) und (24) ergibt sich dann

$$\|m{r}^k\|\leq 2 \|m{A}\| \|m{\hat{v}}^k-m{u}^1\|+|\lambda_1-arrho_k|\leq 2 \|m{A}\| (an arphi_k+ an^2arphi_k),$$

also (25). Falls λ_1 mehrfach ist, kann analog mit γ gemäß (23) geschlossen werden, vgl. 13.3.7 (i).

Zu (i): Für die Residuumsnormen gilt wegen der Optimalität des Rayleigh-Quotienten

$$\begin{aligned} |\mathbf{r}^{k+1}| &= \|(\mathbf{A} - \varrho_{k+1}\mathbf{I})\,\mathbf{v}^{k+1}\| \leq \|(\mathbf{A} - \varrho_k\mathbf{I})\,\mathbf{v}^{k+1}\| = |\mathbf{v}^{k\mathsf{T}}(\mathbf{A} - \varrho_k\mathbf{I})\,\mathbf{v}^{k+1}| \\ &= |\mathbf{v}^{k+1\mathsf{T}}(\mathbf{A} - \varrho_k\mathbf{I})\,\mathbf{v}^k| \leq \|(\mathbf{A} - \varrho_k\mathbf{I})\,\mathbf{v}^k\| = \|\mathbf{r}^k\|. \end{aligned}$$

Das zweite Gleichheitszeichen folgt aus der Tatsache, daß $\boldsymbol{y}^{k} := (\boldsymbol{A} - \varrho_{k}\boldsymbol{I})\boldsymbol{v}^{k+1} = \boldsymbol{v}^{k}/\omega_{k+1}$ gilt, also \boldsymbol{y}^{k} gleich seiner Projektion auf \boldsymbol{v}^{k} ist. Zur Konvergenzuntersuchung beachten wir, daß neben $||\boldsymbol{v}^{k}|| = 1$ noch $|\varrho_{k}| \leq ||\boldsymbol{A}||$ gilt, d. h., die Paare $\{\varrho_{k}, \boldsymbol{v}^{k}\}$ liegen in einer beschränkten und abgeschlossenen Menge. Es gibt dann nach dem Satz von WEIERSTRASS eine konvergente Teilfolge $\{\varrho_{k_{i}}, \boldsymbol{v}^{k_{i}}\}$ mit dem Grenzwert $\{\bar{\varrho}, \bar{\boldsymbol{v}}\}, |\bar{\varrho}| \leq ||\boldsymbol{A}||, ||\bar{\boldsymbol{v}}|| = 1$, und aus Stetigkeitsgründen gilt

$$\lim_{i \to \infty} ||(\boldsymbol{A} - \varrho_{\boldsymbol{k}_l} \boldsymbol{I}) \boldsymbol{v}^{\boldsymbol{k}_l}|| = \lim_{i \to \infty} ||\boldsymbol{r}^{\boldsymbol{k}_l}|| = ||(\boldsymbol{A} - \bar{\varrho} \boldsymbol{I}) \, \bar{\boldsymbol{v}}|| =: \sigma.$$
(28)

Wegen der Monotonie der Residuen existiert nun stets der Grenzwert

$$\lim_{k \to \infty} \|\boldsymbol{r}^k\| = \sigma \ge 0. \tag{29}$$

Falls $\sigma = 0$ ist, folgt $||(A - \bar{\varrho}I)\bar{v}|| = 0$ aus (28) und (29), d. h., $\{\bar{\varrho}, \bar{v}\} = \{\lambda_j, u^j\}$ ist ein Eigenpaar von A. Es gibt dann ein genügend großes $l := k_i$, so daß $\eta_1 \tan^2 \varphi_l \leq q < 1$ gilt. Aus dem Beweis zu (ii) folgt dann

$$an \varphi_{k+1} \leq [\eta_1 \tan^2 \varphi_k] \tan \varphi_k \leq q \tan \varphi_k$$

für alle $k \ge l$, d. h., die gesamte Folge $\{\varrho_k, \vartheta_k v^k\}$ konvergiert gegen $\{\lambda_j, v^j\}$. Der praktisch nicht auftretende Fall $\sigma > 0$ ist komplizierter zu analysieren, wir verweisen auf die Literatur, siehe B 13.6. \Box

13.4.8. Bemerkung. (i) Bei exakter Realisierung von 13.4.6 kann durch Zufall $A - \varrho_k I$ singulär sein. Dann ist ϱ_k ein Eigenwert von A, und ein zugehöriger Eigenvektor läßt sich aus dem homogenen Gleichungssystem $(A - \varrho_k I) v = o$ berechnen. Bei Computerrealisierung wird dies praktisch nicht auftreten. Wenn dabei doch $A - \varrho_k I$ numerisch singulär sein sollte, also kein von 0 verschiedenes Pivot bei der Dreiecksfaktorisierung gefunden werden kann, ist wie in 13.4.5 (i) zu verfahren.

(ii) Bei Computerrechnung wird das durch 13.4.2 beschriebene ideale Konvergenzverhalten durch die Rundungsfehler gestört, wegen 13.4.4 jedoch nicht verschlechtert, sondern eher verbessert. Es ist also auch bei Computerrealisierung praktisch immer mit schneller Konvergenz zu rechnen.

(iii) Im allgemeinen kann nichts darüber ausgesagt werden, gegen welchen Eigenwert λ_i die Rayleigh-Quotienten ϱ_k konvergieren. Insbesondere braucht die Konvergenz nicht gegen denjenigen Eigenwert zu erfolgen, der minimalen Abstand von ϱ_1 hat. \Box

Übungsaufgabe

Ü 13.4.1. Man zeige: Wenn bei der RQ-Iteration $||\boldsymbol{r}^{k+1}|| = ||\boldsymbol{r}^{k}||$ ist, gilt $\varrho_{k} = \varrho_{k+1}$, und \boldsymbol{v}^{k} ist ein Eigenvektor von $(\boldsymbol{A} - \varrho_{k}\boldsymbol{I})^{2}$.

13.5. Orthogonale Ähnlichkeitstransformation auf Tridiagonalform

Die im vorherigen Abschnitt behandelten Varianten der inversen Iteration wurden direkt mit der Matrix $A \in \mathbf{S}^{n,n}$ des zu lösenden Eigenwertproblems $Ax = \lambda x$ durchgeführt. Eine alternative Vorgehensweise besteht darin, das an anderen Stellen bereits praktizierte Prinzip der Transformation des Ausgangsproblems in ein äquivalentes Problem einfacherer Gestalt auch für die Lösung des Eigenwertproblems auszunützen. Als Transformation kommt dabei eine orthogonale Ähnlichkeitstransformation in Frage, bei der Symmetrie und Eigenwerte erhalten bleiben, vgl. 13.1.A.

Da eine Ähnlichkeitstransformation auf Diagonalform in endlich vielen Schritten i. allg. nicht möglich ist — die Diagonalelemente wären dann bereits die gesuchten
Eigenwerte -, bietet sich als Ziel eine Tridiagonalmatrix

$$\boldsymbol{T} = \begin{pmatrix} a_1 & b_2 \\ b_2 & a_2 & b_3 \\ \ddots & \ddots & \ddots \\ b_{n-1} & a_{n-1} & b_n \\ & & b_n & a_n \end{pmatrix} =: \text{trid } (a_1, \dots, a_n, b_2, \dots, b_n)$$

an. Man beachte, daß A durch $\sim n^2/2$ Daten $\{a_{ij}: i \leq j\}$ repräsentiert wird, während T durch nur $\sim 2n$ Daten $\{a_1, \ldots, a_n, b_2, \ldots, b_n\}$ festgelegt ist. Beim Übergang von A zu T tritt daher eine wesentliche Datenreduktion ein. Durch die Transformation

$$A \to T = Q^{\mathsf{T}} A Q, \qquad Q \in \mathsf{R}^{n,n} \text{ orthogonal},$$
 (1)

wird also das Eigenwertproblem einer allgemeinen symmetrischen Matrix A auf das einer symmetrischen Tridiagonalmatrix T zurückgeführt. Das letztere läßt sich jedoch wesentlich effektiver lösen, wie wir in den Abschnitten 13.6 und 13.7 sehen werden.

A. Transformation mittels Householder-Spiegelungen

Im folgenden wird gezeigt, daß die Transformation (1) nach der Vorschrift

$$A^{(1)} := A, \quad A^{(k+1)} := H_k A^{(k)} H_k \quad (k = 1, ..., n-2), \quad T := A^{(n-1)}$$
 (2)

unter Verwendung von Householder-Spiegelungen bzw. Einheitsmatrizen H_k in n-2 Schritten durchgeführt werden kann, wobei im k-ten Schritt die Nullen in der k-ten Spalte und Zeile erzeugt werden. Es sei also $A^{(k)}$ bereits von der Form

$$A^{(k)} = \begin{pmatrix} T^{(k)} & | \\ & | \\ \hline & | \\ a^{k} & | \\ B^{(k)} \end{pmatrix}$$
(3)

 mit

$$T^{(k)} = \text{trid} (a_1, ..., a_k, b_2, ..., b_k) \in S^{k,k}, \quad a^k \in \mathbb{R}^{n-k}, \quad B^{(k)} \in S^{n-k, n-k}.$$

Im Fall $a^k = o$ sind die geforderten Nullen bereits vorhanden, so daß

$$H_k := I \tag{4}$$

gewählt werden kann. Andernfalls setzen wir H_k als Householder-Spiegelung

$$\boldsymbol{H}_{\boldsymbol{k}} := \left(\begin{array}{c|c} \boldsymbol{I}_{\boldsymbol{k}} & | \\ \hline & | & \boldsymbol{\overline{H}}_{\boldsymbol{k}} \end{array} \right) \tag{5}$$

an und bestimmen

$$\overline{H}_{k} := I - \overline{v}^{k} \overline{v}^{k\top} / \overline{\gamma}_{k} \in S^{n-k, n-k}, \qquad \overline{\gamma}_{k} := \overline{v}^{k\top} \overline{v}^{k} / 2$$
(6)

so, daß

$$oldsymbol{a}^{k} = egin{pmatrix} a_{k+1,k}^{(k)} \ a_{k+2,k}^{(k)} \ dots \ a_{nk}^{(k)} \end{pmatrix} ext{ in } oldsymbol{ar{a}}^{k} := oldsymbol{ar{H}}_{k} oldsymbol{a}^{k} = arrho_{k} oldsymbol{e}^{1} = egin{pmatrix} arrho_{k} \ 0 \ dots \ 0 \ dots \ 0 \end{pmatrix}$$

übergeht. Dies ist bei der Festlegung

$$\overline{\boldsymbol{v}}^{k} = (\boldsymbol{v}_{k+1,k}, \dots, \boldsymbol{v}_{nk})^{\mathsf{T}} := \boldsymbol{e}^{1} - \boldsymbol{a}^{k}/\varrho_{k}$$

$$\tag{7}$$

mit

$$\varrho_k := \begin{cases} \|\boldsymbol{a}^k\| & \text{für } \boldsymbol{a}_{k+1,k}^{(k)} \leq 0\\ -\|\boldsymbol{a}^k\| & \text{für } \boldsymbol{a}_{k+1,k}^{(k)} > 0, \end{cases}$$
(8)

der Fall; es gilt dann $\bar{\gamma}_k = v_{k+1,k}$, vgl. 3.3.A und Abschnitt 10.2. Für $A^{(k+1)}$ ergibt sich damit

$$A^{(k+1)} = H_k A^{(k)} H_k = \begin{pmatrix} T^{(k)} & & \\ & & \bar{a}^{k} \\ \hline & & \bar{a}^{k} & \bar{B}^{(k)} \end{pmatrix} = \begin{pmatrix} T^{(k+1)} & & \\ & & a^{k+1} \\ \hline & & a^{k+1} & B^{(k+1)} \end{pmatrix}$$

mit

$$\overline{\boldsymbol{B}}^{(k)} := \overline{\boldsymbol{H}}_k \boldsymbol{B}^{(k)} \overline{\boldsymbol{H}}_k. \tag{9}$$

Im Unterschied zur Householder-Faktorisierung aus 10.2 tritt die Spiegelung \overline{H}_k in (9) auf beiden Seiten von $B^{(k)}$ auf, und $\overline{B}^{(k)}$ ist wie $B^{(k)}$ symmetrisch. Es ist daher nicht zweckmäßig, die transformierte Matrix \overline{B} gemäß $\overline{B} := (\overline{H}B) \overline{H}$ oder $\overline{B} := \overline{H}(B\overline{H})$ zu berechnen, denn $\overline{H}B$ und $B\overline{H}$ sind i. allg. nicht symmetrisch; zur Vereinfachung ignorieren wir den Index k. Die Symmetrie wird dagegen voll ausgenutzt, wenn (9) unter Beachtung von (6) in der Form

$$\overline{B} = \overline{H}B\overline{H} = B - \overline{v}\overline{v}^{\mathsf{T}}B/\overline{\gamma} - B\overline{v}\overline{v}^{\mathsf{T}}/\overline{\gamma} + \overline{v}\overline{v}^{\mathsf{T}}B\overline{v}\overline{v}^{\mathsf{T}}/\overline{\gamma}^{2}
= B - \overline{v}\overline{w}^{\mathsf{T}} - \overline{w}\overline{v}^{\mathsf{T}} + [\overline{v}^{\mathsf{T}}\overline{w}/\overline{\gamma}] \overline{v}\overline{v}^{\mathsf{T}}$$
(10)

 mit

$$\overline{\boldsymbol{w}} := \boldsymbol{B}\overline{\boldsymbol{v}}/\overline{\boldsymbol{\gamma}} \tag{11}$$

geschrieben wird. Die gesonderte Berechnung des dritten dyadischen Produktes in (10) kann schließlich auch noch vermieden werden, indem dieses zur Hälfte den beiden anderen Produkten zugeschlagen wird, was auf

$$\overline{B} := B - (\overline{p}\overline{v}^{\mathsf{T}} + \overline{v}\overline{p}^{\mathsf{T}}) \tag{12}$$

 \mathbf{mit}

$$\overline{p} := \overline{w} - [\overline{v}^{\mathsf{T}} \overline{w} / (2\overline{\gamma})] \overline{v} \tag{13}$$

führt.

Nach n-2 derart ausgeführten Schritten ergibt sich die gesuchte Tridiagonalmatrix

mit $a^{n-1} = a^{n-1\top} =: b_n \in \mathbb{R}$ und $B^{(n-1)} =: a_n \in \mathbb{R}$, und es gilt

$$\boldsymbol{\Gamma} = \boldsymbol{H}_{\boldsymbol{n}-2} \cdots \boldsymbol{H}_{2} \boldsymbol{H}_{1} \boldsymbol{A} \boldsymbol{H}_{1} \boldsymbol{H}_{2} \cdots \boldsymbol{H}_{\boldsymbol{n}-2},$$

also (1) mit der durch

$$\boldsymbol{Q} := \boldsymbol{H}_1 \boldsymbol{H}_2 \cdots \boldsymbol{H}_{n-2} \tag{14}$$

definierten orthogonalen Matrix Q. Zusammenfassend erhalten wir den folgenden Algorithmus.

13.5.1. Tridiagonalisierung mittels Householder-Spiegelungen. Für gegebenes $A \in S^{n,n}$ werde die Transformation (2) mit H_k gemäß (4) oder (5) bis (8) in exakter Arithmetik durchgeführt. Dann gilt

$$T = Q^{\mathsf{T}} A Q$$

mit der Tridiagonalmatrix $T := A^{(u-1)}$ und der Produktdarstellung (14) für die orthogonale Transformationsmatrix Q.

Aufwand: $\sim 2n^3/3$ opms $+ \sim n$ opr für **T**, falls **A** voll besetzt ist und $\overline{B}^{(k)}$ gemäß (11) bis (13) aus **B**^(k) berechnet wird.

13.5.2. Bemerkung. (i) Algorithmus 13.5.1 kann auf dem Platz des unteren Dreiecks von A ausgeführt werden. Falls H_k als Householder-Spiegelung gewählt wird, werden die Elemente a_{ik} mit den Komponenten $\{v_{ik}: i \ge k+1\}$ des Vektors \overline{v}^k überspeichert. Es gilt dann $v_{k+1,k} = \overline{\gamma}_k \ge 1$, vgl. 3.3.5. Falls $H_k = I$ gesetzt wurde, ordnen wir H_k den Vektor $\overline{v}^k = o$ zu. Da dies nur für $a^k = o$ eintreten kann, sind die entsprechenden Nullen bereits vorhanden. Insbesondere ist dann $v_{k+1,k} = 0$, so daß beide Fälle zuverlässig rekonstruiert werden können. Die Subdiagonalelemente $b_{k+1} = \varrho_k$ von T müssen dann in einem gesonderten Feld der Länge n aufgehoben werden. Der Aufwand im k-ten Teilschritt wird für voll besetztes A durch die Berechnung von \overline{w} gemäß (11) und die Bildung von \overline{B} gemäß (12) mit zusammen $\sim 2(n - k)^2$ opms bestimmt; man beachte, daß nur ein Dreieck von \overline{B} berechnet zu werden braucht.

(ii) Die Berechnung von $Q^{\mathsf{T}}b$ bzw. Qf aus der durch die $\{v_{ik}: i \geq k+1\}$ repräsentierten Produktform (14) ist analog zu 10.2.6 mit jeweils $\sim n^2$ opms möglich; dabei sind die Fallunterscheidungen aus (i) zu berücksichtigen. Falls Q explizit benötigt wird, kann die Berechnung mit zusätzlichen $\sim 2n^3/3$ opms analog zu Ü 10.2.2 erfolgen. \Box

Wir geben ohne Beweis die Ergebnisse der Rundungsfehleranalyse an.

13.5.3. Rundungsfehleranalyse. Algorithmus 13.5.1 werde für symmetrisches $A \in \Re^{n,n}$ durchgeführt, und $T \in \Re^{n,n}$ sei die berechnete symmetrische Tridiagonalmatrix. Dann gilt

$$\mathbf{A} + \boldsymbol{\delta} \mathbf{A} = \boldsymbol{Q} \boldsymbol{T} \boldsymbol{Q}^{\mathsf{T}} \tag{15}$$

mit exakt orthogonalem $Q \in \mathbb{R}^{n,n}$ und einer Störung $dA \in \mathbb{S}^{n,n}$, für die

$$\|\delta A\|_F \le \nu F_0 \|A\|_F \quad \text{mit} \quad F_0 := 4.2n^2 \tag{16}$$

und

$$\| \boldsymbol{\delta} \boldsymbol{A} \| \leq \nu F_1 \| \boldsymbol{A} \| \quad ext{mit} \quad F_1 := 0.4 n^{5/2} (1 + 9.5/n^{1/2})$$

gilt.

13.5.4. Bemerkung. (i) Die Matrix Q in (15) ist das exakte Produkt (14) der Matrizen H_k , die den berechneten $A^{(k)}$ in exakter Arithmetik zugeordnet sind. Die exakte Matrix Q wird sowohl durch die mittels der berechneten $\{v_{ik}\}$ definierte Produktform (14) als auch durch die explizit analog zu Ü 10.2.2 berechnete Matrix $Q = \text{fl}(H_1H_2\cdots H_{n-2})$ gut repräsentiert, vgl. 10.2 und 11.1 für analoge Überlegungen.

(ii) Die Matrix T und die Spiegelungen H_k und damit deren Produkt Q sind stetige Funktionen von A, sofern der Fall $H_k = I$ ausgeschlossen wird, also $||a^k|| = |\varrho_k|$ $= |b_{k+1}| \neq 0$ (k = 1, ..., n - 2) gilt. Trotzdem sind die berechneten Faktoren Tund Q oft sehr empfindlich gegenüber Störungen von A und Änderungen der Computerarithmetik, siehe B 13.7 für Hinweise auf eindrucksvolle Beispiele in der Spezialliteratur. Die möglicherweise stark differierenden berechneten Matrizen T sind jedoch alle exakt orthogonal ähnlich zu $A + \delta A$ mit kleinem δA , d. h., es liegt numerische Gutartigkeit vor, und die Eigenwerte von T approximieren die von A gut.

(iii) Ist A eine im Sinne von 13.1.B absteigend gestufte Matrix, also etwa

$$A = \begin{pmatrix} 100 & 10 & 1 \\ 10 & 10 & 1 \\ 1 & 1 & 1 \end{pmatrix},$$

so sind T und δA bei Rechnung nach 13.5.1 meistens in derselben Weise gestuft, so daß auch die beitragskleinen Eigenwerte von A durch T ausreichend repräsentiert werden. Ist A jedoch in inverser Weise gestuft, also etwa

$$A = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 10 & 10 \\ 1 & 10 & 100 \end{pmatrix},$$

so gehen beim Übergang zu T mittels 13.5.1 Informationen über die betragskleinen Eigenwerte verloren. In diesem Fall sollte mit

$$ar{m{A}}:=m{P}m{A}m{P}, \qquad m{P}=m{P}^{\intercal}:=egin{pmatrix} 0&0&1\0&1&0\1&0&0 \end{pmatrix}, ext{ allgemein } m{P}=\sum\limits_{j=1}^nm{e}^{n-j+1}m{e}^{j\intercal}$$

gearbeitet werden, oder 13.5.1 sollte in inverser Reihenfolge realisiert werden, d. h., im ersten Schritt sollten die Nullen in der n-ten Spalte und Zeile erzeugt werden, siehe Ü 13.5.2. □

B. Tridiagonalisierung mittels Givens-Drehungen und Anwendung auf Bandmatrizen

In Analogie zum Vorgehen bei der QR-Faktorisierung in 10.2 kann die Householder-Spiegelung H_k durch ein Produkt von n - k - 1 Givens-Drehungen ersetzt werden: Die Vorschrift $A^{(k+1)} := H_k A^{(k)} H_k$ wird ersetzt durch

$$A^{(k+1)} := G_k^{\mathsf{T}} A^{(k)} G_k := [G_{k+1,k+2} \cdots (G_{n-1,n} A^{(k)} G_{n-1,n}^{\mathsf{T}}) \cdots G_{k+1,k+2}^{\mathsf{T}}]; \qquad (17)$$

selbstverständlich sind auch hier andere Drehungsindizes möglich, siehe Ü 13.5.3. Die Drehung $G_{i-1,i}$ wird dabei so gewählt, daß $G_{i-1,i}$ $(k+2 \le i \le n)$ bei Anwendung auf die aktuelle Transformierte in der Position $\{i, k\}$ die gewünschte Null erzeugt.

13.5.5. Bemerkung. (i) Bei voll besetztem A und expliziter Realisierung der Givens-Drehungen ist der Aufwand $\sim n^3$ opms $+ \sim n^2/2$ opr, also etwa um die Hälfte höher als in 13.5.1, wenn die Bildung von $\bar{A} := G_{i-1,i}AG_{i-1,i}^{\mathsf{T}}$ für das jeweils aktuelle A analog zu Ü 3.3.6 erfolgt. Man beachte, daß Drehungsindizes und Drehungsparameter $\{c, s\}$ beim Jacobi-Verfahren jedoch anders festgelegt sind als hier. Diese auf GIVENS zurückgehende, historisch älteste Version der Tridiagonalisierung - siehe B 13.7 scheidet daher aus Aufwandsgründen aus. Bei impliziter Realisierung sind $\sim 2n^3/3$ opms wie bei 13.5.1 erforderlich, siehe wieder B 13.7 für weiterführende Hinweise zur Implementierung.

(ii) Falls A schwach besetzt ist, zeigen Givens-Drehungen i. allg. deutliche Vorteile gegenüber Householder-Spiegelungen, da nur die Nichtnullelemente (= NNE) annulliert zu werden brauchen, was mittels Givens-Drehungen effektiver möglich ist. Wir demonstrieren dies im folgenden für den wichtigen Sonderfall einer Bandmatrix. Man beachte, daß die Bandgestalt durch 13.5.1 sukzessive zerstört wird. 🗌

Es sei also $A \in S^{n,n}$ eine Bandmatrix des Typs $\{s, s\}$ mit der Bandbreite l = 2s + 1, vgl. 6.4.A, d. h., A habe die Gestalt

$$A = \begin{pmatrix} \ddots & & \\ &$$

Wir erläutern das Vorgehen am Beispiel einer Matrix mit n = 10, s = 3, vgl. das Schema auf S. 403, und betrachten dazu nur den ersten Hauptschritt k = 1 bei insitu-Realisierung. Im ersten Schritt wird G_{34} so festgelegt, daß $A := G_{34}AG_{34}^{\mathsf{T}}$ in der Position {4, 1} (und der jeweils symmetrischen Position, wir erwähnen diese nicht gesondert) eine 0 hat. Dabei wird in der Position $\{7, 3\}$ außerhalb des Bandes das mit $\mathbf{2}$

bezeichnete NNE erzeugt. Um die Bandgestalt wieder herzustellen, wird dieses

Element durch eine zweite Drehung G_{67} nach der Vorschrift $A := G_{67}AG_{67}^{\mathsf{T}}$ anulliert, wobei allerdings in {10, 6} ein zweites NNE entsteht, nämlich 3. Dieses wird durch eine weitere Zusatzdrehung $G_{9,10}$ annulliert, womit die Bandgestalt wieder vorliegt. Im nächsten Teilschritt wird in {3, 1} mittels G_{23} eine 0 erzeugt, und die in {6, 2} bzw. {9, 5} entstehenden NNE außerhalb des Bandes werden mittels G_{56} bzw. G_{89} annulliert usw.



Im allgemeinen Fall lautet eine Grobbeschreibung des Verfahrens wie folgt:

13.5.6. Tridiagonalisierung einer Bandmatrix mittels Givens-Drehungen. Für die Bandmatrix $A \in S^{n,n}$ vom Typ $\{s, s\}$ werde der folgende Algorithmus in exakter Arithmetik ausgeführt:

for k := 1(1) n - 2 do k-ter Hauptschritt: for $i := \min \{k + s, n\} (-1) k + 2$ do S1 (Annullierung von a_{ik}): Lege $G_{i-1,i}$ so fest, daß $(G_{i-1,i}A)_{ik} = 0$ Bilde $A := G_{i-1,i}AG_{i-1,i}^{\mathsf{T}}$ S2 (Wiederherstellung der Bandgestalt): p := ifor p := p + s while $p \le n$ do S3 (Annullierung von $a_{p,p-s-1}$): Lege $G_{p-1,p}$ so fest, daß $(G_{p-1,p}A)_{p,p-s-1} = 0$ Bilde $A := G_{p-1,p}AG_{p-1,p}^{\mathsf{T}}$

Dann ist A mit der Tridiagonalmatrix T = trid $(a_1, \ldots, a_n, b_2, \ldots, b_n)$ überspeichert, und T ist orthogonal ähnlich zu A.

13.5.7. Bemerkung. (i) Bei expliziter Realisierung der Drehungen analog zu Ü 3.3.6 ist der Aufwand $\sim 3(s-1) n^2$ opms $+ \sim [(s-1)/(2s)] n^2$ opr, wenn nur T berechnet

wird. Im Fall $s \leq n/5$ ist 13.5.6 also billiger als 13.5.1. Mittels impliziter Givens-Drehungen läßt sich der Aufwand etwa um ein Drittel reduzieren.

(ii) Da n bei einer Bandmatrix meist groß ist, werden nur wenige Eigenwerte zu berechnen sein. Diese können aus T nach dem Bisektionsverfahren 13.6.8 bestimmt werden. Zugehörige Eigenvektoren können dann mittels inverser Iteration nach 13.4.3 ermittelt werden. Da das Speichern der Transformationsmatrix O aufwendig ist, sollte die inverse Iteration direkt mit der originalen Bandmatrix A unter Verwendung einer nach 6.4.A berechneten spaltenpivotisierten Gauß-Faktorisierung erfolgen, vgl. auch 13.6.11(vi).

C. Eindeutigkeit der Faktoren T und O

Für gewisse Realisierungen des **QR**-Algorithmus aus 13.7 ist die nachfolgende Eindeutigkeitsaussage von Bedeutung.

13.5.8. Satz. Für
$$A \in S^{n,n}$$
 gelte

 $A = QTQ^{\intercal} = \overline{Q}\overline{T}\overline{Q}^{\intercal}$ (19) mit Tridiagonalmatrizen $T = \text{trid}(a_i, b_i), \ \overline{T} = \text{trid}(\overline{a}_i, \overline{b}_i)$ und orthogonalen Matrizen $Q = (q^1, ..., q^n), \ \overline{Q} = (\overline{q}^1, ..., \overline{q}^n).$ Wenn die Bedingung

$$b_i \neq 0$$
 $(i = 2, ..., n)$ (20)

erfüllt ist und $q^1 = Q$ gilt, ist $ar{Q} = Q h$

$$\boldsymbol{q}^1 = \boldsymbol{Q}\boldsymbol{e}^1 = \bar{\boldsymbol{Q}}\boldsymbol{e}^1 = \bar{\boldsymbol{q}}^1 \tag{21}$$

$$\bar{Q} = QD$$
 und $\bar{T} = DTD$ (22)

mit

 $D = ext{diag}(d_i), \quad d_1 = 1, \quad |d_i| = 1 \quad (i = 2, ..., n).$ (23)

Beweis. Aus (19) folgt

$$\overline{T}P = PT$$
 mit orthogonalem $P = (p^1, ..., p^n) := \overline{Q}^{\mathsf{T}}Q$, (24)

und wegen (21) ist $e^1 = \bar{Q}^{\mathsf{T}} Q e^1 = P e^1 = p^1$. Die erste Spalte der Identität (24) lautet daher

$$\overline{T}Pe^{1} = \overline{T}p^{1} = \overline{T}e^{1} = \overline{a}_{1}e^{1} + \overline{b}_{2}e^{2} = PTe^{1} = a_{1}p^{1} + b_{2}p^{2} = a_{1}e^{1} + b_{2}p^{2}.$$
 (25)

Linksmultiplikation mit $e^{1\intercal}$ führt auf $\bar{a}_1 = a_1$, also $\bar{b}_2 e^2 = b_2 p^2$ und damit $|\bar{b}_2| = |b_2|$ sowie $p^2 = d_2 e^2$ mit $d_2 := \bar{b}_2/b_2$, $|d_2| = 1$; man beachte (20). Für die zweite Spalte folgt in analoger Weise

$$\begin{split} \bar{T}Pe^2 &= \bar{T}p^2 = d_2\bar{T}e^2 = d_2(\bar{b}_2e^1 + \bar{a}_2e^2 + \bar{b}_3e^3) \\ &= PTe^2 = b_2p^1 + a_2p^2 + b_3p^3 = b_2e^1 + d_2a_2e^2 + b_3p^3. \end{split}$$

Dies liefert $b_3 p^3 = (d_2 \bar{b}_2 - b_2) e^1 + d_2 (\bar{a}_2 - a_2) e^2 + d_2 \bar{b}_3 e^3$. Die erste Klammer verschwindet wegen der Festlegung von d_2 . Multiplikation mit $e^{2\uparrow}$ liefert $\bar{a}_2 = a_2$, mithin $p^3 = d_3 e^3$, $d_3 := d_2 \bar{b}_3 / b_3$, $|d_3| = 1$. Durch Induktion folgt schließlich P = D mit D aus (23), womit (24) in (22) übergeht.

13.5.9. Bemerkung. (i) Spalten- bzw. elementweise liest sich (22) als

$$\bar{q}^{j} = d_{j}q^{j}, \quad \bar{a}_{j} = a_{j} \ (j = 1, ..., n), \quad \bar{b}_{j} = d_{j-1}d_{j}b_{j} \ (j = 2, ..., n)$$
 (26)

und bedeutet: Die Faktoren Q und T der Faktorisierung $A = QTQ^{\intercal}$ sind unter der Voraussetzung (20) durch die erste Spalte q^1 von Q im wesentlichen festgelegt, nämlich bis auf das Vorzeichen der Spalten von Q und der Nebendiagonalelemente $\{b_i\}$ von T. Eine vergleichbare Aussage über die Eindeutigkeit der QR-Faktorisierung einer spaltenregulären Matrix haben wir in 10.1.3 kennengelernt. Hier sind allerdings mehr Freiheitsgrade vorhanden, so daß q^1 vorgegeben werden kann. Daß dies in der Tat bis auf gewisse Ausnahmesituationen möglich ist, zeigt der Lanczos-Algorithmus aus Ü 13.5.6.

(ii) Die Bedingung (21) kann durch $q^n = \bar{q}^n$ ersetzt werden, d. h., statt q^1 kann q^n vorgegeben werden. Der Beweis ist dann in umgekehrter Richtung — beginnend mit Spalte n von (24) — zu führen. \Box

Übungsaufgaben

Ü 13.5.1. Man schreibe ein Programm zur Berechnung von w := Bv, wenn $B \in S^{n,n}$ durch die Elemente des unteren Dreiecks $\{b_{ij}: i \ge j\}$ repräsentiert ist.

Ü 13.5.2. Man modifiziere 13.5.1 so, daß im ersten Schritt $A^{(2)} := H_1 A^{(1)} H_1$ Nullen in den Positionen $\{1, n\}, \{2, n\}, ..., \{n - 2, n\}$ und $\{n, 1\}, \{n, 2\}, ..., \{n, n - 2\}$ erzeugt werden usw. Welche Gestalt müssen dann H_1 und \overline{v}^1 aus $\overline{H}_1 = I - \overline{v}^1 \overline{v}^{1T} / \overline{\gamma}_1$ haben? Hin weis: Man beachte Ü 3.3.7.

Ü 13.5.3. Man überlege sich, daß die Transformation $A^{(k)} \rightarrow A^{(k+1)}$ auch in der Form

$$egin{aligned} A^{(k+1)} &:= \left\{ m{G}_{k+1,k+2} \cdots \left[\, m{G}_{k+1,n-1} \left(m{G}_{k+1,n} A^{(k)} m{G}_{k+1,n}^{\mathsf{T}}
ight) \, m{G}_{k+1,n-1}^{\mathsf{T}}
ight] \cdots \, m{G}_{k+1,k+2}^{\mathsf{T}}
ight\} \end{aligned}$$

realisiert werden kann.

Ü 13.5.4. Es sei

$$C = \begin{pmatrix} a_1 & c_1 \\ b_2 & a_2 & c_2 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ b_n & a_n \end{pmatrix} \quad \text{mit} \quad c_j * b_{j+1} > 0 \quad (j = 1, ..., n - 1)$$
(27)

eine nicht notwendig symmetrische Tridiagonalmatrix. Man konstruiere eine Diagonalmatrix $D = \text{diag}(d_i), d_i > 0$, so daß $T := D^{-1}CD$ symmetrisch ist. Eine Matrix des Typs (27) heißt quasisymmetrisch; bezüglich des Eigenwertproblems kann sie wie eine symmetrische Matrix behandelt werden.

Ú 13.5.5. Man zeige, daß die durch 13.5.1 erzeugte Matrix Q die erste Spalte $q^1 = e^1$ besitzt.

U 13.5.6. Man schreibe (1) als AQ = QT, mit $Q = (q^1, ..., q^n)$ also spaltenweise als $Aq^k = b_k q^{k-1} + a_k q^k + b_{k+1} q^{k+1}$ bzw.

$$b_{k+1}\boldsymbol{q}^{k+1} = \boldsymbol{A}\boldsymbol{q}^k - a_k\boldsymbol{q}^k - b_k\boldsymbol{q}^{k-1} =: \boldsymbol{p}^{k+1} \qquad (k = 1, ..., n)$$
 (28)

mit $b_1 := b_{n+1} := 0$, $q^0 := q^{n+1} := 0$.

(i) Man überlege sich: Es seien b_k , q^{k-1} und q^k bekannt. Linksmultiplikation von (28) mit $q^{k\uparrow}$ führt auf $a_k = q^{k\uparrow}Aq^k$. Im Fall $p^{k+1} \neq o$ ist dann die Richtung von q^{k+1} durch p^{k+1} fest-gelegt, andernfalls kann q^{k+1} orthogonal zu span $\{q^1, \ldots, q^k\}$ beliebig mit $||q^{k+1}|| = 1$ gewählt werden. Dies führt auf die folgende Grundform des sogenannten Lanczos-Algorithmus:

 $b_1 := 0$, $q^0 := o$, wähle q^1 mit $||q^1|| = 1$ for k := 1(1)n do

$$\begin{array}{l} a_k := {\bm{q}}^{k \intercal} A {\bm{q}}^k, \text{ if } k = n \text{ then stop} \\ {\bm{p}}^{k+1} := A {\bm{q}}^k - a_k {\bm{q}}^k - b_k {\bm{q}}^{k-1}, b_{k+1} := \pm || {\bm{p}}^{k+1} || \\ \text{ if } b_{k+1} \neq 0 \text{ then } {\bm{q}}^{k+1} := {\bm{p}}^{k+1} / b_{k+1} \\ & \quad \text{ else [wähle } {\bm{q}}^{k+1} \bot \text{ span } \{ {\bm{q}}^1, \dots, {\bm{q}}^k \} \text{ mit } || {\bm{q}}^{k+1} || = 1 \end{bmatrix} \end{array}$$

(ii) Man zeige durch Induktion über k, daß

$$\boldsymbol{q^{i\mathsf{T}}}\boldsymbol{q^{k+1}} = \begin{cases} 0 & \text{für } i = 1, ..., k, \\ 1 & \text{für } i = k+1 \end{cases} \quad \boldsymbol{q^{i\mathsf{T}}}\boldsymbol{A}\boldsymbol{q^{k+1}} = \begin{cases} 0 & \text{für } i = 1, ..., k-1, \\ b_{k+1} & \text{für } i = k \end{cases}$$

gilt, d. h., in exakter Arithmetik erzeugt der Lanczos-Algorithmus die Faktorisierung $A = QTQ^{\intercal}$. In der angegebenen Form ist das Verfahren jedoch instabil; die berechneten Vektoren $\{q^k\}$ können stark von der Orthogonalität abweichen.

13.6. Schneiden des Spektrums und Bisektionsverfahren

In diesem Abschnitt behandeln wir Algorithmen, die besonders zur Berechnung einiger weniger Eigenwerte der symmetrischen Tridiagonalmatrix

$$\boldsymbol{T} = \text{trid} (a_1, \ldots, a_n, b_2, \ldots, b_n)$$

geeignet sind, etwa zur Berechnung der k kleinsten oder größten Eigenwerte oder der in einem vorgegebenen Intervall liegenden Eigenwerte. Wir können dabei

$$b_i \neq 0$$
 $(i = 2, ..., n)$ (1)

voraussetzen; die Matrix T heißt dann *nichtzerfallend*. Dies ist keine Einschränkung: Ist etwa $b_{k+1} = 0$ für ein gewisses $k \in \{1, ..., n - 1\}$, so zerfällt T gemäß

in die beiden Blöcke T_1 und T_2 der Dimensionen k bzw. n - k. Jedem Eigenpaar $\{\lambda_1, u^1\}$ von T_1 bzw. $\{\lambda_2, u^2\}$ von T_2 entspricht das Eigenpaar $\{\lambda_1, \begin{pmatrix} u^1 \\ o \end{pmatrix}\}$ bzw. $\{\lambda_2, \begin{pmatrix} o \\ u^2 \end{pmatrix}\}$ von T, d. h., die Eigenpaare von T ergeben sich unmittelbar aus denen von T_1 und T_2 .

13.6.1. Bemerkung. Das durch verschwindende Nebendiagonalelemente bedingte Zerfallen des Eigenwertproblems von T in das mehrerer Probleme kleinerer Dimension mit nichtzerfallenden Blöcken T_i reduziert i. allg. den Gesamtaufwand. Es ist daher zweckmäßig, auch ein von 0 verschiedenes, aber ausreichend kleines b_{k+1} als 0 anzusehen. Als grobes Kriterium kann

$$|b_{k+1}| \leq \nu \|\boldsymbol{T}\| =: \varepsilon \tag{2}$$

genommen werden. Raffinierter ist das Kriterium

$$|b_{k+1}| \le v\{|a_k| + |a_{k+1}|\} =: \varepsilon.$$
(3)

In beiden Fällen entspricht das Nullsetzen von b_{k+1} einer Störung von A (bzw. T) mit $\|\delta A\| \leq \varepsilon$. \Box

Eine Folgerung von (1) wird durch die nachstehende Aussage angegeben.

13.6.2. Aussage. Die Eigenwerte $\lambda_j = \lambda_j[T]$ einer nichtzerfallenden Tridiagonalmatrix $T \in S^{n,n}$ sind sämtlich einfach, d. h., in der natürlichen Ordnung gilt

$$\lambda_1 < \lambda_2 < \dots < \lambda_{n-1} < \lambda_n. \tag{4}$$

Beweis. Für jedes λ sind die ersten n-1 Spalten von $T - \lambda I$ wegen (1) linear unabhängig, d. h., es gilt rang $(T - \lambda I) \ge n - 1$. Im Fall $\lambda = \lambda_j$ ist $T - \lambda_j I$ singulär, so daß dann rang $(T - \lambda_j I) = n - 1$ folgt. Der zu λ_j gehörende Eigenraum $\mathcal{N}(T - \lambda_j I)$ ist also eindimensional, d. h., es gibt nicht mehr als einen linear unabhängigen Eigenvektor zu λ_j . \Box

Aussage 13.6.2 schließt nicht aus, daß auch für deutlich von 0 verschiedene Nebendiagonalelemente gewisse Eigenwerte eng benachbart sein können. Ein viel zitiertes Beispiel dafür stellt die Wilkinson-Matrix

$$\boldsymbol{T} = \text{trid} \ (10, 9, \dots, 1, 0, 1, \dots, 9, 10; 1, \dots, 1) \in \boldsymbol{S}^{21, 21}$$
(5)

mit den Eigenwerten

$$\lambda_{20} = 10.74619\ 41829\ 0339\ \dots, \qquad \lambda_{21} = 10.74619\ 41829\ 0332\ \dots$$

dar, siehe B 13.8.

A. Die symmetrische Faktorisierung und das Schneiden des Spektrums

Die in diesem Abschnitt beschriebenen Algorithmen beruhen auf der Möglichkeit, für gegebenes $\mu \in \mathbf{R}$ die Anzahl $s(\mu) \in \{0, ..., n\}$ derjenigen Eigenwerte λ_j von \mathbf{T} , die kleiner als μ sind, einfach berechnen zu können. Bei der Numerierung (4) ist $s = s(\mu)$ durch

$$\lambda_1 < \lambda_2 < \dots < \lambda_s < \mu \leqq \lambda_{s+1} < \dots < \lambda_n \tag{6}$$

festgelegt. Durch μ wird also das Spektrum in die beiden Bereiche $\{\lambda_1, \ldots, \lambda_s\}$ und $\{\lambda_{s+1}, \ldots, \lambda_n\}$ "zerschnitten", die links bzw. rechts von μ auf der λ -Achse liegen. Offensichtlich ist $s(\mu)$ gleich der Anzahl der negativen Eigenwerte $\lambda_j[\mathbf{T}_{\mu}] = \lambda_j - \mu$

von

$$\boldsymbol{T}_{\boldsymbol{\mu}} := \boldsymbol{T} - \boldsymbol{\mu} \boldsymbol{I}. \tag{7}$$

Um $s(\mu)$ einfach berechnen zu können, nutzen wir diesen Fakt aus und erinnern uns an die folgende Aussage, siehe Ü 6.1.4: Wenn die symmetrische Dreiecksfaktorisierung

$$\boldsymbol{T}_{\mu} = \boldsymbol{L}_{\mu} \boldsymbol{D}_{\mu} \boldsymbol{L}_{\mu}^{\mathsf{T}} \tag{8}$$

mit

$$\boldsymbol{D}_{\mu} = \begin{pmatrix} d_{1} & & \\ & d_{2} & \\ & & d_{3} & \\ & & \ddots & \\ & & & d_{n} \end{pmatrix}, \quad \boldsymbol{L}_{\mu} = \begin{pmatrix} 1 & & & \\ & l_{2} & 1 & & \\ & & l_{3} & 1 & \\ & & \ddots & \ddots & \\ & & & \ddots & \ddots & \\ & & & l_{n} & 1 \end{pmatrix}$$
(9)

existiert, sind T_{μ} und D_{μ} kongruent, und die Anzahl der negativen Diagonalelemente von D_{μ} ist gleich der Anzahl der negativen Eigenwerte von T_{μ} , also gleich $s(\mu)$. Daß L_{μ} bei Existenz der Zerlegung bidiagonal ist, folgt sofort aus der Tridiagonalität von T_{μ} , vgl. 6.4.1.

Elementweise Auswertung von (8) führt auf die Identitäten

$$a_1 - \mu = d_1, \qquad a_k - \mu = d_k + l_k * d_{k-1} * l_k \qquad (k = 2, ..., n),$$
 (10)

$$b_k = l_k * d_{k-1}$$
 $(k = 2, ..., n).$ (11)

Wegen (1) und (11) ist

$$d_{k-1} \neq 0$$
 $(k = 2, ..., n),$ (12)

so daß nach $l_k = b_k/d_{k-1}$ aufgelöst werden kann. Wenn auch (10) nach d_k aufgelöst und rechts $l_k * d_{k-1}$ durch b_k ersetzt wird, ergibt sich der folgende Algorithmus zur Berechnung von $\{d_k, l_k\}$ aus $\{a_k, b_k, \mu\}$:

13.6.3. LDL^T-Faktorisierung von T_{μ} . Für $\mu \in \mathbb{R}$ und nichtzerfallendes T = trid $(a_1, \ldots, a_n, b_2, \ldots, b_n)$ existiere die Zerlegung (8), (9). Dann ist der Algorithmus

$$ie := 0, \quad d_1 := a_1 - \mu$$

for $k := 2(1)n$ do
$$| if d_{k-1} = 0 \text{ then } [ie := -1, \text{ stop}]$$

$$| l_k := b_k/d_{k-1}, \quad d_k := a_k - l_k * b_k - \mu$$

 $|\iota_k := b_k/d_{k-1}, \quad d_k := a_k - l_k * b_k - \mu$ in exakter Arithmetik mit ie = 0 durchführbar und liefert die signifikanten Elemente von D_{μ} und L_{μ} .

Aufwand: ~ 2n opms, ~ 2n S (für $\{d_k, l_k\}$).

Für die Berechnung von $s(\mu)$ — wir sagen dazu auch "Schneiden des Spektrums mit μ " (engl. "slicing the spectrum") — sind nur die Vorzeichen der d_k erforderlich. Es bietet sich daher an, den Ausdruck für l_k direkt in die Formel für d_k einzusetzen, was auf

$$d_k := a_k - (b_k)^2 / d_{k-1} - \mu$$

führt. Wenn alle d_k auf dem Platz von d gespeichert werden, erhalten wir die folgende Modifikation von 13.6.3:

13.6.4. Schneiden des Spektrums von T mit μ . Für $\mu \in \mathbb{R}$ und nichtzerfallendes $T = \text{trid}(a_1, \ldots, a_n, b_2, \ldots, b_n)$ existiere die Zerlegung (8), (9). Dann ist der Algorithmus

S1 (Berechnung der Hilfsgrößen $c_k = (b_k)^2$): for k := 2(1)n do $c_k := (b_k)^2$

S2 (Berechnung von $s = s(\mu)$ aus $\{a_k, c_k, \mu\}$): $ie := 0, d := a_1 - \mu$, if d < 0 then s := 1 else s := 0for k := 2(1)n do $\begin{vmatrix} \text{if } d = 0 \text{ then } [ie := -1, \text{ stop}] \\ d := a_k - c_k/d - \mu \\ \text{if } d < 0 \text{ then } s := s + 1 \end{vmatrix}$

in exakter Arithmetik mit ie = 0 durchführbar und liefert die Zahl $s = s(\mu)$ als Anzahl der Eigenwerte von T, die kleiner als μ sind.

Aufwand: $\sim n \text{ opm} + \sim n \text{ S}$ (für $\{c_k\}$), $\sim n(2 \text{ ops} + 1 \text{ opm})$ (für s)

13.6.5. Bemerkung. (i) Algorithmus 13.6.3 ist identisch mit der symmetrischen Version 6.4.4 (i) des Faktorisierungsalgorithmus 6.4.3 für eine Tridiagonalmatrix. Die Voraussetzung der Durchführbarkeit, d. h. der Existenz der Zerlegung (8), ist äquivalent zu (12) und praktisch nicht einschneidend: Offensichtlich ist $T_{\mu} = L_{\mu}R_{\mu}$, $R_{\mu} := D_{\mu}L_{\mu}^{T}$ die Gauß-Faktorisierung von T_{μ} . Nach Ü 5.2.5 existiert diese genau dann mit (12), wenn

$$\tau_k := \det \left((T_{\mu})_k \right) \neq 0 \qquad (k = 1, ..., n - 1)$$
 (13)

gilt, wobei

$$(T_{\mu})_k := T_k - \mu I_k, \qquad T_k := \text{trid} (a_1, ..., a_k, b_2, ..., b_k)$$
 (14)

ist. Nun ist τ_k der Wert des charakteristischen Polynoms der k-ten Hauptabschnittsmatrix $(\mathbf{T})_k$ von \mathbf{T} und kann an höchstens k verschiedenen Stellen verschwinden, so daß (13) für insgesamt höchstens $\sum_{k=1}^{n-1} k \sim n^2/2$ Werte von μ nicht erfüllt ist. Wenn zufälligerweise doch $d_{k-1} = 0$ ist, sollte mit einem leicht modifizierten Wert von μ neu gestartet werden, oder d_{k-1} sollte durch eine kleine, von 0 verschiedene Zahl ersetzt werden, siehe 13.6.7 (iv).

(ii) Zur Eigenwertberechnung muß 13.6.4 mehrfach für verschiedene μ -Werte ausgeführt werden. Es ist daher sinnvoll, den von μ unabhängigen Schritt S1 nur einmal zu Beginn der Rechnung auszuführen. Man beachte, daß in 13.6.4 nur die Quadrate der b_k vorkommen. Dies ist nicht verwunderlich, denn die Eigenwerte von T sind unabhängig vom Vorzeichen der b_k , siehe Ü 13.6.1. (iii) Nach Ü 5.2.5 gilt

$$d_k = \tau_k / \tau_{k-1}$$
 $(k = 1, ..., n), \quad \tau_0 := 1,$ (15)

so daß $s = s(\mu)$ auch als Anzahl der Vorzeichenwechsel in der Folge $\{\tau_0, \tau_1, ..., \tau_n\}$ berechnet werden kann. Die Zahlen $\{\tau_k\}$ lassen sich mittels einer dreigliedrigen Rekursionsformel einfach ermitteln, siehe Ü 13.6.2.

13.6.6. Rundungsfehleranalyse. Für $\mu \in \Re$ und nichtzerfallendes $T = \text{trid}(a_1, \ldots, a_n, b_2, \ldots, b_n)$ sei Algorithmus 13.6.3 bzw. 13.6.4 mit ie = 0 durchführbar. Dann gilt für die durch die berechneten Ausgangsdaten $\{d_k, l_k\}$ bzw. $\{d_k\}$ festgelegten Matrizen D_{μ} und L_{μ}

$$\boldsymbol{T}_{\mu} + \boldsymbol{\delta} \boldsymbol{T}_{\mu} = \boldsymbol{T} + \boldsymbol{\delta} \boldsymbol{T}_{\mu} - \mu \boldsymbol{I} = \boldsymbol{\hat{L}}_{\mu} \boldsymbol{D}_{\mu} \boldsymbol{\hat{L}}_{\mu}^{\mathsf{T}}$$
(16)

mit einer Dreiecksmatrix

$$m{L}_{\mu} = egin{pmatrix} 1 & & & \ \hat{l}_2 & 1 & & \ & \hat{l}_3 & 1 & & \ & \ddots & \ddots & \ & & \hat{l}_n & 1 \end{pmatrix} \in \mathbf{R}^{n.n}$$

und einer tridiagonalen Störung $\mathbf{\delta T}_{\mu} \in \mathbf{S}^{n,n}$, die den Abschätzungen

$$|\boldsymbol{\delta}\boldsymbol{T}_{\boldsymbol{\mu}}| \leq \boldsymbol{\nu}[2 |\boldsymbol{T}| + |\boldsymbol{\mu}| \boldsymbol{I}]$$
(17)

sowie

$$|\hat{l}_k - l_k| \leq 2\nu |l_k| \qquad (k = 2, ..., n)$$
 (18)

genügen.

13.6.7. Bemerkung. (i) Praktisch werden die angegebenen Algorithmen nur mit μ -Werten, für die

$$|\mu| \leq ||\mathbf{T}||_{\infty} = \max \{|a_k| + |b_k| + |b_{k+1}| : k = 1, ..., n\},$$

$$b_1 := b_{n+1} := 0$$
(19)

gilt, durchgeführt. Aus (17) folgt dann

$$\|\boldsymbol{\delta}\boldsymbol{T}_{\boldsymbol{\mu}}\|_{\infty} \leq \nu \{2 \|\boldsymbol{T}\|_{\infty} + |\boldsymbol{\mu}|\} \leq 3\nu \|\boldsymbol{T}\|_{\infty}, \tag{20}$$

d. h., die berechnete Diagonalmatrix D_{μ} ist exakt kongruent zu $T_{\mu} + \delta T_{\mu}$, und δT_{μ} ist klein gegenüber T. Die berechnete Zahl $s = s(\mu)$ ist also die Anzahl der Eigenwerte von $T + \delta T_{\mu}$, die kleiner als μ sind. Die Berechnung von $s = s(\mu)$ nach 13.6.4 ist daher ein numerisch gutartiger Proze β .

(ii) Wenn $|\mu|$ in der Größenordnung von min $\{|a_k|\}$ liegt, ist δT_{μ} eine elementweise kleine Störung von T. Dies zieht nach sich, daß unter Verwendung von $s(\mu)$ oft auch die betragskleinsten Eigenwerte einer gestuften Tridiagonalmatrix T ausreichend genau lokalisiert werden können. (iii) Aus 13.6.6 folgt

$$\boldsymbol{T}_{\mu} + \boldsymbol{\delta} \hat{\boldsymbol{T}}_{\mu} = \boldsymbol{L}_{\mu} \boldsymbol{D}_{\mu} \boldsymbol{L}_{\mu}^{\mathsf{T}} \quad \text{mit} \quad \boldsymbol{\delta} \hat{\boldsymbol{T}}_{\mu} := \boldsymbol{\delta} \boldsymbol{T}_{\mu} + \boldsymbol{\delta} \tilde{\boldsymbol{T}}_{\mu}, \qquad (21)$$

wobei

$$|\delta \tilde{T}_{\mu}| \leq 2\nu \begin{pmatrix} 0 & |b_{2}| & & \\ |b_{2}| & 2 & |l_{2}b_{2}| & |b_{3}| & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & \ddots & \\ & & & |b_{n}| & 2 & |l_{n}b_{n}| \end{pmatrix}$$
(22)

ist, siehe Ü 13.6.3. Da die Produkte $|l_k b_k|$ groß werden können, ist $\delta \tilde{T}_{\mu}$ nicht notwendig klein gegenüber T_{μ} , d. h., die berechneten Faktoren D_{μ} , L_{μ} sind i. allg. nicht die exakten Faktoren einer leicht gestörten Matrix. Der Faktorisierungsalgorithmus 13.6.3 ist also i. allg. nicht numerisch gutartig.

(iv) Bei der Berechnung von $d_k := \operatorname{fl} (a_k - l_k * b_k - \mu)$ gilt

$$d_{k} = \{ [a_{k} - l_{k} * b_{k}(1 + \varepsilon_{k})] (1 + \delta_{k}) - \mu \} (1 + \xi_{k})$$
(23)

mit $|\varepsilon_k|, |\delta_k|, |\xi_k| \leq \nu$, also

$$d_k - [a_k - l_k * b_k - \mu]| \leq \nu[2 |a_k| + 3 |l_k * b_k| + |\mu|] =: \Delta d_k.$$

Eine Störung des berechneten d_k in der Größenordnung von Δd_k ändert also das Fehlerniveau von d_k nicht wesentlich. Dies kann ausgenutzt werden, um in 13.6.3 vorzeitigen erfolglosen Abbruch wegen $d_{k-1} = 0$ zu vermeiden, etwa durch Verwendung der folgenden Modifikation:

$$\begin{array}{l} ie := 0, \, d_1 := a_1 - \mu + \nu(|a_1| + |\mu|) \\ \text{if } d_1 = 0 \text{ then } [ie := -1, \, \text{stop}] \\ \text{for } k := 2(1)n \text{ do} \\ & \left| \begin{array}{l} l_k := b_k/d_{k-1}, \, p := l_k \ast b_k, \, d_k := a_k - p - \mu \\ & \text{if } d_k = 0 \text{ then } d_k := \nu(|a_k| + |p| + |\mu|) \end{array} \right. \end{array}$$

Bis auf den Ausnahmefall $a_1 = \mu = 0$ ist damit unerwünschter Abbruch ausgeschlossen, und die Abschätzungen in 13.6.6 werden nur unwesentlich schlechter. Analog kann 13.6.4 modifiziert werden; die Korrekturgröße ist dann $d_k := \nu\{|a_k| + |c_k/d_{k-1}| + |\mu|\}$. Die durch die Modifikation hervorgerufene Überlaufgefahr ist gering, wenn $\|T\|_{\infty}$ nicht zu groß ist. \Box

B. Das Bisektionsverfahren

Wir zeigen im folgenden, wie durch das Schneiden des Spektrums gemäß 13.6.4 jeder Eigenwert λ_j von T mit vorgegebenem Index j bezüglich der Numerierung (4) in exakter Arithmetik beliebig genau lokalisiert werden kann. Es sei dazu $[\alpha, \beta]$ ein Intervall mit

$$s(\alpha) < j \leq s(\beta).$$
 (24)

Nach Definition von $s(\mu)$ ist dies zu

$$\alpha \leq \lambda_j < \beta \tag{25}$$

äquivalent, vgl. (6). Um die Lokalisierungsaussage (25) zu verbessern, berechnen wir den Mittelpunkt $\mu := (\alpha + \beta)/2$ des Intervalls $[\alpha, \beta]$ und bestimmen $s = s(\mu)$ nach 13.6.4. Dann können zwei Fälle auftreten:

Fall 1: $s(\mu) < j$. Dann sind weniger als j Eigenwerte kleiner als μ , d. h., es gilt

$$\mu \leq \lambda_j < \beta.$$

Fall 2: $s(\mu) \ge j$. Dann sind mindestens j Eigenwerte kleiner als μ , d. h., es gilt

$$\alpha \leq \lambda_j < \mu.$$

Im ersten Fall wird $[\alpha, \beta]$ durch $[\mu, \beta]$ ersetzt, im zweiten durch $[\alpha, \mu]$, wobei die Intervallänge halbiert wird. Fortsetzung dieses Vorgehens führt auf das im folgenden angegebene Bisektionsverfahren:

13.6.8. Das Bisektionsverfahren zur Lokalisierung von λ_i . Gegeben seien die nichtzerfallende Tridiagonalmatrix $T = \text{trid}(a_1, \ldots, a_n, b_2, \ldots, b_n)$, ein Index $j \in \{1, \ldots, n\}$, Zahlen $\{\alpha, \beta\}$ mit (25) und eine Abbruchschranke $\varepsilon > 0$. Dann werden die Eingangsdaten $\{\alpha, \beta\}$ durch den in exakter Arithmetik ausgeführten Algorithmus

S1 (Bisektionsschritt): $\delta := \beta - \alpha$

if $\delta > \varepsilon$ then

 $\begin{array}{l} \mu := \alpha + \delta/2, \quad \text{berechne } s = s(\mu) \text{ nach } 13.6.4 \\ \text{if } s < j \text{ then } \alpha := \mu \text{ else } \beta := \mu \\ \text{goto } S1 \\ \text{mit Zahlen } \{\alpha, \beta\} \text{ überspeichert, für die } (25) \text{ mit } \beta - \alpha \leq \varepsilon \text{ gilt.} \\ Aufwand : \sim kn(2 \text{ ops } + 1 \text{ opm}), \text{ wobei } k - \text{ die Anzahl der Bisektionsschritte } - \\ \text{is blaister } \pi = \pi \delta =$

die kleinste ganze Zahl größer gleich $\log_2 ((\beta - \alpha)/\varepsilon)$ ist.

13.6.9. Bemerkung. (i) Wenn $[\alpha_1, \beta_1]$ das Startintervall bezeichnet, gilt nach k Bisektionsschritten $\delta_{k+1} = \beta_{k+1} - \alpha_{k+1} = (\beta_1 - \alpha_1)/2^k$, d. h., die Abbruchbedingung $\delta_{k+1} \leq \varepsilon$ ist zu

$$k \ge \log_2\left((\beta - \alpha)/\varepsilon\right) \tag{26}$$

äquivalent. Dies führt auf die oben angegebene Schrittzahl.

(ii) Als Starteinschließung (25) kann die aus dem Satz von GERŠGORIN folgende Einschließung mit

$$\begin{aligned} \alpha &:= G_{\min} := \min \left\{ a_k - |b_k| - |b_{k+1}| : k = 1, ..., n \right\}, \\ \beta &:= G_{\max} := \max \left\{ a_k + |b_k| + |b_{k+1}| : k = 1, ..., n \right\}, \end{aligned}$$
(27)

 $b_1 := b_{n+1} := 0$, verwendet werden, die für jedes j gültig ist, siehe 1.2.10. Im Fall j = n könnte $\lambda_n = G_{\max}$, also $s(G_{\max}) = n - 1$ sein. Dann ist theoretisch G_{\max} durch eine wenig größere Zahl zu ersetzen. Praktisch tritt dieser Fall nicht auf. Für jedes $\mu \in [\alpha, \beta]$ gilt bei der Festlegung (27)

$$|\mu| \leq \max\{|\alpha|, |\beta|\} = G_{\max} = \|\boldsymbol{T}\|_{\infty}, \qquad (28)$$

d. h., (19) ist in allen Schritten erfüllt.

(iii) Mit Hilfe des Bisektionsverfahrens 13.6.8 können in offensichtlicher Weise z. B. die p kleinsten Eigenwerte $\{\lambda_1, \ldots, \lambda_p\}$, die p größten Eigenwerte $\{\lambda_{n-p+1}, \ldots, \lambda_n\}$, die im Intervall $[\alpha, \beta)$ liegenden Eigenwerte $\{\lambda_i, \ldots, \lambda_j\}$ o. ä. berechnet werden. Im letzten Fall ist $i = s(\alpha) + 1$, $j = s(\beta)$. Wenn wie in den angegebenen Fällen mehrere Eigenwerte $\{\lambda_{j_1}, \ldots, \lambda_{j_p}\}$ zu lokalisieren sind und sich die aktuellen Einschließungsintervalle $[\alpha^{(l)}, \beta^{(l)}]$ für λ_{j_1} überlappen, können mit éinem berechneten $s = s(\mu)$ alle Intervalle aufdatiert werden, die μ als inneren Punkt enthalten. Dadurch kann in der Anfangsphase eine wesentliche Aufwandssenkung erzielt werden. In der Schlußphase werden die Einschließungsintervalle in der Regel disjunkt sein, so daß jedes für sich behandelt werden muß. \Box

13.6.10. Rundungsfehleranalyse. Es sei $[\alpha, \beta]$ das durch 13.6.8 berechnete Einschließungsintervall. Dann gibt es eine tridiagonale Störung $\delta T \in S^{n,n}$, so daß

$$\alpha \leq \lambda_j [T + \delta T] < \beta \quad \text{mit} \quad \beta - \alpha \leq \varepsilon$$
(29)

gilt, wobei

$$|\boldsymbol{\delta}\boldsymbol{T}| \leq \nu[2 |\boldsymbol{T}| + \max\{|\boldsymbol{\alpha}|, |\boldsymbol{\beta}|\} \boldsymbol{I}]$$
(30)

ist.

Beweis. Nach 13.6.6 existieren für das Schlußintervall $[\alpha, \beta]$ Störungen $\delta T_{\alpha}, \delta T_{\beta}$ mit

$$\alpha \leq \lambda_{j} [\boldsymbol{T} + \boldsymbol{\delta} \boldsymbol{T}_{a}], \quad |\boldsymbol{\delta} \boldsymbol{T}_{a}| \leq v\{2 |\boldsymbol{T}| + |\boldsymbol{\alpha}| \boldsymbol{I}\},$$

$$\lambda_{j} [\boldsymbol{T} + \boldsymbol{\delta} \boldsymbol{T}_{\beta}] < \beta, \quad |\boldsymbol{\delta} \boldsymbol{T}_{\beta}| \leq v\{2 |\boldsymbol{T}| + |\boldsymbol{\beta}| \boldsymbol{I}\}.$$

$$(31)$$

Im Fall $\alpha \leq \lambda_i[T] < \beta$ gelten (29) und (30) mit $\delta T = 0$. Im Fall $\lambda_i[T] < \alpha$ ist

 $\lambda_{j}[\boldsymbol{T}+0\cdot\boldsymbol{\delta}\boldsymbol{T}_{\alpha}]<\alpha\leq\lambda_{j}[\boldsymbol{T}+1\cdot\boldsymbol{\delta}\boldsymbol{T}_{\alpha}].$

Wegen der stetigen Abhängigkeit der Eigenwerte von der Matrix gibt es daher ein ξ , $0 < \xi \leq 1$, so daß (29) für $\delta T := \xi \cdot \delta T_{\alpha}$ erfüllt ist, und wegen (31) gilt auch (30). Im Fall $\beta \leq \lambda_i[T]$ kann analog geschlossen werden. \Box

13.6.11. Bemerkung. (i) Für die Anfangseinschließung (27) folgt aus (30) wegen (28)

$$\|\boldsymbol{\delta}\boldsymbol{T}\|_{\infty} \leq 3\boldsymbol{v} \|\boldsymbol{T}\|_{\infty},\tag{32}$$

d. h., im berechneten Intervall $[\alpha, \beta]$ liegt der *j*-te Eigenwert der zu T im Sinne von (32) benachbarten Matrix $T + \delta T$. Dies bedeutet die *numerische Gutartigkeit* des Bisektionsverfahrens.

(ii) Aus (29) folgt

$$|\lambda_j[T + \delta T] - \mu| \leq (\beta - \alpha)/2 \leq \varepsilon/2 \quad \text{mit} \quad \mu := (\alpha + \beta)/2.$$

Andererseits gilt nach 13.1.2 und (32)

 $|\lambda_j[T + \delta T] - \lambda_j[T]| \leq \|\delta T\| \leq \|\delta T\| \leq \|\delta T\|_{\infty} \leq 3\nu \|T\|_{\infty},$

für das vorletzte Ungleichheitszeichen siehe Ü 1.2.9, so daß sich insgesamt

$$|\lambda_{j}[\boldsymbol{T}] - \mu| \leq \varepsilon/2 + 3\nu \|\boldsymbol{T}\|_{\infty}$$
(33)

als Abschätzung des Fehlers von μ gegen $\lambda_i[T]$ ergibt.

(iii) Der Wahl $\varepsilon := \nu ||\mathbf{T}||_{\infty}$ entspricht i. allg. die größtmögliche Genauigkeit des Bisektionsverfahrens; aus (33) folgt dann

$$|\lambda_{j}[T] - \mu| \leq 3.5\nu \|T\|_{\infty} \leq 3.5\nu \sqrt{3} \|T\| \approx 6.1\nu \|T\|,$$
(34)

man beachte Ü 13.6.4. Andererseits ergibt sich aus Ü 13.6.4 und 13.1.2 das optimale Fehlerniveau von $\lambda_j[T]$ zu

$$|\lambda_{j}[\boldsymbol{T} + \boldsymbol{\delta}\boldsymbol{T}_{D}] - \lambda_{j}[\boldsymbol{T}]| \leq \|\boldsymbol{\delta}\boldsymbol{T}_{D}\| \leq \sqrt{3} \, \boldsymbol{\nu} \, \|\boldsymbol{T}\| \approx 1.7 \, \|\boldsymbol{T}\|, \tag{35}$$

siehe auch (13.1.22) für den Fall einer allgemeinen Matrix, wo \sqrt{n} statt $\sqrt{3}$ steht. Der Vergleich von (34) und (35) zeigt, daß das Bisektionsverfahren mit der obigen ε -Wahl hinsichtlich der Genauigkeit fast optimal ist. Wenn für die Starteinschließung $\beta - \alpha \approx ||\mathbf{T}||_{\infty}$ gilt, sind jedoch der Größenordnung nach $\log_2 (1/\nu)$ Schritte erforderlich. Selbst mit den in 13.6.9(iii) beschriebenen Methoden zur Effektivitätserhöhung werden dann zur Bestimmung aller Eigenwerte der Größenordnung nach $\sim \log_2 (1/\nu)$ $\times n^2$ opms benötigt. Dies ist deutlich mehr als bei Verwendung der Varianten des **QR**-Algorithmus aus 13.7. In der Literatur wird empfohlen, das Bisektionsverfahren nur zu verwenden, wenn weniger als n/4 Eigenwerte gesucht sind. Andernfalls ist es billiger, alle Eigenwerte mittels des **QR**-Algorithmus zu berechnen.

(iv) Falls die betragskleinen Eigenwerte einer gestuften Tridiagonalmatrix T zu bestimmen sind, sollte $\varepsilon \approx r ||T_0||_{\infty}$ gewählt werden, wobei T_0 den betragskleinen Teil von T bezeichnet.

(v) Da das Bisektionsverfahren i. allg. eine sehr genaue Näherung μ für den gesuchten Eigenwert λ von T liefert, kann eine Näherung v für den zugehörigen Eigenvektor u durch wenige Schritte der inversen Iteration 13.4.3 berechnet werden. Man beachte jedoch, daß die für das Schneiden des Spektrums geeignete LDL^{T} -Faktorisierung aus 13.6.3 für die Gleichungsauflösung bei der inversen Iteration i. allg. ungeeignet ist, denn sie ist nicht numerisch gutartig, siehe 13.6.7 (iii) und Text nach 13.4.5. Es sollte die in 13.4.3 vorgesehene pivotisierte Gauß-Faktorisierung verwendet werden, die nur $\sim n(1 \text{ opm} + 2 \text{ opms})$ kostet, siehe 6.4.1. Da T in der Regel nach dem Tridiagonalisierungsverfahren 13.5.1 gemäß $T = Q^{T}AQ$ aus einer voll besetzten Matrix A entstanden ist, muß der Eigenvektor v von T noch in den Eigenvektor Qv von A rücktransformiert werden. Mit der Produktform von Qerfordert dies $\sim n^{2}$ opms pro Eigenvektor, siehe 13.5.2 (ii).

(vi) Wenn einige Eigenwerte einer Bandmatrix A gesucht sind, braucht A nicht unbedingt erst nach 13.5.6 tridiagonalisiert zu werden, sondern das Schneiden des Spektrums kann direkt mit $A - \mu I$ unter Verwendung der etwa nach 6.1.1 berechneten LDL^{\intercal} -Faktorisierung $A - \mu I = L_{\mu}D_{\mu}L_{\mu}^{\intercal}$ durchgeführt werden. Der Test " $a_{kk} \leq 0$ " ist dabei durch " $a_{kk} = 0$ " zu ersetzen, da $A - \mu I$ i. allg. nicht positiv definit ist. Im Gegensatz zum tridiagonalen Fall kann jedoch die berechnete Diagonalmatrix D_{μ} zu einer Matrix $A + \delta A$ kongruent sein, die sich von A stark unterscheidet. Damit ist insbesondere im Fall großer $|l_{ij}|$ zu rechnen. Wenn $a_{kk} = 0$ ist oder große $|l_{ij}|$ auftreten, sollte die Faktorisierung deshalb mit einem leicht veränderten Wert von μ wiederholt werden. Meist gelingt es auf diese Weise, den gesuchten Eigenwert so genau zu approximieren, daß einige nachfolgende Schritte der inversen oder der RQ-Iteration ein akzeptables Eigenpaar liefern. Für jede der genannten Iterationen sollte jedoch eine numerisch gutartige Faktorisierung etwa die spaltenpivotisierte Gauß-Faktorisierung — benutzt werden, vgl. 13.4.

Übungsaufgaben

Ü 13.6.1. Man konstruiere eine Matrix $D = \text{diag}(d_i)$ mit $|d_i| = 1$, so daß

trid
$$(a_1, ..., a_n, |b_2|, ..., |b_n|) = \boldsymbol{D}[\text{trid } (a_1, ..., a_n, b_2, ..., b_n)] \boldsymbol{D}$$

gilt. Wegen $D = D^{\intercal} = D^{-1}$ sind die beiden Tridiagonalmatrizen also ähnlich und besitzen dieselben Eigenwerte. Letztere hängen daher nur von den Beträgen der Nebendiagonalelemente ab.

Ü 13.6.2. Man zeige durch Entwicklung von det $((T_{\mu})_k)$ nach der letzten Spalte, daß die Hauptabschnittsdeterminanten τ_k von T_{μ} der Rekursion

$$\tau_k = (a_k - \mu) \tau_{k-1} - (b_k)^2 \tau_{k-2} \qquad (k = 2, ..., n)$$
(36)

mit $\tau_0 = 1$ genügen. Die Folge $\{\tau_0, ..., \tau_n\}$ und damit $s = s(\mu)$ als Zahl der Vorzeichenwechsel in dieser Folge läßt sich also mit $\sim n(2 \text{ ops} + 2 \text{ opm})$ bestimmen, falls $c_k = b_k^2$ bekannt ist, allerdings kann für großes n Über- bzw. Unterlauf eintreten. Wenn zu dessen Vermeidung zu den Quotienten $d_k := \tau_k/\tau_{k-1}$ übergegangen wird, geht (36) in die Rekursion aus 13.6.4 über.

Ü 13.6.3. Man beweise die Gültigkeit von (21) mit $\delta \tilde{T}_{\mu} := -[\delta L_{\mu} D_{\mu} L_{\mu}^{\mathsf{T}} + L_{\mu} D_{\mu} \delta L_{\mu}^{\mathsf{T}} + \delta L_{\mu} D_{\mu} \delta L_{\mu}^{\mathsf{T}}], \delta L_{\mu} := \hat{L}_{\mu} - L_{\mu}$ und schätze $|\delta \tilde{T}_{\mu}|$ gemäß (22) ab.

Ü 13.6.4. Man zeige, daß für jede Tridiagonalmatrix $T = trid (a_i, b_i)$

$$\|\boldsymbol{T}\|_{\infty} \leq \sqrt{3} \|\boldsymbol{T}\| \tag{37}$$

gilt, und folgere hieraus die Abschätzung

$$\|\boldsymbol{\partial}\boldsymbol{T}_{D}\| \leq \boldsymbol{\nu}\,\sqrt{3}\,\|\boldsymbol{T}\| \tag{38}$$

für den Darstellungsfehler $\delta T_D = rd(T) - T von T$.

Hinweis: Wenn k die Zeile mit maximaler Betragssumme von T bezeichnet, gilt

 $\|T\|_{\infty} \leq \sqrt{3} \ \sqrt{|b_k|^2 + |a_k|^2 + |b_{k+1}|^2} = \sqrt{3} \ \sqrt{(T^2)_{kk}}.$

Zum Nachweis von (38) beachte man 2.3.15 und Ü 1.2.9.

13.7. Der **QR**-Algorithmus

Der vor etwa 25 Jahren entdeckte QR-Algorithmus — zur Geschichte siehe B 13.9 — ist der effektivste der derzeit bekannten Algorithmen zur Berechnung aller Eigenwerte und gegebenenfalls Eigenvektoren einer symmetrischen Tridiagonalmatrix. In Kombination mit den im Abschnitt 13.5 diskutierten Tridiagonalisierungsverfahren stellt er das Verfahren zur Lösung des vollständigen Eigenwertproblems voll besetzter symmetrischer Matrizen mittlerer Dimension dar.

Da ein großer Teil der Theorie auch für beliebige symmetrische Matrizen gilt, setzen wir zunächst außer der Symmetrie keine weiteren Eigenschaften der Matrix voraus.

A. Motivierung, Grundform und Zusammenhang mit Teilraumiteration

Ausgangspunkt der folgenden Überlegungen ist der Algorithmus 13.3.10 zur Teilraumiteration, der mit p = m = n zur Berechnung aller Eigenvektoren von A verwendet werden kann. Um später mit der beim QR-Algorithmus üblichen Bezeichnung nicht zu kollidieren, schreiben wir 13.3.10 in der Form

Wähle
$$\overline{V}_1$$
, bilde $\overline{W}_{k+1} := A \overline{V}_k$, faktorisiere \overline{W}_{k+1} gemäß $\overline{W}_{k+1} = \overline{V}_{k+1} \overline{R}_k$ (1)

(k = 1, 2, ...); in 13.3.10 steht \overline{R}_{k+1} statt \overline{R}_k . Dabei sind $\{\overline{V}_k\}$ und $\{\overline{R}_k\}$ Folgen orthogonaler bzw. oberer Dreiecksmatrizen.

Wenn die aus A durch orthogonale Ähnlichkeitstransformation mit \overline{V}_k gemäß

$$\bar{A}_{k} = \bar{V}_{k}^{\mathsf{T}} A \bar{V}_{k} \tag{2}$$

entstehenden Matrizen \bar{A}_k eingeführt werden, ergibt sich wegen (1)

$$\bar{\boldsymbol{A}}_{k} = \bar{\boldsymbol{V}}_{k}^{\mathsf{T}} \overline{\boldsymbol{W}}_{k+1} = \bar{\boldsymbol{V}}_{k}^{\mathsf{T}} \overline{\boldsymbol{V}}_{k+1} \overline{\boldsymbol{R}}_{k} = \bar{\boldsymbol{Q}}_{k} \overline{\boldsymbol{R}}_{k}$$
(3)

mit

$$\bar{\boldsymbol{Q}}_{k} := \bar{\boldsymbol{V}}_{k}^{\mathsf{T}} \bar{\boldsymbol{V}}_{k+1}, \quad \text{also} \quad \bar{\boldsymbol{V}}_{k+1} = \bar{\boldsymbol{V}}_{k} \bar{\boldsymbol{Q}}_{k} = \bar{\boldsymbol{V}}_{1} \bar{\boldsymbol{Q}}_{1} \bar{\boldsymbol{Q}}_{2} \cdots \bar{\boldsymbol{Q}}_{k}, \tag{4}$$

d. h., \overline{R}_k ist ein Dreiecksfaktor der **QR**-Faktorisierung (3) von \overline{A}_k , und \overline{V}_{k+1} entsteht gemäß (4) aus \overline{V}_k durch Aufdatierung mit dem zugehörigen Orthogonalfaktor \bar{Q}_k . Mit (4), (3) folgt weiter

$$\bar{A}_{k+1} = \vec{V}_{k+1}^{\mathsf{T}} A \vec{V}_{k+1} = \bar{Q}_{k}^{\mathsf{T}} \bar{A}_{k} \bar{Q}_{k} = \bar{Q}_{k}^{\mathsf{T}} \bar{Q}_{k} \bar{R}_{k} \bar{Q}_{k} = \bar{R}_{k} \bar{Q}_{k}, \qquad (5)$$

d. h., \bar{A}_{k+1} ist das Produkt der Faktoren \bar{Q}_k , \bar{R}_k der Faktorisierung (3) in umgekehrter Reihenfolge. Da die **QR**-Faktorisierung für eine reguläre Matrix nach 10.1.3 bis auf das Vorzeichen der Spalten von O und der zugehörigen Zeilen von R eindeutig ist, könnte man versuchen, die Formeln (3) bis (5) zur Erzeugung der Folgen $\{A_k\}, \{Q_k\}, \{R_k\}, \{V_k\}$ zu verwenden. Um die durch möglicherweise andere Vorzeichenfestlegung bei der **QR**-Faktorisierung hervorgerufenen Unterschiede auszudrücken, bezeichnen wir die derart definierten Matrizen ohne Querstrich. Bei der sich anbietenden Festlegung $V_1 := I$, also $A_1 = A$, führt dies auf die folgende Vorschrift, die die Grundform des **QR**-Algorithmus darstellt:

13.7.1. QR-Algorithmus ohne Verschiebungen.

- S0 (Initialisierung): Setze $A_1 := A$, $V_1 := I$, k := 1S1 (Faktorisierung von A_k): Berechne QR-Faktorisierung $A_k = Q_k R_k$ S2 (Bildung von A_{k+1} , V_{k+1}): Bestimme $A_{k+1} := R_k Q_k$, $V_{k+1} := V_k Q_k$
- S3: Setze k := k + 1, goto S1

Es gilt dann wie in (5)

$$A_{k+1} = V_{k+1}^{\mathsf{T}} A V_{k+1} = Q_k^{\mathsf{T}} A_k Q_k, \tag{6}$$

und die Transformationsmatrix

$$V_{k+1} = V_k Q_k = Q_1 Q_2 \cdots Q_k \tag{7}$$

als Produkt der Q-Faktoren ist orthogonal.

13.7.2. Aussage. Es sei $A \in S^{n,n}$ regulär, d. h., $\lambda = 0$ sei kein Eigenwert von A. Die Folgen $\{\overline{V}_k\}$, $\{\overline{R}_k\}$ seien nach der Teilraumiteration 13.3.10 mit $\overline{V}_1 := I$ erzeugt, und $\{\overline{A}_k\}$, $\{\overline{Q}_k\}$ bezeichne die durch (2) bzw. (4) zugeordneten Matrizen. Die Folgen $\{A_k\}$, $\{Q_k\}$, $\{R_k\}$, $\{V_k\}$ seien durch den QR-Algorithmus 13.7.1 definiert. Dann gilt

$$\overline{V}_k = V_k D_{k-1} \tag{8}$$

und

$$\bar{\boldsymbol{A}}_{k} = \boldsymbol{D}_{k-1} \boldsymbol{A}_{k} \boldsymbol{D}_{k-1}, \qquad \boldsymbol{\bar{Q}}_{k} = \boldsymbol{D}_{k-1} \boldsymbol{Q}_{k} \boldsymbol{D}_{k}, \qquad \boldsymbol{\bar{R}}_{k} = \boldsymbol{D}_{k} \boldsymbol{R}_{k} \boldsymbol{D}_{k-1}$$
(9)

für k = 1, 2, ... mit

$$\boldsymbol{D}_0 := \boldsymbol{I} \quad \text{und} \quad |\boldsymbol{D}_k| = \boldsymbol{I}, \quad \text{d. h.} \quad \boldsymbol{D}_k = \text{diag}(\pm 1).$$
 (10)

Beweis. Angenommen, es gelten (8) bis (10) bis zum Index k. Dann folgt

$$\overline{V}_{k+1} = \overline{V}_k \overline{Q}_k = V_k D_{k-1} D_{k-1} Q_k D_k = V_k Q_k D_k = V_{k+1} D_k$$

also (8) für den Index k + 1. Weiter ist

$$\bar{A}_{k+1} = \bar{Q}_k^{\mathsf{T}} \bar{A}_k \bar{Q}_k = \left(\boldsymbol{D}_k Q_k^{\mathsf{T}} \boldsymbol{D}_{k-1} \right) \left(\boldsymbol{D}_{k-1} A_k \boldsymbol{D}_{k-1} \right) \left(\boldsymbol{D}_{k-1} Q_k \boldsymbol{D}_k \right) = \boldsymbol{D}_k A_{k+1} \boldsymbol{D}_k$$

Mit dieser Beziehung ergibt sich

$$\bar{Q}_{k+1}\bar{R}_{k+1} = \bar{A}_{k+1} = D_k A_{k+1} D_k = (D_k Q_{k+1}) (R_{k+1} D_k)$$

Da \bar{A}_{k+1} zu A ähnlich und folglich regulär ist, gibt es nach 10.1.3 ein D_{k+1} mit $|D_{k+1}| = I$, so daß $\bar{Q}_{k+1} = (D_k Q_{k+1}) D_{k+1}$, $\bar{R}_{k+1} = D_{k+1}(R_{k+1}D_k)$ gilt, denn $D_k Q_{k+1}$ und $R_{k+1}D_k$ sind wie \bar{Q}_{k+1} und \bar{R}_{k+1} Faktoren von \bar{A}_{k+1} . Die Gültigkeit für k = 1 wird analog gezeigt bzw. folgt aus $\bar{V}_1 = V_1 = I$. \Box

13.7.3. Bemerkung. Im Beweis wurde nur ausgenutzt, daß die Folgen $\{\overline{V}_k\}, \{\overline{R}_k\}, \{\overline{A}_k\}, \{\overline{Q}_k\}$ den zu 13.7.1 analogen Gleichungen genügen. Es wurde also eigentlich gezeigt, daß die durch 13.7.1 definierten Matrizenfolgen im wesentlichen — d. h. im Sinne von (8) bis (10) — durch A eindeutig festgelegt sind, sofern A regulär ist. Wird zusätzlich z. B. $(R_k)_{jj} > 0$ (j = 1, ..., n) gefordert, was sich durch gleichzeitiges Vertauschen der Vorzeichen der *j*-ten Spalte von Q_k und der *j*-ten Zeile von R_k stets erreichen ließe, ist $D_k = I$, und die Folgen sind sogar im strengen Sinne eindeutig festgelegt. \Box

Da die Spalten von V_k und \overline{V}_k bis auf das Vorzeichen identisch sind, übertragen sich die Konvergenzaussagen für die Teilraumiteration auf den **QR**-Algorithmus.

13.7.4. Satz. Für die Matrix $A \in S^{n,n}$ mit den Eigenwerten $\{\lambda_j\}$,

$$0 < |\lambda_n| \le |\lambda_{n-1}| \le \dots \le |\lambda_2| \le |\lambda_1|, \tag{11}$$

und zugehörigen orthonormierten Eigenvektoren $\{u^{j}\}$ werde der **QR**-Algorithmus 13.7.1 in exakter Arithmetik ausgeführt. Es gelte

$$|\lambda_{p+1}| < |\lambda_p| \tag{12}$$

und

$$\sigma := \cos\left(\langle (\mathcal{Y}_1, \mathcal{S}_p) \right) > 0 \tag{13}$$

 $_{\rm mit}$

$$\mathcal{Y}_1 := \operatorname{span} \left\{ e^1, ..., e^p \right\}, \qquad \mathcal{S}_p := \operatorname{span} \left\{ u^1, ..., u^p \right\}.$$

Die durch 13.7.1 erzeugten Matrizen seien gemäß

$$egin{aligned} V_k &= (m{v}^{k,1}, \, ..., m{v}^{k,n}) = (V_1^{(k)} \mid V_2^{(k)})\,, \ A_k &= egin{pmatrix} A_{11}^{(k)} \mid A_{12}^{(k)} \ A_{21}^{(k)} \mid A_{22}^{(k)} \end{pmatrix}, \qquad m{Q}_k = egin{pmatrix} Q_{11}^{(k)} \mid Q_{12}^{(k)} \ Q_{21}^{(k)} \mid Q_{22}^{(k)} \end{pmatrix}, \qquad m{R}_k = egin{pmatrix} R_{11}^{(k)} \mid R_{12}^{(k)} \ O \mid R_{22}^{(k)} \end{pmatrix} \end{aligned}$$

partitioniert mit $V_1^{(k)} \in \mathbf{R}^{n,p}, A_{11}^{(k)} \in \mathbf{R}^{p,p}$ usw., und es sei

$$\mathcal{Y}_k := \mathcal{R}(V_1^{(k)}) = \operatorname{span} \left\{ \boldsymbol{r}^{k,1}, \dots, \boldsymbol{r}^{k,p} \right\}.$$
(14)

Dann gilt

$$\tan \varphi_{k-1} \leq \varkappa \tan \varphi_k \leq \varkappa^k \tan \varphi_1 = \varkappa^k \sqrt{1 - \sigma^2} / \sigma =: \varepsilon_{k+1}$$
(15)

mit

$$\varkappa := |\lambda_{p+1}/\lambda_p| < 1 \quad \text{und} \quad \varphi_k := \measuredangle (\mathcal{Y}_k, \mathcal{S}_p) = \measuredangle (\mathcal{Y}_k^{\perp}, \mathcal{Z}_{p+1}); \tag{16}$$

dabei ist

$$\mathcal{Y}_k^{\perp} = \mathcal{R}(V_2^{(k)}) = \operatorname{span} \left\{ \boldsymbol{v}^{k,p+1}, \ldots, \boldsymbol{v}^{k,n} \right\}, \quad \mathcal{X}_{p+1} := \operatorname{span} \left\{ \boldsymbol{u}^{p+1}, \ldots, \boldsymbol{u}^n \right\} = \mathcal{S}_p^{\perp}.$$

Überdies gelten die Abschätzungen

$$\|A_{12}^{(k)}\| = \|A_{21}^{(k)}\| \le (1 + \varkappa) \varepsilon_k \|A\| \le 2\varepsilon_k \|A\|,$$
(17)

$$\|Q_{12}^{(k)}\|, \|Q_{21}^{(k)}\| \le (1+\varkappa) \varepsilon_k \le 2\varepsilon_k, \tag{18}$$

$$|\mathbf{R}_{12}^{(k)}|| \le (1+\varkappa)^2 \,\varepsilon_k \, \|A\| \le 4\varepsilon_k \, \|A\|. \tag{19}$$

Beweis. Die Gültigkeit von (15), (16) folgt sofort aus 13.3.9 und 13.3.11 (ii) unter Beachtung von (13.3.45). Zum Nachweis von (18) gehen wir von $Q_k = V_k^{\mathsf{T}} V_{k+1}$ aus und erhalten

$$Q_{21}^{(k)} = V_2^{(k)\mathsf{T}} V_1^{(k+1)} = \begin{bmatrix} U^{\mathsf{T}} V_2^{(k)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} U^{\mathsf{T}} V_1^{(k+1)} \end{bmatrix}$$
(20)

mit $U = (u^1, ..., u^p | u^{p+1}, ..., u^n) =: (U_1 | U_2)$. Wir setzen jetzt

$$U^{\intercal}V_1^{(k+1)} = \left(rac{U_1^{\intercal}V_1^{(k+1)}}{U_2^{\intercal}V_1^{(k+1)}}
ight) =: \left(rac{F}{G}
ight).$$

419

Wegen $||U_1|| = ||V_1^{(k+1)}|| = 1$ ist $||F|| \leq 1$. Zur Abschätzung von ||G|| beachten wir, daß

$$\|\boldsymbol{G}\| = \cos\left(\langle \boldsymbol{g}_{\min} \left(\boldsymbol{\mathcal{Y}}_{k+1}, \boldsymbol{\mathcal{X}}_{p+1}\right)\right) = \sin\varphi_{k+1}$$

$$\tag{21}$$

gilt, vgl. Ü 13.3.4 und (13.3.46). Mit (15) folgt hieraus $||G|| \leq \varepsilon_{k+1}$. In analoger Weise erhält man für

$$U^{\intercal}V_2^{(k)}=:\left(\!rac{G'}{F'}\!
ight)$$

die Abschätzungen $\|F'\| \leq 1$ und $\|G'\| \leq \epsilon_k$. Damit kann (20) gemäß

$$\|oldsymbol{Q}_{21}^{(k)}\| = \|oldsymbol{G}'^{\intercal}oldsymbol{F} + oldsymbol{F}'^{\intercal}oldsymbol{G}\| \le \|oldsymbol{G}'\| \|oldsymbol{F}\| + \|oldsymbol{F}'\| \|oldsymbol{G}\| \le arepsilon_k + arepsilon_{k+1} \le (1+arphi)arepsilon_k)$$

abgeschätzt werden, und dieselbe Abschätzung gilt für $Q_{12}^{(k)}$. Zum Beweis von (17) beachten wir $A_k = Q_k R_k$, also $A_{21}^{(k)} = Q_{21}^{(k)} R_{11}^{(k)}$. Mit $||R_{11}^{(k)}|| \leq ||R_k|| = ||A||$ und (18) ergibt sich (17). Schließlich folgt aus $R_k = A_{k+1} Q_k^{\mathsf{T}}$ die Darstellung $R_{12}^{(k)} = A_{11}^{(k+1)} Q_{21}^{(k)\mathsf{T}} + A_{12}^{(k-1)} Q_{22}^{(k)\mathsf{T}}$. Mit $||A_{11}^{(k+1)}|| \leq ||A_{k+1}|| = ||A||$, (18), (17) und $||Q_{22}^{(k)}|| \leq 1$ ergibt sich daraus die noch ausstehende Ungleichung (19).

13.7.5. Bemerkung. (i) Die Ungleichung (15) besagt, daß \mathcal{Y}_k gegen \mathcal{S}_p und \mathcal{Y}_k^{\perp} gegen \mathcal{S}_{p+1} linear mit dem Konvergenzfaktor $\varkappa < 1$ konvergieren, und die Konvergenz ist monoton im Sinne von tan $\varphi_{k+1}/\tan \varphi_k \leq \varkappa$. Die Abschätzungen (17) bis (19) zeigen, daß die Außerdiagonalblöcke von A_k , Q_k und R_k mit der Majorante ε_k gegen 0 gehen. Obwohl $\varepsilon_{k+1}/\varepsilon_k = \varkappa$ gilt, braucht $||A_{12}^{(k+1)}||/||A_{12}^{(k)}||$ usw. selbst nicht durch \varkappa beschränkt zu sein, also keine monotone Konvergenz vorzuliegen.

(ii) Meist wird (12) für mehrere Indizes erfüllt sein. Es gelte etwa

$$\cdots \leq |\lambda_{p'+1}| < |\lambda_{p'}| \leq \cdots \leq |\lambda_{p+1}| < |\lambda_p| \leq \cdots,$$

und (13) sei ebenfalls für p und p' erfüllt. Mit der Partitionierung

$$A_k = egin{pmatrix} \displaystyle A_{k}^{(k)} & A_{12}^{(k)} & A_{13}^{(k)} \ \displaystyle rac{A_{21}^{(k)} & A_{22}^{(k)} & A_{23}^{(k)} \ \displaystyle A_{22}^{(k)} & A_{23}^{(k)} \ \displaystyle A_{33}^{(k)} & A_{33}^{(k)} & A_{33}^{(k)} \ \end{pmatrix}_{p'+1}^p$$

folgt dann aus 13.7.4

$$\|A_{21}^{(k)}\|, \|A_{31}^{(k)}\| \leq \left\| \left(\!\! rac{A_{21}^{(k)}}{A_{31}^{(k)}}\!\!
ight\| \leq 2 arepsilon_k \, \|A\|$$

und

$$\|A_{31}^{(k)}\|, \|A_{32}^{(k)}\| \le \|(A_{31}^{(k)} \mid A_{32}^{(k)})\| \le 2\varepsilon'_k \|A\|,$$

mithin

$$\|A_{31}^{(k)}\| \leq 2 \min \{\varepsilon_k, \varepsilon'_k\} \|A\|,$$
(22)

wobei $\varepsilon_{k+1} = \varkappa^k \tan \varphi_1, \ \varepsilon'_{k+1} = (\varkappa')^k \tan \varphi'_1 \ \mathrm{mit} \ \varkappa = |\lambda_{p+1}/\lambda_p|, \ \varkappa' = |\lambda_{p'+1}/\lambda_{p'}| \ \mathrm{ist.}$

(iii) Falls (12), (13) für p = 1, 2, ..., n - 1 erfüllt sind, gilt (22) für jedes p < p'. Da speziell $|a_{p'+1,p}^{(k)}| \leq ||A_{31}^{(k)}||$ ist, folgt daraus mit i := p' + 1, j := p die elementweise Abschätzung

$$|a_{ij}^{(k+1)}| = O((\varkappa_{ij})^k) \quad \text{mit} \quad \varkappa_{ij} := \min\{|\lambda_i/\lambda_{i-1}|, |\lambda_{j+1}/\lambda_j|\}.$$
(23)

Die Außerdiagonalelemente konvergieren also — i. allg. verschieden schnell — gegen 0, und die Diagonalelemente approximieren folglich die gesuchten Eigenwerte mit wachsendem k beliebig genau.

In der bisher betrachteten Grundform 13.7.1 ist der **QR**-Algorithmus wenig attraktiv: Für voll besetztes **A** kostet eine **QR**-Faktorisierung mittels Householder-Spiegelungen als billigster Methode $\sim 2n^3/3$ opms, siehe 10.2.11, und für die Bildung von **R**_k**Q**_k sind nochmals $\sim 2n^3/3$ opms nötig, so daß ein **QR**-Schritt insgesamt $\sim 4n^3/3$ opms erfordert. Die Konvergenz ist nur linear mit dem Faktor $\varkappa = |\lambda_{p+1}/\lambda_p|$, also i. allg. langsam, so daß viele solcher teuren **QR**-Schritte ausgeführt werden müssen.

Im folgenden diskutieren wir eine Reihe von Maßnahmen, die zur Überwindung der genannten Nachteile führen und die hohe Effektivität des entsprechend verfeinerten QR-Algorithmus bewirken.

B. Deflation und Verschiebungen

Unter den Voraussetzungen von 13.7.4 werden die Außerdiagonalblöcke von

$$A_{k} = \left(\frac{A_{11}^{(k)} | A_{12}^{(k)}}{A_{21}^{(k)} | A_{22}^{(k)}}\right) \tag{24}$$

für hinreichend großes k im Sinne von

$$\|A_{12}^{(k)}\| = \|A_{21}^{(k)}\| \le \varepsilon := \nu F \|A\|$$
(25)

mit akzeptablem F genügend klein sein und können im Rahmen der Computergenauigkeit vernachlässigt werden, d. h., A_k kann durch

$$\hat{A}_k = \left(egin{array}{c|c} A_{11}^{(k)} & O \ \hline O & A_{22}^{(k)} \end{array}
ight)$$

ersetzt werden. Dies entspricht einer Störung

$$\delta \hat{A}_k := -\left(egin{array}{c|c} O & A_{12}^{(k)} \ \hline A_{21}^{(k)} & O \end{array}
ight), \qquad \|\delta \hat{A}_k\| \leq arepsilon,$$

von A_k , wegen (6) also einer Störung

$$d ilde{A}_k := V_k d ilde{A}_k V_k^\intercal, \qquad \|d ilde{A}_k\| = \|d ilde{A}_k\| \leq arepsilon,$$

von A. Das Eigenwertproblem der Dimension n zerfällt damit in die zwei kleineren Probleme mit den Diagonalblöcken $A_{11}^{(k)}$ und $A_{22}^{(k)}$ der Dimensionen p bzw. n - p; man spricht auch hier von *Deflation*. Das Verfahren kann dann mit jedem der Diagonalblöcke einzeln und damit insgesamt billiger fortgesetzt werden, wobei V_k entsprechend aufzudatieren ist. Wird etwa mit dem unteren Block weitergearbeitet und ist $A_{22}^{(k)} = Q_{22}^{(k)} R_{22}^{(k)}$ dessen **QR**-Faktorisierung, so ist der nächste Schritt mit

$$\boldsymbol{Q}_{\boldsymbol{k}} := \left(\begin{array}{c|c} \boldsymbol{I} & \boldsymbol{O} \\ \hline \boldsymbol{O} & \boldsymbol{Q}_{22}^{(\boldsymbol{k})} \end{array} \right) \tag{26}$$

durch

$$A_{k+1} = \left(\frac{A_{11}^{(k+1)} \mid O}{O \mid A_{22}^{(k+1)}}\right) = Q_k^{\mathsf{T}} A_k Q_k = \left(\frac{A_{11}^{(k)} \mid O}{O \mid R_{22}^{(k)} Q_{22}^{(k)}}\right)$$
(27)

und

$$V_{k+1} = (V_1^{(k+1)} \mid V_1^{(k+1)}) = V_k Q_k = (V_1^{(k)} \mid V_2^{(k)} Q_{22}^{(k)})$$
(28)

definiert. Dabei wurde wieder A_k statt \hat{A}_k geschrieben.

Je schneller $A_{21}^{(k)}$ gegen O geht, um so weniger QR-Schritte sind bis zur Erfüllung der Annullierungsbedingung (25) nötig und um so geringer wird der Gesamtaufwand. Da die Konvergenzgeschwindigkeit durch $\varkappa = |\lambda_{p+1}/\lambda_p|$ bestimmt wird, versuchen wir in Analogie zum Vorgehen bei der inversen Iteration in 13.4, durch eine geeignete Spektralverschiebung

$$\boldsymbol{A} := \boldsymbol{A} - \boldsymbol{\mu} \boldsymbol{I} \tag{29}$$

einen kleineren Konvergenzfaktor zu erreichen. Wie dort numerieren wir die Eigenwerte $\lambda_j - \mu$ von \bar{A} nach wachsenden Beträgen und setzen

$$0 < |\lambda_1 - \mu| \ll |\lambda_2 - \mu| \le \dots \le |\lambda_n - \mu|$$
(30)

voraus; in der (11) entsprechenden Numerierung ist dann $\lambda_n[\bar{A}] = \lambda_1 - \mu$, $\lambda_{n-1}[\bar{A}] = \lambda_2 - \mu$ usw. Bezüglich \bar{A} ergibt sich damit für p = n - 1 der Konvergenzfaktor

$$\bar{\mathbf{x}} = |\lambda_n[\bar{\mathbf{A}}]/\lambda_{n-1}[\bar{\mathbf{A}}]| = |(\lambda_1 - \mu)/(\lambda_2 - \mu)| \ll 1.$$
(31)

Die nach 13.7.1 aus \bar{A} erzeugte Folge $\{\bar{A}_k\}$ der **QR**-Iterierten ist durch

$$\bar{A}_1 := \bar{A}$$
, faktorisiere $\bar{A}_k = \bar{Q}_k \bar{R}_k$, bilde $\bar{A}_{k+1} := \bar{R}_k \bar{Q}_k$ $(k = 1, 2, ...)$ (32)

definiert. Wenn \bar{A}_k gemäß

$$\bar{A}_{k} = \left(\frac{\bar{A}_{11}^{(k)} \mid a^{k}}{a^{k\top} \mid \bar{a}_{nn}^{(k)}}\right)$$
(33)

partitioniert wird, entspricht $a^k \in \mathbb{R}^{n-1}$ dem Block $A_{12}^{(k)}$ in 13.7.4. Unter der (13) entsprechenden Voraussetzung

$$\cos \bar{\varphi}_1 > 0, \qquad \bar{\varphi}_1 := \measuredangle (\operatorname{span} \{ e^1, \dots, e^{n-1} \}, \mathscr{S}_{n-1}) = \measuredangle (e^n, u^n)$$
(34)

gilt dann (17), also

$$\|\boldsymbol{a}^{k}\| \leq (1+\bar{z})\,\bar{\varepsilon}_{k}\,\|\bar{\boldsymbol{A}}\| \quad \text{mit} \quad \bar{\varepsilon}_{k+1} = \bar{z}\,\bar{\varepsilon}_{k} = \bar{z}^{k}\,\tan\bar{\varphi}_{1}. \tag{35}$$

Sofern tan $\bar{\varphi}_1$ nicht zu groß ist, wird der Block a^k also schnell klein werden, und $\bar{a}_{nn}^{(k)}$ konvergiert gegen $\lambda_n[\bar{A}] = \lambda_1 - \mu$.

Wenn analog zu (29)

$$\bar{A}_k = A_k - \mu I$$

gesetzt wird, schreibt sich (32) als

$$A_{1} := A, \text{ faktorisiere } A_{k} - \mu I = \bar{Q}_{k} \bar{R}_{k}, \text{ bilde } A_{k+1} := \bar{R}_{k} \bar{Q}_{k} + \mu I$$

$$(k = 1, 2, \ldots).$$
(36)

Wegen

$$\boldsymbol{A}_{k+1} = \boldsymbol{\bar{R}}_{k} \boldsymbol{\bar{Q}}_{k} + \mu \boldsymbol{I} = \boldsymbol{\bar{Q}}_{k}^{\mathsf{T}} (\boldsymbol{A}_{k} - \mu \boldsymbol{I}) \, \boldsymbol{\bar{Q}}_{k} + \mu \boldsymbol{I} = \boldsymbol{\bar{Q}}_{k}^{\mathsf{T}} \boldsymbol{A}_{k} \boldsymbol{\bar{Q}}_{k} = \boldsymbol{\bar{V}}_{k+1}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{\bar{V}}_{k+1} \quad (37)$$

sind die A_k orthogonal ähnlich zur Originalmatrix A, die A_k dagegen zu A. Bei Partitionierung von A_k in zu (33) analoger Weise ergibt sich

$$\boldsymbol{A}_{k} = \left(\frac{\boldsymbol{A}_{11}^{(k)} \mid \boldsymbol{a}^{k}}{\boldsymbol{a}^{k\intercal} \mid \boldsymbol{a}_{nn}^{(k)}}\right) = \left(\frac{\boldsymbol{\tilde{A}}_{11}^{(k)} + \mu \boldsymbol{I} \mid \boldsymbol{a}^{k}}{\boldsymbol{a}^{k\intercal} \mid \boldsymbol{\tilde{a}}_{nn}^{(k)} + \mu}\right),\tag{38}$$

d. h., der Außerdiagonalblock a^k ist derselbe wie in \bar{A}_k und geht gemäß (35) gegen o. Sobald $||a^k|| \leq \varepsilon$ mit ε aus (25) gilt, wird zu

$$\hat{A}_{k} = \left(\frac{A_{11}^{(k)} \mid o}{o^{\mathsf{T}} \mid a_{nn}^{(k)}}\right) \tag{39}$$

übergegangen. Das Diagonalelement $a_{nn}^{(k)}$ ist dann eine akzeptable Approximation für λ_1 , und das Verfahren wird mit $A_{11}^{(k)} \in \mathbf{S}^{n-1,n-1}$ fortgesetzt. Die Dimension reduziert sich dabei um 1. Sollten noch andere Außerdiagonalblöcke von A_k im Sinne von (25) vernachlässigbar sein, so zerfällt $A_{11}^{(k)}$ in weitere Diagonalblöcke niedrigerer Dimension, die dann jeder für sich weiter bearbeitet werden können.

Die bisherigen Überlegungen zeigen, daß der QR-Algorithmus mit einer Verschiebung μ , die einen Eigenwert von A möglichst gut approximiert, durchgeführt werden sollte. Dabei bietet es sich an, $\mu = \mu_k$ in der Vorschrift (36) in jedem Schritt auf Grund der aktuellen Information neu festzulegen. Dies führt auf den folgenden Algorithmus:

- 13.7.6. QR-Algorithmus mit Verschiebungen. S0: Setze $A_1 := A, V_1 := I, k := 1$ S1: Wähle Verschiebung μ_k S2: Berechne QR-Faktorisierung $A_k \mu_k I = Q_k R_k$ S3: Bilde $A_{k+1} := R_k Q_k + \mu_k I, V_{k+1} := V_k Q_k$ S4: Setze k := k + 1, goto S1

Auch für die durch 13.7.6 definierten Matrizen gilt (37) mit $\mu = \mu_k$, d. h., die A_k sind orthogonal ähnlich zu A.

13.7.7. Bemerkung. Unter der Voraussetzung

 μ_k ist kein Eigenwert von A

ist $A_k = \mu_k I$ regulär, so daß die Faktoren Q_k, R_k und damit die durch 13.7.6 erzeugten Matrizenfolgen im Sinne von (8) bis (10) im wesentlichen eindeutig durch

(40)

A und $\{\mu_k\}$ festgelegt sind, vgl. 13.7.3. Für die Durchführbarkeit ist (40) jedoch nicht erforderlich, denn die **QR**-Faktorisierung existiert auch für singuläres $A_k - \mu_k I$. Bei spaltenweiser Erzeugung von R_k durch Linksmultiplikation mit elementaren orthogonalen Matrizen entsprechend 10.2 treten dann in einer Spalte j in den Positionen $\{j, j\}$ bis $\{n, j\}$ Nullen auf, so daß sofort zur nächsten Spalte übergegangen werden kann. Bei der für uns interessanten Anwendung auf Tridiagonalmatrizen führt ein singuläres $A_k - \mu_k I$ in exakter Arithmetik zum Abspalten des Eigenwertes μ_k im nächsten Schritt, siehe 13.7.12 (ii). \Box

Wir diskutieren im folgenden zwei Möglichkeiten zur Festlegung von μ_k . Da $a_{nn}^{(k)}$ unter den Voraussetzungen (30), (34) gegen einen Eigenwert von A konvergiert, liegt es nahe, dieses letzte Diagonalelement als Verschiebung zu benutzen.

13.7.8. Rayleigh-Quotienten-Verschiebungen. Spezifiziere S1 in 13.7.6 gemäß S1: Setze $\mu_k := a_{nn}^{(k)}$

Daß die auf die Rayleigh-Quotienten-Iteration 13.4.6 hinweisende Bezeichnung dieser Verschiebung berechtigt ist, zeigt die folgende Aussage:

13.7.9. Aussage. Der **QR**-Algorithmus 13.7.6 werde mit der Verschiebungsstrategie 13.7.8 realisiert, und es gelte (40). Dann sind die *n*-ten Spalten $v^{k,n} := V_k e^n$ der Matrizen V_k in exakter Arithmetik bis auf das Vorzeichen identisch mit den Iterierten v^k , die sich nach der RQ-Iteration 13.4.6 mit dem Startvektor $v^1 := e^n$ ergeben.

Beweis. Wegen (37) gilt

$$\mu_k = a_{nn}^{(k)} = e^{n \mathsf{T}} A_k e^n = e^{n \mathsf{T}} V_k^{\mathsf{T}} A V_k e^n = v^{k,n \mathsf{T}} A v^{k,n} = \varrho(v^{k,n}),$$

d. h., μ_k ist der zu $v^{k,n}$ gehörende Rayleigh-Quotient. Aus der Iterationsvorschrift und (37) folgt weiter

$$(\boldsymbol{A} - \mu_k \boldsymbol{I}) \boldsymbol{v}^{k+1,n} = \boldsymbol{V}_k (\boldsymbol{A}_k - \mu_k \boldsymbol{I})^{\mathsf{T}} \boldsymbol{V}_k^{\mathsf{T}} \boldsymbol{v}^{k+1,n} = \boldsymbol{V}_k (\boldsymbol{Q}_k \boldsymbol{R}_k)^{\mathsf{T}} \boldsymbol{V}_k^{\mathsf{T}} \boldsymbol{v}^{k+1,n}$$
$$= \boldsymbol{V}_k \boldsymbol{R}_k^{\mathsf{T}} \boldsymbol{Q}_k^{\mathsf{T}} \boldsymbol{V}_k^{\mathsf{T}} \boldsymbol{v}^{k+1,n} = \boldsymbol{V}_k \boldsymbol{R}_k^{\mathsf{T}} \boldsymbol{V}_{k+1}^{\mathsf{T}} \boldsymbol{v}^{k+1,n} = \boldsymbol{V}_k \boldsymbol{R}_k^{\mathsf{T}} \boldsymbol{e}^n$$
$$= r_{nn}^{(h)} \boldsymbol{V}_k \boldsymbol{e}^n = r_{nn}^{(k)} \boldsymbol{v}^{k,n}.$$

Wegen (40) muß $r_{nn}^{(k)} \neq 0$ sein, so daß $\boldsymbol{w}^{k+1,n} := \boldsymbol{v}^{k+1,n}/r_{nn}^{(k)}$ der Gleichung $(\boldsymbol{A} - \mu_k \boldsymbol{I}) \boldsymbol{w}^{k+1,n} = \boldsymbol{v}^{k,n}$ genügt. Daraus folgt die Behauptung. \Box

Für das zu $\{\mu_k, v^{k,n}\}$ gehörende Residuum gilt

$$(A - \mu_k I) v^{k,n} = V_k(A_k - \mu_k I) V_k^{\mathsf{T}} v^{k,n} = V_k(A_k - \mu_k I) e^n,$$

also mit den Bezeichnungen aus (38)

$$\mathbf{r}^{k,n} := (\mathbf{A} - \mu_k \mathbf{I}) \, \mathbf{v}^{k,n} = V_k \left(\frac{\mathbf{a}^k}{0} \right), \qquad \|\mathbf{r}^{k,n}\| = \|\mathbf{a}^k\|.$$
 (41)

Falls $\{\mu_k, v^{k,n}\}$ ein Eigenpaar von A ist, ergibt sich $\|r^{k,n}\| = \|a^k\| = 0$, so daß $\mu_k = a_{nn}^{(k)}$ abgespalten werden kann und das Verfahren abbricht. Andernfalls können die Konvergenzaussagen für die RQ-Iteration aus 13.4.7 übernommen werden:

13.7.10. Satz. Für den QR-Algorithmus mit RQ-Verschiebungen fällt ||a^k|| monoton, und fast immer gilt $\|\boldsymbol{a}^{k}\| \to 0$, $a_{nn}^{(k)} \to \lambda_{j}$, $\vartheta_{k}\boldsymbol{v}^{k,n} \to \boldsymbol{u}^{j}$ mit geeignetem $\vartheta_{k} \in \{+1, -1\}$. Für genügend großes k ist dabei $\tan \varphi_{k+1} \leq \eta_{1} (\tan \varphi_{k})^{3}$ (42) mit $\eta_{1} > 0$ und $\varphi_{k} := \langle (\boldsymbol{u}^{j}, \vartheta_{k}\boldsymbol{v}^{k,n}),$ und es gelten die Abschätzungen $|\hat{\boldsymbol{a}}_{k} = q_{k}^{(k)}| \leq 2 \|\boldsymbol{A}\|$ (tan \boldsymbol{w}_{k})² $\|\boldsymbol{a}^{k}\| \leq w$ tan \boldsymbol{w}_{k} mit $\boldsymbol{w}_{k} > 0$ (43)

$$\tan \varphi_{k+1} \le \eta_1 (\tan \varphi_k)^3 \tag{42}$$

$$|\lambda_j - a_{nn}^{(k)}| \leq 2 \|A\| (\tan \varphi_k)^2, \qquad \|a^k\| \leq \eta_2 \tan \varphi_k \quad \text{mit} \quad \eta_2 > 0.$$

$$(43)$$

13.7.8 konvergiert also fast immer und monoton, und die Konvergenz ist asymptotisch kubisch.

Die in A_k enthaltene Information wird durch die folgende, auf WILKINSON zurückgehende Verschiebungsstrategie noch besser ausgenutzt:

13.7.11. Wilkinson-Verschiebungen. Spezifiziere S1 in 13.7.6 wie folgt:

S1: Berechne die Eigenwerte $\{\mu', \mu''\}$ der in der rechten unteren Ecke von A_k stehenden (2, 2)-Matrix

$$\boldsymbol{P}_{k} := \begin{pmatrix} a_{n-1,n-1}^{(k)} & a_{n-1,n}^{(k)} \\ a_{n,n-1}^{(k)} & a_{nn}^{(k)} \end{pmatrix} \in \mathbf{S}^{2,2}$$

$$\tag{44}$$

Ordne
$$\{\mu', \mu''\}$$
 gemäß $|\mu' - a_{nn}^{(k)}| < |\mu'' - a_{nn}^{(k)}|$ bzw. $|\mu'| \leq |\mu''|$, wenn $|\mu' - a_{nn}^{(k)}| = |\mu'' - a_{nn}^{(k)}|$
Setze $\mu_k := \mu'$

Man beachte, daß P_k gerade die zu $Q := (v^{k,n-1}, v^{k,n})$ gehörende Matrix aus dem Rayleigh-Ritz-Algorithmus ist und $\{\mu', \mu''\}$ die entsprechenden Ritzschen Eigenwerte bezüglich span $\{v^{k,n-1}, v^{k,n}\}$ sind, siehe 13.3.C, so daß 13.7.11 als natürliche Erweiterung von 13.7.8 auf den zweidimensionalen Fall angesehen werden kann. Zur numerischen Realisierung siehe Ü 13.7.2. Wir werden später sehen, daß 13.7.11 bei Anwendung auf eine Tridiagonalmatrix sogar globale Konvergenz garantiert.

Bei Verwendung der obigen Verschiebungsstrategien besteht die Tendenz, daß die Eigenwerte in betragsmäßig wachsender Reihenfolge abgespalten werden. Allerdings braucht das nicht immer der Fall zu sein.

C. Der QR-Algorithmus für Tridiagonalmatrizen

Mit den oben angegebenen Verschiebungen erhöht sich die Konvergenzgeschwindigkeit des **OR**-Algorithmus drastisch; in der Regel sind nur einige Schritte bis zur Abspaltung eines Eigenwertes nötig. Selbst wenn pro Eigenwert nur ein solcher Schritt erforderlich wäre, würde jedoch der Aufwand zur Berechnung aller Eigenwerte einer voll besetzten Matrix A bei Berücksichtigung der Dimensionsreduktion

 $\sum\limits_{j=1}^n 4j^3/3 \sim n^4/3$ opms betragen, also um eine n-Potenz höher sein als beim Jacobi-Verfahren, vgl. 13.2. Die Kosten für einen **QR-**Schritt können jedoch signifikant

gesenkt werden, wenn A vor Ausführung des QR-Algorithmus nach den in 13.5

beschriebenen Methoden orthogonal ähnlich auf symmetrische Tridiagonalform

$$\boldsymbol{T} = \operatorname{trid} \left(a_1, \dots, a_n, b_2, \dots, b_n \right) = \boldsymbol{Q}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{Q} \tag{45}$$

transformiert wird. Nach 13.5.1 ist dies mittels Householder-Spiegelungen mit $\sim 2n^3/3$ opms möglich, kostet also nur die Hälfte eines einzelnen **QR**-Schrittes für voll besetztes **A**. Anschließend wird der **QR**-Algorithmus 13.7.6 mit $A_1 := T$ durchgeführt.

Es zeigt sich, daß mit der Startmatrix A_1 auch alle folgenden QR-Iterierten A_k tridiagonal sind: Es sei A_k tridiagonal. Zur Berechnung der QR-Faktorisierung

$$\boldsymbol{A}_{\boldsymbol{k}} - \boldsymbol{\mu}_{\boldsymbol{k}} \boldsymbol{I} = \boldsymbol{Q}_{\boldsymbol{k}} \boldsymbol{R}_{\boldsymbol{k}} \tag{46}$$

bietet es sich an, das in der *j*-ten Spalte stehende Subdiagonalelement mittels einer Givens-Drehung $G_{j,j+1} = G_{j,j+1}^{(k)}$ zu annullieren, was auf

$$G_{n-1,n}\Big(\cdots\Big(G_{23}\big(G_{12}(A_k-\mu_kI)\big)\Big)\cdots\Big)=R_k$$
(47)

führt, vgl. 10.2. Dann ist

$$\boldsymbol{Q}_{k} = \boldsymbol{G}_{12}^{\mathsf{T}} \boldsymbol{G}_{23}^{\mathsf{T}} \cdots \boldsymbol{G}_{n-1,n}^{\mathsf{T}}, \tag{48}$$

und wegen der Tridiagonalform von $A_k - \mu_k I$ hat R_k die Gestalt

ist also eine obere Dreiecksmatrix der Bandbreite 3 vom Typ $\{0, 2\}$; zur Bezeichnung siehe 6.4.A. Bei Bildung von

$$\boldsymbol{A}_{k+1} - \boldsymbol{\mu}_{k}\boldsymbol{I} = \boldsymbol{R}_{k}\boldsymbol{Q}_{k} = \left(\cdots\left((\boldsymbol{R}_{k}\boldsymbol{G}_{12}^{\mathsf{T}})\ \boldsymbol{G}_{23}^{\mathsf{T}}\right)\cdots\right)\boldsymbol{G}_{n-1,n}^{\mathsf{T}}$$
(49)

werden nacheinander die Spalten j und j + 1 kombiniert. Die Matrix $R_k Q_k$ ist daher eine obere Hessenberg-Matrix, wegen der Symmetrie also notwendig tridiagonal.

13.7.12. Bemerkung. (i) Wir haben gezeigt, daß die Tridiagonalform von A_k unter der gemäß (46) bis (49) realisierten **QR**-Transformation

$$oldsymbol{A}_k o oldsymbol{A}_{k+1} = oldsymbol{Q}_k^{\mathsf{T}} oldsymbol{A}_k oldsymbol{Q}_k = oldsymbol{R}_k oldsymbol{Q}_k + \mu_k oldsymbol{I}$$

invariant ist. Unter der zusätzlichen Voraussetzung (40) gilt dies auch für jede andere Realisierung des **QR**-Schrittes, denn die transformierten Matrizen A_{k+1} können sich dann höchstens im Vorzeichen der Elemente unterscheiden, vgl. 13.7.3 und 13.7.7. (ii) Wenn A_k nicht zerfällt, also

$$a_{i,i-1}^{(k)} = b_i^{(k)} \neq 0 \qquad (i = 2, ..., n)$$
(50)

gilt, und (40) nicht erfüllt, also μ_k ein Eigenwert von A ist, hat die letzte Zeile von A_{k+1} bei Rechnung in exakter Arithmetik die Gestalt $(0, \ldots, 0, a_{nn}^{(k+1)})$, siehe Ü 13.7.3. Es ist also $a_{n,n-1}^{(k-1)} = 0$, und $a_{nn}^{(k+1)} = \mu_k$ wird als Eigenwert erkannt und kann abgespalten werden. Die Voraussetzung (50) stellt dabei keine Einschränkung dar, denn andernfalls hätte bereits im vorhergehenden Schritt zu nichtzerfallenden Teilproblemen niedrigerer Dimension übergegangen werden können, vgl. 13.7.B.

(iii) Bei Verwendung der RQ-Verschiebungen besagt 13.7.10, daß $|a_{n,n-1}^{(k)}| = |b_n^{(k)}|$ in exakter Arithmetik fast immer monoton und asymptotisch kubisch gegen 0 geht. Bei Verwendung der Wilkinson-Verschiebungen kann sogar gezeigt werden, daß $b_n^{(k)}$ immer gegen 0 konvergiert, d. h., es liegt globale Konvergenz vor. Die Konvergenz ist dabei mindestens quadratisch im Sinne der Gültigkeit von

$$|b_n^{(k+1)}| \leq C_1 |b_n^{(k)}|^2 \quad \text{mit} \quad C_1 > 0$$

und meistens sogar besser als kubisch. Für eine Formulierung des letztgenannten Sachverhaltes und die Beweise muß auf die Spezialliteratur verwiesen werden, siehe B 13.9.

Für die numerische Realisierung ist wesentlich, daß die ersten j Spalten der in (47) auftretenden intermediären Matrix

$$\boldsymbol{G}_{j,j+1}\left(\boldsymbol{G}_{j-1,j}\left(\cdots\left(\boldsymbol{G}_{12}(\boldsymbol{A}_{k}-\mu_{k}\boldsymbol{I})\right)\right)\right)=:\boldsymbol{A}_{k}^{(j-1)}$$

bei den nachfolgenden Linksmultiplikationen mit $G_{j+1,j+2}$ usw. weder geändert noch benutzt werden. Sie stellen daher bereits die entsprechenden Spalten von \mathbf{R}_k dar, und die auf diese Spalten gemäß (49) von rechts angewendeten Drehungen $G_{12}^{\mathsf{T}}, \ldots,$ $G_{j-1,j}^{\mathsf{T}}$ können bereits ausgeführt werden, ohne daß sie die noch benötigten Spalten $j + 1, \ldots, n$ von $A_k^{(j-1)}$ beeinflussen. Die Matrix \mathbf{R}_k braucht als Ganzes also überhaupt nicht gebildet und gespeichert zu werden, sondern die Transformation kann direkt in der durch die Klammern gekennzeichneten Reihenfolge gemäß

$$A_{k+1} - \mu_k I = R_k Q_k = \cdots \left\{ G_{45} [\left(G_{34} \left\{ \left[G_{23} G_{12} (A_k - \mu_k I) \right] G_{12}^{\mathsf{T}} \right\} \right) G_{23}^{\mathsf{T}} \right] \right\} G_{34}^{\mathsf{T}} \cdots$$
(51)

realisiert werden. Da hier die Verschiebungen explizit ausgeführt werden, spricht man auch vom **QR**-Algorithmus mit explizit realisierten Verschiebungen oder kurz vom *expliziten* **QR**-Algorithmus. Um unnötige Indizierungen zu vermeiden, geben wir einen Schritt bei in-situ-Realisierung auf dem Platz von **A** an, wobei außerdem V_k auf dem Platz von **V** entsprechend (7) aufdatiert wird.

13.7.13. Expliciter **QR**-Schritt. Setze $A := A_k - \mu_k I$, $V := V_k$, $G_{01} := I$ for j := 1(1)n - 1 do | Lege $G_{j,j+1}$ so fest, daß $(G_{j,j+1}A)_{j+1,j} = 0$ gilt | Bilde $A := G_{j,j+1}A$, $A := AG_{j-1,j}^{\mathsf{T}}$, $V := VG_{j,j+1}^{\mathsf{T}}$, Setze $A := AG_{n-1,n}^{\mathsf{T}}$, $A_{k+1} := A + \mu_k I$, $V_{k+1} := V$

Der Teilschritt $j \rightarrow j + 1$ wird für n = 5 und j = 3 durch das folgende Muster illustriert, wobei sich verändernde Elemente wie üblich eingerahmt worden sind.

Die durch die stark gezogenen Linien links bzw. rechts unten abgegrenzten Elemente sind bereits die von $R_k Q_k$ bzw. noch die von $A_k - \mu_k I$.

13.7.14. Bemerkung. (i) Unter Beachtung der Symmetrie läßt sich 13.7.13 bei expliziter Realisierung der Givens-Drehungen mit $\sim n(K_1 \text{ ops} + K_2 \text{ opm} + 1 \text{ opr})$ ausführen. Die genauen Werte der K_i hängen von den Details der Implementierung ab; für die bei WILKINSON/REINSCH [71] enthaltene Version ist $K_1 = 4$ und $K_2 = 13$. Die Berechnung von V_{k+1} erfordert dagegen bei Realisierung gemäß U 3.3.6 weitere $\sim 3n^2$ opms, also we sentlich mehr. Unter Verwendung einiger Hilfsspeicher kann die Transformation $A_k \rightarrow A_{k+1}$ auf den $\sim 2n$ Plätzen eines Dreiecks von T realisiert werden, und V_k kann offensichtlich mit V_{k+1} überspeichert werden.

(ii) Beim expliziten **QR**-Algorithmus ist es zweckmäßig, in Abänderung von 13.7.6 sofort

$$A_{k+1} = R_k Q_k = Q_k^{\mathsf{T}} (A_k - \mu_k I) Q_k$$

zu setzen, also die Addition von $\mu_k I$ bei der Bildung von A_{k+1} wegzulassen. Für die derart definierte Folge gilt

$$A_{k+1} = V_{k+1}^{\mathsf{T}} [A - (\mu_1 + \dots + \mu_k) I] V_{k+1},$$

d. h., $A_{k+1} + \sigma_{k+1}I$ mit $\sigma_{k+1} := \sum_{j=1}^{k} \mu_j$ ist orthogonal ähnlich zu A. Die Verschiebungen müssen dann gemäß $\sigma_1 := 0$, $\sigma_{k+1} := \sigma_k + \mu_k$ (k = 1, 2, ...) gesondert akkumuliert werden.

Eine zu 13.7.13 alternative Realisierung des **QR**-Schrittes beruht auf der Tatsache, daß Q_k und damit A_{k+1} und V_{k+1} im wesentlichen eindeutig durch A_k festgelegt sind, wenn die erste Spalte von Q_k vorgegeben wird und A_k nicht zerfällt, vgl. 13.5.8. Für die durch (48) definierte Matrix Q_k gilt nun

$$Q_k e^1 = G_{12}^{\mathsf{T}} G_{23}^{\mathsf{T}} \cdots G_{n-1,n}^{\mathsf{T}} e^1 = G_{12}^{\mathsf{T}} e^1,$$
(52)

d. h., die erste Spalte von Q_k ist gleich der ersten Spalte von G_{12}^{T} und damit durch die erste Spalte von $A_k - \mu_k I$ festgelegt, denn G_{12} transformiert letztere in ein Vielfaches von e^1 . Jede Matrix

$$\overline{\mathcal{Q}}_k := \mathbf{G}_{12}^{\mathsf{T}} \overline{\mathbf{G}}_{23}^{\mathsf{T}} \overline{\mathbf{G}}_{34}^{\mathsf{T}} \cdots \overline{\mathbf{G}}_{n-1,n}^{\mathsf{T}}$$
(53)

besitzt folglich für beliebige $\bar{G}_{23}, \ldots, \bar{G}_{n-1,n}$ dieselbe erste Spalte wie Q_k . Wenn es gelingt, die Drehungen $\bar{G}_{23}, \ldots, \bar{G}_{n-1,n}$ so festzulegen, daß

$$\bar{\boldsymbol{A}}_{k+1} := \bar{\boldsymbol{Q}}_{k}^{\mathsf{T}} \boldsymbol{A}_{k} \bar{\boldsymbol{Q}}_{k} = \bar{\boldsymbol{G}}_{n-1,n} \cdots \bar{\boldsymbol{G}}_{23} \boldsymbol{G}_{12} \boldsymbol{A}_{k} \boldsymbol{G}_{12}^{\mathsf{T}} \bar{\boldsymbol{G}}_{23}^{\mathsf{T}} \cdots \bar{\boldsymbol{G}}_{n-1,n}^{\mathsf{T}}$$
(54)

tridiagonal ist, muß $\{\bar{A}_{k+1}, \bar{Q}_k\}$ im Sinne von 13.5.8 im wesentlichen mit $\{A_{k+1}, Q_k\}$ übereinstimmen. Nun ist $F_2 := G_{12}A_kG_{12}^{\mathsf{T}}$ von der Form

weicht also nur in den Positionen (3, 1) und (1, 3) von der Tridiagonalform ab. Wenn \overline{G}_{23} so festgelegt wird, daß $(\overline{G}_{23}F_2)_{31} = 0$ gilt, hat $F_3 := \overline{G}_{23}F\overline{G}_{23}^{\mathsf{T}}$ in diesen Positionen Nullen, allerdings sind in $\{4, 2\}$ und $\{2, 4\}$ neue Nichtnullelemente erzeugt worden. Diese werden durch Übergang zu $F_4 := \overline{G}_{34}F_3\overline{G}_{34}^{\mathsf{T}}$ in analoger Weise annulliert usw. Nach n - 2 Schritten ist die Tridiagonalform wieder hergestellt, vgl. 11.1 und 13.5.B für ein ähnliches Vorgehen. Die Reihenfolge, in der Nichtnullelemente außerhalb des tridiagonalen Bandes auftreten und wieder annulliert werden, wird für n = 5 durch das folgende Muster charakterisiert.

$$\begin{pmatrix} \times & \times & 1 \\ \times & \times & \times & 2 \\ 1 & \times & \times & 3 \\ & & 2 & \times & \times \\ & & & 3 & \times & \times \\ & & & & 3 & \times & \times \\ \end{pmatrix}$$

Da beim derart ausgeführten QR-Schritt die Verschiebung μ_k nicht mehr explizit auftritt — sie wird lediglich bei der Festlegung der Drehungsparameter von G_{12} benötigt —, heißt diese Version auch QR-Algorithmus mit impliziter Realisierung der Verschiebungen oder kurz *impliziter* QR-Algorithmus. Eine 13.7.13 entsprechende Grobbeschreibung des Einzelschrittes lautet wie folgt:

13.7.15. Impliziter **QR**-Schritt. Setze $A := A_k$, $V := V_k$ Lege G_{12} so fest, daß $[G_{12}(A - \mu_k I)]_{21} = [G_{12}(a_{11} - \mu_k, a_{21}, 0, ..., 0)^{\mathsf{T}}]_2 = 0$ gilt. Bilde $A := G_{12}AG_{12}^{\mathsf{T}}$, $V := VG_{12}^{\mathsf{T}}$ for i := 3(1)n do | Lege $\bar{G}_{i-1,i}$ so fest, daß $(\bar{G}_{i-1,i}A)_{i,i-2} = 0$ gilt Bilde $A := \bar{G}_{i-1,i}A\bar{G}_{i-1,i}^{\mathsf{T}}$, $V := V\bar{G}_{i-1,i}^{\mathsf{T}}$ Setze $A_{k+1} := A$, $V_{k+1} := V$ 13.7.16. Bemerkung. (i) Von der Implementierung her ist der implizite QR-Algorithmus etwas günstiger als der explizite. Der Aufwand ist nur unwesentlich geringer; für die bei WILKINSON/REINSCH [71] enthaltene Version werden $\sim n(6 \text{ ops} + 11 \text{ opm} + 1 \text{ opr})$ zur Berechnung von A_{k+1} aus A_k benötigt. Die Bildung von V_{k+1} entspricht dem expliziten Verfahren und kostet $\sim 3n^2$ opms. In-situ-Realisierung ist bezüglich A_k und V_k möglich.

(ii) Der eigentliche Vorteil des impliziten **QR**-Algorithmus liegt darin, daß $A_k - \mu_k I$ nicht explizit gebildet und daher ein möglicher Informationsverlust bei der Berechnung von $a_{jj}^{(k)} - \mu_k$ bei betragsmäßig großem μ_k und kleinem $a_j^{(k)}$ vermieden wird. Selbst für ansteigend gestufte Matrizen liefert der implizite Algorithmus daher oft auch die betragskleinen Eigenwerte mit ausreichender Genauigkeit.

(iii) Besonders für den impliziten QR-Algorithmus ist wesentlich, daß nach jedem Schritt *alle* Subdiagonalelemente — nicht nur das letzte! — auf Kleinheit etwa im Sinne von

$$|b_i^{(k)}| \le \nu(|a_{i-1}^{(k)}| + |a_i^{(k)}|), \qquad i \in \{2, ..., n\},$$
(55)

getestet und bei Erfülltsein von (55) gleich 0 gesetzt werden, vgl. 13.6.1. Wenn ein kleines $b_i^{(k)}$ mit $2 \leq i \leq n-1$ nicht zur Aufspaltung in kleinere Teilprobleme genutzt wird, erhöht sich nicht nur wegen der verschenkten Dimensionsreduktion der Aufwand für einen QR-Schritt, sondern das Konvergenzverhalten kann sich wesentlich verschlechtern: Mit $b_i^{(k)}$ wird meist auch $b_i^{(k+1)}$ klein sein. Dies bedeutet, daß die ersten Spalten $\{1, ..., i-1\}$ von \bar{Q}_k – man beachte $A_k = \bar{Q}_k \bar{A}_{k+1} \bar{Q}_k^{\mathsf{T}}$ – nur schwach mit den letzten Spalten $\{i, ..., n\}$ gekoppelt sind, was z. B. aus dem Beweis von 13.5.8 oder aus der Rekursionsformel des Lanczos-Algorithmus aus Ü 13.5.6 zu ersehen ist. Die in der ersten Spalte von \bar{Q}_k enthaltene Information über die Verschiebung μ_k , die mittels der Drehungen $\bar{G}_{i-1,i}$ implizit auf die übrigen Spalten übertragen werden soll, geht dann bei der Durchführung des k-ten QR-Schrittes teilweise oder vollständig verloren; die Drehungen $\overline{G}_{i-1,i}, \ldots, \overline{G}_{n-1,n}$ werden deformiert, vgl. auch 13.5.4 (ii) für ähnliche Überlegungen bei Verwendung von Householder-Spiegelungen. Die durch 13.5.8 in exakter Arithmetik garantierte Äquivalenz der explizit bzw. implizit berechneten Matrizen A_{k+1} bzw. \bar{A}_{k+1} liegt bei Computerrechnung nicht mehr vor, was zum Verlust der schnellen Konvergenz führt.

(iv) Aus den in (iii) dargelegten Gründen ist es zweckmäßig, das in den EISPACK-Routinen verwendete, relativ einfache Kriterium (55) durch ein selektiveres zu ersetzen. Wir empfehlen die Verwendung von

$$|b_i| \leq \nu(\beta_i + \gamma_i) (1 + \alpha_i/\beta_i), \qquad i \in \{2, \dots, n\},$$
(56)

mit

$$lpha_i:=|a_{i-1}|+|a_i|\,,\qquad eta_i:=|b_{i-1}|+|b_i|+|b_{i+1}|\,,\qquad \gamma_i:=|a_{i-1}-a_i|\,;$$

der obere Index k wurde weggelassen. Da dieser Test aufwendig ist, sollte er nur ausgeführt werden, wenn beispielsweise $|b_i| \leq \sqrt{\nu} ||T||_{\infty}$ gilt. \Box

Nach einer nicht zu großen Zahl von Schritten wird (55) oder (56) für alle *i* erfüllt sein, d. h., A_{k+1} ist diagonal. Die Diagonalelemente approximieren dann die gesuchten Eigenwerte, und die Spalten von V_{k+1} sind zugehörige Eigenvektornäherungen zu T.

429

Mit Q aus (45) ergeben sich schließlich die Eigenvektornäherungen zu A als Spalten von QV_{k+1} . Wenn die Eigenvektoren auch gesucht sind, sollte daher sofort QV_{k+1} durch Aufdatierung von QV_k mit Q_k ermittelt werden, d. h., 13.7.6 sollte mit $V_1 := Q$ statt $V_1 := I$ gestartet werden.

Wir geben im folgenden das Grobschema einer in-situ-Realisierung des QR-Algorithmus für eine Tridiagonalmatrix T an, bei der die obigen Überlegungen berücksichtigt werden. Die Indizes p, q mit $1 \leq p \leq q \leq n$ geben dabei die Position des gerade behandelten, untersten nichtzerfallenden Diagonalblockes von T an.

13.7.17. **QR**-Algorithmus für Tridiagonalmatrizen. Gegeben seien die Tridiagonalmatrix $T = \text{trid} (a_1, ..., a_n, b_2, ..., b_n) \in S^{n,n}$ und die orthogonale Matrix $V \in \mathbb{R}^{n,n}$. Algorithmus:

- S0 (Initialisierung): Setze k := 1, q := n
- S1 (Abbruchtest): If q = 1 then stop
- S2 (Nullsetzen kleiner Subdiagonalelemente und Festlegung des untersten nichtzerfallenden Diagonalblockes $T^{(p,q)} := \text{trid} (a_p, \ldots, a_q, b_{p+1}, \ldots, b_q)$ der Dimension m := q - p + 1 von T): for p := q(-1)2 do

 $\begin{vmatrix} \text{ if } [|b_p| \text{ klein im Sinne von (56)}] \\ \text{ then } & b_p := 0 \\ \text{ if } p = q \text{ then } [q := q - 1, \text{ goto S1}] \text{ else goto S3} \\ p := 1 \end{aligned}$

S3 (k-ter **QR**-Schritt): Führe **QR**-Schritt 13.7.13 bzw. 13.7.15 mit der Wilkinson-Verschiebung 13.7.11 für

$$A_{k} := T^{(p,q)} \in \mathsf{S}^{m,m}, \qquad V_{k} := V^{(p,q)} \in \mathsf{R}^{n,m}$$

auf dem Platz von T bzw. V durch, wobei $V^{(p,q)}$ die aus den Spalten p bis q von V bestehende Teilmatrix bezeichnet

S4: Setze k := k + 1, goto S1

13.7.18. Bemerkung. (i) Umfangreiche numerische Experimente zeigen, daß die Eingangsmatrix T im Mittel nach $k \approx 1.7n$ QR-Schritten diagonalisiert ist, d. h., bei Abbruch ist $T = T_{k+1} = \text{diag}(a_j^{(k+1)}) =: \text{diag}(\hat{\mu}_j) =: M$. Für die zuerst berechneten Eigenwerte $\hat{\mu}_n, \hat{\mu}_{n-1}, \ldots$ werden dabei i. allg. etwas mehr — etwa 4 bis 6 — Schritte benötigt, die zuletzt berechneten Eigenwerte $\ldots, \hat{\mu}_3, \hat{\mu}_2, \hat{\mu}_1$ können meist schon nach einem, höchstens zwei Schritten abgespalten werden. Wenn pro Eigenwert zwei QR-Schritte veranschlagt werden und ein Schritt der Dimension m ohne Aufdatierung von V mit $\sim m(10 \text{ opms} + 1 \text{ opr})$ gezählt wird, vgl. 13.7.14 und 13.7.16, kostet die Berechnung aller Eigenwerte von T bei dieser Modellvorstellung $\sim n^2(10 \text{ opms} + 1 \text{ opr})$, ist also extrem billig. Wegen der zusätzlich auftretenden Deflation außerhalb der letzten Zeile — in 13.7.17 durch p > 1 charakterisiert — ist der Aufwand praktisch sogar meist geringer. Die zugehörigen Eigenvektor-approximationen ergeben sich mit der Eingangsbelegung V := I durch Aufdatierung von V mit maximal $\sim 3n^3$ opms.

(ii) Zur Berechnung aller Eigenwerte der voll besetzten Matrix A muß zunächst $T = Q^{\intercal}AQ$ mit $\sim 2n^3/3$ opms nach 13.5.1 berechnet werden, so daß insgesamt $\sim n^2[(2n/3 + 10) \text{ opms} + 1 \text{ opr}] \sim 2n^3/3$ opms erforderlich sind. Wenn auch die Eigenvektoren gefordert werden, wird Q mit $\sim 2n^3/3$ opms explizit gebildet — siehe 13.5.2 — und in 13.7.17 als Eingangsbelegung V := Q verwendet. Mit den zur Aufdatierung benötigten $\sim 3n^3$ opms kostet daher die Berechnung der aus den Eigenvektorapproximationen zu A gebildeten Ausgangsmatrix $V = QV_{k+1}$ maximal $\sim 11n^3/3$ opms, ist also wesentlich teurer als die Berechnung der Eigenwerte allein.

(iii) Sind nur wenige Eigenvektoren gesucht, so sollten die entsprechenden Eigenvektoren von T mittels inverser Iteration mit jeweils $\sim K_3 n$ opms berechnet werden, siehe 13.4.3. Unter Verwendung der Produktform von Q sind die Eigenvektoren v^j von T dann noch in Eigenvektoren Qv^j von A zu transformieren, was jeweils $\sim n^2$ opms kostet. Insbesondere bei mehrfachen und dicht benachbarten Eigenwerten ist dabei i. allg. eine Re-Orthogonalisierung erforderlich, vgl. 13.4.5 (iv).

(iv) Wenn A eine Bandmatrix ist, kann die Tridiagonalisierung nach 13.5.6 erfolgen. Eine alternative Möglichkeit besteht darin, den QR-Algorithmus direkt mit $A_1 := A$ durchzuführen, denn die Bandgestalt ist invariant unter der QR-Transformation. Bei impliziter Realisierung wird dabei wie in 13.5.B vorgegangen. Ist die Dimension n groß und sind nur wenige Eigenwerte und gegebenenfalls Eigenvektoren gesucht, so empfiehlt sich das Bisektionsverfahren bzw. die inverse Iteration direkt mit der Bandmatrix A, vgl. 13.6.11(vi) und 13.5.7(ii).

(v) Man kann zeigen, daß die nach 13.7.17 bei Abbruch nach k Schritten vorliegende, berechnete Diagonalmatrix $M = T_{k+1}$ der Gleichung

$$A + \delta A = VMV^{\intercal} \quad \text{mit} \quad \delta A \in \mathbf{S}^{n,n}, \qquad \|\delta A\| \leq \nu F \|A\|, \tag{57}$$

und akzeptablem F genügt, wobei V eine gewisse exakt orthogonale Matrix ist. Die Berechnung der Eigenwerte von A nach dem QR-Algorithmus ist also ein numerisch gutartiger Proze β . Das exakt orthogonale V wird durch die mit $V = V_1 := Q$, Q aus 13.5.1, in 13.7.17 berechnete Matrix $\tilde{V} = V_{k+1}$ im Sinne von

$$\|V - \tilde{V}\| \le \nu F_1, \qquad F_1 \text{ akzeptabel}, \tag{58}$$

ausreichend genau approximiert. Die berechneten Eigenvektornäherungen sind daher ausreichend orthogonal und akzeptabel im Sinne kleiner Residuen.

(vi) Wenn T eine ansteigend gestufte Matrix ist, werden die Verschiebungen am Anfang groß sein. Die im betragskleinen oberen Teil enthaltene Information über die kleinen Eigenwerte wird daher — beim expliziten QR-Algorithmus in katastrophaler Weise, beim impliziten weniger stark — verlorengehen. Um dies zu vermeiden, kann der QR-Algorithmus in inverser Reihenfolge von rechts nach links durchgeführt und die Verschiebung durch die links oben stehende (2,2)-Matrix festgelegt, also die QR-Faktorisierung von $A_k - \mu_k I$ durch die analog gebildete QL-Faktorisierung $A_k - \mu_k I = Q_k L_k$ mit einer unteren Dreiecksmatrix L_k ersetzt werden. Diese Version wird QL-Algorithmus genannt und heute in den Programmpaketen vorwiegend verwendet, und zwar in impliziter Realisierung als *impliziter QL*-Algo*rithmus.* Wenn A bereits ansteigend gestuft ist, sollte die Tridiagonalisierung ebenfalls in inverser Reihenfolge von rechts nach links erfolgen, siehe 13.5.4(iii) und Ü 13.5.2.

Übungsaufgaben

Ü 13.7.1. Man zeige, daß für die Grundform 13.7.1 des QR-Algorithmus $R_k = V_{k+1}^{\mathsf{T}} A V_k$ gilt, und folgere hieraus die Gültigkeit von

$$V_{k+1}(R_kR_{k-1}\cdots R_2R_1)=A^k.$$

Ú 13.7.2. Gegeben sei die Matrix $\mathbf{P} = \begin{pmatrix} \alpha & \beta \\ \beta & \gamma \end{pmatrix}$ mit $\beta \neq 0$, und es sei $\delta := \frac{\gamma - \alpha}{2\beta}$. Man zeige: (i) Ein Schritt des Jacobi-Verfahrens aus 13.2 liefert bei Anwendung auf \mathbf{P} die Eigenwerte $\mu' = \gamma + \beta t, \ \mu'' = \gamma - \beta/t$, wobei

$$t := \mathrm{sgn} \left(\delta \right) / \left(\left| \delta \right| + \sqrt{1 + \delta^2} \right) \quad \mathrm{im \ Fall} \quad \delta \neq 0,$$

$$t:=1$$
 im Fall $\delta=0$

ist.

(ii) Es gilt

 $ert \mu' - \gamma ert < ert \mu'' - \gamma ert$ im Fall $\delta \neq 0$, $ert \mu' - \gamma ert = ert \mu'' - \gamma ert$ im Fall $\delta = 0$,

so daß

$$\mu := \begin{cases} \gamma + \mathrm{sgn} (\delta) \beta / (|\delta| + \sqrt{1 + \delta^2}) & \text{für } \delta \neq 0, \\ \gamma - \mathrm{sgn} (\gamma) |\beta| & \text{für } \delta = 0, \gamma \neq 0 \\ \beta & \text{für } \delta = 0, \gamma = 0 \end{cases}$$

die zu P gehörende Wilkinson-Verschiebung gemäß 13.7.11 ist.

Ü 13.7.3. Es sei $T - \mu I = QR$ eine QR-Faktorisierung von $T - \mu I$, und T sei tridiagonal und nichtzerfallend. Man zeige, daß dann $r_{jj} \neq 0$ (j = 1, ..., n - 1) gilt und daß $T - \mu I$ genau dann singulär ist, wenn $r_{nn} = 0$ gilt. Die letzte Bedingung ist gleichwertig dazu, daß die *n*-te Zeile von RQ verschwindet, also $RQ + \mu I$ die *n*-te Zeile $(0, ..., 0, \mu)$ hat.

Bemerkungen zum Kapitel 13

B 13.1. Die im Abschnitt 13.1 zusammengestellten Ergebnisse gehören zum Standardwissen der Eigenwerttheorie, siehe etwa WILKINSON [65] und PARLETT [80a], wo auch historische Kommentare zu finden sind. Der Rayleigh-Quotient wurde durch Lord RAYLEIGH [1899] eingeführt.

B 13.2. Das klassische Jacobi-Verfahren geht auf JACOBI [1846] zurück. Es wurde 100 Jahre später zu Beginn des Computer-Zeitalters durch BARGMANN/MONTGOMERY/VON NEUMANN [1946] wiederentdeckt; vgl. auch GOLDSTINE/MURRAY/VON NEUMANN [59]. Schwellwert-strategien sind von POPE/TOMPKINS [57], RUTISHAUSER [66] u. a. vorgeschlagen und analysiert worden. Zum Nachweis der asymptotisch quadratischen Konvergenz des zyklischen Jacobi-Verfahrens sei auf SCHÖNHAGE [64] und WILKINSON [62] verwiesen. Eine ausgefeilte

433

Implementierung hat RUTISHAUSER [66/71] publiziert; dort sind alle im Text beschriebenen Tricks zur Verringerung des Rundungsfehlereinflusses zu finden. Die zweite Formel (13.2.23), mit der sich eine Multiplikation einsparen läßt, scheint nicht allgemein bekannt zu sein. Zur Fehleranalyse siehe WILKINSON [65]; die dort angegebene Schranke (13.2.36) für F bezieht sich auf eine zyklische Realisierung nach den Formeln (13.2.10) und (13.2.11) und ist daher i. allg. zu pessimistisch.

B 13.3. Die einfache Vektoriteration wird seit Jahrzehnten in jedem Lehrbuch der Numerischen Mathematik beschrieben; der Name von-Mises-Iteration geht auf die Arbeit von MISES/ POLLACZEK-GEIRINGER [29] zurück. Die bei der Vektoriteration erzeugte Vektorfolge $\{v, Av, A^2v, \ldots\}$ bildet die Grundlage zur Betrachtung der sog. Krylov-Teilräume \mathcal{K}_p := span $\{v, Av, \ldots, A^{p-1}v\}$, mit deren Hilfe ebenfalls die dominanten Eigenwerte und zugehörigen Eigenvektoren von A approximiert werden können, z. B. nach den verschiedenen Varianten des sog. Lanczos-Algorithmus, siehe PARLETT [80a] für Details. Diese Technik ist allerdings nur für große schwach besetzte Matrizen zweckmäßig und wird deshalb hier nicht untersucht.

B 13.4. Wesentliche Arbeiten zur Teilraumiteration gehen auf BAUER [57] — dort simultane Iteration genannt — und RUTISHAUSER [69, 70/71] zurück, siehe auch STEWART [69, 75]. Auf Grund der bei der Teilraumiteration vorgenommenen Orthogonalisierung der Basis besteht ein enger Zusammenhang zum QR-Algorithmus, siehe 13.7.2. Die zur Teilraumapproximation benötigten Begriffe sind bei DAVIS/KAHAN [70], STEWART [73a] und PARLETT [80a] zu finden. Der Rayleigh-Ritz-Algorithmus zur Beschaffung von Eigenwert/Eigenvektor-Näherungen aus approximierenden Teilräumen ist eine numerische Standardtechnik; die Ursprünge liegen bei RAYLEIGH [1899] und RITZ [1909]. Eine sehr ausgefeilte Prozedur ritzit zur Kombination von Teilraumiteration und RR-Algorithmus, in die noch weitere Verbesserungen aufgenommen worden sind, hat RUTISHAUSER [70] angegeben. Die dabei benutzte Version des RR-Algorithmus nach 13.3.17 (vi) ist von REINSCH vorgeschlagen worden; zur Analyse siehe wieder PARLETT [80a]. Alternative Varianten sind für allgemeinere Aufgabenklassen bei McCORMICK/NOE [77] zu finden.

B 13.5. Spezielle Verfahren der Teilraumiteration beruhen auf einer Maximierung des Funktionals tr (P) = tr ($Q^T A Q$) über alle $Q \in \mathbb{R}^{n,p}$ mit $Q^T Q = I$. Hierzu gehören die Verfahren der simultanen Koordinatenrelaxation — siehe etwa SCHWARZ [77] —, simultane Gradientenverfahren und simultane Verfahren der konjugierten Gradienten — siehe GERADIN [71], LONGSINE/MCCORMICK [80] und DÖHLER [82]. Diese Verfahren sind in der Regel für das allgemeine symmetrische Eigenwertproblem $Ax = \lambda Bx$ entwickelt und formuliert worden, dazu siehe Kapitel 14. Eine zusammenfassende Darstellung gibt MEYER [87].

B 13.6. Die inverse Iteration geht auf WIELANDT [44] zurück. Modifikationen, bei denen neben v auch μ verbessert wird, haben UNGER [50] und viele andere vorgeschlagen; diese Verfahren hängen eng mit dem Newton-Verfahren zusammen, siehe etwa PETERS/WILKINSON [79]. Die sogar kubisch konvergente Rayleigh-Quotienten-Iteration geht nicht etwa auf Lord RAYLEIGH, sondern auf OSTROWSKI [58, 59] zurück. Das monotone Fallen der Residuen bemerkte KAHAN, siehe PARLETT [74, 80a]. Der Beweis der globalen Konvergenz ist von PARLETT/ KAHAN [69] geführt worden, siehe wieder PARLETT [80a]. Zwischen der RQ-Iteration und dem **QR**-Algorithmus mit speziellen Verschiebungen besteht ein enger Zusammenhang, siehe 13.7.9. Während noch vor 25 Jahren angenommen wurde, daß die Fastsingularität von $A - \mu I$ sich nachteilig auf die inverse Iteration auswirkt, hat später WILKINSON in verschiedenen Arbeiten gezeigt, daß sie sich sogar sehr vorteilhaft auf die berechneten Eigenvektornäherungen auswirkt. Eine detaillierte Rundungsfehleranalyse ist bei WILKINSON [77] zu finden.

B 13.7. Die Transformation einer symmetrischen Matrix mittels Givens-Drehungen auf Tridiagonalform wurde erstmals von GIVENS [54] beschrieben. HOUSEHOLDER/BAUER [59] verwenden später für dieselbe Aufgabe Householder-Spiegelungen. Daß Q und T empfindlich von A und der zur Berechnung genutzten Arithmetik abhängen, zeigen instruktive Bei-
spiele bei WILKINSON [65] und PARLETT [80a]. Zur Anwendung impliziter Givens-Drehungen sei auf RATH [82] verwiesen. Eine alternative, besonders für große schwach besetzte Matrizen geeignete Tridiagonalisierungsmethode stellt der sog. Lanczos-Algorithmus dar, siehe Ü 13.5.6, der auf LANCZOS [50] zurückgeht und detailliert von PAIGE [76], PARLETT [80a], GOLUB/ VAN LOAN [83] und CULLUM/WILLOUGHBY [85] untersucht wird.

B 13.8. Das Bisektionsverfahren mit Berechnung von $s(\mu)$ über die Folge { $\tau_0, ..., \tau_n$ } gemäß Ü 13.6.2 wurde erstmals von GIVENS [53, 54] vorgeschlagen. Da die Hauptabschnittsdeterminanten $\tau_k = \tau_k(\mu)$ eine sog. Sturmsche Kette bilden, spricht man auch von der *Methode der Sturmschen Ketten*. Wesentliche Verbesserungen und eine Fehleranalyse gehen auf WILKINson [65] zurück; dort ist auch die Matrix (13.6.5) zu finden. Bereits bei der Implementierung in WILKINSON/REINSCH [71] wird Über- und Unterlauf durch Übergang zu $d_k = \tau_k/\tau_{k-1}$ vermieden, vgl. Ü 13.6.2. Die natürlichere Motivierung ohne Sturmsche Ketten mittels des Trägheitsgesetzes für quadratische Formen kongruenter Matrizen wird von verschiedenen Autoren bevorzugt, z. B. von PARLETT [80a].

B 13.9. Der *QR*-Algorithmus ist das Analogon des auf RUTISHAUSER [58] zurückgehenden sog. *LR*-Algorithmus und wurde unabhängig und etwa gleichzeitig von FRANCIS [61] und KUBLANOVSKAJA [61] gefunden. Ausführungen zur historischen Entwicklung können bei PARLETT [64] nachgelesen werden. Bei der Motivierung der Grundform des *QR*-Algorithmus über die Teilraumiteration folgen wir methodisch der Darstellung bei WATKINS [82]. Der Zusammenhang zwischen dem *QR*-Algorithmus mit den Verschiebungen $\mu_k = a_{nn}^{(k)}$ und der RQ-Iteration wurde von WILKINSON und PARLETT bemerkt, siehe PARLETT [80a]. Die globale Konvergenz des tridiagonalen *QR*-Algorithmus mit den nach ihm benannten Verschiebungen hat WILKINSON [68] bewiesen, siehe auch HOFFMAN/PARLETT [78] und PARLETT [80a], wo auch Aussagen über die in der Regel schnellere als kubische Konvergenz zu finden sind. Alternative Verschiebungen haben z. B. REINSCH/BAUER [71] verwendet. Bei PARLETT [80a] wird auch auf die Vermeidung von Quadratwurzeln eingegangen, siehe dazu ebenfalls RATH [82], wo implizite Givens-Drehungen herangezogen werden. Das Kriterium (13.7.56) für das Nullsetzen von Subdiagonalelementen ist eine leichte Verbesserung eines auf KAHAN zurückgehenden Kriteriums, das bei PARLETT [80a] angegeben und diskutiert wird.

14. Das allgemeine symmetrische Eigenwertproblem

14.1. Grundlegende Eigenschaften

Das allgemeine symmetrische Eigenwertproblem lautet wie folgt: Für gegebene Matrizen $A, B \in S^{n,n}$ sind Zahlen $\lambda_i = \lambda_i [A, B]$, für die

$$A z^{j} = \lambda_{j} B z^{j}$$
 mit $z^{j} \neq o$ (1)

gilt, und gegebenenfalls die Vektoren z^j zu bestimmen. Wie beim speziellen symmetrischen Eigenwertproblem, das sich aus (1) für B = I ergibt, werden die λ_j Eigenwerte, die z^j zugehörige Eigenvektoren und die Paare $\{\lambda_j, z^j\}$ Eigenpaare von $\{A, B\}$ genannt. Die Matrizenschar $A - \lambda B$ heißt das durch $\{A, B\}$ erzeugte Büschel (engl. "pencil", russ. "пучок"); häufig wird das Paar $\{A, B\}$ selbst Büschel genannt. Man beachte, daß $\{\lambda, z\}$ genau dann ein Eigenpaar von $\{A, B\}$ ist, wenn

$$(\boldsymbol{A} - \boldsymbol{\lambda}\boldsymbol{B})\,\boldsymbol{z} = \boldsymbol{o} \tag{2}$$

435

gilt, also

$$p(\lambda) := \det \left(\boldsymbol{A} - \lambda \boldsymbol{B} \right) = 0 \quad \text{und} \quad \boldsymbol{z} \in \mathcal{N}(\boldsymbol{A} - \lambda \boldsymbol{B}), \qquad \boldsymbol{z} \neq \boldsymbol{o}, \tag{3}$$

ist. Die Funktion $p(\lambda)$ ist ein Polynom vom Höchstgrad n, das charakteristische Polynom von $\{A, B\}$. Im Unterschied zum speziellen Eigenwertproblem kann jedoch der Grad von $p(\lambda)$ kleiner als n sein. In diesem Fall hat das aus (1) durch die Substitution $\lambda = 1/\omega$ entstehende Eigenwertproblem

$$\omega A \boldsymbol{z} = \boldsymbol{B} \boldsymbol{z} \tag{4}$$

den Eigenwert $\omega = 0$, siehe Ü 14.1.1. Wir sagen dann, daß (1) den Eigenwert

 $= 1/0 = \infty$ hat. Es kann insbesondere eintreten, daß $p(\lambda)$ identisch verschwindet, d. h., daß jedes λ Eigenwert ist. Außerdem können auch komplexe Eigenwerte auftreten. In Ü 14.1.2 sieht man, daß die beschriebenen ungewöhnlichen Fälle bereits für n = 2 möglich sind.

Ein wesentliches theoretisches wie praktisches Hilfsmittel bei der Behandlung des speziellen symmetrischen Eigenwertproblems $Ax = \lambda x$ war der Übergang von $A - \lambda I$ zu

$$\hat{A} - \lambda \hat{I} = E(A - \lambda I) F$$
(5)

mittels regulärer Transformationsmatrizen E, F. Damit das transformierte Problem wieder ein spezielles Eigenwertproblem wird, muß $\hat{I} = EIF = I$, also $E = F^{-1}$ gefordert werden, und \hat{A} ergibt sich dann zu $\hat{A} = F^{-1}AF$, d. h., \hat{A} entsteht durch Ähnlichkeitstransformation mit F aus A. Damit \hat{A} mit A symmetrisch ist, haben wir $F^{-1} = F^{T}$, also die Orthogonalität der Transformationsmatrix F gefordert, vgl. 1.2.B und 13.1.A.

Für das allgemeine Eigenwertproblem $Az = \lambda Bz$ ist die Situation etwas anders. Die zu (5) analoge Transformation führt auf

$$\hat{A} - \lambda \hat{B} = E(A - \lambda B) F$$
 mit $\hat{A} = EAF$, $\hat{B} = EBF$. (6)

Damit die Klasse der allgemeinen Eigenwertprobleme beim Übergang von $\{A, B\}$ zu $\{\hat{A}, \hat{B}\}$ nicht verlassen wird, braucht außer der Regularität keine weitere Forderung an E und F gestellt zu werden. Büschel, die durch die simultanen Äquivalenztransformationen (6) auseinander hervorgehen, heißen *äquivalent*. Äquivalente Büschel haben dieselben Eigenwerte, und \hat{z} ist Eigenvektor von $\{\hat{A}, \hat{B}\}$ genau dann, wenn z mit

$$\boldsymbol{z} = \boldsymbol{F} \hat{\boldsymbol{z}} \tag{7}$$

Eigenvektor von $\{A, B\}$ zum selben Eigenwert ist. Man beachte jedoch, daß bei der Transformation (6) die Symmetrie von $\{A, B\}$ i. allg. verlorengeht. Damit sich die Symmetrie von $\{A, B\}$ auf $\{\hat{A}, \hat{B}\}$ überträgt, muß $E = F^{\intercal}$ gefordert werden, d. h., die Äquivalenztransformationen (6) müssen auf die symmetrieerhaltenden simultanen Kongruenztransformationen

$$\hat{A} - \lambda \hat{B} = F^{\mathsf{T}}(A - \lambda B) F$$
 mit $\hat{A} = F^{\mathsf{T}}AF$, $\hat{B} = F^{\mathsf{T}}BF$ (8)

eingeschränkt werden; die Relation (7) zwischen den Eigenvektoren ändert sich dabei nicht. Büschel, die gemäß (8) auseinander hervorgehen, heißen kongruent. Der Übergang zu geeigneten einfacheren kongruenten Büscheln, bei dem die Symmetrie erhalten bleibt, ist die für das allgemeine symmetrische Eigenwertproblem angemessene Transformation. Wir weisen darauf hin, daß in (8) nicht die Orthogonalität von \mathbf{F} gefordert wird.

Es zeigt sich, daß die eingangs beschriebenen pathologischen Situationen bei einem symmetrischen Büschel nicht auftreten können, wenn B zusätzlich positiv definit ist.

14.1.1. Satz. Gegeben sei das symmetrische Büschel $\{A, B\}$, und **B** sei positiv definit. Dann gibt es eine Matrix **Z** mit

$$\boldsymbol{Z}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{Z} = \boldsymbol{\Lambda} = \operatorname{diag}\left(\lambda_{i}\right), \quad \lambda_{i} \in \mathbf{R}, \tag{9}$$

und

$$\mathbf{Z}^{\mathsf{T}}\mathbf{B}\mathbf{Z}=\mathbf{I},\tag{10}$$

d. h., $\{A, B\}$ ist kongruent zum diagonalen Büschel $\{A, I\}$, und die Transformationsmatrix ist Z.

Beweis. Wir konstruieren die gesuchte Transformation in zwei Schritten. Es sei zunächst

$$Q^{\mathsf{T}}BQ = \varDelta = \operatorname{diag}(\delta_i), \quad \delta_1 \ge \cdots \ge \delta_n > 0, \quad Q \text{ orthogonal},$$
 (11)

die Eigenwertzerlegung von B; wegen der positiven Definitheit von B sind alle Eigenwerte δ_j positiv. Mit

gilt dann

$$\boldsymbol{F}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{F} = \varDelta^{-1/2}\boldsymbol{Q}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{Q}\varDelta^{-1/2} =: C \quad \text{und} \quad \boldsymbol{F}^{\mathsf{T}}\boldsymbol{B}\boldsymbol{F} = \varDelta^{-1/2}\boldsymbol{Q}^{\mathsf{T}}(\boldsymbol{Q}\varDelta\boldsymbol{Q}^{\mathsf{T}}) \boldsymbol{Q}\varDelta^{-1/2} = \boldsymbol{I},$$

d. h., $\{A, B\}$ ist kongruent zu $\{C, I\}$ mit F, wobei $\varDelta^{-1/2} := (\varDelta^{1/2})^{-1}$ gesetzt ist. Die Matrix C ist symmetrisch und besitzt ihrerseits die Eigenwertzerlegung

$$U^{\mathsf{T}}CU = \Lambda = \operatorname{diag}(\lambda_i), \quad \lambda_i \in \mathsf{R}, \quad U \text{ orthogonal}, \quad \mathrm{d. h.} \quad U^{\mathsf{T}}IU = I.$$
 (12)

Die Gleichungen (12) besagen, daß $\{C, I\}$ mit U kongruent zu $\{.1, I\}$ ist, d. h., $\{A, B\}$ ist mit dem Produkt Z := FU der einzelnen Transformationsmatrizen kongruent zu $\{.1, I\}$. \Box

Die $\{\lambda_j, e^j\}$ sind offensichtlich die Eigenpaare von $\{A, I\}$, so daß die $\{\lambda_j, z^j\}$, wobei die $z^j = Ze^j$ die Spalten von $Z = (z^1, ..., z^n)$ bezeichnen, die Eigenpaare von $\{A, B\}$ sind; man lese (7) mit Z statt F oder schreibe (9) unter Beachtung von (10) in der Form AZ = BZA, was zu (1) äquivalent ist. Die Eigenschaft (10) von Z wird B-Orthogonalität genannt; in elementweiser Lesart bedeutet sie

$$(\boldsymbol{z}^{i})^{\mathsf{T}} \boldsymbol{B} \boldsymbol{z}^{j} = \begin{cases} 1 & \text{für } i = j, \\ 0 & \text{für } i \neq j \end{cases} \quad (i, j = 1, ..., n).$$
(13)

Die Eigenvektoren $\{z^{j}\}$ heißen deswegen auch **B**-orthonormal, vgl. Ü 14.1.3. Offenbar sind **B**-orthonormale Vektoren stets linear unabhängig. 14.1.1 bedeutet daher: Jedes symmetrische Büschel $\{A, B\}$ mit positiv definitem **B** hat n reelle Eigenwerte $\{\lambda_j\}$, und es gibt n zugehörige **B**-orthonormale Eigenvektoren $\{z^i\}$. Wie im Fall des speziellen symmetrischen Eigenwertproblems kann diese Aussage noch verschärft werden: Zu verschiedenen Eigenwerten $\lambda_i \neq \lambda_j$ gehörende Eigenvektoren z^i und z^j sind notwendig **B**-orthogonal im Sinne von $[z^i, z^j]_B = (z^i)^{\mathsf{T}} B z^j = 0$. Falls λ eine α -fache Nullstelle von $p(\lambda)$ ist, gibt es α linear unabhängige Eigenvektoren, die den zugehörigen Eigenraum $\mathcal{N}(A - \lambda B)$ aufspannen, und diese Eigenvektoren können **B**-orthonormal gewählt werden. Im folgenden setzen wir stets voraus, daß $\{A, B\}$ symmetrisch und **B** positiv definit ist.

Wir bemerken an dieser Stelle, daß auch für das allgemeine symmetrische Eigenwertproblem eine Reihe von nützlichen Residualkriterien und Störungsaussagen existiert, siehe B 14.1. Aus Platzgründen verzichten wir jedoch auf eine Wiedergabe.

Übungsaufgaben

Ü 14.1.1. Man überlege sich, daß $p(\lambda) = \det (A - \lambda B)$ und $q(\omega) = \det (\omega A - B)$ den Beziehungen

$$q(\omega) = p(1/\omega) \; \omega^{n} \quad ext{und} \quad p(\lambda) = q(1/\lambda) \; \lambda^{n}$$

genügen, und folgere hieraus, daß $\omega = 0$ Eigenwert von $\{B, A\}$, also $\lambda = 1/0 = \infty$ Eigenwert von $\{A, B\}$ genau dann ist, wenn der Grad von $p(\lambda)$ kleiner als n ist.

Ü14.1.2. Man zeige (PARLETT [80a]):

(i) Für
$$\mathbf{A} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$$
, $\mathbf{B} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$ sind $\{1, e^1\}$ und $\{\lambda, e^2\}$, λ beliebig, Eigenpaare.

(ii) Für
$$\mathbf{A} = \begin{pmatrix} \mathbf{1} & 0 \\ 0 & 0 \end{pmatrix}, \mathbf{B} = \begin{pmatrix} 0 & 0 \\ 0 & \mathbf{1} \end{pmatrix}$$
 sind $\{\infty, \mathbf{e}^1\}$ und $\{0, \mathbf{e}^2\}$ Eigenpaare.

(iii) Für $\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$, $\mathbf{B} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}$ sind {i, (i, -1)^T} und {-i, (i, 1)^T} Eigenpaare, wobei i mit i² = -1 die imaginäre Einheit bedeutet.

Ü 14.1.3. Mit der positiv definiten Matrix $B \in S^{n,n}$ werde

$$[\boldsymbol{x}, \boldsymbol{y}] := [\boldsymbol{x}, \boldsymbol{y}]_{\boldsymbol{B}} := \boldsymbol{x}^{\mathsf{T}} \boldsymbol{B} \boldsymbol{y} \tag{14}$$

definiert. Man zeige, daß die derart erklärte Abbildung $(x, y) \rightarrow [x, y]_B$ von $\mathbb{R}^n \times \mathbb{R}^n$ in \mathbb{R} die Eigenschaften

 $(S1) [\boldsymbol{x}, \boldsymbol{y}] = [\boldsymbol{y}, \boldsymbol{x}],$

 $(\mathbf{S2}) \quad [\lambda_1 \boldsymbol{x}^1 + \lambda_2 \boldsymbol{x}^2, \boldsymbol{y}] = \lambda_1 [\boldsymbol{x}^1, \boldsymbol{y}] + \lambda_2 [\boldsymbol{x}^2, \boldsymbol{y}].$

(S3) $[x, x] \ge 0$, [x, x] = 0 genau für x = ohat, d. h., [x, y] hat dieselben Eigenschaften wie das gewöhnliche Skalarprodukt $[x, y]_I = x^{\mathsf{T}}y$ und wird daher das *durch* **B** erzeugte Skalarprodukt oder kurz **B**-Skalarprodukt genannt. Durch die Vorschrift

$$\|\boldsymbol{x}\|_{\boldsymbol{B}} := \sqrt[]{[\boldsymbol{x}, \boldsymbol{x}]_{\boldsymbol{B}}} = \sqrt[]{\boldsymbol{x}^{\mathsf{T}} \boldsymbol{B} \boldsymbol{x}}$$
(15)

wird auf \mathbb{R}^n eine Norm definiert, die sog. durch **B** erzeugte (energetische) Norm oder kurz **B**-Norm. Für **B**-Skalarprodukt und **B**-Norm gilt die (verallgemeinerte) Schwarzsche Ungleichung

$$|[\boldsymbol{x},\boldsymbol{y}]_{\boldsymbol{B}}| \leq ||\boldsymbol{x}||_{\boldsymbol{B}} ||\boldsymbol{y}||_{\boldsymbol{B}}.$$
(16)

Ü 14.1.4. Betrachtet wird das quadratische Eigenwertproblem

$$(\boldsymbol{A} - \boldsymbol{\lambda}\boldsymbol{B} - \boldsymbol{\lambda}^2 \mathbf{C})\,\boldsymbol{x} = 0 \tag{17}$$

mit $A, B, C \in S^{n,n}$ und positiv definiten Matrizen A, C. Man zeige: Wenn $\{\lambda, x\}$ ein Eigenpaar von (17) ist, ist $\{\lambda, z\}$ mit $z := \left(\frac{x}{\lambda x}\right)$ ein Eigenpaar von

$$Nz = \lambda Mz \tag{18}$$

mit

$$N := \left(\frac{O \mid A}{A \mid -B}\right), \qquad M := \left(\frac{A \mid O}{O \mid C}\right) \in \mathbf{S}^{2n, 2n}.$$
(19)

Ist umgekehrt $\{\lambda, z\}$ mit $z = \left(\frac{x}{y}\right)$ ein Eigenpaar von (18), so gilt $y = \lambda x$, und $\{\lambda, x\}$ ist ein Eigenpaar von (17). Das quadratische symmetrische Eigenwertproblem (17) kann also auf das lineare allgemeine symmetrische Eigenwertproblem (18) mit positiv definitem M zurückgeführt werden, wobei sich die Dimension verdoppelt.

14.2. Explizite Reduktion auf ein spezielles Eigenwertproblem

Wir betrachten in diesem Abschnitt Methoden zur Lösung des allgemeinen symmetrischen Eigenwertproblems

$$A\boldsymbol{z} = \lambda \boldsymbol{B}\boldsymbol{z},\tag{1}$$

die auf der Transformation von $\{A, B\}$ in ein äquivalentes bzw. kongruentes Büschel $\{\hat{A}, \hat{B}\} = \{C, I\}$ beruhen, d. h., die (1) auf das spezielle Eigenwertproblem

$$\boldsymbol{C}\boldsymbol{u} = \lambda \boldsymbol{u} \tag{2}$$

zurückführen. Werden äquivalente Büschel mit $\hat{A} = EAF$, $\hat{B} = EBF$ zugelassen, so muß also

$$\hat{B} = EBF = I \tag{3}$$

gelten, was am einfachsten durch die Festlegung

$$E := E_1 := B^{-1}, \quad F := F_1 := I, \text{ also } C := C_1 := B^{-1}A$$
 (4)

erreicht werden kann. Dies entspricht dem Übergang von (1) zu

$$\boldsymbol{B}^{-1}\boldsymbol{A}\boldsymbol{z}=\boldsymbol{\lambda}\boldsymbol{z}.$$

Da die entstehende Matrix $C_1 = B^{-1}A$ nur in trivialen Ausnahmefällen symmetrisch ist, kann dieses naheliegende Vorgehen jedoch nicht empfohlen werden. Wir beschränken uns daher auf kongruente Büschel, bei denen die Symmetrie erhalten bleibt. Die Bedingung (3) geht dabei in

$$\boldsymbol{F}^{\mathsf{T}}\boldsymbol{B}\boldsymbol{F} = \boldsymbol{I}, \quad \text{also} \quad \boldsymbol{B} = \boldsymbol{F}^{-\mathsf{T}}\boldsymbol{F}^{-1} \tag{6}$$

über, d. h., es wird eine symmetrische Faktorisierung von B benötigt. Eine Möglichkeit zur Festlegung von F haben wir bereits im ersten Schritt des Beweises zu 14.1.1 kennengelernt:

14.2.1. Reduktion mittels der Eigenwertzerlegung von B. Es sei

 $B = Q \varDelta Q^{\intercal}$ mit $\varDelta = \text{diag}(\delta_j), \quad \delta_1 \ge \cdots \ge \delta_n > 0, \quad Q \text{ orthogonal}, \quad (7)$ die Eigenwertzerlegung von B, und mit $\varDelta^{1/2} := \text{diag}(\sqrt[]{\delta_j})$ werde

$$\boldsymbol{F}_{2} := \boldsymbol{O} \boldsymbol{\varDelta}^{-1/2} \tag{8}$$

gesetzt. Dann gilt

$$\boldsymbol{F}_{2}^{\mathsf{T}}(\boldsymbol{A}-\boldsymbol{\lambda}\boldsymbol{B})\,\boldsymbol{F}_{2}=\boldsymbol{C}_{2}-\boldsymbol{\lambda}\boldsymbol{I} \tag{9}$$

mit

$$C_2 := F_2^{\mathsf{T}} A F_2 = A^{-1/2} Q^{\mathsf{T}} A Q A^{-1/2}.$$
(10)

14.2.2. Bemerkung. (i) Für voll besetztes **B** sollte der **QR**-Algorithmus — d. h. die Kombination von 13.5.1 und 13.7.17 — zur Berechnung der Eigenwertzerlegung verwendet werden. Man beachte, daß bei der Reduktion die Matrix **Q** als Produkt der Householder-Spiegelungen aus 13.5.1 und der Givens-Drehungen aus 13.7.17 überhaupt nicht explizit gebildet zu werden braucht: Von $\{A, B\}$ wird direkt zu $\{\tilde{A}, \tilde{B}\}$ mit $\tilde{A} := \mathbf{Q}^{\mathsf{T}} A \mathbf{Q}$ und $\tilde{B} := \mathbf{Q}^{\mathsf{T}} B \mathbf{Q} = \Delta$ übergegangen, indem die einzelnen Faktoren von **Q** nicht nur auf **B**, sondern sofort bei ihrem Auftreten auch auf **A** angewendet werden. Das kostet insgesamt $\sim Kn^3$ opms mit $K \leq 4$, vgl. 13.7.18(ii), denn der Aufwand zur Berechnung von \tilde{A} ist derselbe wie zur expliziten Berechnung von **Q**, und die Berechnung von A aus der Tridiagonalform kostet nur $O(n^2)$ opms, also relativ fast nichts. Zum Schluß wird $C_2 := A^{-1/2} \tilde{A} A^{-1/2}$ mit $\sim n^2$ opm gebildet. Für die so berechnete Matrix C_2 gilt dann

$$C_2 = A^{-1/2} Q^{\intercal} (A + \delta A) Q A^{-1/2}$$
 mit kleinem $\delta A \in S^{n,n}$,

und C_2 enthält fast die volle Information über das Originalproblem. Falls die Eigenvektoren auch gesucht sind, muß Q für die Rücktransformation $z = F_2 u = Q \Lambda^{-1/2} u$ allerdings bereitgestellt werden, was den Aufwand verdoppelt.

(ii) Im Fall $\delta_1 \gg \delta_n$ ist C_2 ansteigend gestuft. Es sind dann die in 13.5.4 und 13.7.18 beschriebenen Vorsichtsmaßnahmen durchzuführen, damit eventuell auch die betragskleinen Eigenwerte ausreichend genau berechnet werden. Für die Behandlung des semidefiniten Falles $\delta_n = 0$ siehe B 14.2.

(iii) Falls A, B Bandmatrizen sind, kann das Büschel $\{A, B\}$ in Analogie zu 13.5.B ohne Vergrößerung der Bandbreite kongruent in ein Büschel $\{T, I\}$ mit tridiagonalem T transformiert werden, siehe B 14.3.

Eine billigere Reduktion als in 14.2.1 ist möglich, wenn (6) mittels der LDL^{T} -Faktorisierung von **B** realisiert wird. Zur Verbesserung der Stabilität verwenden wir dazu die mit Diagonalpivotisierung gemäß 6.1.5(iii) berechnete Faktorisierung, bei der die Diagonalelemente monoton fallen.

14.2.3. Reduktion mittels LDL^T-Faktorisierung von B. Es sei

$$\boldsymbol{P}\boldsymbol{B}\boldsymbol{P}^{\mathsf{T}} = \boldsymbol{L}\boldsymbol{D}\boldsymbol{L}^{\mathsf{T}} \quad \text{mit} \quad \boldsymbol{D} = \text{diag}\,(d_{i}), \qquad d_{1} \geq \cdots \geq d_{n} > 0, \tag{11}$$

einer Permutationsmatrix P und einer unteren Einsdreiecksmatrix L die mittels Diagonalpivotisierung gemäß 6.1.1/6.1.5 berechnete LDL^{\intercal} -Faktorisierung von B, und mit $D^{1/2} := \text{diag}\left(\sqrt{d_j}\right)$ werde

$$\boldsymbol{F}_{3} := \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L}^{-\mathsf{T}} \boldsymbol{D}^{-1/2} \tag{12}$$

gesetzt. Dann gilt

$$\boldsymbol{F}_{3}^{\mathsf{T}}(\boldsymbol{A}-\boldsymbol{\lambda}\boldsymbol{B})\,\boldsymbol{F}_{3}=\boldsymbol{C}_{3}-\boldsymbol{\lambda}\boldsymbol{I} \tag{13}$$

mit

$$\boldsymbol{C}_3 := \boldsymbol{F}_3^{\mathsf{T}} \boldsymbol{A} \boldsymbol{F}_3 = \boldsymbol{D}^{-1/2} \boldsymbol{L}^{-1} \boldsymbol{P} \boldsymbol{A} \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L}^{-\mathsf{T}} \boldsymbol{D}^{-1/2}. \tag{14}$$

Beweis. Der Vergleich von $B = P^{\mathsf{T}}LDL^{\mathsf{T}}P$ mit (6) zeigt, daß $F^{-1} = D^{1/2}L^{\mathsf{T}}P$ gesetzt werden muß. \Box

14.2.4. Bemerkung. (i) Bei der Realisierung von 14.2.3 sollte die LDL^{T} -Faktorisierung in der (n - 1)-stufigen Version analog zu 6.1.B durchgeführt werden, was der Produktdarstellung

$$\tilde{\boldsymbol{B}} := \boldsymbol{L}^{-1} \boldsymbol{P} \boldsymbol{B} \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L}^{-\mathsf{T}} = \boldsymbol{D} = \boldsymbol{L}_{n-1} \cdots \boldsymbol{L}_1 \boldsymbol{T}_{n-1} \cdots \boldsymbol{T}_1 \boldsymbol{B} \boldsymbol{T}_1 \cdots \boldsymbol{T}_{n-1} \boldsymbol{L}_1^{\mathsf{T}} \cdots \boldsymbol{L}_{n-1}^{\mathsf{T}} \quad (15)$$

mit

$$\boldsymbol{L}^{-\mathsf{T}} = \boldsymbol{L}_1^{\mathsf{T}} \cdots \boldsymbol{L}_{n-1}^{\mathsf{T}} \tag{16}$$

entspricht, wobei die Matrizen $L_k = L_k(-l^k)$ die zur Elimination verwendeten nichtorthogonalen elementaren Transformationsmatrizen bezeichnen. Die Matrix L^{-1} wird nicht explizit gebildet, sondern $\tilde{A} := L^{-1}PAP^{\mathsf{T}}L^{-\mathsf{T}}$ wird analog zu (15) aus Aerzeugt, indem dieselben Transformationen wie mit B auch mit A ausgeführt werden. In dieser Form kostet die Transformation von $\{A, B\}$ in $\{\tilde{A}, \tilde{B}\}$ insgesamt $\sim 2n^3/3$ opms, nämlich $\sim n^3/2$ opms für \tilde{A} und $\sim n^3/6$ für $\tilde{B} = D$, ist also um den Faktor 6 billiger als die entsprechende Transformation in 14.2.1. Aus \tilde{A} wird $C_3 = D^{-1/2}\tilde{A}D^{-1/2}$ mit $\sim n^2$ opm berechnet. Es gilt dann

$$C_3 = \boldsymbol{D}^{-1/2} \boldsymbol{L}^{-1} \boldsymbol{P} (A + \boldsymbol{\sigma} A) \boldsymbol{P}^{\mathsf{T}} \boldsymbol{L}^{-\mathsf{T}} \boldsymbol{D}^{-1/2}$$

aber die Störung $dA \in S^{n,n}$ ist a priori nur gemäß

$$\| \delta A \|_F \leq \nu F \text{ [cond (L)]}^2 \| A \|$$

mit akzeptablem F beschränkt. Auf Grund der Diagonalpivotisierung wird cond (L)meist nicht sehr groß sein, so daß δA klein wird und C_3 ausreichende Information über $\{A, B\}$ enthält. Falls cond (L) jedoch groß ist — dies wird selten eintreten, theoretisch kann cond (L) aber exponentiell mit n wachsen, vgl. Ü 11.2.2 —, entsteht beim Übergang zu C_3 ein großer Informationsverlust. Dann wird i. allg. auch \tilde{A} groß sein, so daß $\|\tilde{A}\|/\|A\|$ als Kriterium für die Zulässigkeit genommen werden kann. Ist dieser Wert groß, so sollte besser die teurere, aber stabilere Reduktion 14.2.1 verwendet werden. Wegen (16) gilt außerdem $z = F_3 u = P^{\intercal} L_1^{\intercal} \cdots L_{n-1}^{\intercal} D^{-1/2} u$, d. h., L^{-1} wird auch bei der Rücktransformation der Eigenvektoren von C_3 in solche von $\{A, B\}$ nicht explizit benötigt. Diese Transformation kostet $\sim n^2/2$ opms pro Vektor, ist also viel billiger als in 14.2.1. (ii) Wenn auf die Pivotisierung verzichtet wird, vergrößert sich das Risiko, ein großes cond (L) zu erhalten, wesentlich.

(iii) Im Fall $d_1 \gg d_n$ ist C_3 ansteigend gestuft, so daß wie in 14.2.2(ii) vorzugehen ist. \Box

14.2.5. Bemerkung. (i) Die explizite Reduktion nach 14.2.1 oder 14.2.3 ist nur für kleine bzw. mittlere Dimensionen n zweckmäßig; eine eventuelle schwache Besetztheit von A, B wird beim Übergang zu C zerstört.

(ii) Wenn mehr als n/4 Eigenwerte bzw. Eigenpaare gesucht sind, sollten alle Eigenwerte bzw. Eigenpaare von C nach dem QR-Algorithmus bestimmt werden. Interessieren nur wenige Eigenwerte bzw. Eigenpaare, so empfiehlt sich das Bisektionsverfahren aus 13.6 bzw. die inverse Iteration aus 13.4. In beiden Fällen müssen die Eigenvektoren u von C gemäß z = Fu in solche von $\{A, B\}$ transformiert werden.

(iii) Falls A, B Bandgestalt haben, kann das Bisektionsverfahren direkt mit $A - \mu B$ durchgeführt werden: Existiert die Faktorisierung $A - \mu B = L_{\mu}D_{\mu}L_{\mu}^{\mathsf{T}}$, so ist die Zahl der negativen Diagonalelemente von D_{μ} gleich der Zahl $s(\mu)$ der Eigenwerte von $\{A, B\}$, die kleiner als μ sind, vgl. Ü 14.2.1. Für weiterführende Hinweise zur Behandlung von Bandmatrizen siehe B 14.3, vgl. auch 13.6.11 (vi).

Übungsaufgabe

Ü 14.2.1. Man beweise die Aussage von 14.2.5 (iii). Hinweis: Man zeige mit 14.1.1, daß $A - \mu I$ und D_{μ} kongruent sind, und beachte 1.2.13.

14.3. Vektor- und Teilraumiteration

Wenn die Dimension n groß ist, sind die expliziten Reduktionsverfahren aus dem vorangegangenen Abschnitt aus Aufwandsgründen nicht mehr anwendbar. Meist sind dann jedoch nur einige Eigenwerte bzw. Eigenpaare von Interesse — etwa die pbetragsgrößten oder betragskleinsten —, so daß sich die Verfahren der direkten bzw. inversen Vektor- und Teilraumsteration — siehe 13.3 bzw. 13.4 — als Alternative anbieten. Da wir diese Verfahren für das spezielle symmetrische Eigenwertproblem formuliert haben, setzen wir zunächst voraus, daß das Büschel $\{A, B\}$ bereits mittels einer regulären Matrix F kongruent in das Büschel $\{C, I\}$ transformiert worden ist, d. h., es gelte

$$C = F^{\mathsf{T}} A F \quad \text{und} \quad I = F^{\mathsf{T}} B F, \quad \text{also} \quad B = F^{-\mathsf{T}} F^{-1}, \tag{1}$$

mithin

$$\boldsymbol{C} - \lambda \boldsymbol{I} = \boldsymbol{F}^{\mathsf{T}} (\boldsymbol{A} - \lambda \boldsymbol{B}) \, \boldsymbol{F}. \tag{2}$$

Die Vektoriteration 13.3.1 für C lautet dann

$$\hat{\boldsymbol{v}}^{k+1} := C \hat{\boldsymbol{v}}^k / \hat{\boldsymbol{\omega}}_{k+1} \qquad (k = 1, 2, ...),$$
(3)

wobei $\|\hat{v}^1\| = 1$ und $\hat{\omega}_{k+1}$ so gewählt wird, daß $\|\hat{v}^{k+1}\| = 1$ gilt. Falls λ_1 ein dominanter Eigenwert von C bzw. $\{A, B\}$ ist, konvergiert $\{\pm \hat{v}^k\}$ gegen den zu λ_1 gehörenden

Eigenvektor u^1 von C, vgl. 13.3.A. Da die Eigenvektoren u von C und z von $\{A, B\}$ durch

$$\boldsymbol{z} = \boldsymbol{F}\boldsymbol{u} \tag{4}$$

verknüpft sind, siehe (2), konvergieren die i. allg. nicht normierten transformierten Iterierten

$$\boldsymbol{v}^{\boldsymbol{k}} := \boldsymbol{F} \hat{\boldsymbol{v}}^{\boldsymbol{k}} \tag{5}$$

bei geeigneter Vorzeichenfestlegung gegen den u^1 entsprechenden Eigenvektor z^1 von $\{A, B\}$. Aus (3), (5) und (1) folgt nun

$$\boldsymbol{v}^{k+1} = \boldsymbol{F}\hat{\boldsymbol{v}}^{k+1} = \boldsymbol{F}C\hat{\boldsymbol{v}}^{k}/\hat{\boldsymbol{\omega}}_{k+1} = (\boldsymbol{F}\boldsymbol{F}^{\mathsf{T}}) \boldsymbol{A}(\boldsymbol{F}\hat{\boldsymbol{v}}^{k})/\hat{\boldsymbol{\omega}}_{k+1} = \boldsymbol{B}^{-1}\boldsymbol{A}\boldsymbol{v}^{k}/\hat{\boldsymbol{\omega}}_{k+1}, \quad (6)$$

d. h., die eigentlich interessierenden, rücktransformierten Iterierten $\{v^k\}$ sind — abgesehen von der für die Approximation unwesentlichen Normierung — gerade die Iterierten der Vektoriteration für die i. allg. nichtsymmetrische Matrix $B^{-1}A$, und die Konvergenzaussagen für $\{\hat{v}^k\}$ übertragen sich sinngemäß auf $\{v^k\}$. Die oben vorausgesetzte Reduktion auf $\{C, I\}$ braucht also überhaupt nicht ausgeführt zu werden; das Verfahren läßt sich allein mit den Originaldaten $\{A, B\}$ formulieren. Zur Berechnung von $w^{k+1} := B^{-1}Av^k$ braucht $B^{-1}A$ — die Matrix C_1 aus 14.2 — auch nicht explizit gebildet zu werden, sondern w^{k+1} wird als Lösung des Gleichungssystems $Bw^{k+1} = Av^k$ etwa unter Verwendung einer symmetrischen Faktorisierung von Bbestimmt. Bei Normierung von v^k auf $||v^k||_B = 1$ — siehe Ü 14.1.3 — ergibt sich damit die folgende implizite Realisierung der Vektoriteration für C:

14.3.1. Implizite Vektoriteration.

- S0: Berechne Faktorisierung $B = LDL^{\intercal}$, wähle v^1 mit $||v^1||_B = 1$, setze k := 1
- S1: Berechne w^{k+1} unter Verwendung der LDL^{\intercal} -Faktorisierung von B aus dem Gleichungssystem

$$\boldsymbol{B}\boldsymbol{w}^{k+1} = \boldsymbol{A}\boldsymbol{v}^k \tag{7}$$

S2: Bestimme

$$\omega_{k+1} := \|\boldsymbol{w}^{k+1}\|_{\boldsymbol{B}} = \sqrt{(\boldsymbol{w}^{k+1})^{\mathsf{T}} \boldsymbol{B} \boldsymbol{w}^{k+1}} = \sqrt{(\boldsymbol{w}^{k+1})^{\mathsf{T}} (\boldsymbol{A} \boldsymbol{v}^{k})},$$
(8)

setze

$$oldsymbol{v}^{k+1}:=oldsymbol{w}^{k+1}/\omega_{k+1}$$

S3: Setze k := k + 1, goto S1

Aus den v^k können analog zu 13.3.A Näherungen für λ_1 ermittelt werden, insbesondere kann dazu der Rayleigh-Quotient benutzt werden. Der Vektor $\hat{v}^k = F^{-1}v^k$ ist eine i. allg. nicht normierte Eigenvektorapproximation für C, und der zugehörige Rayleigh-Quotient ergibt sich wegen (1) zu

$$\varrho_k = \varrho_C(\hat{\boldsymbol{v}}^k) = \frac{(\hat{\boldsymbol{v}}^k)^{\mathsf{T}} \, C\hat{\boldsymbol{v}}^k}{(\hat{\boldsymbol{v}}^k)^{\mathsf{T}} \, \hat{\boldsymbol{v}}^k} = \frac{(\boldsymbol{v}^k)^{\mathsf{T}} \, \boldsymbol{F}^{-\mathsf{T}}(\boldsymbol{F}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{F}) \, \boldsymbol{F}^{-1}\boldsymbol{v}^k}{(\boldsymbol{v}^k)^{\mathsf{T}} \, \boldsymbol{F}^{-\mathsf{T}}\boldsymbol{F}^{-1}\boldsymbol{v}^k} = \frac{(\boldsymbol{v}^k)^{\mathsf{T}} \, \boldsymbol{A}\boldsymbol{v}^k}{(\boldsymbol{v}^k)^{\mathsf{T}} \, \boldsymbol{B}\boldsymbol{v}^k},\tag{9}$$

vgl. (13.1.44).

In den Anwendungen ist meist nicht der betragsgrößte, sondern der betragskleinste oder der am dichtesten an einer vorgegebenen Zahl μ gelegene Eigenwert gesucht. Dann muß statt der Vektoriteration die inverse Iteration aus 13.4 verwendet werden, d. h., statt C wird die Matrix

$$\hat{C} := (C - \mu I)^{-1} \quad \text{mit den Eigenwerten } \hat{\lambda}_j = 1/(\lambda_j - \mu) \tag{10}$$

zur Iteration verwendet. Wir numerieren die Eigenwerte λ_j von C bzw. $\{A, B\}$ wie in 13.4 gemäß

$$0 \leq |\lambda_1 - \mu| \leq |\lambda_2 - \mu| \leq \dots \leq |\lambda_n - \mu|$$
(11)

und setzen $0 < |\lambda_1 - \mu| < |\lambda_2 - \mu|$ voraus, d. h., \hat{C} existiere, und $\hat{\lambda}_1$ sei ein dominanter Eigenwert von \hat{C} . Die Iterierten $\{\pm \hat{v}^k\}$ der inversen Iteration

$$\hat{\boldsymbol{v}}^{k+1} := \hat{\boldsymbol{C}}\hat{\boldsymbol{v}}^{k}/\hat{\omega}_{k+1} = (\boldsymbol{C} - \boldsymbol{\mu}\boldsymbol{I})^{-1}\,\hat{\boldsymbol{v}}^{k}/\hat{\omega}_{k+1}$$
(12)

konvergieren dann gegen den zu λ_1 gehörenden Eigenvektor u^1 von C. Nun folgt aus (2) mit μ statt λ unter Beachtung von (1)

$$(C - \mu I)^{-1} = F^{-1}(A - \mu B)^{-1} BF$$

so daß (12) mit den rücktransformierten Größen \boldsymbol{v}^k in der Form

$$\boldsymbol{v}^{k+1} = \boldsymbol{F} \hat{\boldsymbol{v}}^{k+1} = (\boldsymbol{A} - \boldsymbol{\mu} \boldsymbol{B})^{-1} \, \boldsymbol{B} \boldsymbol{v}^k / \hat{\boldsymbol{\omega}}_{k+1} \tag{13}$$

geschrieben werden kann. Man beachte, daß die i. allg. nichtsymmetrische Matrix $(\mathbf{A} - \mu \mathbf{B})^{-1} \mathbf{B}$ ähnlich zu $(\mathbf{C} - \mu \mathbf{I})^{-1}$ ist und folglich dieselben Eigenwerte $1/(\lambda_j - \mu)$ wie die letztere hat. Wie die direkte Vektoriteration kann daher auch die inverse Iteration allein mit den Originaldaten $\{\mathbf{A}, \mathbf{B}\}$ durchgeführt werden:

14.3.2. Implizite inverse Iteration.

- S0: Wähle $\mu \in \mathbf{R}$, berechne numerisch gutartige Faktorisierung von $A \mu B$, wähle v^1 mit $||v^1||_B = 1$, setze k := 1
- S1: Berechne w^{k+1} unter Verwendung der in S0 berechneten Faktorisierung aus dem Gleichungssystem

$$(\boldsymbol{A} - \boldsymbol{\mu}\boldsymbol{B})\,\boldsymbol{w}^{k+1} = \boldsymbol{B}\boldsymbol{v}^k \tag{14}$$

S2: Setze
$$\boldsymbol{v}^{k-1} := \boldsymbol{w}^{k+1} / \| \boldsymbol{w}^{k+1} \|_{\boldsymbol{B}}$$

S3: Setze
$$k := k + 1$$
, goto S1

14.3.3. Bemerkung. (i) Im Unterschied zur Situation beim speziellen Eigenwertproblem ist sowohl in 14.3.1 als auch in 14.3.2 pro Schritt ein lineares Gleichungssystem zu lösen. Bei der direkten Vektoriteration ist die Koeffizientenmatrix **B** positiv definit, so daß die **LDL**^T-Faktorisierung mit positiven Diagonalelementen — und gegebenenfalls mit Diagonalpivotisierung — berechnet werden kann, vgl. 6.1.A. Die bei der inversen Iteration auftretende Matrix $A - \mu B$ ist jedoch i. allg. indefinit, so daß eine auch dann numerisch gutartige Faktorisierung — etwa die spaltenpivotisierte Gauß-Faktorisierung oder die BKP-Faktorisierung — verwendet werden sollte, vgl. 13.4.2 (i). (ii) Bei der Faktorisierung von **B** bzw. $A - \mu B$ kann eine für großes *n* fast immer vorliegende schwache Besetztheit von A und B nach den in 6.4 beschriebenen bzw. erwähnten Methoden ausgenutzt werden. Sollte *n* so groß sein, daß auch dann die Faktorisierung wegen zu großer Speicherplatzforderungen nicht möglich ist, müssen die Systeme (7) bzw. (14) durch ein inneres Iterationsverfahren gelöst werden.

(iii) Wenn in S1 die feste Verschiebung μ durch den zu v^k gehörenden Rayleigh-Quotienten $\mu = g_k$ gemäß (9) ersetzt wird, geht 14.3.2 in die implizit realisierte Rayleigh-Quotienten-Iteration 13.4.6 für C über. Allerdings muß dann als Preis für die kubische Konvergenz — siehe 13.4.7 — pro Schritt eine neue Faktorisierung von $A - g_k B$ berechnet werden. \Box

Wenn nicht nur der erste, sondern p dominante Eigenvektoren von $\{A, B\}$ gesucht sind, wird zur Teilraumiteration übergegangen. Wir diskutieren zunächst, wie die Rayleigh-Ritz-Prozedur 13.3.13 für das allgemeine Eigenwertproblem zu realisieren ist. Ausgangspunkt ist wieder die durch (1), (2) beschriebene Reduktion. Der zu $\{A, B\}$ gehörende approximierende *m*-dimensionale Teilraum $\mathcal{R}(W)$ mit $p \leq m \leq n$ werde durch die spaltenreguläre Matrix $W \in \mathbb{R}^{n,m}$ repräsentiert. Wegen (2), (4) ist dann

$$\mathcal{R}(\hat{W})$$
 mit $\hat{W} := F^{-1}W$ (15)

der entsprechende Teilraum für C. Um 13.3.12 mit C statt A anwenden zu können, benötigen wir eine durch die Spalten einer Matrix \hat{Q} repräsentierte orthonormale Basis von $\mathcal{R}(\hat{W})$; wie bisher werden die zu $\{C, I\}$ gehörenden Größen durch ein " \wedge " gekennzeichnet. Die Forderung $\mathcal{R}(\hat{Q}) = \mathcal{R}(\hat{W})$ ist gleichbedeutend mit

$$\hat{Q} = \hat{W}G$$
 mit regulärem $G \in \mathbb{R}^{m,m}$, (16)

und die Bedingung der Spaltenorthonormalität führt auf $\hat{Q}^{\mathsf{T}}\hat{Q} = G^{\mathsf{T}}\hat{W}^{\mathsf{T}}\hat{W}G = I$, also

$$\hat{W}^{\mathsf{T}}\hat{W} = W^{\mathsf{T}}F^{-\mathsf{T}}F^{-1}W = W^{\mathsf{T}}BW = G^{-\mathsf{T}}G^{-1},\tag{17}$$

man beachte (15) und (1). Da $W^{\intercal}BW$ positiv definit ist, existiert die Cholesky-Faktorisierung

$$W^{\mathsf{T}}BW = LL^{\mathsf{T}},\tag{18}$$

und (17) ist mit $G := L^{-\intercal}$ erfüllt, was auf

$$\hat{\boldsymbol{Q}} = \hat{\boldsymbol{W}} \boldsymbol{L}^{-\mathsf{T}} = \boldsymbol{F}^{-1} \boldsymbol{W} \boldsymbol{L}^{-\mathsf{T}} \tag{19}$$

führt. Die in 13.3.12 benötigte Matrix \hat{P} ist dann wegen (1)

$$\hat{\boldsymbol{P}} := \hat{\boldsymbol{Q}}^{\mathsf{T}} \boldsymbol{C} \hat{\boldsymbol{Q}} = \boldsymbol{L}^{-1} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{F}^{-\mathsf{T}} (\boldsymbol{F}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{F}) \, \boldsymbol{F}^{-1} \boldsymbol{W} \boldsymbol{L}^{-\mathsf{T}} = \boldsymbol{L}^{-1} \boldsymbol{W}^{\mathsf{T}} \boldsymbol{A} \boldsymbol{W} \boldsymbol{L}^{-\mathsf{T}}, \tag{20}$$

und die zugehörige Eigenwertzerlegung sei

$$\hat{P} = \hat{X}M\hat{X}^{\mathsf{T}}$$
 mit $M = \text{diag}(\mu_i)$ und orthogonalem \hat{X} . (21)

Die mit den Spalten von $\hat{\mathbf{l}} := (\hat{\mathbf{r}}^1, ..., \hat{\mathbf{r}}^m) := \hat{\mathbf{Q}}\hat{\mathbf{X}}$ gebildeten Ritzschen Näherungen $\{\mu_j, \hat{\mathbf{r}}^j\}$ für C bezüglich $\mathcal{R}(\hat{\mathbf{Q}}) = \mathcal{R}(\hat{\mathbf{M}})$ minimieren

$$\begin{aligned} \|\hat{R}\|_{F}^{2} &= \operatorname{tr} \left(\hat{R}^{\mathsf{T}}\hat{R}\right) = \|C\hat{V} - \hat{V}M\|_{F}^{2} = \sum_{j=1}^{m} \|C\hat{v}^{j} - \mu_{j}\hat{v}^{j}\|^{2}, \\ \hat{R} &:= C\hat{V} - \hat{V}M, \end{aligned}$$
(22)

über alle $M = \text{diag}(\mu_j)$ und alle spaltenorthonormalen \hat{V} mit $\mathcal{R}(\hat{V}) = \mathcal{R}(\hat{Q})$. Die entsprechenden Paare $\{\mu_j, v^j\}$ bezüglich $\{A, B\}$ entstehen dann wegen (2), (15) durch die Rücktransformation

$$\boldsymbol{v}^{j} := \boldsymbol{F} \hat{\boldsymbol{v}}^{j} \qquad (j = 1, \dots, m), \tag{23}$$

und für die zugehörige Matrix $V := (v^1, ..., v^m)$ gilt wegen (19) und (15)

$$V = F\hat{V} = F\hat{Q}\hat{X} = F\hat{W}L^{-\intercal}\hat{X} = WL^{-\intercal}\hat{X} =: WX$$

mit $X := (x^1, ..., x^m) := L^{-\intercal}\hat{X}.$ (24)

Zufolge (18) gilt dabei

$$X^{\mathsf{T}}(W^{\mathsf{T}}BW) X = V^{\mathsf{T}}BV = \hat{X}^{\mathsf{T}}L^{-1}W^{\mathsf{T}}BWL^{-\mathsf{T}}\hat{X} = \hat{X}^{\mathsf{T}}\hat{X} = I, \qquad (25)$$

d. h., die $\{x^j\}$ sind $W^{\mathsf{T}}BW$ -orthonormal, die $\{v^j\}$ dagegen *B*-orthonormal. Aus (20), (21) ergibt sich außerdem

$$X^{\mathsf{T}}(W^{\mathsf{T}}AW) X = \hat{X}^{\mathsf{T}}L^{-1}W^{\mathsf{T}}BWL^{-\mathsf{T}}\hat{X} = \hat{X}^{\mathsf{T}}\hat{P}\hat{X} = M,$$
(26)

d. h., die Matrix X transformiert $\{W^{T}AW, W^{T}BW\}$ kongruent in $\{M, I\}$. Schließlich kann \hat{R} in (22) mit (1) und (24) gemäß

$$\hat{\boldsymbol{R}} = (\boldsymbol{F}^{\mathsf{T}}\boldsymbol{A}\boldsymbol{F})\,\boldsymbol{F}^{-1}\boldsymbol{V} - \boldsymbol{F}^{-1}\boldsymbol{V}\boldsymbol{M} = \boldsymbol{F}^{\mathsf{T}}(\boldsymbol{A}\boldsymbol{V} - \boldsymbol{B}\boldsymbol{V}\boldsymbol{M}) =: \boldsymbol{F}^{\mathsf{T}}\boldsymbol{R}$$
(27)

mit $\mathbf{R} := (\mathbf{r}^1, \dots, \mathbf{r}^m) := \mathbf{AV} - \mathbf{BVM}$ dargestellt werden. Wieder mit (1) folgt $\|\mathbf{F}^{\mathsf{T}}\mathbf{r}\|^2 = \mathbf{r}^{\mathsf{T}}\mathbf{F}\mathbf{F}^{\mathsf{T}}\mathbf{r} = \mathbf{r}^{\mathsf{T}}\mathbf{B}^{-1}\mathbf{r} = \|\mathbf{r}\|_{\mathbf{B}^{-1}}^2$, so daß (22) in

$$\|\hat{\boldsymbol{R}}\|_{F}^{2} = \operatorname{tr}\left(\boldsymbol{R}^{\mathsf{T}}\boldsymbol{B}^{-1}\boldsymbol{R}\right) = \sum_{j=1}^{m} \|\boldsymbol{F}^{\mathsf{T}}\boldsymbol{r}^{j}\|^{2} = \sum_{j=1}^{m} \|\boldsymbol{r}^{j}\|_{\boldsymbol{B}^{-1}}^{2} = \sum_{j=1}^{m} \|\boldsymbol{A}\boldsymbol{v}^{j} - \mu_{j}\boldsymbol{B}\boldsymbol{v}^{j}\|_{\boldsymbol{B}^{-1}}^{2} \quad (28)$$

übergeht. Die Paare $\{\mu_j, v^j\}$ minimieren also die mit B^{-1} gewichteten Residuumsnormen. Als Zusammenfassung erhalten wir das folgende Analogon zu 13.3.12:

14.3.4. Satz. Gegeben seien das symmetrische Büschel $\{A, B\}$ mit positiv definitem **B** und die reguläre Matrix $W \in \mathbb{R}^{n,m}$, $m \leq n$. Es sei $X = (x^1, \ldots, x^m) \in \mathbb{R}^{m,m}$ diejenige Matrix, die $\{\bar{A}, \bar{B}\}$ mit $\bar{A} := W^{\mathsf{T}}AW$, $\bar{B} := W^{\mathsf{T}}BW$ kongruent in $\{M, I\}$ mit $M = \text{diag}(\mu_i)$ transformiert, d. h., es gelte

$$ar{m{A}}m{x^j} = \mu_j ar{m{B}}m{x^j} \qquad (j=1,...,m)$$

mit \overline{B} -orthonormalen Eigenvektoren $\{x^{j}\}$. Aus W und X werde

$$V := (\boldsymbol{v}^1, \ldots, \boldsymbol{v}^m) := WX$$

gebildet. Dann sind die Paare $\{\mu_i, v^j\}$ die Ritzschen Eigenpaare für das durch $\{A, B\}$ definierte allgemeine Eigenwertproblem bezüglich des Teilraumes $\mathcal{R}(W)$, d. h., sie minimieren (28) über alle diagonalen $M \in \mathbb{R}^{m,m}$ und alle $V \in \mathbb{R}^{n,m}$ mit $V^{\mathsf{T}}BV = I$ und $\mathcal{R}(V) = \mathcal{R}(W)$.

Wir geben abschließend eine zu 13.3.17(i) analoge Kombination der Teilraumiteration und des RR-Algorithmus aus 14.3.4 an. Aus den bereits diskutierten Gründen beziehen wir uns dabei nicht auf C, sondern auf $(C - \mu I)^{-1}$, wobei die Realisierung wie in 14.3.2 implizit mit den Originaldaten und rücktransformierten Teilräumen erfolgt.

14.3.5. Implizite inverse Teilraumiteration mit RR-Algorithmus.

- S0: Wähle $\mu \in \mathbf{R}$ und berechne numerisch gutartige Faktorisierung von $\mathbf{A} \mu \mathbf{B}$, wähle m mit $1 \leq m \leq n$ und $V_1 \in \mathbf{R}^{n,m}$ mit $V_1^{\mathsf{T}} \boldsymbol{B} V_1 = \boldsymbol{I}$, setze k := 1
- S1: Berechne $W_{k+1} = (w^{k+1,1}, \dots, w^{k+1,m}) \in \mathbb{R}^{n,m}$ unter Verwendung der in S0 berechneten Faktorisierung von $A \mu B$ als Lösung von

$$(\boldsymbol{A} - \boldsymbol{\mu}\boldsymbol{B}) \boldsymbol{W}_{k+1} = \boldsymbol{B}\boldsymbol{V}_k \tag{29}$$

S2: Bilde $\bar{A}_{k+1} := W_{k+1}^{\mathsf{T}} A W_{k+1}, \ \bar{B}_{k+1} := W_{k-1}^{\mathsf{T}} B W_{k+1}.$ Bestimme die Eigenwerte $\{\mu_i\}$ und die zugehörigen \overline{B}_{k+1} -orthonormalen Eigenvektoren $\{x^i\}$ des allgemeinen symmetrischen Eigenwertproblems

$$ar{oldsymbol{A}}_{k+1}oldsymbol{x}^j = \mu_iar{oldsymbol{B}}_{k+1}oldsymbol{x}^j \qquad (j=1,...,m)$$

$$\tag{30}$$

 $A_{k+1}x' = \mu_j B_{k+1}x^j$ (j = 1, ..., m)S3: Bilde $V_{k+1} := (v^{k+1,1}, ..., v^{k+1,m}) := W_{k+1}X_{k+1}$ mit $X_{k+1} := (x^1, ..., x^m)$ S4: Setze k := k + 1, goto S1

14.3.6. Bemerkung. (i) Unter geeigneten Voraussetzungen approximieren die Ritzschen Eigenpaare $\{\mu_i, v^j\}$ die in der Numerierung (11) zugehörigen Eigenpaare $\{\lambda_i, \boldsymbol{z}^j\}$ von $\{\boldsymbol{A}, \boldsymbol{B}\}$, siehe 13.3.B,C.

(ii) Die Matrixgleichung (29) entspricht den m linearen Gleichungen

 $(A - \mu B) w^{k+1,j} = B v^{k,j}$ (j = 1, ..., m)

für die Spalten von W_{k+1} ; die Koeffizientenmatrix ist für alle *j* dieselbe. Für ihre Lösung gilt 14.3.3 sinngemäß.

(iii) In der Regel ist $m \ll n$, d. h., das Eigenwertproblem (30) ist von kleiner Dimension und i. allg. voll besetzt. Zur Lösung empfiehlt sich daher die explizite Reduktion auf ein spezielles Eigenwertproblem nach den Methoden aus 14.2.

14.3.7. Bemerkung. Alternative Verfahren zur Lösung des allgemeinen symmetrischen Eigenwertproblems sind einmal die Relaxations- und Gradientenverfahren, die auf einer Minimierung des Rayleigh-Quotienten $o = v^{\mathsf{T}} A v / v^{\mathsf{T}} B v$ bzw. einer geeigneten, auf Teilräume definierten Verallgemeinerung von diesem beruhen, siehe B 14.5 für Hinweise. Zum anderen können die bereits in 13.3 erwähnten Verfahren vom Lanczos-Typ mit C bzw. $(C - \mu I)^{-1}$, aber auch direkt mit $B^{-1}A$ bzw. $(A - \mu B)^{-1} B$ angewendet werden; man beachte, daß die beiden letztgenannten Matrizen bezüglich des durch B definierten Skalarproduktes selbstadjungiert sind, siehe Ü 14.3.1. Die Lanczos-Algorithmen, die in den letzten Jahren sehr intensiv untersucht worden sind, konvergieren i. allg. deutlich besser als die Teilraumiteration. Wir weisen auch hier auf die Spezialliteratur hin, siehe wieder B 14.5.

Übungsaufgabe

Ü 14.3.1. Es sei **B** symmetrisch und positiv definit. Dann heißt die Matrix **M** selbstadjungiert bezüglich $[x, y]_B$ oder kurz **B**-selbstadjungiert, wenn

 $[\boldsymbol{x}, \boldsymbol{M}\boldsymbol{y}]_{\boldsymbol{B}} = [\boldsymbol{M}\boldsymbol{x}, \boldsymbol{y}]_{\boldsymbol{B}}$

für alle $\boldsymbol{x}, \boldsymbol{y}$ gilt. Man zeige:

(i) *M* ist *I*-selbstadjungiert genau dann, wenn *M* symmetrisch ist.

(ii) Es sei A symmetrisch und μ sei kein Eigenwert von $\{A, B\}$. Dann sind die Matrizen $B^{-1}A$ und $(A - \mu B)^{-1} B$ beide **B**-selbstadjungiert.

Bemerkungen zum Kapitel 14

B 14.1. Eine detaillierte Diskussion des allgemeinen symmetrischen Eigenwertproblems ist bei PARLETT [80a] zu finden, wo viele nützliche Fakten und Hinweise gegeben werden. Zur Störungstheorie sei auf STEWART [72, 78, 79a], ELSNER/SUN [82], SUN [83] verwiesen. Neuere Resultate und weiterführende Literaturhinweise über Probleme, die mit Matrixbüscheln zusammenhängen, sind in dem von Kågström/Ruhe [83] herausgegebenen Sammelband enthalten.

B 14.2. Die in 14.2 angegebenen Methoden zur expliziten Reduktion des allgemeinen Eigenwertproblems auf ein spezielles gehören zu den Standardtechniken der numerischen linearen Algebra, siehe WILKINSON [65], PETERS/WILKINSON [70b] und PARLETT [80a]. Sie sind an die positive Definitheit von \boldsymbol{B} gebunden. Falls \boldsymbol{B} nur positiv semidefinit ist, also $\delta_n = 0$ Eigenwert von \boldsymbol{B} ist, kann die nach FIX/HEIBERGER [72] benannte Reduktion vorgenommen werden; sie wird auch bei PARLETT [80a] beschrieben. Numerisch unterscheidet sich der Fall einer semidefiniten Matrix \boldsymbol{B} nicht von dem einer positiv definiten mit schlechter Kondition.

B 14.3. Die in 14.2.2(iii) erwähnte Transformation eines Büschels $\{A, B\}$ aus Bandmatrizen in ein kongruentes Büschel $\{T, I\}$ mit tridiagolem T geht auf CRAWFORD [73] zurück. Andere Verfahren für Bandmatrizen werden von PETERS/WILKINSON [69], SCHWARZ [77, 80] und WALDVOGEL [82] beschrieben.

B 14.4. Ein alternatives Verfahren zur Lösung des nicht notwendig symmetrischen allgemeinen Eigenwertproblems stellt der QZ-Algorithmus von MOLER/STEWART [73] dar, bei dem $\{A, B\}$ mittels einer Folge orthogonaler Äquivalenztransformationen simultan in ein äquivalentes Büschel $\{R_A, R_B\}$ mit oberen Dreiecksmatrizen R_A, R_B transformiert wird. Die Quotienten $(R_A)_{ij}/(R_B)_{ij}$ sind dann die Eigenwerte von $\{A, B\}$. Falls B regulär ist, ist der QZ-Algorithmus identisch mit einer impliziten Realisierung des QR-Algorithmus für AB^{-1} . Wir gehen auf diesen Algorithmus nicht weiter ein, da bei ihm die Symmetrie von $\{A, B\}$ zerstört wird; der Aufwand für symmetrische Büschel ist derselbe wie für nichtsymmetrische. **B 14.5.** Die Verfahren der direkten und inversen Vektor- und Teilraumiteration für das allgemeine Eigenwertproblem werden ausführlich bei PARLETT [80a] behandelt. Stärker auf die Belange der Ingenieurwissenschaften eingegangen wird in den Büchern von BATHÉ/ WILSON [76] und JENNINGS [77]. Algorithmus 14.3.5 ist z. B. bei PARLETT [80b] zu finden. Eine Variante der RQ-Iteration für das quadratische Eigenwertproblem (14.1.17) ohne Reduktion auf ein lineares der doppelten Dimension beschreiben SCOTT/WARD [82]. Zur Literatur über Relaxations- und Gradientenverfahren siehe B 13.5. Zur Klasse der Lanczos-Algorithmen sei wieder auf PARLETT [80a] und CULLUM/WILLOUGHBY [85] verwiesen. Hier sind auch in nächster Zeit noch interessante neue Ergebnisse zu erwarten.

15. Zusammenfassung zum Teil IV

Das Vorgehen zur Lösung des speziellen symmetrischen Eigenwertproblems $Ax = \lambda x$ hängt wesentlich von der Struktur der Matrix A und der Anzahl der gesuchten Eigenwerte und gegebenenfalls Eigenvektoren ab.

Wenn A voll besetzt ist und mehr als ca. 25% der Eigenwerte bzw. Eigenvektoren berechnet werden sollen, ist die Tridiagonalisierung von A nach 13.5.A und die vollständige Lösung des entstehenden tridiagonalen Eigenwertproblems nach dem QR-Algorithmus aus 13.7.C – etwa in der impliziten QL-Version – der effektivste Weg. Von der Genauigkeit vergleichbar, aber vom Aufwand her ungünstiger ist das Jacobi-Verfahren aus 13.2. Es sollte nur dann angewendet werden, wenn A wenig von einer Diagonalmatrix abweicht, eine geringe Genauigkeit gefordert ist oder Wert auf ein kurzes, kompaktes Programm gelegt wird. Sind nur wenige Eigenwerte gesucht, sollten diese durch Schneiden des Spektrums der Tridiagonalmatrix nach dem Bisektionsverfahren aus 13.6.B ermittelt werden. Zugehörige Eigenvektoren können mittels inverser Iteration nach 13.4 bestimmt werden; bei mehrfachen Eigenwerten bzw. Eigenwerthaufen sind dann Zusatzmaßnahmen erforderlich, um die Orthogonalität der berechneten Eigenvektoren zu sichern.

In der Regel werden die betragskleinen Eigenwerte von A mit geringerer relativer Genauigkeit als die betragsgroßen berechnet, siehe 13.1.B. Für gestuftes A können jedoch auch die betragskleinen Eigenwerte bei sachgemäßem Vorgehen ausreichend genau berechnet werden. Für ansteigend gestuftes A muß dazu die Tridiagonalisierung von rechts nach links erfolgen, und der QR-Algorithmus ist in der QL-Variante zu verwenden, siehe 13.5.A und 13.7.C.

Wenn A eine Bandmatrix ist, kann die Tridiagonalisierung mittels Givens-Drehungen bei Erhalt der Bandbreite vorgenommen werden, siehe 13.5.B. Da dann i. allg. nur wenige Eigenwerte bzw. Eigenvektoren gesucht sind, sollten diese durch Schneiden des Spektrums bzw. inverse Iteration bestimmt werden, wobei unter Umständen günstiger direkt mit $A - \mu I$ statt mit der Tridiagonalmatrix $T - \mu I$ gearbeitet wird.

Wenn A schwach besetzt und von hoher Dimension ist, sind die Verfahren der Vektor- und Teilraumiteration zu empfehlen. Falls der betragsgrößte oder die pbetragsgrößten Eigenwerte und zugehörigen Eigenvektoren gesucht sind, kann die direkte Vektor- oder Teilraumiteration aus 13.3 verwendet werden, bei der nur Ausdrücke der Form $\boldsymbol{w} = A\boldsymbol{v}$ für gegebenes \boldsymbol{v} ausgewertet werden müssen. Sind dagegen die betragskleinsten oder allgemeiner die nahe an einer vorgegebenen Zahl μ gelegenen Eigenwerte gesucht, so kommt die inverse Vektor- oder Teilraumiteration in Frage, siehe 13.4. Pro Schritt sind dann ein oder mehrere Gleichungssysteme $(\boldsymbol{A} - \mu \boldsymbol{I}) \boldsymbol{w} = \boldsymbol{v}$ für gegebenes \boldsymbol{v} zu lösen, wobei die schwache Besetztheit nach den in 6.4 angegebenen Methoden ausgenutzt werden kann. Bei der Kopplung mit dem RR-Algorithmus kommt pro Schritt die Lösung eines speziellen Eigenwertproblems niedriger Dimension hinzu. Alternative Methoden sind Relaxations- und Gradientenverfahren und die Verfahren vom Lanczos-Typ, siehe dazu B 13.5 und B 13.7.

Die berechneten Eigenwerte und Eigenvektoren können nach DONGARRA/MOLER/ WILKINSON [83] verbessert werden.

Für das allgemeine symmetrische Eigenwertproblem $Az = \lambda Bz$ mit positiv definitem **B** ist im Fall voll besetzter Matrizen **A**, **B** die explizite Reduktion auf ein spezielles symmetrisches Eigenwertproblem nach 14.2 der Standardweg. In den meisten Fällen kann dazu die möglichst mit Diagonalpivotisierung berechnete **LDL**^T-Faktorisierung von **B** verwendet werden; andernfalls sollte mit der stabileren Eigenwertzerlegung gearbeitet werden. Falls **B** semidefinit oder sehr schlecht konditioniert ist, ist die Reduktion nach FIX/HEIBERGER — siehe B 14.2 — zu empfehlen.

Wenn A, B Bandmatrizen hoher Dimension sind, kann das Büschel $\{A, B\}$ unter Erhalt der Bandbreite kongruent in $\{T, I\}$ mit tridiagonalem T transformiert werden, siehe B 14.3. Da dann i. allg. nur wenige Eigenwerte gesucht sind, sollten diese durch Schneiden des Spektrums bestimmt werden, wobei auch direkt auf $A - \mu B$ zurückgegriffen werden kann. Letzteres trifft auch auf die inverse Iteration zur Ermittlung zugehöriger Eigenvektoren zu.

Für schwach besetzte Matrizen A, B hoher Dimension gelten die Ausführungen für das spezielle Eigenwertproblem sinngemäß. Die direkte bzw. inverse Vektorund Teilraumiteration kann direkt mit den Originaldaten durchgeführt werden, wobei bei der direkten Iteration Gleichungssysteme des Typs Bw = Av, bei der inversen solche des Typs $(A - \mu B)w = Av$ für gegebenes v zu lösen sind. Bei Kombination mit dem RR-Algorithmus tritt pro Schritt zusätzlich ein allgemeines symmetrisches Eigenwertproblem niedriger Dimension auf, siehe 14.3. Auch hier stellen Relaxations- und Gradientenverfahren sowie die Verfahren vom Lanczos-Typ alternative Möglichkeiten dar, siehe B 13.5, B 13.7 und B 14.5.

Den Problemkreis der Rang-1-Modifikation symmetrischer Eigenwertaufgaben haben wir aus Platzgründen nicht diskutiert; siehe dazu GOLUB [73] und BUNCH/ NIELSEN/SORENSEN [78].

In manchen Anwendungen spielen sog. inverse Eigenwertprobleme eine Rolle, bei denen Matrizen A bzw. $\{A, B\}$ innerhalb gewisser Klassen so festzulegen sind, daß sie vorgegebene Eigenwerte besitzen. Wir weisen dazu auf GOLUB [73], FRIEDLAND [75], DE BOOR/GOLUB [78], WANG/GARBOW [83], aber auch auf die entsprechenden Abschnitte bei PARLETT [80a] und GOLUB/VAN LOAN [83] hin.

Eigenwertzerlegungen können auch zur Berechnung von Matrixfunktionen f(A)– etwa e^A – herangezogen werden; siehe GOLUB/VAN LOAN [83] und speziell MOLER/VAN LOAN [78]. Eigenwertprobleme für nichtsymmetrische Matrizen sind von der Theorie wie von der Praxis wesentlich komplizierter als solche für symmetrische; wir verweisen auf WILKINSON [65] und GOLUB/VAN LOAN [83].

16. Software für Aufgabenklassen der numerischen linearen Algebra

Wir beginnen mit der — möglicherweise überspitzten — Feststellung, daß das Hauptziel der numerischen linearen Algebra die Schaffung von qualitativ hochwertiger numerischer Software ist. Unter *numerischer Software* soll dabei die Gesamtheit der Computerprogramme verstanden werden, in denen Algorithmen zur Lösung numerischer Probleme aus gewissen Problemklassen dem fortgeschrittenen Entwicklungsstand entsprechend realisiert worden sind. In diesem Sinne impliziert die Bezeichnung "numerische Software" einen Qualitätsanspruch, der natürlich relativ ist: Er hängt vom jeweiligen theoretischen Erkenntnisstand über die Problemklasse und die zu ihrer Lösung bekannten Algorithmen, von den Möglichkeiten der Computer, den verwendeten Programmiersprachen und Betriebssystemen, vom Stand der Programmiertechnologie und schließlich auch vom ins Auge gefaßten Nutzerkreis ab.

Die auf die Software wirkenden Faktoren weisen eine große Variationsbreite auf: Die zu einer Klasse gehörenden Probleme können in "kleine" und "große" eingeteilt werden. Dabei versteht man unter kleinen Problemen in der Regel solche, bei denen alle Eingangsdaten und Zwischenergebnisse im Hauptspeicher mit schnellem Zugriff Platz finden und wo eine mögliche schwache Besetztheit nicht systematisch ausgenutzt wird; die Einteilung hängt also vom Computer ab. Das Spektrum der Computer selbst reicht vom Ein-Chip-Rechner bis zu den imposanten Supercomputern mit mehreren 100 Millionen Gleitpunktoperationen pro Sekunde. Die Operationsgeschwindigkeiten können sich um Faktoren in der Größenordnung 10⁴ bis 10⁵ unterscheiden; dasselbe trifft in etwa auch für die Preise zu.

Die Mehrzahl der Nutzer numerischer Software sind Ingenieure, Naturwissenschaftler u. a., die bei der Behandlung ihrer fachspezifischen Aufgaben auf numerische Probleme treffen und diese möglichst einfach gelöst haben wollen. Sie verfügen in der Regel nicht über Spezialkenntnisse in Numerischer Mathematik und sehen die Software als "black box" an, die aus den Eingangsdaten in hoffentlich vernünftiger, sie im einzelnen jedoch nicht interessierender Weise die gewünschten Ausgangsdaten erzeugen soll. Der Kreis dieser Nutzer vergrößert sich in dem Maße, wie leistungsfähige Rechentechnik am Arbeitsplatz verfügbar gemacht und die Mathematik in den Fachdisziplinen mehr und mehr angewendet wird. Es ist klar, daß dieser große Nutzerkreis auf entsprechend konzipierte Software angewiesen ist. Den Gegenpol bildet die kleine Schar der hochspezialisierten Nutzer, die wesentlich komplexere Aufgaben bearbeiten und z. T. die derzeit verfügbaren Erkenntnisse der Theorie wie auch die Möglichkeiten der Computer bis zur Grenze ausschöpfen. Für das Gebiet der numerischen linearen Algebra muß außerdem beachtet werden, daß die einzelnen Probleme i. allg. nicht eigenständig, sondern als Bestandteil umfangreicher technisch-wissenschaftlicher u. ä. Berechnungen auftreten. Die Eingangsdaten sind das Ergebnis vorhergehender Rechnungen, und die Ausgangsdaten werden im nachfolgenden Schritt als Eingangsdaten benutzt, ohne daß sie etwa über einen Drucker ausgegeben oder über einen Bildschirm angezeigt werden. Wegen dieses sog. verdeckten Datenflusses ist es zweckmäßig — wenn auch nicht immer nötig —, von den Ausgangsdaten die höchste erreichbare Genauigkeit zu verlangen.

Die bisherigen Ausführungen motivieren die folgenden Forderungen an numerische Software:

- Zuverlässigkeit: Alle Probleme, die zur Problemklasse gehören, für die der implementierte Algorithmus geeignet ist und die im Rahmen der Computerarithmetik auch gelöst werden können, sollen mit maximal möglicher Genauigkeit gelöst werden. Wünschenswert ist eine Information über die erreichte Genauigkeit. Die verwendeten Algorithmen sollten numerisch stabil und möglichst sogar numerisch gutartig sein; gegebenenfalls sollte ein Stabilitätsverlust durch Beobachtung der dafür charakteristischen Größen signalisiert werden. Sonderfälle sollten sachgemäß behandelt und eine erfolglose Bearbeitung mit Angabe von Gründen mitgeteilt werden.
- Effektivität: Die numerische Lösung soll mit möglichst geringer Rechenzeit und möglichst wenig zusätzlichem Speicherplatz berechnet werden. Fast immer sind Zuverlässigkeit und Effektivität gegensätzliche Forderungen: Die zuverlässigsten Algorithmen sind meist nicht die billigsten, und die korrekte Behandlung aller denkbaren Sonderfälle erfordert zusätzlichen Aufwand. Für die große Klasse der Nutzer, die nur wenige nicht zu große Probleme lösen wollen, spielt die Effektivität gegenüber der Zuverlässigkeit i. allg. eine untergeordnete Rolle. Sind dagegen sehr viele und/oder sehr große Probleme zu lösen, so wird die Forderung nach Effektivität dominieren.
- Portabilität: Die Software soll möglichst ohne bzw. mit wenigen Änderungen auf unterschiedlichen Computern und unter unterschiedlichen Betriebssystemen verwendbar sein. Da numerische Software heute fast ausschließlich in der Programmiersprache FORTRAN geschrieben wird, sollte eine solche Teilmenge dieser Sprache verwendet werden, die auf der Mehrzahl und insbesondere auch auf kleineren Computern verfügbar ist. Die Software sollte keine Ein- und Ausgabeoperationen enthalten; wenn solche nicht vermieden werden können, sollten sie über Unterprogramme realisiert werden, die dann i. allg. vom Computer und der genutzten Peripherie abhängen. Die Datenübermittlung sollte möglichst über Parameterlisten und nicht über COMMON-Bereiche erfolgen. Computerspezifische Größen wie v, MAX, MIN usw. sollten nicht explizit verwendet werden. Der Kleinheitstest

$$\mathbf{if} \; (|\alpha| \le \nu * |\beta|) \tag{1}$$

kann z. B. etwa gleichwertig in der Form

$$\mathbf{if} \left(|\alpha| + |\beta| = |\beta| \right) \tag{2}$$

ohne explizite Verwendung von ν realisiert werden. Wenn doch gewisse Computerkonstanten benutzt werden müssen, sollten sie am Programmanfang durch eine DATA-Anweisung und damit leicht erkennbar gemacht werden, oder sie sollten durch Unterprogramme vermittelt werden, vgl. auch B 2.5.

Nutzerfreundlichkeit: Dazu gehört eine verständliche und übersichtliche Dokumentation, und zwar möglichst auf zwei Niveaus: eins für den black-box-Nutzer, ein zweites für den erfahrenen Nutzer, der Anpassungen, Erweiterungen o. ä. vornehmen und daher im Detail informiert sein will. Letzteres wird durch eine ausreichende Erläuterung der Quelltexte durch Kommentare gefördert. Sofern die Auswahl der Verfahren nicht automatisch erfolgt, sollten Hinweise über die Leistungsfähigkeit und Anwendungsgebiete der implementierten Algorithmen gegeben werden. Programme, die für eine oder mehrere benachbarte Problemklassen gedacht sind, sollten nach einheitlichen Gesichtspunkten implementiert und dokumentiert werden, in verschiedenen Programmen auftretende Größen mit derselben oder ähnlicher Bedeutung sollten dieselben Bezeichnungen haben, die Parameterlisten sollten vereinheitlicht und nicht zu lang sein, die Bezeichnung der Teilprogramme und Parameter sollte nach einem durchschaubaren, einheitlichen Prinzip erfolgen.

Eine nach den letztgenannten Gesichtspunkten erarbeitete Zusammenstellung von aufeinander abgestimmten Programmen für die Lösung von Aufgaben aus gewissen Aufgabenklassen wird *Programmpaket* genannt. Programmpakete stellen die heute übliche Form dar, in der numerische Software ausgearbeitet, angeboten und genutzt wird. Wenn der Automatisierungsgrad durch Ermöglichen der anwendernahen Formulierung des Problems in einer angepaßten Spezialsprache, automatische Algorithmenauswahl, Übernahme der Datenorganisation und der Ein- und Ausgabe u. ä. weiter erhöht wird, spricht man von *Programmsystemen*. Die Erarbeitung solcher Programmsysteme ist allerdings außerordentlich aufwendig, so daß sie sich nur für komplexere und häufig auftretende Problemklassen wie etwa Schaltkreisentwurf, mechanische Berechnungen u. ä. lohnt.

Als erstes maßstabsetzendes Programmpaket zur numerischen linearen Algebra müssen die ALGOL-Prozeduren aus dem von WILKINSON/REINSCH [71] herausgegebenen und wesentlich mitgetragenen Band "Linear Algebra" des "Handbook for Automatic Computation" angesehen werden. Mit diesen Programmen wurde ein bis zu diesem Zeitpunkt nicht erreichter Stand demonstriert, der die weitere Entwicklung im Bereich der numerischen Software wesentlich beeinflußt hat. Auf der Grundlage der Eigenwertprogramme des "Handbooks" entstanden in den darauffolgenden Jahren die weiterentwickelten FORTRAN-Pakete EISPACK von SMTTH et al. [74] und GABBOW et al. [77]; eine überarbeitete dritte, allerdings nicht in Buchform veröffentlichte Version liegt inzwischen auch vor. Die EISPACK-Routinen sind für das symmetrische und nichtsymmetrische Eigenwertproblem konzipiert. Den Problemklassen der linearen Gleichungssysteme und Quadratmittelprobleme ist das von DONGABRA et al. [79] erarbeitete Programmpaket LINPACK gewidmet. Hinsichtlich Sorgfalt der Implementierung, Berücksichtigung theoretischer Erkenntnisse und Programmierstil können die in den genannten Paketen enthaltenen FORTRAN-Routinen als Musterbeispiel hochwertiger numerischer Software dienen und jedem, der sich mit der Implementierung numerischer Algorithmen befaßt, zum intensiven Studium empfohlen werden. Auch in der Ausbildung der Studenten ist es viel vernünftiger, sie etwa ein LINPACK-Programm verstehen und testen zu lassen als zu fordern, daß sie ein — dann in der Regel schlechtes — Programm zur Lösung des entsprechenden Problems selbst schreiben.

Eine Besonderheit des letztgenannten Paketes LINPACK ist die durchgängige Verwendung der sog. "Basic Linear Algebra Subprograms for FORTRAN Usage" - kurz: BLAS -, siehe B 16.2. Diese Unterprogramme realisieren die in der numerischen linearen Algebra häufigsten Grundoperationen wie Addition des Vielfachen eines Vektors zu einem anderen, Skalarproduktberechnung usw., wobei die Vektoren sowohl Spalten als auch Zeilen von Matrizen sein können, die in zweidimensionalen Feldern gespeichert sind. Der Ersatz der sonst benötigten FORTRAN-Schleifen durch die BLAS-Aufrufe führt zu einer Verkürzung der Programme. Außerdem besteht die Möglichkeit, die BLAS in der Maschinensprache zu kodieren und dabei Besonderheiten des Computers - etwa vorhandene Vektorprozessoren - auszunutzen, womit die Effektivität erhöht werden kann. Beim Übergang zu anderen Computern brauchen dann nur die BLAS neu geschrieben zu werden. Als Nachteil steht dem der durch den Organisationsaufwand bei den BLAS-Aufrufen bedingte Effektivitätsverlust bei der Behandlung von Problemen niedriger Dimension gegenüber, vgl. die entsprechenden Angaben in der LINPACK-Dokumentation.

Besondere Aufmerksamkeit muß bei Verwendung von FORTRAN den in der linearen Algebra typischen zweidimensionalen Feldern gewidmet werden. Da solche Felder spaltenweise abgespeichert werden, sollten die Algorithmen stets spaltenweise orientiert werden, d. h., der innerste Schleifenindex sollte in einer Spalte laufen. Ist die (M,N)-Matrix A z. B. in dem gleichnamigen Feld A der Dimension (M,N) gespeichert worden, so entspricht dem zweifach indizierten Element A(I,J) in der im Speicher vorhandenen linearen Anordnung das einfach indizierte Element A(L) mit

$$L = (J-1)*M + I,$$
 (3)

vgl. 1.1.C. Bei Zugriff auf die J-te Spalte werden die aufeinanderfolgenden Elemente A(L) mit L zwischen (J-1)*M + 1 und J*M benötigt; das Inkrement von L ist 1. Bei Zugriff auf die I-te Zeile wächst L dagegen von I bis (N-1)*M + I um das Inkrement M. Bei Betriebssystemen mit virtuellem Speicher kann dies zu häufigem Seitenwechsel und daher zu großem Effektivitätsverlust führen.

Das zweite Problem entsteht, wenn mit einem Programm mehrere Probleme mit Matrizen unterschiedlicher Dimensionen bearbeitet werden sollen, was in den Anwendungen meist der Fall ist. Da die Vereinbarung von Feldern variabler Dimension in FORTRAN nicht zulässig ist, muß am Anfang des Hauptprogramms ein Feld mit maximaler Dimension vereinbart werden, das alle aktuellen Matrizen aufnimmt, etwa durch die Vereinbarung

$$DIMENSION A(100,20). \tag{4}$$

In diesem Feld A können (M,N)-Matrizen — sie mögen ebenfalls mit A bezeichnet sein — mit $M \leq 100$, $N \leq 20$ gespeichert werden. Als Beispiel betrachten wir eine (50, 10)-Matrix A, die im Hauptprogramm durch die Anweisungen



erzeugt werden möge. Im danebenstehenden Schema sind die dadurch belegten Teile des Feldes A wie üblich gekennzeichnet worden. In linearer Reihenfolge stehen also am Anfang die 50 Elemente der ersten Spalte, dann kommen 50 nicht belegte Plätze, dann die 50 Elemente der zweiten Spalte, wieder 50 "Lücken" usw. Einem Unterprogramm – etwa einer SUBROUTINE QRFACT zur *QR*-Faktorisierung der Matrix *A* nach HOUSEHOLDER – muß dann neben der aktuellen Zeilenzahl M=50 und Spaltenzahl N=10 noch die erste Dimension (engl. "leading dimension") LDA=100 des tatsächlich gemäß (4) vereinbarten Feldes mitgeteilt werden. Das Unterprogramm müßte also mit

SUBROUTINE QRFACT(A,LDA,M,N,...)
$$(5)$$

DIMENSION
$$A(LDA,1)$$
 (6)

beginnen, und im Hauptprogramm muß dem aktuellen Parameter LDA vor Aufruf von QRFACT der Wert 100 durch die Anweisung

$$LDA = 100$$
 bzw. DATA $LDA/100/$ (7)

zugewiesen werden. Man beachte, daß durch (6) kein Feld vereinbart wird, sondern dem Compiler wird lediglich mitgeteilt, daß A ein zweidimensionales Feld der Zeilenzahl LDA ist, so daß im Unterprogramm der Index von A(I,J) gemäß (3) mit LDA statt M richtig ausgewertet wird. Da die zweite Dimension — die Spaltenzahl des vereinbarten Feldes — in (3) überhaupt nicht vorkommt und daher für die Indexauswertung nicht benötigt wird, kann in (6) für sie 1 gesetzt werden.

Bei IBM- und ESER-Computern wird statt (6) häufig

REAL A(LDA,1)

geschrieben; bei doppelter Genauigkeit ist

```
DOUBLE PRECISION A(LDA,1) bzw. REAL*8 A(LDA,1)
```

zu verwenden.

Würde das Unterprogramm in naiver Weise in der Form

```
SUBROUTINE QRFACT(A,M,N,...)
DIMENSION A(M,1)
```

programmiert werden, so müßten die Elemente der (50,10)-Matrix A spaltenweise nacheinander auf den ersten 500 Plätzen des Feldes A stehen, d. h., das Feld müßte nach dem Schema



belegt werden. Im Hauptprogramm könnte dann nicht mit zweifacher Indizierung gearbeitet werden, da die dort gültige Auswertungsformel (3) mit LDA = 100 statt M nicht der oben angegebenen Belegung entspricht. Zum Beispiel wird das im Unterprogramm mit A(1,2) erreichbare Element im Hauptprogramm als A(51,1) geführt. Die Indexrechnung im Hauptprogramm müßte also vom Bearbeiter selbst in eindimensionaler Weise organisiert werden, was umständlich ist, zu unübersichtlichen Programmen führt und besonders für Anfänger eine beliebte Fehlerquelle darstellt.

Wir bemerken abschließend, daß die oben zitierten Programmpakete für Probleme mit voll besetzten Matrizen — bei Gleichungssystemen und Eigenwertaufgaben auch für solche mit Bandmatrizen — konzipiert wurden. Zur Software für Probleme mit großen und schwach besetzten Matrizen sei auf B 16.4 verwiesen.

Bemerkungen zum Kapitel 16

B 16.1. Es gibt inzwischen eine umfangreiche Literatur über numerische Software und speziell über Software zur numerischen linearen Algebra. Wir zitieren lediglich die Sammelbände von RICE [77], JACOBS [78], FOSDICK [79] und COWELL [84] sowie das informative, auch für Anfänger geeignete Buch von RICE [81]. Zum Problem der Portabilität — und andere interessante Fragen — sei auf Cowell [76] verwiesen. Auf die in den letzten Jahren sehr stark in den Blickpunkt des Interesses getretenen Vektor- und Parallelcomputer gehen u. a. HELLER [78], RODRIGUE [82], DONGARRA et al. [84] und KOTOV/MIKLOŠKO [84] ein. In jüngster Zeit werden auch sog. "array processors" und "systolic arrays" untersucht; als Beispiele seien etwa die Arbeiten von HELLER/IFSEN [83] und BOJANCZYK/BRENT/KUNG [84] genannt.

B 16.2. FORTRAN-Programme für Quadratmittelprobleme mit voll besetzten Matrizen sind außer im LINPACK auch bei LAWSON/HANSON [74] zu finden. Für Eigenwertaufgaben mit voll besetzten und Bandmatrizen wurden von LANG et al. [79] auf der Grundlage von EISPACK das weiterentwickelte Programmsystem MEIWEP erarbeitet; siehe auch LANG [82]. Die BLAS gehen auf LAWSON et al. [79] zurück; sie sind auch als Anhang im LINPACK aufgelistet.

B 16.3. Eine Vorstellung über die Zentralprozessorzeit, die zur Lösung eines Gleichungssystems der Dimension 100 mittels der LINPACK-Implementierung des Gaußschen Algorithmus erforderlich ist, vermitteln die von DONGARBA [84] zusammengestellten Tabellen, wo eine Vielzahl von Computern vom Apple II bis zur CRAY-XM erfaßt worden ist. Instruktive Zeitvergleiche sind auch in den LINPACK- und EISPACK-Dokumentationen enthalten.

B 16.4. Eine Übersicht über Software für Probleme mit schwach besetzten Matrizen gibt DUFF [84]. Zu den am meisten verwendeten Paketen zur Lösung von Gleichungssystemen mit schwach besetzten, symmetrischen und positiv definiten Matrizen gehören das von GEORGE/ LIU [81] entwickelte SPARSPAK und das auf EISENSTAT et al. [77] zurückgehende Yale Sparse Matrix Package. Falls die Matrix indefinit ist und aus Stabilitätsgründen auch (2, 2)-Pivots zugelassen werden müssen — vgl. 6.1.B —, sind die Programme von DUFF/REID [82] zu empfehlen. Für nichtsymmetrische Systeme erwähnen wir wieder das Yale Sparse Matrix Package, die von DUFF/REID [77] erarbeiteten Routinen und die darauf aufbauenden Entwicklungen von ZLATEV/WASNIEWSKI/SCHAUMBURG [81] und ØSTERBY/ZLATEV [83]. Zur Software für große Eigenwertaufgaben sei außer auf DUFF [84] noch auf die Übersichten von STEWART [76a] und PARLETT [84] sowie auf CULLUM/WILLOUGHBY [85] und die dort zitierte Literatur verwiesen, für große Quadratmittelprobleme auf die am Schluß des Kapitels 12 erwähnten Arbeiten.

Literatur

- AASEN, J. O. (71). On the reduction of a symmetric matrix to tridiagonal form. BIT 11, 233 to 242.
- BABUŠKA, I. (69). Numerical stability in numerical analysis, in: Information processing 68. North-Holland, Amsterdam, 11-23.
- BABUŠKA, I. (72). Numerical stability in problems of linear algebra. SIAM J. Numer. Anal. 9. 53-77.
- BARD, Y. (73). Nonlinear parameter estimation. Academic Press, New York. Russ. Übers.: Nelineinoe ocenivanie parametrov. Statistika, Moskva 1979.
- BARGMANN, V.; MONTGOMERY, C.; VON NEUMANN, J. (46). Solution of linear systems of high order. Institute for Advanced Study, Princeton.
- BARTELS, R. H.; STOER, J.; ZENGER, C. (71). A realization of the simplex method based on triangular decompositions, in: WILKINSON, J. H.; REINSCH, C. (Hrsg.), 152-190.
- BARWELL, V.; GEORGE, A. (76). A comparison of algorithms for solving symmetric indefinite systems of linear equations. ACM Trans. Math. Software 2, 242-251.
- BATHÉ, K. J.; WILSON, E. (76). Numerical methods in finite element analysis. Prentice-Hall, Englewood Cliffs, N. J.
- BAUER, F. L. (57). Das Verfahren der Treppeniteration und verwandte Verfahren zur Lösung algebraischer Eigenwertprobleme. Z. Angew. Math. Phys. 8, 214-235.
- BAUER, F. L. (63). Optimally scaled matrices. Numer. Math. 5, 73-87.
- BAUER, F. L. (66). Genauigkeitsfragen bei der Lösung linearer Gleichungssysteme. Z. Angew. Math. Mech. 46, 409-421.
- BAUER, F. L. (69). Remarks on optimally scaled matrices. Numer. Math. 13, 1-3.
- BAUER, F. L. (74). Computational graphs and rounding errors. SIAM J. Numer. Anal. 11, 87-96.
- BENNETT, J. M. (65). Triangular factors of modified matrices. Numer. Math. 7, 217 to 221.
- BERSENEV, S. M. (79). O peresčete faktorizacii Holeckogo. Ž. Vyčisl. Mat. i Mat. Fiz. 19, 1318-1319.
- BJÖRCK, A. (67a). Solving least squares problems by Gram-Schmidt orthogonalization. BIT 7. 1-21.
- BJÖRCK, A. (67b, 68). Iterative refinement of linear least squares solutions I, II. BIT 7, 257-278; 8, 8-30.
- BJÖRCK, A. (76). Methods for sparse least squares problems, in: BUNCH, J. R.; ROSE, D. J. (Hrsg.), 177-199.
- BJÖRCK, A. (78). Comment on the iterative refinement of least squares solutions. J. Amer. Statist. Assoc. 73, 161-166.
- BJÖRCK, A. (81). Least squares methods in physics and engineering. Report CERN 81-16, Geneva.

- BJÖRCK, A.; GOLUB, G. H. (67). Iterative refinement of linear least squares solutions by Householder transformations. BIT 7, 322-337.
- BLUE, J. L. (78). A portable FORTRAN program to find the Euclidean norm of a vector. ACM Trans. Math. Software 4, 15-23.
- BOJANCZYK, A.; BRENT, R. P.; KUNG, H. T. (84). Numerically stable solution of dense systems of linear equations using mesh-connected processors. SIAM J. Sci. Statist. Comput. 5, 95-104.
- BROWN, W. S. (81). A simple but realistic model of floating-point computation. ACM Trans. Math. Software 7, 445-480.
- BUNCH, J. R. (71). Analysis of the diagonal pivoting method. SIAM J. Numer. Anal. 8, 656-680.
- BUNCH, J. R.; KAUFMAN, L. (77). Some stable methods for calculating inertia and solving symmetric linear systems. Math. Comp. 31, 162-179.
- BUNCH, J. R.; KAUFMAN, L.; PARLETT, B. N. (76). Decomposition of a symmetric matrix. Numer. Math. 27, 95-109.
- BUNCH, J. R.; NIELSEN, C. P.; SORENSEN, D. C. (78). Rank-one modification of the symmetric eigenproblem. Numer. Math. 31, 31-48.
- BUNCH, J. R.; PARLETT, B. N. (71). Direct methods for solving symmetric indefinite systems of linear equations. SIAM J. Numer. Anal. 8, 639-655.
- BUNCH, J. R.; ROSE, D. J. (Hrsg.) (76). Sparse matrix computations. Academic Press, New York.
- BUNSE, W.; BUNSE-GERSTNER, A. (85). Numerische lineare Algebra. Teubner, Stuttgart.
- BURMEISTER, W. (76). Ein Verfahren zur Modifizierung von LU-Zerlegungen bei Spaltenersetzung. Unveröffentl. Manuskript.
- BUSINGER, P. (68). Matrices which can be optimally scaled. Numer. Math. 12, 346-348.
- BUSINGER, P.; GOLUB, G. H. (65). Linear least squares solutions by Householder transformations. Numer. Math. 7, 269-276. Auch in: WILKINSON, J. H.; REINSCH, C. (Hrsg.). 111-118.
- CHAN, T. F. (82). An improved algorithm for computing the singular value decomposition. ACM Trans. Math. Software 8, 72-83.
- CLINE, A. K.; CONN, A. R.; VAN LOAN, C. (82). Generalizing the LINPACK condition estimator, in: HENNART, J. P. (Hrsg.), Numerical analysis. Lecture Notes in Mathematics, Vol. 909. Springer, Berlin, 73-83.
- CLINE, A. K.; MOLER, C. B.; STEWART, G. W.; WILKINSON, J. H. (79). An estimate for the condition number of a matrix. SIAM J. Numer. Anal. 16, 368-375.
- CLINE, A. K.; REW, R. K. (83). A set of counter-examples to three condition number estimators. SIAM J. Sci. Statist. Comput 4, 602-611.
- COWELL, W. R. (Hrsg.) (76). Portability of numerical software. Lecture Notes in Computer Science, Vol. 57. Springer, Berlin.
- COWELL, W. R. (Hrsg.) (84). Sources and development of mathematical software. Prentice-Hall, Englewood Cliffs, N. J.
- CRAWFORD, C. R. (73). Reduction of a band-symmetric generalized eigenvalue problem. Comm. ACM 16, 41-44.
- CULLUM, J. K.; WILLOUGHBY, R. A. (85). Lanczos algorithms for large symmetric eigenvalue computations. Vol. I: Theory, Vol. II: Programs. Progress in Scientific Computing, Vol. 3, 4. Birkhäuser, Basel.
- DANIEL, J.; GRAGG, W. B.; KAUFMAN, L.; STEWART, G. W. (76). Reorthogonalization and stable algorithms for updating the Gram-Schmidt QR factorization, Math. Comp. 30, 772-795.
- DAVIS, C.; KAHAN, W. M. (70). The rotation of eigenvectors by a perturbation. III. SIAM J. Numer. Anal. 7, 1-46.
- DAX, A. (83). A diagonal modification for the downdating algorithm. SIAM J. Sci. Statist. Comput. 4, 85-93.

- DE BOOR, C.; GOLUB, G. H. (78). The numerically stable reconstruction of a Jacobi matrix from spectral data. Linear Algebra Appl. 21, 245-260.
- Döhler, B. (82). Ein neues Gradientenverfahren zur simultanen Berechnung der kleinsten und größten Eigenwerte des allgemeinen Eigenwertproblems. Numer. Math. 40, 79-91.
- DONGARRA, J. J. (84). Performance of various computers using standard linear equations software in a FORTRAN environment. Technical Memorandum MCS-TM-23, Argonne National Laboratory, Argonne.
- DONGARRA, J. J.; BUNCH, J. R.; MOLER, C. B.; STEWART, G. W. (79). LINPACK. Users guide. SIAM Publications, Philadelphia.
- DONGARRA, J. J.; GUSTAVSON, F. G.; KARP, A. (84). Implementing linear algebra algorithms for dense matrices on a vector pipeline machine. SIAM Rev. 26, 91-112.
- DONGARRA, J. J.; MOLER, C. B.; WILKINSON, J. H. (83). Improving the accuracy of computed eigenvalues and eigenvectors. SIAM J. Numer. Anal. 20, 23-45.
- DRAPER, N. R.; SMITH, H. (66). Applied regression analysis. J. Wiley, New York.
- DUFF, I. S. (Hrsg.) (81). Sparse matrices and their use. Academic Press, New York.
- DUFF, I. S. (84). A survey of sparse matrix software, in: COWELL, W. R. (Hrsg.).
- DUFF, I. S.; REID, J. K. (76). A comparison of some methods for the solution of sparse overdetermined systems of linear equations. J. Inst. Math. Appl. 17, 267-280.
- DUFF, I. S.; REID, J. K. (77). MA 28 A set of FORTRAN-subroutines for sparse unsymmetric matrices. Report R-8730, Computer Science and Systems Division, AERE Harwell.
- DUFF, I. S.; REID, J. K. (82). MA 27 A set of FORTRAN-subroutines for solving sparse symmetric sets of linear equations. Report R-10533, Computer Science and Systems Division, AERE Harwell.
- DUFF, I. S.; STEWART, G. W. (Hrsg.) (79). Sparse matrix proceedings 1978. SIAM Publishers, Philadelphia.
- EFROYMSON, M. A. (60). Multiple regression analysis, in: RALSTON, A.; WILF, H. (Hrsg.), Mathematical methods for digital computers, Vol. 1, J. Wiley, New York. Dtsch. Übers.: Mathematische Methoden für Digitalrechner, Bd. 1. Oldenbourg, München 1963.
- EISENSTAT, S. C.; GURSKY, M. C.; SCHULTZ, M. H.; SHERMAN, A. H. (77). Yale sparse matrix package. I. The symmetric codes. Research Report no. 112, II. The nonsymmetric codes. Research Report no. 114. Department of Computer Science, Yale University, New Haven.
- ELDÉN, L. (72). Stepwise regression analysis with orthogonal transformations. Report 2-1972. Department of Applied Mathematics, Linköping Institute of Technology, Linköping.
- ELDÉN, L. (77). Algorithms for the regularization of illconditioned least squares problems. BIT 17, 134-145.
- ELSNER, L.; SUN, J. (82). Perturbation theorems for the generalized eigenvalue problem. Linear Algebra Appl. 48, 341-357.
- FADDEEV, D. K.; FADDEEVA, V. N. (63). Vyčislitel'nye metody lineinoi algebry. Fizmatgiz, Moskva. Dtsch. Übers.: Numerische Methoden der linearen Algebra. 5. Auflage, VEB Deutscher Verlag der Wissenschaften, Berlin; Oldenbourg, München-Wien 1978.
- FADDEEV, D. K.; FADDEEVA, V. N. (69). Stability in linear algebra problems, in: Information processing 68. North-Holland, Amsterdam, 33-39.
- FADDEEV, D. K.; FADDEEVA, V. N. (70). Natural norms in algebraic processes. SIAM J. Numer. Anal. 7, 520-531.
- FADDEEV, D. K.; KUBLANOVSKAJA, V. N.; FADDEEVA, V. N. (68). Sur les Systèmes Linéaires Algébriques de Matrices Rectangularies et Mal-Conditionnées, in: Programmation en Mathématiques Numériques VII. Editions Centre Nat. Recherche Sci., Paris, 161-170.
- FADDEEVA, V. N.; IKRAMOV, H. D. et al. (82). Vyčislitel'nye metody lineňnoĭ algebry. Bibliografičeskiĭ ukasatel' 1975-1980. Leningradskoe Otdelenie Mat. Inst. Steklova, Leningrad.
- FADDEEVA, V. N.; KUZNECOV, JU. A. et al. (76). Vyčislitel'nye metody lineĭnoĭ algebry. Bibliografičeskiĭ ukasatel' 1828—1974. Leningradskoe Otdelenie Mat. Inst. Steklova, Vyčisl. Centr Sibirskogo Otdelenija Akad. Nauk SSSR, Novosibirsk.

- FENNER, T. I.; LOIZOU, G. (77). Optimally scalable matrices. Philos. Trans. Roy. Soc. London Ser. A 287, no. 1345, 307-349.
- FIX, G.; HEIBERGER, R. (72). An algorithm for the ill-conditioned generalized eigenvalue problem. SIAM J. Numer. Anal. 9, 78-88.
- FLETCHER, R.; MATTHEWS, S. P. J. (83). A stable algorithm for updating triangular factors under a rank one change. Report NA/69. Department of Mathematical Sciences, University of Dundee, Dundee.
- FLETCHER, R.; POWELL, M. J. D. (74). On the modification of LDL^{\intercal} factorizations. Math. Comp. 29, 1067-1087.
- FORSYTHE, G. E.; MOLER, C. B. (67). Computer solutions of linear algebraic systems. Prentice Hall, Englewood Cliffs, N. J. Dtsch. Übers.: Computer-Verfahren für lineare algebraische Systeme. Oldenbourg, München 1971. Russ. Übers.: Čislennoe rešenie sistem lineinyh algebraičeskyh uravnenii. Mir, Moskva, 1969.
- FORSYTHE, G. E.; STRAUS, E. G. (55). On best conditioned matrices. Proc. Amer. Math. Soc. 6, 340-345.
- FOSDICK, L. (Hrsg.) (79). Performance evaluation of numerical software. North-Holland, Amsterdam.
- FRANCIS, J. G. F. (61). The QR transformation a unitary analogue to the LR transformation. Comput. J. 4, 265–271, 332–345.
- FRIEDLAND, S. (75). On inverse multiplicative eigenvalue problems for matrices. Linear Algebra Appl. 12, 127-138.
- FRIEDRICH, V.; HOFMANN, B.; TAUTENHAHN, U. (79). Möglichkeiten der Regularisierung bei der Auswertung von Meßdaten. Schriftenreihe der TH Karl-Marx-Stadt 10, 1-42.
- GANTMAHER, F. R. (66). Teorija matric. Izd. vtoroe. Nauka, Moskva. Dtsch. Übers.: Matrizentheorie. VEB Deutscher Verlag der Wissenschaften, Berlin 1986.
- GARBOW, B. S.; BOYLE, J. M.; DONGARRA, J. J.; MOLER, C. B. (77). Matrix eigensystem routines — EISPACK guide extension. Lecture Notes in Computer Science, Vol. 51. Springer, Berlin.
- GENTLEMAN, W. M. (73). Least squares computations by Givens transformations without square roots. J. Inst. Math. Appl. 12, 329-336.
- GENTLEMAN, W. M. (75). Error analysis of QR decompositions by Givens transformations. Linear Algebra Appl. 10, 189-197.
- GENTLEMAN, W. M.; MAROVICH, S. B. (74). More on algorithms that reveal properties of floating point arithmetic units. Comm. ACM 17, 276-277.
- GEORGE, A.; LIU, J. W. H. (81). Computer solution of large positive definite systems. Prentice Hall, Englewood Cliffs, N. J. Russ. Übers.: Čislennoe rešenie bol'ših razrežennyh sistem uravnenii. Mir, Moskva 1984.
- GEORGE, A.; HEATH, M. T.; NG, E. (83). A comparison of some methods for solving sparse linear least squares problems. SIAM J. Sci. Statist. Comput. 4, 177-187.
- GEORGE, A.; HEATH, M. T.; PLEMMONS, R. J. (81). Solution of large-scale sparse least squares problems using auxiliary storage. SIAM J. Sci. Statist. Comput. 2, 416-429.
- GERADIN, M. (71). The computational efficiency of a new minimization algorithm for eigenvalue analysis. J. Sound Vibration 19, 319-331.
- GILL, P. E.; GOLUB, G. H.; MURRAY, W.; SAUNDERS, M. A. (74). Methods for modifying matrix factorizations. Math. Comp. 28, 505-535.
- GILL, P. E.; MURRAY, W. (77). Modification of matrix factorizations after a rank-one change, in: JACOBS, D. E. (Hrsg.), 55-83.
- GILL, P. E.; MURRAY, W.; SAUNDERS, M. A. (75). Methods for computing and modifying the LDV factors of a matrix. Math. Comp. 29, 1051-1077.
- GILL, P. E.; MURRAY, W.; WRIGHT, M. H. (82). Practical optimization. Academic Press, New York.
- GIVENS, W. (53). A method for computing eigenvalues and eigenvectors suggested by classical results on symmetric matrices. Nat. Bur. Standards Appl. Math. Ser. 29, 117-122.
- GIVENS, W. (54). Numerical computation of the characteristic values of a real symmetric matrix. Report ORNL-1574, Oak Ridge National Laboratory, Oak Ridge.

- GIVENS, W. (58). Computation of plane unitary rotations transforming a general matrix to triangular form. SIAM J. Appl. Math. 6, 26-50.
- GOLDSTINE, H. H.; MURRAY, F. J.; VON NEUMANN, J. (59). The Jacobi method for real symmetric matrices. J. Assoc. Comput. Mach. 6, 59-96.
- GOLUB, G. H. (65). Numerical methods for solving linear least squares problems. Numer. Math. 7, 206-216.
- GOLUB, G. H. (69). Matrix decompositions and statistical computation, in: MILTON, R. C.; NELDER, J. A. (Hrsg.), Statistical computation. Academic Press, New York.
- GOLUB, G. H. (73). Some modified matrix eigenvalue problems. SIAM Rev. 15, 318-334.
- GOLUB, G. H.; HEATH, M.; WAHBA, G. (79). Generalized cross-validation as a method for choosing a good ridge parameter. Technometrics 21, 215-223.
- GOLUB, G. H.; KAHAN, W. (65). Calculating the singular values and pseudoinverse of a matrix. SIAM J. Numer. Anal. 2, 205-224.
- GOLUB, G. H.; KAHAN, W. (68). Calculating the singular values and pseudoinverse of a matrix. Appl. Mat. 13, 44-51.
- GOLUB, G. H.; REINSCH, C. (70). Singular value decomposition and least squares solutions. Numer. Math. 14, 403-420. Auch in: WILKINSON, J. H.; REINSCH, C. (Hrsg.), 134-151.
- GOLUB, G. H.; STYAN, G. P. (73). Numerical computations for univariate linear models. J. Statist. Comput. Simulation 2, 253-274.
- GOLUB, G. H.; VAN LOAN, C. F. (80). An analysis of the total least squares problem. SIAM J. Numer. Anal. 17, 883-893.
- GOLUB, G. H.; VAN LOAN, C. F. (83). Matrix computations. North Oxford Academic, Oxford.
- GOLUB, G. H.; WILKINSON, J. H. (66). Note on the iterative refinement of least squares solution. Numer. Math. 9, 139-148.
- GREGORY, R. T.; KRISHNAMURTHY, E. V. (84). Methods and applications of error-free computation. Springer, Berlin.
- GRIMES, R. G.; LEWIS, J. G. (81). Condition number estimation for sparse matrices. SIAM J. Sci. Statist Comput. 2, 384-388.
- HACKBUSCH, W.; TROTTENBERG, U. (Hrsg.) (82). Multigrid methods. Lecture Notes in Mathematics, Vol. 960. Springer, Berlin.
- HAGEMAN, L. A.; YOUNG, D. M. (81). Applied iterative methods. Academic Press, New York. Russ. Übers.: Prikladnye iteracionnye metody. Mir, Moskva 1986.
- HAMMARLING, S. (74). A note on modifications of Givens plane rotation. J. Inst. Math. Appl. 13, 215-218.
- HASKELL, K. H.; HANSON, R. J. (81). An algorithm for linear least squares problems with equality and nonnegativity constraints. Math. Programming 21, 98-118.
- HEATH, M. T. (83). Numerical methods for large sparse linear least squares problems. Report ORNL/CSD-114, Oak Ridge National Laboratory, Oak Ridge.
- HEINRICH, H. (63). Bemerkungen zu einem Konditionsmaß für lineare Gleichungssysteme. Z. Angew. Math. Mech. 43, 568.
- HELLER, D. E.; IPSEN, I. C. F. (83). Systolic networks for orthogonal decompositions. SIAM J. Sci. Statist. Comput. 4, 261-269.
- HELLER, P. (78). A survey of parallel algorithms in numerical linear algebra. SIAM Rev. 20, 740-777.
- HENRICI, P. (58). On the speed of convergence of cyclic and quasicyclic Jacobi methods for computing eigenvalues of Hermitian matrices. J. Soc. Industr. Appl. Math. 6, 144-162.
- HOFFMANN, W.; PARLETT, B. N. (78). A new proof of global convergence for the tridiagonal QL algorithm. SIAM J. Numer. Anal. 15, 929-937.
- HOFMANN, B. (86). Regularization for applied inverse and ill-posed problems. Teubner, Leipzig.
- HOUSEHOLDER, A. S. (58). Unitary triangularization of a nonsymmetric matrix. J. Assoc. Comput. Mach. 5, 339-342.
- HOUSEHOLDER, A. S. (64). The theory of matrices in numerical analysis. Blaisdell, New York. HOUSEHOLDER, A. S.; BAUER, F. L. (59). On certain methods for expanding the characteristic
 - polynomial. Numer. Math. 1, 29-37.

- HUMAK, K. M. S. (77/83). Statistische Methoden der Modellbildung. Bd. I/II. Akademie-Verlag, Berlin.
- IKRAMOV, H. D. (82). Razrežennye matricy. Itogi Nauki i Tehniki. VINITI, Ser. Mat. Analiz, 179-257.
- IKRAMOV, H. D. (84). Čislennoe rešenie matričnyh uravneniĭ. Ortogonal'nye metody. Nauka, Moskva.
- JACOBI, C. G. J. (1846). Über ein leichtes Verfahren, die in der Theorie der Säcularstörungen vorkommenden Gleichungen numerisch aufzulösen. J. Reine Angew. Math. 30, 51–95.
- JACOBS, D. E. (Hrsg.) (77). The state of the art in numerical analysis. Academic Press, London.
- JACOBS, D. A. H. (Hrsg.) (78). Numerical software needs and availability. Academic Press, London.
- JANKOWSKI, M.; WOŹNIAKOWSKI, H. (77). Iterative refinement implies numerical stability. BIT 17, 303-311.
- JENNINGS, A. (77). Matrix computations for engineers and scientists. J. Wiley, London.
- KÅGSTRÖM, B.; RUHE, A. (Hrsg.) (83). Matrix pencils. Lecture Notes in Mathematics, Vol. 973. Springer, Berlin.
- KAHAN, W. (66). Numerical linear algebra. Canad. Math. Bull. 9, 757-801.
- KAHAN, W. (71). A survey of error analysis, in: Information processing 71. North-Holland, Amsterdam, 1214-1239.
- KIEŁBASIŃSKI, A. (78). Basic concepts in numerical error analysis, in: Mathematical models and numerical methods. Banach Center Publications. Vol. 3. PWN, Warszawa.
- Korov, V. J.; MIKLOŠKO, J. (Hrsg.) (84). Algorithms, software, and hardware of parallel computers, VEDA, Bratislava; Springer, Berlin.
- KUBLANOVSKAJA, V. N. (61). O nekotoryh algoritmah dlja rešenija polnoĭ problemy sobstvennyh značenił, Ž. Vyčisl. Mat. i Mat. Fiz. 1, 555-570.
- KUBLANOVSKAJA, V. N. (78a). K rešeniju spektral'nož zadači dlja singuljarnogo pučka matric. Ž. Vyčisl. Mat. i Mat. Fiz. 18, 1056-1060.
- KUBLANOVSKAJA, V. N. (78b). Primenenie normalizovannogo processa k postroeniju algoritmov rešenija spektral'nyh zadač dlja pučkov matric, in: Proceedings of the Fourth Sympos. on Basic Problems of Numerical Mathematics. Prague.
- KUZNECOV, JU. A. (83). Metod soprjažennyh gradientov, ego obobščenija i primenenija, in: MARČUK, G. I. (Hrsg.), 267-301.
- LANCZOS, C. (50). An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. J. Res. Nat. Bur. Standards, Sect. B 45, 255-282.
- LANG, W. et al. (79). Dokumentation zum Programmsystem MEIWEP. Teil 1: Dokumentation des Gesamtsystems. Teil 2: Dokumentation der Moduln (Heft 1: Symmetrische Matrizen, Heft 2: Reell nichtsymmetrische Matrizen, Heft 3: Komplexe Matrizen). TH Karl-Marx-Stadt, Sektion Mathematik, Karl-Marx-Stadt.
- LANG, W. (82). Ein Programmsystem für Matrizeneigenwertaufgaben. Akademie-Verlag, Berlin.
- LAWSON, C. L.; HANSON, R. J. (74). Solving least squares problems. Prentice-Hall, Englewood Cliffs, N. J. Russ. Übers.: Čislennoe rešenie zadač metoda naimen'ših kvadratov. Mir, Moskva 1986.
- LAWSON, C. L.; HANSON, R. J.; KINCAID, D. R.; KROGH, F. T. (79). Basic linear algebra subprograms for FORTRAN usage. ACM Trans. Math. Software 5, 308-323.
- LEMEIRE, F. (75). Bounds for condition numbers of triangular and trapezoid matrices. BIT 15, 58-64.
- LEVENBERG, K. (44). A method for the solution of certain nonlinear problems in least squares. Quart. Appl. Math. 2, 164-168.
- LINNIK, I. V. (62). Metod naimen'ših kvadratov i osnovy matematiko-statističeskoĭ teorii

obrabotki nabljudeniĭ. Fizmatgiz, Moskva. Dtsch. Übers. d. 1. Auflage: Die Methode der kleinsten Quadrate in moderner Darstellung. Deutscher Verlag der Wissenschaften, Berlin 1961.

- LÖTSTEDT, P. (84). Solving the minimal least squares problem subject to bounds on the variables. BIT 24, 206-224.
- LONGSINE, D. E.; MCCORMICK, S. F. (80). Simultaneous Rayleigh-quotient minimization methods for $Ax = \lambda Bx$. Linear Algebra Appl. 34, 195–234.
- MAESS, G. (84). Vorlesungen über Numerische Mathematik. Bd. 1. Akademie-Verlag, Berlin.
- MALCOLM, M. (72). Algorithms to reveal properties of floating-point arithmetic. Comm. ACM 15, 949-951.
- MARČUK, G. I. (Hrsg.) (83). Vyčislitel'nye processy i sistemy. Vyp. 1. Nauka, Moskva.
- MARKOWITZ, H. M. (57). The elimination form of the inverse and its application to linear programming. Management Sci. 3, 255-269.
- MCCORMICK, S. F.; NOE, T. (77). Simultaneous iteration for the matrix eigenvalue problem. Linear Algebra Appl. 16, 43-56.
- MEYER, A. (87). Modern algorithms for large sparse eigenvalues problems. Akademie-Verlag, Berlin.
- MILLER, W. (75). Computational complexity and numerical stability. SIAM J. Comput. 4, 97-107.
- MILLER, W.; WRATHALL, C. (80). Software for roundoff analysis of matrix algorithms. Academic Press, New York.
- MOLER, C. B. (67). Iterative refinement in floating point. J. Assoc. Comput. Mach. 14, 316-371.
- MOLER, C. B.; STEWART, G. W. (73). An algorithm for generalized matrix eigenvalue problems. SIAM J. Numer. Anal. 10, 241-256.
- MOLER, C. B.; VAN LOAN, C. F. (78). Nineteen dubious ways to compute the exponential of a matrix. SIAM Rev. 20, 801-836.
- MOORE, R. E. (79). Methods and applications of interval analysis. SIAM Publications, Philadelphia.
- Morozov, V. A. (66). O reguljarizacii nekorrektno postavlennyh zadač i vybore parametra reguljarizacii. Ž. Vyčisl. Mat. i Mat. Fiz. 6, 170–175.
- Morozov, V. A. (73). O principe nevjaski pri rešenii nesovmestnyh uravnenii metodom reguljarizacii A. N. Tyhonova. Ž. Vyčisl. Mat. i Mat. Fiz. 13, 1099-1111.
- MOROZOV, V. A. (84). Methods for solving incorrectly posed problems. Springer, Berlin.
- NICKEL, K. (Hrsg.) (80). Interval mathematics 1980. Academic Press, New York.
- O'LEARY, D. P. (80). Estimating matrix condition numbers. SIAM J. Sci. Statist. Comput. 1, 205-209.
- OSBORNE, M. R. (76). On the computation of stepwise regression. Austral. Comput. J. 8, 61-68.
- ØSTERBY, O.; ZLATEV, Z. (83). Direct methods for sparse matrices. Lecture Notes in Computer Science, Vol. 157. Springer, Berlin. Russ. Übers.: Prjamye metody dlja razrežennyh matric. Mir, Moskva 1987.
- OSTROWSKI, A. M. (58/59). On the convergence of the Rayleigh quotient iteration for the computation of characteristic roots and vectors. I-VI. Arch. Rat. Mech. Anal. 1 (1958), 233-241; 2 (1959), 423-428; 3 (1959), 325-340, 341-347, 472-481; 4 (1959), 153 to 165.
- PAIGE, C. C. (76). Error analysis of the Lanczos algorithm for tridiagonalizing a symmetric matrix. J. Inst. Math. Appl. 18, 341-349.
- PAIGE, C. C. (79a). Computer solution and perturbation analysis of generalized linear least squares problems. Math. Comp. 33, 171-184.
- PAIGE, C. C. (79b). Fast numerically stable computations for generalized linear least squares problems. SIAM J. Numer. Anal. 16, 165-171.

- PAIGE, C. C. (80). Error analysis of some techniques for updating orthogonal decompositions. Math. Comp. 34, 465-471.
- PARLETT, B. N. (64). The origin and development of methods of LR type. SIAM Rev. 6, 275 to 295.
- PARLETT, B. N. (74). The Rayleigh quotient iteration and some generalizations for nonnormal matrices. Math. Comp. 28, 679-693.
- PARLETT, B. N. (80a). The symmetric eigenvalue problem. Prentice-Hall, Englewood Cliffs, N. J. Russ. Übers.: Simmetričnaja problema sobstvennyh značeniĭ. Čislennye metody. Mir, Moskva 1983.
- PARLETT, B. N. (80b). How to solve $(K \lambda M) z = 0$ for large K and M. Numer. Methods Engrg. 1, 97-106.
- PARLETT, B. N. (84). The software scene in the extraction of eigenvalues from sparse matrices. SIAM J. Sci. Statist. Comput. 5, 590-604.
- PARLETT, B. N.; KAHAN, W. (69). On the convergence of a practical QR algorithm, in: Information processing 68. North-Holland, Amsterdam, 114-118.
- PETERS, G.; WILKINSON, J. H. (69). Eigenvalues of $Ax = \lambda Bx$ with band symmetric A and B. Comput. J. 12, 398-404.
- PETERS, G.; WILKINSON, J. H. (70a). The least squares problem and pseudoinverses. Comput. J. 13, 309-316.
- PETERS, G.; WILKINSON, J. H. (70b). $Ax = \lambda Bx$ and the generalized eigenproblem. SIAM J. Numer. Anal. 7, 479-492.
- PETERS, G.; WILKINSON, J. H. (75). On the stability of Gauss-Jordan elimination with pivoting. Comm. ACM 18, 20-24.
- PETERS, G.; WILKINSON, J. H. (79). Inverse iteration, ill-conditioned equations, and Newton's method. SIAM Rev. 21, 339-360.
- PISSANETZKY, S. (84). Sparse matrix technology. Academic Press, New York.
- POPE, D. A.; TOMPKINS, C. (57). Maximizing functions of rotations experiments concerning speed of diagonalization of symmetric matrices using Jacobi's method. J. Assoc. Comput. Mach. 4, 459-466.
- RATH, W. (82). Fast Givens rotations for orthogonal similarity transformations. Numer. Math. 40, 47-56.
- RICE, J. R. (66). Experiments on Gram-Schmidt orthogonalization. Math. Comp. 20, 325-328.
- RICE, J. R. (Hrsg.) (77). Mathematical software III. Academic Press, New York.
- RICE, J. R. (81). Matrix computations and mathematical software. McGraw-Hill, New York. Russ. Übers.: Matričnye vyčislenija i matematičeskoe obespečenie. Mir, Moskva 1984.
- RODRIGUE, G. (Hrsg.) (82). Parallel computations. Academic Press, New York.
- RUTISHAUSER, H. (58). Solution of eigenvalue problems with the LR transformation. Nat. Bur. Standards Appl. Math. Ser. 49, 47-81.
- RUTISHAUSER, H. (66). The Jacobi method for real symmetric matrices. Numer. Math. 9, 1-10. Auch in: WILKINSON, J. H.; REINSCH, C. (Hrsg.), 202-211.
- RUTISHAUSER, H. (69). Computational aspects of F. L. Bauer's simultaneous iteration method. Numer. Math. 13, 4-13.
- RUTISHAUSER, H. (70). Simultaneous iteration method for symmetric matrices. Numer. Math. 16, 205-223. Auch in: WILKINSON, J. H.; REINSCH, C. (Hrsg.), 202-211.
- SAMARSKII, A. A.; NIKOLAEV, E. S. (78). Metody rešenija setočnyh uravnenil. Nauka, Moskva.
- SCHITTKOWSKI, K.; STOER, J. (79). A factorization method for the solution of constrained linear least squares problems allowing for subsequent data changes. Numer. Math. 31, 431 to 463.
- SCHÖNHAGE, A. (64). Zur quadratischen Konvergenz des Jacobi-Verfahrens. Numer. Math. 6, 410-412.
- SCHWARZ, H. R. (77). Two algorithms for treating $Ax = \lambda Bx$. Comput. Methods Appl. Mech. Engrg. 12, 181–199.

- SCHWARZ, H. R. (79). Zur Eigenwertaufgabe $Ax = \lambda Bx$, in: Numerische Behandlung von Eigenwertaufgaben. ISNM 43. Birkhäuser, Basel, 161–175.
- SCHWARZ, H. R. (80). Methode der finiten Elemente. Teubner, Stuttgart.
- SCHWETLICK, H. (79). Numerische Lösung nichtlinearer Gleichungen. VEB Deutscher Verlag der Wissenschaften, Berlin; Oldenbourg, München.
- SCHWETLICK, H.; TILLER, V. (85). Numerical methods for estimating parameters in nonlinear models with errors in the variables. Technometrics 27, 17-24.
- SCOTT, D. S.; WARD, R. C. (82). Solving symmetric-definite quadratic λ -matrix problems without factorization. SIAM J. Sci. Statist. Comput. 3, 58-67.
- SKEEL, R. D. (79). Scaling for numerical stability in Gaussian elimination. J. Assoc. Comput. Mach. 26, 494-526.
- SKEEL, R. D. (80). Iterative refinement implies numerical stability for Gaussian elimination. Math. Comp. 35, 817-832.
- SKEEL, R. D. (81). Effect of equilibration on residual size for partial pivoting. SIAM J. Numer. Anal. 18, 449-454.
- SMITH, B. T.; BOYLE, J. M.; DONGARRA, J. J.; GARBOW, B. S.; IKEBE, Y.; KLEMA, V. C.; MOLER, C. B. (74). Matrix eigensystem routines — EISPACK guide. Lecture Notes in Computer Science, Vol. 6. Springer, Berlin. 2nd ed. 1976.
- SORENSEN, D. C. (77). Updating the symmetric indefinite factorization with applications in a modified Newton's method. Report ANL-77-49, Argonne National Laboratory. Argonne.
- STERBENZ, P. H. (74). Floating-point computation. Prentice-Hall, Englewood Cliffs, N. J.
- STEWART, G. W. (69). Accelerating the orthogonal iteration for the eigenvalues of a Hermitian matrix. Numer. Math. 13, 362-376.
- STEWART, G. W. (72). On the sensitivity of the eigenvalue problem $Ax = \lambda B\dot{x}$. SIAM J. Numer. Anal. 9, 669–686.
- STEWART, G. W. (73a). Error and perturbation bounds for subspaces associated with certain eigenvalue problems. SIAM Rev. 15, 727-764.
- STEWART, G. W. (73b). Introduction to matrix computations. Academic Press, New York.
- STEWART, G. W. (75). Methods of simultaneous iteration for calculating eigenvectors of matrices, in: MILLER, J. H. (Hrsg.), Topics in numerical analysis II. Academic Press, New York, 185-196.
- STEWART, G. W. (76a). A bibliographic tour of the large, sparse generalized eigenvalue problem, in: BUNCH, J. R.; ROSE, D. J. (Hrsg.), 113-130.
- STEWART, G. W. (76b). The economical storage of plane rotations. Numer. Math. 25, 137 to 138.
- STEWART, G. W. (77). On the perturbation of pseudo-inverses, projections, and linear least squares problems. SIAM Rev. 19, 634-662.
- STEWART, G. W. (78). Perturbation theory for the generalized eigenvalue problem, in: DE BOOR, C.; GOLUB, G. H. (Hrsg.), Recent advances in numerical analysis. Academic Press, New York, 69-86.
- STEWART, G. W. (79a). Perturbation bounds for the definite generalized eigenvalue problem. Linear Algebra Appl. 23, 69-86.
- STEWART, G. W. (79b). The effects of rounding error on an algorithm for downdating a Cholesky factorization. J. Inst. Math. Appl. 23, 203-213.
- STEWART, G. W. (80). The efficient generation of random orthogonal matrices with an application to condition estimators. SIAM J. Numer. Anal. 17, 403-409.
- STIEFEL, E. (61). Einführung in die Numerische Mathematik. Teubner, Stuttgart.
- STOER, J. (71). On the numerical solution of constrained least squares problems. SIAM J. Numer. Anal. 8, 382-411.
- STOER, J. (83). Einführung in die Numerische Mathematik. Bd. 1. 4., verb. Auflage. Springer, Berlin.
- STRASSEN, V. (69). Gaussian elimination is not optimal. Numer. Math. 13, 354-356.
- SUN, J. (83). Perturbation analysis for the generalized eigenvalue and the generalized singular value problem, in: KÅGSTRÖM, B.; RUHE, A. (Hrsg.), 221-224.

- TURING, A. M. (48). Rounding-off errors in matrix processes. Quart. J. Mech. Appl. Math. 1, 287-308.
- TYHONOV, A. N. (63). O reguljarizacii nekorrektno postavlennyh zadač. Dokl. Akad. Nauk SSSR 153, 49-52.
- TYHONOV, A. N.; ARSENIN, V. JA. (79). Metody rešenija nekorrektnyh zadač. Izd. vtoroe. Nauka, Moskva.
- TYHONOV, A. N.; GONČARSKII, A. V.; STEPANOV, V. V.; JAGOLA, A. G. (83). Reguljarizirujuščie algoritmy i apriornaja informacija. Nauka, Moskva.
- UNGER, H. (50). Nichtlineare Behandlung von Eigenwertaufgaben. Z. Angew. Math. Mech. 30, 281-282.
- VAN DER SLUIS, A. (69). Condition numbers and equilibration of matrices. Numer. Math. 14, 14-23.

VAN DER SLUIS, A. (70). Condition, equilibration, and pivoting in linear algebraic systems. Numer. Math. 15, 74-86.

VARGA, R. S. (62). Matrix iterative analysis. Prentice-Hall, Englewood Cliffs, N. J.

VOEVODIN, V. V. (69a). O metode reguljarizacii. Ž. Vyčisl. Mat. i Mat. Fiz. 9, 671-673.

VOEVODIN, V. V. (69b). Ošibki okruglenija i ustoičivost' v prjamyh metodah lineinoi algebry. Izd. MGU, Moskva.

VOEVODIN, V. V. (77). Vyčislitel'nye osnovy lineinoi algebry. Nauka, Moskva.

- VOEVODIN, V. V. (80). Lineinaja algebra. Nauka, Moskva. Engl. Übers.: Linear algebra. Mir, Moscow 1983.
- VOEVODIN, V. V.; KUZNECOV, JU. A. (84). Matricy i vyčislenija. Nauka, Moskva.
- VOEVODIN, V. V.; TYRTYŠNIKOV, E. E. (83). Vyčislenija s teplicevymi matricami, in: MARČUK, G. I. (Hrsg.), 124–266.
- von MISES, R.; POLLACZEK-GEIRINGER, H. (29). Praktische Verfahren der Gleichungsauflösung. Z. Angew. Math. Mech. 9, 152-164.
- VON NEUMANN, J.; GOLDSTINE, H. H. (47). Numerical inverting of matrices of high order. Bull. Amer. Math. Soc. 53, 1021-1099.
- WALDVOGEL, P. (82). Bisection for $Ax = \lambda Bx$ with matrices of variable band width. Computing 28, 171-180.
- WANG, J. Y.; GARBOW, B. S. (83). A numerical method for solving inverse real symmetric eigenvalue problems. SIAM J. Sci. Statist. Comput. 4, 45-51.
- WATKINS, D. S. (82). Understanding the QR algorithm. SIAM Rev. 24, 427-440.
- WEDIN, P. A. (73). Perturbation theory for pseudo-inverses. BIT 13, 217-232.
- WILKINSON, J. H. (61). Error analysis of direct methods of matrix inversion. J. Assoc. Comput. Mach. 8, 281-330.
- WILKINSON, J. H. (62). Note on the quadratic convergence of the cyclic Jacobi process. Numer. Math. 4, 296-300.
- WILKINSON, J. H. (63). Rounding errors in algebraic processes. Prentice-Hall, Englewood Cliffs, N. J. Dtsch. Übers.: Rundungsfehler. Springer, Berlin 1969. Russ. Übers.: Algebraičeskaja problema sobstvennyh značenii. Nauka, Moskva 1970.
- WILKINSON, J. H. (65). The algebraic eigenvalue problem. Clarendon Press, Oxford.
- WILKINSON, J. H. (68). Global convergence of tridiagonal QR algorithm with origin shifts. Linear Algebra Appl. 1, 409-420.
- WILKINSON, J. H. (71). Modern error analysis. SIAM Rev. 13, 548-568.
- WILKINSON, J. H. (77). Some recent advances in numerical linear algebra, in: JACOBS, D. E. (Hrsg.), 3-53.
- WILKINSON, J. H.; REINSCH, C. (Hrsg.) (71). Linear algebra. Springer, Berlin. Russ. Übers.: Lineĭnaja algebra. Spravočnik algoritmov na jazyke ALGOL. Izd. Mašinostroenie, Moskva 1976.

YOUNG, D. M. (71). Iterative solution of large linear systems. Academic Press, New York.

- ZIELKE, G. (70). Numerische Berechnung von benachbarten inversen Matrizen und linearen Gleichungssystemen. Vieweg, Braunschweig.
- ZIELKE, G. (78). Zur historischen Entwicklung von verallgemeinerten inversen Matrizen. Wiss. Z. Martin-Luther-Univ. Halle-Wittenberg, Math.-Natur. Reihe 4, 109-118.
- ZIELKE, G. (86). Report on test matrices for generalized inverses. Computing 36, 105-162.
- ZLATEV, Z.; WASNIEWSKI, J.; SCHAUMBURG, K. (81). Y 12 M. Solution of large and sparse systems of linear algebraic equations. Documentation of subroutines. Lecture Notes in Computer Science, Vol. 121. Springer, Berlin.

Sachverzeichnis

Abbildung 20 -, identische 22 -, inverse 21 -. lineare 21 absolute Norm 30 Ähnlichkeitstransformation 35, 354 Akkumulation in höherer Genauigkeit 84 Äquilibriertheit 134 Äquivalenz numerischer Probleme 96 Äquivalenztransformation 41 Aufdatierung einer LR-Faktorisierung 214 - einer LL^T-Faktorisierung 217 f. - einer inversen Matrix 235f. - einer **QR**-Faktorisierung 301 ff. Auffüllung 227 Aufwand 74 Ausgangsdaten 50 Auslöschung 70

Bandbreite 15, 221 Bandmatrix 15, 221 Basis einer Computerzahl 64 - eines Teilraumes 14 Bauer-Fike-Theorem 356 Betragsmatrix 29 Betragsvektor 29 Bidiagonalisierung 316ff. -, modifizierte, für $m \gg n$ 324 **Bidiagonalmatrix 15** Bisektionsverfahren 411 ff. Bit 65 BKP-Faktorisierung \rightarrow Bunch-Kaufman-Parlett-Faktorisierung BLAS - Basic Linear Algebra Subprograms 453 Blockelimination 191 ff. Blockmatrix 12 **B**-Norm 437 **B**-Orthogonalität 436

B-Skalarprodukt 437 Bunch-Kaufman-Parlett-Faktorisierung 197 Büschel 434 -, äquivalente 435 –, kongruente 436 Byte 65 Charakteristisches Polynom eines Büschels 435 — — einer Matrix 32, 353 Cholesky-Faktorisierung $\rightarrow LL^{T}$ -Faktorisierung Computerarithmetik 69ff. Computerzahlen 64ff. Darstellungsfehler einer Matrix 85 — eines Vektors 85 - einer Zahl 67 Deflation bei **QR**-Algorithmus 420 - bei SVD-Berechnung 320 - bei Vektoriteration 381 Determinanten 22f. -; Berechnung 47, 160, 163 -; Entwicklungssatz 23 Diagonalmatrix 15 Diagonalpivotisierung 168, 189 Dimension einer linearen Mannigfaltigkeit 242 — eines Teilraumes 14 direkte Verfahren 96, 142 dominanter Eigenwert 374, 441, 444 Teilraum 374 doppelte Genauigkeit 65 Drehungsmatrix 112 $Dreiecksfaktorisierung \rightarrow LR$ -Faktorisierung Dreiecksmatrix 15 Dreiecksungleichung 25 dyadisches Produkt 19

B-Selbstadjungiertheit 447

Eigenwerte 31, 434

Eigenwerte, konjugiert-komplexe 33 Eigenpaar 32, 434 Eigenvektor 32, 434 Eigenwerthaufen 358, 372, 387, 394t. Eigenwertproblem, allgemeines symmetrisches 49, 434 -. inverses 449 -, spezielles nichtsymmetrisches 49, 54 f. -, - symmetrisches 49, 353 **Eigenwertzerlegung 36** einfache Genauigkeit 65 Eingangsdaten 50 Einhüllende einer Matrix 226 Einsdreiecksmatrix 145 Eliminationsschritt mittels Givens-Drehungen 114 - mittels impliziter Givens-Drehungen 121 - mittels Householder-Spiegelung 110 - mittels LNT-Matrix 103 mittels stabilisierter LNT-Matrix 105 Euklidische Norm 25 Exponent 63ff. Exponentenüberlauf 66 Exponentenunterlauf 66 Fehler, absoluter 51

-, relativer 56
-, unvermeidlicher 53, 91
Fehlerabschätzung bei linearen Gleichungssystemen 174 ff.
Fehleranalyse 79, 90 ff.
Fehlerkumulationskonstante 79, 90 ff.
Fehlerniveau 51 ff.
-, optimales 52, 91
Frobeniusnorm 27

Gauß-Jordan-Verfahren 234 Gaußsche Dreiecksfaktorisierung 158ff. $\rightarrow LR$ -Faktorisierung -r Algorithmus 143, 153ff. -r -, verketteter 203 Givens-Drehungen 111ff. -, implizite 116ff. Givens-Orthogonalisierung 290ff. -, implizite 291 Givens-Spiegelungen 125 Gleitpunktzahlen 64 -, normalisierte 64 Gradientenverfahren 446 Gram-Schmidt-Orthogonalisierung 277 ff. -; Instabilität 279 -, modifizierte 281 ff.

Hauptachsentransformation 40 Hessenbergmatrix 16 Hölderstetigkeit 54 Householder-Orthogonalisierung 288ff. – mit Pivotisierung 292 Householder-Spiegelungen 106ff.

Inkorrekt gestelltes Problem 59 inverse Iteration 389ff., 443ff. - Matrix 21 --; Berechnung 46f., 229ff. iterative Verbesserung bei Gleichungssystemen 176ff. - - bei Quadratmittelproblemen 275ff., 299ff., 301, 315 – Verfahren 96 - - für Eigenwertprobleme 373ff., 441ff. — – für lineare Gleichungssysteme 239 Jacobi-Verfahren 365ff. Jordanblock 37 Jordansche Normalform 37 Kaskadensummation 93 Kompaktspeicherung 228 Konditionsschätzer 178ff. Konditionszahl, absolute 56 einer Matrix 130 -, partielle 58 -, relative 57 Kongruenztransformation 40 Konsistenz 241 Koordinatenvektor 13 korrekt gestelltes Problem 53 Kreisesatz von Geršgorin 39

Lanczos-Algorithmus 405f., 434, 447 *LDL*^T-Faktorisierung, direkte 205 — indefiniter Matrizen 190ff.

- positiv definiter Matrizen 185

- zum Schneiden des Spektrums 407 ff.

- lineare Abhängigkeit 14
- -- Gleichungssysteme mit Dreiecksmatrizen 145 ff.
- -, inkonsistente 241
- -, konsistente 241
- --- mit orthogonalen Matrizen 148ff.
- -, quadratische reguläre 45, 142ff.
- --, rechteckige 240ff.
- -, schwach besetzte 46
- --. überbestimmte 47, 240 ff.
- --, unterbestimmte 301
- Mannigfaltigkeit 242
- Quadratmittelprobleme 47 f., 243 ff.
- -. gewichtete 266f., 268, 274
- - mit Nebenbedingungen 269
- —, rangdefiziente 253, 264, 283, 315ff.
lineare Quadratmittelprobleme, spaltenreguläre 264, 270ff., 277ff. - -, steife 267, 268, 274, 292 - -. totale 269 Unabhängigkeit 14 -s Ausgleichsproblem 243 Linearkombination 12 Lipschitzkonstante 54 -. lokale 53 -, relative 56 Lipschitzstetigkeit 54 -, lokale 53 LL^T-Faktorisierung 188ff. -; Aufdatierung 217ff. - von Bandmatrizen 224, 226ff. direkte 206 f. LNT-Matrix 101 Lösbarkeitsbedingung 241 Lösungsmenge eines Gleichungssystems 242 - eines Quadratmittelproblems 243, 244 LR-Faktorisierung 143, 158, 162 - von Bandmatrizen 221 ff. -, direkte 201 ff. von diagonaldominanten Matrizen 165ff. - von positiv definiten Matrizen 166ff., 204 ff. - von schwach besetzten Matrizen 220, 226 ff. - von symmetrischen Matrizen 226 von Tridiagonalmatrizen 225 Mantisse 63, 64 Matrix 14 -. defektive 34 -, diagonalähnliche 37 -, dominante 165 -, gestufte 359 -, indefinite 22 -, inverse 21 -, orthogonale 24 -, positiv definite 22 -, - semidefinite 22 -, quadratische 18 -, range fiziente 21 -, reguläre 21 -, schiefsymmetrische 30 -, schwach besetzte 46 -, singuläre 21 -, spaltenorthonormale 24 -, spaltenreguläre 21 -, symmetrische 17 -, transponierte 16 -, zeilenorthonormale 24 Matrixnormen 27 ff. Matrizen 90

Matrizen, ähnliche 35 -, äquivalente 41 -, kongruente 40 -, vertauschbare 19 Meßfehler 137 Methode der kleinsten Quadrate 243 MGS-Orthogonalisierung 281ff. - mit Pivotisierung 282f. mit Re-Orthogonalisierung 284 modifizierte Gram-Schmidt-Orthogonalisierung \rightarrow MGS-Orthogonalisierung monotone Norm 30 Moore-Penrose-Inverse \rightarrow Pseudoinverse Nichtorthogonale elementare Transformationsmatrix 101 - - -, stabilisierte 105 NNE = Nichtnullelemente 220 ff.Normalgleichungen 247 Normalgleichungsverfahren 270ff. mit höherer Genauigkeit 274 -: Instabilität 273 mit iterativer Verbesserung 275 Normallösung 248ff., 301, 315ff. NT-Matrix 101 Nullmatrix 16 Nullraum einer Matrix 21 numerische Gutartigkeit 79, 90 - Probleme 50 - Problemklassen 50 Stabilität 79, 91 -r Algorithmus 73 -r -; Computerrealisierung 73 -s Verfahren \rightarrow numerischer Algorithmus Orthogonale Faktorisierung $\rightarrow QR$ -Faktorisierung Projektion 244 -r Projektor 245, 250 -s Komplement 246 Orthogonalität 24 Orthogonalisierungsverfahren 143, 277 ff., 316ff. Orthonormalität 24 Penrose-Bedingungen 249 Permutation 97 Permutationsmatrix 98 Pivotelement 153 Pivotisierung beim Gaußschen Algorithmus 153ff., 161ff. - bei impliziten Givens-Drehungen 123, 292 beim Jacobi-Verfahren 369ff. bei der QR-Faktorisierung 282 f., 292

Potenzmethode \rightarrow Vektoriteration

Profil einer Matrix 226 Programmpaket 452 Programmsystem 452 $Projection \rightarrow orthogonale Projection$ $Projektor \rightarrow orthogonaler Projektor$ Pseudoinverse 48, 249 Pseudorang 327, 332 OL-Algorithmus 431 f. **OR**-Algorithmus 415ff. -, expliziter 426 -; Grundform 416 -, impliziter 428 -, tridiagonaler 424 ff., 430 **OR**-Faktorisierung 143 -; Eindeutigkeit 279 - nach GIVENS 290ff. - nach GRAM-SCHMIDT 277 ff. - nach HOUSEHOLDER 285ff. quadratische Form 22 Quadratmittellösung 243 **QZ**-Algorithmus 447 Rang einer Matrix 21 Rang-1-Modifikation 208, 234f., 275, 301ff., 449 rangerhöhende Störungen 253f., 326, 340 Rayleigh-Quotient 362, 385, 442

Rayleigh-Quotienten-Iteration 395, 444
Rayleigh-Quotienten-Verschiebung 423
Rayleigh-Ritz-Algorithmus 384ff., 444ff.
Reduktion eines allgemeinen Eigenwertproblems auf ein spezielles, explizite 438ff.
- - - - - - , implizite 442ff.
Regularisierung 59
, diskrete 325 ff.
-. kontinuierliche 338 ff.
Regularisierungsparameter 327, 338
relative Maschinengenauigkeit 67
-s Rundungsfehlerniveau 67

- Re-Orthogonalisierung bei Aufdatierung von **QR**-Faktorisierungen 313
- bei Gram-Schmidt-Orthogonalisierung 284
 bei inverser Iteration 394
- Residualkriterien bei Eigenwertproblemen 360ff.
- bei Gleichungssystemen 139ff.
- bei Quadratmittelproblemen 261 ff.
- Residuum bezüglich eines allgemeinen Eigenwertproblems 445
- bezüglich eines speziellen Eigenwertproblems 360
- bezüglich eines Gleichungssystems 139
- bezüglich eines Quadratmittelproblems 261, 299

Reskalierung 292 Ritzsche Eigenwerte und Eigenvektoren 385 RQ-Iteration \rightarrow Rayleigh-Quotienten-Iteration RR-Algorithmus \rightarrow Rayleigh-Ritz-Algorithmus Rücksubstitution 159 Rückwärtsanalyse 79 Rundung 66 -, symmetrische 66 -, unsymmetrische 66 Rundungsfehler 67, 91 Satz von Schur 38 Schneiden des Spektrums 408ff. schrittweise Regression 266, 269 Schwarzsche Ungleichung 25 Sherman-Morrison-Formel 235 simultane Iteration 387 Singulärvektoren 42 Singulärwerte 41 Singulärwertzerlegung 41 -; Berechnung 316ff., 321 Skalarprodukt 19 Skalierung bei linearen Gleichungssystemen 135---- Quadratmittelproblemen 261, 266, 291 SLNT-Matrix 105 Software 450ff. Spaltenpivotisierung nach Bunch-Kaufman-Parlett 197 beim Gaußschen Algorithmus 161, 203 - bei Givens-Drehungen 123 bei QR-Faktorisierung 292 Spaltenraum einer Matrix 21 Spektralnorm 27 Spektralverschiebung 355 - bei inverser Iteration 389ff. bei OR-Algorithmus 420ff. bei SVD 319f. Spektrumslücke 358 Spiegelungsmatrizen 107ff. Störungslemma 128 Störungstheorie von Eigenwertproblemen 355ff. von Gleichungssystemen 129 - von iterativen Matrizen 128 - von Pseudoinversen 251 ff. von Quadratmittelproblemen 257 ff. Sturmsche Ketten 434 $SVD \rightarrow Singulärwertzerlegung$

symmetrische Blockfaktorisierung 193

Teilraum 13

—, invarianter 374

Teilraumiteration 381 ff., 441 ff.
Trägheitsgesetz für quadratische Formen 40
Transformation auf Tridiagonalform 397 ff.
Trapezmatrix 15
Tridiagonalmatrix 224, 398

; Dreiecksfaktorisierung 225
, nichtzerfallende 406
; QR-Algorithmus 424 ff.

Vektoren 11

, linear abhängige/unabhängige 14
, orthogonale/orthonormale 24
Vektornormen 25, 26

Vektoriteration 375ff., 441ff. verallgemeinerte Inverse → Pseudoinverse Verschiebung → Spektralverschiebung Vertauschungsmatrix 99 Vielfachheit, algebraische 33 --, geometrische 34 Vollrang 21 Vollrangmatrix 21 vollständige Pivotisierung 165, 197 Vorwärtsanalyse 79 Vorwärtselimination 159

Wertebereich einer Matrix 21 Wielandt-Hoffman-Theorem 356 Wilkinson-Verschiebung 424, 432 Winkel zweier Vektoren 26 - zweier Teilräume 381

Zeilenpivotisierung beim Gaußschen Algorithmus 164