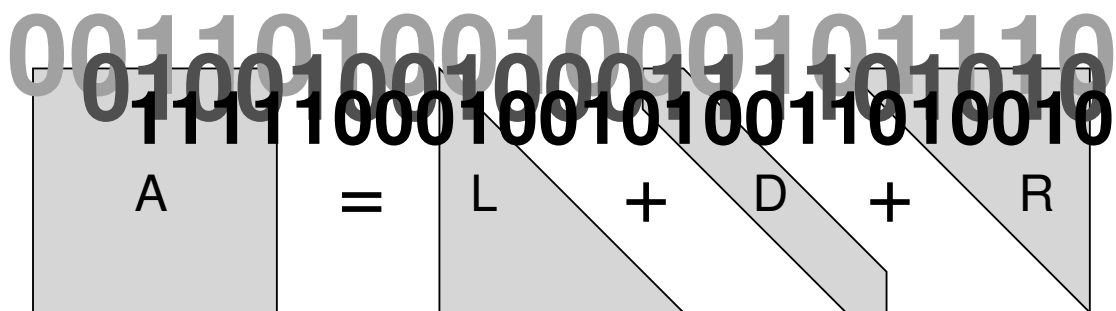

Stefan Funken

Numerik I

(Einführung in die Numerische Lineare Algebra)



Danksagung. Frau Kristin Kirchner und Herrn Moritz Reinhard bin ich für das sorgfältige Lesen des Manuskripts zu besonderem Dank verpflichtet. Ihre zahlreichen Kommentare, Vorschläge, Korrekturen und Hinweise haben die Qualität des Textes wesentlich verbessert. Insbesondere möchte ich Herrn Reinhard für die Ausarbeitung der Aufgaben und zugehöriger Lösungen danken.

Ganz besonderer Dank gebührt auch Frau Petra Hildebrand, die meine handschriftlichen Aufzeichnungen in \LaTeX umgesetzt und zahlreiche Grafiken erstellt hat. Frau Brandner, Frau Serbine und Herrn Weithmann möchte ich für das \LaTeX 'en der Lösungen danken.

Die aufmerksame Leserin oder Leser sei ermuntert, mir Hinweise auf Druckfehler, sprachliche Unzulänglichkeiten oder inhaltliche Flüchtigkeiten per Email zukommen zu lassen (`stefan.funken@uni-ulm.de`).

Copyright. Alle Rechte, insbesondere das Recht auf Vervielfältigung und Verbreitung sowie der Übersetzung, vorbehalten. Kein Teil des Werkes darf in irgendeiner Form ohne schriftliche Genehmigung des Autors reproduziert oder unter Verwendung elektronischer Systeme oder auf anderen Wegen verarbeitet, vervielfältigt oder verbreitet werden.

Stand. Ulm, Oktober 2010.

Inhaltsverzeichnis

1	Lineare Gleichungssysteme	1
1.1	Einführung, Cramersche Regel	1
1.2	Gestaffelte Systeme	4
1.3	Gaußsche Eliminationsmethode	7
1.4	Pivot-Strategien	11
1.5	Nachiteration	17
1.6	Cholesky-Verfahren	17
1.7	Bandgleichungen	21
1.8	Vandermond-Matrizen	23
2	Zahlendarstellung, Fehlerarten und Kondition	27
2.1	Zahlendarstellung	27
2.2	Rundung und Maschinengenauigkeit	36
2.3	Gleitkommaarithmetik	40
2.4	Kondition eines Problems	41
2.5	Numerische Stabilität eines Algorithmus	42
3	Fehlerabschätzung mittels Kondition	47
3.1	Absolute und relative Fehler	47
3.2	Fehlerabschätzung & Kondition	47
4	Iterative Lösung linearer Gleichungssysteme	53
4.1	Folgen von Iterationsmatrizen	55
4.2	Konvergenz des Jacobi-Verfahrens	57
4.3	Konvergenz des Gauß-Seidel-Verfahrens	60
4.4	Abbruch-Kriterien	64
4.5	Gradienten-Verfahren	66
4.6	Verfahren der konjugierten Gradienten	70
4.7	Präkonditionierung, das pcg-Verfahren	77
5	Ausgleichsprobleme	81
5.1	Die Methode der kleinsten Quadrate	81
5.2	Die <i>QR</i> -Zerlegung	83
5.3	<i>Givens</i> -Rotationen	85
5.4	Update einer <i>QR</i> -Zerlegung	87
5.5	<i>Householder</i> -Spiegelungen	90
A	Normen	97
A.1	Vektornorm-Eigenschaften	97
A.2	Matrixnormen	98

B	Einführung in MATLAB	101
	B.1 Grundlegende MATLAB-Befehle	101
	B.2 Mathematik mit Matrizen	106
	B.3 Datenverwaltung	110
	B.4 Ausgabe von Text	111
	B.5 Kontrollbefehle	112
	B.6 Graphische Darstellung	114
	B.7 Fortgeschrittenes	117
C	Speicherformate für schwach besetzte Matrizen	121
	C.1 Koordinatenformat (COO-Format)	121
	C.2 Komprimierte Zeilenspeicherung (CRS-Format)	126
	C.3 Modifizierte komprimierte Zeilenspeicherung (MRS-Format)	128
	C.4 Harwell-Boeing-Format (CCS-Format)	132
D	Aufgaben und Lösungen	133
E	Matlab-Programmieraufgaben	165
	Literaturverzeichnis	174
	Stichwortverzeichnis	176

1 LINEARE GLEICHUNGSSYSTEME

1.1 EINFÜHRUNG, CRAMERSCHE REGEL

Wir beginnen mit dem klassischen Verfahren zur Lösung eines linearen Gleichungssystems (LGS), der Gaußschen¹ Eliminationsmethode.

Zu lösen ist ein System von n linearen Gleichungen mit n Unbekannten $x_1, \dots, x_n \in \mathbb{R}$

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \dots + a_{1n}x_n &= b_1 \\ a_{21}x_1 + a_{22}x_2 + \dots + a_{2n}x_n &= b_2 \\ \vdots & \\ a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n &= b_n \end{aligned} \quad (1.1)$$

oder kurz

$$Ax = b, \quad (1.2)$$

wobei $A \in \mathbb{R}^{n \times n}$ eine reelle $n \times n$ -Matrix ist und $b, x \in \mathbb{R}^n$ reelle (Spalten-)Vektoren sind.

Wann ist ein lineares Gleichungssystem überhaupt lösbar? Aus der Linearen Algebra (siehe z.B. [Wille]) kennen wir das folgende Resultat, das die Lösbarkeit mit Hilfe der Determinante der Matrix A charakterisiert.

Satz 1.1.1 (Lösbarkeit) Sei $A \in \mathbb{R}^{n \times n}$ mit $\det A \neq 0$ und $b \in \mathbb{R}^n$. Dann existiert genau ein $x \in \mathbb{R}^n$, so dass $Ax = b$.

Falls $\det A \neq 0$, so lässt sich die Lösung $x = A^{-1}b$ mit der **Cramerschen Regel** berechnen, d.h.

$$x_i = \frac{1}{\det A} \begin{vmatrix} a_{11} & \dots & b_1 & \dots & a_{1n} \\ a_{21} & \dots & b_2 & \dots & a_{2n} \\ \vdots & & \vdots & & \vdots \\ a_{n1} & \dots & b_n & \dots & a_{nn} \end{vmatrix} = \frac{D_i}{D} \quad (i = 1, \dots, n).$$

Dabei geht die im Zähler stehende Determinante D_i dadurch aus $D = \det A$ hervor, dass man deren i -te Spalte durch den Vektor b der rechten Seite ersetzt.

Man beachte hier die Verbindung von Existenz- und Eindeutigkeitsaussage mit dem Rechenverfahren, was einen „guten“ Algorithmus ausmacht. Dieser braucht dabei nicht unbedingt optimal zu sein!

Die Determinanten von $n \times n$ -Matrizen lassen sich mittels der **Leibnizschen Darstellung** berechnen, d.h.

$$\det A := \sum_{\pi} (-1)^{j(\pi)} a_{1i_1} a_{2i_2} \dots a_{ni_n},$$

wobei die Summe über alle möglichen $n!$ Permutationen π der Zahlen $1, 2, \dots, n$ zu berechnen ist. Der Wert des Ausdrucks $(-1)^{j(\pi)}$ ergibt sich aus der Anzahl $j(\pi)$ der Inversionen der Permutation $\pi = \begin{pmatrix} 1 & 2 & \dots & n \\ i_1 & i_2 & \dots & i_n \end{pmatrix}$.

¹Carl Friedrich Gauß, 1777 - 1855. Lagrange hatte 1759 die Methode schon vorweggenommen und in China war sie schon vor dem ersten Jahrhundert bekannt. Näheres zu Gauß, Lagrange und weiteren Mathematikern findet man im Internet unter www-groups.dcs.st-andrews.ac.uk/~history.

Bemerkung 1.1.2 (Rechenoperationen) *Im Folgenden werden die Verknüpfungen Multiplikation, Addition, Division und Subtraktion in ihrem Rechenaufwand nicht unterschieden und unter dem Begriff **Gleitkommaoperation** zusammengefasst (1 Gleitkommaoperation $\simeq 1 \text{ FLOP}^2$). Systematische Multiplikationen mit ± 1 bleiben im Allgemeinen unberücksichtigt. Anderweitige Operationen wie z.B. Wurzelziehen werden gesondert betrachtet.*

Satz 1.1.3 (Aufwand Leibnizsche Darstellung) *Sei $A \in \mathbb{R}^{n \times n}$. Der Aufwand zur Berechnung von $\det A$ mit der Leibnizschen Darstellung, d.h. als Summe über alle Permutationen der Menge $\{1, \dots, n\}$, beträgt*

$$\text{FLOP}(\det A) = n n! - 1.$$

Beweis. Der Aufwand zur Berechnung von $\det A$ für $A \in \mathbb{R}^{n \times n}$ in der Leibnizschen Darstellung (unter Vernachlässigung der Multiplikation mit $(-1)^{j(\pi)}$) ergibt sich wie folgt:

$$\begin{aligned} \text{FLOP}(\det A) &= \text{'' } (n! - 1) \text{ Additionen } + n! \text{ Produkte mit } n \text{ Faktoren ''} \\ &= n! - 1 + n!(n - 1) = n n! - 1. \end{aligned}$$

Damit ist die Aussage des Satzes bewiesen. □

Alternativ lässt sich die Determinante rekursiv mit Hilfe des **Laplaceschen Entwicklungssatzes** berechnen. Dieser lautet für eine quadratische Matrix $A \in \mathbb{R}^{n \times n}$

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

wobei A_{1j} diejenige Matrix ist, die aus A durch Streichen der ersten Zeile und der j -ten Spalte entsteht.

Satz 1.1.4 (Aufwand Laplacescher Entwicklungssatz) *Sei $A \in \mathbb{R}^{n \times n}$ ($n \geq 2$). Der Aufwand zur Berechnung von $\det A$ mit dem Laplaceschen Entwicklungssatz beträgt*

$$\text{FLOP}(\det A) = 2 n! - 1.$$

Beweis. Wir zeigen zuerst induktiv, dass für $n \geq 2$

$$\sum_{k=1}^{n-1} \frac{k}{(k+1)!} = 1 - \frac{1}{n!} \tag{1.3}$$

gilt. Der Induktionsanfang für $n = 2$ ist offensichtlich erfüllt. Der Induktionsschritt $n \rightarrow n + 1$ ergibt sich aus

$$\sum_{k=1}^n \frac{k}{(k+1)!} = \sum_{k=1}^{n-1} \frac{k}{(k+1)!} + \frac{n}{(n+1)!} = \frac{(n! - 1)}{n!} + \frac{n}{(n+1)!} = 1 - \frac{1}{(n+1)!}.$$

Nun zeigen wir weiter, dass der Aufwand zur Berechnung der Determinante einer $n \times n$ -Matrix mit dem Laplaceschen Entwicklungssatz

$$\text{FLOP}(\det(\mathbb{R}^{n \times n})) = n! \left(1 + \sum_{k=1}^{n-1} \frac{k}{(k+1)!} \right) \tag{1.4}$$

²**FLOP** ist nicht zu verwechseln mit **FLOP/s** oder **flops** (floating point operations per second), welches als Maßeinheit für die Geschwindigkeit von Computersystemen verwendet wird.

beträgt.

Induktionsanfang ($n = 2$): Die Determinante der Matrix $A = (a_{ij}) \in \mathbb{R}^{2 \times 2}$ lautet $a_{11}a_{22} - a_{21}a_{12}$, d.h. der Aufwand beträgt 2 Multiplikationen und 1 Addition, also

$$\text{FLOP}(\det(\mathbb{R}^{2 \times 2})) = 3 = 2! \left(1 + \frac{1}{2!}\right).$$

Induktionsschritt ($n \rightarrow n + 1$): Für $n \geq 2$ gilt

$$\begin{aligned} \text{FLOP}(\det(\mathbb{R}^{(n+1) \times (n+1)})) &= \text{„}n \text{ Additionen und } (n+1) \text{ Multiplikationen mit } \det(\mathbb{R}^{n \times n})\text{“} \\ &= n + (n+1) \cdot \text{FLOP}(\det(\mathbb{R}^{n \times n})) \\ &= \frac{(n+1)!n}{(n+1)!} + (n+1)! \left(1 + \sum_{k=1}^{n-1} \frac{k}{(k+1)!}\right) \\ &= (n+1)! \left(1 + \sum_{k=1}^{n-1} \frac{k}{(k+1)!} + \frac{n}{(n+1)!}\right) \\ &= (n+1)! \left(1 + \sum_{k=1}^n \frac{k}{(k+1)!}\right). \end{aligned}$$

Aus den Gleichungen (1.3) und (1.4) ergibt sich die Behauptung des Satzes. \square

Beispiel 1.1.5 Ein einfaches Beispiel soll zeigen, dass man die Determinante überhaupt nur für sehr kleine allgemeine Matrizen $n \ll 23$ mit dem Laplaceschen Entwicklungssatz berechnen kann.

Ein Jahr hat $365 \cdot 24 \cdot 60 \cdot 60 \approx 3 \cdot 10^7$ Sekunden. Geht man nun von einem schnellen Rechner³ mit 180 TFlops aus, so benötigt man zur Berechnung der Determinante einer 20×20 -Matrix mit dem Laplaceschen Entwicklungssatz

$$\frac{\text{Anzahl der Operationen}}{\text{Gleitkommaoperationen pro Sekunde}} [s] = \frac{2 \cdot 20! - 1}{180 \cdot 10^{12}} [s] \approx 27032 [s] \approx 7.5 [h]$$

bzw. für eine 23×23 -Matrix

$$\frac{2 \cdot 23! - 1}{180 \cdot 10^{12}} [s] \approx 2.872 \cdot 10^8 [s] \approx 9.1 \text{ Jahre}.$$

Bemerkung 1.1.6 Bei der Cramerschen Regel ist zum Lösen eines linearen Gleichungssystems $Ax = b$ mit $A \in \mathbb{R}^{n \times n}$ und $b \in \mathbb{R}^n$ die Berechnung von $n + 1$ Determinanten und n Quotienten notwendig. Der Aufwand zur Lösung eines linearen Gleichungssystems lässt sich somit zusammenfassen, wobei wir auf die Komplexität des Gauß-Verfahrens (siehe Seite 9) erst später eingehen werden.

Cramersche Regel		Gauß-Elimination
Leibniz	Laplace	
$n(n+1)! - 1$	$2(n+1)! - 1$	$\frac{4n^3 + 9n^2 - n}{6}$

³Der weltweit auf Platz sechs rangierende und zugleich schnellste europäische Supercomputer JUGENE am nordrhein-westfälischen Forschungszentrum Jülich hat 180 TFlops (Stand Okt. 2008). Ein aktueller PC hat eine Leistung von 1-10 GFlops. T = Tera = 10^{12} , G = Giga = 10^9 .



Bemerkung 1.1.7 Für das Lösen eines LGS mit 20 Unbekannten mit der Cramerschen Regel und dem Laplaceschen Entwicklungssatz benötigt man auf einem der schnellsten Rechner etwa eine knappe Woche, d.h. $21 \cdot 7.5[h] = 157.5[h] = 6.5625[d]$ (20 verschiedene Determinanten im Zähler und eine im Nenner, vgl. Beispiel 1.1.5). Ein System mit 23 Unbekannten ist mit einem heutigen Supercomputer und der Cramerschen Regel nicht in einem Menschenleben zu lösen.

Beispiel 1.1.8 (Vergleich Rechenaufwand) Aufwand zum Lösen eines linearen Gleichungssystems $Ax = b$ via Cramerscher Regel und Gauß-Verfahren.

Für einige n sei die Anzahl der notwendigen FLOPs wiedergegeben.

n	Cramer/Leibniz $= n(n+1)! - 1$	Cramer/Laplace $= 2(n+1)! - 1$	Gauß $= (4n^3 + 9n^2 - n)/6$
$n = 2$	11	11	11
$n = 3$	71	47	31
$n = 4$	479	239	66
$n = 5$	3599	1439	120
$n = 8$	2903039	725759	436
$n = 10$	399167999	79833599	815

Im Folgenden werden wir nun Verfahren beschreiben, die bereits für $n \geq 3$ effektiver als die Cramersche Regel sind. Mit Hilfe dieser Verfahren lassen sich auch Determinanten in polynomialer Zeit bestimmen. Daher wird die Cramersche Regel im Allgemeinen nur für $n = 2, 3$ verwendet. Der Vorteil der Cramerschen Regel ist jedoch die explizite Schreibweise, d.h. sie ist eine Formel für alle Fälle. Man spart sich bei ähnlichem Rechenaufwand (d.h. $n = 2, 3$) aufwendige Fallunterscheidungen (z.B. Pivotwahl) im Programm.

1.2 GESTAFFELTE SYSTEME

Betrachten wir zuerst besonders einfach zu lösende Spezialfälle. Am einfachsten ist sicherlich der Fall einer diagonalen Matrix A .

Diagonalmatrix

Man bezeichnet eine Diagonalmatrix $A = (a_{ij})$ als Diagonalmatrix, falls $a_{ij} = 0$ für $i \neq j$ gilt. Häufig schreibt man eine Diagonalmatrix A mit Diagonaleinträgen a_{11}, \dots, a_{nn} auch einfach $A = \text{diag}(a_{11}, \dots, a_{nn})$. Die Anwendung der Inversen einer Diagonalmatrix A mit Diagonalelementen $a_{ii} \neq 0$ auf einen Vektor b lässt sich als Pseudocode wie folgt schreiben:

Algorithmus 1.2.1: Lösen von $Ax = b$ mit Diagonalmatrix A

Sei $A \in \mathbb{R}^{n \times n}$ eine invertierbare Diagonalmatrix und $b \in \mathbb{R}^n$

Input $A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$

for $j = 1, \dots, n$

$$x_j = b_j / a_{jj}$$

end

Output $x = (x_1, \dots, x_n)^T$

Eine einfache Matlab Realisierung ist im Folgenden dargestellt.

MATLAB-Beispiel:

Man löse $Ax = b$ mit

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 4 \end{pmatrix}, b = \begin{pmatrix} 5 \\ 3 \\ 2 \end{pmatrix}.$$

```
>> A=diag([1,2,4]);
>> b=[5;3;2];
>> for j=1:3, x(j)=b(j)/A(j,j); end
>> x
x =
    5.0000    1.5000    0.5000
```

Dreiecksmatrix

Der nächstschwierigere Fall ist der einer **Dreiecksmatrix** A . Man spricht von einer oberen Dreiecksmatrix $A = (a_{ij})$, falls $a_{ij} = 0$ für $i > j$, und von einer unteren Dreiecksmatrix $A = (a_{ij})$, falls $a_{ij} = 0$ für $i < j$ gilt. Häufig verwenden wir auch R für eine obere Dreiecksmatrix und L für eine untere Dreiecksmatrix.

Betrachten wir nun das „gestaffelte“ Gleichungssystem

$$\begin{aligned} r_{11}x_1 + r_{12}x_2 + \cdots + r_{1n}x_n &= z_1 \\ r_{22}x_2 + \cdots + r_{2n}x_n &= z_2 \\ &\vdots \\ r_{nn}x_n &= z_n \end{aligned} \quad (1.5)$$

oder in Matrix-Vektor-Schreibweise

$$Rx = z. \quad (1.6)$$

Offenbar erhalten wir x durch sukzessive Auflösung, beginnend mit der n -ten Zeile:

$$\begin{aligned} x_n &:= z_n/r_{nn} && \text{, falls } r_{nn} \neq 0 \\ x_{n-1} &:= (z_{n-1} - r_{n-1,n}x_n)/r_{n-1,n-1} && \text{, falls } r_{n-1,n-1} \neq 0 \\ &\vdots && \vdots \\ x_k &:= (z_k - \sum_{i=k+1}^n r_{ki}x_i)/r_{kk} && \text{, falls } r_{kk} \neq 0 \\ &\vdots && \vdots \\ x_1 &:= (z_1 - r_{12}x_2 - \dots - r_{1n}x_n)/r_{11} && \text{, falls } r_{11} \neq 0. \end{aligned} \quad (1.7)$$

Für die obere Dreiecksmatrix R gilt, dass

$$\det R = r_{11} \cdot r_{22} \cdot \dots \cdot r_{nn} \quad (1.8)$$

und daher

$$\det R \neq 0 \iff r_{ii} \neq 0 \quad (i = 1, \dots, n). \quad (1.9)$$

Der angegebene Algorithmus ist also wiederum genau dann anwendbar, wenn $\det R \neq 0$, also unter der Bedingung des Existenz- und Eindeutigkeitsatzes (Satz 1.1.1).

Analog zum obigen Vorgehen lässt sich auch ein gestaffeltes lineares Gleichungssystem der Form

$$Lx = z$$

mit einer unteren Dreiecksmatrix $L \in \mathbb{R}^{n \times n}$ lösen. In diesem Fall beginnt man in der ersten Zeile mit der Berechnung von x_1 und arbeitet sich dann bis zur letzten Zeile zur Bestimmung von x_n vor.



Bemerkung 1.2.1 (Vorwärts-, Rückwärtssubstitution) Das Lösen eines LGS mit oberer Dreiecksmatrix nennt man auch Rückwärtseinsetzen/-substitution (Index läuft „rückwärts“ von n nach 1), bzw. mit einer unteren Dreiecksmatrix auch Vorwärtseinsetzen/-substitution (Index läuft „vorwärts“ von 1 nach n).



Satz 1.2.2 (Rechenaufwand Ab und $A^{-1}b$ (Dreiecksmatrix)) Sei $A \in \mathbb{R}^{n \times n}$ eine reguläre obere oder untere Dreiecksmatrix und $b \in \mathbb{R}^n$. Dann gilt

$$\text{FLOP}(Ab) = n^2 \quad \text{und} \quad \text{FLOP}(A^{-1}b) = n^2.$$

Bei der Berechnung von $A^{-1}b$ ist im Gegensatz zu Ab die Reihenfolge, in der die Zeilen „abgearbeitet“ werden, festgelegt.

Beweis. Es genügt den Fall einer regulären oberen Dreiecksmatrix $R \in \mathbb{R}^{n \times n}$ zu betrachten. Für den Rechenaufwand zur Lösung von $Rx = b$ (oder die Anwendung von R^{-1} auf einen Vektor b) ergibt sich:

- i) für die i -te Zeile: je $(n - i)$ Additionen und Multiplikationen sowie eine Division
- ii) insgesamt für die Zeilen n bis 1 : Betrachten wir zuerst die Summenformel

$$\sum_{i=1}^n (n - i) = \sum_{i=1}^{n-1} i = \frac{n(n-1)}{2} = \frac{n^2 - n}{2}.$$

Mit i) erhält man insgesamt einen Aufwand von $2 \frac{n^2 - n}{2} + n = n^2$ Operationen.

Der zweite Teil der Aussage bleibt Ihnen als Übungsaufgabe überlassen. □

Aufgabe 1.2.3 Es sei $A \in \mathbb{R}^{n \times n}$ eine obere (oder untere) Dreiecksmatrix und $b \in \mathbb{R}^n$. Man zeige, dass das Matrix-Vektor-Produkt Ab mit n^2 Operationen berechnet werden kann.

Eine Matlab Realisierung zur Bestimmung von $R^{-1}b$, bei der die Matrix R zeilenweise durchlaufen wird, und ein Anwendungsbeispiel sind im Folgenden dargestellt.

MATLAB-Funktion: Rinvb.m

```

1 function x = Rinvb(R,b)
2 % compute solution of R * x = b with upper triangular matrix R
3 n = size(R,1);
4 x = zeros(n,1);
5 for j = n:-1:1
6     for k=j+1:n
7         b(j) = b(j) - R(j,k) * x(k);
8     end
9     x(j) = b(j) / R(j,j);
10 end

```

MATLAB-Beispiel:

Man löse $Rx = b$ mit

$$R = \begin{pmatrix} 4 & 1 & 1 \\ 0 & 3 & 2 \\ 0 & 0 & 1 \end{pmatrix}, b = \begin{pmatrix} 9 \\ 12 \\ 3 \end{pmatrix}.$$

```
>> R = [4,1,1;0,3,2;0,0,1];
>> b = [9;12;3];
>> x = R\invb(R,b);
>> x'
```

ans =

1	2	3
---	---	---

Eine alternative Realisierung, bei der die Matrix R spaltenweise durchlaufen wird, ist mit `R\invb2.m` gegeben.

MATLAB-Funktion: R\invb2.m

```
1 function x = R\invb2(R,b)
2 % compute solution of R * x = b with upper triangular matrix R
3 n = size(R,1);
4 x = zeros(n,1);
5 for j = n:-1:1
6     x(j) = b(j) / R(j,j);
7     for k=1:j-1
8         b(k) = b(k) - R(k,j) * x(j);
9     end
10 end
```

Was ist der Unterschied? Machen Sie sich dies klar!

1.3 GAUSSSCHE ELIMINATIONSMETHODE

Gäbe es nun zu einer beliebigen Matrix A eine Zerlegung

$$A = L \cdot R, \quad (1.10)$$

so könnte ein beliebiges lineares Gleichungssystem $Ax = b$ durch eine Rückwärts- und Vorwärts-substitution mittels der beiden Schritte

i) löse $Lz = b$,

ii) löse $Rx = z$,

gelöst werden. Die folgende Gaußsche Eliminationsmethode liefert gerade eine solche Zerlegung. Betrachten wir ein allgemeines lineares Gleichungssystem $Ax = b$ ($A \in \mathbb{R}^{n \times n}$ regulär, $b \in \mathbb{R}^n$)

$$\begin{array}{cccccc} a_{11}x_1 & + & a_{12}x_2 & + & \cdots & + & a_{1n}x_n & = & b_1 \\ a_{21}x_1 & + & a_{22}x_2 & + & \cdots & + & a_{2n}x_n & = & b_2 \\ \vdots & & \vdots & & & & \vdots & & \vdots \\ a_{n1}x_1 & + & a_{n2}x_2 & + & \cdots & + & a_{nn}x_n & = & b_n \end{array} \quad (1.11)$$

und versuchen dies in ein gestaffeltes System umzuformen.

Durch die folgenden drei Äquivalenzoperationen

1. Vertauschung von Zeilen,
2. Multiplikation einer Zeile mit einem Skalar $\neq 0$,
3. Addition eines Vielfachen einer Zeile zu einer anderen,

wird die Lösungsmenge des Gleichungssystems (1.11) nicht verändert.

Wir setzen zunächst $a_{11} \neq 0$ voraus. Um (1.11) nun in ein gestaffeltes System umzuformen, muss die erste Zeile nicht verändert werden. Die restlichen Zeilen werden so modifiziert, dass die Koeffizienten vor x_1 verschwinden, d.h. die Variable x_1 aus den Gleichungen in den Zeilen 2 bis n eliminiert wird.

So entsteht ein System der Art

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1 \\ a'_{22}x_2 + \cdots + a'_{2n}x_n &= b'_2 \\ &\vdots \\ a'_{n2}x_2 + \cdots + a'_{nn}x_n &= b'_n. \end{aligned} \quad (1.12)$$

Haben wir dies erreicht, so können wir dasselbe Verfahren auf die letzten $(n-1)$ Zeilen anwenden und so rekursiv ein gestaffeltes System erhalten. Mit dem Quotienten

$$l_{i1} = a_{i1}/a_{11} \quad (i = 2, 3, \dots, n) \quad (1.13)$$

sind die Elemente in (1.12) gegeben durch

$$\begin{aligned} a'_{ik} &= a_{ik} - l_{i1}a_{1k}, \\ b'_i &= b_i - l_{i1}b_1. \end{aligned} \quad (1.14)$$

Damit ist der erste Eliminationsschritt unter der Annahme $a_{11} \neq 0$ ausführbar.

Wenden wir auf diese Restmatrix die Eliminationsvorschrift erneut an, so erhalten wir eine Folge

$$A = A^{(1)} \rightarrow A^{(2)} \rightarrow \dots A^{(n)} =: R \quad (1.15)$$

von Matrizen der speziellen Gestalt

$$A^{(k)} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & & a_{1n}^{(1)} \\ & a_{22}^{(2)} & \cdots & & a_{2n}^{(2)} \\ & & \ddots & & \\ & & & a_{kk}^{(k)} & \cdots & a_{kn}^{(k)} \\ & & & \vdots & & \vdots \\ & & & a_{nk}^{(k)} & \cdots & a_{nn}^{(k)} \end{pmatrix} \quad (1.16)$$

mit einer $(n-k+1, n-k+1)$ -Restmatrix, auf die wir den Eliminationsschritt

$\begin{aligned} l_{ik} &:= a_{ik}^{(k)} / a_{kk}^{(k)} && \text{für } i = k+1, \dots, n \\ a_{ij}^{(k+1)} &:= a_{ij}^{(k)} - l_{ik}a_{kj}^{(k)} && \text{für } i, j = k+1, \dots, n \\ b_i^{(k+1)} &:= b_i^{(k)} - l_{ik}b_k^{(k)} && \text{für } i = k+1, \dots, n \end{aligned} \quad (1.17)$
--

ausführen können, wenn das **Pivotelement** (Dreh- und Angelpunkt) $a_{kk}^{(k)}$ nicht verschwindet. Da jeder Eliminationsschritt eine lineare Operation auf den Zeilen von A ist, lässt sich der Übergang

von $A^{(k)}$ und $b^{(k)}$ zu $A^{(k+1)}$ und $b^{(k+1)}$ als Multiplikation mit einer Matrix $L_k \in \mathbb{R}^{n \times n}$ von links darstellen, d.h.

$$A^{(k+1)} = L_k A^{(k)}, \quad b^{(k+1)} = L_k b^{(k)}. \tag{1.18}$$

Die Matrix

$$L_k = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -l_{k+1,k} & 1 & \\ & & \vdots & & \ddots \\ & & -l_{n,k} & & & 1 \end{pmatrix} = I - \begin{pmatrix} 0 \\ \vdots \\ 0 \\ l_{k+1,k} \\ \vdots \\ l_{n,k} \end{pmatrix} \cdot e_k^T =: I - \ell_k \cdot e_k^T, \tag{1.19}$$

wobei e_k der k -te Einheitsvektor sei, hat die Eigenschaft, dass die Inverse L_k^{-1} aus L_k durch einen Vorzeichenwechsel in den Einträgen l_{ik} ($i > k$) entsteht und für das Produkt der L_k^{-1} gilt

$$L := L_1^{-1} \cdot \dots \cdot L_{n-1}^{-1} = \begin{pmatrix} 1 & 0 & \dots & \dots & 0 \\ l_{21} & 1 & \ddots & & \vdots \\ l_{31} & l_{32} & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ l_{n1} & \dots & & l_{n,n-1} & 1 \end{pmatrix}. \tag{1.20}$$

Zusammengefasst erhalten wir auf diese Weise das zu $Ax = b$ äquivalente gestaffelte System $Rx = z$ mit der oberen Dreiecksmatrix

$$R = L^{-1}A \quad \text{und} \quad z = L^{-1}b. \tag{1.21}$$

Das **Gaußsche Eliminationsverfahren** schreibt sich dann wie folgt:
 Sei $A \in \mathbb{R}^{n \times n}$ und $b \in \mathbb{R}^n$. Berechne nacheinander

- i) $L \cdot R = A$ (Zerlegung mit R obere und L untere Dreiecksmatrix),
- ii) $Lz = b$ (Vorwärtseinsetzen bzw. -substitution),
- iii) $Rx = z$ (Rückwärtseinsetzen bzw. -substitution).

Definition 1.3.1 (unipotente Matrix) Eine untere oder obere Dreiecksmatrix, deren Diagonalelemente alle gleich eins sind, heißt unipotent.

Definition 1.3.2 (Gaußsche Dreieckszerlegung, LR-Zerlegung) Die obige Darstellung $A = L \cdot R$ der Matrix A als Produkt einer unipotenten unteren Dreiecksmatrix L und einer oberen Dreiecksmatrix R heißt **Gaußsche Dreieckszerlegung** oder **LR-Zerlegung** von A .

Satz 1.3.3 Existiert die LR-Zerlegung einer quadratischen Matrix A , so sind L und R eindeutig bestimmt.

Beweis. Siehe Hausübung. □

Satz 1.3.4 (Rechenaufwand Gaußsche Dreieckszerlegung, Gaußsche Elimination)

Sei $A \in \mathbb{R}^{n \times n}$ regulär. Existiert die LR-Zerlegung, so gilt

$$\text{FLOP}(LR - \text{Zerlegung}) = \frac{4n^3 - 3n^2 - n}{6}.$$

Des Weiteren sei $b \in \mathbb{R}^n$. Dann gilt für die Gaußsche Elimination

$$\text{FLOP}(A^{-1}b) = \frac{4n^3 + 9n^2 - n}{6},$$

d.h. die Lösung von $Ax = b$ kann mit $(4n^3 + 9n^2 - n)/6$ FLOPs berechnet werden.



Bemerkung 1.3.5 Die Aufwandsberechnung unterscheidet sich teilweise in der Literatur. Neu- erdings werden $+$ und \cdot Operationen nicht mehr unterschieden!

Beweis von Satz 1.3.4. Für den Rechenaufwand der Gaußschen Dreieckszerlegung gilt Folgendes:

- i) für den i -ten Eliminationsschritt:
 je $(n - i)$ Divisionen zur Berchnung der l_{ik} und je $(n - i)^2$ Multiplikationen und Subtrak- tionen zur Berechnung der $a_{jk}^{(i+1)}$, d.h. $2(n - i)^2 + (n - i)$ Operationen
- ii) für alle $n - 1$ Eliminationschritte erhalten wir durch Umindizierung und Summenformeln aus Analysis 1 (Aufgabe 1.4.23 in [Analysis I])

$$\begin{aligned} & \sum_{i=1}^{n-1} (n - i) + 2(n - i)^2 = \sum_{i=1}^{n-1} i + 2i^2 = \frac{n(n - 1)}{2} + \frac{(2n - 1)n(n - 1)}{3} \\ & = \frac{2(2n - 1)(n - 1)n + 3n(n - 1)}{6} = \frac{n(n - 1)(4n - 2 + 3)}{6} \\ & = \frac{n(n - 1)(4n + 1)}{6} = \frac{4n^3 - 3n^2 - n}{6} \end{aligned}$$

- iii) insgesamt benötigt man zum Lösen von $Ax = b$ mittels Gauß-Elimination, d.h. LR - Zerlegung und Vorwärts-, Rückwärtssubstitution (siehe Satz 1.2.2)

$$\frac{4n^3 - 3n^2 - n}{6} + 2n^2 = \frac{4n^3 + 9n^2 - n}{6} \text{ Operationen.}$$

Somit sind die Behauptungen des Satzes bewiesen. □

Das Speicherschema für die Gauß-Elimination orientiert sich an der Darstellung von $A^{(k)}$. In die erzeugten Nullen des Eliminationsschritts können die l_{ik} eingetragen werden. Da die Elemente mit den Werten 0 oder 1 natürlich nicht eigens gespeichert werden müssen, kann die LR -Zerlegung mit dem Speicherplatz der Matrix A bewerkstelligt werden:

$$\begin{aligned} A & \rightarrow \left(\begin{array}{cccc} a_{11}^{(1)} & \dots & \dots & a_{1n}^{(1)} \\ l_{21} & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ \vdots & \vdots & & \\ l_{n1} & a_{n2}^{(2)} & & a_{nn}^{(2)} \end{array} \right) \rightarrow \dots \rightarrow \left(\begin{array}{cccc} a_{11}^{(1)} & \dots & & a_{1n}^{(1)} \\ l_{21} & a_{22}^{(2)} & & a_{2n}^{(2)} \\ \vdots & \ddots & & \\ \vdots & & \ddots & \\ l_{n1} & \dots & \dots & l_{n,n-1} & a_{nn}^{(n)} \end{array} \right) \\ \Rightarrow L & = \left(\begin{array}{ccccc} 1 & 0 & \dots & \dots & 0 \\ l_{21} & \ddots & & & \vdots \\ \vdots & & \ddots & & \vdots \\ \vdots & & & \ddots & 0 \\ l_{n1} & \dots & \dots & l_{n,n-1} & 1 \end{array} \right), \quad R = \left(\begin{array}{cccc} a_{11}^{(1)} & \dots & \dots & a_{1n}^{(1)} \\ & \ddots & & \vdots \\ & & \ddots & \vdots \\ & & & \ddots & \vdots \\ & & & & a_{nn}^{(n)} \end{array} \right). \end{aligned}$$

1.4 PIVOT-STRATEGIEN

Bisher haben wir uns noch nicht mit der Frage beschäftigt, ob eine LR -Zerlegung immer existiert und ob eine Vertauschung von Zeilen bei $Ax = b$ zu einem anderen Ergebnis führt, wenn der Algorithmus mit endlicher Rechengenauigkeit durchgeführt wird. (Bei exakter Arithmetik sind beide Ergebnisse gleich, bzw. Vertauschung von Zeilen in A führt zur gleichen Lösung). Betrachten wir hierzu zwei Beispiele:

Beispiel 1.4.1 i) **Permutation liefert LR -Zerlegung:** Sei $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$. Wie wir leicht sehen, gilt für die Eigenwerte $\lambda_{1,2} = \pm 1$ und somit $\det A \neq 0$. Falls eine Zerlegung existieren würde, gäbe es $a, b, c, d \in \mathbb{R}$ mit

$$A = \begin{pmatrix} 1 & 0 \\ a & 1 \end{pmatrix} \begin{pmatrix} b & c \\ 0 & d \end{pmatrix} = \begin{pmatrix} b & c \\ ab & ac + d \end{pmatrix}.$$

Ein einfacher Vergleich der Koeffizienten zeigt jedoch (zuerst nach b und dann nach a), dass keine LR -Zerlegung existiert, obwohl nach Vertauschen der Zeilen trivialerweise eine LR -Zerlegung sofort gegeben ist. Die Frage, die sich nun stellt, lautet: Kann man für jede reguläre Matrix A durch Vertauschen ihrer Zeilen eine Matrix finden, für die dann eine LR -Zerlegung existiert?

ii) **Permutation liefert höhere Genauigkeit:** Berechnen wir die Lösung des folgenden linearen Gleichungssystems ($\epsilon > 0$)

$$\epsilon x_1 + x_2 = 1 \quad (1.22)$$

$$x_1 + x_2 = 2. \quad (1.23)$$

Bei exakter Arithmetik erhalten wir

$$\frac{1}{\epsilon-1} \begin{pmatrix} 1 & -1 \\ -1 & \epsilon \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \frac{1}{1-\epsilon} \begin{pmatrix} 1 \\ 1-2\epsilon \end{pmatrix}$$

d.h. $x_1 = \frac{1}{1-\epsilon}$, $x_2 = \frac{1-2\epsilon}{1-\epsilon}$. Für $2\epsilon < \text{Rechengenauigkeit}$ (d.h. $1 \pm 2\epsilon$ und 1 werden im Rechner durch dieselbe Zahl dargestellt) gilt nun

$$x_1 = 1, \quad x_2 = 1.$$

Für die Gauß-Elimination erhalten wir, wenn wir den Koeffizienten vor x_1 in (1.23) eliminieren, d.h. das $1/\epsilon$ -fache von (1.22) subtrahieren

$$\begin{aligned} \epsilon x_1 + x_2 &= 1 \\ \left(1 - \frac{1}{\epsilon}\right) x_2 &= 2 - \frac{1}{\epsilon}. \end{aligned}$$

Somit ergibt sich für $2\epsilon < \text{Rechengenauigkeit}$ aus der letzten Gleichung $x_2 = 1$ und damit aus der ersten Gleichung $x_1 = 0$. Vertauschen wir jedoch 1. und 2. Zeile (bzw. eliminieren den Koeffizienten vor x_1 in (1.22)) so erhalten wir

$$\begin{aligned} x_1 + x_2 &= 2 \\ (1 - \epsilon)x_2 &= 1 - 2\epsilon. \end{aligned}$$

Dies liefert $x_1 = x_2 = 1$ bei $2\epsilon < \text{Rechengenauigkeit}$.

MATLAB-Beispiel:

Dieses scheinbar theoretische Ergebnis aus Beispiel 1.4.1 können wir direkt in Matlab umsetzen. Da die Recheneinheiten heutiger Rechner meist 2 Stellen genauer rechnen als unsere bisherigen theoretischen Untersuchungen vermuten lassen, muss hier $\epsilon \leq \text{Maschinengenauigkeit}/8$ gelten.

```
>> e = eps/8; % eps/2 genügt nicht!
>> x(2) = (2 - 1/e) / (1 - 1/e);
>> x(1) = (1 - x(2)) / e
x =
    0    1
>> x(2) = (1 - 2*e) / (1 - e);
>> x(1) = 2 - x(2)
x =
    1    1
```

Das Vertauschen kann folglich nötig sein. Zum einen, damit eine LR -Zerlegung überhaupt existiert, zum anderen, um den „Rechenfehler“ bedingt durch Rundungsfehler möglichst klein zu halten. Daraus leiten sich besondere Strategien zur Wahl des Pivotelements ab.

Betrachten wir nun die Gauß-Elimination mit **Spaltenpivotstrategie**, welche theoretisch nur dann nicht zielführend sein kann, falls die Matrix A singular ist.

Algorithmus 1.4.1: Gauß-Elimination mit Spaltenpivotstrategie

- a) Wähle im Eliminationsschritt $A^{(k)} \rightarrow A^{(k+1)}$ ein $p \in \{k, \dots, n\}$, so dass

$$|a_{pk}^{(k)}| \geq |a_{jk}^{(k)}| \quad \text{für } j = k, \dots, n.$$

Die Zeile p soll die Pivotzeile werden.

- b) Vertausche die Zeilen p und k

$$A^{(k)} \rightarrow \tilde{A}^{(k)} \quad \text{mit} \quad \tilde{a}_{ij}^{(k)} = \begin{cases} a_{kj}^{(k)} & \text{falls } i = p \\ a_{pj}^{(k)} & \text{falls } i = k \\ a_{ij}^{(k)} & \text{sonst.} \end{cases}$$

Nun gilt

$$|\tilde{l}_{ik}| = \left| \frac{\tilde{a}_{ik}^{(k)}}{\tilde{a}_{kk}^{(k)}} \right| = \left| \frac{\tilde{a}_{ik}^{(k)}}{a_{pk}^{(k)}} \right| \leq 1.$$

- c) Führe den nächsten Eliminationsschritt angewandt auf $\tilde{A}^{(k)}$ aus,

$$\tilde{A}^{(k)} \rightarrow A^{(k+1)}.$$

Bemerkung 1.4.2 Anstelle der Spaltenpivotstrategie mit Zeilentausch kann man auch eine Zeilenpivotstrategie mit Spaltentausch durchführen. Beide Strategien benötigen im schlimmsten Fall $\mathcal{O}(n^2)$ zusätzliche Operationen. In der Praxis werden im Gegensatz zum oben genannten Algorithmus die Zeile bzw. Spalte nicht umgespeichert, sondern besondere Indexvektoren verwendet (siehe Matlabfunktion `mylu.m`). Die Kombination von Spalten- und Zeilenpivotstrategie führt zur vollständigen Pivotsuche, bei der die gesamte Restmatrix nach dem betragsgrößten Eintrag durchsucht wird. Wegen des Aufwands $\mathcal{O}(n^3)$ wird dies jedoch so gut wie nie angewandt.



MATLAB-Funktion: mylu.m

```

1 function [L,R,P] = mylu(A)
2 n = size(A,1); % get leading dimension of A
3 p = 1 : n; % pivot element vector
4 for k = 1 : n-1 % consider k-th column
5 [m,mptr] = max(abs(A(p(k:end),k))); % find pivot element
6 tmp = p(k); % interchange in vector p
7 p(k) = p(k-1+mptr);
8 p(k-1+mptr) = tmp;
9
10 for j = k+1 : n % modify entries in
11 A(p(j),k) = A(p(j),k)/A(p(k),k); % compute l_jk, store in A
12 for i = k+1 : n % (n-k-1)*(n-k-1) submatrix
13 A(p(j),i) = A(p(j),i) ...
14 - A(p(j),k)*A(p(k),i);
15 end
16 end
17 end
18 L = tril(A(p,:),-1)+eye(n); % these lines could be
19 R = triu(A(p,:)); % neglected, all information
20 P = eye(n); % already in A and p
21 P(:,p) = P;

```

Satz 1.4.3 Es sei $A \in \mathbb{R}^{n \times n}$ regulär. Dann existiert vor dem k -ten Eliminationsschritt des Gauß-Algorithmus stets eine Zeilen-/Spaltenpermutation derart, dass das k -te Diagonalelement von Null verschieden ist. Bei Zeilenpermutation, d.h. Spaltenpivotisierung, sind alle Einträge von L vom Betrag kleiner oder gleich eins.

Beweis. Nach Voraussetzung gilt $\det A \neq 0$. Angenommen, es sei $a_{11} = 0$. Dann existiert eine Zeilenvertauschung, so dass das erste Pivotelement $a_{11}^{(1)}$ von Null verschieden und das betragsgrößte Element in der Spalte ist, d.h.

$$0 \neq |a_{11}^{(1)}| \geq |a_{i1}^{(1)}| \quad \text{für } i = 1, \dots, n,$$

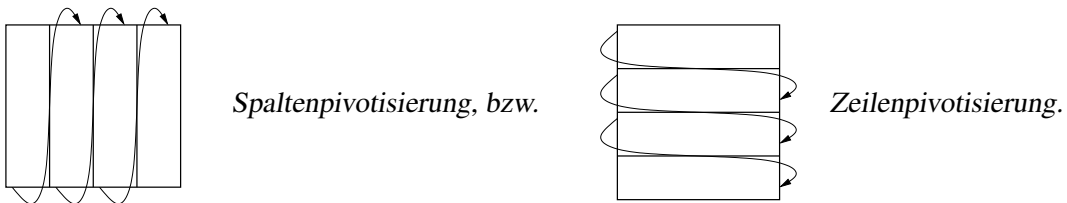
denn andernfalls wäre die Determinante von A in Widerspruch zur Voraussetzung gleich Null. Die Zeilenvertauschung hat dabei nur einen Vorzeichenwechsel in der Determinante zur Folge. Die Überlegungen für den ersten Eliminationsschritt übertragen sich sinngemäß auf die folgenden, reduzierten Systeme, bzw. ihre zugehörigen Determinanten. Diese Schlussfolgerungen gelten analog bei Zeilenpivotisierung. \square

Eine unmittelbare Folge aus der Beweisführung des letzten Satzes ist das folgende Lemma.

Lemma 1.4.4 *Erfolgen im Verlauf des Gauß-Algorithmus total m Zeilenvertauschungen (Spaltenvertauschungen), so ist die Determinante von A gegeben durch*

$$\det A = (-1)^m \prod_{k=1}^n r_{kk}.$$

Bemerkung 1.4.5 *Numerisch sind Spalten- und Zeilenpivotisierungen äquivalent. Die Auswahl hängt von der Speichermethode der Matrix A ab, d.h. man favorisiert die Spaltenpivotisierung, wenn die Einträge der Matrix spaltenweise im Speicher abgelegt sind, z.B. bei den Programmiersprachen Matlab, Fortran und die Zeilenpivotisierung u.a. bei C, denn hier sind die Matrizen als Folge von Zeilenvektoren abgelegt.*



Die genannten Pivotisierungen bzw. Zeilenvertauschungen lassen sich formal durch eine geeignete Permutationsmatrix P beschreiben. Diese ist eine quadratische Matrix, welche in jeder Zeile und in jeder Spalte genau eine Eins und sonst Nullen enthält. Die Determinante ist $\det P = \pm 1$. Aus dem letzten Satz ergibt sich folgendes Resultat.

Satz 1.4.6 (Existenz einer Zerlegung $PA = LR$) *Zu jeder regulären Matrix A existiert eine Permutationsmatrix P , so dass $P \cdot A$ in ein Produkt $L \cdot R$ zerlegbar ist, d.h.*

$$P \cdot A = L \cdot R.$$

Ist nun immer eine Pivotisierung notwendig? Wir nennen im Folgenden ein Kriterium, welches leicht zu überprüfen ist.

Definition 1.4.7 (Diagonaldominanz) *Eine Matrix heißt diagonaldominant, falls gilt*

$$|a_{ii}| \geq \sum_{k=1, k \neq i}^n |a_{ik}| \quad \forall i = 1, \dots, n.$$

Sie heißt strikt diagonaldominant, falls für alle i „>“ anstatt „ \geq “ gilt.

Satz 1.4.8 *Es sei $A \in \mathbb{R}^{n \times n}$ eine diagonaldominante und reguläre Matrix. Dann existiert eine Zerlegung*

$$A = L \cdot R.$$

Beweis. Einen Beweis dieser Aussage (in leicht abgewandelter Form) findet man z.B. bei [Schwarz]. Nach Voraussetzung ist entweder $|a_{11}| \geq \sum_{k=2}^n |a_{1k}| > 0$ oder es gilt $\sum_{k=2}^n |a_{1k}| = 0$ und $|a_{11}| > 0$; dies folgt aus der Regularität von A . Somit ist $a_{11} \neq 0$ ein zulässiges Pivotelement für den ersten Eliminationsschritt. Man muss nun zeigen, dass sich die Eigenschaft der diagonalen Dominanz auf das reduzierte Gleichungssystem überträgt. \square

Zum Schluss dieses Abschnitts betrachten wir noch folgendes Beispiel zur konkreten Berechnung einer LR -Zerlegung einer Matrix A .

Beispiel 1.4.9 (Berechnung von $PA = LR$) Sei

$$A = \begin{pmatrix} 0 & 2 & -1 & -2 \\ 2 & -2 & 4 & -1 \\ 1 & 1 & 1 & 1 \\ -2 & 1 & -2 & 1 \end{pmatrix}.$$

Diese Matrix ist offensichtlich nicht diagonaldominant, also wenden wir die Gauß-Elimination mit Spaltenpivotisierung an. Hierzu sei $P(i, j)$ diejenige Permutationsmatrix, die aus der Einheitsmatrix durch Vertauschen der i -ten und j -ten Zeile entsteht. Das Vertauschen der ersten mit der zweiten Zeile von A entspricht der Multiplikation der Matrix $P_1 := P(1, 2)$ mit A von links, also

$$\tilde{A}^{(1)} = P_1 A = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 1 & 1 & 1 & 1 \\ -2 & 1 & -2 & 1 \end{pmatrix}.$$

Der erste Schritt der LR -Zerlegung liefert

$$L_1 = \begin{pmatrix} 1 & & & \\ 0 & 1 & & \\ -\frac{1}{2} & & 1 & \\ 1 & & & 1 \end{pmatrix}, \quad L_1 P_1 A = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 2 & -1 & \frac{3}{2} \\ 0 & -1 & 2 & 0 \end{pmatrix}.$$

Bei der nächsten Spaltenpivotisierung sind keine Zeilenvertauschungen notwendig, also $P_2 = I$. Es folgt

$$L_2 = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & -1 & 1 & \\ & \frac{1}{2} & & 1 \end{pmatrix}, \quad L_2 P_2 L_1 P_1 A = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & 0 & \frac{7}{2} \\ 0 & 0 & \frac{3}{2} & -1 \end{pmatrix}.$$

Um nun auf eine obere Dreiecksgestalt zu kommen, müssen lediglich noch die dritte und vierte Zeile vertauscht werden, also formal $P_3 = P(3, 4)$, $L_3 = I$ und damit

$$L_3 P_3 L_2 P_2 L_1 P_1 A = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & \frac{3}{2} & -1 \\ 0 & 0 & 0 & \frac{7}{2} \end{pmatrix} =: R.$$

Aber wie sehen nun P und L aus, sodass

$$P \cdot A = L \cdot R \quad ?$$

Wir betrachten hierzu im allgemeinen Fall mit $A \in \mathbb{R}^{n \times n}$ invertierbar für $k = 1, \dots, n-1$ die Matrizen

$$\tilde{L}_k := P_n \cdot \dots \cdot P_{k+1} L_k P_{k+1} \cdot \dots \cdot P_n,$$

wobei $P_n := I$ gesetzt wird und P_1, \dots, P_{n-1} wie in Beispiel 1.4.9 gewählt werden. Dann gilt

$$\begin{aligned} \tilde{L}_k &= P_n \cdot \dots \cdot P_{k+1} L_k P_{k+1} \cdot \dots \cdot P_n \\ &\stackrel{(1.19)}{=} P_n \cdot \dots \cdot P_{k+1} (I - \ell_k e_k^T) P_{k+1} \cdot \dots \cdot P_n \\ &= I - \underbrace{P_n \cdot \dots \cdot P_{k+1} \ell_k}_{=: \tilde{\ell}_k} \underbrace{e_k^T P_{k+1} \cdot \dots \cdot P_n}_{=: e_k^T} \\ &= I - \tilde{\ell}_k e_k^T. \end{aligned}$$

Da die Multiplikation mit den Matrizen P_{k+1}, \dots, P_n von links bzw. rechts nur Zeilen- bzw. Spaltenvertauschungen innerhalb der Zeilen bzw. Spalten $k+1, \dots, n$ bewirkt, besitzt die Matrix \tilde{L}_k dieselbe Struktur wie L_k . Aus der Definition der \tilde{L}_k folgt nun

$$\tilde{L}_{n-1} \cdot \dots \cdot \tilde{L}_1 P_{n-1} \cdot \dots \cdot P_1 = L_{n-1} P_{n-1} \cdot \dots \cdot L_2 P_2 L_1 P_1.$$

Nach Konstruktion der L_k und P_k gilt schließlich

$$\tilde{L}_{n-1} \cdot \dots \cdot \tilde{L}_1 P_{n-1} \cdot \dots \cdot P_1 A = R.$$

Somit ergibt sich mit

$$L := \left(\tilde{L}_{n-1} \cdot \dots \cdot \tilde{L}_1 \right)^{-1} = \tilde{L}_1^{-1} \cdot \dots \cdot \tilde{L}_{n-1}^{-1}$$

und

$$P := P_{n-1} \cdot \dots \cdot P_1$$

eine Zerlegung $PA = LR$.

Bemerkung 1.4.10 Man beachte, dass mit den obigen Überlegungen ein alternativer, konstruktiver Beweis des Satzes 1.4.6 vorliegt.

Führen wir nun Beispiel 1.4.9 weiter fort:

Beispiel 1.4.11 Zu der Matrix A aus Beispiel 1.4.9 sollen die Matrizen L und P der Zerlegung $PA = LR$ bestimmt werden. Wir erhalten

$$P = P_3 P_2 P_1 = P(1, 2) I P(3, 4) = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{pmatrix},$$

$$\tilde{L}_1 = P_3 P_2 L_1 P_2 P_3 = P(3, 4) \begin{pmatrix} 1 & & & \\ 0 & 1 & & \\ -\frac{1}{2} & & 1 & \\ 1 & & & 1 \end{pmatrix} P(3, 4) = \begin{pmatrix} 1 & & & \\ 0 & 1 & & \\ 1 & & 1 & \\ -\frac{1}{2} & & & 1 \end{pmatrix},$$

$$\tilde{L}_2 = P_3 L_2 P_3 = P(3, 4) \begin{pmatrix} 1 & & & \\ & 1 & & \\ -1 & & 1 & \\ & \frac{1}{2} & & 1 \end{pmatrix} P(3, 4) = \begin{pmatrix} 1 & & & \\ & \frac{1}{2} & & \\ -1 & & 1 & \\ & & & 1 \end{pmatrix},$$

$$\tilde{L}_3 = P_4 L_3 P_4 = I,$$

$$L = \tilde{L}_1^{-1} \tilde{L}_2^{-1} \tilde{L}_3^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -1 & -\frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 1 & 0 & 1 \end{pmatrix}$$

und somit

$$PA = \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ -2 & 1 & -2 & 1 \\ 1 & 1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ -1 & -\frac{1}{2} & 1 & 0 \\ \frac{1}{2} & 1 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & -2 & 4 & -1 \\ 0 & 2 & -1 & -2 \\ 0 & 0 & \frac{3}{2} & -1 \\ 0 & 0 & 0 & \frac{7}{2} \end{pmatrix} = LR.$$

1.5 NACHITERATION

Die oben diskutierten Pivotstrategien schließen offenbar nicht aus, dass die so berechnete Lösung x immer noch „ziemlich ungenau“ ist. Wie kann man nun x ohne großen Aufwand verbessern?

Häufig lässt sich dies durch eine **Nachiteration** mit Hilfe einer expliziten Auswertung des **Residuums** $r := b - A\tilde{x}$ erreichen. Dabei bezeichne \tilde{x} die durch ein numerisches Verfahren berechnete Näherung an x . Ausgehend von \tilde{x} soll die exakte Lösung mittels des Korrekturansatzes

$$x = \tilde{x} + s$$

ermittelt werden. Der Korrekturvektor ist so zu bestimmen, dass die Gleichungen erfüllt sind, d.h.

$$Ax - b = A(\tilde{x} + s) - b = A\tilde{x} + As - b = 0.$$

Der Korrekturvektor s ergibt sich somit als Lösung von

$$As = b - A\tilde{x} = r.$$

Bei der numerischen Lösung dieser Korrekturgleichung erhalten wir im Allgemeinen eine wiederum fehlerhafte Korrektur $\tilde{s} \neq s$. Trotzdem erwarten wir, dass die Näherungslösung

$$\tilde{x} + \tilde{s}$$

„besser“ ist als \tilde{x} .

Bemerkung 1.5.1 *Wie wir in Kapitel 3 zeigen werden, genügt bei Spalten- oder Zeilenpivotisierung eine einzige Nachiteration, um eine Lösung auf eine dem Problem angepasste Genauigkeit zu erhalten.*

1.6 CHOLESKY-VERFAHREN

Wir wollen nun die Gauß-Elimination auf die eingeschränkte Klasse von Gleichungssystemen mit symmetrisch positiv definiten Matrizen anwenden. Es wird sich herausstellen, dass die Dreieckszerlegung in diesem Fall stark vereinfacht werden kann und dass dieser vereinfachte Algorithmus für große Matrizen nur den halben Aufwand an Operationen benötigt. Wir erinnern hier nochmals an die

Definition 1.6.1 *Eine symmetrische Matrix $A \in \mathbb{R}^{n \times n}$ heißt genau dann positiv definit, wenn $x^T Ax > 0$ für alle $x \in \mathbb{R}^n$ mit $x \neq 0$ ist.*

Positiv definite Matrizen haben folgende Eigenschaften.

Satz 1.6.2 (Eigenschaften positiv definiter Matrizen) *Für jede positiv definite Matrix $A \in \mathbb{R}^{n \times n}$ gilt:*

- i) A ist invertierbar,
- ii) $a_{ii} > 0$ für $i = 1, \dots, n$,
- iii) $\max_{i,j=1,\dots,n} |a_{ij}| = \max_{i=1,\dots,n} a_{ii}$,
- iv) Bei der Gauß-Elimination ohne Pivotsuche ist jede Restmatrix wiederum positiv definit.

Beweis. Wäre A nicht invertierbar, gäbe es einen Eigenwert $\lambda = 0$, da $\det A = 0$ nach Annahme. Dies steht aber im Widerspruch zur Voraussetzung, dass A positiv definit ist, da es ansonsten einen Eigenvektor $x \neq 0$ zu $\lambda = 0$ mit $Ax = \lambda x$ gäbe und somit dann auch $x^T Ax = 0$ gälte. Somit ist i) bewiesen. Setzt man für x einen kanonischen Einheitsvektor e_i ein, so folgt gerade $a_{ii} = e_i^T A e_i > 0$ und daher gilt die Aussage ii). Behauptung iii) sei als Hausübung überlassen. Um die verbleibende Behauptung iv) zu beweisen, schreiben wir $A = A^{(1)}$ in der Form

$$A^{(1)} = \left(\begin{array}{c|c} a_{11} & z^T \\ \hline z & B^{(1)} \end{array} \right),$$

wobei $z = (a_{12}, \dots, a_{1n})^T$ sei. Nach einem Eliminationsschritt ergibt sich

$$A^{(2)} = L_1 A^{(1)} = \left(\begin{array}{c|c} a_{11} & z^T \\ \hline 0 & B^{(2)} \end{array} \right) \quad \text{mit} \quad L_1 = \left(\begin{array}{cccc} 1 & & & \\ -l_{21} & 1 & & \\ \vdots & & \ddots & \\ -l_{n1} & & & 1 \end{array} \right).$$

Multipliziert man nun $A^{(2)}$ von rechts mit L_1^T , so wird auch z^T in der ersten Zeile eliminiert und die Teilmatrix $B^{(2)}$ bleibt unverändert, d.h.

$$L_1 A^{(1)} L_1^T = \left(\begin{array}{c|ccc} a_{11} & 0 & \cdots & 0 \\ \hline 0 & & & \\ \vdots & & B^{(2)} & \\ 0 & & & \end{array} \right).$$

Damit ist bewiesen, dass die Restmatrix $B^{(2)}$ symmetrisch ist. Können wir nun noch zeigen, dass $B^{(2)}$ auch wiederum positiv definit ist, so können wir unsere Argumentation sukzessive fortsetzen.

Es sei $y \in \mathbb{R}^{n-1}$ mit $y \neq 0$. Da L_1 regulär ist, gilt $x^T := (0 : y^T) L_1 \neq 0$. Nach Voraussetzung ist A positiv definit, also gilt

$$0 < x^T Ax = (0 : y^T) L_1 A L_1^T \begin{pmatrix} 0 \\ \cdots \\ y \end{pmatrix} = (0 : y^T) \left(\begin{array}{c|c} a_{11} & 0 \\ \hline 0 & B^{(2)} \end{array} \right) \begin{pmatrix} 0 \\ \cdots \\ y \end{pmatrix} = y^T B^{(2)} y.$$

Somit haben wir gezeigt, dass auch $B^{(2)}$ wieder positiv definit ist. □

Aufgabe 1.6.3 Man beweise Aussage iii) in Satz 1.6.2.

Mit Hilfe des letzten Satzes über die LR -Zerlegung können wir jetzt die **rationale** Cholesky-Zerlegung für symmetrische, positiv definite Matrizen herleiten.

Satz 1.6.4 (rationale Cholesky-Zerlegung) Für jede **symmetrische**, **positiv definite** Matrix A existiert eine **eindeutig bestimmte** Zerlegung der Form

$$A = LDL^T,$$

wobei L eine unipotente untere Dreiecksmatrix und D eine positive Diagonalmatrix ist.

Beweis. Wir setzen die Konstruktion im Beweis von Satz 1.6.2.iv für $k = 2, \dots, n - 1$ fort und erhalten so unmittelbar L als Produkt der $L_1^{-1}, \dots, L_{n-1}^{-1}$ und D als Diagonalmatrix der Pivotelemente. \square

Bemerkung 1.6.5 Man beachte die Eigenschaften der L_k bei der wirklichen Berechnung von L , d.h. das Produkt $L_1^{-1} \cdot \dots \cdot L_{n-1}^{-1}$ ist nicht wirklich durch Multiplikation auszurechnen (vgl. (1.20)).

Bemerkung 1.6.6 Da alle Einträge der Diagonalmatrix D positiv sind, existiert $D^{1/2} = \text{diag}(\pm\sqrt{d_i})$ und daher die Cholesky-Zerlegung

$$A = \bar{L} \bar{L}^T,$$

wobei \bar{L} die untere Dreiecksmatrix $\bar{L} := LD^{1/2}$ ist.

Bemerkung 1.6.7 Die Cholesky-Zerlegung ist nicht eindeutig, da $D^{1/2}$ nur bis auf Vorzeichen der Elemente in der Diagonalen eindeutig bestimmt ist und somit auch $\bar{L} = LD^{1/2}$ nicht eindeutig ist.

Zur Herleitung des Algorithmus zur Berechnung einer Cholesky-Zerlegung betrachten wir einfach die Gleichung

$$\begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & & \ddots & \\ l_{n1} & \dots & \dots & l_{nn} \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & \dots & l_{n1} \\ & l_{22} & & \vdots \\ & & \ddots & \vdots \\ & & & l_{nn} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & & & \vdots \\ \vdots & & & \vdots \\ a_{n1} & \dots & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} l_{11}^2 & l_{11}l_{21} & l_{11}l_{31} & \dots & l_{11}l_{n1} \\ l_{11}l_{21} & l_{21}^2 + l_{22}^2 & l_{21}l_{31} + l_{22}l_{32} & \dots & l_{21}l_{n1} + l_{22}l_{n2} \\ l_{11}l_{31} & l_{21}l_{31} + l_{22}l_{32} & l_{31}^2 + l_{32}^2 + l_{33}^2 & \dots & l_{31}l_{n1} + l_{32}l_{n2} + l_{33}l_{n3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{11}l_{n1} & l_{21}l_{n1} + l_{22}l_{n2} & l_{31}l_{n1} + l_{32}l_{n2} + l_{33}l_{n3} & \dots & \sum_{k=1}^n l_{nk}^2 \end{pmatrix}$$

d.h.

$$\text{für } i = k \quad \text{gilt } a_{kk} = \sum_{j=1}^k l_{kj}^2 \quad \text{und}$$

$$\text{für } i \neq k \quad \text{gilt } a_{ik} = \sum_{j=1}^k l_{ij} \cdot l_{kj} \quad (k + 1 \leq i \leq n)$$

und werten diese in einer geschickten Reihenfolge aus.

Cholesky-Verfahren zur Berechnung von L mit $A = LL^T$

Spaltenweise berechnet man für $k = 1, 2, \dots, n$

$$l_{kk} = \left(a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2 \right)^{1/2}$$

$$l_{ik} = \frac{1}{l_{kk}} \left(a_{ik} - \sum_{j=1}^{k-1} l_{ij} \cdot l_{kj} \right) \quad (k + 1 \leq i \leq n).$$

Bemerkung 1.6.8 i) Aus der Cholesky-Zerlegung $A = L \cdot L^T$ ergibt sich für alle $1 \leq k \leq n$ die Abschätzung

$$\sum_{j=1}^k l_{kj}^2 \leq \max_{1 \leq j \leq n} |a_{jj}|.$$

Folglich sind alle Elemente der Matrix L betragsweise durch $\max_{1 \leq j \leq n} \sqrt{|a_{jj}|}$ beschränkt. Die Elemente können damit nicht allzu stark anwachsen, was sich günstig auf die Stabilität des Verfahrens auswirkt.

ii) Da A symmetrisch ist, wird nur Information oberhalb und einschließlich der Hauptdiagonalen benötigt. Wenn man die Diagonalelemente l_{kk} separat in einem Vektor der Länge n speichert, und die Elemente l_{jk} , $k < j$ unterhalb der Diagonale, so kann man die Information der Matrix A bewahren.

iii) Bei der algorithmischen Durchführung der Cholesky-Zerlegung liefert das Verfahren auch die Information, ob die Matrix positiv definit ist. Man mache sich dies als Übungsaufgabe klar!

Aufgabe 1.6.9 Man beweise Aussage iii) in Bemerkung 1.6.8.

Satz 1.6.10 (Rechenaufwand Cholesky-Zerlegung) Sei $A \in \mathbb{R}^{n \times n}$ positiv definit. Dann gilt

$$\text{FLOP}(\text{Cholesky-Zerl. von } A) = \frac{2n^3 + 3n^2 - 5n}{6} + n \text{ Wurzelberechnung} = \frac{1}{3}n^3 + \mathcal{O}(n^2).$$

Beweis. Untersuchen wir nun die Komplexität der Cholesky-Zerlegung zuerst für einen k -ten Zerlegungsschritt und bestimmen dann den Gesamtaufwand.

i) Für den k -ten Zerlegungsschritt, d.h. die Berechnung der Elemente l_{ik} für festen Spaltenindex, sind jeweils $(k-1)$ Multiplikationen, $(k-1)$ Subtraktionen und eine Wurzelberechnung zur Berechnung des Diagonalelements nötig. Für jedes Nebendiagonalelement werden $(k-1)$ Multiplikationen, $(k-1)$ Subtraktionen und eine Division benötigt, d.h. bei $(n-k)$ Elementen unterhalb des k -ten Diagonalelements sind dies in der Summe

$$\begin{aligned} & (2(k-1) + (n-k)(2k-1)) \text{FLOP und 1 Wurzelberechnung} \\ & = (-2k^2 + (3+2n)k - (n+2)) \text{FLOP und 1 Wurzelberechnung} \end{aligned}$$

ii) Insgesamt ergibt sich dann für die Summe über alle Zerlegungsschritte

$$\begin{aligned} & n \text{ Wurzelberechnungen} + \sum_{k=1}^n (-2k^2 + (3+2n)k - (n+2)) \\ & = n \text{ Wurzelberechnungen} + \left(-2 \frac{(2n+1)(n+1)n}{6} + (3+2n) \frac{(n+1)n}{2} - (n+2)n \right) \\ & = n \text{ Wurzelberechnungen} + \left(\frac{-(4n^3 + 6n^2 + 2n) + (6n^3 + 15n^2 + 9n) - (6n^2 + 12n)}{6} \right) \\ & = n \text{ Wurzelberechnungen} + \frac{2n^3 + 3n^2 - 5n}{6}. \end{aligned}$$

Da eine einzelne Wurzelberechnung an Aufwand ein konstantes Vielfaches einer Gleitkommaoperation benötigt, kann man den Term n Wurzelberechnungen $+ (3n^2 - 5n)/6$ zu $\mathcal{O}(n^2)$ Operationen zusammenfassen. \square

Bemerkung 1.6.11 Für große n ist der Aufwand des Cholesky-Verfahrens im Vergleich zum Gauß-Algorithmus ungefähr halb so groß.



1.7 BANDGLEICHUNGEN

Eine weitere Klasse spezieller Matrizen, welche beim Lösen linearer Gleichungssysteme eine besondere Rolle spielen, sind Bandmatrizen. Man spricht von einer Bandmatrix A , falls alle von Null verschiedenen Elemente a_{ik} in der Diagonale und in einigen dazu benachbarten Nebendiagonalen stehen. Für die Anwendungen sind insbesondere die symmetrischen, positiv definiten Bandmatrizen wichtig.

Definition 1.7.1 Unter der **Bandbreite** m einer symmetrischen Matrix $A \in \mathbb{R}^{n \times n}$ versteht man die kleinste natürliche Zahl $m < n$, so dass gilt

$$a_{ik} = 0 \quad \text{für alle } i \text{ und } k \text{ mit } |i - k| > m.$$

Bemerkung 1.7.2 Die Bandbreite m gibt somit die Anzahl der Nebendiagonalen unterhalb bzw. oberhalb der Diagonalen an, welche die im Allgemeinen von Null verschiedenen Matrixelemente enthalten.

Bemerkung 1.7.3 In verallgemeinerter Form spricht man auch von unterer und oberer Bandbreite, welche sich auf suggestive Weise definieren.

Satz 1.7.4 Die untere Dreiecksmatrix L der Cholesky-Zerlegung $A = LL^T$ einer symmetrischen, positiv definiten Bandmatrix mit der Bandbreite m besitzt dieselbe Bandstruktur, denn es gilt

$$l_{ik} = 0 \quad \text{für alle } i, k \text{ mit } i - k > m.$$

Beweis. Ergibt sich direkt aus einer Hausübung zur Hülle einer Matrix. □

Aufgabe 1.7.5 Man beweise Satz 1.7.4.

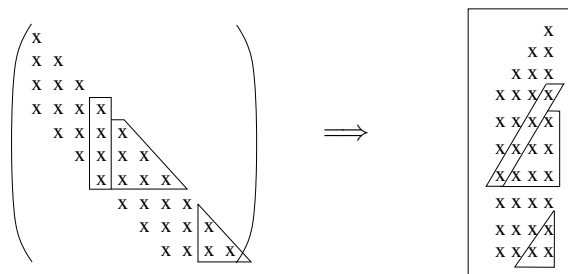


Abb. 1.1: Reduktion und Speicherung einer symmetrischen, positiv definiten Bandmatrix.

Analog zur Herleitung des Cholesky-Verfahrens auf Seite 19 erhalten wir die Rechenvorschrift des Cholesky-Verfahrens für Bandmatrizen wobei nun die Koeffizienten zeilenweise bestimmt werden. Für $n = 5$ und $m = 2$ sieht die Situation wie folgt aus:

$$\begin{pmatrix} l_{11} & 0 & 0 & 0 & 0 \\ l_{21} & l_{22} & 0 & 0 & 0 \\ l_{31} & l_{32} & l_{33} & 0 & 0 \\ 0 & l_{42} & l_{43} & l_{44} & 0 \\ 0 & 0 & l_{53} & l_{54} & l_{55} \end{pmatrix} = \begin{pmatrix} l_{11} & l_{21} & l_{31} & 0 & 0 \\ 0 & l_{22} & l_{32} & l_{42} & 0 \\ 0 & 0 & l_{33} & l_{43} & l_{53} \\ 0 & 0 & 0 & l_{44} & l_{54} \\ 0 & 0 & 0 & 0 & l_{55} \end{pmatrix} = \begin{pmatrix} l_{11}^2 & & & & \\ l_{21}l_{11} & l_{21}^2 + l_{22}^2 & & & \\ l_{31}l_{11} & l_{31}l_{21} + l_{32}l_{22} & l_{31}^2 + l_{32}^2 + l_{33}^2 & & \\ & l_{42}l_{22} & l_{42}l_{32} + l_{43}l_{33} & l_{42}^2 + l_{43}^2 + l_{44}^2 & \\ & & l_{53}l_{33} & l_{53}l_{43} + l_{54}l_{44} & l_{53}^2 + l_{54}^2 + l_{55}^2 \end{pmatrix}.$$

Cholesky-Verfahren zur Berechnung von L mit $A = LL^T$, wobei $A \in \mathbb{R}^{n \times n}$ mit Bandbreite $0 \leq m < n$.

Zeilenweise berechnet man für $k = 1, 2, \dots, n$

$$l_{ki} = \frac{1}{l_{ii}} \left(a_{ki} - \sum_{j=\max\{1, k-m\}}^{i-1} l_{kj} \cdot l_{ij} \right) \quad (\max\{1, k-m\} \leq i \leq k-1)$$

$$l_{kk} = \left(a_{kk} - \sum_{j=\max\{1, k-m\}}^{k-1} l_{kj}^2 \right)^{1/2}$$

Satz 1.7.6 Es sei $A \in \mathbb{R}^{n \times n}$ eine positiv definite Bandmatrix mit Bandbreite $1 \leq m \leq n$. Dann beträgt der Aufwand zur Berechnung der Cholesky-Zerlegung von A

$$\text{FLOP}(\text{Cholesky-Zerl. von } A) = n(m^2 + 2m) - \frac{4m^3 + 9m^2 + 5m}{6} + n \text{ Wurzelberechnungen.}$$

Beweis. Wendet man die Cholesky-Zerlegung auf eine positiv definite Bandmatrix A an, so ist die resultierende untere Dreiecksmatrix L mit $A = LL^T$ wieder eine Bandmatrix mit der gleichen Bandbreite, wie A sie hat.

- i) Gehen wir davon aus, dass die Matrix L zeilenweise berechnet wird und betrachten zuerst den Aufwand für die ersten m Zeilen. Hier kann man das Ergebnis aus Satz 1.6.10 verwenden. Es sind hierfür somit

$$m \text{ Wurzelberechnungen} + \frac{2m^3 + 3m^2 - 5m}{6}$$

Operationen notwendig.

- ii) In den verbleibenden $n - m$ Zeilen, in denen jeweils $m + 1$ Koeffizienten zu bestimmen sind, ergibt sich Folgendes:

Um den j -ten von Null verschiedenen Nebendiagonaleintrag von L in der k -ten Zeile $m + 1 \leq k \leq n$ zu berechnen, sind jeweils $j - 1$ Multiplikationen, $j - 1$ Subtraktionen und eine Division notwendig. Zur Berechnung des Diagonalelements werden m Multiplikationen, m Subtraktionen und eine Wurzelberechnung benötigt. Bei m Nebendiagonalelementen sind dies in der Summe

$$1 \text{ Wurzelberechnung} + 2m + \sum_{j=1}^m (2j - 1) = 1 \sqrt{\cdot} + m + 2 \sum_{j=1}^m j$$

$$= 1 \sqrt{\cdot} + m + m(m + 1) = 1 \sqrt{\cdot} + m(m + 2).$$

- iii) Der gesamte Aufwand beträgt also

$$n \text{ Wurzelberechnungen} + \frac{2m^3 + 3m^2 - 5m}{6} + (n - m)m(m + 2)$$

$$= n \text{ Wurzelberechnungen} + n(m^2 + 2m) - \frac{4m^3 + 9m^2 + 5m}{6}.$$

□

Bemerkung 1.7.7 Ist für eine Klasse von positiv definiten Bandmatrizen aus $\mathbb{R}^{n \times n}$ für beliebiges n die Bandbreite beschränkt, so wächst der Aufwand zur Berechnung der Cholesky-Zerlegung asymptotisch nur **linear** in n .

1.8 VANDERMOND-MATRIZEN

Bei vielen Approximations- und Interpolationsproblemen treten sogenannte Vandermond-Matrizen auf, d.h. Matrizen $V \in \mathbb{R}^{(n+1) \times (n+1)}$ von der Form

$$V = V(x_0, \dots, x_n) = \begin{pmatrix} 1 & 1 & \dots & 1 \\ x_0 & x_1 & \dots & x_n \\ \dots & \dots & \dots & \dots \\ x_0^n & x_1^n & \dots & x_n^n \end{pmatrix}.$$

Sei z.B. zu $(n + 1)$ -Daten (x_i, f_i) ein Polynom n -ten Grades $p \in \mathbb{P}_n$ gesucht mit der Eigenschaft $p(x_i) = f_i$ ($i = 0, \dots, n$). Wählt man den Ansatz

$$p(x) = \sum_{j=0}^n a_j x^j,$$

so führt dies zu einem LGS

$$V^T a = f$$

mit $a = (a_0, \dots, a_n)^T$ und $f = (f_0, \dots, f_n)^T$. Man mache sich dies an einem einfachen Beispiel klar!

Diesen Zusammenhang zwischen dem Interpolationsproblem und dem System $V^T a = f$ macht man sich auch beim Lösen zu Nutze. Die Bestimmung des Koeffizientenvektors a bzw. das Lösen des Gleichungssystems gliedert sich dabei in zwei Teile.

- Zuerst bestimmt man das Interpolationspolynom p in der Newtondarstellung, d.h.

$$\begin{aligned} p(x) &= \sum_{k=0}^n c_k \left(\prod_{j=0}^{k-1} (x - x_j) \right) \\ &= c_0 + c_1(x - x_0) + c_2(x - x_0)(x - x_1) + \dots + c_n(x - x_0) \dots (x - x_{n-1}) \end{aligned}$$

Die Konstanten c_k werden als dividierte Differenzen bezeichnet, sie lösen das lineare Gleichungssystem

$$\begin{pmatrix} 1 & & & & \\ 1 & (x_1 - x_0) & & & \\ 1 & (x_2 - x_0) & (x_2 - x_0)(x_2 - x_1) & & \\ \vdots & & & \ddots & \\ 1 & (x_n - x_0) & (x_n - x_0)(x_n - x_1) & \dots & \prod_{j=0}^{n-1} (x_n - x_j) \end{pmatrix} \begin{pmatrix} c_0 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} f_0 \\ \vdots \\ f_n \end{pmatrix}$$

und können wie folgt bestimmt werden.

```

c(0:n) = f(0:n)
for k = 0 : n - 1
    for i = n : -1 : k + 1
        c_i = (c_i - c_{i-1}) / (x_i - x_{i-k-1})
    end
end
    
```

(1.24)

- Im zweiten Schritt gewinnt man nun die a_i 's aus den c_i 's durch eine Basistransformation von $(1, (x - x_0), (x - x_0)(x - x_1), \dots, \prod_{j=0}^{n-1} (x - x_j))$ nach $(1, x, x^2, \dots, x^n)$. Die Transformationsmatrix U ergibt sich aus entsprechendem Koeffizientenvergleich

$$\begin{pmatrix} 1 & -x_0 & x_0x_1 & \dots & \prod_{j=0}^{n-1}(-x_j) \\ & 1 & -(x_0 + x_1) & & \\ & & 1 & & \vdots \\ & & & \ddots & \\ & & & & 1 \end{pmatrix} c = a. \quad (1.25)$$

Anstatt die Matrix U zu bestimmen, kann a wie folgt bestimmt werden.

```

a(0:n) = c(0:n)
for k = n - 1 : -1 : 0
    for i = k : n - 1
        a_i = a_i - x_k a_{i+1}
    end
end
end

```

(1.26)

Die Kombination beider Schritte liefert den folgenden Algorithmus.

Algorithmus 1.8.1: Polynominterpolation: $V^T a = f$

Gegeben seien $x(0 : n) \in \mathbb{R}^{n+1}$ mit paarweise verschiedenen Einträgen und $f = f(0 : n) \in \mathbb{R}^{n+1}$. Der folgende Algorithmus überschreibt f mit der Lösung $a = a(0 : n)$ zu dem Vandermondschen System

$$V(x_0, \dots, x_n)^T a = f$$

```

for k = 0 : n - 1
    for i = n : -1 : k + 1
        f(i) = (f(i) - f(i - 1)) / (x(i) - x(i - k - 1))
    end
end
for k = n - 1 : -1 : 0
    for i = k : n - 1
        f(i) = f(i) - f(i + 1)x(k)
    end
end
end

```

Satz 1.8.1 Das Verfahren 1.8.1 benötigt zum Lösen des Vandermondschen System

$$V(x_0, \dots, x_n)^T a = f$$

nur $5(n^2 + n)/2$ Operationen.

Beweis. Die Anweisung in der ersten inneren FOR-Schleife von Algorithmus 1.8.1 wird für ein $k \in \{0, \dots, n-1\}$ $(n-k)$ -mal ausgeführt, insgesamt also $\sum_{k=0}^{n-1} (n-k) = \sum_{k=1}^n k = n(n+1)/2$

-mal. Hierbei werden jedesmal 3 Gleitkommaoperationen benötigt. Für die zweite Doppelschleife ergeben sich, analog zu der obigen Betrachtung, $n(n+1)/2$ Aufrufe der inneren Anweisung, für die jeweils 2 Gleitkommaoperationen benötigt werden. \square

Bemerkung 1.8.2 Man beachte, dass für spezielle vollbesetzte Matrizen Verfahren existieren, die weniger als $\mathcal{O}(n^3)$ Operationen zur Berechnung der Lösung benötigen. Auch die Matrix braucht nicht vollständig gespeichert werden.

1.8.1 Das System $Vz = b$

Der Vollständigkeit wegen sei hier auch ein Verfahren für $Vz = b$ aufgeführt. Es sei $L_k(\alpha) \in \mathbb{R}^{(n+1) \times (n+1)}$ eine untere bidiagonale Matrix definiert durch

$$L_k(\alpha) = \left(\begin{array}{c|cccc} I_k & & & & 0 \\ \hline & 1 & & & 0 \\ & -\alpha & 1 & & \\ & & & \ddots & \ddots \\ 0 & \vdots & & \ddots & \ddots & \vdots \\ & & & & \ddots & 1 \\ & 0 & & \dots & -\alpha & 1 \end{array} \right)$$

und die Diagonalmatrix D_k sei definiert durch

$$D_k = \text{diag}(\underbrace{1, \dots, 1}_{(k+1)\text{-mal}}, x_{k+1} - x_0, \dots, x_n - x_{n-k-1}).$$

Mit diesen Definitionen kann man leicht aus (1.24) nachvollziehen, falls $f = f(0:n)$ und $c = c(0:n)$ der Vektor der dividierten Differenzen ist, dass $c = Lf$ gilt, wobei L die Matrix sei, die wie folgt definiert ist:

$$L = D_{n-1}^{-1} L_{n-1}(1) \cdots D_0^{-1} L_0(1).$$

Ähnlich erhält man aus (1.26), dass

$$a = Uc \tag{1.27}$$

gilt, wobei die durch (1.27) definierte Matrix U sich auch definieren lässt durch

$$U = L_0(x_0)^T \cdots L_{n-1}(x_{n-1})^T.$$

Mit diesen Definitionen ergibt sich sofort die Aussage $a = ULf$ mit $V^{-T} = UL$. Die Lösung von $Vz = b$ ergibt sich somit durch

$$\begin{aligned} z &= V^{-1}b = L^T(U^T b) \\ &= (L_0(1)^T D_0^{-1} \cdots L_{n-1}(1)^T D_{n-1}^{-1}) (L_{n-1}(x_{n-1}) \cdots L_0(x_0)) b. \end{aligned}$$

Diese Erkenntnis führt zu dem folgenden Algorithmus.

Algorithmus 1.8.2: Vandermondsches System: $Vz = b$

Gegeben seien $x(0:n) \in \mathbb{R}^{n+1}$ mit paarweise verschiedenen Einträgen und $b = b(0:n) \in \mathbb{R}^{n+1}$. Der folgende Algorithmus überschreibt b mit der Lösung $z = z(0:n)$ zu dem Vandermondschen System

$$V(x_0, \dots, x_n)z = f.$$

```

for k = 0 : n - 1
  for i = n : -1 : k + 1
    b(i) = b(i) - x(k)b(i - 1)
  end
end
for k = n - 1 : -1 : 0
  for i = k + 1 : n
    b(i) = b(i)/(x(i) - x(i - k - 1))
  end
  for i = k + 1 : n
    b(i) = b(i) - b(i + 1))
  end
end
end

```

Satz 1.8.3 *Das Verfahren 1.8.2 benötigt zum Lösen des Vandermondschen System*

$$V(x_0, \dots, x_n)a = f$$

nur $5(n^2 + n)/2$ Operationen.

Beweis. Analog zu den Betrachtungen im Beweis von Satz 1.8.1 erkennt man, dass die Anweisungen in den inneren FOR-Schleifen von Algorithmus 1.8.2 jeweils $n(n+1)/2$ -mal aufgerufen werden. Jedesmal werden 2 bzw. 1 Gleitkommaoperation benötigt. \square

2 ZAHLENDARSTELLUNG, FEHLERARTEN UND KONDITION

Im letzten Kapitel haben wir aus den Eingabegrößen (A, b) das Resultat $f(A, b) := A^{-1}b$ berechnet. Etwas abstrakter formuliert, wir haben eine Abbildung $f : U \subset X \rightarrow Y$ an einer Stelle $x \in U$ ausgewertet.



Bei exakter Arithmetik würden wir mit dem Gauß-Verfahren ein exaktes Ergebnis erhalten. Im Allgemeinen treten jedoch sowohl Fehler in der Eingabe als auch im Algorithmus auf, z.B. lässt sich die Zeit nur auf eine Genauigkeit von $10^{-14}s$ und die Masse auf $10^{-8}g$ genau bestimmen¹, Zahlen nur näherungsweise im Rechner als Gleitkommazahl (auch Fließkommazahl genannt, engl. floating point number) darstellen und z. B. auch die Darstellung

$$\sin(x) = \sum_{k=0}^{\infty} \frac{(-1)^k x^{2k+1}}{(2k+1)!} \quad (2.1)$$

nur näherungsweise auswerten. Welchen Einfluss haben nun beide Fehlerarten

- Eingabefehler (lassen sich nicht vermeiden) und
- Fehler im Algorithmus (wobei keine algorithmischen Fehler gemeint sind)

auf den Fehler im Resultat? Die Unterscheidung beider Arten von Fehlern führt uns zu den Begriffen

Kondition eines Problems und

Stabilität eines Algorithmus.

2.1 ZAHLENDARSTELLUNG

Ein erstes Beispiel soll motivieren, wieso wir uns mit der Frage, wie Zahlen im Rechner dargestellt werden, beschäftigen sollten.

Beispiel 2.1.1 Eine einfache Überlegung zeigt, dass das Ergebnis von $4 \cdot 13860^4 - 19601^4 + 2 \cdot 19601^2$ ungerade ist, eine erste Rechnung, dass an der letzten Stelle eine 1 oder 9 stehen muss. (Dabei könnte sich die 9 z.B. durch $-1 \cdot 10 + 1$ ergeben.)

$$\begin{aligned} 4 \cdot 13860^4 - 19601^4 + 2 \cdot 19601^2 &= 4(1386 \cdot 10)^4 - (1960 \cdot 10 + 1)^4 + 2(1960 \cdot 10 + 1)^2 \\ &= 4 \cdot 1386^4 \cdot 10^4 - (1960^4 \cdot 10^4 + 4 \cdot 1960^3 \cdot 10^3 + 6 \cdot 1960^2 \cdot 10^2 + 4 \cdot 1960 \cdot 10 + 1) + 2(1960^2 \cdot 10^2 + 2 \cdot 1960 \cdot 10 + 1) \\ &= 10^4(\dots) + 10^3(\dots) + 10^2(\dots) + 10^1(\dots) + 10^0(-1 + 2) \end{aligned}$$

Rechnet man auch die anderen Stellen aus (per Hand oder mit Maple) so ergibt sich

$$4 \cdot 13860^4 - 19601^4 + 2 \cdot 19601^2 = 1.$$

¹siehe www.ptb.de Physikalisch-technische Bundesanstalt

MATLAB-Beispiel:

Überraschenderweise liefert Matlab aber ein anderes Resultat.

```
>> myfun = @(x,y) 4*x^4-y^4+2*y^2;
>> myfun(13860,19601)
ans =
     2
```

Auch die folgende Realisierung mittels quadratischer Ergänzung liefert ein (anderes) falsches Ergebnis.

```
>> x = 13860; y = 19601;
>> u = 4 * x * x * x * x + 1;
>> v = y * y - 1;
>> u - v * v
ans =
     0
```



Und Vorsicht beim Umgang mit Maple! Schreibt man dort versehentlich nach einer Zahl noch einen Punkt, so ergibt es das scheinbar überraschende folgende Ergebnis. (Man beachte dabei, dass Maple, wenn man den Punkt hinzufügt, nicht mehr mit exakter Arithmetik rechnet, sondern Gleitkommazahlen verwendet, deren Genauigkeit man einstellen kann. Die Anweisung `Digits:=18` führt dazu, dass Maple Gleitkommazahlen mit 18 Stellen verwendet.

```
> restart:
4 *13860^4-19601^4+2*19601^2;

4. *13860^4-19601^4+2*19601^2;
Digits:=15:
4. *13860^4-19601^4+2*19601^2;
Digits:=16:
4. *13860^4-19601^4+2*19601^2;
Digits:=17:
4. *13860^4-19601^4+2*19601^2;
Digits:=18:
4. *13860^4-19601^4+2*19601^2;

1
68398402.
402.
2.
2.
1.
```

MATLAB-Beispiel:

Aus dem Topf der mathematischen Mirakel ziehen wir noch folgendes Beispiel:

```
>> floor(114)
ans =
    114
>> floor(1.14*100)
ans =
    113
```

Die Routine `floor` rundet dabei auf die nächstkleinere ganze Zahl. 114 auf die nächstkleinere ganze Zahl gerundet sollte also in beiden Fällen 114 ergeben!

Ein Rechner verwendet als einfachste Form der Zahlendarstellung Festkommazahlen. Dabei wird eine (meistens) begrenzte Ziffernfolge gespeichert und an festgelegter Stelle das Komma angenommen. Bei vielen Anwendungen reicht der darstellbare Zahlenbereich jedoch nicht aus, so dass es naheliegend ist, bei jedem Wert die genaue Stelle des Kommas zusätzlich zu speichern. Das bedeutet mathematisch nichts anderes als die Darstellung der Zahl x mit zwei Werten, der Mantisse m und dem Exponenten e . Wir bekommen somit die Darstellung $x = m \cdot 10^e$. Diese Darstellung ist jedoch nicht eindeutig, siehe z.B. $1014 = 1.014 \cdot 10^3 = 10140 \cdot 10^{-1}$. Legt man jedoch den Wertebereich der Mantisse m geschickt fest, z. B. $1 \leq m < 10$, so erhält man eine eindeutige Darstellung (wenn man endlich viele Stellen hat). Diesen Schritt bezeichnet man als **Normalisierung** der Mantisse. Wir haben soeben stillschweigend vorausgesetzt, dass wir die Basis 10 verwenden¹. Darauf ist man jedoch nicht festgelegt. Pascal² erkannte, dass man jede Basis $b \geq 2$ verwenden kann. Heutige Rechner verwenden meist binäre Zahlendarstellungen, d.h. zur Basis 2. Betrachten wir nun die Zahlendarstellung noch etwas allgemeiner.

Definition 2.1.2 (*b*-adischer Bruch) Es sei $b \geq 2$ eine natürliche Zahl, $n \in \mathbb{N}_0$, $a_k \in \mathbb{N}_0$ mit $0 \leq a_k < b$ und $a_n \neq 0$. Eine Reihe der Gestalt

$$\pm \sum_{k=-\infty}^n a_k b^k$$

wird als *b*-adischer Bruch bezeichnet. Falls die Basis festgelegt ist, kann man einen *b*-adischen Bruch auch durch die Reihung der Ziffern a_k angeben:

$$\pm a_n a_{n-1} \cdots a_1 a_0 . a_{-1} a_{-2} a_{-3} \cdots$$

Bemerkung 2.1.3 i.) Für $b = 10$ spricht man von **Dezimalbrüchen**, für $b = 2$ von **dyadischen Brüchen** und die Babylonier³ verwendeten ein System zur Basis $b = 60$, was sich heute noch im Verhältnis von Sekunde, Minute und Stunde wiederfindet.

ii.) Die Basis wird auch als Grundzahl bezeichnet, weswegen man auch die Bezeichnung **g-adisch** findet.

iii.) Vorsicht, nicht **b-adisch** mit **p-adisch** verwechseln ($p \sim$ Primzahl)!



iv.) Der Null kommt immer eine Sonderrolle zu. Sie wird gesondert betrachtet, um die Darstellung einfacher zu halten. Ohne den Sonderfall der Null jeweils aufzuführen, schließen wir ihn indirekt mit ein.

Die nächsten beiden Sätze besagen, dass sich jede reelle Zahl durch einen *b*-adischen Bruch darstellen lässt und umgekehrt.

Satz 2.1.4 (Konvergenz *b*-adischer Bruch) Jeder *b*-adische Bruch stellt eine Cauchy-Folge dar, d.h. dass jeder *b*-adische Bruch gegen eine reelle Zahl konvergiert.

¹**Dezimalsystem** Unter dem Einfluss der Schrift „De Thiende“ (Leiden 1585) von Simon Stevin (1548-1620) setzte sich im 16. Jahrhundert das Rechnen mit Dezimalbrüchen in Westeuropa durch.

²**Blaise Pascal**, 1623-1662

³**Sexagesimalsystem** Zählt man mit den Fingern der linken Hand wie üblich von 1 bis 5 und registriert die an der linken Hand gezählten 5-er Gruppen durch das Drücken des Daumens der rechten Hand auf das passende Fingerglied an der rechten Hand, so kommt man zum $5 \times 12 = 60$ -System. Diese Zählweise wurde schon von den Vorfahren der Sumerer im 4 Jahrtausend v. Chr. angewandt. Die Babylonier (1750 v. Chr.) verwendeten dann das Sexagesimalsystem in dem Sinne, dass sie 60 Herzschläge zu einer Minute zusammenfassten und 60 Minuten zu einer Stunde.

Beweis. Wir beschränken uns auf den Fall eines nicht-negativen b -adischen Bruchs $\sum_{k=-\infty}^n a_k b^k$. Wir definieren für $-\ell \leq n$ die Partialsummen

$$s_\ell := \sum_{k=-\ell}^n a_k b^k.$$

Es ist nun zu zeigen, dass die Folge $(s_\ell)_{-\ell \leq n}$ eine Cauchy-Folge ist. Sei $\varepsilon > 0$ und $L \in \mathbb{N}$ so groß, dass $b^{-L} < \varepsilon$ gelte. Dann gilt für $\ell \geq m \geq L$

$$\begin{aligned} |s_\ell - s_m| &= \sum_{k=-\ell}^n a_k b^k - \sum_{k=-m}^n a_k b^k = \sum_{k=-\ell}^{-m-1} a_k b^k \leq \sum_{k=-\ell}^{-m-1} (b-1)b^k \\ &= (b-1)b^{-m-1} \sum_{k=-\ell+m+1}^0 b^k = (b-1)b^{-m-1} \sum_{k=0}^{\ell-m-1} b^{-k} \\ &< (b-1)b^{-m-1} \frac{1}{1-b^{-1}} = b^{-m} \leq b^{-L} < \varepsilon. \end{aligned}$$

□

Satz 2.1.5 (Reelle Zahl als b -adischer Bruch) Sei $b \geq 2$ eine natürliche Zahl. Dann lässt sich jede reelle von Null verschiedene Zahl in einen b -adischen Bruch entwickeln.

Beweis. Es genügt wieder, nur den Fall $x \geq 0$ zu zeigen. Zuerst zeigen wir, dass es zu jeder Basis $b \geq 2$ mindestens eine Zahl $m \in \mathbb{N}$ gibt mit $x < b^m$. Nach dem Archimedischen Axiom gibt es ein $m \in \mathbb{N}$ mit $m > (x-1)/(b-1) \in \mathbb{R}$. Mit diesem m und der Bernoullischen Ungleichung erhalten wir

$$b^m = (1 + (b-1))^m \geq 1 + m(b-1) > 1 + x - 1 = x.$$

Sei $n := \min\{k \in \mathbb{N}_0 \mid 0 \leq x < b^{k+1}\}$. Durch vollständige Induktion konstruieren wir nun eine Folge $(a_{-\ell})_{\ell \geq -n}$ natürlicher Zahlen $0 \leq a_{-\ell} < b$, so dass für alle $m \geq -n$ gilt

$$x = \sum_{\ell=-m}^n a_\ell b^\ell + \nu_{-m} \quad \text{mit } 0 \leq \nu_{-m} < b^{-m}.$$

Da $\lim_{m \rightarrow \infty} \nu_{-m} = 0$ gilt, folgt $x = \sum_{\ell=-\infty}^n a_\ell b^\ell$ und somit die Behauptung.

Induktionsanfang $-m = n$: Es gilt $0 \leq x b^{-n} < b$. Somit existiert ein $a_n \in \{0, 1, \dots, b-1\}$ und ein $\delta \in \mathbb{R}$ mit $0 \leq \delta < 1$, so dass $x b^{-n} = a_n + \delta$ gilt. Mit $\nu_n := \delta b^n$ gewinnt man

$$x = a_n b^n + \nu_n \quad \text{mit } 0 \leq \nu_n < b^n.$$

Induktionsschritt $(-m) \rightarrow (-m-1)$: Es gilt $0 \leq \nu_{-m} b^{m+1} < b$, also gibt es ein $a_{-m-1} \in \{0, 1, \dots, b-1\}$ und ein $\delta \in \mathbb{R}$ mit $0 \leq \delta < 1$, so dass $\nu_{-m} b^{m+1} = a_{-m-1} + \delta$ gilt. Mit $\nu_{-m-1} := \delta b^{-m-1}$ gewinnt man

$$x = \sum_{\ell=-m}^n a_\ell b^\ell + \nu_{-m} = \sum_{\ell=-m}^n a_\ell b^\ell + (a_{-m-1} + \delta) b^{-m-1} = \sum_{\ell=-m-1}^n a_\ell b^\ell + \nu_{-m-1},$$

wobei $0 \leq \nu_{-m-1} < b^{-m-1}$ gilt. □

Bemerkung 2.1.6 Die b -adische Darstellung ist in der gewählten Form nicht eindeutig, siehe z. B. $1 = 0.\bar{9} = 0.999\dots$, falls $b = 10$ gewählt ist. Durch die zusätzliche Bedingung „ $a_k < b-1$ für unendlich viele $k \leq n$ “ erhält man eine eindeutige Darstellung. Somit wäre z.B. für $b = 10$ die Darstellung $0.\bar{9} = 0.999\dots$ nicht zugelassen.

Bemerkung 2.1.7 Da es unmöglich ist, reelle Zahlen als unendliche b -adische Brüche zu speichern, werden reelle Zahlen auf Rechnern durch eine endliche Entwicklung approximiert.



Definition 2.1.8 (Festpunktzahl) Zu gegebener Basis b ist

$$F(\ell, n) := \left\{ \pm \sum_{k=-\ell}^n a_k b^k = \pm (a_n \dots a_1 a_0 \cdot a_{-1} \dots a_{-\ell})_b \right\}$$

die Menge der Festpunktzahlen mit $(n+1)$ Vor- und ℓ Nachkommastellen.

Bemerkung 2.1.9 i.) Typischerweise werden ganze Zahlen (engl. integers) auf Computern dargestellt mit $\ell = 0$ und $b = 2$.

ii.) In modernen Rechnern werden negative Zahlen häufig als 2-er-Komplement⁴ abgespeichert.

Beispiel 2.1.10 Binärdarstellung ($b = 2$) von $x = (14.625)_{10} = (1110.101)_2$

$$\begin{aligned} 14 &= 7 \cdot 2 + 0 && \Rightarrow a_0 = 0 \\ 7 &= 3 \cdot 2 + 1 && \Rightarrow a_1 = 1 \\ 3 &= 1 \cdot 2 + 1 && \Rightarrow a_2 = 1 \\ 1 &= 0 \cdot 2 + 1 && \Rightarrow a_3 = 1 \end{aligned}$$

Im Folgenden bezeichnet $\lfloor \cdot \rfloor$ die untere Gaußklammer, d.h.

$$\lfloor x \rfloor := \max \{j \in \mathbb{Z} : j \leq x\}.$$

$$\begin{aligned} a_{-1} &= \lfloor 2 * 0.625 \rfloor = \lfloor 1.25 \rfloor = 1 \\ a_{-2} &= \lfloor 2 * 0.25 \rfloor = \lfloor 0.5 \rfloor = 0 \\ a_{-3} &= \lfloor 2 * 0.5 \rfloor = \lfloor 1 \rfloor = 1 \end{aligned}$$

Bemerkung 2.1.11 Da im 2-er-Komplement der Wert 0 den positiven Zahlen zugeordnet wird, erhält man für 8-Bit⁵, 16-Bit und 32-Bit, bzw. 1-, 2- und 4-Byte⁶ Darstellungen die Zahlenbereiche

$$\begin{aligned} 8 - \text{bit}, n = 7 : & 0, \dots, 255 \text{ ohne Vorzeichen} \\ & \text{bzw. } -128, \dots, 127 \text{ (8 Stellen)} \end{aligned}$$

$$\begin{aligned} 16 - \text{bit}, n = 15 : & 0, \dots, 65535 \text{ ohne Vorzeichen} \\ & \text{bzw. } -32768, \dots, 32767 \text{ (16 Stellen)} \end{aligned}$$

$$\begin{aligned} 32 - \text{bit}, n = 31 : & 0, \dots, 4.294.967.295, \dots, \text{ ohne Vorzeichen} \\ & \text{bzw. } -2.147.483.648, \dots, 2.147.483.647 \text{ (32 Stellen)} \end{aligned}$$

⁴**Zweierkomplement** ist eine Möglichkeit, negative Zahlen im Dualsystem darzustellen. Dabei werden keine zusätzlichen Symbole wie + und - benötigt. Da im Zweierkomplement der Wert 0 den positiven Zahlen zugerechnet wird, umfasst der Wertebereich bei n binären Stellen allgemein den Bereich $-2^{n-1}, \dots, 0, \dots, 2^{n-1} - 1$. Zum Beispiel bei 4 Bit von -8 bis 7. Negative Zahlen gewinnt man aus positiven Zahlen, indem sämtliche binären Stellen negiert werden und zu dem Ergebnis der Wert 1 addiert wird. Somit wird aus $(4)_{10} = (0100)_2$ durch Invertieren $\text{Not}(0100) = 1011$ und Addition der 1 $(1011)_2 + (0001)_2 = (1100)_2$. Addiert man nun 4 und -4 so gilt $(4)_{10} + (-4)_{10} = (0100)_2 + (1100)_2 = (0000)_2 = (0)_{10}$, wobei die vorderste Stelle verworfen wurde.

⁵**Bit** ~ (engl.: binary digit, lat. digitus ~ Finger)

⁶**Byte** ~ Maßeinheit für eine Datenmenge von 8 Bit

Bemerkung 2.1.12 Als Problem stellt sich dabei heraus, dass damit nur ein kleiner Zahlenbereich darstellbar ist.

Definition 2.1.13 (Gleitpunkt-Darstellung) Die normalisierte Gleitpunkt-Darstellung einer reellen Zahl $x \in \mathbb{R}$ hat die Form $x = f \cdot b^\ell$, wobei

- die **Mantisse** f ist eine feste Zahl, die m Stellen hat und

$$b^{-1} \leq |f| < 1, \text{ falls } x \neq 0 \text{ und } f = 0 \text{ falls, } x = 0$$

genügt, d.h.

$$f = \begin{cases} \pm (0.d_1 \dots d_m)_b = \pm \sum_{j=1}^m d_j b^{-j} & , d_1 \neq 0 \\ 0 & , d_1 = 0 \end{cases}$$

- der **Exponent** ℓ ist eine ganze Zahl innerhalb fester Schranken

$$r \leq \ell \leq R,$$

welche sich in der Form

$$\ell = \pm (\ell_{n-1} \dots \ell_0)_b = \pm \sum_{j=0}^{n-1} \ell_j b^j$$

darstellen lässt. (Für 8-stelligen binären Exponenten wird z.B. $r = -128$ und $R = 127$ gewählt.)

Die Menge der Maschinenzahlen mit Basis b , m -stelliger Mantisse und n -stelligem Exponenten wird als $\mathbb{M}(b, m, n)$ bezeichnet.

Beispiel 2.1.14 $\mathbb{M}(2, 3, 1)$ besteht aus den Zahlen

$$x = \pm(d_{-1} \cdot 2^{-1} + d_{-2} \cdot 2^{-2} + d_{-3} \cdot 2^{-3}) \cdot 2^{\pm \ell_0},$$

wobei $d_{-1} \neq 0 \vee x = d_{-1} = d_{-2} = d_{-3} = 0$ gilt. Somit

$$\mathbb{M}(2, 3, 1) = \left\{ 0, \pm \frac{1}{4}, \pm \frac{5}{16}, \pm \frac{3}{8}, \pm \frac{7}{16}, \pm \frac{1}{2}, \pm \frac{5}{8}, \pm \frac{3}{4}, \pm \frac{7}{8}, \pm 1, \pm \frac{5}{4}, \pm \frac{3}{2}, \pm \frac{7}{4} \right\}.$$

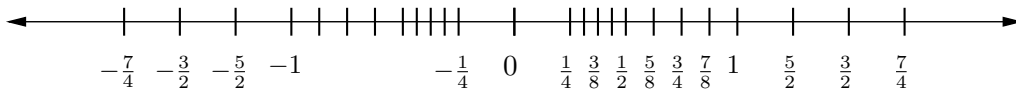


Abb. 2.1: Graphische Darstellung von $\mathbb{M}(2, 3, 1)$.

Man beachte den kleiner werdenden Abstand zur Null hin und den Abstand zwischen den betragsmäßig kleinsten Zahlen und der Null.

Bemerkung 2.1.15 Die betragsmäßig kleinste und größte Zahl in $\mathbb{M}(b, m, n)$ bezeichnen wir mit x_{\min} bzw. x_{\max} ,

$$\begin{aligned} x_{\min} &:= \min \{ |z| \in \mathbb{M}(b, m, n) \}, & |x| < x_{\min} &\Rightarrow x = 0 \\ x_{\max} &:= \max \{ |z| \in \mathbb{M}(b, m, n) \}, & |x| > x_{\max} &\Rightarrow \text{overflow} \end{aligned} \tag{2.2}$$

Beispiel 2.1.16

$$\begin{array}{lll}
 \mathbb{M}(10, 4, 2), & x = -51.34 & \rightsquigarrow -0.5134 \cdot 10^2 \\
 \mathbb{M}(2, 4, 2), & x = 3.25 = (11.01)_2 & \rightsquigarrow (0.1101)_2 \cdot 2^{(10)}_2 \\
 \mathbb{M}(2, 4, 2), & x = 14.625 = (1110.101)_2 & \rightsquigarrow (0.1110101)_2 \cdot 2^{(100)}_2 \notin \mathbb{M}(2, 4, 2)
 \end{array}$$

Bemerkung 2.1.17 Üblicherweise werden 8 Byte (64 Bit) Speicherplatz für eine Gleipunktzahl wie folgt aufgeteilt,

- 1 Bit für das Vorzeichen,
- 11 Bit für den Exponenten, d.h. Werte sind darstellbar im Bereich $-1022 \dots 1023$ ($1 \dots 2046$ - Bias⁷-Wert 1023). Die Werte -1023 und 1024 , bzw. 0 und $2047 = 2^{11} - 1$ als Charakteristik sind reserviert für die speziellen Zahlenwerte „Null“, „Unendlich“ und „NaN“,
- 52 Bit für die Mantisse. Dies sind eigentlich Werte zwischen $(0 \dots 01)_2$ und $(1 \dots 1)_2$, da die Null gesondert abgespeichert wird. Da die erste Null jedoch redundant ist, gewinnt man noch eine Stelle hinzu, wenn man in der normalisierten Darstellung ist. Also muss auch noch irgendwie festgehalten werden, ob man in der normalisierten Darstellung ist oder nicht.

Exponent e	Mantisse	Bedeutung	Bezeichnung
$e = 0$	$m = 0$	± 0	Null
$e = 0$	$m > 0$	$\pm 0, m \times 2^{1-Bias}$	denormalisierte Zahl
$0 < e < e_{max}$	$m \geq 0$	$\pm 1, m \times 2^{e-Bias}$	normalisierte Zahl
$e = e_{max}$	$m = 0$	$\pm \infty$	Unendlich (Inf)
$e = e_{max}$	$m > 0$	keine Zahl	Not a Number (NaN)

Tab. 2.1: Interpretation des Zahlenformats nach IEEE-754

Typ	Größe	Exponent	Mantisse	Werte des Exponenten (e)
single	32 Bit	8 Bit	23 Bit	$-126 \leq e \leq 127$
double	64 Bit	11 Bit	52 Bit	$-1022 \leq e \leq 1023$
	rel. Abstand zweier Zahlen	Dezimalstellen	betragsmäßig kleinste Zahl	größte Zahl
single	$2^{-(23+1)}$	7-8	$2^{-23} \cdot 2^{-126}$	$(1 - 2^{-24})2^{128}$
double	$2^{-(52+1)}$	15-16	$2^{-52} \cdot 2^{-1022}$	$(1 - 2^{-53})2^{1024}$

Tab. 2.2: Zahlenformat nach IEEE-754

Bemerkung 2.1.18 Der Standardtyp in Matlab, wenn nichts Weiteres angegeben wird, ist double.

⁷(oder auch Offset) hier zu verstehen als ein konstanter Wert, der zu etwas hinzuaddiert bzw. subtrahiert wird

Mit Hilfe der folgenden mex-Routine `digits.c` lassen sich auf einem 32-Bit Rechner unter Matlab die einzelnen Bits einer double precision-Variablen auslesen. (Kompilieren der c-Routine auf der Kommandozeilenebene von Matlab mit `mex digits.c.`)

C-Funktion: `digits.c`

```

1 #include "mex.h"
2 bool get_bit_at(double y[], int pos)
3 {
4     unsigned long bitmask;
5     unsigned long *cy = (unsigned long *) y;
6
7     if (pos >= 0) {
8         if (pos < 32) {
9             bitmask = 1 << pos; // bitmask setzen
10            return (cy[0] & bitmask) != 0;
11        } else if (pos < 64) {
12            bitmask = 1 << (pos-32); // bitmask setzen
13            // unsigned long ist nur
14            // 4 Byte = 32 Bit lang
15            return (cy[1] & bitmask) != 0;
16        }
17    }
18    mexErrMsgTxt("No valid position for Bit.");
19 }
20
21 void mexFunction( int nlhs, mxArray *plhs[],
22                  int nrhs, const mxArray *prhs[] )
23 {
24     double *x,*y;
25     mwSize mrows,ncols;
26     int k;
27     /* Check for proper number of arguments. */
28     if(nrhs!=1) {
29         mexErrMsgTxt("One input required.");
30     } else if(nlhs>1) {
31         mexErrMsgTxt("Too many output arguments");
32     }
33     /* The input must be a noncomplex scalar double.*/
34     mrows = mxGetM(prhs[0]);
35     ncols = mxGetN(prhs[0]);
36     if( !mxIsDouble(prhs[0]) || mxIsComplex(prhs[0]) ||
37         !(mrows==1 && ncols==1) ) {
38         mexErrMsgTxt("Input must be a noncomplex scalar double.");
39     }
40     /* Create matrix for the return argument. */
41     plhs[0] = mxCreateDoubleMatrix(1,64, mxREAL);
42     /* Assign pointers to each input and output. */
43     x = mxGetPr(prhs[0]);
44     y = mxGetPr(plhs[0]);
45     /* Call the digit subroutine. */
46     for (k=0; k<64; k++) y[63-k] = (double) get_bit_at(x,k);
47 }

```


Versuchen Sie sich die Ergebnisse der folgenden Matlab-Ausgaben zu erklären bzw. testen Sie es auf Ihrem eigenen Rechner.

```
>> A=digits(-1);A(1)
ans =
    1
>> A=digits(1);A(1)
ans =
    0
>> A=digits(2^(1023)*(1.9999999999999998));
>> A(1),A(2:12),[A(13:29); A(30:46); A(47:63)],A(64)
ans =
    0
ans =
    1  1  1  1  1  1  1  1  1  1  1  0
ans =
    1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1
    1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1
    1  1  1  1  1  1  1  1  1  1  1  1  1  1  1  1
ans =
    1
>> A=digits(2^(-971)*(2-1.9999999999999995));
>> A(1),A(2:12),[A(13:29); A(30:46); A(47:63)],A(64)
ans =
    0
ans =
    0  0  0  0  0  0  0  0  0  0  0  1
ans =
    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
ans =
    0
>> A=digits(2^(-972)*(2-1.9999999999999995));
>> A(1),A(2:12),[A(13:29); A(30:46); A(47:63)],A(64)
ans =
    0
ans =
    0  0  0  0  0  0  0  0  0  0  0
ans
    1  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
    0  0  0  0  0  0  0  0  0  0  0  0  0  0  0  0
ans =
    0
>> A=digits(realmax)           % try by yourself
>> A=digits(realmin)
>> A=digits(realmin/2^52)
>> A=digits(2^0+2^(-1)); A(13:29)
>> A=digits(2^0+2^(-1)+2^(-2)); A(13:29)
>> A=digits(2^0+2^(-1)+2^(-2)+2^(-3)); A(13:29)
>> A=digits(2^0+2^(-1)+2^(-2)+2^(-3)+2^(-17))A(13:29)1
```

Bemerkung 2.1.19 i) Die Werte von b, m, ℓ hängen in der Regel vom Typ des jeweiligen Computers ab, ebenso r, R . Der IEEE⁸-Standard versucht dies zu vereinheitlichen.

ii) Bei der Umwandlung „exakter“ Eingaben in Maschinenzahlen können Fehler auftreten, z.B.

$$x = 0.1 = (0.000\overline{1100})_2$$

2.2 RUNDUNG UND MASCHINENGENAUIGKEIT

Aufgabe: Approximiere $x \in \mathbb{R}$ durch $fl(x) \in \mathbb{M}(b, m, \ell)$ („float“)

$$fl : \mathbb{R} \rightarrow \mathbb{M}(b, m, \ell)$$

bzw.

$$fl : [-x_{\max}, x_{\max}] \rightarrow \mathbb{M}(b, m, \ell)$$

Dies wird durch geeignete Rundungsstrategien gewährleistet. Wir betrachten zwei konkrete Beispiele:

2.2.1 Standardrundung, kaufmännische Rundung

(entspricht „normalem Runden“ $1.4 \rightarrow 1, 1.6 \rightarrow 2$)

Die **Standardrundung** ist definiert durch

$$fl(x; b, m, \ell) = fl(x) := \pm \sum_{j=1}^m d_j b^{\ell_j} + \begin{cases} 0 & , \text{ falls } d_{m+1} < \frac{b}{2} \\ b^{\ell_m} & , \text{ falls } d_{m+1} \geq \frac{b}{2} \end{cases}$$

und

$$|x| < x_{\min} \Rightarrow fl(x) := 0$$

$$|x| > x_{\max} \Rightarrow \text{overflow} .$$

Bemerkung 2.2.1 Man mache sich klar, dass die Rundung im Falle $d_{m+1} = b/2$ (b gerade, $b/2$ ungerade) rein willkürlich ist.

2.2.2 „Round to even“ (oder auch Banker’s Rule oder mathematisch unverzerrte Rundung)

Die unverzerrte Rundung (englisch „round to even“), auch als Banker’s Rule bekannt, ist eine Rundungsregel, welche sich dann von der Standardrundung unterscheidet, wenn der zu rundende Rest genau zwischen zwei Zahlen mit der gewählten Anzahl von Ziffern liegt. Voraussetzung zur Anwendung der Banker’s Rule ist, dass die Basis b gerade und $b/2$ ungerade ist, dies ist z.B. für die interessanten Fälle $b = 2, 10$ erfüllt. Die Standardrundung erzeugt statistische „Fehler“, da das Aufrunden um $b/2$ vorkommt, das Abrunden um $b/2$ jedoch nie. Rundet man z.B. bei großen Mengen von Daten, bei denen der Trend untersucht werden soll, immer auf, falls der zu rundende Rest genau $b/2$ ist, so verschiebt sich der Trend aufwärts. Außerdem ist das Verhalten bei positiven und negativen Zahlen unterschiedlich, wenn die zu rundende Ziffer $b/2$ ist. Das „Round to Even“ vermeidet dieses Verhalten. Es rundet von der genauen Mitte zwischen zwei Ziffern immer zur nächsten geraden Zahl auf oder ab. Ansonsten entspricht es der Standardrundung.

⁸IEEE~ Institute of Electrical and Electronics Engineers

Bemerkung 2.2.2 Man könnte genauso gut auch zur nächsten ungeraden Zahl runden, dann würde man jedoch nie auf Null runden.

Banker's Rule: Voraussetzung: b gerade, $b/2$ ungerade. Es sei m die vorgegebene Mantissenlänge und $x = \pm(0.d_1d_2d_3\dots)_b b^{e_1+1}$ ($d_1 \neq 0, d_k < b - 1$ für unendlich viele $k \in \mathbb{N}$). Dann ist die Banker's Rule definiert durch

$$fl(x) := \pm \sum_{j=1}^m d_j b^{e_j} + \begin{cases} 0 & , \text{ falls } d_{m+1} < b/2 \\ b^{e_m} & , \text{ falls } d_{m+1} > b/2 \text{ oder} \\ & d_{m+1} = b/2 \text{ und es ex. ein } k > m + 1 \text{ mit } d_k \neq 0 \\ 0 & \sum_{k=m+1}^{\infty} d_k b^{e_k} = b^{e_m}/2 \text{ und } d_m \text{ gerade} \\ b^{e_m} & \sum_{k=m+1}^{\infty} d_k b^{e_k} = b^{e_m}/2 \text{ und } d_m \text{ ungerade} \end{cases}$$

Beispiel 2.2.3 Es sei $b = 10$. Dann erhält man mit der Banker's Rule folgende Ergebnisse:

aus 2.33 wird 2.3,	aus -2.33 wird -2.3 ,
aus 2.35 wird 2.4,	aus -2.45 (exakt) wird -2.4 ,
aus 2.45 (exakt) wird 2.4,	aus -2.4500001 wird -2.5 ,
aus 2.4500001 wird 2.5 und	aus -2.449 wird -2.4 .

Bemerkung 2.2.4 i.) Der Vorteil der Banker's Rule liegt in einer „besseren“ statistischen Verteilung der Rundungsfehler.

ii.) Zur Basis 10 lautet die Regel: Folgt auf die letzte beizubehaltende Ziffer lediglich eine 5 (oder eine 5, auf die nur Nullen folgen), so wird derart gerundet, dass die letzte beizubehaltende Ziffer gerade wird.

Beispiel 2.2.5 (Wiederholtes Runden)

Standard	2.445 \rightarrow 2.45 \rightarrow 2.5 \rightarrow 3	(Fehler 0.555)
	1.445 \rightarrow 1.45 \rightarrow 1.5 \rightarrow 2	(Fehler 0.555)
Banker's Rule	2.445 \rightarrow 2.44 \rightarrow 2.4 \rightarrow 2	(Fehler 0.445)
	1.445 \rightarrow 1.44 \rightarrow 1.4 \rightarrow 1	(Fehler 0.445)

Definition 2.2.6 (absoluter/relativer Fehler) Es sei $x \in \mathbb{R}$ und $\tilde{x} \in \mathbb{R}$ eine Approximation.

(i) Die Größe $|x - \tilde{x}|$ heißt **absoluter Fehler**.

(ii) Die Größe $\frac{|x - \tilde{x}|}{|x|}$ heißt **relativer Fehler**.

Lemma 2.2.7 Die Basis b sei gerade. Es sei $x = \pm b^\ell \sum_{k=1}^{\infty} d_k b^{-k}$ ($0 < d_1 < b$) und $fl(x)$ die Standardrundung von x auf m Stellen. Dann gilt

$$|fl(x) - x| \leq \frac{b^{1-m}}{2} \cdot b^{\ell-1}$$

und

$$\frac{|fl(x) - x|}{|x|} \leq \frac{b^{1-m}}{2}.$$

Beweis. O.B.d.A sei x positiv. Sei $n \in \mathbb{N}$. Für $b > 1$ und $0 \leq d_k \leq b - 1$ ($k = n, n + 1, \dots$) gilt

$$\sum_{k=n}^{\infty} d_k b^{-k} \leq \sum_{k=n}^{\infty} (b - 1)b^{-k} = \sum_{k=n-1}^{\infty} b^{-k} - \sum_{k=n}^{\infty} b^{-k} = b^{-(n-1)} \quad (2.3)$$

Betrachten wir zuerst die Situation des Abrundens bei der Standardrundung, was der Situation $d_{m+1} < b/2$ entspricht. Da b nach Voraussetzung gerade ist, gilt somit $b/2 \in \mathbb{Z}$. Somit folgt aus $d_{m+1} < b/2$, dass auch $d_{m+1} + 1 \leq b/2$ gilt. Da $0 < d_1 < b$ vorausgesetzt wurde (was der normierten Gleitpunktdarstellung 2.1.13 entspricht), hat man $b^{-1} \leq \sum_{k=1}^{\infty} d_k b^{-k} \leq 1$ und somit $b^{\ell-1} \leq x \leq b^\ell$. Beachtet man $x - fl(x) = b^\ell \sum_{m+1}^{\infty} d_k b^{-k}$, so schließt man nun für $d_{m+1} < b/2$ mit (2.3)

$$\begin{aligned} |fl(x) - x| &= b^\ell \left[d_{m+1} b^{-(m+1)} + \sum_{k=m+2}^{\infty} d_k b^{-k} \right] \leq b^\ell \left[d_{m+1} b^{-(m+1)} + b^{-(m+1)} \right] \\ &\leq b^{\ell-(m+1)} (d_{m+1} + 1) \leq b^{\ell-(m+1)} \frac{b}{2} = \frac{b^{1-m}}{2} b^{\ell-1} \end{aligned} \quad (2.4)$$

Für das Aufrunden, d.h. für den Fall $d_{m+1} \geq \frac{b}{2}$ gilt nun $d_{m+1} b^{-1} \geq \frac{1}{2}$ und analog zu (2.4)

$$\begin{aligned} |fl(x) - x| &= b^\ell \left[b^{-m} - \sum_{k=m+1}^{\infty} d_k b^{-k} \right] \leq b^\ell \left[b^{-m} - d_{m+1} b^{-(m+1)} \right] \\ &= b^{\ell-m} \underbrace{\left| 1 - \underbrace{d_{m+1} b^{-1}}_{\substack{\geq 1/2 \\ \leq 1/2}} \right|}_{\leq 1/2} \leq \frac{b^{1-m}}{2} b^{\ell-1}. \end{aligned}$$

Damit ist die Abschätzung für den absoluten Fehler bewiesen und die Abschätzung für den relativen Fehler erhält man sofort mit $b^{\ell-1} \leq x$. \square

Definition 2.2.8 Die Zahl $eps := \frac{b^{1-m}}{2}$ heißt (relative) **Maschinengenauigkeit**.

Definition 2.2.9 Wir definieren

$$\mathbb{M}' := \{x \in \mathbb{R} : x_{\min} \leq |x| \leq x_{\max}\} \cup \{0\}$$

Bemerkung 2.2.10 (i) Es gilt $eps = \inf \{\delta > 0 : fl(1 + \delta) > 1\}$

(ii) In Matlab liefert die Anweisung `eps` die Maschinengenauigkeit, die sich aber auch durch folgende Routine `mineps.m` bestimmen lässt. Unter anderen Programmiersprachen geht dies ähnlich, man muss sich jedoch vergewissern, dass der Compiler die Anweisung `while 1 < 1 + value` nicht zu `while 0 < value` optimiert!

(iii) $\forall x \in \mathbb{M}', x \neq 0 \exists \varepsilon \in \mathbb{R}$ mit $|\varepsilon| < eps$ und $\frac{|fl(x) - x|}{|x|} = \varepsilon$, d.h.

$$fl(x) = x(1 + \varepsilon) \quad (2.5)$$

MATLAB-Funktion: mineps.m

```
1 function value = mineps
2 value = 1;
3 while 1 < 1 + value
4     value = value / 2;
5 end
6 value = value * 2;
```

MATLAB-Beispiel:

Die Funktion `mineps` liefert das gleiche Ergebnis wie die Matlab-Konstante `eps`. Ergänzend sei noch die Matlab-Konstante `realmin` wiedergegeben. Sie ist die kleinste positive darstellbare Maschinenzahl.

```
>> eps
ans =
    2.220446049250313e-016
>> mineps
ans =
    2.220446049250313e-016
>> realmin
ans =
    2.225073858507201e-308
```

Voraussetzung 2.2.11 Wir unterstellen im Folgenden, dass für die Addition \oplus im Rechner gilt

$$x \oplus y = fl(x + y) \quad \text{für alle } x, y \in \mathbb{M} \text{ mit } |x + y| \in \mathbb{M}' \quad (2.6)$$

und entsprechendes auch für die anderen Grundrechenarten. Man sagt \oplus ist **exakt gerundet**.

Folgerung 2.2.12 Für alle $x, y \in \mathbb{M}$ mit $x + y \in \mathbb{M}'$ gilt

$$x \oplus y = (x + y)(1 + \varepsilon) \quad (2.7)$$

für ein geeignetes $|\varepsilon| \leq eps$, und entsprechendes auch für die anderen Grundrechenarten.

Bemerkung 2.2.13 Die Maschinenaddition und -multiplikation sind nicht assoziativ, z.B. gilt für $b = 10$ und $m = 2$

$$(100 \oplus 4) \oplus 4 = 100 \neq 110 = 100 \oplus (4 \oplus 4). \quad (2.8)$$

Man beachte dabei $100 \oplus 4 = 100$, da $100 + 4$ auf 2 Stellen gerundet 100 ergibt, und dass $100 + 8$ auf 2 Stellen gerundet 110 ergibt.

Der unangenehmste Effekt der Gleitpunktarithmetik ist die sogenannte Auslöschung.

Beispiel 2.2.14 (Auslöschung) Wir betrachten die Rechnung

$$0.1236 + 1.234 - 1.356 = 0.0016 = 0.1600 \cdot 10^{-2}. \quad (2.9)$$

Gleitpunktrechnung mit $b = 10$, $m = 4$ liefert

$$(0.1236 \oplus 1.234) \ominus 1.356 = 1.358 \ominus 1.356 = 0.2000 \cdot 10^{-2}, \quad (2.10)$$

d.h. schon die erste Stelle des berechneten Ergebnisses ist falsch! Der Rundungsfehler der ersten Addition wird durch die nachfolgende Subtraktion extrem verstärkt. (Dieser Effekt tritt nicht bei Multiplikation und Division auf.)

Bemerkung 2.2.15 Auch wenn alle Summanden positiv sind, kann etwas schiefgehen: Sei wie im vorigen Beispiel $b = 10$ und $m = 4$. Wir wollen berechnen

$$\sum_{n=0}^N a_n, \quad a_0 = 1. \quad (2.11)$$

Falls $0 \leq a_n < 10^{-4}$ für alle $n > 0$, und falls wir der Reihe nach summieren, ist das berechnete Ergebnis 1, egal wie groß die Summe wirklich ist. Eine erste Abhilfe ist es, der Größe nach zu summieren (zuerst die kleinste Zahl). In unserem Beispiel (2.11) hilft das allerdings nicht wirklich, das berechnete Ergebnis ist dann höchstens 2. Besser ist jedoch folgendes Vorgehen: Man ersetzt die zwei kleinsten Zahlen durch ihre Summe und wiederholt diesen Vorgang so lange, bis nur noch eine Zahl, nämlich die gesuchte Summe, übrigbleibt.

2.3 GLEITKOMMAARITHMETIK

Beispiel 2.3.1 (Guard Digit) . Sei $\mathbb{M} = \mathbb{M}(10, 3, 1)$ und betrachte \ominus . Sei weiter $x = 0.215 \cdot 10^9$ und $y = 0.125 \cdot 10^{-9}$. Naive Realisierung von $x \ominus y = \text{rd}(x - y)$ erfordert **schieben** von y auf den größeren Exponenten $y = 0.125 \cdot 10^{-18} \cdot 10^9$ und **subtrahieren** der Mantissen:

$$\begin{array}{r} x = 0.21500000000000000000 \cdot 10^9 \\ y = 0.000000000000000000125 \cdot 10^9 \\ \hline x - y = 0.21499999999999999875 \cdot 10^9 \end{array} \quad (2.12)$$

Runden auf drei Stellen liefert dann $x \ominus y = 0.215 \cdot 10^9$. Dies erfordert einen Addierer mit $2(b^\ell - 1) + m = 21$ Stellen. In diesem Fall hätten wir das Ergebnis auch durch die Abfolge **Schieben**, **Runde** y , **Subtrahiere** erhalten. Im Allgemeinen ist das aber nicht gut wie folgendes Beispiel zeigt:

$$\begin{array}{r} x = 0.101 \cdot 10^1 \\ y = 0.993 \cdot 10^0 \end{array} \longrightarrow \begin{array}{r} x = 0.101 \cdot 10^1 \\ y = 0.099 \cdot 10^1 \\ \hline x \ominus y = 0.002 \cdot 10^1 \end{array} .$$

Für den relativen Fehler im Ergebnis gilt dann

$$\frac{(x \ominus y) - (x - y)}{(x - y)} = \frac{0.02 - 0.017}{0.017} \approx 0.176 \approx 35\text{eps},$$

wobei $\text{eps} = 1/2 \cdot 10^{-3+1} = 0.005$ sei.

Verwenden wir nun einen $m + 1$ -stelligen Addierer, dann erhalten wir

$$\begin{array}{r} x = 0.1010 \cdot 10^1 \\ y = 0.0993 \cdot 10^1 \\ \hline x - y = 0.0017 \cdot 10^1 \end{array} .$$



Das Ergebnis $x \ominus y = 1.7 \cdot 10^{-2}$ ist exakt!

Bemerkung 2.3.2 Allgemein kann man zeigen (siehe [Knuth]): Mit einer zusätzlichen Stelle (so genannter **Schutzziffer** oder **Guard Digit**) gilt

$$\frac{(x \ominus y) - (x - y)}{(x - y)} \leq 2 \text{eps} .$$

Mit nur 2 Guard Digits erhält man sogar genau das gerundete Ergebnis der exakten Rechnung, d.h. das gerundete Ergebnis aus (2.12) ohne einen $2(b^\ell - 1) + m$ -stelligen Addierer.

Völlig unkalkulierbar sind Gleitkommaoperationen aber nicht. Es gilt z.B. folgendes Resultat: Sind u, v Gleitkommazahlen und

$$\begin{array}{ll} u' = (u \oplus v) \ominus v, & v' = (u \oplus v) \ominus u, \\ u'' = (u \oplus v) \ominus v', & v'' = (u \oplus v) \ominus u', \end{array}$$

dann folgt

$$u + v = (u \oplus v) + ((u \ominus u') \oplus (v \ominus v')) .$$

Dies erlaubt eine Berechnung des Fehlers mittels Gleitkommaarithmetik. (Siehe [Knuth, 4.2.2, Theorem B].)

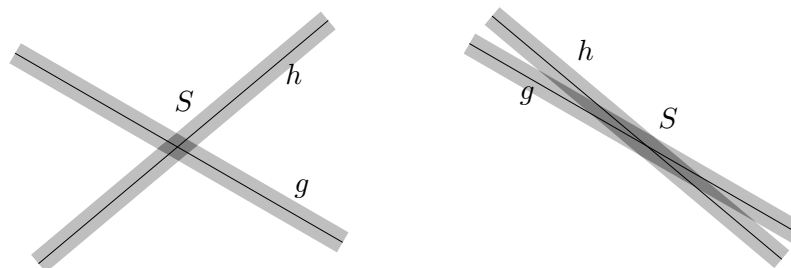


Abb. 2.2: Schnittpunkt S zweier Geraden g und h . Links gut konditioniert, rechts schlecht konditioniert.

2.4 KONDITION EINES PROBLEMS

Diskutieren wir am Anfang zuerst folgendes geometrische Problem: die zeichnerische Bestimmung des Schnittpunkts S zweier Geraden g und h in der Ebene. Schon beim Zeichnen haben wir Schwierigkeiten, die Geraden ohne Fehler darzustellen. Die Frage, die wir näher untersuchen wollen, lautet, wie stark der Schnittpunkt S (Output) von den Zeichenfehlern (Fehler im Input) abhängt.

Wie wir der Grafik direkt entnehmen können, hängt der Fehler in der Ausgabe stark davon ab, in welchem Winkel $\angle(g, h)$ sich die Geraden schneiden. Stehen g und h annähernd senkrecht aufeinander, so variiert der Schnittpunkt S etwa im gleichen Maße wie der Fehler beim Zeichnen von g und h . Man bezeichnet das Problem, den Schnittpunkt S zeichnerisch zu bestimmen, als gut konditioniert. Ist jedoch der Winkel $\angle(g, h)$ sehr klein (g und h sind fast parallel), so kann man schon mit dem Auge keinen genauen Schnittpunkt S ausmachen, und eine kleine Lageänderung von g oder h liefert einen gänzlich anderen Schnittpunkt. Man spricht hier von einem schlecht konditionierten Problem.

Kommen wir nun zu einer mathematischen Präzisierung des Konditionsbegriffs.

Problem 2.4.1 Seien X, Y Mengen und $\varphi : X \rightarrow Y$. Wir betrachten das Problem:

$$\text{Gegeben sei } x \in X, \text{ gesucht } y = \varphi(x).$$

Wir untersuchen, wie sich Störungen in den Daten x auf das Ergebnis y auswirken. Wir betrachten den Spezialfall $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\varphi(x) = \langle a, x \rangle + b, \quad a \in \mathbb{R}^n, \quad b \in \mathbb{R}, \tag{2.13}$$

wobei $\langle \cdot, \cdot \rangle$ das (euklidische) Skalarprodukt bezeichne. Seien $x, \tilde{x} \in \mathbb{R}^n$. Der relative Fehler von $\varphi(\tilde{x})$ bezüglich $\varphi(x)$ lässt sich folgendermaßen abschätzen (falls $\varphi(x) \neq 0$ und $x_j \neq 0$ für alle j):

$$\left| \frac{\varphi(x) - \varphi(\tilde{x})}{\varphi(x)} \right| = \frac{|\langle a, x - \tilde{x} \rangle|}{|\varphi(x)|} \leq \frac{\sum_{j=1}^n |a_j| |x_j - \tilde{x}_j|}{|\varphi(x)|} = \sum_{j=1}^n \frac{|x_j|}{|\varphi(x)|} |a_j| \cdot \frac{|x_j - \tilde{x}_j|}{|x_j|}. \tag{2.14}$$

Setzen wir $\tilde{x}_j = fl(x_j)$ in (2.14) ein, so ergibt sich mit Lemma 2.2.7, dass die Zahl

$$eps \cdot \sum_{j=1}^n \frac{|x_j|}{|\varphi(x)|} |a_j| \tag{2.15}$$

den unvermeidlichen Fehler (d.h. selbst bei exakter Rechnung) bei der Berechnung von φ repräsentiert. Sei nun $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}$ differenzierbar in x . Dann ist

$$\varphi(x) - \varphi(\tilde{x}) = \langle \nabla \varphi(x), x - \tilde{x} \rangle + o(\|x - \tilde{x}\|), \tag{2.16}$$

und (2.14) wird zu

$$\left| \frac{\varphi(x) - \varphi(\tilde{x})}{|\varphi(x)|} \right| \leq \sum_{j=1}^n \frac{|x_j|}{|\varphi(x)|} \left| \frac{\partial \varphi}{\partial x_j}(x) \right| \cdot \frac{|x_j - \tilde{x}_j|}{|x_j|} + o(\|x - \tilde{x}\|). \quad (2.17)$$

Ist schließlich $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$, so können wir die bisherigen Überlegungen auf jede Komponente von φ einzeln anwenden.

Definition 2.4.2 (Konditionszahlen) Sei $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ differenzierbar in $x \in \mathbb{R}^n$, sei $\varphi_i(x) \neq 0$, $1 \leq i \leq m$. Die Zahlen

$$k_{ij}(x) = \frac{|x_j|}{|\varphi_i(x)|} \left| \frac{\partial \varphi_i}{\partial x_j}(x) \right|, \quad 1 \leq i \leq m, \quad 1 \leq j \leq n, \quad (2.18)$$

heißen die Konditionszahlen von φ in x .

Beispiel 2.4.3 1. Multiplikation: $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\varphi(x_1, x_2) = x_1 \cdot x_2$,

$$k_1(x) = \frac{|x_1|}{|x_1 x_2|} \left| \frac{\partial \varphi}{\partial x_1}(x) \right| = 1, \quad k_2(x) = 1. \quad (2.19)$$

Keine Verstärkung des relativen Fehlers; die Multiplikation ist „gut konditioniert“.

2. Addition: $\varphi : \mathbb{R}^2 \rightarrow \mathbb{R}$, $\varphi(x_1, x_2) = x_1 + x_2$,

$$k_1(x) = \frac{|x_1|}{|x_1 + x_2|} \left| \frac{\partial \varphi}{\partial x_1}(x) \right| = \frac{|x_1|}{|x_1 + x_2|}, \quad k_2(x) = \frac{|x_2|}{|x_1 + x_2|}. \quad (2.20)$$

Im Falle der Auslöschung, d.h. wenn $|x_j| \gg |x_1 + x_2|$, große Verstärkung des relativen Fehlers; die Addition ist in diesem Fall „schlecht konditioniert“.

3. Lösen der quadratischen Gleichung $x^2 + 2px - q = 0$ im Fall $p, q > 0$, Berechnung der größeren der beiden Nullstellen:

$$\varphi(p, q) = -p + \sqrt{p^2 + q}, \quad \varphi : \mathbb{R}^2 \rightarrow \mathbb{R}, \quad (2.21)$$

Eine Rechnung ergibt

$$k_p = \frac{p}{\sqrt{p^2 + q}}, \quad k_q = \frac{p + \sqrt{p^2 + q}}{2\sqrt{p^2 + q}}, \quad (2.22)$$

also $k_p \leq 1$, $k_q \leq 1$; das Problem ist ebenfalls gut konditioniert.

2.5 NUMERISCHE STABILITÄT EINES ALGORITHMUS

Jeder Algorithmus zur Lösung von Problem 2.4.1 lässt sich auffassen als eine Abbildung $\tilde{\varphi} : X \rightarrow Y$. Von einem guten Algorithmus wird man erwarten, dass er die relativen Fehler nur unwesentlich mehr verstärkt als die Konditionszahlen $k_{ij}(x)$ des Problems φ es erwarten lassen. Für $\tilde{\varphi} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ sollte also eine Abschätzung der Form

$$\left| \frac{\tilde{\varphi}_i(\tilde{x}) - \tilde{\varphi}_i(x)}{\tilde{\varphi}_i(x)} \right| \leq C_1 \underbrace{\sum_{j=1}^n k_{ij}(x) \frac{|\tilde{x}_j - x_j|}{|x_j|}}_{\geq \frac{|\varphi(\tilde{x}) - \varphi(x)|}{|\varphi(x)|} - o(\|\tilde{x} - x\|)} + C_2 n \text{ eps}, \quad (2.23)$$

(oder so ähnlich) gelten mit Konstanten C_i , welche nicht viel größer als 1 sind.

Definition 2.5.1 (Numerische Stabilität eines Algorithmus) Ein Algorithmus $\tilde{\varphi} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ zur Lösung eines Problems $\varphi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ heißt numerisch stabil (oder gut konditioniert), falls (2.23) oder etwas Ähnliches gilt mit vernünftigen Konstanten C_i . Andernfalls heißt der Algorithmus numerisch instabil (oder schlecht konditioniert).

Bemerkung 2.5.2 (Konsistenz eines Algorithmus) Aus (2.23) ist nicht ersichtlich, was das exakte Datum $\varphi(x)$ mit $\tilde{\varphi}(x)$ zu tun hat. (Z.B. ist $\tilde{\varphi}(x) = 0$ ein stabiles Verfahren, aber im Allgemeinen nutzlos.) Für diesen Zusammenhang hat man in der Numerik den Begriff der **Konsistenz**. Dies ist die Eigenschaft eines numerischen Verfahrens, die bedeutet, dass der Algorithmus (in einer noch näher zu bestimmenden Art und Weise) tatsächlich das gegebene Problem löst und nicht ein anderes. Wir werden dies später noch insbesondere im Zusammenhang mit der numerischen Lösung von gewöhnlichen Differentialgleichungen untersuchen, aber man mag schon jetzt festhalten, dass Stabilität und Konsistenz für die Konvergenz eines Verfahrens wichtig sind.

Bemerkung 2.5.3 (Standardfehler) Typisch für das Entstehen numerischer Instabilität ist, dass man das Ausgangsproblem φ in zwei Teilschritte $\varphi = \varphi^{(2)} \circ \varphi^{(1)}$ zerlegt, von denen einer erheblich schlechter konditioniert ist als das Ausgangsproblem.

Beispiel 2.5.4 (Quadratische Gleichung) Wir untersuchen zwei verschiedene Verfahren zur Lösung von (2.21) in Beispiel 2.4.3

$$\varphi(p, q) = -p + \sqrt{p^2 + q}, \quad x^2 + 2px - q = 0$$

Methode 2.5.5

$$u = \sqrt{p^2 + q}, \quad y = \chi_1(p, u) = -p + u. \quad (2.24)$$

Falls $u \approx p$, d.h. falls $p \gg q$, sind die Konditionszahlen von χ_1 erheblich größer als 1 (Auslöschung).

Methode 2.5.6

$$u = \sqrt{p^2 + q}, \quad y = \chi_2(p, q, u) = \frac{q}{p + u}. \quad (2.25)$$

Die Konditionszahlen von χ_2 sind kleiner als 1. (Das Verfahren beruht darauf, dass $-p - u$ die andere Lösung und das Produkt der beiden Lösungen gleich $-q$ ist (Satz von Viëta).)

Es stellt sich heraus, dass Algorithmus 2.5.5 numerisch instabil, aber Algorithmus 2.5.6 numerisch stabil ist.

Bemerkung 2.5.7 Man beachte beim Algorithmus 2.5.5 in Beispiel 2.5.4, dass die numerische Auswertung der Funktion χ_1 für sich genommen (trotz Auslöschung) nicht numerisch instabil ist! Schlecht konditioniert ist hier das Problem, $\chi_1(p, \sqrt{p^2 + q})$ zu berechnen. Umgekehrt kann man daran sehen, dass das Zusammensetzen zweier gut konditionierter Algorithmen sehr wohl einen schlecht konditionierten Algorithmus für das Gesamtproblem ergeben kann. Diese Tatsache steht in unangenehmem, aber unvermeidlichem Kontrast zu dem Wunsch, ein Problem in unabhängig voneinander zu bearbeitende Teilprobleme zu zerlegen.

2.5.1 Vorwärtsanalyse

Die sogenannte Vorwärtsanalyse besteht darin, die Rundungsfehler im Laufe einer Rechnung durch Anwendung der Formeln (2.5) und (2.7) (und weiterer entsprechender) für jeden einzelnen Rechenschritt abzuschätzen. Man erhält dadurch eine obere Schranke für den gesamten durch die Rundungen bedingten Fehler. Kennt man außerdem die Konditionszahlen des Problems, so kann man auf diese Weise (theoretisch wenigstens) feststellen, ob der Algorithmus gut oder schlecht konditioniert ist.

Beispiel 2.5.8 Berechnung des Skalarprodukts $s = \langle x, y \rangle$ für $x, y \in \mathbb{R}^n$ mit dem Verfahren

$$s := 0; \text{ for } k := 1 \text{ to } n \text{ do } s := s + x_k \cdot y_k; \quad (2.26)$$

Notation 2.5.9 Für $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ definieren wir $|x| \in \mathbb{R}^n$ durch

$$|x| = (|x_1|, \dots, |x_n|). \quad (2.27)$$

Für die Matrix $A \in \mathbb{R}^{m \times n}$ definieren wir die Matrix $|A| \in \mathbb{R}^{m \times n}$ durch

$$|A|_{ij} = |a_{ij}|, \quad 1 \leq j \leq n. \quad (2.28)$$

Satz 2.5.10 (Vorwärtsanalyse des Rundungsfehlers für das Skalarprodukt)

Sei $0 < eps < 1/n$. Dann gilt

$$|s - \langle x, y \rangle| \leq \frac{n \, eps}{1 - n \, eps} \langle |x|, |y| \rangle. \quad (2.29)$$

Beweis. Wegen (2.5) lassen sich die Zwischenergebnisse s_k schreiben als

$$s_1 = x_1 y_1 (1 + \delta_1) \quad , \quad |\delta_1| \leq eps, \quad (2.30)$$

$$s_k = [s_{k-1} + x_k y_k (1 + \delta_k)] (1 + \varepsilon_k) \quad , \quad |\delta_k|, |\varepsilon_k| \leq eps, \quad (2.31)$$

($x_i \odot y_i = x_i y_i (1 + \delta_i)$, $x_i \oplus y_i = (x_i + y_i)(1 + \varepsilon_i)$) also gilt für $s = s_n$

$$s - \langle x, y \rangle = \sum_{k=1}^n x_k y_k [((1 + \delta_k) \prod_{j=k}^n (1 + \varepsilon_j)) - 1], \quad \varepsilon_1 = 0. \quad (2.32)$$

Die Behauptung folgt aus (2.32) mit der Bernoullischen-Ungleichung

$$\begin{aligned} |((1 + \delta_k) \prod_{j=k}^n (1 + \varepsilon_j)) - 1| &\leq |(1 + eps)^n - 1| \leq \left| \frac{1}{(1 - eps)^n} - 1 \right| \\ &\leq \left| \frac{1}{1 - n \, eps} - 1 \right| \leq \left| \frac{n \, eps}{1 - n \, eps} \right|. \end{aligned} \quad (2.33)$$

□

Aus der Formel (2.29) ergibt sich

$$\left| \frac{s - \langle x, y \rangle}{\langle x, y \rangle} \right| \leq \frac{\langle |x|, |y| \rangle}{|\langle x, y \rangle|} \left| \frac{n \, eps}{1 - n \, eps} \right|, \quad (2.34)$$

d.h. bei gegebenem $n \, eps$ wird der relative Fehler groß, falls der Winkel zwischen x und y nahe bei 90 Grad ist.

Bemerkung 2.5.11 (Intervallarithmetik) In der Intervallarithmetik wird jede Zahl $x \in \mathbb{M}^l$ durch ein (mittels zweier Gleitpunktzahlen) definiertes Intervall dargestellt, welches x einschließt. Die arithmetischen Operationen werden für solche Intervalle so definiert, dass das resultierende Intervall das exakte Ergebnis der Operation einschließt (und zwar für jede Kombination von Operationen in den Ausgangsintervallen). Falls die Hardware die Intervallarithmetik nicht unterstützt, muss sie softwareseitig simuliert werden (was die Rechnung natürlich stark verlangsamt). In der eben beschriebenen einfachen Form ist sie allerdings für die meisten Probleme der Praxis unbrauchbar, da die Intervalle unrealistisch groß werden, d.h. das mit normaler Gleitpunktarithmetik berechnete Ergebnis ist erheblich genauer.

Bemerkung 2.5.12 (Hochgenauigkeitsarithmetik) Die Grundlage der sogenannten Hochgenauigkeitsarithmetik besteht darin, das Skalarprodukt bis auf Maschinengenauigkeit auszurechnen, d.h. zu erreichen, dass der berechnete Wert s des Skalarprodukts die Eigenschaft

$$s = fl(\langle x, y \rangle) \quad (2.35)$$

hat. Hierzu wird hardwareseitig ein langer Akkumulator benötigt, d.h. ein spezielles Register, welches einen festen Dezimalpunkt und so viele Stellen hat, dass jede Gleitkommazahl in \mathbb{M} dargestellt werden kann. Hält man die Zwischensummen s aus (2.26) in diesem Akkumulator, so kann $\langle x, y \rangle$ exakt berechnet werden, falls $x_i, y_i \in \mathbb{M}$. Der Rundungsfehler entsteht dann beim Abspeichern von s als Gleitkommazahl.

In Verbindung mit Intervallarithmetik lassen sich Algorithmen konstruieren, die bereits für einen recht breiten Anwendungsbereich sehr gute Einschließungen und damit genaue und sichere numerische Ergebnisse liefern. Ein wesentliches Problem ist die benötigte Hardwareunterstützung. Die prinzipiellen Schwierigkeiten schlecht konditionierter Probleme werden natürlich nicht beseitigt, insbesondere auch nicht die Auswirkungen des „Standardfehlers“ (Zerlegung eines gut konditionierten Problems in schlecht konditionierte Teilschritte).

2.5.2 Rückwärtsanalyse

Während in der Vorwärtsanalyse die Kondition eines Algorithmus durch Vergleich der Fehlerverstärkung des Algorithmus (durch Abschätzung der akkumulierten Rundungsfehler) mit den Konditionszahlen des Problems ermittelt wird, geht man in der Rückwärtsanalyse folgendermaßen vor: Sei $\varphi : X \rightarrow Y$ ein Problem und $\tilde{\varphi} : X \rightarrow Y$ ein Algorithmus. Zu gegebenen $x \in X$ sucht man ein $\tilde{x} \in X$ mit

$$\tilde{\varphi}(x) = \varphi(\tilde{x}). \quad (2.36)$$

Gilt dann

$$\left| \frac{\tilde{x} - x}{x} \right| \leq Cn \text{ eps}, \quad (2.37)$$

mit einer kleinen Konstante C (n repräsentiert in einem geeigneten Sinn die Problemgröße), so ist man zufrieden, da sich das berechnete Ergebnis $\tilde{\varphi}(x)$ interpretieren lässt als das Ergebnis einer exakten Rechnung auf mit kleinem relativen Fehler behafteten Daten.

Definition 2.5.13 (Alternative Definition der numerischen Stabilität) Ein Algorithmus heißt gut konditioniert (oder numerisch stabil), falls er (2.36), (2.37) erfüllt mit einer kleinen Konstanten C .

Beispiel 2.5.14 (Rückwärtsanalyse des Skalarprodukts) Die Rechnung im Beweis von Satz 2.5.10 zeigt, dass für das mit dem Verfahren (2.26) berechnete Skalarprodukt s gilt

$$s = \langle x, \tilde{y} \rangle, \quad \tilde{y}_k = y_k(1 + \delta_k) \prod_{j=k}^n (1 + \varepsilon_j), \quad (2.38)$$

also mit (2.33)

$$\left| \frac{\tilde{y}_k - y_k}{y_k} \right| \leq \frac{n \text{ eps}}{1 - n \text{ eps}}, \quad \text{falls } \text{eps} \leq 1/n. \quad (2.39)$$

Bemerkung 2.5.15 i.) Also ist das Verfahren (2.26) gut konditioniert (nochmal: das Problem, das Skalarprodukt zu berechnen, kann schlecht konditioniert sein, wenn die Konditionszahlen des Skalarprodukts groß sind).

ii.) Man beachte, dass die Rückwärtsanalyse prinzipiell nichts über die Kondition des Problems aussagt.

iii.) Die Rückwärtsanalyse hat sich als recht brauchbares Werkzeug der Rundungsfehleranalyse erwiesen, und zwar vor allem bei numerischen Verfahren der Linearen Algebra.

3 FEHLERABSCHÄTZUNG MITTELS KONDITION

3.1 ABSOLUTE UND RELATIVE FEHLER

Es sei $\tilde{x} \in \mathbb{R}^n$ eine Approximation von $x \in \mathbb{R}^n$. Als **Fehler** bezeichnet man $e = x - \tilde{x}$ und zu einer gegebenen Vektornorm $\|\cdot\|$ ist

$$\|x - \tilde{x}\|$$

der **absolute Fehler** von \tilde{x} . Es sei $x \neq 0$, dann beschreibt

$$\frac{\|x - \tilde{x}\|}{\|x\|}$$

den **relativen Fehler** von \tilde{x} . Relative Fehler in der ∞ -Norm können durch eine Aussage über die Anzahl korrekter Stellen von \tilde{x} ausgedrückt werden, d.h.

$$\frac{\|x - \tilde{x}\|_\infty}{\|x\|_\infty} \approx 10^{-p},$$

falls die betragsgrößte Komponente von \tilde{x} näherungsweise p korrekte signifikante Stellen hat.

3.2 FEHLERABSCHÄTZUNG & KONDITION

Wir wollen nun zwei wichtige Fragestellungen untersuchen, welche bei dem numerischen Lösen von linearen Gleichungssystemen (LGS) auftreten. Zum einen betrachten wir eine Näherungslösung \tilde{x} zu der Lösung x von $Ax = b$. Welche Rückschlüsse kann man von der Größe des **Residuums** $r := b - A\tilde{x}$ auf den Fehler $e = x - \tilde{x}$ ziehen? Des Weiteren rechnet man im Allgemeinen mit endlicher Genauigkeit auf einem Rechner. Welche Einflüsse haben dabei Störungen der Ausgangsdaten auf die Lösung x ?

Sei $\|A\|$ eine beliebige Matrixnorm und $\|x\|$ eine dazu verträgliche Vektornorm (d.h. $\|Ax\| \leq \|A\| \|x\|$). Nach Definition des Fehlers e und Residuums r gilt

$$Ae = Ax - A\tilde{x} = b - A\tilde{x} = r.$$

Daraus folgt

$$\|e\| = \|A^{-1}r\| \leq \|A^{-1}\| \|r\|.$$

Mit

$$\|b\| = \|Ax\| \leq \|A\| \|x\|$$

folgt für den relativen Fehler die Abschätzung

$$\frac{\|e\|}{\|x\|} \leq \frac{\|A^{-1}\| \|r\|}{\|b\| \cdot \frac{1}{\|A\|}} = \|A^{-1}\| \|A\| \cdot \frac{\|r\|}{\|b\|}. \quad (3.1)$$

Dies motiviert den Begriff der Konditionszahl.

Definition 3.2.1 (Konditionszahl) Man bezeichnet

$$\kappa(A) = \|A\| \|A^{-1}\| \quad (3.2)$$

als die **Konditionszahl** der Matrix A bezüglich der verwendeten Matrixnorm.

Beispiel 3.2.2 Man betrachte das LGS $Ax = b$ mit

$$A = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} \\ \frac{1}{5} & \frac{1}{6} & \frac{1}{7} & \frac{1}{8} & \frac{1}{9} \end{pmatrix} \quad b := (b_j) \quad b_j := \sum_{k=1}^5 \frac{1}{k+j-1} \quad (1 \leq j \leq 5).$$

Für die Matrixnormen und die Konditionszahlen erhält man in den Fällen $p = 1, 2, \infty$

$$\begin{array}{lll} \|A\|_1 & = & 2.28\bar{3}, & \|A^{-1}\|_1 & = & 413280, & \kappa_1(A) & = & 943656, \\ \|A\|_2 & = & 1.567051, & \|A^{-1}\|_2 & = & 304142.84, & \kappa_2(A) & = & 476607.34, \\ \|A\|_\infty & = & 2.28\bar{3}, & \|A^{-1}\|_\infty & = & 413280, & \kappa_\infty(A) & = & \kappa_1(A). \end{array}$$

Sei die exakte Lösung $x = (1, \dots, 1)^T$. Ist nun $\tilde{x} = (1 - \varepsilon)x$, dann gilt

$$\frac{\|e\|}{\|x\|} = \frac{\|\varepsilon x\|}{\|x\|} = |\varepsilon|$$

und

$$r = b - A\tilde{x} = Ax - A\tilde{x} = A(x - (1 - \varepsilon)x) = \varepsilon Ax = \varepsilon b,$$

d.h.

$$\frac{\|r\|}{\|b\|} = |\varepsilon|.$$

Obwohl die Konditionszahl sehr groß ist, verhält sich bei dieser Störung $\frac{\|e\|}{\|x\|}$ wie $\frac{\|r\|}{\|b\|}$. Obige Abschätzung ist offensichtlich eine „worst case“ Abschätzung. Betrachten wir nun



$$\tilde{x} = (0.993826, 1.116692, 0.493836, 1.767191, 0.623754)^T. \quad (3.3)$$

(A habe die Eigenwerte $0 \leq \lambda_1 < \dots < \lambda_5$ mit den zugehörigen Eigenvektoren $\varphi_1, \dots, \varphi_5$ ($\|\varphi_j\|_2 = 1$). Man beachte, dass \tilde{x} in (3.3) so gewählt ist, dass $x - \tilde{x} \approx \varphi_1$ gilt und x näherungsweise φ_5 entspricht.) Dann erhalten wir

$$\frac{\|e\|_1 \|b\|_1}{\|x\|_1 \|r\|_1} \approx 0.41 \kappa_1(A), \quad \frac{\|e\|_2 \|b\|_2}{\|x\|_2 \|r\|_2} \approx 0.89 \kappa_2(A), \quad \frac{\|e\|_\infty \|b\|_\infty}{\|x\|_\infty \|r\|_\infty} \approx 0.74 \kappa_\infty(A)$$

und schätzen dann in der richtigen Größenordnung ab!

Aufgabe 3.2.3 Es sei $A \in \mathbb{R}^{n \times n}$ eine positiv definite Matrix. Die Eigenwerte von Matrix A seien $0 < \lambda_1 \leq \dots \leq \lambda_n$ mit den zugehörigen Eigenvektoren $\varphi_1, \dots, \varphi_n$ ($\|\varphi_j\|_2 = 1$).

i) Es sei $x \in \mathbb{R}^n \setminus \{0\}$, $b := Ax$ und $\tilde{x} = cx$ ($0 \neq c \in \mathbb{R}$). Man zeige

$$\frac{\|x - \tilde{x}\|}{\|x\|} \frac{\|b\|}{\|b - A\tilde{x}\|} = 1.$$

ii) Man zeige

$$\kappa_2(A) = \lambda_n / \lambda_1.$$

iii) Es sei $x = \varphi_n$, $b := Ax$ und $\tilde{x} = x + \varphi_1$. Man zeige

$$\frac{\|x - \tilde{x}\|}{\|x\|} \frac{\|b\|}{\|b - A\tilde{x}\|} = \kappa_2(A).$$

Untersuchen wir nun, welchen Einfluss kleine Störungen in den Ausgangsdaten A, b auf die Lösung x des linearen Gleichungssystems haben können, d.h. wir sind interessiert an der **Empfindlichkeit** der Lösung x auf Störungen in den Koeffizienten. Die genaue Frage lautet: Wie groß kann die Änderung δx der Lösung x von $Ax = b$ sein, falls die Matrix um δA und b durch δb gestört sind. Dabei seien δA und δb kleine Störungen, so daß $A + \delta A$ immer noch regulär ist. Es sei $x + \delta x$ die Lösung zu

$$(A + \delta A)(x + \delta x) = (b + \delta b).$$

Teilweises Ausmultiplizieren liefert

$$Ax + \delta Ax + (A + \delta A)\delta x = b + \delta b$$

und somit mit $Ax = b$

$$\begin{aligned} \delta x &= (A + \delta A)^{-1}(\delta b - \delta Ax) \\ &= (I + A^{-1}\delta A)^{-1}A^{-1}(\delta b - \delta Ax). \end{aligned}$$

Für eine submultiplikative Matrixnorm und eine dazu verträgliche Vektornorm ergibt sich somit

$$\|\delta x\| \leq \|(I + A^{-1}\delta A)^{-1}\| \|A^{-1}\| (\|\delta b\| + \|\delta A\| \|x\|)$$

und wir erhalten die Abschätzung

$$\frac{\|\delta x\|}{\|x\|} \leq \|(I + A^{-1}\delta A)^{-1}\| \|A^{-1}\| \left(\frac{\|\delta b\|}{\|x\|} + \|\delta A\| \right). \quad (3.4)$$

Wir zeigen in Lemma 3.2.6, dass für $\|B\| < 1$ die Abschätzung

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|} \quad \text{existiert.}$$

Setzen wir nun $B = A^{-1}\delta A$ und $\|A^{-1}\| \|\delta A\| < 1$ (d.h. δA ist eine kleine Störung von A) voraus, erweitern auf der rechten Seite mit $\|A\|$, so erhalten wir unter Ausnutzung der Verträglichkeit der Normen und $Ax = b$

$$\begin{aligned} \frac{\|\delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1}\delta A\|} \left(\frac{\|\delta b\|}{\|A\| \|x\|} + \frac{\|\delta A\|}{\|A\|} \right) \\ &\leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1}\| \|\delta A\|} \left(\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right). \end{aligned}$$

Wir halten das Ergebnis nun in einem Satz fest.

Satz 3.2.4 *Es sei $A \in \mathbb{R}^{n \times n}$ regulär und $x \in \mathbb{R}^n$ die exakte Lösung von $Ax = b$. Die rechte Seite b sei um δb gestört und für die Störung von A gelte $\kappa(A) \|\delta A\| / \|A\| = \|A^{-1}\| \|\delta A\| < 1$ und $A + \delta A$ ist immer noch regulär. Dann gilt für die Lösung \tilde{x} des gestörten Systems mit Koeffizientenmatrix $A + \delta A$ und rechter Seite $b + \delta b$*

$$\frac{\|x - \tilde{x}\|}{\|x\|} \leq \frac{\kappa(A)}{1 - \kappa(A) \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta b\|}{\|b\|} + \frac{\|\delta A\|}{\|A\|} \right). \quad (3.5)$$



Bemerkung 3.2.5 Die Konditionszahl $\kappa(A)$ der Koeffizientenmatrix A ist folglich die entscheidende Größe, welche die Empfindlichkeit der Lösung bzgl. der Störungen δA und δb beschreibt.

Es bleibt folgendes Lemma zu beweisen:

Lemma 3.2.6 (Neumann-Reihe) Sei $A \in \mathbb{R}^{n \times n}$, $\|\cdot\|$ eine submultiplikative Matrixnorm und $\|A\| < 1$. Dann ist $(I - A)$ regulär und

$$(I - A)^{-1} = \sum_{k=0}^{\infty} A^k \quad (\text{Neumann-Reihe}) \quad (3.6)$$

mit

$$\|(I - A)^{-1}\| \leq \frac{1}{1 - \|A\|}. \quad (3.7)$$

Beweis. Mit $\|A\| < 1$, der Submultiplikativität und der Summenformel für die geometrische Reihe erhält man

$$\sum_{k=0}^m \|A^k\| \leq \sum_{k=0}^m \|A\|^k \leq \sum_{k=0}^{\infty} \|A\|^k = \frac{1}{1 - \|A\|} < \infty \quad (m \in \mathbb{N}). \quad (3.8)$$

Der Raum $\mathbb{R}^{n \times n}$ ist isomorph zu \mathbb{R}^{n^2} (siehe Anhang A.2). In [Analysis II, Beispiel 8.4.6] wurde gezeigt, dass \mathbb{R}^{n^2} vollständig ist bezüglich irgendeiner Norm auf \mathbb{R}^{n^2} . Somit ist $\mathbb{R}^{n \times n}$ vollständig bezüglich $\|\cdot\|$ und aus der absoluten Konvergenz folgt die Konvergenz von $\sum_{k=0}^{\infty} A^k$. Ebenso folgt aus $\|A^k\| \leq \|A\|^k \rightarrow 0$ für $k \rightarrow \infty$ die Konvergenz von $\lim_{k \rightarrow \infty} A^k = 0$. Weiter gilt die Gleichung („Teleskopsumme“)

$$\left(\sum_{k=0}^m A^k \right) (I - A) = I - A^{m+1} \quad (m \in \mathbb{N}). \quad (3.9)$$

Der Grenzübergang von (3.9) führt zur Gleichung

$$\left(\sum_{k=0}^{\infty} A^k \right) (I - A) = I. \quad (3.10)$$

Das bedeutet, $I - A$ ist regulär und $(I - A)^{-1} = \sum_{k=0}^{\infty} A^k$. Zusammen mit 3.8 wäre damit die Behauptung bewiesen. Mit der Summenformel für die geometrische Reihe erhält man nun

$$\|(I - A)^{-1}\| \leq \lim_{N \rightarrow \infty} \sum_{k=0}^N \|A^k\| \leq \lim_{N \rightarrow \infty} \sum_{k=0}^N \|A\|^k = \frac{1}{1 - \|A\|}, \quad (3.11)$$

womit der Satz bewiesen ist. □

Zum Ende des Kapitels wollen wir nun den praktischen Nutzen der Abschätzung (3.5) in einer Faustregel festhalten.

Bei einer d -stelligen dezimalen Gleitkommarechnung können die relativen Fehler der Ausgangsgrößen für beliebige, kompatible Normen von der Größenordnung

$$\frac{\|\delta A\|}{\|A\|} \approx 5 \cdot 10^{-d} \quad \frac{\|\delta b\|}{\|b\|} \approx 5 \cdot 10^{-d}$$

sein. Ist die Konditionszahl $\kappa(A) \approx 10^\alpha$ mit $5 \cdot 10^{\alpha-d} \ll 1$, so ergibt die Abschätzung

$$\frac{\|\delta x\|}{\|x\|} \leq 10^{\alpha-d+1}$$

Mit der Aussage zu (3.1) besagt die Schätzung, dass $\|\delta x\|$ maximal in der Größenordnung der $(d - \alpha - 1)$ -ten Dezimalstelle von $\|x\|$ liegen kann und dies motiviert folgende Daumenregel:

Bemerkung 3.2.7 (Daumenregel zur Genauigkeit) Wird $Ax = b$ mit d -stelliger dezimaler Gleitkommarechnung gelöst, und beträgt die Konditionszahl $\kappa(A) \approx 10^\alpha$, so sind, bezogen auf die betragsgrößte Komponente, nur $(d - \alpha - 1)$ Dezimalstellen sicher.



MATLAB-Beispiel:

Es ist bekannt, dass die Gauß-Elimination selbst mit Spaltenpivotstrategie zu überraschend ungenauen Ergebnissen beim Lösen von linearen Gleichungssystemen führen kann, obwohl die Matrix gut konditioniert ist.

Betrachten wir hierzu die von Wilkinson angegebene pathologische Matrix

$$A = \begin{pmatrix} 1 & & & 1 \\ -1 & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ -1 & \cdots & -1 & 1 \end{pmatrix}.$$

```
>> A=toeplitz([1,-ones(1,59)], ...
              [1,zeros(1,59)]);
>> A(:,60)=1;
>> cond(A)
ans =
    26.8035 % rel. gut konditioniert
>> randn('state', 3383)
>> x=randn(60,1);
>> b=A*x;
>> x1=A\b;
>> norm(x-x1)/norm(x)
ans =
    0.3402 % großer rel. Fehler
```

Bemerkung 3.2.8 Das Beispiel lässt vermuten, dass das Gauß-Verfahren über die ganze Menge der invertierbaren Matrizen betrachtet nicht stabil ist. Für die in der Praxis auftretenden Matrizen, ist das Gauß-Verfahren mit Spaltenpivotierung jedoch „in der Regel“ stabil. Für eine weitere Stabilitätsanalyse des Gauß-Verfahrens sei auf [Deuffhard/Hohmann], die grundlegenden Artikel von Wilkinson [Wilkinson65, Wilkinson69] sowie auf die Übersichtsartikel [Higham, Discroll/Maki] verwiesen.

4 ITERATIVE LÖSUNG LINEARER GLEICHUNGSSYSTEME

Gegeben sei eine reguläre Matrix $A \in \mathbb{R}^{n \times n}$ und ein lineares Gleichungssystem

$$Ax = b$$

mit der exakten Lösung x . Wir betrachten Iterationsverfahren der Form

$$x^{(i+1)} = \phi(x^{(i)}), \quad i = 0, 1, \dots$$

Mit Hilfe einer beliebigen regulären Matrix $B \in \mathbb{R}^{n \times n}$ erhält man solche Iterationsvorschriften aus der Gleichung

$$Bx + (A - B)x = b,$$

indem man

$$Bx^{(i+1)} + (A - B)x^{(i)} = b$$

setzt und nach $x^{(i+1)}$ auflöst

$$\begin{aligned} x^{(i+1)} &= x^{(i)} - B^{-1}(Ax^{(i)} - b) \\ &= (I - B^{-1}A)x^{(i)} + B^{-1}b. \end{aligned} \quad (4.1)$$

Jede Wahl einer nichtsingulären Matrix B führt zu einem möglichen **Iterationsverfahren**. Es wird umso brauchbarer, je besser B die folgenden Kriterien erfüllt

- i) B ist leicht zu invertieren (**einfache Realisierbarkeit**)
- ii) die Eigenwerte von $(I - B^{-1}A)$ sollen möglichst kleine Beträge haben. (**Konvergenzeigenschaft**)

Wir wollen hier nun einige Beispiele angeben. Dazu verwenden wir folgende (additive) Standardzerlegung

$$A = L + D + R,$$

wobei D eine **Diagonalmatrix**, L strikte **untere Dreiecksmatrix** und R strikte **obere Dreiecksmatrix** seien. Die Wahl

- i) $B = \gamma I$ liefert das **Richardson-Verfahren**
- ii) $B = D$ liefert das **Jacobi-Verfahren (Gesamtschrittverfahren)**
- iii) $B = L + D$ oder $B = D + R$ liefert das **Gauß-Seidel-Verfahren (Einzelschrittverfahren)**

Was zeichnet nun die einzelnen Verfahren aus? Betrachten wir dazu ein Beispiel.

Beispiel 4.0.1 Zu gegebenem $n \in \mathbb{N}$ und

$$A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{n \times n}, b = \vec{1} \in \mathbb{R}^n, x^{(0)} = \vec{0} \in \mathbb{R}^n$$

bestimmen wir die **Konvergenzrate**

$$c = \max_k \frac{\|x - x^{(k+1)}\|_2}{\|x - x^{(k)}\|_2},$$

für das Jacobi- und das Gauß-Seidel-Verfahren, d.h. den Faktor, um den der **Fehler** in der 2-Norm in jedem Iterationsschritt **mindestens reduziert** wird.

MATLAB-Funktion: runKonvergenz.m

```

1 n = 10;
2 e = ones(n,1);
3 A = spdiags([e -2*e e], -1:1, n, n);
4 x_ex = rand(n,1); % exakte Loesung
5 b = A * x_ex; % exakte rechte Seite
6 x{1} = rand(n,1); % zufaelliger Startv.
7 x{2}=x{1};
8 W{1} = triu(A); % Gauss-Seidel
9 W{2} = diag(diag(A)); % Jacobi
10 for j = 1:length(x)
11 error_old = norm(x{j}-x_ex);
12 for k = 1 : 20
13 x{j} = x{j} + W{j} \ (b-A*x{j});
14 error_new = norm(x{j}-x_ex);
15 quot{j}(k) = error_new/error_old;
16 error_old = error_new;
17 end
18 end
19 plot(1:20, quot{1}, 'm-s', 1:20, quot{2}, 'k:*');
20 xlabel('Anzahl der Iterationen'), ylabel('Kontraktion')
21 legend({'Gauss-Seidel-Verf.', 'Jacobi-Verf.'}, 4)

```

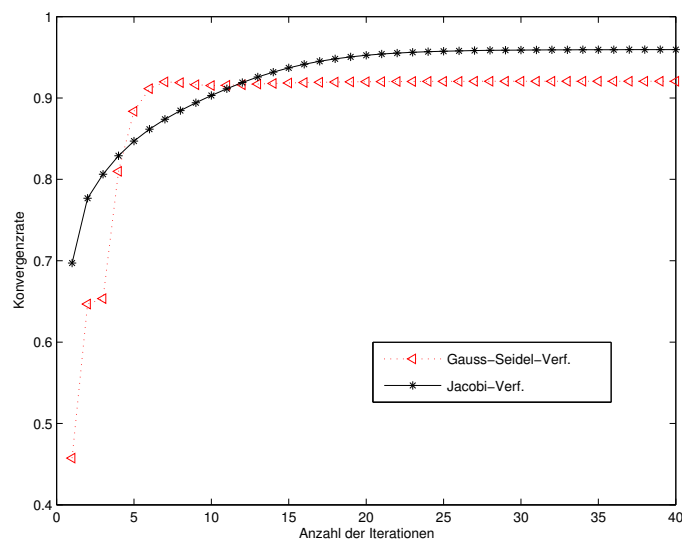


Abb. 4.1: Kontraktionszahlen für Gesamt- und Einzelschrittverfahren.

Dem numerischen Experiment kann man entnehmen, dass beide Verfahren konvergieren, da die Konvergenzrate jeweils kleiner 1 ist, und dass das Gauß-Seidel-Verfahren schneller konvergiert, da die Konvergenzrate hier kleiner ist.

Bemerkung 4.0.2 *Iterationsverfahren finden besonders Anwendung bei schwachbesetzten Matrizen, d.h. die Matrix hat „viele Nulleinträge“. Ist die Hülle jedoch groß und in jeder Zeile treten nur wenige Nichtnulleinträge auf, dann ist ein Bandlöser somit „teuer“ und die folgenden Verfahren liefern gute Alternativen.*

Bei dem Gauß-Seidel-Verfahren verwendet man in der praktischen Anwendung folgende Formulierung um die Anzahl der Operatoren zu reduzieren

$$x^{(i+1)} = (L + D)^{-1}(b - Rx^{(i)})$$

da

$$\begin{aligned} x^{(i+1)} &= x^{(i)} - B^{-1}(Ax^{(i)} - b) \\ &= x^{(i)} - B^{-1}([(A - B) + B]x^{(i)} - b) \\ &= x^{(i)} - B^{-1}(A - B)x^{(i)} - B^{-1}Bx^{(i)} + B^{-1}b \\ &= B^{-1}(b - (A - B)x^{(i)}) \end{aligned}$$

gilt und mit $B = L + D$ und $A - B = R$ folgt

$$x^{(i+1)} = (L + D)^{-1}(b - Rx^{(i)}).$$

Ein Schritt eines Gauß-Seidel-Verfahrens ist also etwa so aufwendig wie eine Matrix-Vektor-Multiplikation. Ähnlich kann man auch bei dem Jacobi-Verfahren vorgehen, d.h.

$$x^{(i+1)} = D^{-1}(b - (L + R)x^{(i)})$$

und das Richardson-Verfahren vereinfacht man zu

$$x^{(i+1)} = x^{(i)} + \frac{1}{\gamma}(b - Ax^{(i)}).$$

4.1 FOLGEN VON ITERATIONSMATRIZEN

Es sei x Lösung von $Ax = b$. Mit (4.1) erhalten wir

$$\begin{aligned} x - x^{(k)} &= x - B^{-1}b - (I - B^{-1}A)x^{(k-1)} \\ &= Ix - B^{-1}Ax - (I - B^{-1}A)x^{(k-1)} \\ &= (I - B^{-1}A)(x - x^{(k-1)}) = \dots = (I - B^{-1}A)^k(x - x^{(0)}). \end{aligned}$$

Die Konvergenz des dargestellten Verfahrens hängt also nur von den Eigenschaften der Iterationsmatrix $I - B^{-1}A$ ab. Es sei C eine beliebige komplexwertige $(n \times n)$ -Matrix, $\lambda_i := \lambda_i(C)$ ($i = 1, \dots, n$) seien die Eigenwerte von C . Dann bezeichnen wir mit

$$\rho(C) := \max_{1 \leq i \leq n} \{|\lambda_i(C)|\}$$

den **Spektralradius** von C . Bevor wir ein Konvergenzkriterium angeben, bereiten wir noch den Begriff der Jordanschen Normalform vor.

Definition 4.1.1 (Jordan-, bzw. Elementarmatrix) Eine Matrix $E_k(\lambda) \in \mathbb{C}^{k \times k}$ heißt Jordanmatrix (oder Elementarmatrix) zum Eigenwert λ , wenn

$$E_k(\lambda) = \begin{pmatrix} \lambda & 1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda \end{pmatrix}. \quad (4.2)$$

Satz 4.1.2 (Jordansche Normalform (siehe z.B. [Fischer])) Zu jeder Matrix $A \in \mathbb{C}^{n \times n}$ existiert eine reguläre Matrix $T \in \mathbb{C}^{n \times n}$, so dass

$$A = T^{-1}JT,$$

wobei J , die durch die Paare $(\lambda_1, n_1), \dots, (\lambda_k, n_k)$ mit $\lambda_i \in \mathbb{C}$, $n_i \geq 1$ (eindeutig bis auf die Reihenfolge) bestimmte Jordansche Normalform

$$J = \begin{pmatrix} E_{n_1}(\lambda_1) & & 0 \\ & \ddots & \\ 0 & & E_{n_k}(\lambda_k) \end{pmatrix}$$

von A ist.

Satz 4.1.3 (Konvergenzkriterium) Es sei $C \in \mathbb{C}^{n \times n}$. Die Folge $(C^k)_{k \in \mathbb{N}}$ ist genau dann eine Nullfolge, wenn $\varrho(C) < 1$ gilt.

Beweis. Sei zunächst $\varrho(C) \geq 1$. Dann gibt es einen Eigenwert λ mit $|\lambda| \geq 1$ und einen Vektor $x \neq 0$ mit $Cx = \lambda x$. Wegen $C^k x = \lambda^k x$ und $\lim_{k \rightarrow \infty} \lambda^k \neq 0$ kann folglich $(C^k)_k$ keine Nullfolge sein. Die Bedingung $\varrho(C) < 1$ ist somit notwendig.

Sei nun $\varrho(C) < 1$. Weil $(TCT^{-1})^k = TC^kT^{-1}$ für jede Ähnlichkeitstransformation T gilt, reicht es, $\lim_{k \rightarrow \infty} (TCT^{-1})^k = 0$ zu zeigen. Die Matrix C lässt sich durch Ähnlichkeitstransformation auf die Jordansche Normalform J transformieren. Wir zeigen, daß $\lim_{k \rightarrow \infty} J^k = 0$ gilt, wenn alle Eigenwerte $\lambda_1, \dots, \lambda_n$ dem Betrag nach kleiner Eins sind. Dazu sei

$$E_\mu = E_{n_\mu}(\lambda_\mu) = \begin{pmatrix} \lambda_\mu & 1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda_\mu \end{pmatrix} \in \mathbb{C}^{n_\mu \times n_\mu}$$

eine Elementarmatrix zum Eigenwert λ_μ der Jordanschen Normalform J von C . Da offenbar

$$J^k = \begin{pmatrix} E_1^k & & & \\ & E_2^k & & \\ & & \ddots & \\ & & & E_\ell^k \end{pmatrix}$$

mit $1 \leq \ell \leq n$ gilt, genügt es, das Konvergenzverhalten einer Jordanmatrix E_μ zu untersuchen. Wir schreiben E_μ in der Form $E_\mu = \lambda_\mu I + S$ mit

$$S = \begin{pmatrix} 0 & 1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ 0 & & & 0 \end{pmatrix}$$

und bilden $E_\mu^k = (\lambda_\mu I + S)^k$. Nach Anwendung der binomischen Formel und unter Beachtung von $S^n = 0$ erhält man die Beziehung

$$E_\mu^k = \sum_{\nu=0}^{\min\{k, n-1\}} \binom{k}{\nu} \lambda_\mu^{k-\nu} S^\nu.$$

Für feste ν hat man mit

$$\binom{k}{\nu} = \frac{k!}{\nu!(k-\nu)!} = \frac{k(k-1)\cdots(k-\nu+1)}{1\cdots\nu} \leq k^\nu$$

die Abschätzung

$$\left| \binom{k}{\nu} \lambda_\mu^{k-\nu} \right| \leq |\lambda_\mu^{k-\nu} k^\nu|.$$

Da $|\lambda_\mu| < 1$, strebt $k \log |\lambda_\mu| + \nu \log(k) \rightarrow -\infty$ für $k \rightarrow \infty$ und mit

$$|\lambda_\mu^{k-\nu} k^\nu| \leq \exp((k-\nu) \log |\lambda_\mu| + \nu \log k)$$

folgt die Konvergenz $\lim_{k \rightarrow \infty} \left| \binom{k}{\nu} \lambda_\mu^{k-\nu} \right| = 0$. Damit ist gezeigt, dass $(E_\mu^k)_k$ eine Nullfolge ist und somit auch die Behauptung. \square

Um die Konvergenz der Richardson-Iteration mathematisch zu beweisen, muss der Spektralradius der Iterationsmatrix bestimmt werden. Dies ist im Allgemeinen nicht analytisch möglich und numerisch sehr aufwendig. Es ist daher das Ziel der nächsten beiden Abschnitte, aus einfachen algebraischen Eigenschaften der Matrix A auf die Konvergenz der Jacobi- bzw. Gauß-Seidel-Iteration zu schließen. Anders als in Satz 4.1.3 sind die folgenden Ergebnisse hinreichende Konvergenzkriterien, im Allgemeinen aber nicht notwendig.

4.2 KONVERGENZ DES JACOBI-VERFAHRENS

Bei dem **Jacobi-Verfahren** (auch Gesamtschrittverfahren genannt) werden alle Komponenten des Lösungsvektors in einem Iterationsschritt gleichzeitig korrigiert. Die Iterationsvorschrift lautet

$$x^{(i+1)} = D^{-1}(b - (L + R)x^{(i)}),$$

d.h. die Iterationsmatrix ist $C_J = I - D^{-1}A = -D^{-1}(L + R)$.

Satz 4.2.1 (Starkes Zeilen- und Spaltensummenkriterium) *Es sei $A \in \mathbb{C}^{n \times n}$. Das Jacobi-Verfahren konvergiert für jeden Startvektor $x^{(0)} \in \mathbb{C}^n$, wenn für die Matrix A das*

i) (starke Zeilensummenkriterium)

$$|a_{ii}| > \sum_{\substack{k=1 \\ k \neq i}}^n |a_{ik}|, \quad \text{für } i = 1, 2, \dots, n,$$

d.h. A ist strikt diagonaldominant, oder das

ii) (starke Spaltensummenkriterium)

$$|a_{kk}| > \sum_{\substack{i=1 \\ i \neq k}}^n |a_{ik}|, \quad \text{für } k = 1, 2, \dots, n,$$

d.h. A^T ist strikt diagonaldominant, erfüllt ist.

Für den Beweis von Satz 4.2.1 benötigen wir das folgende Lemma.

Lemma 4.2.2 Es sei $A \in \mathbb{C}^{n \times n}$. Dann gilt für jede natürliche p -Matrixnorm $\varrho(A) \leq \|A\|_p$.

Beweis. Jeder Eigenwert λ von A mit zugehörigem Eigenvektor x genügt für jede natürliche p -Matrixnorm $\|\cdot\|_p$ der Beziehung

$$\frac{\|Ax\|_p}{\|x\|_p} = |\lambda|$$

und damit der Abschätzung $\|A\|_p \geq |\lambda|$. □

Beweis von Satz 4.2.1. i): Die Iterationsmatrix des Jacobi-Verfahrens bezeichnen wir mit $C_J = -D^{-1}(L + R)$. Wenn das starke Zeilensummenkriterium erfüllt ist, gilt die Abschätzung

$$\|C_J\|_\infty = \left\| \begin{pmatrix} 0 & -\frac{a_{1,2}}{a_{1,1}} & \dots & \dots & -\frac{a_{1,n}}{a_{1,1}} \\ -\frac{a_{2,1}}{a_{2,2}} & 0 & -\frac{a_{2,3}}{a_{2,2}} & \dots & -\frac{a_{2,n}}{a_{2,2}} \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \dots & \dots & \ddots & \vdots \\ -\frac{a_{n,1}}{a_{n,n}} & \dots & \dots & -\frac{a_{n,n-1}}{a_{n,n}} & 0 \end{pmatrix} \right\|_\infty = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} < 1.$$

Lemma 4.2.2 liefert dann die Behauptung i).

ii) Ist für A das starke Spaltensummenkriterium (ii) erfüllt, so gilt (i) für A^T . Also konvergiert das Jacobi-Verfahren für A^T und es ist daher wegen Lemma 4.2.2, $\varrho(C) < 1$ für $C = I - D^{-1}A^T$. Nun hat C die gleichen Eigenwerte wie C^T und $D^{-1}C^T D = I - D^{-1}A$. Also ist auch $\varrho(I - D^{-1}A) < 1$, d.h. das Jacobi-Verfahren ist auch für die Matrix A konvergent. □

Definition 4.2.3 Eine Matrix $A \in \mathbb{R}^{n \times n}$, heißt zerlegbar, wenn es nichtleere Teilmengen N_1 und N_2 der Indexmenge $N := \{1, 2, \dots, n\}$ gibt mit den Eigenschaften

- i) $N_1 \cap N_2 = \emptyset$
- ii) $N_1 \cup N_2 = N$
- iii) $a_{ij} = 0$ für alle $i \in N_1$ und $j \in N_2$

Eine Matrix, die nicht zerlegbar ist, heißt unzerlegbar (irreduzibel).

Beispiel 4.2.4 i)

$$A = \begin{pmatrix} a_{11} & \dots & a_{1k} & 0 & \dots & 0 \\ \vdots & & \vdots & \vdots & & \vdots \\ a_{N1} & \dots & a_{NN} & 0 & \dots & 0 \\ a_{N+1,1} & \dots & a_{N+1,N} & a_{N+1,N+1} & \dots & a_{N+1,2N} \\ \vdots & & \vdots & \vdots & & \vdots \\ a_{2N,1} & \dots & a_{2N,N} & a_{2N,N+1} & \dots & a_{2N,2N} \end{pmatrix}$$

Die Teilmengen $N_1 = \{1, 2, \dots, N\}$, $N_2 = \{N + 1, \dots, 2N\}$ haben alle in der Definition geforderten Eigenschaften. Somit ist A zerlegbar.

- ii) Eine Tridiagonalmatrix mit nicht verschwindenden Nebendiagonal- und Diagonalelementen ist unzerlegbar.

$$A = \begin{pmatrix} 2 & -1 & 0 \\ 1 & 4 & 0 \\ 0 & -1 & 3 \end{pmatrix}, \quad G(A) : \begin{array}{ccc} & \curvearrowright & \\ \curvearrowleft & K_1 & \curvearrowright \\ & & K_2 & \curvearrowright \\ & & & K_3 \end{array}$$

Abb. 4.2: Beispiel einer zerlegbaren Matrix A und ihres Graphen $G(A)$.

Bemerkung 4.2.5 Dass eine Matrix A unzerlegbar (irreduzibel) ist, kann man häufig leicht mit Hilfe des der Matrix A zugeordneten (gerichteten) Graphen $G(A)$ zeigen. Wenn A eine $n \times n$ -Matrix ist, so besteht $G(A)$ aus n Knoten K_1, \dots, K_n und es gibt eine gerichtete Kante $K_i \rightarrow K_j$ in $G(A)$ genau dann, wenn $a_{ij} \neq 0$. Man zeigt leicht, dass A genau dann unzerlegbar ist, falls der Graph $G(A)$ in dem Sinne zusammenhängend ist, dass es für jedes Knotenpaar (K_i, K_j) in $G(A)$ einen gerichteten Weg K_i nach K_j gibt.

Definition 4.2.6 (Schwachtes Zeilensummenkriterium) Eine Matrix $A \in \mathbb{R}^{n \times n}$ erfüllt das schwache Zeilensummenkriterium, falls

$$\sum_{\substack{\nu=1 \\ \nu \neq \mu}}^n |a_{\mu\nu}| \leq |a_{\mu\mu}|$$

für alle Zeilen $j = 1, \dots, n$ gilt, d.h. A ist diagonaldominant, und

$$\sum_{\substack{\nu=1 \\ \nu \neq \mu_0}}^n |a_{\mu_0\nu}| < |a_{\mu_0\mu_0}|$$

für mindestens einen Index $\mu_0 \in \{1, \dots, n\}$ erfüllt ist.

Eine Matrix $A \in \mathbb{R}^{n \times n}$ erfüllt das schwache Spaltensummenkriterium, wenn A^T das schwache Zeilensummenkriterium erfüllt.

Satz 4.2.7 (Schwachtes Zeilensummenkriterium) Es sei $A \in \mathbb{R}^{n \times n}$ eine unzerlegbare Matrix, die das schwache Zeilensummenkriterium erfüllt. Dann ist das Jacobi-Verfahren konvergent.

Zum Beweis von Satz 4.2.7 werden wir direkt den Spektralradius der Iterationsmatrix abschätzen. Die wesentliche Beobachtung dabei ist, dass jede irreduzible Matrix, die das schwache Zeilensummenkriterium erfüllt, bereits regulär ist.

Lemma 4.2.8 Jede irreduzible Matrix $M \in \mathbb{R}^{n \times n}$, die das schwache Zeilensummenkriterium erfüllt, ist regulär und für die Diagonalelemente gilt $m_{jj} \neq 0$ ($j = 1, \dots, n$).

Beweis. Wir nehmen an, M sei nicht regulär, d.h. es existiert ein $x \in \mathbb{K}^n \setminus \{0\}$ mit $Mx = 0$. Insbesondere folgt aus der Dreiecksungleichung

$$|m_{jj}| |x_j| \leq \sum_{\substack{\ell=1 \\ \ell \neq j}}^n |m_{j\ell}| |x_\ell| \quad \text{für alle } j = 1, \dots, n. \tag{4.3}$$

Wir definieren die Indextmengen $J := \{j : |x_j| = \|x\|_\infty\}$ und $K := \{k : |x_k| < \|x\|_\infty\}$. Offensichtlich gilt $J \cap K = \emptyset$, $J \cup K = \{1, \dots, n\}$ und $J \neq \emptyset$. Wäre $K = \emptyset$, so könnte man in (4.3) die x_j - und x_ℓ -Terme herauskürzen und erhielte einen Widerspruch zum schwachen

Zeilensummenkriterium. Also gilt auch $K \neq \emptyset$, und aufgrund der Irreduzibilität von M existieren Indizes $j \in J$ und $k \in K$ mit $m_{jk} \neq 0$. Mit diesem ergibt sich

$$|m_{jj}| \leq \sum_{\substack{\ell=1 \\ \ell \neq j}}^n |m_{j\ell}| \frac{|x_\ell|}{|x_j|} < \sum_{\substack{\ell=1 \\ \ell \neq j}}^n |m_{j\ell}|,$$

denn der Quotient ist stets ≤ 1 wegen $|x_j| = \|x\|_\infty$ und er ist < 1 für $\ell = k \in K$. Also erhalten wir einen Widerspruch zum schwachen Zeilensummenkriterium von M , d.h. M ist regulär. Gäbe es schließlich ein triviales Diagonalelement $m_{jj} = 0$, so folgte aus dem schwachen Zeilensummenkriterium, dass bereits die j -te Zeile die Nullzeile wäre. Da M regulär ist, folgt insbesondere $m_{jj} \neq 0$ für alle $j = 1, \dots, n$. \square

Beweis von Satz 4.2.7. Wegen $a_{jj} \neq 0$ für alle $j = 1, \dots, n$ ist $C_J = -D^{-1}(A - D)$ wohldefiniert. Um $\rho(C_J) < 1$ zu zeigen, beweisen wir, dass $M := C_J - \lambda I$ für $\lambda \in \mathbb{C}$ mit $|\lambda| \geq 1$ regulär ist. Da A irreduzibel ist, ist auch $A - D$ irreduzibel, denn es wird lediglich die Diagonale verändert. C_J entsteht durch zeilenweise Multiplikation von $A - D$ mit Werten $\neq 0$. Deshalb ist auch C_J irreduzibel. Da M und C_J sich nur auf der Diagonale unterscheiden, ist M irreduzibel. Aufgrund des schwachen Zeilensummenkriteriums von A gilt

$$\sum_{\substack{k=1 \\ k \neq j}}^n |m_{jk}| = \sum_{\substack{k=1 \\ k \neq j}}^n |c_{jk}^{(J)}| = \sum_{\substack{k=1 \\ k \neq j}}^n \frac{|a_{jk}|}{|a_{jj}|} \leq 1 \leq |\lambda| = |m_{jj}| \quad \text{für alle } j = 1, \dots, n.$$

und für mindestens einen Index j gilt diese Ungleichung strikt. Also erfüllt M auch das schwache Zeilensummenkriterium und ist nach Lemma 4.2.8 insgesamt regulär. \square

Aufgabe 4.2.9 Man zeige, dass das Jacobi-Verfahren auch konvergent ist unter der Annahme, dass $A \in \mathbb{R}^{n \times n}$ irreduzibel ist und das schwache Spaltensummenkriterium erfüllt.

4.3 KONVERGENZ DES GAUSS-SEIDEL-VERFAHRENS

Die Iterationsvorschrift des **Gauß-Seidel-Verfahrens** (auch Einzelschrittverfahren genannt) lautet

$$x^{(i+1)} = (L + D)^{-1}(b - Rx^{(i)}),$$

d.h. die Iterationsmatrix ist $C_{GS} = -(L + D)^{-1}R$.

Satz 4.3.1 (Konvergenzsatz) *Es sei $A \in \mathbb{R}^{n \times n}$ eine reguläre Matrix, die das starke Zeilensummenkriterium erfüllt. Dann ist das Gauß-Seidel-Verfahren zur Lösung des linearen Gleichungssystems $Ax = b$ konvergent. Genauer gilt in diesem Fall die Abschätzung*

$$\|C_{GS}\|_\infty \leq \|C_J\|_\infty < 1.$$

Beweis. Sei das starke Zeilensummenkriterium erfüllt. Die Iterationsmatrizen des Gauß-Seidel-Verfahrens bzw. des Jacobi-Verfahrens werden mit $C_{GS} := -(L + D)^{-1}R$ bzw. $C_J := -D^{-1}(L + R)$ bezeichnet. Wenn das starke Zeilensummenkriterium erfüllt ist, gilt die Abschätzung

$$\|C_J\|_\infty = \max_{1 \leq i \leq n} \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} < 1.$$

Es sei jetzt $y \in \mathbb{R}^n$ beliebig und $z = C_{GS} y$. Durch vollständige Induktion beweisen wir, dass alle Komponenten z_i des Vektors z der Abschätzung

$$|z_i| \leq \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} \|y\|_\infty$$

genügen. Dazu schreiben wir die Gleichung $z = C_{GS} y$ um in $-(L + D)z = Ry$ und schätzen ab:

$$|z_1| \leq \sum_{j=2}^n \frac{|a_{1j}|}{|a_{11}|} |y_j| \leq \sum_{j=2}^n \frac{|a_{1j}|}{|a_{11}|} \|y\|_\infty.$$

Schreiben wir das Gauß-Seidel-Verfahren mit $B = L + D$ in der Form

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{j=1}^{i-1} a_{ij} x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij} x_j^{(k)} \right) \quad (1 \leq i \leq n),$$

so folgt daraus dann mit der Induktionsvoraussetzung

$$\begin{aligned} |z_i| &\leq \frac{1}{|a_{ii}|} \left(\sum_{j=1}^{i-1} |a_{ij}| |z_j| + \sum_{j=i+1}^n |a_{ij}| |y_j| \right) \\ &\leq \frac{1}{|a_{ii}|} \left(\sum_{j=1}^{i-1} |a_{ij}| \|C_J\|_\infty + \sum_{j=i+1}^n |a_{ij}| \right) \|y\|_\infty \leq \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} \|y\|_\infty. \end{aligned}$$


Hiermit erhält man die Abschätzung

$$\|C_{GS} y\|_\infty = \|z\|_\infty \leq \|C_J\|_\infty \|y\|_\infty \quad \forall y \in \mathbb{R}^n$$

und somit

$$\|C_{GS}\|_\infty \leq \|C_J\|_\infty < 1. \tag{4.4}$$

Da $\varrho(C_{GS}) \leq \|C_{GS}\|_\infty$, folgt daraus die Konvergenz des Verfahrens. □

Bemerkung 4.3.2 Häufig verleitet (4.4) zu der falschen Schlußfolgerung, dass das Gauß-Seidel-Verfahren schneller als das Jacobi-Verfahren konvergiert, wenn die Matrix diagonaldominant ist. 

Beispiel 4.3.3 Dass bei strikter Diagonaldominanz einer Matrix $A \in \mathbb{R}^{n \times n}$ aus $\|C_{GS}\|_\infty \leq \|C_J\|_\infty < 1$ nicht $\varrho(C_{GS}) \leq \varrho(C_J)$ folgen muss, sieht man, wenn man die Matrix A folgendermaßen wählt:

$$A = \begin{pmatrix} 50 & -10 & -20 \\ -20 & 49 & -20 \\ 20 & -10 & 49 \end{pmatrix}.$$

Dann besitzen die zugehörigen Iterationsmatrizen

$$C_{GS} = -(L + D)^{-1}R = \frac{1}{12005} \begin{pmatrix} 0 & 2401 & 4802 \\ 0 & 980 & 6860 \\ 0 & -780 & -560 \end{pmatrix} \quad C_J = \frac{1}{245} \begin{pmatrix} 0 & 49 & 98 \\ 100 & 0 & 100 \\ -100 & 50 & 0 \end{pmatrix}$$

nämlich die Spektralradien $\varrho(C_{GS}) = (4\sqrt{5})/49 \approx 0.1825$ und $\varrho(C_J) = 2/49 \approx 0.04082$ sowie die Maximumnormen $\|C_{GS}\|_\infty = 32/49 \approx 0.6531$ und $\|C_J\|_\infty = 40/49 \approx 0.8163$.

Bemerkung 4.3.4 Mit Bezug auf das letzte Beispiel halten wir fest: Ist eine Matrix A strikt diagonaldominant, gilt für die Iterationsmatrizen $\|C_{GS}\|_\infty \leq \|C_J\|_\infty < 1$, es folgt im Allgemeinen aber nicht $\varrho(C_{GS}) \leq \varrho(C_J)$.

Satz 4.3.5 (Schwachtes Zeilensummenkriterium) Ist $A \in \mathbb{R}^{n \times n}$ irreduzibel und erfüllt das schwache Zeilensummenkriterium, so ist das Gauß-Seidel-Verfahren konvergent.

Beweis. Die Wohldefiniertheit von $C_{GS} = -(L + D)^{-1}R$ ist wieder klar. Wir betrachten $M := C_{GS} - \lambda I$ für $\lambda \in \mathbb{C}$ mit $|\lambda| \geq 1$. Durch Multiplikation mit $-(L + D)$ sieht man, dass M genau dann regulär ist, wenn $\widetilde{M} := R + \lambda L + \lambda D$ regulär ist. Offensichtlich erbt \widetilde{M} die Irreduzibilität von $A = D + L + R$. Ferner erfüllt \widetilde{M} das schwache Zeilensummenkriterium, denn es gilt

$$\sum_{\substack{k=1 \\ k \neq j}}^n |\widetilde{m}_{jk}| = |\lambda| \sum_{k=1}^{j-1} |a_{jk}| + \sum_{k=j+1}^n |a_{jk}| \leq |\lambda| \sum_{\substack{k=1 \\ k \neq j}}^n |a_{jk}| \leq |\lambda| |a_{jj}| = |\widetilde{m}_{jj}| \quad \text{für } j = 1, \dots, n$$

mit strikter Ungleichung für mindestens einen Index j . Nach Lemma 4.2.8 ist \widetilde{M} regulär. Insgesamt erhalten wir wie zuvor $\varrho(C_{GS}) < 1$. \square

Beispiel 4.3.6 Es sei $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ mit $a, b, c, d \in \mathbb{C}$. Die zugehörige Iterationsmatrix zum Jacobi-Verfahren lautet somit

$$C_J := -D^{-1}(L + R) = - \begin{pmatrix} a^{-1} & 0 \\ 0 & d^{-1} \end{pmatrix} \begin{pmatrix} 0 & b \\ c & 0 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{b}{a} \\ -\frac{c}{d} & 0 \end{pmatrix}.$$

Das charakteristische Polynom hierzu lautet $p(\lambda) = \lambda^2 - \frac{bc}{ad}$ und hat Nullstellen $\lambda_{1,2} = \pm \sqrt{\frac{bc}{ad}}$. Entsprechend erhält man für das Gauß-Seidel-Verfahren

$$C_{GS} := -(L + D)^{-1}R = -\frac{1}{ad} \begin{pmatrix} d & 0 \\ -c & a \end{pmatrix} \begin{pmatrix} 0 & b \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0 & -\frac{bd}{ad} \\ 0 & \frac{bc}{ad} \end{pmatrix}.$$

Hier lautet das charakteristische Polynom $p(\lambda) = \lambda(\lambda - bc/ad)$ und hat Nullstellen

$$\lambda_1 = 0, \quad \lambda_2 = \frac{bc}{ad}.$$

womit

$$\varrho(C_J) = \sqrt{\frac{|bc|}{|ad|}} \quad \text{und} \quad \varrho(C_{GS}) = \frac{|bc|}{|ad|}$$

gilt. Man beachte

$$\|C_J\|_1 = \|C_J\|_\infty = \max\{|b/a|, |c/d|\}, \quad \|C_{GS}\|_1 = \frac{|b|(|c| + |d|)}{|ad|}, \quad \|C_{GS}\|_\infty = \frac{|b| \max\{|c|, |d|\}}{|ad|}.$$

Wir können somit für A aus $\mathbb{R}^{2 \times 2}$ festhalten: Konvergiert das Jacobi- oder Gauß-Seidel-Verfahren, so konvergiert auch das jeweilige andere. Und im Falle der Konvergenz, konvergiert das Gauß-Seidel doppelt so schnell wie das Jacobi-Verfahren. Gilt dies immer, bzw. kann dies ggf. einfach charakterisiert werden?

Aufgabe 4.3.7 Man nutze das letzte Beispiel, um zu zeigen, dass aus $\|C_J\|_1 < 1$ im Allgemeinen

$$\|C_{GS}\|_1 \leq \|C_J\|_1 < 1$$

nicht folgt. (Z.B. $a = d = 1$, $b = c = 2/3$ in der Matrix von Beispiel 4.3.6 liefert ein Gegenbeispiel).

Definition 4.3.8 Eine Matrix $A \in \mathbb{R}^{m \times n}$ heißt nichtnegativ, wenn alle Koeffizienten a_{ij} von A nichtnegativ sind.

Satz 4.3.9 (von Stein und Rosenberg [Hämmerlin/Hoffmann]) Die Iterationsmatrix $C_J \in \mathbb{R}^{n \times n}$ des Jacobi-Verfahrens sei nichtnegativ. Dann gilt genau eine der folgenden Aussagen:

- i) $\rho(C_J) = \rho(C_{GS}) = 0$
- ii) $\rho(C_J) = \rho(C_{GS}) = 1$
- iii) $0 < \rho(C_{GS}) < \rho(C_J) < 1$
- iv) $1 < \rho(C_J) < \rho(C_{GS})$

Beispiel 4.3.10

$$A = \begin{pmatrix} 2 & -1 & 0 & 0 \\ -1 & 2 & -1 & 0 \\ 0 & -1 & 2 & -1 \\ 0 & 0 & -1 & 2 \end{pmatrix} \Rightarrow C_J = -D^{-1}(L + R) = \begin{pmatrix} 0 & \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & 0 \end{pmatrix}$$

Die Iterationsmatrix ist nichtnegativ und somit folgt, dass das Gauß-Seidel-Verfahren schneller ist als das Jacobi-Verfahren!

Bemerkung 4.3.11 (Jacobi- immer schlechter als Gauß-Seidel-Verfahren?) Dass das Gauß-Seidel-Verfahren nicht immer besser sein muss als das Jacobi-Verfahren oder aus der Divergenz des Jacobi-Verfahrens nicht auch die Divergenz des Gauß-Seidel-Verfahrens folgen muss, zeigen die folgenden beiden Beispiele.

i.) Die Iterationsmatrizen C_J bzw. C_{GS} zur Matrix

$$A = \begin{pmatrix} 1 & 2 & -2 & 2 \\ 1 & 1 & 1 & 0 \\ 2 & 2 & 1 & 2 \\ -1 & -2 & 1 & 1 \end{pmatrix}$$

lauten

$$C_J = \begin{pmatrix} 0 & -2 & 2 & -2 \\ -1 & 0 & -1 & 0 \\ -2 & -2 & 0 & -2 \\ 1 & 2 & -1 & 0 \end{pmatrix} \quad \text{bzw.} \quad C_{GS} = -(L + D)^{-1}R = \begin{pmatrix} 0 & -2 & 2 & -2 \\ 0 & 2 & -3 & 2 \\ 0 & 0 & 2 & -2 \\ 0 & 2 & -6 & 4 \end{pmatrix}$$

mit den Spektralradien $\rho(C_J) = 0$ und $\rho(C_{GS}) \approx 7.3850$.

ii.) Die Iterationsmatrizen C_J bzw. C_{GS} zur Matrix

$$A = \frac{1}{3} \begin{pmatrix} 3 & -2 & -1 & 1 \\ 1 & 2 & -2 & 1 \\ 1 & -2 & 2 & 2 \\ 1 & -2 & 1 & 1 \end{pmatrix}$$

lauten

$$C_J = \frac{1}{6} \begin{pmatrix} 0 & 4 & 2 & -2 \\ -3 & 0 & 6 & -3 \\ -3 & 6 & 0 & -6 \\ -6 & 12 & -6 & 0 \end{pmatrix} \quad \text{bzw.} \quad C_{GS} = -(L + D)^{-1}R = \frac{1}{6} \begin{pmatrix} 0 & 4 & 2 & -2 \\ 0 & -2 & 5 & -2 \\ 0 & -4 & 4 & -7 \\ 0 & -4 & 4 & 5 \end{pmatrix}$$

mit den Spektralradien $\rho(C_J) \approx 1.4527$ und $\rho(C_{GS}) \approx 0.9287 < 1$.

Bemerkungen 4.3.12 i) Die obigen Iterationsverfahren ließen sich in der Form $x^{(k+1)} = B^{-1}(B - A)x^{(k)} + B^{-1}b = Cx^{(k)} + d$ schreiben. Da die Iterationsmatrix C für alle k konstant ist, spricht man von stationären Iterationsverfahren.

ii) Das quantitative Verhalten solcher stationärer Verfahren lässt sich durch die Einführung eines (Relaxations-) Parameters ω verbessern:

$$x^{(k+1)} = \omega(Cx^{(k)} + d) + (1 - \omega)x^{(k)}.$$

Für $0 < \omega < 1$ spricht man von einem **Unterrelaxationsverfahren** und für $\omega > 1$ von einem **Überrelaxationsverfahren**. Man kann zeigen, dass der optimale Parameter für eine positiv definite Matrix A beim gedämpften Jacobiverfahren

$$\omega_{\text{opt}} = \frac{2}{\lambda_{\min}(D^{-1}A) + \lambda_{\max}(D^{-1}A)}$$

lautet und für das überrelaxierte Gauß-Seidel-Verfahren (SOR = successive overrelaxation method) angewandt auf eine positiv definite Matrix $A = L + D + L^T$ ergibt sich der optimale Parameter zu

$$\omega_{\text{opt}} = \frac{2}{1 + \sqrt{\lambda_{\min}(D^{-1}A) + \lambda_{\max}((D + 2L)D^{-1}(D + 2L^T)A^{-1})}}.$$

Ergebnisse für allgemeinere Fälle findet man z.B. bei [Niethammer].

iii) Die Bedeutung der oben genannten Iterationsverfahren liegt heute weniger in deren unmittelbarem Einsatz zur Lösung von $Ax = b$, sondern auf Grund deren „**Glättungseigenschaften**“ als Beschleuniger anderer moderner Verfahren (vorkonditioniertes konjugiertes Gradienten-Verfahren, Mehrgitter).

4.4 ABBRUCH-KRITERIEN

Da ein Iterationsverfahren aufeinanderfolgende Näherungen der Lösung liefert, ist ein praktischer Test notwendig, um das Verfahren zu stoppen, wenn die gewonnene Approximation genau genug ist. Da es nicht möglich ist, den Fehler $x - x^{(k)}$ zu bestimmen, d.h. den Abstand zur eigentlichen (gesuchten) Lösung zu bestimmen, müssen andere Quantitäten gewonnen werden, die meist auf dem Residuum $r = b - Ax$ basieren.

Die hier vorgestellten Verfahren liefern eine Folge $(x^{(i)})$ von Vektoren die gegen den Vektor x streben, welcher Lösung des linearen Gleichungssystems $Ax = b$ ist. Um effizient zu sein, muss die Methode wissen, wann sie abbrechen soll. Eine gute Methode sollte

- i) feststellen, ob der Fehler $e^{(i)} := x - x^{(i)}$ klein genug ist,
- ii) abbrechen, falls der Fehler nicht weiter kleiner wird oder nur noch sehr langsam, und
- iii) den maximalen Aufwand, der zur Iteration verwendet wird, beschränken.

Das folgende Abbruchkriterium ist eine einfache aber häufig genügende Variante. Man muss hierzu die Quantitäten \maxit , $\|b\|$, tol und wenn möglich auch $\|A\|$ (und $\|A^{-1}\|$) zur Verfügung stellen. Dabei ist

- die natürliche Zahl \maxit die maximale Anzahl an möglichen Iterationen des Verfahrens,
- die reelle Zahl $\|A\|$ eine Norm von A , (Jede einigermaßen vernünftige Approximation des betragsgrößten Eintrags in A genügt schon.)

- die reelle Zahl $\|b\|$ eine Norm von b , (Auch hier genügt eine einigermaßen vernünftige Approximation des betragsgrößten Eintrags in b .)
- die reelle Zahl tol eine Schranke für die Größe des Residuums bzw. des Fehlers.

MATLAB-Beispiel: Beispiel eines vernünftigen Abbruchkriteriums

```

k = 0;
while 1
    k = k + 1;
    % Berechne die Approximation x^(k)
    % Berechne das Residuum r^(k) = b - A x^(k)
    % Berechne norm_ak = || A * x^(k) ||, norm_rk = || r^(k) ||
    % und norm_b = || b ||
    if (k >= maxit) | ( norm_rk <= tol * max( norm_ak, norm_b ) )
        break
    end
end
end

```

Da sich nach einigen Iterationen der Term $\|Ax^{(k)}\|$ nicht mehr groß ändert, braucht man diesen nicht immer neu zu bestimmen. Zu bestimmen ist eigentlich $\|e^{(k)}\| = \|A^{-1}r^{(k)}\| \leq \|A^{-1}\| \|r^{(k)}\|$. Man beachte, dass man $\|A^{-1}B\|$ bei den bisherigen Verfahren mit Lemma 3.2.6 zur Neumannschen-Reihe abschätzen kann. Es gilt $x^{(k+1)} = B^{-1}(B - A)x^{(k)} + B^{-1}b = Cx^{(k)} + d$ und

$$B^{-1}A = I - B^{-1}(B - A) = I - C \quad \text{sowie} \quad \|A^{-1}B\| = \|(I - C)^{-1}\| \leq \frac{1}{1 - \|C\|}.$$

4.5 GRADIENTEN–VERFAHREN

Im Folgenden nehmen wir stets an, dass

$$A \in \mathbb{R}^{n \times n} \quad \text{positiv definit (s.p.d.) ist.} \quad (4.5)$$

Wir ordnen nun dem Gleichungssystem $Ax = b$ die Funktion $f(x) := \frac{1}{2}x^T Ax - b^T x$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$ zu. Der Gradient von f ist $f'(x) = \frac{1}{2}(A + A^T)x - b$. Da $A = A^T$ nach Voraussetzung (4.5), lautet die Ableitung

$$f'(x) = \nabla f(x) = Ax - b.$$

Notwendig für ein Minimum von f ist das Verschwinden des Gradienten, d.h. $Ax = b$. Da die Hesse-Matrix $f''(x) = A$ positiv definit ist, liegt für die Lösung von $Ax = b$ tatsächlich ein Minimum vor.

Mit $\arg \min_{y \in \mathbb{R}^n} f(y)$ bezeichnen wir denjenigen Wert aus \mathbb{R}^n , der den Term f minimiert, d.h. $f(x) = \min_{y \in \mathbb{R}^n} f(y)$, wenn $x = \arg \min_{y \in \mathbb{R}^n} f(y)$ ist.

Idee:

$$Ax = b \iff x = \arg \min_{y \in \mathbb{R}^n} f(y)$$

mit $f(y) := \frac{1}{2}y^T Ay - b^T y$, $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

Bewiesen haben wir soeben das folgende Lemma.

Lemma 4.5.1 *Es sei $A \in \mathbb{R}^{n \times n}$ positiv definit, dann gilt*

$$Ax = b \iff x = \arg \min_{y \in \mathbb{R}^n} f(y).$$

Alternativer Beweis zu Lemma 4.5.1. Ein zweiter Beweis folgt aus der Darstellung

$$f(x) = f(x^*) + \frac{1}{2}(x - x^*)^T A(x - x^*) \quad \text{mit } x^* := A^{-1}b. \quad (4.6)$$

Hieraus folgt $f(x) > f(x^*)$ für $x \neq x^*$, d.h. $x^* := A^{-1}b$ ist das eindeutige Minimum von f . Man beachte dabei, dass (4.6) ein Sonderfall der folgenden Entwicklung von f um einen beliebigen Wert $\tilde{x} \in \mathbb{R}^n$ ist, welche sich durch ausmultiplizieren zeigen lässt:

$$f(x) = f(\tilde{x}) + \langle A\tilde{x} - b, x - \tilde{x} \rangle + \frac{1}{2}\langle A(x - \tilde{x}), x - \tilde{x} \rangle$$

□

Folgerung: Man kann also Verfahren zur numerischen Optimierung/Minimierung verwenden. Der **Gradient** ist die Richtung des **steilsten Anstiegs**, also kann man $-\nabla f$ als Abstiegsrichtung wählen und entlang dieser Geraden ein Minimum suchen.

Gradientenverfahren (allgemein):

Es sei $\Omega \subseteq \mathbb{R}^n$, $f: \Omega \rightarrow \mathbb{R}$, $x_0 \in \Omega$ Startwert, für $k = 1, 2, 3, \dots$

- 1) *Bestimmung der Abstiegsrichtung* $d_k = -\nabla f(x_k)$
- 2) *Liniensuche:* suche auf der Geraden $\{x_k + td_k : t \geq 0\} \cap \Omega$ ein Minimum, d.h. bestimme $\alpha_k \geq 0 : f(x_k + \alpha_k d_k) \leq f(x_k)$

$$x_{k+1} = x_k + \alpha_k d_k.$$

Daraus folgt $f(x_0) \geq f(x_1) \geq f(x_2) \geq \dots$

Für die quadratische Funktionen $f(x) = \frac{1}{2}x^T Ax - b^T x$ und $\Omega \equiv \mathbb{R}^n$ kann man 1) und 2) leicht berechnen: Da $\nabla f(x) = Ax - b$, ergibt sich $d_k = -\nabla f(x_k) = b - Ax_k$. Sei $p \in \mathbb{R}^n \setminus \{0\}$ und $F(\lambda) := f(x + \lambda p)$, dann gilt für die Liniensuche

$$\begin{aligned} F(\lambda) &= f(x + \lambda p) \\ &= \frac{1}{2}\langle x + \lambda p, A(x + \lambda p) \rangle - \langle b, x + \lambda p \rangle \\ &= \frac{1}{2}\langle x, Ax \rangle - \langle b, x \rangle + \lambda\langle p, Ax - b \rangle + \frac{1}{2}\lambda^2\langle p, Ap \rangle \\ &= f(x) + \lambda\langle p, Ax - b \rangle + \frac{1}{2}\lambda^2\langle p, Ap \rangle \end{aligned} \tag{4.7}$$

Da $p \neq 0$ nach Voraussetzung, folgt $\langle p, Ap \rangle > 0$. F ist also eine quadratische Funktion mit positivem führenden Koeffizienten. Somit folgt aus

$$0 \stackrel{!}{=} F'(\lambda) = \langle p, Ax - b \rangle + \lambda\langle p, Ap \rangle, \tag{4.8}$$

dass der Parameter

$$\lambda_{opt}(x, p) = \frac{\langle p, b - Ax \rangle}{\langle p, Ap \rangle}.$$

zu gegebenem Vektor $p \in \mathbb{R}^n \setminus \{0\}$ das Funktional $F(\lambda) := f(x + \lambda p)$ minimiert.

Für allgemeine Funktionen wird die Liniensuche angenähert, z.B. mit der **Schrittweitenregel von Armijo**

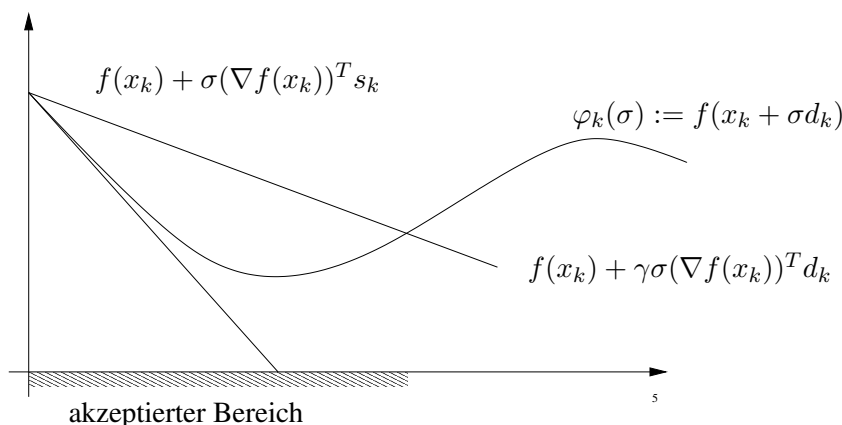


Abb. 4.3: Optimale Wahl des Dämpfungsparameters .

Schrittweitenregel von Armijo:

Wähle $\beta \in (0, 1)$ (z.B. $\beta = \frac{1}{2}$) und $\gamma \in (0, 1)$ (z.B. $\gamma \in [10^{-3}, 10^{-2}]$)

Bestimme die größte Schrittweite $\sigma_k \in \{1, \beta, \beta^2, \beta^3, \dots\}$ mit

$$f(x_k) - f(x_k - \sigma_k d_k) \geq -\gamma\sigma_k \nabla f(x_k)^T d_k,$$

d.h.

$$f(x_k - \sigma_k d_k) \leq f(x_k) + \gamma\sigma_k \nabla f(x_k)^T d_k$$

Wir untersuchen nun die Konvergenz des Verfahrens für quadratische Funktionen. Hierzu bietet sich die sogenannte *Energienorm* an

$$\|x\| := \sqrt{x^T Ax}, \quad (A \in \mathbb{R}^{n \times n})$$

(beachte: alle Normen auf dem $\mathbb{R}^{n \times n}$ sind äquivalent).

Algorithmus 4.5.1: Gradientenverfahren: $Ax = b$

Input: Initial guess x^0

$$r^0 := b - Ax^0$$

Iteration: $k = 0, 1, \dots$

$$a^k := Ar^k$$

$$\lambda_{opt} := \langle r^k, r^k \rangle / \langle r^k, a^k \rangle$$

$$x^{k+1} := x^k + \lambda_{opt} r^k$$

$$r^{k+1} := r^k - \lambda_{opt} a^k$$

Man beachte $r^{k+1} = b - Ax^{k+1} = b - A(x^k + \lambda_{opt} r^k) = r^k - \lambda_{opt} Ar^k$

Lemma 4.5.2 *Es sei $A \in \mathbb{R}^{n \times n}$ positiv definit und $x^* \in \mathbb{R}^n$ erfülle $Ax^* = b$. Dann gilt*

$$\|x_{k+1} - x^*\|_A^2 \leq \|x_k - x^*\|_A^2 \left(1 - \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle \langle r_k, A^{-1}r_k \rangle} \right).$$

Beweis. Es gilt

$$\begin{aligned} f(x_{k+1}) &= f(x_k + \lambda_{opt} r_k) = \frac{1}{2} \lambda_{opt}^2 \langle r_k, Ar_k \rangle - \lambda_{opt} \langle r_k, r_k \rangle + f(x_k) \\ &= f(x_k) + \frac{1}{2} \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle} - \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle} = f(x_k) - \frac{1}{2} \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle} \end{aligned} \quad (4.9)$$

Für die exakte Lösung x^* von $Ax = b$ gilt für $x \in \mathbb{R}^n$

$$\begin{aligned} f(x^*) + \frac{1}{2} \|x - x^*\|_A^2 &= \frac{1}{2} \langle x^*, Ax^* \rangle - \langle b, x^* \rangle + \frac{1}{2} \langle (x - x^*), A(x - x^*) \rangle \\ &= \frac{1}{2} \langle x^*, Ax^* \rangle - \langle x^*, b \rangle + \frac{1}{2} \langle x, Ax \rangle - \langle x, Ax^* \rangle + \frac{1}{2} \langle x^*, Ax^* \rangle \\ &= \langle x^*, Ax^* - b \rangle + \frac{1}{2} \langle x, Ax \rangle - \langle x, b \rangle \\ &= f(x), \quad (\text{d.h. } \|x - x^*\|_A^2 = 2(f(x) - f(x^*))) \end{aligned}$$

also mit (4.9)

$$\begin{aligned} \|x_{k+1} - x^*\|_A^2 &= 2f(x_{k+1}) - 2f(x^*) = 2f(x_{k+1}) - 2f(x_k) + \|x_k - x^*\|_A^2 \\ &= \|x_k - x^*\|_A^2 - \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle} \end{aligned}$$

Mit $r_k = b - Ax_k = A(x^* - x_k)$ folgt wegen

$$\begin{aligned} \|x_k - x^*\|_A^2 &= \|x^* - x_k\|_A^2 = \langle (x^* - x_k), A(x^* - x_k) \rangle \\ &= \langle A^{-1}A(x^* - x_k), r_k \rangle = \langle r_k, A^{-1}r_k \rangle \end{aligned}$$

die Behauptung. □

Frage: Was sagt Lemma 4.5.2 bzgl. der Konvergenz und Kondition aus?

Lemma 4.5.3 (Kantorowitsch-Ungleichung)

Es sei $A \in \mathbb{R}^{n \times n}$ positiv definit und $\kappa := \kappa_2(A) := \|A\|_2 \|A^{-1}\|_2$. Dann gilt für alle $x \in \mathbb{R} \setminus \{0\}$

$$\frac{\langle x, Ax \rangle \langle x, A^{-1}x \rangle}{\langle x, x \rangle^2} \leq \frac{1}{4} \left(\sqrt{\kappa} + \sqrt{\kappa^{-1}} \right)^2$$

Beweis:

Beweis. Die Eigenwerte von A seien geordnet gemäß

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n, \quad \kappa = \frac{\lambda_n}{\lambda_1}$$

Da A symmetrisch ist, existiert eine orthonormale Matrix Q mit $Q^T A Q = \Lambda = \text{diag}(\lambda_i)$. Für $y = Q^T x$ gilt dann

$$x^T A x = x^T Q \Lambda Q^T x = y^T \Lambda y = \sum_{i=1}^n \lambda_i y_i^2, \quad x^T A^{-1} x = x^T Q \Lambda^{-1} Q^T x = \sum_{i=1}^n \lambda_i^{-1} y_i^2$$

sowie $x^T x = x^T Q Q^T x = y^T y$, also

$$\frac{\langle x, Ax \rangle \langle x, A^{-1}x \rangle}{\langle x, x \rangle^2} = \frac{\left(\sum_i \lambda_i y_i^2 \right) \left(\sum_i \lambda_i^{-1} y_i^2 \right)}{\|y\|_2^4} = \left(\sum_{i=1}^n \lambda_i z_i \right) \left(\sum_{i=1}^n \lambda_i^{-1} z_i \right) \quad (4.10)$$

mit $z_i := \frac{y_i^2}{\|y\|_2^2}$. Man beachte $\sum_{i=1}^n z_i = 1$. Für $\lambda_1 \leq x \leq \lambda_n$ folgt

$$0 \geq (x - \lambda_1)(x - \lambda_n) = x^2 - x(\lambda_1 + \lambda_n) + \lambda_1 \lambda_n$$

und somit

$$\lambda_1 \lambda_n + \lambda_k^2 \leq \lambda_k (\lambda_1 + \lambda_n) \quad \Rightarrow \quad \frac{\lambda_1 \lambda_n}{\lambda_k} + \lambda_k \leq \lambda_1 + \lambda_n \quad (k = 1, \dots, n).$$

Es gilt nun

$$\lambda_1 \lambda_n \sum_{i=1}^n \lambda_i^{-1} z_i + \sum_{i=1}^n \lambda_i z_i = \left(\frac{\lambda_1 \lambda_n}{\lambda_1} + \lambda_1 \right) z_1 + \left(\frac{\lambda_1 \lambda_n}{\lambda_2} + \lambda_2 \right) z_2 + \dots + \left(\frac{\lambda_1 \lambda_n}{\lambda_n} + \lambda_n \right) z_n \leq \lambda_1 + \lambda_n,$$

d.h.

$$\sum_{i=1}^n \lambda_i^{-1} z_i \leq \frac{\lambda_1 + \lambda_n - \lambda}{\lambda_1 \lambda_n}$$

mit $\lambda := \sum_{i=1}^n \lambda_i z_i$. Somit lässt sich (4.10) abschätzen durch

$$\frac{\langle x, Ax \rangle \langle x, A^{-1}x \rangle}{\langle x, x \rangle^2} \leq \lambda \cdot \underbrace{\frac{\lambda_1 + \lambda_n - \lambda}{\lambda_1 \lambda_n}}_{=: h(\lambda)}$$

Für welches λ wird nun das Polynom h maximal?

$$\begin{aligned} h'(\lambda) &= \frac{\lambda_1 + \lambda_n - \lambda}{\lambda_1 \lambda_n} - \frac{\lambda}{\lambda_1 \lambda_n} = \frac{\lambda_1 + \lambda_n}{\lambda_1 \lambda_n} - \lambda \frac{2}{\lambda_1 \lambda_n} \stackrel{!}{=} 0 \quad \Rightarrow \quad \lambda^* = \frac{\lambda_1 + \lambda_n}{2} \\ h''(\lambda) &= -\frac{2}{\lambda_1 \lambda_n} < 0 \quad \Rightarrow \quad \lambda^* \text{ maximiert } h \end{aligned}$$

d.h.

$$\max_{\lambda \in [\lambda_1, \lambda_n]} h(\lambda) = h(\lambda^*) = \frac{(\lambda_1 + \lambda_n)^2}{4\lambda_1\lambda_n} = \frac{1}{4} \left(\sqrt{\frac{\lambda_1}{\lambda_n}} + \sqrt{\frac{\lambda_n}{\lambda_1}} \right)^2 = \frac{1}{4} \left(\sqrt{\kappa} + \sqrt{\kappa^{-1}} \right)^2.$$

□

Satz 4.5.4 Es sei $A \in \mathbb{R}^{n \times n}$ positiv definit. Dann gilt für das Gradientenverfahren

$$\|x_k - x^*\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^k \|x_0 - x^*\|_A$$

Beweis. Lemma 4.5.2 liefert die Abschätzung

$$\|x_{k+1} - x^*\|_A^2 \leq \|x_k - x^*\|_A^2 \left(1 - \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle \langle r_k, A^{-1}r_k \rangle} \right).$$

und mit Lemma 4.5.3 ergibt sich

$$\begin{aligned} 1 - \frac{\langle r_k, r_k \rangle^2}{\langle r_k, Ar_k \rangle \langle r_k, A^{-1}r_k \rangle} &\leq 1 - 4(\sqrt{\kappa} + \sqrt{\kappa^{-1}})^{-2} = \frac{(\sqrt{\kappa} + \sqrt{\kappa^{-1}})^2 - 4}{(\sqrt{\kappa} + \sqrt{\kappa^{-1}})^2} \\ &= \frac{\kappa + 2 + \kappa^{-1} - 4}{\kappa + 2 + \kappa^{-1}} = \frac{\kappa^2 - 2\kappa + 1}{\kappa^2 + 2\kappa + 1} = \left(\frac{\kappa - 1}{\kappa + 1} \right)^2. \end{aligned}$$

□

Bemerkungen:

(a) Für große κ gilt

$$\frac{\kappa - 1}{\kappa + 1} = \underbrace{\frac{\kappa}{\kappa + 1}}_{\approx 1} - \frac{1}{\kappa + 1} \approx 1 - \frac{1}{\kappa + 1},$$

also sehr nahe bei 1, d.h. geringe Konvergenzgeschwindigkeit!

(b) Dies tritt auch schon bei einfachen Beispielen auf:

$$A = \begin{pmatrix} 1 & 0 \\ 0 & a \end{pmatrix}, \quad a \gg 1, \quad b = \begin{pmatrix} 0 \\ a \end{pmatrix} \quad \text{und} \quad x_0 = \begin{pmatrix} a \\ 1 \end{pmatrix}$$

daraus folgt (*Übung*).

$$\begin{pmatrix} x_{k+1} \\ y_{k+1} \end{pmatrix} = \rho \begin{pmatrix} x_k \\ -y_k \end{pmatrix} \quad \text{mit} \quad \rho = \frac{a - 1}{a + 1}$$

wegen $a = \kappa_2(A)$ ist das genau die Rate aus Satz 4.5.4!

(c) Anschaulich sieht man ein „Zick-Zack-Verhalten“ der Iteration.

4.6 VERFAHREN DER KONJUGIERTEN GRADIENTEN

- entstanden 1952 von Hestenes und Stiefel
- enormer Gewinn an Bedeutung 1971 durch Vorkonditionierung
- gehören zu den schnellsten Verfahren, werden heute sehr oft verwendet

- für die 2D-Standard-Matrix ab Systemgröße 2000–4000 besser als Gauß bei zusätzlich erheblich geringerem Speicherbedarf.

Idee: Vermeide Zick-Zack-Verhalten durch Verwendung von Orthogonalität bzgl.

$$(x, y)_A = x^T A y,$$

daraus folgt „konjugierte Gradienten“: *cg-Verfahren* (conjugate gradient)

Bemerkungen 4.6.1

- (a) Zwei Vektoren $x, y \in \mathbb{R}^n$ heißen **konjugiert** oder **A-orthogonal**, falls $(x, y)_A = 0$.
- (b) $\{x_1, \dots, x_k\}$ paarweise konjugiert, d.h. $(x_i, x_j)_A = \delta_{ij} \|x_i\|_A^2$, $x_i \neq 0$ ($i, j \in \{1, \dots, k\}$)
daraus folgt $\{x_1, \dots, x_k\}$ linear unabhängig.
- (c) Jeder Vektor $x \in \mathbb{R}^n$ besitzt eine eindeutige Entwicklung

$$x = \sum_{k=1}^n \alpha_k d_k \tag{4.11}$$

bezüglich konjugierter Richtungen $\{d_1, \dots, d_n\}$. Aus (4.11) folgt

$$(x, d_i)_A = \sum_{k=1}^n \alpha_k \underbrace{(d_k, d_i)_A}_{=\delta_{ik} \|d_i\|_A^2} = \alpha_i \|d_i\|_A^2,$$

also

$$\alpha_k = \frac{d_k^T A x}{d_k^T A d_k} \quad (k = 1, \dots, n). \tag{4.12}$$

- (d) Für die Lösung x^* von $Ax = b$ gilt offenbar

$$\alpha_i = \frac{d_i^T b}{d_i^T A d_i}$$

Lemma 4.6.2 Seien $\{d_1, \dots, d_n\}$ konjugierte Richtungen. Für jedes $x_0 \in \mathbb{R}^n$ und

$$x_k = x_{k-1} + \alpha_k d_k, \quad \alpha_k = \frac{r_{k-1}^T d_k}{d_k^T A d_k}, \quad r_k := b - A x_k \quad (k \geq 1) \tag{4.13}$$

gilt nach (höchstens) n Schritten $x_n = A^{-1}b$.

Beweis. Aus (4.12), (4.13) folgt für $x^* = A^{-1}b$,

$$x^* - x_0 = \sum_{k=1}^n \tilde{\alpha}_k d_k \quad \text{mit} \quad \tilde{\alpha}_k = \frac{d_k^T A(x^* - x_0)}{d_k^T A d_k} = \frac{d_k^T (b - A x_0)}{d_k^T A d_k}.$$

Da d_k zu $d_i, i \neq k$ konjugiert ist, gilt

$$d_k^T A(x_{k-1} - x_0) = d_k^T A \left(\sum_{i=1}^{k-1} \hat{\alpha}_i d_i \right) = \sum_{i=1}^{k-1} \hat{\alpha}_i \underbrace{d_k^T A d_i}_{=0, \text{ da } i \neq k} = 0,$$

also

$$\begin{aligned} d_k^T A(x^* - x_0) &= d_k^T A(x^* - x_{k-1}) + \underbrace{d_k^T A(x_{k-1} - x_0)}_{=0} \\ &= d_k^T (b - A x_{k-1}) = d_k^T r_{k-1} \Rightarrow \tilde{\alpha}_k = \alpha_k. \end{aligned}$$

□

Bemerkungen 4.6.3

- (a) $r_k = b - Ax_k$ wird als **Residuum** von $Ax = b$ bzgl x_k bezeichnet.
 (b) Lemma 4.6.2 besagt, dass das Verfahren ein direktes Verfahren ist, welches nach n Iterationen konvergiert.

$$A \text{ sparse} \Rightarrow O(n^2)$$

(dies ist nicht optimal).

- (c) Wie in (4.8) gilt

$$\begin{aligned} \frac{d}{d\lambda} f(x_k + \lambda d_k) &= \lambda \langle d_k, Ad_k \rangle + \langle d_k, (Ax_k - b) \rangle \\ f(x_k) = f(x_{k-1} + \alpha_k d_k) &= \min_{\lambda \in \mathbb{R}} : f(x_{k-1} + \lambda d_k), \end{aligned}$$

dann

$$\lambda_{opt} = -\frac{\langle d_k, (Ax_{k-1} - b) \rangle}{\langle d_k, Ad_k \rangle} = \frac{\langle r_{k-1}, d_k \rangle}{\langle d_k, Ad_k \rangle} = \alpha_k$$

Satz 4.6.4 Seien $\{d_1, \dots, d_n\}$ konjugierte Richtungen und r_k ($k = 0, \dots, n-1$), die durch (4.13) definierten Residuen. Dann gilt

$$r_k^T d_j = 0 \text{ bzw. } r_k \perp U_k := \text{span}\{d_1, \dots, d_k\} \quad (1 \leq k \leq n, 1 \leq j \leq k). \quad (4.14)$$

Beweis. Nach (4.13) gilt für $k \in \{1, \dots, n\}$

$$r_k = b - Ax_k = r_{k-1} - \alpha_k Ad_k = r_{k-2} - \alpha_{k-1} Ad_{k-1} - \alpha_k Ad_k = \dots = r_0 - \sum_{j=1}^k \alpha_j Ad_j.$$

Daraus folgt nun für $1 \leq j \leq k$

$$\begin{aligned} r_k^T d_j &= r_0^T d_j - \sum_{\ell=1}^k \alpha_\ell d_\ell^T Ad_j \\ &= r_0^T d_j - \alpha_j d_j^T Ad_j = r_0^T d_j - \frac{r_{j-1}^T d_j}{d_j^T Ad_j} d_j^T Ad_j \\ &= (r_0^T - r_{j-1}^T) d_j = \sum_{\ell=1}^{j-1} \alpha_\ell d_\ell^T Ad_j = 0 \end{aligned}$$

womit der Satz bewiesen ist. □

Frage: Wie sind nun die d_k und damit erzeugten Teilräume $\text{span}\{d_1, \dots, d_k\}$ zu wählen?

Falls $r_0 \neq 0$ gilt (sonst ist x_0 schon die gesuchte Lösung) setzt man $d_1 = r_0$.

Für $k = 1, 2, 3, \dots$ verfahren wir wie folgt. Falls $r_k \neq 0$ (sonst ist x_k schon die gesuchte Lösung), gewinnt man formal d_{k+1} mittels des Gram-Schmidtschen-Orthogonalisierungsverfahren aus r_k und den schon bestimmten konjugierten Richtungen d_1, \dots, d_k , d.h.

$$d_{k+1} = r_k - \sum_{j=1}^k \frac{\langle r_k, Ad_j \rangle}{\langle d_j, Ad_j \rangle} d_j. \quad (4.15)$$

Aufgabe 4.6.5 Man zeige induktiv, dass die so erzeugten A -orthogonalen Richtungen

$$d_{k+1} \notin \text{span}\{r_0, \dots, r_{k-1}\} \quad \text{und} \quad d_{k+1} \in \text{span}\{r_0, \dots, r_k\}$$

erfüllen für $0 \leq k \leq n-1$ mit $r_k \neq 0$.

Damit das ganze Verfahren effizient wird, benötigen wir noch folgende Eigenschaft

$$Ad_k \in U_{k+1} := \text{span}\{d_1, \dots, d_{k+1}\} = \text{span}\{r_0, \dots, r_k\},$$

wenn $r_k \neq 0$ gilt. Denn daraus ergibt sich $\langle r_k, Ad_j \rangle = 0$ für $1 \leq j \leq k-1$ und (4.15) verkürzt sich zu

$$d_{k+1} = r_k - \underbrace{\frac{\langle r_k, Ad_k \rangle}{\langle d_k, Ad_k \rangle}}_{=: \beta_{k+1}} d_k. \quad (4.16)$$

(Beweis: Aus der Definition von $r_k = r_{k-1} - \alpha_k Ad_k$ folgt sofort $Ad_k \in \{r_0, \dots, r_k\}$, da $r_{k-1} \in U_{k-1}$ und $\alpha_k \neq 0$ gilt. Die Gleichheit von U_k und $\{r_0, \dots, r_k\}$ ergibt sich aus Aufgabe 4.6.5.)

Für den Algorithmus schreiben wir nur noch α_k, β_k um: $\alpha_k = \frac{r_{k-1}^T d_k}{d_k^T Ad_k}$ und

$$\begin{aligned} r_{k-1}^T d_k &= r_{k-1}^T r_{k-1} - \beta_k \underbrace{r_{k-1}^T d_{k-1}}_{=0} \quad \text{also} \\ \alpha_k &= \frac{r_{k-1}^T r_{k-1}}{d_k^T Ad_k} \end{aligned} \quad (4.17)$$

und wegen $\alpha_k r_k^T Ad_k = (r_{k-1} - r_k)^T r_k = -r_k^T r_k$

$$\beta_{k+1} = \frac{r_k^T Ad_k}{d_k^T Ad_k} = -\frac{r_k^T r_k}{\alpha_k d_k^T Ad_k} = -\frac{r_k^T r_k}{r_{k-1}^T r_{k-1}}. \quad (4.18)$$

Bemerkung 4.6.6 Die Ausdrücke (4.17), (4.18) haben sich als numerisch stabiler und Speicherplatz-optimal herausgestellt.

Algorithmus 4.6.1: Konjugiertes Gradientenverfahren (cg-Verfahren): $Ax = b$ **Input:** Initial guess x_0

$$r_0 := b - Ax_0$$

$$\rho_0 := \langle r_0, r_0 \rangle$$

$$d_1 := r_0$$

Iteration: $k = 1, 2, \dots$ as long as $k \leq n$ and $r_k \neq 0$

$$a_k := Ad_k$$

$$\alpha_k := \rho_{k-1} / \langle d_k, a_k \rangle$$

$$x_k := x_{k-1} + \alpha_k d_k$$

$$r_k := r_{k-1} - \alpha_k a_k$$

$$\rho_k := \langle r_k, r_k \rangle$$

$$d_{k+1} := r_k + \frac{\rho_k}{\rho_{k-1}} d_k$$

Satz 4.6.7 Es sei A positiv definit. Für das cg-Verfahren gilt folgende Abschätzung

$$\|x^* - x_k\|_A \leq 2 \left(\frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x^* - x_0\|_A.$$

Beweis. Der Beweis gliedert sich in 4 Schritte. Es sei $e_k := x^* - x_k$.(1) Zuerst zeige man induktiv, dass Polynome $p_k \in \mathbb{P}_k$, $k = 0, \dots, n-1$, existieren mit

$$e_k = p_k(A)e_0, \quad p_k(0) = 1. \quad (4.19)$$

 $k = 0$: klar! $k-1 \Rightarrow k$: Aus der Definition (4.13) folgt

$$\alpha_k d_k = x_k - x_{k-1} = x^* - x_{k-1} - (x^* - x_k) = e_{k-1} - e_k, \text{ bzw. } e_k = e_{k-1} - \alpha_k d_k.$$

Mit (4.16) ergibt sich nun

$$\begin{aligned} d_k &= r_{k-1} - \beta_k d_{k-1} = b - Ax_{k-1} - \beta_k \frac{1}{\alpha_{k-1}} (e_{k-2} - e_{k-1}) \\ &= A(x^* - x_{k-1}) + \frac{\beta_k}{\alpha_{k-1}} (e_{k-1} - e_{k-2}) \\ &= \left(\frac{\beta_k}{\alpha_{k-1}} I + A \right) e_{k-1} - \frac{\beta_k}{\alpha_{k-1}} e_{k-2} \\ &= \underbrace{\left\{ \left(\frac{\beta_k}{\alpha_{k-1}} I + A \right) p_{k-1}(A) - \frac{\beta_k}{\alpha_{k-1}} p_{k-2}(A) \right\}}_{=: \tilde{q}_k(A), \tilde{q}_k \in \mathbb{P}_k} e_0, \end{aligned}$$

also

$$e_k = \underbrace{\left[\left(\frac{\alpha_k \beta_k}{\alpha_{k-1}} I + A \right) p_{k-1}(A) - \frac{\alpha_k \beta_k}{\alpha_{k-1}} p_{k-2}(A) + p_{k-1}(A) \right]}_{=: (\alpha_k \tilde{q}_k(A) + p_{k-1}(A)) =: p_k(A)} e_0$$

- (2) Nun zeige man, dass für alle $q_k \in \mathbb{P}_k$ mit $q_k(0) = 1$ die Ungleichung $\|e_k\|_A \leq \|q_k(A)e_0\|_A$ gilt.

Zuerst halten wir $r_k = b - Ax_k = A(x^* - x_k) = Ae_k$ fest. Des weiteren gilt

$$\begin{aligned} \|e_k\|_A^2 &:= e_k^T Ae_k = r_k^T e_k \\ &= r_k^T \left(e_k - \sum_{j=0}^{k-1} \sigma_j r_j \right) \quad \text{für beliebige } \sigma_0, \dots, \sigma_{k-1} \in \mathbb{R} \quad (r_j^T r_k = \delta_{jk}) \\ &= r_k^T \left(e_k + \sum_{j=0}^{k-1} \sigma_j Ae_j \right) \\ &= (Ae_k)^T \left(p_k(A) + \sum_{j=0}^{k-1} \sigma_j Ap_j(A) \right) e_0. \end{aligned}$$

Sei $q_k \in \mathbb{P}_k$ mit $q_k(0) = 1$ beliebig. Da die $\{p_0, \dots, p_{k-1}\}$ linear unabhängig sind, folgt aus der letzten Umformung, dass es eindeutig bestimmte $\tilde{\sigma}_0, \dots, \tilde{\sigma}_{k-1}$ existieren mit

$$q_k(t) = p_k(t) + \sum_{j=0}^{k-1} \tilde{\sigma}_j t p_j(t).$$

Mit $\sigma_i = \tilde{\sigma}_i$ folgt also

$$e_k^T Ae_k = e_k^T A^{1/2} A^{1/2} q_k(A) e_0 \stackrel{CSU}{=} e_k^T Ae_k (q_k(A) e_0)^T A q_k(A) e_0$$

Wir haben somit gezeigt, dass

$$\|e_k\|_A \leq \|q_k(A)e_0\|_A \quad \text{für alle } q_k \in \mathbb{P}_k, q_k(0) = 1 \tag{4.20}$$

gilt.

- (3) Die Eigenwerte von A seien $0 < \lambda_{\min} = \lambda_1 \leq \dots \leq \lambda_n = \lambda_{\max}$. Es gelte $A = Q\Lambda Q^T$ mit $QQ^T = Q^TQ = I$ und $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$. Beachtet man $q_k(A) = q_k(Q\Lambda Q^T) = Qq_k(\Lambda)Q^T$, so schätzt man wie folgt ab

$$\begin{aligned} \|q_k(A)e_0\|_A^2 &= e_0^T Q q(\Lambda) Q^T Q \Lambda Q^T Q q(\Lambda) Q^T e_0 = e_0^T Q q(\Lambda \Lambda q(\Lambda) Q^T e_0 \tag{4.21} \\ &\leq \max_{\lambda \in \{\lambda_1, \dots, \lambda_n\}} \{q(\lambda)^2\} e_0^T Q \Lambda Q^T e_0 \leq \left(\max_{\lambda \in [\lambda_{\min}, \lambda_{\max}]} |q(\lambda)| \right)^2 e_0^T Ae_0, \end{aligned}$$

d.h. es gilt

$$\|q_k(A)e_0\|_A \leq \|e_0\|_A \cdot \max_{\lambda \in [\lambda_1, \lambda_n]} |q_k(\lambda)|.$$

- (4) Man sucht nun Polynome p_k , die die rechte Seite minimieren. Man kann zeigen, dass die *Tschebyscheff-Polynome* $T_k : \mathbb{C} \rightarrow \mathbb{C}$ ($k = 0, 1, \dots$) definiert als

$$T_k(x) := \frac{1}{2} \left[(x + \sqrt{x^2 - 1})^k + (x + \sqrt{x^2 - 1})^{-k} \right]$$

auf $[-1, 1]$ folgende Eigenschaften haben

$$T_k(1) = 1, \quad |T_k(x)| \leq 1 \quad \forall -1 \leq x \leq 1$$

und minimal bzgl. $\|\cdot\|_\infty$ unter allen Polynomen $p \in \mathbb{P}_k$ mit $p(1) = 1$ sind (siehe Aufgabe 4.6.9). Gesucht ist nun eine Transformation, die $[\lambda_{\min}, \lambda_{\max}]$ auf $[-1, 1]$ abbildet. Somit

ist $\tilde{T}_k(z) := T_k\left(\frac{\lambda_{\max} + \lambda_{\min} - 2z}{\lambda_{\max} - \lambda_{\min}}\right)$ minimal bzgl. $\|\cdot\|_\infty$ auf $[\lambda_{\min}, \lambda_{\max}]$ unter allen Polynomen aus $p \in \mathbb{P}_k$ mit $p(\lambda_{\max}) = 1$. Man wählt also

$$q_k(z) := T_k\left(\frac{\lambda_{\max} + \lambda_{\min} - 2z}{\lambda_{\min} - \lambda_{\max}}\right) / T_k\left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\min} - \lambda_{\max}}\right)$$

auf $[0, \lambda_{\max}]$, $q_k(0) = 1$
Daraus folgt

$$\min_{q_k(0)=1} \max_{\lambda \in [\lambda_1, \lambda_n]} |q_k(\lambda)| \leq \frac{1}{T_k\left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}}\right)}$$

und

$$T_k\left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\min} - \lambda_{\max}}\right) = T_k\left(\frac{\lambda_{\max} + \lambda_{\min}}{\lambda_{\max} - \lambda_{\min}}\right) = T_k\left(\frac{\kappa + 1}{\kappa - 1}\right) \geq \frac{1}{2}(z + \sqrt{z^2 - 1})^k$$

mit $z = \frac{\kappa+1}{\kappa-1}$, also

$$\begin{aligned} z + \sqrt{z^2 - 1} &= \frac{\kappa + 1}{\kappa - 1} + \sqrt{\frac{\kappa^2 + 2\kappa + 1 - \kappa^2 + 2\kappa - 1}{(\kappa - 1)^2}} \\ &= \frac{1}{\kappa - 1}(\kappa + 1 + 2\sqrt{\kappa}) = \frac{(\sqrt{\kappa} + 1)^2}{(\sqrt{\kappa} + 1)(\sqrt{\kappa} - 1)} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1}. \end{aligned}$$

□

Bemerkung 4.6.8 Im letzten Beweis wurde unter (3) (siehe (4.21)) gezeigt, dass

$$\|q_k(A)e_0\|_A \leq \max_{\lambda \in \{\lambda_1, \dots, \lambda_n\}} |q(\lambda)| \|e_0\|_A$$

gilt. Daraus folgt sofort, dass das cg-Verfahren nach m -Schritten terminiert, wenn das Spektrum von A nur m verschiedene Eigenwerte besitzt.

Aufgabe 4.6.9 (Tschebyscheff-Polynome) Die Tschebyscheff-Polynome T_n sind orthogonal bezüglich des Skalarprodukts

$$(f, g)_\omega := \int_{-1}^1 \frac{f(x)g(x)}{\sqrt{1-x^2}} dx$$

und werden standardisiert durch $T_n(1) = 1$.

Man zeige, dass die Tschebyscheff-Polynome folgende Eigenschaften besitzen:

- (i) Sie haben stets ganzzahlige Koeffizienten
- (ii) Der höchste Koeffizient von T_n ist $a_n = 2^{n-1}$
- (iii) T_n ist stets eine gerade Funktion, falls n gerade und eine ungerade, falls n ungerade ist
- (iv) $T_n(1) = 1, T_n(-1) = (-1)^n$ (vgl. hierzu auch Beispiel 3)
- (v) $|T_n(x)| \leq 1$ für $x \in [-1, 1]$
- (vi) Die Nullstellen von $T_n(x)$ sind

$$x_k := \cos\left(\frac{2k-1}{2n}\pi\right), \quad (k = 1, \dots, n)$$

(vii)

$$T_k(x) = \begin{cases} \cos(k \cdot \arccos(x)), & -1 \leq x \leq 1; \\ \cosh(k \cdot \operatorname{Arccosh}(x)), & x \geq 1; \\ (-1)^k \cosh(k \cdot \operatorname{Arccosh}(-x)), & x \leq -1. \end{cases}$$

(viii) Die *Tschebyscheff*-Polynome besitzen die globale Darstellung

$$T_k(x) = \frac{1}{2} \left((x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right), \quad \text{wobei } x \in \mathbb{R}$$

(ix) $|T_n(x)|$ nimmt seinen maximalen Wert im Intervall $[-1, 1]$ an den sogenannten *Tschebyscheff*-Abszissen $\bar{x}_k = \cos\left(\frac{k\pi}{n}\right)$ für $k = 0, \dots, n$ an, d.h.

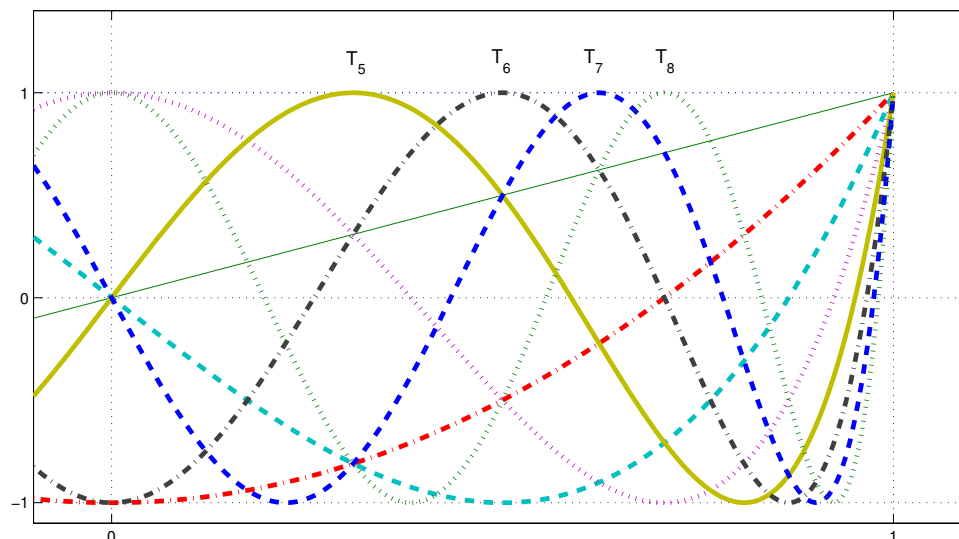
$$|T_n(x)| = 1 \Leftrightarrow x = \bar{x}_k = \cos\left(\frac{k\pi}{n}\right) \quad \text{mit } k = 0, \dots, n.$$

(x) Für $x \in \mathbb{R}$ erfüllen die *Tschebyscheff*-Polynome die folgende Drei-Term-Rekursion

$$T_0(x) = 1, \quad T_1(x) = x, \quad T_k(x) = 2xT_{k-1}(x) - T_{k-2}(x), \quad k \geq 2$$

Ergänzend seien für $n = 0, \dots, 5$ die T_n explizit angeben und diese grafisch dargestellt:

$$\begin{array}{ll} T_0 = 1 & T_3 = 4x^3 - 3x \\ T_1 = x & T_4 = 8x^4 - 8x^2 + 1 \\ T_2 = 2x^2 - 1 & T_5 = 16x^5 - 20x^3 + 5x \end{array}$$



4.7 PRÄKONDITIONIERUNG, DAS PCG-VERFAHREN

Die Abschätzung der Konvergenzgeschwindigkeit des cg-Verfahrens hängt gemäß Satz 4.6.7 monoton von der Kondition κ der Matrix A ab. Ziel dieses Abschnittes ist es daher, das Problem

$Ax = b$ so zu transformieren, dass die entstehende Matrix möglichst „gut konditioniert“ (d.h. κ möglichst klein) ist.

Idee: Betrachte anstatt

$$Ax = b \text{ mit } A \text{ s.p.d.} \quad (4.22)$$

$$\bar{A}\bar{x} = b \text{ mit } \bar{x} = B^{-1}x \text{ und } \bar{A} = AB \quad (4.23)$$

mit einer invertierbaren Matrix $B \in \mathbb{R}^{n \times n}$ und wende hierauf ein iteratives Lösungsverfahren an. Um die Symmetrie des Problems zu erhalten, betrachtet man im Falle, dass B ebenfalls s.p.d. ist, anstatt des Euklidischen Skalarproduktes $\langle \cdot, \cdot \rangle$ das von der Matrix B iduzierte Skalarprodukt $\langle \cdot, \cdot \rangle_B := \langle \cdot, B \cdot \rangle$. Aufgrund der folgenden Gleichung

$$\langle x, \bar{A}y \rangle_B = \langle x, AB y \rangle_B = \langle x, BAB y \rangle = \langle ABx, By \rangle = \langle ABx, y \rangle_B = \langle \bar{A}x, y \rangle_B$$

ist $\bar{A} = AB$ bezüglich $\langle \cdot, \cdot \rangle_B$ selbstadjungiert. Man kann daher zur Lösung von $\bar{A}\bar{x} = b$ das cg-Verfahren anwenden, vorausgesetzt man berücksichtigt, dass hierbei an die Stelle des Euklidischen Skalarproduktes die Form $\langle \cdot, \cdot \rangle_B$ rückt. Da man jedoch nicht an dem Vektor \bar{x} , sondern an der eigentlichen Lösung $x = B\bar{x}$ von $Ax = b$ interessiert ist, ersetzt man die so bestimmten Iterierten \bar{x}_k durch $x_k = B\bar{x}_k$. Hieraus ergibt sich unmittelbar der Algorithmus 4.7.1.

Algorithmus 4.7.1: Präkonditioniertes Konjugiertes Gradientenverfahren (pcg-Verfahren)

Input: Initial guess x_0

$$r_0 := b - Ax_0$$

$$\rho_0 := \langle r_0, Br_0 \rangle$$

$$d_1 := Br_0$$

Iteration: $k = 1, 2, \dots$ as long as $k \leq n$ and $r_k \neq 0$

$$a_k := Ad_k$$

$$\alpha_k := \rho_{k-1} / \langle d_k, a_k \rangle$$

$$x_k := x_{k-1} + \alpha_k d_k$$

$$r_k := r_{k-1} - \alpha_k a_k$$

$$\rho_k := \langle r_k, Br_k \rangle$$

$$d_{k+1} := Br_k + \frac{\rho_k}{\rho_{k-1}} d_k$$

Pro Iterationsschritt benötigt dieser Algorithmus gegenüber dem cg-Verfahren nur eine Multiplikation mit der Matrix B mehr. Doch dieser zusätzliche Aufwand rentiert sich, sofern sich die Kondition $\kappa(\bar{A})$ der Matrix $\bar{A} = AB$ gegenüber der Konditionszahl $\kappa(A)$ des nicht präkonditionierten Systems „verbessert“, d.h. sie nimmt ab.

Beispiel 4.7.1 Eine sehr einfache, aber häufig schon wirkungsvolle Vorkonditionierung einer s.p.d. Matrix A mit nichtverschwindenden Diagonalelementen liefert die Wahl $B := D^{-1}$, also der Inversen der Diagonale von A . Man spricht in diesem Fall von *diagonaler Präkonditionierung*.

Bemerkung 4.7.2 *Es ist klar, dass die Wahl des Prädiktionierers B elementar von den Eigenschaften der Matrix A abhängt. Die Wahl eines geeigneten Prädiktionierers ist daher oft nicht leicht. Hat man aber einen gefunden, so erhöht sich die Konvergenzgeschwindigkeit beträchtlich, was besonders für große lineare Gleichungssysteme von Bedeutung ist.*

5 AUSGLEICHSPROBLEME

Motivation Gesucht ist eine Gerade $g(x) = mx + b$, deren y -Werte den kleinsten Quadratsummenabstand zu den vorgegebenen Daten $(x_i, f(x_i))$ haben, die sogenannte *Ausgleichsgerade*. Die geometrische Veranschaulichung dazu ist in der nachfolgenden Abbildung dargestellt.

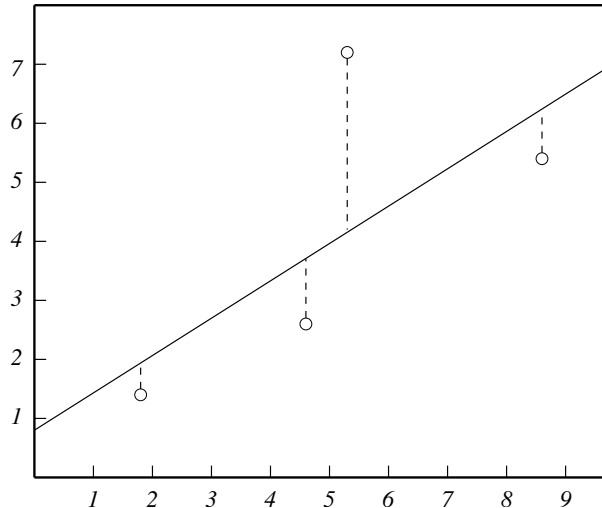


Abb. 5.1: Die Ausgleichsgerade, die die Quadratsumme der Abstände minimiert

5.1 DIE METHODE DER KLEINSTEN QUADRATE

Gemäß obiger Motivation gilt es zu gegebenen Punktepaaren $(x_1, f_1), (x_2, f_2), \dots, (x_n, f_n)$ die in der Geradengleichung $g(x) = mx + b$ freien Parameter m und b so zu bestimmen, dass gilt

$$F(m, b) := \sum_{i=1}^n (g(x_i) - f_i)^2 \rightarrow \min,$$

oder in alternativer Formulierung

$$\min_{(b,m) \in \mathbb{R}^2} \|A \begin{pmatrix} b \\ m \end{pmatrix} - d\|_2^2,$$

wobei wir nachfolgende Definitionen verwenden

$$A = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad d = \begin{pmatrix} f_1 \\ \vdots \\ f_n \end{pmatrix}.$$

Hierzu bilden wir alle partiellen Ableitungen bzgl. m und b , woraus sich dann die kritischen Punkte der Funktion $F(m, b)$ wie folgt berechnen

$$\nabla F(m, b) = \left(\frac{\partial}{\partial b} F, \frac{\partial}{\partial m} F \right) = \vec{0}.$$

Im vorliegenden Fall erhalten wir damit

$$\begin{aligned}\frac{\partial}{\partial b} F(m, b) &= 2 \sum_{i=1}^n (mx_i + b - f_i) \cdot 1 \stackrel{!}{=} 0 \\ \frac{\partial}{\partial m} F(m, b) &= 2 \sum_{i=1}^n (mx_i + b - f_i) \cdot x_i \stackrel{!}{=} 0.\end{aligned}$$

Diese beiden Bestimmungsgleichungen können auch wie folgt geschrieben werden

$$\begin{aligned}1^T \left(A \begin{pmatrix} b \\ m \end{pmatrix} - d \right) &= 0 \\ x^T \left(A \begin{pmatrix} b \\ m \end{pmatrix} - d \right) &= 0,\end{aligned}$$

wobei $1 = (1, \dots, 1)^T \in \mathbb{R}^n$ und $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ seien.

Es gilt also

$$\begin{aligned}\nabla F(m, b) = 0 &\Rightarrow 1^T \left(A \begin{pmatrix} b \\ m \end{pmatrix} - d \right) = 0 \\ &\quad x^T \left(A \begin{pmatrix} b \\ m \end{pmatrix} - d \right) = 0.\end{aligned}$$

Zusammengefasst folgt für potentielle Minimalstellen

$$\begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}^T \left(A \begin{pmatrix} b \\ m \end{pmatrix} - d \right) = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \Leftrightarrow \boxed{A^T A \begin{pmatrix} b \\ m \end{pmatrix} = A^T d}$$

Bemerkung 5.1.1 Falls in einem linearen Gleichungssystem die Anzahl der Gleichungen die der Unbekannten übersteigen sollte (wie es im vorliegenden Fall bei $Ax = d$ ist), so nennen wir das System **überbestimmt** und können es in der Regel nicht exakt lösen. Deshalb betrachten wir sinnvollerweise das Ausgleichsproblem der **Gaußschen Normalgleichung**

$$A^T Ax = A^T d \tag{5.1}$$

und erhalten dafür genau dann eine eindeutige Lösung, falls A vollen Rang besitzt.

Bemerkung 5.1.2 Als einführendes Beispiel haben wir uns auf die Betrachtung von linear unabhängigen Basisvektoren des \mathbb{P}_1 in der Darstellung von g beschränkt, d.h. $1, x$. Als Nächstes werden wir zu allgemeineren Funktionen in der Funktionsvorschrift von g übergehen, wie etwa

$$g(x) = a \exp(x) + b \exp(-x) + c \log(x).$$

Damit erhalten wir folgendes Minimierungsproblem bei bekannten Daten (x_i, f_i) ($i = 1, \dots, n$)

$$\sum_{i=1}^n (g(x_i) - f_i)^2 \rightarrow \min,$$

welches wir völlig analog zu obigem Verfahren als Normalgleichung in der Form

$$A^T A \begin{pmatrix} a \\ b \\ c \end{pmatrix} = A^T d, \text{ mit } A = \begin{pmatrix} e^{x_1} & e^{-x_1} & \log x_1 \\ \vdots & \vdots & \vdots \\ e^{x_n} & e^{-x_n} & \log x_n \end{pmatrix} \text{ und } d = \begin{pmatrix} f_1 \\ \vdots \\ f_n \end{pmatrix}$$

interpretieren können. Je komplexer also die Funktionsvorschrift von g , desto aufwendiger wird auch das Lösen der Gaußschen Normalgleichung (5.1). Im Folgenden werden wir ein Verfahren bereitstellen, das eine numerisch stabile und effiziente Lösung des Ausgleichsproblems garantiert.

5.2 DIE QR-ZERLEGUNG

Definition 5.2.1 Eine quadratische Matrix $A \in \mathbb{R}^{n \times n}$ heißt orthogonal, falls gilt:

$$AA^T = A^T A = I.$$

Weiterhin bezeichne $\|x\|$ die vom Euklidischen Skalarprodukt im \mathbb{R}^n induzierte Norm in der Form

$$\|x\| := \sqrt{\sum_{i=1}^n x_i^2} = \sqrt{x^T x}.$$

Darüber hinaus halten wir noch fest, dass orthogonale Abbildungen stets längen- und winkeltreu sind, dass also gilt

$$\|Qx\|^2 = (Qx)^T Qx = x^T Q^T Qx = x^T x = \|x\|^2, \text{ d.h. } \|Qx\| = \|x\|.$$

Motivation Zur Lösung des Minimierungsproblems $\min_{x \in \mathbb{R}^n} \|Ax - b\|^2$ setzen wir voraus, dass zu $A \in \mathbb{R}^{m \times n}$ ($\text{rg } A = n < m$) stets eine Zerlegung $A = QR$ existiert, wobei $Q \in \mathbb{R}^{m \times m}$ orthogonal ist und R folgende Gestalt besitzt

$$R = \begin{pmatrix} * & * & * \\ & * & * \\ & & * \\ 0 & & \end{pmatrix} = \begin{pmatrix} \widehat{R} \\ \hline 0 \end{pmatrix},$$

\widehat{R} stellt dabei eine reguläre obere Dreiecksmatrix dar. Damit können wir unser bisheriges Minimierungsproblem wie folgt modifizieren (beachte die Faktorisierung $A = QR$):

$$\|Ax - b\|^2 = \min \Leftrightarrow \|Q^T(Ax - b)\|^2 = \min \Leftrightarrow \|Rx - Q^T b\|^2 = \min,$$

hierbei weisen wir vor allen Dingen auf die Gestalt von R hin, als

$$Rx = \begin{pmatrix} * & * & * \\ & * & * \\ & & * \\ 0 & & \end{pmatrix} \cdot x = \begin{pmatrix} b_1 \\ \hline b_2 \end{pmatrix} = Q^T b$$

Dieser Faktorisierungsansatz führt uns nun zu folgenden Fragen:

- Zu welchen Matrizen $A \in \mathbb{R}^{m \times n}$ existiert eine derartige QR-Zerlegung?
- Wie bestimmen wir eine solche Zerlegung möglichst effizient?

Satz 5.2.2 Zu jeder Matrix $A \in \mathbb{R}^{m \times n}$ mit Maximalrang $n < m$ existiert eine orthogonale Matrix $Q \in \mathbb{R}^{m \times m}$ derart, dass

$$A = QR \text{ mit } R = \begin{pmatrix} \widehat{R} \\ 0 \end{pmatrix}, \widehat{R} \in \mathbb{R}^{n \times n}$$

gilt, wobei \widehat{R} eine reguläre obere Dreiecksmatrix darstellt.

Der soeben erbrachte Existenzbeweis ist rein konstruktiv und liefert damit eine erste Möglichkeit eine Zerlegung QR zu bestimmen. Wir unterscheiden die folgenden beiden Vorgehensweisen:

- Givens¹-Rotation
- Householder²-Spiegelung

5.3 Givens-ROTATIONEN

Es bezeichne im Folgenden stets $c = \cos \varphi$ und $s = \sin \varphi$. Die Grundidee der Givens-Rotation besteht darin durch die Multiplikation einer geeigneten orthogonalen (Dreh-)Matrix einen Vektor $x \in \mathbb{R}^n$ so zu drehen, dass möglichst viele seiner Komponenten verschwinden.

Im Falle $(a, b) \in \mathbb{R}^2$ werden wir also $c, s \in \mathbb{R}$ so bestimmen, dass gilt

$$\begin{pmatrix} c & -s \\ s & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}.$$

Setzen wir zuerst $a, b \neq 0$ voraus. Wir erhalten damit folgende Gleichung

$$as = -bc \tag{5.3}$$

und mit $s^2 + c^2 = 1$ gilt

$$r^2 = (ac - bs)^2 = a^2c^2 - abcs - abcs + b^2s^2 \stackrel{(5.3)}{=} a^2c^2 + a^2s^2 + b^2c^2 + b^2s^2 = a^2 + b^2.$$

Weiter ergibt sich für $a \neq 0$

$$\begin{aligned} ac - bs = r &\Rightarrow c = \frac{1}{a}(r + bs) \Rightarrow as + \frac{b}{a}(r + bs) = 0 \\ &\Leftrightarrow a^2s + br + b^2s = 0 \\ &\Leftrightarrow br = -(a^2 + b^2)s \\ &\Leftrightarrow s = \frac{-br}{a^2 + b^2} \stackrel{(r^2=a^2+b^2)}{=} \frac{-b}{r}. \end{aligned}$$

Und damit erhalten wir dann auch für $b \neq 0$

$$-bc = as \Rightarrow c = -\frac{a}{b}s = \frac{a}{r}.$$

Offen ist noch die Wahl des Vorzeichens des Vorzeichens $r = \pm\sqrt{a^2 + b^2}$. Setzen wir $r = \text{sign}(a)\sqrt{a^2 + b^2}$, so folgt $c > 0$, falls $a \neq 0$, und für $a = 0$ setzen wir $c = 0$ und $s = 1$. Somit sind auch s und c für die anderen Spezialfälle $b = 0$ und $a = b = 0$ wohldefiniert. Damit gilt (obwohl der Winkel bisher gar nicht explizit genannt wurde) $\varphi \in (-\frac{\pi}{2}, \frac{\pi}{2}]$ und somit lässt sich $\cos \varphi$ eindeutig aus $\sin \varphi$ mit $\cos \varphi = \sqrt{1 - \sin^2 \varphi}$ bestimmen.

Seien p und q zwei Zeilenindizes von A . Zur Eliminierung des Matrixelements a_{qj} gehen wir daher wie folgt vor:

$$\begin{aligned} \text{Falls } a_{pj} \neq 0: \quad \cos \varphi &= \frac{|a_{pj}|}{\sqrt{a_{pj}^2 + a_{qj}^2}}, \\ \sin \varphi &= \frac{-\text{sign}(a_{pj}) a_{qj}}{\sqrt{a_{pj}^2 + a_{qj}^2}}, \\ a_{pj} = 0: \quad \cos \varphi &= 0, \quad \sin \varphi = 1. \end{aligned}$$

¹Givens, W.

²Householder, A.

Da sich $\cos \varphi$ aus $\sin \varphi$ bestimmen lässt, ermöglicht das Verfahren eine effiziente Speicherung, denn es genügt $\sin \varphi$ an den entsprechenden freiwerdenden Stellen in A abzulegen:

$$\begin{aligned}
 A &= \begin{pmatrix} a_{11}^{(0)} & \cdots & a_{1n}^{(0)} \\ \vdots & & \vdots \\ a_{m1}^{(0)} & \cdots & a_{mn}^{(0)} \end{pmatrix} \rightsquigarrow \begin{pmatrix} a_{11}^{(1)} & \cdots & \cdots & a_{1n}^{(1)} \\ \sin \varphi_{21} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ a_{31}^{(0)} & a_{32}^{(0)} & \cdots & a_{3n}^{(0)} \\ a_{41}^{(0)} & a_{42}^{(0)} & \cdots & a_{4n}^{(0)} \\ \vdots & & & \vdots \\ a_{m1}^{(0)} & \cdots & \cdots & a_{mn}^{(0)} \end{pmatrix} \rightsquigarrow \begin{pmatrix} a_{11}^{(2)} & \cdots & \cdots & a_{1n}^{(2)} \\ \sin \varphi_{21} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \sin \varphi_{31} & a_{32}^{(2)} & \cdots & a_{3n}^{(2)} \\ a_{41}^{(0)} & a_{42}^{(0)} & \cdots & a_{4n}^{(0)} \\ \vdots & & & \vdots \\ a_{m1}^{(0)} & \cdots & \cdots & a_{mn}^{(0)} \end{pmatrix} \\
 &\rightsquigarrow \begin{pmatrix} a_{11}^{(3)} & \cdots & \cdots & a_{1n}^{(3)} \\ \sin \varphi_{21} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \sin \varphi_{31} & a_{32}^{(2)} & \cdots & a_{3n}^{(2)} \\ \sin \varphi_{41} & a_{42}^{(3)} & \cdots & a_{4n}^{(3)} \\ \vdots & & & \vdots \\ a_{m1}^{(0)} & \cdots & \cdots & a_{mn}^{(0)} \end{pmatrix} \rightsquigarrow \begin{pmatrix} a_{11}^{(m-1)} & \cdots & \cdots & a_{1n}^{(m-1)} \\ \sin \varphi_{21} & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \sin \varphi_{31} & a_{32}^{(2)} & \cdots & a_{3n}^{(2)} \\ \sin \varphi_{41} & a_{42}^{(3)} & \cdots & a_{4n}^{(3)} \\ \vdots & & & \vdots \\ \sin \varphi_{m1} & a_{m2}^{(m-1)} & \cdots & a_{mn}^{(m-1)} \end{pmatrix} \rightsquigarrow \begin{pmatrix} \ddots & & & \widehat{R} \\ & \ddots & & \\ \text{„Q“} & & \ddots & \\ & & & \ddots \end{pmatrix}.
 \end{aligned}$$

Die QR -Zerlegung liefert also QR in kompakter Form, gespeichert in A . Der Aufwand dieser Vorgehensweise liegt bei ungefähr $\mathcal{O}(n^3)$ Operationen.

Bemerkung 5.3.1 Die Multiplikation eines gegebenen Vektors $x \in \mathbb{R}^n$ mit der Drehmatrix $\mathcal{U}(p, q; \varphi)$ ist äquivalent zur Drehung von x um einen Winkel φ entgegen dem Uhrzeigersinn in der Koordinatenebene. Vergleiche hierzu die nachfolgende Abbildung, die diesen geometrischen Sachverhalt verdeutlicht:

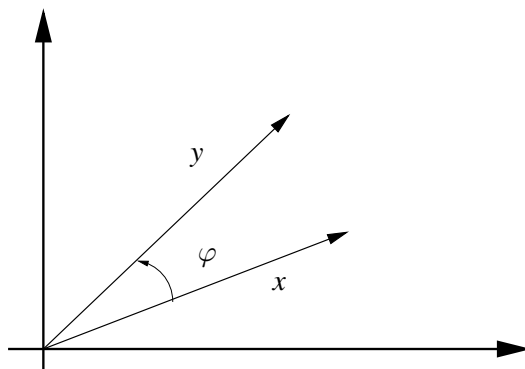


Abb. 5.2: Drehung um einen Winkel φ in der Ebene

MATLAB-Funktion: GivensRotMat.m

```

1 function rot = GivensRotMat(ap, aq)
2 if ap == 0
3     s = 1; c = 0;
4 else
5     dist2 = sqrt(ap^2+aq^2);
6     c = abs(ap)/dist2;
7     s = -sign(ap)*aq/dist2;
8 end
9 rot = [c, s; -s, c];

```

MATLAB-Funktion: Givens.m

```

1 function [Q,A] = Givens(A)
2 Q = eye(size(A,1));
3 for j=1:size(A,2)
4     for i=j+1:size(A,1)
5         rot = GivensRotMat(A(j,j),A(i,j));
6         A([j,i],j:end) = rot' * A([j,i],j:end);
7         Q(:, [j,i]) = Q(:, [j,i]) * rot;
8     end
9     A(j+1:end, j) = 0;
10 end

```

MATLAB-Beispiel:

Alternativ zu der Matlab-Funktion `qr` können wir auch die Routine `Givens` verwenden, die eine QR -Zerlegung, wie oben diskutiert, berechnet.

```

>> A = [1,2;3,4;5,6];
> [Q,R] = Givens(A);
>> Q*R
ans =
    1.0000    2.0000
    3.0000    4.0000
    5.0000    6.0000

```

5.4 UPDATE EINER QR -ZERLEGUNG

Bei verschiedenen numerischen Verfahren müssen wiederholt QR -Zerlegungen bestimmt werden, wobei sich die Matrizen nur wenig voneinander unterscheiden. Dies ist z.B. der Fall beim Broyden-Verfahren (siehe Numerik II, nichtlineare Gleichungen). Sei die Zerlegung $A_k = Q_k R_k$ gegeben und die QR -Zerlegung von $A_{k+1} = A_k + st^T$ zu bestimmen. Wir werden zeigen, dass

anstatt jeweils $\mathcal{O}(n^3)$ Operationen für eine vollständige QR -Zerlegung zu verwenden, für jeden einzelnen Updateschritt $\mathcal{O}(n^2)$ Operationen genügen.

Wir benötigen noch einige Hilfsmittel, um diese Aussage zu zeigen.

Definition 5.4.1 (Hessenberg-Matrix) Eine Matrix $A \in \mathbb{R}^{n \times n}$ der Form

$$A = \begin{pmatrix} * & \cdots & \cdots & \cdots & * \\ * & & & & \vdots \\ 0 & \cdot & & & \vdots \\ \vdots & \cdot & \cdot & & \vdots \\ 0 & \cdots & 0 & * & * \end{pmatrix},$$

d.h. $a_{ij} = 0$ für $i - j \geq 2$, wird obere **Hessenberg-Matrizen** genannt. Gilt $a_{ij} = 0$ für $j - i \geq 2$, so hat die Matrix A untere Hessenberg-Gestalt.

Lemma 5.4.2 Es seien $u, v \in \mathbb{R}^n$ und $R \in \mathbb{R}^{n \times n}$ eine obere Dreiecksmatrix. Dann benötigt man $n - 2$ Givens-Rotationen, um $R + uv^T$ auf obere Hessenberg-Gestalt zu transformieren.

Beweis. Multipliziert man R sukzessive von links mit Drehmatrizen $\mathcal{U}^T(p, q)$ (siehe (5.2)) zu Indexpaaren $(n, n - 1), \dots, (4, 3), (3, 2)$ so hat QR obere Hessenberg-Gestalt, wobei $Q := \mathcal{U}_{32} \cdots \mathcal{U}_{n, n-1}$ ist. Dabei lassen sich die Drehmatrizen so wählen, dass $Qu = (*, *, 0, \dots, 0)^T$ gilt. Beachtet man nun, dass $(Qu)v^T$ und somit auch $Q(R + uv^T)$ obere Hessenberg-Form hat, so ist die Aussage bewiesen. \square

Bemerkung 5.4.3 Beachtet man, dass sich jede Givens-Rotation angewandt auf $A \in \mathbb{R}^{n \times n}$ durch $\mathcal{O}(n)$ Operationen realisieren lässt, so genügen $\mathcal{O}(n^2)$ Operationen um eine orthogonale Matrix \tilde{Q} und eine Hessenberg-Matrix \tilde{H} mit $\tilde{Q}\tilde{H} = R + uv^T$ zu bestimmen.

Bemerkung 5.4.4 Es genügen $n - 1$ Givens-Rotation um eine Hessenberg-Matrix $A \in \mathbb{R}^{n \times n}$ auf Dreiecksgestalt zu transformieren.

Kommen wir nun zu unserer ursprünglichen Fragestellung zurück, d.h. der effizienten Update-Berechnung. Der Ansatz

$$A_{k+1} = Q_{k+1}R_{k+1} = Q_k(R_k + uv^T) = A_k + Q_k uv^T = A_k + st^T$$

liefert $v = t$ und $s = Q_k u$ bzw. $u = Q_k^T s$. Das Update besteht nun aus zwei Schritten.

1. Mittels $n - 2$ Givens-Rotationen wird $R_k + uv^T$ auf eine obere Hessenbergform gebracht, d.h. wir bestimmen eine orthogonale Matrix Q'_{k+1} und eine Hessenberg-Matrix H_{k+1} mit $Q'_{k+1}H_{k+1} = R_k + uv^T = R_k + (Q_k^T s)t^T$.
2. Mittels $n - 1$ Givens-Rotationen wird diese obere Hessenbergform zur neuen Matrix R_{k+1} transformiert, d.h. wir bestimmen eine orthogonale Matrix Q''_{k+1} und eine obere Dreiecksmatrix R_{k+1} mit $Q''_{k+1}R_{k+1} = H_{k+1}$. Die Matrix Q_{k+1} ergibt sich dann als Produkt $Q_k \cdot Q'_{k+1} \cdot Q''_{k+1}$.

MATLAB-Funktion: QR2Hessenberg.m

```

1 function [Q,R] = QR2Hessenberg(Q,R,s,t)
2 % computes QH-decomposition of Q(R+s*t), H upper Hessenberg form
3 for j=size(R,1):-1:3
4     rot = GivensRotMat(s(j),s(j-1));
5     s([j,j-1]) = rot * s([j,j-1]);
6     Q(:,[j,j-1]) = Q(:,[j,j-1]) * rot';
7     R([j,j-1],:) = rot * R([j,j-1],:);
8 end
9 R(1:2,:) = R(1:2,:)+ s(1:2) * t';

```

MATLAB-Funktion: GivensHessenberg.m

```

1 function [Q,H] = GivensHessenberg(Q,H)
2 % computes QR-decomposition of QH, H Hessenberg matrix
3 for j=1:min(size(H,1)-1,size(H,2))
4     rot = GivensRotMat(H(j+1,1),H(j,1));
5     H([j+1,j],j:end) = rot * H([j+1,j],j:end);
6     Q(:,[j+1,j]) = Q(:,[j+1,j]) * rot';
7     H(j+1,j) = 0;
8 end

```

MATLAB-Beispiel:

Ein einfaches Beispiel zeigt, dass man das $\mathcal{O}(n^2)$ -Verhalten des Updates auch im Vergleich zur optimierten Matlab-Routine `qr` sieht und nutzen kann. Verdoppelt man z.B. n auf $n = 4000$ im nebenstehenden Listing, so benötigt das Update nur $2^2 = 4$ -mal länger, die interne Funktion jedoch $2^3 = 8$ -mal. Testen Sie es einmal!

```

>> n = 2000;
>> A=rand(n); s=rand(n,1); t=rand(n,1);
>> An = A + s * t';
>> [Q,R] = qr(A);
>> tic
>> [Qtmp,H] = QR2Hessenberg(Q,R,Q'*s,t);
>> [Qn,Rn] = GivensHessenberg(Qtmp,H);
>> toc
Elapsed time is 1.078000 seconds.
>> tic, [Qn,Rn] = qr(An); toc
Elapsed time is 5.766000 seconds.

```

5.5 Householder-SPIEGELUNGEN

Definition 5.5.1 Sei $\omega \in \mathbb{R}^n$. Die $n \times n$ -Matrix

$$P = I - 2 \frac{\omega \omega^T}{\omega^T \omega} \quad (5.4)$$

wird als Householder-Transformation und ω als Householder-Vektor bezeichnet.

Satz 5.5.2 Eine Householder-Transformation $P = I - 2(\omega \omega^T)/(\omega^T \omega)$ ist symmetrisch und orthogonal, d.h. $P^{-1} = P^T$.

Beweis. Da $\omega \omega^T$ symmetrisch ist ($(\omega \omega^T)^T = \omega \omega^T$), folgt

$$P^T = \left(I - 2 \frac{\omega \omega^T}{\omega^T \omega} \right)^T = I - 2 \frac{\omega \omega^T}{\omega^T \omega} = P.$$

Des Weiteren ergibt sich

$$P P^T = \left(I - 2 \frac{\omega \omega^T}{\omega^T \omega} \right) \left(I - 2 \frac{\omega \omega^T}{\omega^T \omega} \right) = I - 2 \frac{\omega \omega^T}{\omega^T \omega} - 2 \frac{\omega \omega^T}{\omega^T \omega} + 4 \frac{\omega \omega^T \omega \omega^T}{(\omega^T \omega)^2}$$

und mit $\omega \omega^T \omega \omega^T = (\omega^T \omega) (\omega \omega^T)$ somit die Behauptung. □

Betrachten wir nun folgende Frage: Sei $x \in \mathbb{R}^n$. Wie müsste ein $\omega \in \mathbb{R}^n$ aussehen, so dass $Px = \alpha e_1$ gelte ($\alpha \in \mathbb{R}$, e_1 erster kanonischer Einheitsvektor)?

Mit $Px = x - 2 \frac{\omega \omega^T}{\omega^T \omega} x = x - \lambda \omega \stackrel{!}{=} \alpha e_1$ folgt $\omega \in \text{span}\{x - \alpha e_1\}$.

Wir setzen nun $\omega = x - \alpha e_1$ in $Px = \alpha e_1$ ein und erhalten

$$Px = x - 2 \frac{\omega^T x}{\omega^T \omega} \omega = \left(1 - \frac{2(x - \alpha e_1)^T x}{\|x - \alpha e_1\|^2} \right) x + \alpha \frac{2(x - \alpha e_1)^T x}{\|x - \alpha e_1\|^2} e_1 = \alpha e_1.$$

Damit in der letzten Gleichung der Faktor vor x verschwindet, muss

$$1 = \frac{2(x - \alpha e_1)^T x}{\|x - \alpha e_1\|^2} \Leftrightarrow (x - \alpha e_1)^T (x - \alpha e_1) = 2x^T x - 2\alpha x_1 \Leftrightarrow \alpha = \pm \sqrt{x^T x}$$

gelten.

Wie ist nun das Vorzeichen zu wählen, $\alpha = \pm \sqrt{x^T x}$? Die Wahl $\alpha = \|x\|_2$ hat die schöne Eigenschaft, dass Px ein positives Vielfaches von e_1 ist. Aber das Rezept ist gefährlich, wenn x annähernd ein positives Vielfaches von e_1 ist, da Auslöschungen auftreten können. Berechnet man ω_1 mittels $\omega_1 = x_1 - \|x\|_2$ treten für $x_1 \leq 0$ keine Auslöschungen auf und mit der Umformung

$$\omega_1 = x_1 - \|x\|_2 = \frac{x_1^2 - \|x\|_2^2}{x_1 + \|x\|_2} = \frac{-(x_2^2 + \dots + x_n^2)}{x_1 + \|x\|_2}$$

treten auch für $x_1 > 0$ keine Auslöschungen auf. Man beachte außerdem $(\omega_2, \dots, \omega_n) = (x_2, \dots, x_n)$.

Normiert man den Vektor ω so, dass $\omega_1 = 1$ gilt, d.h.

$$\omega := \frac{x - \alpha e_1}{x_1 - \alpha} \quad \text{mit } \alpha^2 = \|x\|^2 \quad (\omega^T e_1 = 1),$$

dann folgt

$$\omega^T \omega = \frac{(x - \alpha e_1)^T (x - \alpha e_1)}{(x_1 - \alpha)^2} = \frac{\|x\|^2 - 2\alpha x^T e_1 + \alpha^2}{(x_1 - \alpha)^2} = \frac{2\alpha(\alpha - x_1)}{(x_1 - \alpha)^2} = \frac{2\alpha}{\alpha - x_1}.$$

Des Weiteren bleibt festzuhalten, dass bei der Speicherung der Matrix $P = I - 2(\omega\omega^T)/(\omega^T\omega) \in \mathbb{R}^{n \times n}$ auch nur der Vektor $(\omega_2, \dots, \omega_n)^T$ zu berücksichtigen ist, der sich auf den $n - 1$ freigewordenen Einträgen speichern lässt.

Algorithmus 5.5.1: Berechnung Householder-Vektor:

Zu gegebenem $x \in \mathbb{R}^n$ berechnet die Funktion ein $\omega \in \mathbb{R}^n$ mit $\omega(1) = 1$ und $\beta \in \mathbb{R}$, so dass $P = I_n - \beta\omega\omega^T$ orthogonal ist und $Px = \|x\|_2 e_1$.

```

function:  $[\omega, \beta] = \text{housevector}(x)$ 
 $n = \text{length}(x)$ 
 $\sigma = x(2:n)^T x(2:n)$ 
 $\omega = \begin{bmatrix} 1 \\ x(2:n) \end{bmatrix}$ 
if  $\sigma = 0$ 
     $\beta = 0$ 
else
     $\mu = \sqrt{x(1)^2 + \sigma}$ 
    if  $x(1) \leq 0$ 
         $\omega(1) = x(1) - \mu$ 
    else
         $\omega(1) = -\sigma / (x(1) + \mu)$ 
    end
     $\beta = 2\omega(1)^2 / (\sigma + \omega(1)^2)$ 
     $\omega = \omega / \omega(1)$ 
end

```

MATLAB-Funktion: HouseholderVektor.m

```

1 function [v,beta] = HouseholderVector(x)
2 n = length(x);
3 if n>1
4     sigma = x(2:end)' * x(2:end);
5     if sigma==0
6         beta = 0;
7     else
8         mu = sqrt(x(1)^2+sigma);
9         if x(1)<=0
10            tmp = x(1) - mu;
11        else
12            tmp = -sigma / (x(1) + mu);
13        end
14        beta = 2*tmp^2/(sigma + tmp^2);
15        x(2:end) = x(2:end)/tmp;
16    end
17    v = [1;x(2:end)];
18 else
19     beta = 0;
20     v = 1;
21 end

```

Wir haben gerade konstruktiv gezeigt, dass sich mit Hilfe der *Householder*-Transformation eine Matrix $A \in \mathbb{R}^{m \times n}$ in eine Matrix der folgenden Form transformieren können:

$$H_1 A = \begin{pmatrix} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \end{pmatrix}$$

Gehen wir nun davon aus, dass wir nach einigen Schritten die Ausgangsmatrix A auf folgende Gestalt gebracht haben:

$$H_2 H_1 A = \begin{pmatrix} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & 0 & \boxplus & * & * \\ 0 & 0 & \boxplus & * & * \\ 0 & 0 & \boxplus & * & * \\ 0 & 0 & \boxplus & * & * \end{pmatrix}$$

Im nächsten Schritt soll nun die nächste Subspalte unterhalb der Diagonalen eliminiert werden. Es ist also eine Householder-Matrix zu bestimmen, so dass

$$\tilde{H}_3 \begin{pmatrix} \boxplus \\ \boxplus \\ \boxplus \\ \boxplus \end{pmatrix} = \begin{pmatrix} * \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

gilt. Definiert man H_3 als Blockdiagonalmatrix

$$H_3 := \begin{pmatrix} I & 0 \\ 0 & \tilde{H}_3 \end{pmatrix}$$

so erhalten wir

$$H_3 H_2 H_1 A = \begin{pmatrix} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & * & * \end{pmatrix}$$

Mit dieser Vorgehensweise erhalten wir folgende Faktorisierung:

$$R = H_{n-1} H_{n-2} \cdots H_1 A \Leftrightarrow A = (H_1 \cdots H_{n-1}) R \rightsquigarrow Q = H_1 \cdots H_{n-1}$$

Stellen wir nun dar, wie A überschrieben wird. Es sei

$$\omega^{(j)} = \underbrace{(0, \dots, 0)}_{j-1}, 1, \omega_{j+1}^{(j)}, \dots, \omega_m^{(j)} \Big)^T$$

der j -te Householder-Vektor, dann erhält man nach Durchführung der QR -Zerlegung

$$A = \begin{pmatrix} r_{11} & r_{12} & r_{13} & r_{14} & r_{15} \\ \omega_2^{(1)} & r_{22} & r_{23} & r_{24} & r_{25} \\ \omega_3^{(1)} & \omega_3^{(2)} & r_{33} & r_{34} & r_{35} \\ \omega_4^{(1)} & \omega_4^{(2)} & \omega_4^{(3)} & r_{44} & r_{45} \\ \omega_5^{(1)} & \omega_5^{(2)} & \omega_5^{(3)} & \omega_5^{(4)} & r_{55} \\ \omega_6^{(1)} & \omega_6^{(2)} & \omega_6^{(3)} & \omega_6^{(4)} & \omega_6^{(5)} \end{pmatrix}$$

Algorithmus 5.5.2: Berechnung Householder- QR :

Gegeben sei $A \in \mathbb{R}^{m \times n}$ mit $m \geq n$. Der folgende Algorithmus bestimmt die Householder-Matrizen H_1, \dots, H_n , so dass mit $Q = H_1 \cdots H_n$ die Matrix $R = Q^T A$ eine obere Dreiecksmatrix ist. Der obere Dreiecksteil von A wird mit R überschrieben und unterhalb der Diagonalen werden die nichttrivialen Komponenten der Householder-Vektoren gespeichert.

```

function:  $A = \text{housematrix}(A)$ 
for  $j = 1 : n$ 
   $[\omega, \beta] = \text{housevector}(A(j : m, j))$ 
   $A(j : m, j : n) = (I_{m-j+1} - \beta\omega\omega^T)A(j : m, j : n)$ 
  if  $j < m$ 
     $A(j + 1 : m, j) = \omega(2 : m - j + 1)$ 
  end
end
end

```

Bemerkung 5.5.3 Analog zur Givens-Rotation möchten wir auch an dieser Stelle kurz die geometrische Anschauung der Householder-Spiegelung anführen. Für einen gegebenen Vektor $x \in \mathbb{R}^n$ ist der Vektor $y = Px$ die Spiegelung von x an der Ebene $\text{span}\{\omega\}^\perp$. Hierzu nachfolgende Abbildung:

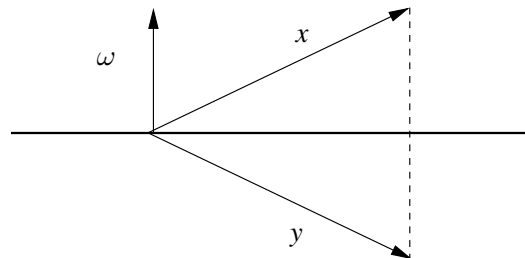


Abb. 5.3: Spiegelung an der zu ω orthogonalen Ebene

Aufgabe 5.5.4 Man zeige, dass die QR -Zerlegung mittels Givens-Rotation $\approx 50\%$ mehr Rechenoperationen benötigt als mittels Householder-Transformation.

MATLAB-Funktion: Householder.m

```

1 function [A,R] = Householder(A)
2 for j = 1:size(A,2)
3     [v,beta(j)] = HouseholderVector(A(j:end,j));
4     A(j:end,j:end) = A(j:end,j:end) - v * (beta(j) * v' * A(j:end,j:end)
5         );
6     if j < size(A,1)
7         A(j+1:end,j) = v(2:end);
8     end
9 end
10 if nargin == 2
11     R = A;
12     A = eye(size(A,1));
13     for j = size(R,2):-1:1
14         v = [1;R(j+1:end,j)];
15         A(j:end,j:end) = A(j:end,j:end) - v*(beta(j)*v'*A(j:end,j:end))
16         ;
17     end
18     R = triu(R);
19 end

```

MATLAB-Funktion: HouseholderMult.m

```

1 function y = prod_rx(A,x)
2 % extract R from A and multiply with x
3 y = triu(A) * x;
4
5 function y = prod_qx(A,x)
6 y = x;
7 for j = size(A,2)-1:-1:1
8     v = [1;A(j+1:end,j)];
9     beta = 2/(v'*v);
10    y(j:end) = y(j:end) - v*(beta*v'*y(j:end));
11 end
12
13 function y = prod_qtx(A,x)
14 y = x;
15 for j = 1:size(A,2)-1
16     v = [1;A(j+1:end,j)];
17     beta = 2/(v'*v);
18     y(j:end) = y(j:end) - v*(beta*v'*y(j:end));
19 end

```

MATLAB-Beispiel:

Mit der Funktion `Householder` wird zum einen die Matrix A kompakt überschrieben, d.h. der Rückgabewert benötigt nur den Speicherplatz der ursprünglichen Matrix; zum anderen kann man sich auch Q und R mit $A = QR$ ausgeben lassen. Die Berechnung von Rx , Qx und $Q^T x$, wenn die kompakte Form vorliegt, ist in den Routinen `prod_qx`, etc. realisiert.

```
>> A = rand(3);
>> x = rand(3,1);
>> [Q,R] = Householder(A);
>> A-Q*R
ans =
  1.0e-015 *
   -0.1110         0   -0.1110
   -0.0035   -0.1110         0
   -0.1110         0   0.2220
>> C = Householder(A);
>> Q'*x - prod_qtx(C,x)
ans =
  1.0e-015 *
   0.2220
   0.0555
   0.5551
```

A NORMEN

Normen erfüllen den gleichen Zweck auf Vektorräumen, den Beträge auf der reellen Achse erfüllen. Genauer gesagt, \mathbb{R}^n mit einer Norm auf \mathbb{R}^n liefert einen metrischen Raum. Daraus ergeben sich bekannte Begriffe wie Umgebung, offene Menge, Konvergenz und Stetigkeit für Vektoren und vektorwertige Funktionen.

Definition A.0.1 (Vektornorm) Eine Vektornorm auf \mathbb{R}^n ist eine Funktion $\|\cdot\| : \mathbb{R}^n \rightarrow \mathbb{R}$ mit den Eigenschaften

$$\begin{aligned} \|x\| &\geq 0 & x \in \mathbb{R}^n & \quad (\|x\| = 0, \text{ g.d.w. } x = 0), \\ \|x + y\| &\leq \|x\| + \|y\| & x, y \in \mathbb{R}^n, \\ \|\alpha x\| &= |\alpha| \|x\| & \alpha \in \mathbb{R}, x \in \mathbb{R}^n. \end{aligned} \tag{A.1}$$

Eine nützliche Klasse von Vektornormen sind die p -Normen, definiert durch

$$\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{\frac{1}{p}} \quad p \geq 1. \tag{A.2}$$

Von diesen sind besonders die 1, 2 und ∞ interessant:

$$\begin{aligned} \|x\|_1 &= |x_1| + \dots + |x_n| \\ \|x\|_2 &= (|x_1|^2 + \dots + |x_n|^2)^{\frac{1}{2}} = (x^T x)^{\frac{1}{2}} \\ \|x\|_\infty &= \max_{1 \leq j \leq n} |x_j| \end{aligned} \tag{A.3}$$

Ein **Einheitsvektor** bzgl. der Norm $\|\cdot\|$ ist ein Vektor x mit $\|x\| = 1$.

A.1 VEKTORNORM-EIGENSCHAFTEN

Ein klassisches Ergebnis bzgl. p -Normen ist die **Hölder-Ungleichung**

$$|x^T y| \leq \|x\|_p \|y\|_q \quad \text{mit } \frac{1}{p} + \frac{1}{q} = 1. \tag{A.4}$$

Spezialfall der Hölder-Ungleichung ist die **Cauchy-Schwarz-Ungleichung**

$$|x^T y| \leq \|x\|_2 \|y\|_2. \tag{A.5}$$

Allen Normen auf \mathbb{R}^n sind äquivalent, d.h. es seien $\|\cdot\|_\alpha$ und $\|\cdot\|_\beta$ Normen auf \mathbb{R}^n , dann existieren positive Konstanten c_1, c_2 , so dass

$$c_1 \|x\|_\alpha \leq \|x\|_\beta \leq c_2 \|x\|_\alpha \quad (x \in \mathbb{R}^n). \tag{A.6}$$

Zum Beispiel, es sei $x \in \mathbb{R}^n$, dann gilt

$$\begin{aligned} \|x\|_2 &\leq \|x\|_1 \leq \sqrt{n} \|x\|_2 \\ \|x\|_\infty &\leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty \\ \|x\|_\infty &\leq \|x\|_1 \leq n \|x\|_\infty. \end{aligned} \tag{A.7}$$

Bemerkung A.1.1 Der Beweis der Ungleichungen (A.4)–(A.7) sei hier zum Selbststudium überlassen.

A.2 MATRIXNORMEN

Da $\mathbb{R}^{m \times n}$ isomorph (d.h. es existiert eine bijektive Abb. mit $\varphi(\lambda A + \mu B) = \lambda\varphi(A) + \mu\varphi(B)$) zu $\mathbb{R}^{m \cdot n}$ ist, sollte die Definition einer Matrixnorm äquivalent sein zur Definition einer Vektornorm.

Definition A.2.1 (Matrixnorm) Es sei $\|\cdot\| : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$. Dann ist $\|\cdot\|$ eine Matrixnorm, wenn folgende Eigenschaften gelten:

$$\begin{aligned} \|A\| &\geq 0 & A &\in \mathbb{R}^{m \times n} & (\|A\| = 0, \text{ g.d.w. } A = 0) \\ \|A + B\| &\leq \|A\| + \|B\| & A, B &\in \mathbb{R}^{m \times n} \\ \|\alpha A\| &= |\alpha| \|A\| & \alpha &\in \mathbb{R}, A \in \mathbb{R}^{m \times n}. \end{aligned} \quad (\text{A.8})$$

Bemerkung A.2.2 Die mit am häufigsten verwendeten Normen in der Numerischen Linearen Algebra sind die Frobenius-Norm

$$\|A\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} \quad (\text{A.9})$$

und die p -Normen

$$\|A\|_p := \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p}. \quad (\text{A.10})$$

Man beachte folgende äquivalente Definition

$$\|A\|_p = \sup_{x \neq 0} \left\| A \begin{pmatrix} x \\ \|x\| \end{pmatrix} \right\|_p = \sup_{\|x\|_p=1} \|Ax\|_p. \quad (\text{A.11})$$

Bemerkung A.2.3 (Submultiplikativität) Die Frobenius-Norm und die p -Normen erfüllen zusätzlich auch noch die Eigenschaft der Submultiplikativität, d.h.

$$\|A \cdot B\| \leq \|A\| \cdot \|B\|. \quad (\text{A.12})$$

Bemerkung A.2.4 Es sei y ein Vektor mit $\|y\|_p = 1$, dann gilt

$$\begin{aligned} \|A \cdot B\|_p &= \max_{x \neq 0} \frac{\|ABx\|_p}{\|x\|_p} = \max_{x \neq 0} \frac{\|ABx\|_p \|Bx\|_p}{\|Bx\|_p \|x\|_p} \\ &\leq \max_{x \neq 0} \frac{\|ABx\|_p}{\|Bx\|_p} \cdot \max_{x \neq 0} \frac{\|Bx\|_p}{\|x\|_p} \\ &= \max_{y \neq 0} \frac{\|Ay\|_p}{\|y\|_p} \cdot \max_{x \neq 0} \frac{\|Bx\|_p}{\|x\|_p} = \|A\|_p \|B\|_p. \end{aligned}$$

Bemerkung A.2.5 • Nicht alle Matrix-Normen erfüllen diese Eigenschaft, z.B. gilt für

$$\|A\|_A := \max |a_{ij}| \text{ und } A = B = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \text{ die Ungleichung } \|A \cdot B\|_A = 2 > \|A\|_A \cdot \|B\|_A.$$

• Für die p -Normen haben wir die wichtige Eigenschaft, dass für alle $A \in \mathbb{R}^{m \times n}$ und $x \in \mathbb{R}^n$

$$\|Ax\|_p \leq \|A\|_p \cdot \|x\|_p. \quad (\text{A.13})$$

A.2.1 Einige Eigenschaften der Matrixnorm

Die Frobenius und p -Normen (speziell $p = 1, 2, \infty$) erfüllen gewisse Ungleichungen, welche in der Analysis von Matrixberechnungen verwendet werden. Es sei $A \in \mathbb{R}^{m \times n}$, dann gilt

$$\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2 \quad (\text{A.14})$$

$$\max_{i,j} |a_{ij}| \leq \|A\|_2 \leq \sqrt{m \cdot n} \max_{i,j} |a_{ij}| \quad (\text{A.15})$$

$$\frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{m} \|A\|_\infty \quad (\text{A.16})$$

$$\frac{1}{\sqrt{m}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1 \quad (\text{A.17})$$

Des Weiteren gilt für $A \in \mathbb{R}^{m \times n}$

$$\begin{aligned} \|A\|_\infty &= \max_{\|x\|_\infty=1} \|Ax\|_\infty = \max_{\|x\|_\infty=1} \left\{ \max_{1 \leq j \leq m} \left| \sum_{k=1}^n a_{jk} x_k \right| \right\} \\ &= \max_{1 \leq j \leq m} \left\{ \max_{\|x\|_\infty=1} \left| \sum_{k=1}^n a_{jk} x_k \right| \right\} = \max_{1 \leq j \leq m} \sum_{k=1}^n |a_{jk}|. \end{aligned}$$

Aufgrund dieser Gleichung bezeichnet man $\|\cdot\|_\infty$ auch als **Zeilensummennorm**. Diese ist mittels der letzten Gleichung auch leicht zu bestimmen, d.h.

$$\|A\|_\infty = \max_{1 \leq j \leq m} \sum_{k=1}^n |a_{jk}| \quad (\text{A.18})$$

Analoges erhält man für die 1-Norm:

$$\begin{aligned} \|A\|_1 &= \max_{\|x\|_1=1} \|Ax\|_1 = \max_{\|x\|_1=1} \sum_{j=1}^m \left| \sum_{k=1}^n a_{jk} x_k \right| = \max_{\|x\|_1=1} \sum_{j=1}^m \sum_{k=1}^n |a_{jk}| |x_k| \text{sign}(a_{jk} x_k) \\ &= \max_{\|x\|_1=1} \sum_{k=1}^n |x_k| \sum_{j=1}^m |a_{jk}| \text{sign}(a_{jk} x_k) = \max_{\|x\|_1=1} \sum_{k=1}^n |x_k| \sum_{j=1}^m |a_{jk}| = \max_{1 \leq k \leq n} \sum_{j=1}^m |a_{jk}|. \end{aligned}$$

Also gilt:

$$\|A\|_1 = \max_{1 \leq k \leq n} \sum_{j=1}^m |a_{jk}| \quad (\text{A.19})$$

Diese Norm bezeichnet man auch als **Spaltensummennorm**.

Bemerkung A.2.6 Einfache Eselbrücke: $\boxed{1}$ -Spaltensummen, $\boxed{\infty}$ -Zeilensummen.

Auch für die 2-Norm läßt sich eine äquivalente Formulierung finden.

Satz A.2.7 Es sei $A \in \mathbb{R}^{m \times n}$. Dann existiert ein Vektor $z \in \mathbb{R}^n$ mit $\|z\|_2 = 1$, so dass $A^T A z = \mu^2 z$ mit $\mu = \|A\|_2$.

Bemerkung A.2.8 Der Satz impliziert, daß $\|A\|_2^2$ eine Nullstelle des Polynoms $p(z) = \det(A^T A - \lambda I)$ ist. Genauer betrachtet, ist die 2-Norm von A die Wurzel des größten Eigenwerts von $A^T A$.

Beweis. Es sei $z \in \mathbb{R}^n$ mit $\|z\|_2 = 1$ und $\|Az\|_2 = \|A\|_2$. Da z die Funktion

$$g(x) = \frac{1}{2} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \frac{1}{2} \frac{x^T A^T A x}{x^T x} \quad (\text{A.20})$$

maximiert, folgt daraus, daß sie $\nabla g(x) = 0$ erfüllt, wobei ∇g der Gradient von g ist ($\nabla := \left(\frac{\partial}{\partial x_1}, \frac{\partial}{\partial x_2}, \dots, \frac{\partial}{\partial x_n} \right)^T$). Eine einfache Ableitung zeigt für $i = 1, \dots, n$

$$\frac{\partial g(z)}{\partial z_i} = \left[z^T z \sum_{j=1}^n (A^T A)_{ij} z_j - (z^T A^T A z) z_i \right] / (z^T z)^2. \quad (\text{A.21})$$

In Vektornotation bedeutet dies $A^T A z = (z^T A^T A z) z$. Setzt man nun $\mu = \|A\|_2$ so ergibt sich daraus die Behauptung. \square

Lemma A.2.9 *Es sei $A \in \mathbb{R}^{m \times n}$. Dann gilt*

$$\|A\|_2 \leq \sqrt{\|A\|_1 \|A\|_\infty}. \quad (\text{A.22})$$

Beweis. Es sei $z \in \mathbb{R}^n$ mit $A^T A z = \mu^2 z$ und $\mu = \|A\|_2$, d.h. es sei $z \neq 0$ ein Vektor, für den das Maximum von $\frac{\|Ax\|_2}{\|x\|_2}$ angenommen wird. Dann gilt

$$\mu^2 \|z\|_1 = \|A^T A z\|_1 \leq \|A^T\|_1 \|A z\|_1 \leq \|A\|_\infty \|A\|_1 \|z\|_1. \quad (\text{A.23})$$

Kondensation von $\|z\|_1$ und Wurzelziehen liefert das gewünschte Ergebnis. \square

Definition A.2.10 (verträgliche Vektornorm) *Eine Matrixnorm $\|A\|$ heißt kompatibel oder verträglich mit der Vektornorm $\|x\|$, falls die Ungleichung gilt*

$$\|Ax\| \leq \|A\| \|x\| \quad (x \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}) \quad (\text{A.24})$$

Kombinationen von verträglichen Normen sind etwa

$\|A\|_G, \|A\|_\infty$ sind verträglich mit $\|x\|_\infty$;

$\|A\|_G, \|A\|_1$ sind verträglich mit $\|x\|_1$;

$\|A\|_G, \|A\|_F, \|A\|_2$ sind verträglich mit $\|x\|_2$, mit

$$\|A\|_G := n \max_{i,k} |a_{ik}| \quad (A \in \mathbb{R}^{m \times n})$$

Bemerkung A.2.11 *Man verifiziere selbständig an einigen Beispielen die Verträglichkeit von o.g. Normenpaaren.*

Abschließend erhalten wir am Ende des Kapitels.

Satz A.2.12 *Die Matrix- p -Normen sind unter allen mit der Vektornorm $\|x\|_p$ verträglichen Matrixnormen die kleinsten.*

B EINFÜHRUNG IN MATLAB

B.1 GRUNDLEGENDE MATLAB-BEFEHLE

Ruft man das Programm MATLAB mit dem Befehl `matlab` auf, so erscheinen auf dem Monitor einige Fenster. Auf den Linux-Rechnern des KIZ müssen Sie vorher die notwendigen Pfade ergänzen. Geben Sie dazu in einem Konsole-Fenster den Befehl `option matlab` ein. Von diesen ist das Befehl-Fenster der primäre Ort um Befehle einzugeben und Fehler oder Ergebnisse abzulesen! Das Prompt-Zeichen `>>` ist im Befehl-Fenster dargestellt und dort findet man üblicherweise einen blinkenden Cursor. Der blinkende Cursor und der MATLAB-Prompt zeigen einem, daß MATLAB eine Eingabe erwartet.

B.1.1 Einfache mathematische Operationen

Genauso wie mit einem simplen Taschenrechner kann man auch mit MATLAB einfache mathematische Operationen ausführen, z.B. ergibt die Eingabe

```
>> 3 + 4
```

die Ausgabe

```
ans =  
    7
```

Man beachte, daß MATLAB im allgemeinen keine Zwischenräume benötigt, um Befehle eindeutig zu verstehen. Alternativ zu dem obigen Beispiel können in MATLAB auch Variablen verwendet werden.

```
>> a = 3  
a =  
    3  
>> b = 4  
b =  
    4  
>> c = a + b  
c =  
    7
```

MATLAB besitzt folgende einfache arithmetische Operationen

Operation	Symbol	Beispiel
Addition, $a + b$	+	$5 + 3$
Subtraktion, $a - b$	-	$23 - 12$
Multiplikation, $a \cdot b$	*	$13.3 * 63.13$
Division, $a \div b$	/ or \	$17/4 = 4 \setminus 17$
Potenz, a^b	^	3^4

Die Reihenfolge, in der eine Folge von Operationen abgearbeitet wird, läßt sich wie folgt beschreiben. Ausdrücke werden von links nach rechts ausgeführt, wobei die Potenzierung die höchste Priorität besitzt gefolgt von Punktoperation, sprich Multiplikation und Division. Die geringste Priorität haben Addition und Subtraktion. Mit Hilfe von Klammern kann diese Vorgehensweise geändert werden, wobei innere Klammern vor äußeren Klammern berechnet werden.



B.1.2 Variablen



Wie in anderen Programmiersprachen hat auch MATLAB Regeln für Variablennamen. Eine Variable repräsentiert ein Datenelement, dessen Wert während der Programmausführung – gegebenenfalls mehrfach – geändert werden kann. Variablen werden anhand ihrer „Namen“ identifiziert. Namen bestehen aus ein bis neunzehn Buchstaben, Ziffern oder Unterstrichen, wobei das erste Zeichen ein Buchstabe sein muß. Man beachte, daß MATLAB Groß- und Kleinschreibung unterscheidet. (Windows ist im Gegensatz zu Linux nicht so restriktiv, da Sie jedoch Programme austauschen wollen, sollten Windows-Benutzer besondere Aufmerksamkeit walten lassen)

Einer Variablen ist Speicherplatz zugeordnet. Wenn man eine Variable verwendet, dann meint man damit entweder den zugeordneten Speicherplatz oder den Wert, der dort augenblicklich abgespeichert ist. Einen Überblick über alle Variablen erhält man mit dem Befehl `who` oder `whos`, wobei letzterer die Angabe des benutzten Speicherplatzes beinhaltet.

Zusätzlich zu selbstdefinierten Variablen gibt es in MATLAB verschiedene spezielle Variablen. Diese lauten

spezielle Variablen	Wert
<code>ans</code>	standard Variablenname benutzt für Ergebnisse
<code>pi</code>	3.1415...
<code>eps</code>	Maschinengenauigkeit
<code>flops</code>	Zähler für die Anzahl der Fließkommaoperationen
<code>inf</code>	steht für Unendlich (eng. infinity). z.B. $1/0$
<code>NaN</code>	eng. Not a Number, z.B. $0/0$
<code>i (und) j</code>	$i = j = \sqrt{-1}$

In MATLAB kann der Speicherplatz, der durch Variablen belegt ist, durch den Befehl `clear` wieder freigegeben werden, z.B.

```
>> clear a b c
```

B.1.3 Kommentare und Punktion

Der Text, der nach einem Prozentzeichen `%` folgt, wird in MATLAB als Kommentar verstanden

```
>> dummy = 4 % Wert von dummy
dummy =
    4
```

Mehrere Befehle können in eine Zeile geschrieben werden, wenn sie durch Kommata oder Semikola getrennt werden. Kommata veranlassen MATLAB, die Ergebnisse anzuzeigen. Bei einem Semikolon wird die Ausgabe unterdrückt. Durch eine Sequenz von drei Punkten kann man einen Befehl in der folgenden Zeile fortsetzen.

B.1.4 Spezielle Funktionen

Eine unvollständige Liste von Funktionen, die MATLAB bereitstellt, ist im folgenden dargestellt. Die meisten Funktionen sind so definiert, wie man sie üblicherweise benutzt.

```
>> y = cos(pi)
y =
   -1
```



MATLAB bezieht sich im Zusammenhang mit Winkelfunktionen auf das Bogenmaß.

Funktion	Bedeutung
abs(x)	Absolutbetrag
cos(x)	Kosinus
exp(x)	Exponentialfunktion: e^x
fix(x)	rundet auf die nächste, vom Betrag her kleinere ganze Zahl
floor(x)	rundet auf die nächste, kleinere ganze Zahl
gcd(x,y)	größter gemeinsamer Teiler von x und y
lcm(x,y)	kleinstes gemeinsames Vielfaches von x und y
log(x)	natürlicher Logarithmus
rem(x,y)	Modulo (eng. remainder of division), z.B. rem(5,2)=1
sign(x)	Signum Funktion, z.B. sign(2.3) = 1, sign(0) = 0, sign(-.3) = -1
sin(x)	Sinus
sqrt(x)	Quadratwurzel
tan(x)	Tangens

B.1.5 Skript-Dateien

Für einfache Probleme ist es schnell und effizient, die Befehle am MATLAB-Prompt einzugeben. Für größere und umfangreichere Aufgabenstellungen bietet MATLAB die Möglichkeit, sogenannte Skript-Dateien zu verwenden, in denen die Befehle in Form einer Textdatei aufgeschrieben sind und die man am Prompt übergibt. MATLAB öffnet dann diese Dateien und führt die Befehle so aus, als hätte man sie am Prompt eingegeben. Die Datei nennt man Skript-Datei oder M-Datei, wobei der Ausdruck M-Datei daher rührt, daß diese Dateien das Suffix `.m` haben, z.B. `newton.m`. Um eine M-Datei zu erstellen, ruft man einen Editor auf und speichert die Datei in dem Verzeichnis, von dem aus man MATLAB gestartet hat oder starten wird. Die Datei, z.B. `newton.m`, wird in MATLAB dann durch Eingabe von `newton` am Prompt aufgerufen.

Für die Benutzung von Skript-Dateien hat MATLAB unter anderem folgende hilfreichen Befehle

M-Datei-Funktionen	
disp(ans)	zeigt den Wert der Variablen <code>ans</code> , ohne ihren Namen auszugeben
input	erwartet vom Benutzer eine Eingabe
keyboard	übergibt zeitweise die Kontrolle an die Tastatur
pause	hält das Programm an, bis eine Taste betätigt wird

Die folgende Skript-Datei `beispiel1.m`

```
% beispiel1.m
% Beispiel fuer eine Skript-Datei
tmp = input(' Geben Sie bitte eine Zahl an >' );
3 * tmp;
```

führt zu der Ausgabe

```

>> beispie11
Geben Sie bitte eine Zahl an > 6
ans =
    18

```

B.1.6 Dateiverwaltung

MATLAB unterstützt eine Vielzahl von Dateiverwaltungsbefehlen, welche es einem ermöglichen Dateien zu listen, Skript-Dateien anzusehen oder zu löschen und Verzeichnisse zu wechseln.

Datei-Management-Funktionen	
<code>cd path</code>	wechselt in das Verzeichnis <code>path</code>
<code>delete beispiel</code>	löscht die Datei <code>beispiel.m</code>
<code>ls</code>	zeigt alle Dateien im aktuellen Verzeichnis an
<code>pwd</code>	zeigt den aktuellen Verzeichnispfad an
<code>type beispiel</code>	zeigt den Inhalt der Datei <code>beispiel.m</code> im Befehl-Fenster
<code>what</code>	zeigt alle M-Dateien und MAT-Dateien im aktuellen Verzeichnis an

B.1.7 Hilfe

Online-Hilfe: Da sich nicht jeder Benutzer alle MATLAB-Befehle merken kann oder auch von einigen auch nur die Syntax unklar ist, bietet MATLAB die Möglichkeit der Online-Hilfe. Dabei gibt es prinzipiell mehrere Möglichkeiten. Ist einem ein Befehl bekannt und man sucht Informationen über die Syntax, so gibt es den Befehl `help`.

```

>> help sqrt

SQRT    Square root.
        SQRT(X) is the square root of the elements of X. Complex
        results are produced if X is not positive.

        See also SQRTM.

```

Als Beispiel haben wir uns hier die Hilfe zu dem Befehl `sqrt` ausgeben lassen.

Die andere Möglichkeit der von MATLAB gelieferten Hilfe ist durch den Befehl `lookfor` gegeben. Hier durchsucht das Programm alle ersten Zeilen der MATLAB Hilfe-Kennwörter und Skript-Dateien die im MATLAB Suchpfad zu finden sind. Das Bemerkenswerte dabei ist, daß dieser Begriff kein Befehl zu sein braucht.

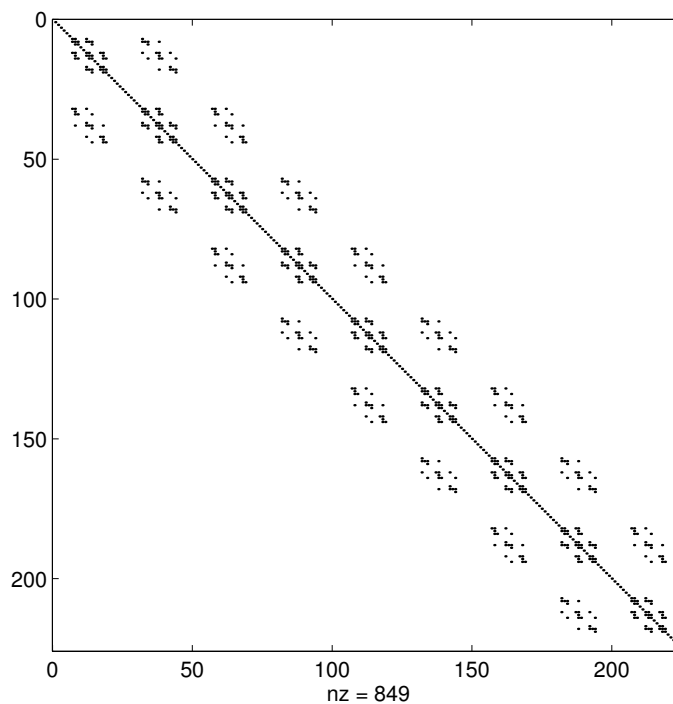
```
>> lookfor cholesky
    CHOL    Cholesky factorization
>> CHOL    Cholesky factorization.

CHOL(X) uses only the diagonal and upper triangle of X.
The lower triangular is assumed to be the (complex conjugate)
transpose of the upper.  If X is positive definite, then
R = CHOL(X) produces an upper triangular R so that R'*R = X.
If X is not positive definite, an error message is printed.

With two output arguments, [R,p] = CHOL(X) never produces an
error message.  If X is positive definite, then p is 0 and R
is the same as above.  But if X is not positive definite, then
p is a positive integer and R is an upper triangular matrix of
order q = p-1 so that R'*R = X(1:q,1:q).

>> lookfor factorization
    CHOL    Cholesky factorization.
    QRDELETE Delete a column from the QR factorization.
    QRINSERT Insert a column in the QR factorization.
    SYMBFACT Symbolic factorization analysis.
```

Eine weitere Möglichkeit, sich Hilfe zu verschaffen, besteht darin, das Helpdesk aufzurufen. Wenn Sie helpdesk am Prompt eingeben, öffnet sich die folgende Hilfsumgebung



B.2 MATHEMATIK MIT MATRIZEN

B.2.1 Matrixkonstruktion und Adressierung

einfache Matrix Konstruktionen	
$x=[1\ 4\ 2*\pi\ 4]$	erstelle einen Zeilenvektor x mit genannten Einträgen
$x=anfang:ende$	erstelle einen Zeilenvektor x beginnend mit $anfang$, Inkrement 1 und endend mit $ende$
$x=anfang:inkrement:ende$	ähnliches wie oben mit dem Inkrement $inkrement$
$x=linspace(anfang,ende,n)$	erzeugt einen Zeilenvektor der Dimension n mit $x(i) = \frac{(n-i) \cdot anfang + (i-1) \cdot ende}{n-1}$

Im folgenden sind einige charakteristische Beispiele aufgeführt.

```
>> B = [1 2 3 4; 5 6 7 8]
B =
     1  2  3  4
     5  6  7  8
```

Der Operator $'$ liefert für reelle Matrizen die Transponierte.

```
>> C = B'
C =
     1  5
     2  6
     3  7
     4  8
```

Der Doppelpunkt $:$ in der zweiten Komponente spricht alle vorhandenen Spalten an, d.h. er ist ein zu $1:4$ äquivalenter Ausdruck.

```
>> C = B(1,:)
C =
     1  2  3  4
>> C = B(:,3)'
C =
     3  7
```

Es lassen sich auch einzelne Komponenten neu definieren.

```
>> A = [1 2 3; 4 5 6; 7 8 9]
A =
     1  2  3
     4  5  6
     7  8  9
>> A(1,3) = 9
A =
     1  2  9
     4  5  6
     7  8  9
```

Ist ein Eintrag noch nicht definiert, so verwendet MATLAB die minimale Erweiterung dieser Matrix und setzt undefinierte Einträge zu Null.


```

>> A(2,5) = 4
A =
    1  2  9  0  0
    4  5  6  0  4
    7  8  9  0  0

```

Im folgenden werden die Vektoren $(3, 2, 1)$ und $(2, 1, 3, 1, 5, 2, 4)$ dazu verwendet, die Matrix C zu indizieren, d.h. C hat die Struktur

```

A(3,2) A(3,1) A(3,3) A(3,1) A(3,5) A(3,2) A(3,4)
A(2,2) A(2,1) A(2,3) A(2,1) A(2,5) A(2,2) A(2,4)
A(1,2) A(1,1) A(1,3) A(1,1) A(1,5) A(1,2) A(1,4)

```

In MATLAB erhält man nun

```

>> C=A(3:-1:1,[2 1 3 1 5 2 4])
C =
    8  7  9  7  0  8  0
    5  4  6  4  4  5  0
    2  1  9  1  0  2  0

```

Ein weiteres Beispiel für Indizierung ist

```

>> C=C(1:2,2:3)
C =
    7  9
    4  6

```

Im nächsten Beispiel wird ein Spaltenvektor dadurch konstruiert, daß alle Elemente aus der Matrix C hintereinander gehängt werden. Dabei wird spaltenweise vorgegangen.

```

>> b=C(:) '
b =
    7  4  9  6

```

Das Löschen einer ganzen Zeile oder Spalte kann durch das Umdefinieren in eine 0×0 -Matrix geschehen, z.B.

```

>> C(2,:)=[ ]
C =
    7  9

```

B.2.2 Skalar-Matrix-Operationen

In MATLAB sind Skalar-Matrix-Operationen in dem Sinne definiert, daß Addition, Subtraktion, Division und Multiplikation mit einem Skalar elementweise durchgeführt werden. Es folgen zwei erklärende Beispiele.

```

>> B - 1
ans =
    0  1  2  3
    4  5  6  7
>> 9 + 3 * B
ans =
   12  15  18  21
   24  27  30  33

```

B.2.3 Matrix-Matrix-Operationen



Die Operationen zwischen Matrizen sind nicht so kanonisch zu definieren wie die zwischen Skalar und Matrix, insbesondere sind Operationen zwischen Matrizen unterschiedlicher Dimension schwer zu definieren. Desweiteren sind die Operationen $*$ und $.*$, bzw. $/$ und $./$ sowie \backslash und $.\backslash$ zu unterscheiden. In nachfolgender Tabelle sind die Matrixoperationen beschrieben.

komponentenweise Matrixoperationen	
Beispieldaten	$a = [a_1, a_2, \dots, a_n], b = [b_1, b_2, \dots, b_n], c$ ein Skalar
komp. Addition	$a + c = [a_1 + c \ a_2 + c \ \dots \ a_n + c]$
komp. Multiplikation	$a * c = [a_1 \cdot c \ a_2 \cdot c \ \dots \ a_n \cdot c]$
Matrix-Addition	$a + b = [a_1 + b_1 \ a_2 + b_2 \ \dots \ a_n + b_n]$
komp. Matrix-Multiplikationen	$a .* b = [a_1 \cdot b_1 \ a_2 \cdot b_2 \ \dots \ a_n \cdot b_n]$
komp. Matrix-Div. von rechts	$a ./ b = [a_1/b_1 \ a_2/b_2 \ \dots \ a_n/b_n]$
komp. Matrix-Div. von links	$a .\backslash b = [b_1/a_1 \ b_2/a_2 \ \dots \ b_n/a_n]$
komp. Matrix-Potenz	$a.^c = [a_1^c \ a_2^c \ \dots \ a_n^c]$ $c.^a = [c^{a_1} \ c^{a_2} \ \dots \ c^{a_n}]$ $a.^b = [a_1^{b_1} \ a_2^{b_2} \ \dots \ a_n^{b_n}]$

Es folgen nun einige Beispiele zu Matrixoperationen

```

>> g=[1 2 3; 4 5 6]; % zwei neue Matrizen
>> h=[2 2 2; 3 3 3];
>> g+h % addiere g und h komponentenweise
ans =
     3     4     5
     7     8     9
>> ans-g % subtrahiere g von der vorherigen Antwort
ans =
     2     2     2
     3     3     3
>> h.*g % multipliziere g mit h komponentenweise
ans =
     2     4     6
    12    15    18
>> g*h' % multipliziere g mit h'
ans =
    12    18
    30    45

```

B.2.4 Matrix-Operationen und -Funktionen

Matrixfunktionen	
<code>reshape(A,m,n)</code>	erzeugt aus den Eintägen der Matrix A eine $m \times n$ -Matrix, wobei die Einträge spaltenweise aus A gelesen werden.
<code>diag(A)</code>	ergibt die Diagonale von A als Spaltenvektor
<code>diag(v)</code>	erzeugt eine Diagonalmatrix mit dem Vektor v in der Diagonalen
<code>tril(A)</code>	extrahiert den unteren Dreiecksanteil der Matrix A
<code>triu(A)</code>	extrahiert den oberen Dreiecksanteil der Matrix A

Es folgen einige Beispiele

```

>> g=linspace(1,9,9)    % ein neuer Zeilenvektor
g =
     1     2     3     4     5     6     7     8     9
>> B=reshape(g,3,3)    % macht aus g eine 3 x 3 Matrix
B =
     1     4     7
     2     5     8
     3     6     9
>> tril(B)
ans =
     1     0     0
     2     5     0
     3     6     9
    
```

Funktion	Bedeutung
R=chol(A)	Choleskyzerlegung
cond(A)	Konditionszahl der Matrix A
d=eig(A)	Eigenwerte und -vektoren
[V,d]=eig(A)	
det(A)	Determinante
hess(A)	Hessenbergform
inv(A)	Inverse
[L,U]=lu(A)	Zerlegung gegeben durch Gauss-Algorithmus
norm(A)	euklidische-Norm
rank(A)	Rang der Matrix A

Bemerkung : Der \backslash Operator ist auch für Matrizen definiert und liefert in Kombination mit Vektoren für reguläre Matrizen ihre Inverse, d.h. $A^{-1}x=A\backslash x$.



B.2.5 Spezielle Matrizen

spezielle Matrizen	
eye(n)	erzeugt eine Einheitsmatrix der Dimension n
ones(m,n)	erzeugt eine $m \times n$ -Matrix mit den Einträgen 1
zeros(m,n)	erzeugt eine $m \times n$ -Matrix mit den Einträgen 0

B.2.6 Spezielle Funktionen für schwachbesetzte Matrizen

Bei vielen numerischen Anwendungen treten schwachbesetzte Matrizen auf. MATLAB hat für solche Matrizen besondere Sparse-Funktionen, die dieser Eigenschaft Rechnung tragen.



Funktion	Bedeutung
find(A)	findet Indizes von Nichtnulleinträgen
nnz(A)	Anzahl an Nichtnulleinträgen
spdiags(v)	erzeugt eine Sparse-Diagonalmatrix mit dem Vektor v als Diagonale
speye(n)	erzeugt eine Sparse-Einheitsmatrix
spy(A)	visualisiert die Struktur der Matrix A

Kurze Illustration der Funktionsweise obiger Befehle anhand einiger Beispiele.

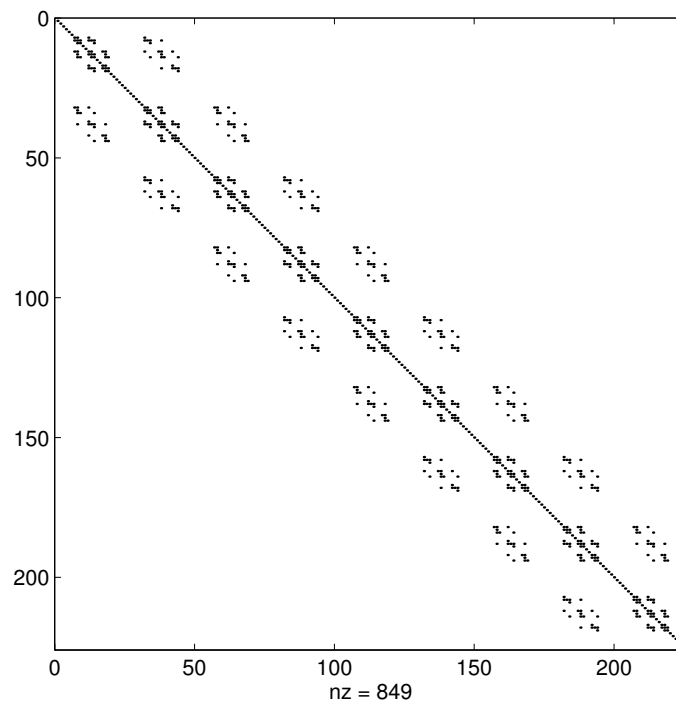
```

>> E=eye(100); % vollbesetzte 100 x 100 Einheitsmatrix
>> Es=sparse(E); % Sparse-Version von E
>> whos
      Name      Size      Elements  Bytes  Density  Comple
      E  100 by 100      10000  80000    Full    No
      Es 100 by 100         100   1600  0.0100    No
Grand total is 10100 elements using 81600 bytes

>> A=spdiags([7*ones(4,1),ones(4,1),2*ones(4,1)],[-1,0,1],4,4);
>> nnz(A)
ans =
      10
>> full(A)
ans =
      1   2   0   0
      7   1   2   0
      0   7   1   2
      0   0   7   1

```

Als Beispiel für `spy` sei hier die Besetzungsstruktur einer 3D-FEM Steifigkeitsmatrix gezeigt.



B.3 DATENVERWALTUNG

Für die meisten Anwendungen genügt es, Datenfelder in einem Format abzuspeichern und wieder laden zu können. Die Befehle `load` und `save` setzen voraus, daß die Daten in einem System unabhängigen, binären Format in einer Datei mit dem Suffix `.mat` gespeichert sind oder in einem einfachen ASCII-Format vorliegen.

B.3.1 Daten speichern

Im folgenden wird eine 3×5 -Matrix im binär-Format in der Datei `A.mat` gespeichert. Diese Daten sind sehr kompakt gespeichert.

```
>> A=zeros(3,5);
>> save A
```

Gibt man sich aber den Inhalt dieser Datei auf dem Bildschirm aus, so gibt er wenig Sinn. Möchte man sich also z.B. einen Lösungsvektor sichern um ihn später „per Hand zu analysieren“ so speichere man die Daten als ASCII-Datei, dabei kann man wählen zwischen einer 8-stelligen oder 16-stelligen Abspeicherung.

```
>> save mat1.dat A -ascii           % 8-stellige Speicherung
>> save mat2.dat A -ascii -double  % 16-stellige Speicherung
```

Hier wurde die Matrix mit 8-stelligem Format in der Datei `mat1.dat` gespeichert, bzw. 16-stellig in der Datei `mat2.dat`.

B.3.2 Daten laden

Mit dem Befehl `load A` versucht MATLAB, die in `A.mat` gespeicherten Daten in einem Datenfeld `A` zu speichern. Auch ASCII-Dateien kann MATLAB lesen. Da es hier jedoch keine Standardendung gibt, ist die Datei inklusive Endung anzugeben. Es ist darauf zu achten, daß ein rechteckiges Feld an Daten vorliegt, d.h. daß m Zeilen mit jeweils n numerischen Werten vorliegen. MATLAB erstellt dann eine $m \times n$ -Matrix mit dem Namen der Datei ohne Suffix.

```
>> load mat1.dat
>> whos

      Name      Size      Elements  Bytes  Density  Complex
      mat1      3 by 5          15     120     Full     No

Grand total is 15 elements using 120 bytes
```

B.4 AUSGABE VON TEXT

Mit dem Befehl `fprintf` lassen sich Strings, d.h. Zeichenfolgen, auf dem Bildschirm ausgeben.

```
>> fprintf('\n Hello world %12.3e\n',4);
      Hello world      4.000e+00
```

Man sieht, daß der auszugebende Text zusätzliche Zeichen enthält, die nicht mit ausgedruckt werden. Diese Zeichen nennt man Escape-Sequenzen. In obigem Beispiel ist die Sequenz `\n` eingebaut. `\n` steht für „newline“ und sorgt dafür, daß bei der Textausgabe an diesen Stellen eine neue Zeile begonnen wird. Der Ausdruck `%12.3e` dient als Platzhalter für einen reellen Wert, der durch Komma getrennt hinter der Zeichenkette folgt. Dabei sei die Zahl in Exponentialdarstellung auszugeben, wofür 12 Stellen mit 3 Nachkommastellen bereitgestellt. Im Beispiel ist dies der Wert 4. Anstatt eines expliziten Wertes können auch Variablen oder Ausdrücke, z.B. `3 * 4` stehen. Ein Platzhalter kann mehrfach in einem `printf`-Befehl vorkommen. In diesem Fall müssen hinter der Zeichenkette genau so viele Werte folgen, wie Platzhalter angegeben sind. Die Reihenfolge der Werte muß mit der Reihenfolge der Platzhalter übereinstimmen, da die Ausdrücke von links nach rechts bewertet werden. Die Escape-Sequenzen dürfen im Text an beliebiger Stelle stehen.

Ausgabeformate	
Befehl	Ausgabe
<code>fprintf('%0e\n',1.234567)</code>	1e00+
<code>fprintf('%2e\n',1.234567)</code>	1.23e00+
<code>fprintf('%5e\n',1.234567)</code>	1.23456e00+
<code>fprintf('%10.0e\n',1.234567)</code>	1e00+
<code>fprintf('%10.2e\n',1.234567)</code>	1.23e00+
<code>fprintf('%10.5e\n',1.234567)</code>	1.2346e00+
<code>fprintf('%10.2f\n',1.234567)</code>	1.23
<code>fprintf('%10.5f\n',1.234567)</code>	1.23457

MATLAB rundet numerische Werte bei der Ausgabe, wenn nötig!

B.5 KONTROLLBEFEHLE

B.5.1 For-Schleifen

1. Mit jeglicher gültigen Matrix-Darstellung läßt sich eine FOR-Schleife definieren, z.B.

```

>> data = [1 7 3 2; 5 4 7 2]
data =
     1  7  3  2
     5  4  7  2
>> for n=data
     x=n(1)-n(2)
end
x =
    -4
x =
     3
x =
    -4
x =
     0

```

2. FOR-Schleifen können nach Belieben geschachtelt werden.

```

>> for k=3:5
     for l=4:-1:2
         A(k,l)=k^2-l;
     end
end
>> A
A =
     0     0     0     0
     0     0     0     0
     0     7     6     5
     0    14    13    12
     0    23    22    21

```

3. FOR-Schleifen sollten vermieden werden, wenn immer sie durch eine äquivalente Matrix-Darstellung ersetzt werden können. Der folgende Ausdruck ist unten optimierter Version aufgeführt.



```

>> for n=1:10
    x(n)=sin(n*pi/10);
end
>>x
x =
    Columns 1 through 7
    0.3090  0.5878  0.8090  0.9511  1.0000  0.9511  0.8090
    Columns 8 through 10
    0.5878  0.3090  0.0000

>> n=1:10;
>> x=sin(n*pi/10);
>> x
x =
    Columns 1 through 7
    0.3090  0.5878  0.8090  0.9511  1.0000  0.9511  0.8090
    Columns 8 through 10
    0.5878  0.3090  0.0000

```

4. Um die Ausführungsgeschwindigkeit zu maximieren, sollte `b` benötigter Speicherplatz vor Ausführung der FOR-Schleife alloziert werden.



```

>> x=zeros(1,10);
>> x(1)=0.5;
>> for n=2:10
    x(n)=x(n-1)*sin(n*pi/10);
end
>> x
x =
    Columns 1 through 7
    0.5000  0.2939  0.2378  0.2261  0.2261  0.2151  0.1740
    Columns 8 through 10
    0.1023  0.0316  0.0000

```

B.5.2 WHILE-Schleifen

WHILE-Schleifen sind wie folgt aufgebaut.

```

while Aussage
    Anweisungen
end

```

Die Anweisungen zwischen `while` und `end` werden so lange ausgeführt, wie Aussage wahr ist, z.B.

```

>> num=0; EPS=1;
>> while (1+EPS) > 1
    EPS=EPS/2;
    um=num+1;
end
>> num
num =
    53

```

B.5.3 IF-ELSE-END Konstrukte

Eine IF-ELSE-END-Schleife enthält nach dem IF eine Aussage, die daraufhin überprüft wird, ob sie wahr oder falsch ist. Ist sie wahr, so werden die in den folgenden Zeilen stehenden Anweisungen ausgeführt und die Schleife beendet. Ist sie falsch, so erfolgen die Anweisungen, die dem ELSE folgen (ELSE ist optional). Solch ein Konstrukt kann erweitert werden um beliebig viele ELSIF Befehle, die dieselbe Funktion haben wie der am Anfang stehende IF Befehl, aber nur beachtet werden, falls alle vorher überprüften Aussagen falsch sind.

```

if Aussage1
    Anweisungen, wenn Aussage1 wahr
elseif Aussage2
    Anweisungen, wenn Aussage1 falsch und Aussage2 wahr
elseif
    Anweisungen, wenn Aussage1 und Aussage2 falsch
end

```

B.5.4 Relationen und logische Operatoren

Relationen

<	kleiner als
<=	kleiner als oder gleich
>	größer als
>=	größer als oder gleich
==	gleich
~=	ungleich

logische Operatoren

&	UND
	ODER
~	NICHT

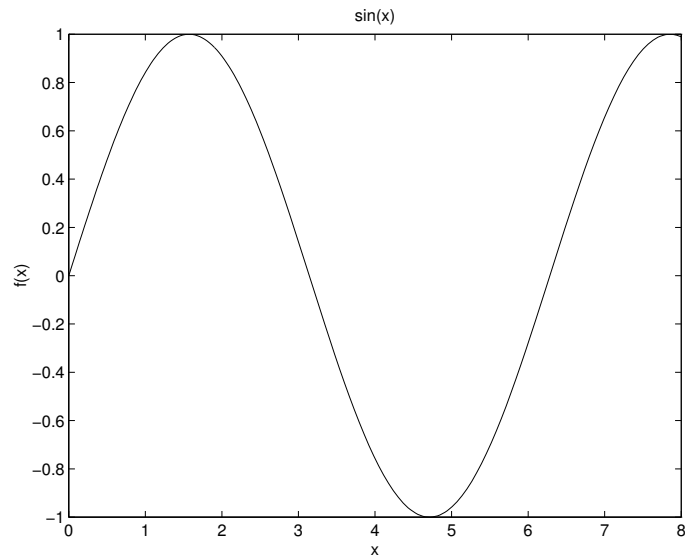
B.6 GRAPHISCHE DARSTELLUNG

B.6.1 Zweidimensionale Graphiken

```

>> f='sin(x)';
>> fplot(f,[0 8]);
>> title(f),xlabel('x'),ylabel('f(x)');

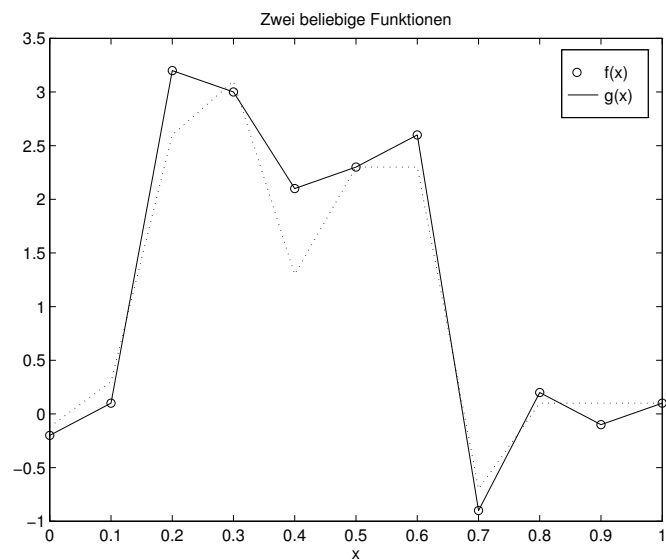
```

```

>> x=0:0.1:1;
>> y=[-0.2 0.1 3.2 3 2.1 2.3 2.6 -0.9 0.2 -1 .1];
>> z=[-0.12 0.3 2.6 3.1 1.3 2.3 2.3 -0.7 0.1 .1 .1];
>> plot(x,y,'o',x,y,x,z,':');
>> title('Zwei beliebige Funktionen'),xlabel('x');
>> legend('f(x)', 'g(x)')

```



Linientypen und Farben

Symbol	Farbe	Symbol	Linientyp
y	gelb	.	Punkt
m	magenta	o	Kreis
c	cyan	x	x-Markierung
r	rot	+	+ -Markierung
g	grün	*	Sternchen
b	blau	—	durchgezogene Linie
w	weiß	:	gepunktete Linie
k	schwarz	—.	Strichpunkt-Linie
		--	gestrichelte Linie

2-D Graphikanweisung

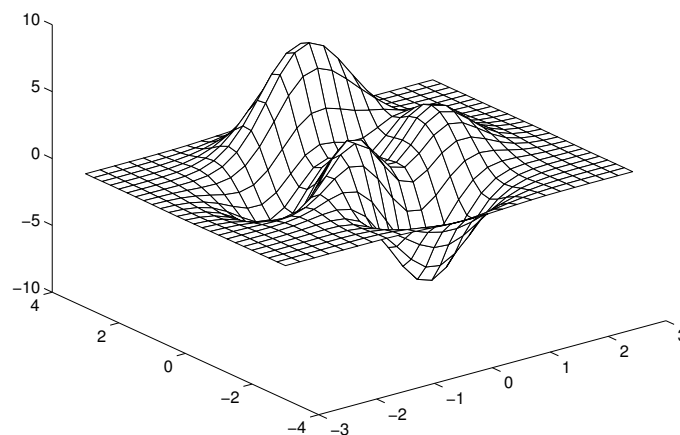
axis	modifiziert die Axen-Proportionen
clf	löscht die Graphik im Graphik-Fenster
close	schließt das Graphik-Fenster
grid	erzeugt ein achsenparalleles Gitter
hold	ermöglicht das Überlagern von Graphiken
subplot	erstellt mehrere Teilgraphiken in einem Fenster
text	gibt Text an vorgegebener Stelle aus
title	zeigt einen Titel an
xlabel	beschriftet die x-Achse
ylabel	beschriftet die y-Achse
colormap(white)	wechselt die Farbtabelle, für S/W-Monitore

B.6.2 Dreidimensionale Graphiken

```

>> [X, Y, Z] = peaks(25)
>> mesh(X, Y, Z)

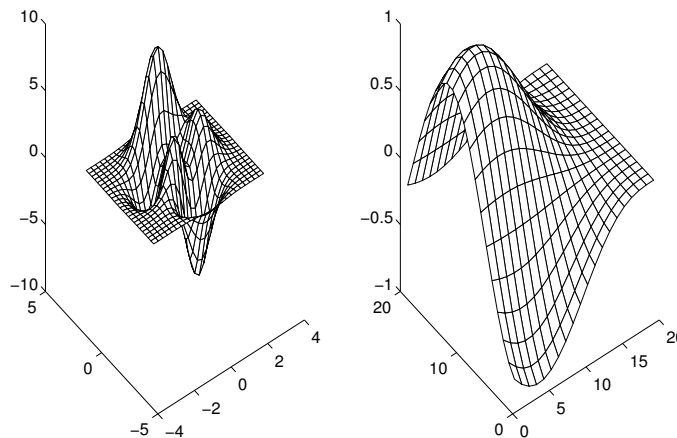
```



```

>> subplot(1,2,1);
>> [X,Y,Z] = peaks(25);
>> mesh(X,Y,Z);
>> subplot(1,2,2);
>> X=1:20;
>> Y=1:20;
>> Z(X,Y) = -(cos(X/4))' * (sin((20-Y)/10).^3);
>> mesh(X,Y,Z);

```



B.6.3 Graphiken drucken

Graphik-Druckbefehl		
print [-dAusgabetyyp] [-Optionen] [Dateiname]		
Ausgabetyyp	-dps	Postscript für Schwarzweißdrucker
	-dpsc	Postscript für Farbdrucker
	-deps	Encapsulated Postscript
	-depvc	Encapsulated Color Postscript
Optionen	-P<Drucker>	Spezifiziert den zu benutzenden Drucker

Die Eingabe

```
>> print fig4 -deps
```

erzeugt die Datei fig4.eps.

B.7 FORTGESCHRITTENES

B.7.1 MATLAB-Skripte

Wie schon unter B Skript-Dateien erwähnt, läßt sich eine Abfolge von MATLAB-Befehlen auch in einer Datei speichern. Diese läßt sich dann am Befehl-Fenster aufrufen und die gespeicherte Folge von Befehlen wird dann ausgeführt. Solche Dateien werden als MATLAB-Skripte oder M-Files bezeichnet. Alle o.g. Befehle lassen sich in einer Datei z.B. mit dem Namen `abc.m` zusammenstellen. Für die Wahl des Dateinamens gelten dabei die folgenden Regeln:

- das erste Zeichen ist ein Buchstabe und
- die Datei hat die Endung `.m`.

Um ein solches M-File zu erstellen, ruft man den Editor auf, der es erlaubt Text einzugeben und diesen als Datei zu speichern. Wenn man den Editor gestartet hat, gebe man die folgenden Befehle ein und speichere die Datei unter dem Namen `abc.m`

```
a = 1;
b = 3;
c = -5;
x1 = (-b + sqrt(b^2 - 4*a*c)) / (2*a)
x2 = (-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

Um nun die Befehle aus der Datei `abc.m` auszuführen, gibt man im MATLAB-Befehlsfenster

```
>> abc
```

ein. MATLAB sucht im aktuellen Pfad nach der Datei `abc.m` und führt die darin enthaltenen Befehle aus. Das aktuelle Verzeichnis wird angezeigt, wenn man

```
>> pwd
```

eingibt (`pwd`, engl. `print working directory`).

B.7.2 Erstellen eigener Funktionen

Es wird sicherlich etwas komfortabler sein, eine eigene Funktion zu haben, der man a, b und c als Argument übergibt und die beiden Wurzeln als Ergebnis erhält, als jedesmal erneut ein eigenes M-File anzufertigen. Ein solches Programm könnte z.B. die folgende Datei `root.m` sein

```
%-----
modified "abc.m"
%-----

a = input('Enter a: ');
b = input('Enter b: ');
c = input('Enter c: ');

x1 = (-b + sqrt(b^2 - 4*a*c)) / (2*a)
x2 = (-b - sqrt(b^2 - 4*a*c)) / (2*a)
```

Die Prozentzeichen in den ersten Zeilen dienen der Dokumentation. Alles was einem solchen Zeichen folgt, wird von MATLAB nicht ausgewertet, sprich interpretiert. Die Argumente werden in dieser Funktion jedoch nur bedingt übergeben und ausgegeben. Eine Funktion die diese Aufgabe erfüllt ist

```
%-----  
modified "abc.m"  
%-----  
  
function [x1, x2] = quadroot(a,b,c)  
  
radical = sqrt(b^2 - 4*a*c);  
  
x1 = (-b + radical)/(2*a)  
x2 = (-b - radical)/(2*a)
```

Wenn eine Funktion keine Rückgabewerte hat, können die eckigen Klammern mit den Variablen und das Gleichheitszeichen fehlen. Fehlen die Eingabeparameter, so kann auch der Klammerausdruck nach `quadroot` fehlen. Auf die in der Funktion verwendeten Variablen, hier `radical`, kann man vom Befehlsfenster nicht zugreifen. Ist die Funktion in der Datei `quadroot.m` gespeichert, so erhält man die Wurzeln, indem man

```
[x1, x2] = quadroot(1, 3, -5);
```

eingibt. Es kann vom Befehlsfenster oder einer anderen Datei immer nur die erste Funktion in einer Datei aufgerufen werden. Dies bedeutet, in einer Datei können mehrere Funktionen stehen, aber nur die erste Funktion kann extern aufgerufen werden. Alle weiteren Funktionen können nur von Funktionen in der gleichen Datei aufgerufen werden. Für den Anfang ist es einfacher, wenn jede Datei nur eine Funktion enthält. Entscheidend für den Funktionsnamen ist der Name der Datei.

C SPEICHERFORMATE FÜR SCHWACH BESETZTE MATRIZEN

Im Folgenden werden wir einige Speicherformate für schwach besetzte Matrizen vorstellen, welche speziell im Zusammenhang mit Iterationsverfahren von Interesse sind. Den Speicherformaten für schwach besetzte Systeme kommen aus zweierlei Gründen eine besondere Bedeutung zu:

- Die Wahl der Speicherform hat direkte Auswirkungen auf die verwendeten Lösungsverfahren (spezielle Matrixform, z.B. Bandmatrix, symmetrische Matrix, etc.). Unter anderen wegen ihres Auffüllverhaltens (Erzeugung neuer Nichtnullelemente) der verschiedenen Algorithmen, der Zugriffshäufigkeit auf spezielle Einträge (z.B. der Diagonaleinträge) und somit ihrer Konsequenzen auf die Speicherstruktur.
- In den meisten höheren Programmiersprachen C/C++, Java, Fortran, etc. stehen für die meisten Standardprobleme der numerischen Linearen Algebra verschiedene Programmpakete (Lapack, Linpack, Eispack u.v.a.) mit fertig implementierten numerischen Algorithmen zur Verfügung. Diesen Bibliotheken müssen die Daten aber in bestimmten Datenstrukturen übergeben werden.

C.1 KOORDINATENFORMAT (COO-FORMAT)

Das einfachste Format zur Speicherung von schwach besetzten Matrizen für iterative Verfahren ist das Koordinatenformat (Coordinate Storage, COO-Format).

Es verwendet 3 Vektoren der Länge $\text{nnz}A$ ($\text{nnz} \simeq$ number of nonzero elements):

- ein Vektor „Wert“ vom Datentyp der Einträge a_{ij} , in dem die Nichtnullelemente von A in beliebiger Reihenfolge stehen,
- ein Index-Vektor „Zeilenindex“ und
- ein Index-Vektor „Spaltenindex“, in denen zu jedem Element von „Wert“ die entsprechenden Koordinaten, d.h. Zeilen- und Spaltenindices, aus der Matrix A eingetragen werden. Meist sind dies Vektoren vom `unsigned (long) integer`.

Dieses Format wird häufig wegen seiner Einfachheit verwendet (dient oft als eine von vielen möglichen Schnittstelle zu Bibliotheken). Der Umstand, daß die Nichtnull-Matrixelemente in beliebiger Reihenfolge im Feld „Wert“ stehen dürfen, ist sowohl ein Vorteil als auch ein Nachteil des COO-Formats. Einerseits können die bei der Lösung des Gleichungssystems neu auftretenden Nichtnullelemente einfach am Ende der Felder ohne irgendwelchen Umordnungsaufwand angehängt werden. Andererseits ist der Suchaufwand nach einem bestimmten Matrixelement enorm, da ja keinerlei Strukturinformation der Matrix A in das Speicherformat selbst eingeht (ist das gesuchte Element gleich Null, müssen die Felder komplett durchsucht werden, bis diese Information vorliegt). Ein weiterer Nachteil ist auch der immer noch recht hohe Speicherbedarf.

Ein Beispiel soll die Datenspeicherung illustrieren.

Beispiel C.1.1

$$A = \begin{pmatrix} 10 & 0 & 0 & 22 & 31 & 0 \\ -4 & 8 & 0 & 0 & 0 & 0 \\ 0 & 0 & -6 & 9 & 0 & 0 \\ 0 & 0 & 0 & 7 & 0 & 0 \\ 0 & 0 & 5 & 0 & 20 & 11 \\ 12 & 0 & 0 & 14 & 0 & 44 \end{pmatrix}$$

geht über zu

Wert a_{ij}	44	-6	20	12	14	31	5	-4	10	8	7	9	22	11
Zeilenindex i	6	3	5	6	6	1	5	2	1	2	4	3	1	5
Spaltenindex j	6	3	5	1	4	5	3	1	1	2	4	4	4	6

In Matlab kann man einen eigenen Datentypen COO-Matrix definieren und sowohl die Operatoren $+$, $*$, $-$, $/$ überladen, als auch die Ein- und Ausgabe dem bekannten Verhalten anpassen. Im Folgenden beschränken wir uns auf die Wiedergabe einiger Methoden, die den Datentypen `coomatrix` definieren, Ein- und Ausgabe sowie Multiplikation ermöglichen, so dass man u.a. das CG-Verfahren testen könnte. Die Routinen bzw. Methoden zu einer Klasse stehen in Matlab üblicherweise in einem Unterverzeichnis gleichen Namens, d.h. die Methoden `coomatrix.m` etc. stehen im Verzeichnis `@coomatrix`. Für einen ersten Einstieg in Klassendefinitionen in Matlab, sei auf das Beispiel `polynom` in der Matlab-Dokumentation verwiesen.

MATLAB-Funktion: @coomatrix/coomatrix.m

```

1 function p = coomatrix(i, j, a)
2 %COOMATRIX Coordinate format class constructor.
3 %   A = COOMATRIX(i, j, a) creates a compressed matrix object from
4 %   the vectors i, j, a, containing the entries of a matrix A.
5 if nargin == 0
6     p.a = []; p.row = []; p.col = [];
7     p.m = 0;
8     p.n = 0;
9     p = class(p, 'coomatrix');
10 elseif nargin == 1 && isa(i, 'coomatrix')
11     p = i;
12 elseif nargin == 1 && isnumeric(i)
13     [p.row, p.col, p.a] = find(i);
14     p.row = p.row'; p.col = p.col'; p.a = p.a';
15     p.m = max(p.row); p.n = max(p.col);
16     p = class(p, 'coomatrix');
17 else
18     p.row = i(:)';
19     p.col = j(:)';
20     p.a = a(:)';
21     p.m = max(p.row);
22     p.n = max(p.col);
23     p = class(p, 'coomatrix');
24 end

```


MATLAB-Funktion: @coomatrix/display.m

```

1 function display(p)
2 % COOMATRIX/DISPLAY Command window display of coomatrix
3 formattedspace
4 disp([inputname(1),' = '])
5 formattedspace
6 short = ~(max(p.row)>9999 || max(p.col)>9999);
7 for k=1:length(p.row)
8     i = num2str(p.row(k)); j = num2str(p.col(k));
9     if short
10        txt = '          ,          ';
11        txt(5-length(i):5) = ['(',i];
12        txt(7:7+length(j)) = [j,')'];
13    else
14        txt = '          ,          ';
15        txt(11-length(i):11) = ['(',i];
16        txt(13:13+length(j)) = [j,')'];
17    end
18    disp([txt, sprintf('%13.4f',p.a(k))])
19 end
20 formattedspace
21
22 function formattedspace
23 switch get(0,'FormatSpacing')
24     case 'loose'
25         disp(' ');
26 end

```

MATLAB-Funktion: @coomatrix/subsasgn.m

```

1 function p = subsasgn(p,s,b)
2 % SUBSASGN
3 switch s.type
4     case '()'
5         if size(s.subs,1)~=1 || size(s.subs,2)~=2
6             error('??? Index exceeds matrix dimensions.')
7         end
8         ii = s.subs{1,1}; jj = s.subs{1,2};
9         for i = 1:length(ii)
10            for j = 1:length(jj)
11                p.row = [p.row,ii(i)]; p.col = [p.col,jj(j)];
12                p.a = [p.a,b(i,j)];
13            end
14        end
15        p.m = max(p.row); p.n = max(p.col);
16    otherwise
17        error('Specify value for i and j as A(i,j)')
18 end

```

MATLAB-Beispiel:

Anlegen einer Matrix im COO-Format funktioniert mit der Anweisung <code>coomatrix</code>	<pre>>> A=coomatrix([1,2],[2,3],[1.2,1.4]) A = (1,2) 1.2000 (2,3) 1.4000</pre>
Durch die Funktion <code>subsref</code> können einzelne Einträge zu der Matrix hinzugefügt. Indexpaare (i,j) können dabei doppelt auftreten.	<pre>>> A(1,4)=2; >> A(1,4)=3 A = (1,2) 1.2000 (2,3) 1.4000 (1,4) 2.0000 (1,4) 3.0000</pre>

MATLAB-Funktion: @coomatrix/subsref.m

```

1 function b = subsref(a,s)
2 % SUBSREF
3 switch s.type
4 case '()'
5     if size(s.subs,1)~=1 || size(s.subs,2)~=2
6         error('??? Index exceeds matrix dimensions.')
7     end
8     ii = s.subs{1,1};
9     jj = s.subs{1,2};
10    b = zeros(length(ii),length(jj));
11    for i = 1:length(ii)
12        if ii(i) > a.m, error('??? Index exceeds matrix dimensions.')
13            ; end
14        for j = 1:length(jj)
15            if jj(j) > a.n, error('??? Index exceeds matrix dimensions.
16                '); end
17            b(i,j) = get_index(a,ii(i),jj(j));
18        end
19    end
20 otherwise
21    error('Specify value for i and j as A(i,j)')
22 end
23
24 function aij = get_index(mat,i,j)
25 aij = 0;
26 ii = find(mat.row==i);
27 if ~isempty(ii)
28     jj = find(mat.col(ii)==j);
29     if ~isempty(jj)
30         for k=1:length(jj)
31             aij = aij + mat.a(ii(jj(k)));
32         end
33     end
34 end

```

MATLAB-Funktion: @coomatrix/mtimes.m

```

1 function r = mtimes(A,q)
2 % COOMATRIX/MTIMES Implement A * q for coordinate matrix,
3 % q vector or scalar
4 if isa(A,'coomatrix') && isnumeric(q)
5     if size(q,1) * size(q,2) == 1
6         A.a = q * A.a;
7         r = A;
8         return
9     elseif A.n ~= size(q,1),
10        error('Inner matrix dimensions must agree.')
11    end
12    r = zeros(A.m, size(q,2));
13    for j = 1:size(q,2)
14        for k = 1:length(A.a)
15            r(A.row(k)) = r(A.row(k)) + A.a(k) * q(A.col(k),j);
16        end
17    end
18 else
19    error('Multiplication not defined for these types.')
20 end

```

MATLAB-Beispiel:

Durch die Funktion `subsasgn` funktioniert auch der Zugriff auf einzelne Einträge der Matrix wie aus Matlab bekannt. Da ein der Index (i, j) ggf. nicht eindeutig ist, muss die ganze Matrix durchsucht werden.

```

A(1,2)
ans =
    1.2000
>> A(1,3)
ans =
    0

```

Multiplikation einer Matrix im COO-Format mit einem Vektor lässt sich mittels des Überladens des `*`-Operators mit der Funktion `mtimes.m` wie gehabt realisieren.

```

>> A=coomatrix([1,2],[2,3],[1.2,1.4]);
>> b = A * [2;1;3]
b =
    1.2000
    4.2000

```

Aufgabe C.1.2 Realisieren Sie eine Methode `ctranspose`, welches die Verwendung des Operators Transponiert-Operators `'` in Matlab wie bekannt ermöglicht.

C.2 KOMPRIMIERTE ZEILENSPEICHERUNG (CRS-FORMAT)

Das CRS-Format (Compressed Row Storage-Format, Komprimierte Zeilenspeicherung) ist bereits eine relativ kompakte Speicherform, die aber nicht von einer bestimmten Struktur der Matrix A ausgeht. Sinn dieses Formates ist es, sehr schnell alle Nichtnullelemente einer bestimmten Zeile zu finden.

Es sei $A \in \mathbb{R}^{m \times n}$, dann benötigt das Format 3 Vektoren, 2 der Länge $\text{nnz}(A)$ und 1 mit der Länge $m + 1$:

- ein Vektor „Wert“, in dem die Nichtnullelemente von A zeilenweise, von links nach rechts, stehen;
- einen Index-Vektor „Spaltenindex“, in dem zu jedem Element von „Wert“ der entsprechende Spaltenindex aus der Matrix A eingetragen wird;
- einen Index-Vektor „Zeilenzeiger“, in dem der Beginn jeder Matrixzeile im Feld „Wert“ gespeichert wird. Bei Index $m + 1$ wird der Wert $\text{nnz}(A) + 1$ abgelegt.

Das CRS-Format benötigt gegenüber dem COO-Format deutlich weniger Speicher, kein Eintrag kommt mehrfach vor, ein Index-Vektor ist i. Allg. deutlich kürzer. Auch das Suchen nach einem bestimmten Matrixelement geht hier deutlich schneller. Werden vom Lösungsalgorithmus aber neue Nichtnullelemente erzeugt, dann ist die Einordnung doch komplizierter und auch mit mehr Aufwand verbunden als beim COO-Format.

Beispiel C.2.1

Die Matrix

$$A = \begin{pmatrix} 10 & 0 & 0 & 22 & 31 & 0 \\ -4 & 8 & 0 & 0 & 0 & 0 \\ 0 & 0 & -6 & 9 & 0 & 0 \\ 0 & 0 & 0 & 7 & 0 & 0 \\ 0 & 0 & 5 & 0 & 20 & 11 \\ 12 & 0 & 0 & 14 & 0 & 44 \end{pmatrix}$$

geht über zu

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
Wert a_{ij}	10	22	31	-4	8	-6	9	7	5	20	11	12	14	44
Spaltenindex j	1	4	5	1	2	3	4	7	3	5	6	1	4	6
Zeilenzeiger	1	4	6	8	9	12	15							

MATLAB-Funktion: @crsmatrix/crsmatrix.m

```

1  function p = crsmatrix(i,j,a)
2  %CRSMATRIX Compressed row storage format class constructor.
3  %   A = CRSMATRIX(i,j,a) creates a compressed matrix object from
4  %   the vectors i,j,a, containing the entries of a matrix A.
5  if nargin == 0
6      p.a = []; p.col = []; p.ptr = [];
7      p.m = 0; p.n = 0;
8      p = class(p,'crsmatrix');
9  elseif nargin == 1 && isa(i,'crsmatrix')
10     p = i;
11 elseif nargin == 1 && isnumeric(i)
12     [p.col,row,p.a] = find(i');
13     p.n = max(p.col); p.m = max(row);
14     p.col = p.col'; p.a = p.a';
15     p.ptr = find([row(:);row(end)+1]-[0;row(:)])';
16     p = class(p,'crsmatrix');
17 else
18     p = crsmatrix(sparse(i,j,a));
19 end

```

MATLAB-Funktion: @crsmatrix/display.m

```

1  function display(p)
2  % CRSMATRIX/DISPLAY Command window display of crs-matrix
3  formattedspace, disp([inputname(1),' = ']), formattedspace
4  short = ~(max(p.col)>9999 || p.ptr(end)>9999+1);
5  for k = 2:length(p.ptr)
6      for j = p.ptr(k-1):p.ptr(k)-1
7          ii = num2str(k-1); jj = num2str(p.col(j));
8          if short
9              txt = '          ,          ';
10             txt(5-length(i):5) = ['(',ii];
11             txt(7:7+length(j)) = [jj,')'];
12         else
13             txt = '          ,          ';
14             txt(11-length(i):11) = ['(',ii];
15             txt(13:13+length(j)) = [jj,')'];
16         end
17         disp([txt, sprintf('%.13.4f',p.a(j))])
18     end
19 end
20 formattedspace
21
22 function formattedspace
23 switch get(0,'FormatSpacing')
24     case 'loose'
25         disp(' ');
26 end

```

MATLAB-Funktion: @crsmatrix/mtimes.m

```

1 function r = mtimes(A,q)
2 % CRSMATRIX/MTIMES   Implement A * q for crs-format matrix,
3 % q full matrix of matching size
4 if isa(A,'crsmatrix') && isnumeric(q)
5     if A.n ~= size(q,1),
6         error('Inner matrix dimensions must agree.')

```

MATLAB-Beispiel:

Anlegen einer Matrix im COO-Format funktioniert mit der Anwendung `coomatrix`

```

>> A = crsmatrix([3,2,1,2,3,1], ...
>>                [2,2,4,3,1,1], ...
>>                [1,2,3,4,5,6])
A =
    (1,1)    6.0000
    (1,4)    3.0000
    (2,2)    2.0000
    (2,3)    4.0000
    (3,1)    5.0000
    (3,2)    1.0000

```

Aufgabe C.2.2 Realisieren Sie die Methoden `subsref`, `ctranspose`, welche den Zugriff auf einzelne Einträge und das Transponieren in Matlab wie bekannt ermöglicht.

C.3 MODIFIZIERTE KOMPRIMIERTE ZEILENSPEICHERUNG (MRS-FORMAT)

Beim MRS-Format (**M**odified **C**ompressed **R**ow **S**torage-Format, Modifiziertes CRS-Format) wird nochmals versucht, den ohnehin schon geringen Speicherplatzbedarf des CRS-Formats noch weiter zu senken und direkten Zugriff auf die Diagonalelemente zu haben. denken Sie z.B. an das Jacobi-Verfahren, bei dem man nicht nur eine Matrix-Vektor-Multiplikation mit der Koeffizientenmatrix sondern auch mit der Inversen ihrer Diagonalen benötigt.

Beim MRS-Format geht man davon aus, dass die Hauptdiagonale voll (oder zumindest überdurchschnittlich hoch) besetzt ist, eine Bedingung, die fast alle in der Praxis vorkommenden Matrizen

erfüllen. Speichert man nun die Hauptdiagonale extra ab, so wird keinerlei zusätzliche Koordinateninformation dafür benötigt, denn man weiß ja, daß z.B. das dritte Hauptdiagonalelement in der dritten Zeile und der dritten Spalte steht.

Diese Grundidee wird im MRS-Format folgendermaßen umgesetzt:

Es sei $A \in \mathbb{R}^{m \times n}$. In dem Vektor „Wert“ werden zunächst sämtliche Hauptdiagonalelemente eingetragen, danach wird ein Feld freigelassen. Ab Index $m + 2$ werden nun wie beim CRS-Format zeilenweise die Nichtnullelemente, natürlich ohne die Hauptdiagonalelemente, abgespeichert. Die ersten $m + 1$ Elemente des Vektors „Index“ ersetzen das beim CRS-Format separate Feld „Zeilenzeiger“. Ab Index $m + 2$ werden dann wieder die Spaltenindizes der entsprechenden Matrixelemente des Feldes „Wert“ abgespeichert.

Vor- und Nachteile dieser Speicherform sind ähnlich dem CRS-Format. Das MRS-Format ist relativ beliebt.

Beispiel C.3.1

$$A = \begin{pmatrix} 10 & 0 & 0 & 22 & 31 & 0 \\ -4 & 8 & 0 & 0 & 0 & 0 \\ 0 & 0 & -6 & 9 & 0 & 0 \\ 0 & 0 & 0 & 7 & 0 & 0 \\ 0 & 0 & 5 & 0 & 20 & 11 \\ 12 & 0 & 0 & 14 & 0 & 44 \end{pmatrix}$$

geht über zu

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Wert a_{ij}	10	8	-6	7	20	44		22	31	-4	9	5	11	12	14
Index	8	10	11	12	12	14	16	4	5	1	4	3	5	1	4

MATLAB-Funktion: @mrsmatrix/mrsmatrix.m

```

1 function p = coomatrix(i, j, a)
2 %COOMATRIX Coordinate format class constructor.
3 %   A = COOMATRIX(i, j, a) creates a compressed matrix object from
4 %   the vectors i, j, a, containing the entries of a matrix A.
5 if nargin == 0
6     p.a = []; p.ind = [];
7     p.m = 0; p.n = 0;
8     p = class(p, 'mrsmatrix');
9 elseif nargin == 1 && isa(i, 'mrsmatrix')
10    p = i;
11 elseif nargin == 1 && isnumeric(i)
12    % dirty, at least one off-diagonal entry necessary
13    [col, row, a] = find(i');
14    p.n = max(col); p.m = max(row);
15    ptr = find(col==row);
16    p.a = zeros(1, min(p.m, p.n));
17    p.a(col(ptr)) = a(ptr);

```

```

18 ptr = find(col~=row); row = row(ptr);
19 p.ind = 1+min(p.m,p.n)+find([row(:);row(end)+1]-[0;row(:)])';
20 p.ind = [p.ind,col(ptr)]';
21 p.a = [p.a,0,a(ptr)]';
22 p = class(p,'mrsmatrix');
23 else
24 p = mrsmatrix(sparse(i,j,a));
25 end

```

MATLAB-Funktion: @mrsmatrix/display.m

```

1 function display(p)
2 % MRSMATRIX/DISPLAY Command window display of mrs-matrix
3 formattedspace
4 disp([inputname(1),' = '])
5 formattedspace
6 short = ~(p.m>9999 || p.n>9999);
7 for k = 1:min(p.m,p.n)
8     if p.a(k)
9         printentries(k,k,p.a(k),short)
10    end
11 end
12 for k = 1:p.m
13     for j = p.ind(k):p.ind(k+1)-1
14         printentries(k,p.ind(j),p.a(j),short)
15     end
16 end
17 formattedspace
18
19 function formattedspace
20 switch get(0,'FormatSpacing')
21     case 'loose'
22         disp(' ');
23 end
24
25 function printentries(ii,jj,aij,short)
26 i = num2str(ii); j = num2str(jj);
27 if short
28     txt = '          ,          ';
29     txt(5-length(i):5) = ['(',i];
30     txt(7:7+length(j)) = [j,')']';
31 else
32     txt = '          ,          ';
33     txt(11-length(i):11) = ['(',i];
34     txt(13:13+length(j)) = [j,')']';
35 end
36 disp([txt,sprintf('%13.4f',aij)])

```


MATLAB-Funktion: @mrsmatrix/mtimes.m

```

1 function r = mtimes(A,q)
2 % MRSMATRIX/MTIMES   Implement A * q for mrs-format,
3 % q vector or matrix of matching size
4 if isa(A,'mrsmatrix') && isnumeric(q)
5     if A.n ~= size(q,1),
6         error('Inner matrix dimensions must agree.')

```

MATLAB-Beispiel:

Anlegen einer Matrix im MRS-Format funktioniert z.B. mit der Anweisung `mrsmatrix`. Auch die Multiplikation lässt sich kanonisch realisieren.

```

>> A = mrsmatrix([3,2,1,2,3,1], ...
>>                [2,2,4,3,1,1], ...
>>                [1,2,3,4,5,6])
A =
(1,1)          6.0000
(2,2)          2.0000
(1,4)          3.0000
(2,3)          4.0000
(3,1)          5.0000
(3,2)          1.0000
>> x = [0,1;1,1;2,1;2 0];
>> A*x
ans =
     6     6
    10     6
     1     6

```

Aufgabe C.3.2 Realisieren Sie die Methoden `subsref`, `ctranspose`, welche den Zugriff auf einzelne Einträge und das Transponieren in Matlab wie bekannt ermöglicht.

Aufgabe C.3.3 Schreiben Sie eine Funktion, die das Jacobi-Verfahren zur Lösung von linearen Gleichungssystemen für in CRS-Format gespeicherte Matrizen umsetzt. (Tipp: Neben $A \cdot x$ benötigen Sie noch eine Routine, die $\text{diag}(\text{diag}(A)) \setminus x$ effizient umsetzt.

C.4 HARWELL-BOEING-FORMAT (CCS-FORMAT)

Das Harwell-Boeing-Format oder auch CCS (Compressed-Column-Storage-Format, Komprimierte Spaltenspeicherung) genanntes Format entspricht im wesentlichen dem CRS-Format, nur dass man hier nicht schnell auf eine Zeile, sondern eine komplette Spalte zugreifen möchte. Die Speicherung zeilenweise von links nach rechts wird durch spaltenweise Speicherung von oben nach unten ersetzt.

Literaturverzeichnis

- [Deuffhard/Hohmann] P. DEUFLHARD, A. HOHMANN, Numerische Mathematik 1, de Gruyter-Verlag, Berlin, 1993.
- [Fischer] G. FISCHER, Lineare Algebra, Vieweg-Verlag, Wiesbaden, 1989.
- [Analysis I] S. FUNKEN, Skript zur Vorlesung „Analysis I“, gehalten im Wintersemester 07/08 an der Universität Ulm.
- [Analysis II] S. FUNKEN, Skript zur Vorlesung „Analysis II“, gehalten im Sommersemester 2008 an der Universität Ulm.
- [Discroll/Maki] T. A. DISCROLL, K. L. MAKI: Searching for Rare Growth Factors Using Multicanonical Monte Carlo Methods. SIAM Review. Vol. 49, No. 4, pp. 673-692. 2008.
- [Golub/Loan] G. H. GOLUB, C. F. VAN LOAN, Matrix Computations, 3. ed., Hopkins Univ. Press, 1996.
- [Hackbusch] W. HACKBUSCH, Iterative Lösung großer schwachbesetzter Gleichungssysteme, 2. Auflage, Teubner-Verlag, Stuttgart, 1993.
- [Hämmerlin/Hoffmann] G. HÄMMERLIN, K.-H. HOFFMANN, Numerische Mathematik, 4. Auflage, Springer-Verlag, Berlin, 1994.
- [Hanke] M. HANKE-BOURGEOIS, Grundlagen der Numerischen Mathematik und des wissenschaftlichen Rechnens, 1. Auflage, Teubner-Verlag, Stuttgart, 2002.
- [Higham] N. J. HIGHAM, How accurate is Gaussian elimination? Numerical Analysis 1989, Proceedings of the 13th Dundee Conference, Vol. 228 of Pitman Research Notes in Mathematics, 137–154.
- [Knuth] D.E. KNUTH The Art of Computer Programming, Band 2, Addison-Wesley, 3. Auflage, 1998.
- [Niethammer] W. NIETHAMMER, Relaxation bei nichtsymmetrischen Matrizen. Math. Zeitschr. **85** 319-327, 1964.
- [Plato] R. PLATO, Numerische Mathematik kompakt, Vieweg-Verlag.
- [Quateroni et. al.] A. QUATERONI, R. SACCO, F. SALERI, Numerische Mathematik, Band 1 & 2, Springer-Verlag, Berlin, 2002.
- [Schwarz] H. R. SCHWARZ, Numerische Mathematik, 4. Auflage, Teubner-Verlag, Stuttgart, 1997.
- [Stör/Bulirsch] J. STÖR, R. BULIRSCH, Numerische Mathematik, Band 1 & 2, Springer-Verlag, Berlin, 1994.
- [Törnig/Spellucci] W. TÖRNIG, P. SPELLUCCI, Numerische Mathematik für Ingenieure und Physiker, Band 1 & 2, Springer-Verlag, Berlin, 1988.
- [Wilkinson65] J. H. WILKINSON, The Algebraic Eigenvalue Problem, Oxford University Press, 1965.

- [Wilkinson69] J. H. WILKINSON, Rundungsfehler, Springer-Verlag, Berlin, Heidelberg, New York, 1969.
- [Wille] D. WILLE, Repetitorium der Linearen Algebra, Teil 1, 1. Auflage, Feldmann-Verlag, Springe, 1989.

Index

- A-orthogonal, 71
- Abbruch-Kriterien, 64
- absoluter Fehler, 47
- absoluter/relativer Fehler, 37
- Abstiegsrichtung, 66
- Äquivalenzoperationen, 8
- Arithmetik
 - Hochgenauigkeitsarithmetik, 45
 - Intervallarithmetik, 44
- Armijo, Schrittweitenregel, 67
- Ausgleichsgerade, 81
- Auslöschung, 39

- b*-adisch, 29
- b*-adischer Bruch, 29
- Bandbreite, 21
- Bandgleichungen, 21
- Banker's Rule, 36
- Bruch, *b*-adischer, 29

- Cauchy-Schwarz-Ungleichung, 97
- cg-Verfahren, 71, 74
- Cholesky, 17
 - Verfahren, 19, 22
 - Zerlegung, 20
 - Zerlegung, rationale, 18
- Cramersche Regel, 1

- Darstellung, Gleitpunkt-, 32
- Darstellung, Leibnizsche, 1
- Daumenregel zur Genauigkeit, 51
- diagonaldominant, 14
- Diagonalmatrix, 4
- Drei-Term-Rekursion, 77
- Dreiecksmatrix, 5
- Dreieckszerlegung, Gaussche, 9

- einfache Realisierbarkeit, 53
- Einzelschrittverfahren, 53
- Eliminationsmethode, Gaussche, 7
- Eliminationsschritt, 8
- Energienorm, 67

- Fehler, 47
 - absoluter, 47
 - absoluter/relativer, 37
 - relativer, 47
- Fehlerarten, 27
- Festpunktzahl], 31
- Frobenius-Norm, 98

- Gaußsches-Eliminations-Verfahren, 9
- Gauß-Seidel, 53
- Gaus-Elimination mit Spaltenpivotstrategie, 12
- Gaussche Dreieckszerlegung, 9
- Gaussche Eliminationsmethode, 7
- Gaußschen Normalengleichung, 82
- Genauigkeit, Daumenregel zur, 51
- Gesamtschrittverfahren, 53
- Gestaffelte Systeme, 4
- Givens-Rotationen, 85
- Glättungseigenschaften, 64
- Gleitkommaarithmetik, 40
- Gleitkommaoperation, 2
- Gleitkommazahl, 27
- Gleitpunkt-Darstellung, 32
- Gradientenverfahren, 66
- Gradientenverfahren, Konjugiertes, 74
- Gradientenverfahren, Präkonditioniertes Konjugiertes, 78
- Gradientenverfahren:, 68
- Guard Digit, 40

- Hochgenauigkeitsarithmetik, 45
- Holder-Ungleichung, 97
- Householder-Spiegelungen, 90
- Householder-Transformation, 90
- Householder-Vektor, 90

- Interpolationspolynom, 23
- Intervallarithmetik, 44
- irreduzibel, 58
- Iterationsverfahren, 53

- Jacobi, 53
- Jacobi-Verfahren, 57
- Jordansche Normalform, 56

- Kantorowitsch-Ungleichung, 69
- kaufmännische Rundung, 36

- Kondition, 27, 41
- Konditionszahl, 48
- Konditionszahlen, 42
- konjugiert, 71
- Konjugiertes Gradientenverfahren, 74
- Konsistenz, 43
- Konvergenzeigenschaft, 53
- Konvergenzkriterium, 56
- Konvergenzrate, 54

- Laplacescher Entwicklungssatz, 2
- Leibnizsche Darstellung, 1
- Liniensuche, 66
- LR*-Zerlegung, 9

- Maschinengenauigkeit, 38
- mathematisch unverzerrte Rundung, 36
- MATLAB, 101
- Matrix
 - Diagonalmatrix, 4
 - Dreiecksmatrix, 5
 - Elementarmatrix, 56
 - Jordanmatrix, 56
 - positiv definite, 17
 - Rotationsmatrizen, 84
 - unipotente, 9
 - Vandermond, 23
- Matrixnorm, 98
- Methode der kleinsten Quadrate, 81

- Nachiteration, 17
- Neumann-Reihe, 50
- Newtondarstellung, 23
- Normalengleichung, Gaußschen, 82
- Normalisierung, 29

- Operationen
 - Aquivalenzoperationen, 8
 - Gleitkommaoperation, 2
 - Rechenoperationen, 2
- orthogonal, 83

- p*-adisch, 29
- pcg-Verfahren, 78
- Permutation, 11
- Pivot-Strategien, 11
- Pivotelement, 8
- Polynominterpolation, 24
- positiv definite Matrix, 17
- Präkonditioniertes Konjugiertes Gradientenverfahren, 78

- QR-Zerlegung, 83

- Rückwärtsanalyse, 45
- rationale Cholesky-Zerlegung, 18
- Realisierbarkeit, einfache Realisierbarkeit, 53
- Rechenoperationen, 2
- Regel, Cramersche, 1
- Reihe
 - Neumann-, 50
- relativer Fehler, 47
- Residuum, 17, 47, 72
- Richardson, 53
- Rosenberg, 63
- Rotationsmatrizen, 84
- Round to even, 36
- Rundung, 36
 - kaufmännische, 36
 - mathematisch unverzerrte, 36
 - Standardrundung, 36

- Satz
 - Laplacescher Entwicklungssatz, 2
- Schrittweitenregel von Armijo, 67
- Schutzziffer, 40
- schwachbesetzt, 55
- Schwaches Zeilensummenkriterium, 59
- Spaltenpivotisierung, 14
- Spaltenpivotstrategie, 12
- Spaltenpivotstrategie, Gaus-Elimination mit, 12
- Spaltensummenkriterium, starkes , 57
- Spaltensummennorm, 99
- Spektralradius, 55
- Stabilität, 27, 43
- Standardfehler, 43
- Starkes Spaltensummenkriterium, 57
- Starkes Zeilensummenkriterium, 57
- starkes Zeilensummenkriterium, 60
- Strategien, Pivot-, 11
- Submultiplikativität, 98
- Substitution, Vorwärts- Rückwärts-, 6
- Systeme, Gestaffelte, 4

- Tschebyscheff-Polynome, 75, 76

- Überrelaxationsverfahren, 64
- Ungleichung, Cauchy-Schwarz-, 97
- Ungleichung, Holder-, 97
- unipotente Matrix, 9
- Unterrelaxationsverfahren, 64

- Vandermond-Matrizen, 23
- Vektornorm, 97
- Vektornorm, vertragliche, 100
- Verfahren

Überrelaxationsverfahren, 64
cg-Verfahren, 71, 74
Cholesky, 19
Cholesky-, 22
Einzelschrittverfahren, 53
Gaußsches-Eliminations-, 9
Gesamtschrittverfahren, 53
Gradientenverfahren, 66
Gradientenverfahren:, 68
Iterationsverfahren, 53
Jacobi, 57
Konjugiertes Gradientenverfahren, 74
pcg-Verfahren, 78
Präkonditioniertes Konjugiertes Gradientenverfahren, 78
Unterrelaxationsverfahren, 64
vertragliche Vektornorm, 100
von Stein, 63
Vorkonditionierung, 70
Vorwärts- Rückwärtssubstitution, 6
Vorwärtsanalyse, 43

Zahlendarstellung, 27
Zahlenformat, 33
Zeilenpivotisierung, 14
Zeilensummenkriterium, schwaches, 59
Zeilensummenkriterium, starkes, 57, 60
Zeilensummennorm, 99
zerlegbar, 58
Zerlegung, 14
Zerlegung, Cholesky-, 20
Zerlegung, Cholesky-, rationale, 18