

Numerik I

Vorlesungsskript

TU Chemnitz
Fakultät für Mathematik
Volker Mehrmann



Literatur:

- (a) J.H. Wilkinson
Rundungsfehler, Springer
- (b) G.H. Golub/C.F. Van Loan
Matrix Computations, Johns Hopkins University Press
- (c) A. Frommer
Lösung linearer Gleichungssysteme auf Parallelrechnern, Vieweg
- (d) J. Stoer
Numerische Mathematik I, Springer
- e) H.R. Schwarz
Numerische Mathematik, Teubner

u.v.a.m.

Skript und weitere Informationen

<http://www.tu-chemnitz.de/~mehrmann/nla/>

Inhaltsverzeichnis

1 Grundlagen	1
1.1 Was ist numerische Mathematik?	1
1.2 Zahlendarstellung auf dem Digitalrechner	2
1.2.1 Rundungsfehler (Reduktionsfehler)	4
1.3 Fehlerfortpflanzung	7
1.4 Moderne Rechnerarchitekturen	11
1.5 Die '·' Notation, flops und Normen	13
1.6 Übungen zu Kapitel 1	16
2 Lösung linearer Gleichungssysteme	18
2.1 Gauß–Elimination	18
2.1.1 Lösung von Dreieckssystemen	18
2.1.2 LR –Zerlegung	19
2.1.3 Fehleranalyse der Gauß–Elimination	25
2.1.4 Partielle Pivotisierung (Spaltenpivotisierung)	29
2.1.5 Vollständige Pivotisierung	32
2.1.6 Abschätzung der Genauigkeit	33
2.1.7 Skalierung	34
2.1.8 Iterative Verbesserung	35
2.1.9 Übungen zu Kapitel 2.1	35
2.2 Der Cholesky–Algorithmus, Bandmatrizen	39
2.2.1 Verallgemeinerung für nicht positiv definite Systeme.	41
2.2.2 Bandsysteme	41
2.2.3 Übungen zu Kapitel 2.2	43
2.3 Householder und Gram–Schmidt–Orthogonalisierung	45

2.3.1	Householder Matrizen(Spiegelungen)	45
2.3.2	Berechnung der Householder-Matrix	46
2.3.3	Anwendung von Householder-Matrizen	46
2.3.4	Die QR -Zerlegung.	47
2.3.5	Eigenschaften der QR -Zerlegung	49
2.3.6	Gram-Schmidt Orthogonalisierung	49
2.3.7	Übungen zu Kapitel 2.3	50
2.4	Ausgleichsprobleme	54
2.5	Iterative Verfahren	56
2.5.1	Splitting-Verfahren	57
2.5.2	Konvergenzbeschleunigung	59
2.5.3	Konjugierte Gradienten	59
2.5.4	Übungen zu Kapitel 2.5	61
3	Lösung nichtlinearer Gleichungssysteme	64
3.1	Fixpunktverfahren	64
3.2	Das Newtonverfahren	68
3.2.1	Anwendung auf modifiziertes Newtonverfahren	75
3.2.2	Praktische Realisierung des modifizierten Newton-Verfahrens	77
3.3	Nullstellenbestimmung für Polynome	78
3.4	Übungen zu Kapitel 3	80
4	Interpolation	84
4.1	Polynominterpolation	85
4.2	Trigonometrische Interpolation	90
4.3	Spline Interpolation	95
4.3.1	Allgemeine Eigenschaften	95
4.3.2	Berechnung der Splines	97
4.3.3	Fehlerabschätzung und Konvergenz bei Interpolation mit Splines	100
4.4	Übungen zu Kapitel 4	103
5	Numerische Integration	108
5.1	Newton-Cotes Formeln	108
5.2	Extrapolation	111

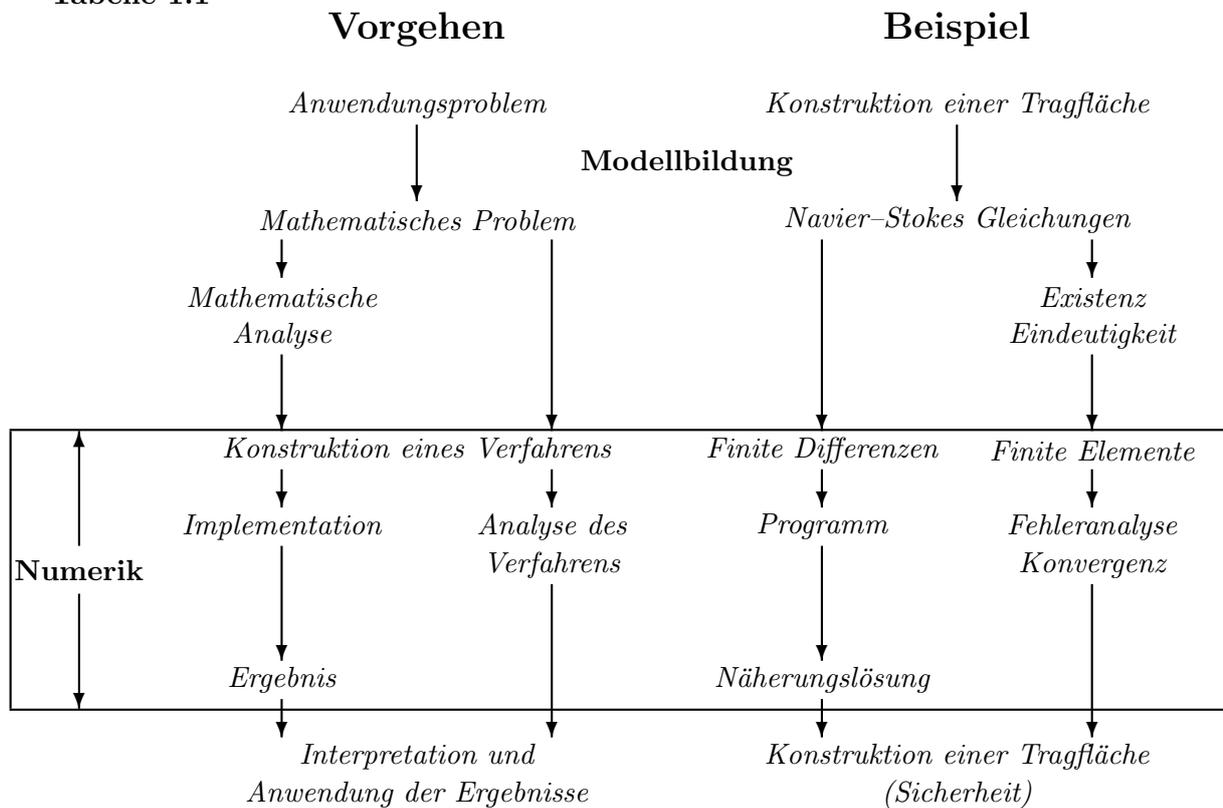
5.3	Übungen zu Kapitel 5	116
6	Eigenwertprobleme	118
6.1	Grundlagen (Wiederholung)	118
6.2	Der QR -Algorithmus für unsymmetrische Matrizen	121
6.3	Der QR -Algorithmus für symmetrische Matrizen	127
6.4	Die Singulärwertzerlegung	127
6.4.1	Die Berechnung der Singulärwertzerlegung	128
6.5	Der Bisektionsverfahren zur Berechnung der Eigenwerte symmetrischer tridia- gonaler Matrizen	135
6.6	Übungen zu Kapitel 6	136

Kapitel 1

Grundlagen

1.1 Was ist numerische Mathematik?

Tabelle 1.1



Gesichtspunkte:

- Genauigkeit des Verfahrens
- Verlässlichkeit der Ergebnisse
- Effizienz des Verfahrens
- Laufzeiten
- Black-box, special purpose,...

1.2 Zahlendarstellung auf dem Digitalrechner

Die Zahlendarstellung auf dem Rechner verwendet den folgenden Satz:

Satz 1.2 (*p*-adische Entwicklung)

Ist $r \in \mathbb{R}$, $p \in \mathbb{N}$, $p > 1$, so gibt es ein eindeutiges $j \in \{0, 1\}$, ein eindeutiges $l \in \mathbb{Z}$ und für alle $k \in \mathbb{Z}$ mit $k \leq l$, eindeutige $\gamma_k \in \mathbb{Z}$, so dass

$$r = (-1)^j \sum_{k=-\infty}^l \gamma_k p^k \quad (1.3)$$

wobei für $r \neq 0$, $\gamma_l \neq 0$ und $j = l = 0$ für $r = 0$, $0 \leq \gamma_k < p$ für alle $k \leq l$, $\gamma_k < p - 1$ für unendliche viele $k \leq l$.

Bemerkung 1.4 Die letzte Bedingung schließt Zahlen wie $3.99999\bar{9}$ aus.

Beispiel 1.5 $r = 18.5$ $p = 2$

$$\begin{aligned} r &= (-1)^0 (1 * 2^{-1} + 0 * 2^0 + 1 * 2^1 + 0 * 2^2 + 0 * 2^3 + 1 * 2^4) \\ &= 10010.1 \text{ Dualzahldarstellung oder Binärdarstellung} \\ p &= 16 \\ r &= (-1)^0 (8 * 16^{-1} + 2 * 16^0 + 1 * 16) \\ &= 12.8 \text{ Hexadezimaldarstellung} \\ &\quad \text{Ziffern (0123456789ABCDEF)}. \end{aligned}$$

→ Übung 1.43.

Die Konvertierung zwischen verschiedenen Darstellungen geschieht durch Division durch p mit Rest für ganzzahligen Anteil und Multiplikation mit p für Nachpunktanteil. Nach Satz 1.2 ist

$$r = \left[\underbrace{(-1)^j \sum_{k=-\infty}^l (\gamma_k p^{k-l-1})}_a \right] \cdot \underbrace{p^{l+1}}_b \quad (1.6)$$

wobei $\frac{1}{p} < |a| < 1$, da $\gamma_l \neq 0$.

Definition 1.7 Die Darstellung $r = a \cdot p^b$ heisst normalisierte Gleitpunktdarstellung von r bezüglich p . a heisst Mantisse, b der Exponent von r .

$$a = V_M \sum_{i=1}^{\infty} \alpha_i p^{-i} \quad \text{mit } \alpha_1 \neq 0$$

$$b = V_E \sum_{i=1}^t \beta_i p^{t-i}, \quad V_M, V_E \text{ Vorzeichen.}$$

! Auf einem Rechner ist für die Darstellung (im allgemeinen, in normierter Sprache) nur eine feste Anzahl von n Stellen möglich. Diese n Stellen werden aufgeteilt in l Stellen für die Mantisse und m Stellen für den Exponenten. Die Zahl

$$V_M(\alpha_1 p^{-1} + \alpha_2 p^{-2} + \dots + \alpha_l p^{-l}) p^{V_E(\beta_1 p^{m-1} + \dots + \beta_m p^0)} \quad (1.8)$$

mit $\alpha_1 \neq 0$ wird durch die Angabe $V_M \alpha_1 \alpha_2 \dots \alpha_l V_E \beta_1 \dots \beta_m$ codiert.

Beispiel 1.9

$$r_1 = -1234.567890, r_2 = 0.01234567890$$

$$p = 10, l = 9, m = 2$$

normalisierte Gleitpunktdarstellung

$$r_1 = -0.123456789 * 10^4$$

$$= -123456789 + 04$$

$$r_2 = -0.123456789 * 10^{-1}$$

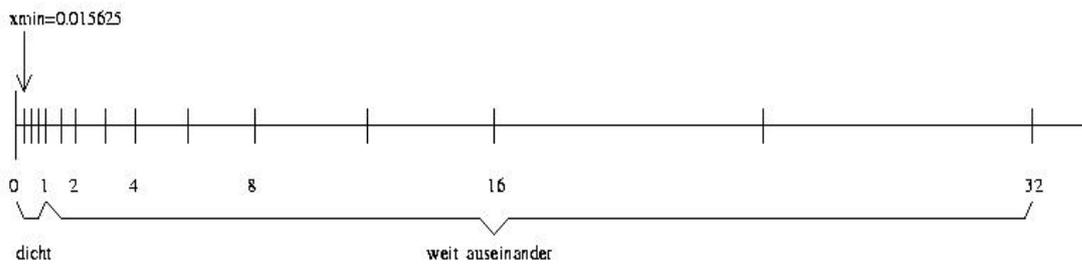
$$= -123456789 - 01$$

Definition 1.10 $\mathbb{M}(p, l, m)$ bezeichnet die Menge der in normalisierter Gleitpunktdarstellung codierbaren (darstellbaren) Zahlen, (auch Maschinenzahlen genannt) mit Basis p , l -Stellen für die Mantisse und m -Stellen für den Exponenten (jeweils zur Basis p).

Beispiel 1.11 $\mathbb{M}(2, 3, 3)$, d.h. $p = 2, l = m = 3$

verfügbare Mantissen	verfügbare Exponenten
	$\pm 000 \hat{=} \pm 0 \quad 1, 1$
$\pm 100 \hat{=} \pm \frac{1}{2}$	$\pm 001 \hat{=} \pm 1 \quad 2, \frac{1}{2}$
	$\pm 010 \hat{=} \pm 2 \quad 4, \frac{1}{4}$
$\pm 101 \hat{=} \pm \frac{5}{8}$	$\pm 011 \hat{=} \pm 3 \quad 8, \frac{1}{8}$
	$\pm 100 \hat{=} \pm 4 \quad 16, \frac{1}{16}$
$\pm 110 \hat{=} \pm \frac{3}{4}$	$\pm 101 \hat{=} \pm 5 \quad 32, \frac{1}{32}$
	$\pm 110 \hat{=} \pm 6 \quad 64, \frac{1}{64}$
$\pm 111 \hat{=} \pm \frac{7}{8}$	$\pm 111 \hat{=} \pm 7 \quad 128, \frac{1}{128}$
kleinste positive Zahl	$x_{\min} = \frac{1}{2} 2^{-7} = 2^{-8}$
größte positive Zahl	$x_{\max} = \frac{7}{8} 2^7 = 112$

Maschinenzahlen sind wie folgt angeordnet



1.2.1 Rundungsfehler (Reduktionsfehler)

Von nun an sei p entweder 10 oder 2^k . Um eine beliebige reelle Zahl darzustellen wird sie gerundet in die Menge der Maschinenzahlen abgebildet,

$$\gamma : \mathbb{R} \longrightarrow \mathbb{M}(p, l, m).$$

Sei $x = \pm (\sum_{i=1}^{\infty} \alpha_i p^{-i}) \cdot p^t, \alpha_1 \neq 0, x \in Z = [x_{\min}, x_{\max}] \cup [-x_{\max}, -x_{\min}]$

Übliche Rundung auf l Stellen

$$\gamma(x) = \begin{cases} \pm \left(\sum_{i=1}^l \alpha_i p^{-i} \right) \cdot p^t & \text{falls } \alpha_{l+1} < \frac{p}{2} \\ \pm \left(\sum_{i=1}^l (\alpha_i p^{-i}) + p^{-l} \right) \cdot p^t & \text{falls } \alpha_{l+1} \geq \frac{p}{2}. \end{cases} \quad (1.12)$$

Bei Zahlen ausserhalb von Z gibt es rechnerabhängige Vorgehensweisen z.B.

→ Übung 1.45

$$|x| < x_{\min} \begin{cases} \gamma(x) = 0 \\ \gamma(x) = \text{sign}(x)x_{\min} \end{cases} \quad \text{'underflow'}$$

$$|x| > x_{\max} \begin{cases} \gamma(x) = \text{sign}(x)x_{\max} \text{ mit Meldung 'overflow'} \\ \gamma(x) = \text{sign}(x)x_{\max} \end{cases} \quad \text{'overflow'}$$

→ Übung 1.45: bestimme x_{\min}, x_{\max}

Absoluter Rundungsfehler

$$|\gamma(x) - x| \leq \frac{p^{-l}}{2} p^t \quad (1.13)$$

Beweis:

abrunden:

$$\begin{aligned} |\gamma(x) - x| &= \left| \left(\sum_{i=1}^{\infty} \alpha_i p^{-i} \right) p^t - \left(\sum_{i=1}^l \alpha_i p^{-i} \right) p^t \right| \\ &= p^t \cdot \left| \sum_{i=l+1}^{\infty} \alpha_i p^{-i} \right| = p^t p^{-(l+1)} \left| \underbrace{\sum_{i=0}^{\infty} \alpha_{l+i+1} p^{-i}}_{\leq \frac{p}{2}} \right| \\ &\leq p^t \frac{p^{-l}}{2}. \end{aligned}$$

aufzurunden:

$$\begin{aligned}
 |\gamma(x) - x| &= \left| \left(\sum_{i=1}^{\infty} \alpha_i p^{-i} \right) p^t - \left(\sum_{i=1}^l (\alpha_i p^{-i}) + p^{-l} \right) p^t \right| \\
 &= \left| \sum_{i=l+1}^{\infty} \alpha_i p^{-i} - p^{-l} \right| p^t \\
 &= p^t \left| -p^{-l} + \alpha_{l+1} p^{-(l+1)} + \dots \right| \\
 &= p^t p^{-l} \left| -1 + \alpha_{l+1} p^{-1} + \alpha_{l+2} p^{-2} + \dots \right| \\
 &\leq p^t p^{-l} \left(-\frac{p}{2} p^{-1} + 1 \right) = p^t \frac{p^{-l}}{2}.
 \end{aligned}$$

□

Relativer Rundungsfehler

$$\frac{|\gamma(x) - x|}{|x|} \leq \frac{\frac{p^{-l}}{2} p^t}{M_x p^t}$$

wobei die Mantisse M_x von x die Ungleichung $M_x \geq \frac{1}{p}$ erfüllt. Bei $M_x = \frac{1}{p}$, d.h. $x = \frac{p^t}{p}$ gibt es keinen Fehler. Also

$$\left| \frac{\gamma(x) - x}{x} \right| \geq \frac{\frac{p^{-l}}{2}}{\frac{1}{p}} = \frac{p}{2} p^{-l} \quad (1.14)$$

Die Zahl

$$\text{eps} := \frac{p}{2} p^{-l} =: (\text{macheps}) \quad (1.15)$$

heißt Maschinengenauigkeit (roundoff unit).

→ Übung: 1.45 Bestimme eps für Deine Maschine.

Es gilt immer $\text{eps} = \min_{\delta \in [x_{\min}, x_{\max}]} |\gamma(1 + \delta)| > 1$.

$$1 = 1 * p^{-1} p^1.$$

Bemerkung 1.16 *Es gilt also stets*

$$\begin{aligned}
 \gamma(x) &= x \cdot (1 + \varepsilon) \text{ mit } |\varepsilon| \leq \text{eps}, \\
 \varepsilon &= \frac{\gamma(x) - x}{x}.
 \end{aligned}$$

Beispiel 1.17

0.2 dezimal $l = 6, p = 2$

ergibt $0.\overline{0011}$ in dualer Darstellung

normalisiert $0.11\overline{0011} * 2^{-2}$

Rundung $l = 6$ ergibt $0.110011 * 2^{-2}$

Rückkonvertiert $\frac{51}{256} = 0.19921875$

Rundungsfehler entstehen beim Einlesen, Speichern, Konvertieren, Rechnen.

Man hat auf dem Rechner nur eine Pseudoarithmetik.

M ist nicht abgeschlossen bzgl. $+, -, *, \div$

Beispiel 1.18 $\mathbb{M}(10, 5, 1)$

$$\begin{aligned}
x &= +25684 + 1 \hat{=} 2.5684 \\
y &= +32791 - 2 \hat{=} 0.0032791 \\
x + y &= \underbrace{2.5716}_{5}781 \notin \mathbb{M}, x * y = 0.00 \underbrace{842204}_{5}4044 \notin \mathbb{M} \\
x/y &= \underbrace{783.26}_{5}37004 \notin \mathbb{M}.
\end{aligned}$$

Die übliche Arithmetik in \mathbb{R} muss so realisiert werden, dass das Ergebnis wieder in \mathbb{M} ist.

→ Übung

Im allgemeinen gilt:

Das Ergebnis der Pseudooperation $\nabla, \nabla \in \{+, -, *, \div\}$ ist festgelegt durch

$$gl(x \nabla y) \equiv x \nabla y = \gamma(x \nabla y), \text{ für } x, y \in \mathbb{M} \quad (1.19)$$

Hardwaremäßig wird üblicherweise mit längerer Mantisse gearbeitet als im Ergebnis, dies normalisiert und dann gerundet.

! Vorsicht: Nicht auf allen Rechnern, man kann böse Überraschungen erleben. (z.B. auf Rechnern die nicht mit IEEE-Standard-Arithmetik arbeiten z.B. Cray.)

Unter der Annahme (1.19) gilt:

Falls $|x \nabla y|$ in $[x_{\min}, x_{\max}]$ liegt, so gilt für den relativen Fehler

$$\left| \frac{x \nabla y - x \nabla y}{x \nabla y} \right| = \left| \frac{\gamma(x \nabla y) - x \nabla y}{x \nabla y} \right| < \text{eps} \quad (1.20)$$

In \mathbb{R} gelten Assoziativ-, Kommutativ- und Distributivgesetz, in der Rechnerarithmetik in \mathbb{M} gilt im allgemeinen nur die Kommutativität für Addition und Multiplikation.

Beispiel 1.21 $\mathbb{M}(10, 5, 1)$ Assoziativität

$$\begin{aligned}
&0.98765 + 0.012424 - 0.0065432 = \underline{0.9935308} \\
&0.98765 \oplus (0.012424 \ominus 0.0065432) = \\
&0.98765 \oplus 0.00588(08) = \underline{0.99353(08)} \\
&(0.98765 \oplus 0.012424) \ominus 0.0065432 = \\
&\underbrace{1.000074}_{1.0001} \ominus 0.0065432 = \underline{0.99351}
\end{aligned}$$

Distributivität

$$\begin{aligned}
&(4.2832 - 4.2821) * 5.7632 = \underline{0.00633952} \\
&(4.2832 \ominus 4.2821) \otimes 5.7632 = \\
&0.0011 \otimes 5.7632 = \underline{0.0063395(2)} \\
&(4.2832 \otimes 5.7632) \ominus (4.2821 \otimes 5.7632) = \\
&\quad \underbrace{24.684(93824)}_{24.685 \ominus 24.679 = 0.0060000} \ominus 24.678(59872)
\end{aligned}$$

! Folgerung: Mathematisch äquivalente Algorithmen zur Auswertung eines rationalen Ausdrucks können bei Durchführung auf dem Rechner zu wesentlich verschiedenen Ergebnissen führen, selbst wenn die Eingabedaten Maschinenzahlen sind. Nichtrationale Operationen wie

$\sqrt{}$, \sin , \exp sind (nicht immer gut) softwaremäßig realisiert. Auch hierfür gilt im allgemeinen: Das Maschinenergebnis ist Rundung des exakten Ergebnisses.

→ Übung.

Auf einigen Maschinen ist selbst \div nicht fest realisiert (z.B. Cray) oder man hat die Wahl zwischen (billig und schlecht oder teuer und gut).

1.3 Fehlerfortpflanzung

Unglücklicherweise pflanzen sich einmal bereits gemachte Fehler (z.B. bei Rundung) auch noch fort. Seien x, y mit Fehlern $\Delta x, \Delta y$ behaftet $\left(\left|\frac{\Delta x}{x}\right|, \left|\frac{\Delta y}{y}\right| \ll 1\right)$. Was passiert bei exakter Durchführung einer elementaren Operation mit diesen Fehlern, d.h. für $\nabla \in \{+, -, *, \div\}$ sei $x \nabla y \in \mathbb{M}$

Fortpflanzung des relativen Fehlers:

$$(x + \Delta x) \nabla (y + \Delta y) = x \nabla y + \underbrace{\Delta(x \nabla y)}_{\text{abs. Fehler des Ergebnisses}}$$

$$\pm : x \pm y + \Delta x \pm \Delta y :$$

$$\frac{\Delta(x \pm y)}{x \pm y} = \frac{\Delta x \pm \Delta y}{x \pm y} = \frac{x}{x \pm y} \frac{\Delta x}{x} \pm \frac{y}{x \pm y} \frac{\Delta y}{y} \quad (1.22)$$

Bemerkung 1.23 Ist $|x \pm y|$ sehr klein gegenüber $|x|$ oder $|y|$, so kann der relative Fehler ausserordentlich verstärkt werden, z.B. wenn zwei annähernd gleich große Zahlen subtrahiert werden. Dieser Effekt heisst *Auslöschung*. Man muss bei der Aufstellung der Algorithmen darauf achten, dass dies soweit wie möglich vermieden wird.

→ Übung

$$* : xy + \underbrace{x \Delta y + y \Delta x + \Delta y \Delta x}_{\Delta(x*y)} :$$

$$\frac{\Delta(x * y)}{x * y} = \frac{\Delta y}{y} + \frac{\Delta x}{x} + \frac{\Delta y \Delta x}{xy} \approx \frac{\Delta y}{y} + \frac{\Delta x}{x} \quad (1.24)$$

Die relativen Fehler addieren sich also im wesentlichen.

$$\div : \frac{x + \Delta x}{y + \Delta y} = \frac{x}{y} + \frac{y \Delta x - x \Delta y}{y(y + \Delta y)} :$$

$$\frac{\Delta(x \div y)}{x \div y} = \frac{1}{x \div y} \left(\frac{y \Delta x - x \Delta y}{y(y + \Delta y)} \right) \quad (1.25)$$

$$= \frac{y \Delta x - x \Delta y}{x(y + \Delta y)} = \frac{\Delta x}{x} \frac{y}{y + \Delta y} - \frac{\Delta y}{y + \Delta y}$$

$$\approx \frac{\Delta x}{x} - \frac{\Delta y}{y}$$

Die relativen Fehler subtrahieren sich im wesentlichen. Bei $*$, \div tritt also keine wesentliche Verstärkung des Fehlers ein.

Beispiel 1.26 Seien a_1, \dots, a_n Maschinenzahlen. Berechne $\sum_{i=1}^n a_i$ auf Rechner mit Maschinengenauigkeit eps . Verfolge die Fehler und deren Fortpflanzung.

→ Übung 1.44

$$\left| gl \left(\sum_{i=1}^n a_i \right) - \sum_{i=1}^n a_i \right| \leq \text{eps} \left(\sum_{i=1}^n (n - i + 1) |a_i| \right)$$

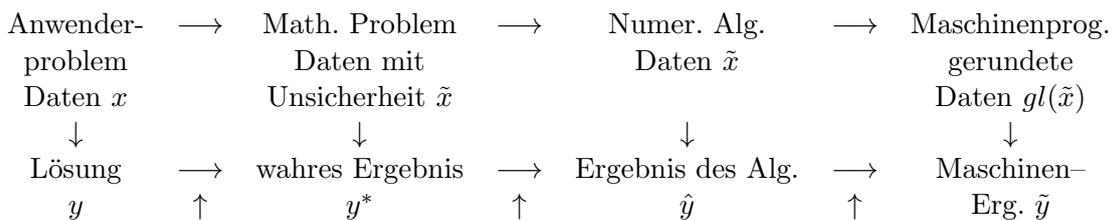
Wie beurteilt man nun numerische Verfahren?

- a) Genauigkeit und Verlässlichkeit des Ergebnisses.
- b) Zeitbedarf, Speicherbedarf, Programmlänge, Kosten, (neuerdings Vektorisierbarkeit + Parallelisierbarkeit).
- b) Hängt von der Architektur ab, bzw. von der Systemumgebung, i.a. kann man dieses gut abschätzen oder testen. Siehe auch 1.5.
- a) Die Verlässlichkeit gibt an, wie ein Versagen eines Teilprogramms den Gesamtprozess beeinflusst. Bei der Genauigkeit hat man **alle möglichen auftretenden** Fehler zu berücksichtigen und ihren Einfluss auf das Ergebnis.
Fehler sind hier nicht Programmier- oder Hardwarefehler, sondern die Abweichung, eines nach richtig angewendeter Vorschrift erzeugten Resultats, vom gewünschten Ergebnis.

Es gibt verschiedene Fehlerarten:

1. Rundungs- oder Reduktionsfehler, s.o.
2. Datenfehler: Üblicherweise sind die Eingangsdaten ungenau, z. B. Messdaten mit Messfehlern. Das Lösen eines Problems entspricht dem Auswerten einer Funktion $f : D \rightarrow W$ für Datenwerte $x \in D$. Anstelle von x hat man gestörten Wert \tilde{x} . Der Datenfehler ist dann $f(x) - f(\tilde{x})$. Die Empfindlichkeit der Lösung des Problems (Auswertung von f , gegenüber kleinen Änderungen in den Daten (in x)) heisst Kondition des Problems. Die gute oder schlechte Kondition ist eine Eigenschaft des Problems und NICHT des Verfahrens. Später mehr.
3. Verfahrensfehler: Exakte Verfahren enden nach endlich vielen Operationen (bei exakter Rechnung) mit dem exakten Ergebnis $y = f(x)$. Näherungsverfahren (z.B. Iterationen) enden in Abhängigkeit von bestimmten Kriterien mit einer Näherung \tilde{y} für die Lösung y .
Verfahrensfehler: $\tilde{y} - y$

!



Datenfehler

Verfahrensfehler

Rundungsfehler

Gegeben sei eine Funktion $f(a_1, \dots, a_n)$ und eine durch Fehler gestörte Funktion $\tilde{f}(a_1, \dots, a_n)$.

Fehleranalyse: Verfolgung der Auswirkung aller Fehler, die in den einzelnen Schritten bei der Berechnung von f auftreten können.

Vorwärtsanalyse: Verfolge den Fehler von Schritt zu Schritt und schätze für jedes Teilergebnis den akkumulierten Fehler ab, d.h. verfolge den Fehler, so dass das Ergebnis die Form $\tilde{f}(a_1, \dots, a_n) = f(a_1, \dots, a_n) + \delta$ hat, wobei δ akkumulierter Fehler ist.

Rückwärtsanalyse: Die Verfolgung des Fehlers geschieht so, dass jedes Zwischenergebnis als exakt berechnetes Ergebnis für gestörte Daten interpretiert wird, d.h. der akkumulierte Fehler im Teilergebnis wird als Datenfehlereffekt interpretiert. Also: $\tilde{f}(a_1, \dots, a_n) = f(\tilde{a}_1, \dots, \tilde{a}_n)$ mit gestörten Daten $\tilde{a}_1, \dots, \tilde{a}_n$. Über den Fehler im Endergebnis erhält man so zunächst nur die Aussage, dass er einem Datenfehler bestimmter Größe entspricht (äquivalenter Datenfehler). Die Größe der äquivalenten Datenfehler dient dann zur Beurteilung der Algorithmen.

triviales Beispiel:

$$f(a_1, a_2) = a_1 + a_2, \quad \tilde{f}(a_1, a_2) = a_1 \oplus a_2.$$

$$\begin{array}{l}
 a_1 \oplus a_2 \\
 \text{Vorwärtsanalyse} \\
 = \\
 \underbrace{a_1 + a_2}_{f(a_1, a_2)} + \underbrace{\varepsilon(a_1 + a_2)}_{\delta} \quad |\varepsilon| \leq \text{eps}, |\delta| \leq \text{eps} (|a_1| + |a_2|) \\
 \text{Rückwärtsanalyse} \\
 = \\
 \underbrace{a_1(1 + \varepsilon) + a_2(1 + \varepsilon)}_{f(\tilde{a}_1, \tilde{a}_2)} \quad |\varepsilon| \leq \text{eps}
 \end{array}$$

Vorwärtsanalyse ist im allgemeinen kaum durchführbar, für die meisten Verfahren ist, wenn überhaupt, nur Rückwärtsanalyse bekannt. Ein idealer Algorithmus im Sinne der Rückwärtsanalyse hat einen äquivalenten Datenfehler von der Größenordnung der Rundungsfehler oder Eingangsfehler.

Aussagen über den Gesamtfehler erhält man über sogenannte Störungssätze.

Vorteil dieser Vorgehensweise: Auswirkungen der Arithmetik und des Verfahrens getrennt von der Kondition des Problems.

Beispiel 1.27 $x^2 - 2a_1x + a_2 = 0 \quad 0 < a_2 \ll 1 < a_1$.

Aufgabe: Finde die kleine Lösung

$$x^* = f(a_1, a_2) = a_1 - \sqrt{a_1^2 - a_2}$$

Kondition des Problems: Setze $a = (a_1, a_2)$ und betrachte Taylorentwicklung von

$$f(a + \Delta a) = \tilde{x}^*$$

$$\begin{aligned}
 f(a + \Delta a) &= f(a) + \frac{\partial f}{\partial a_1}(a) \Delta a_1 + \frac{\partial f}{\partial a_2}(a) \Delta a_2 \\
 &+ \text{Terme höherer Ordnung in } \Delta a_1, \Delta a_2.
 \end{aligned} \tag{1.28}$$

$$\frac{\tilde{x}^* - x^*}{x^*} = \underbrace{\frac{\partial f(a)}{\partial a_1} \frac{a_1}{x^*}}_{\text{Kondition}} \frac{\Delta a_1}{a_1} + \underbrace{\frac{\partial f(a)}{\partial a_2} \frac{a_2}{x^*}}_{\text{Rundungsfehler}} \frac{\Delta a_2}{a_2} + \text{T.h.o.}$$

Verstärkungsfaktoren für
relativen Fehler in Daten.

$$\begin{aligned}\frac{\partial f}{\partial a_1}(a) &= 1 - \frac{1}{2} \frac{2a_1}{\sqrt{a_1^2 - a_2}} = \frac{\sqrt{a_1^2 - a_2} - a_1}{\sqrt{a_1^2 - a_2}} = \frac{-x^*}{\sqrt{a_1^2 - a_2}} \\ \frac{\partial f}{\partial a_2}(a) &= \frac{1}{2} \frac{1}{\sqrt{a_1^2 - a_2}} \\ \frac{\partial f}{\partial a_1}(a) \frac{a_1}{x^*} &= \frac{-x^*}{\sqrt{a_1^2 - a_2}} \frac{a_1}{x^*} \approx -1 \quad \text{da } a_2 \ll 1 < a_1 \\ \frac{\partial f}{\partial a_2}(a) \frac{a_2}{x^*} &= \frac{a_2}{2\sqrt{a_1^2 - a_2} x^*} = \frac{a_2}{\frac{2a_2}{\frac{a_1}{\sqrt{a_1^2 - a_2}} + 1}} \approx 1\end{aligned}$$

Also ist das Problem gut konditioniert. Erwarten würden wir bei einem guten Algorithmus einen relativen Fehler von $c \cdot \text{eps}$, mit kleinem c .

<p>Algorithmus:</p> <p>$y_1 := a_1 * a_1$</p> <p>$y_2 := y_1 - a_2$</p> <p>$y_3 := \sqrt{y_2}$</p> <p>$x := y_4 = a_1 - y_3$</p>	<p>$\mathbb{M}(10, 5, 1) a_1 = 6.002, a_2 = 0.01$</p> <p>$y_1 = 36.024$</p> <p>$y_2 = 36.014$</p> <p>$y_3 = 6.0012$</p> <p>$\tilde{x} = y_4 = 0.00080000$</p> <p>exakt 0.00083311</p> <p>$\frac{\Delta x}{x} = 0.039747$</p>
--	--

Rückwärtsanalyse:

$$\begin{aligned}& \left[a_1 - \sqrt{(a_1^2(1 + \varepsilon_1) - a_2)(1 + \varepsilon_2)(1 + \varepsilon_3)} \right] (1 + \varepsilon_4), |\varepsilon_i| < \text{eps} . \\ &= a_1(1 + \varepsilon_4) - \sqrt{a_1^2(1 + \varepsilon_4)^2 \underbrace{(1 + \varepsilon_1)(1 + \varepsilon_2)(1 + \varepsilon_3)}_{\substack{1 + \varepsilon_5 \\ |\varepsilon_5| < 4 \text{ eps} .}} - a_2 \underbrace{(1 + \varepsilon_2)(1 + \varepsilon_3)^2(1 + \varepsilon_4)^2}_{\substack{1 + \varepsilon_6 \\ |\varepsilon_6| < 5 \text{ eps} .}}} \\ &= a_1(1 + \varepsilon_4) - \sqrt{(a_1(1 + \varepsilon_4))^2 - a_2 \underbrace{\left((1 + \varepsilon_6) - \frac{a_1^2(1 + \varepsilon_4)^2 \varepsilon_5}{a_2} \right)}_{\substack{1 + \varepsilon_7 \\ |\varepsilon_7| < \text{eps} (5 + 4 \frac{a_1^2}{|a_2|})}})}\end{aligned}$$

Der Rechner liefert die exakte Lösung von $x^2 - 2a_1(1 + \varepsilon_4)x + a_2(1 + \varepsilon_7) = 0$ mit $|\varepsilon_4| < \text{eps}$ und $|\varepsilon_7| < \text{eps} (5 + 4 \frac{a_1^2}{|a_2|})$, d.h. der Algorithmus ist ungeeignet.

$\underbrace{\hspace{10em}}_{\text{groß}}$

→ Übung 1.46: besserer Alg.

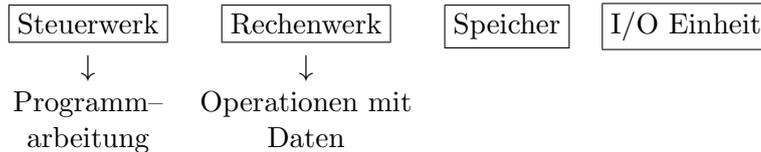
Definition 1.29 Ein Algorithmus für den die Rückwärtsanalyse liefert, dass das berechnete Ergebnis das exakte Ergebnis eines gestörten Problems mit Störungen der Größenordnung $c \cdot \text{eps}$ ist, heißt numerisch rückwärts stabil. Ein Algorithmus für die Vorwärtsanalyse liefert, dass das berechnete Resultat einen relativen Fehler der Größenordnung $c \cdot \text{eps}$ hat,

heisst numerisch vorwärts stabil. Ansonsten ist der Algorithmus nicht numerisch stabil. (Andere Stabilitätsbegriffe.)

! Faustregel: Gut konditioniertes Problem und stabiler Algorithmus \implies gute Ergebnisse.
Schlecht konditioniertes Problem oder instabiler Algorithmus \implies ? Ergebnisse.

1.4 Moderne Rechnerarchitekturen

a) Klassischer von Neumann-Rechner:



Alle vorkommenden Operationen wie Ein/Ausgabe, Speicherzugriff, arithmetische Operationen, etc. werden seriell nacheinander ausgeführt. In modernerer Form können Speicherzugriff und arithmetische Operationen gleichzeitig ausgeführt werden. Dies hat jedoch noch wenig Einfluss auf numerische Methoden.

b) Vektorrechner arbeiten im allgemeinen nach dem SIMD Prinzip. (single instruction, multiple data). In ihren Prozessoren ist ein Pipeline-Prinzip realisiert, d.h. verschiedene Phasen einer Operation können mit verschiedenen Operanden gleichzeitig durchgeführt werden.

Beispiel 1.30 Addition zweier Dualzahlen

$x = a \cdot 2^p, y = b \cdot 2^q, q \geq p$ in normalisierter Darstellung. 4 Teilschritte

1. Exponentenvergleich: $q - p$.
2. Verschieben der Mantisse x um $q - p$ Stellen, so dass $x = \tilde{a}2^q$.
3. Addition der Mantissen $\tilde{a} + b$.
4. Normalisierte Darstellung von $x + y$.

Bei einem seriellen Rechner werden alle 4 Operationen nacheinander ausgeführt, erst wenn Teil 4) für ein paar Operanden fertig ist, wird mit Teil 1 für das nächste Paar begonnen. In einem Vektorprozessor wird das nächste Paar in Teil 1 bearbeitet, sofort wenn ein Paar in Teil 2 ist. Sei τ die Zeit für einen Zyklus

$$\begin{array}{l}
 t = 0 \quad \left[\begin{array}{c|c|c|c} x_1 & & & \\ \hline y_1 & & & \end{array} \right] \quad \xrightarrow{x_2, y_2} \quad \left[\begin{array}{c|c|c|c} x_1 & & & \\ \hline y_1 & & & \end{array} \right] \\
 t = \tau \quad \left[\begin{array}{c|c|c|c} & x_1 & & \\ \hline & y_1 & & \end{array} \right] \quad \xrightarrow{x_3, y_3} \quad \left[\begin{array}{c|c|c|c} x_2 & x_1 & & \\ \hline y_2 & y_1 & & \end{array} \right] \\
 t = 2\tau \quad \left[\begin{array}{c|c|c|c} & & x_1 & \\ \hline & & y_1 & \end{array} \right] \quad \xrightarrow{x_4, y_4} \quad \left[\begin{array}{c|c|c|c} x_3 & x_2 & x_1 & \\ \hline y_3 & y_2 & y_1 & \end{array} \right] \\
 t = 3\tau \quad \xrightarrow{x_2, y_2} \left[\begin{array}{c|c|c|c} & & & x_1 \\ \hline & & & y_1 \end{array} \right] \quad \xrightarrow{x_5, y_5} \quad \left[\begin{array}{c|c|c|c} x_4 & x_3 & x_2 & x_1 \\ \hline y_4 & y_3 & y_2 & y_1 \end{array} \right] \\
 t = 4\tau \quad \left[\begin{array}{c|c|c|c} x_2 & & & \\ \hline y_2 & & & \end{array} \right] \xrightarrow{x_1+y_1} \quad \left[\begin{array}{c|c|c|c} x_5 & x_4 & x_3 & x_2 \\ \hline y_5 & y_4 & y_3 & y_2 \end{array} \right] \xrightarrow{x_1+y_1} \\
 t = 5\tau \quad \left[\begin{array}{c|c|c|c} & x_2 & & \\ \hline & y_2 & & \end{array} \right] \quad \xrightarrow{x_7, y_7} \quad \left[\begin{array}{c|c|c|c} x_6 & x_5 & x_4 & x_3 \\ \hline y_6 & y_5 & y_4 & y_3 \end{array} \right] \xrightarrow{x_2+y_2}
 \end{array}$$

Nach Start-up Phase, je Takt ein Ergebnis, gut falls vielfach dieselben Operationen gemacht werden. Vektoraddition z.B.

Hohe Leistung beim Vektorrechner nur wenn der Datenfluss zum Prozessor gut funktioniert.

Register-Register-Maschinen

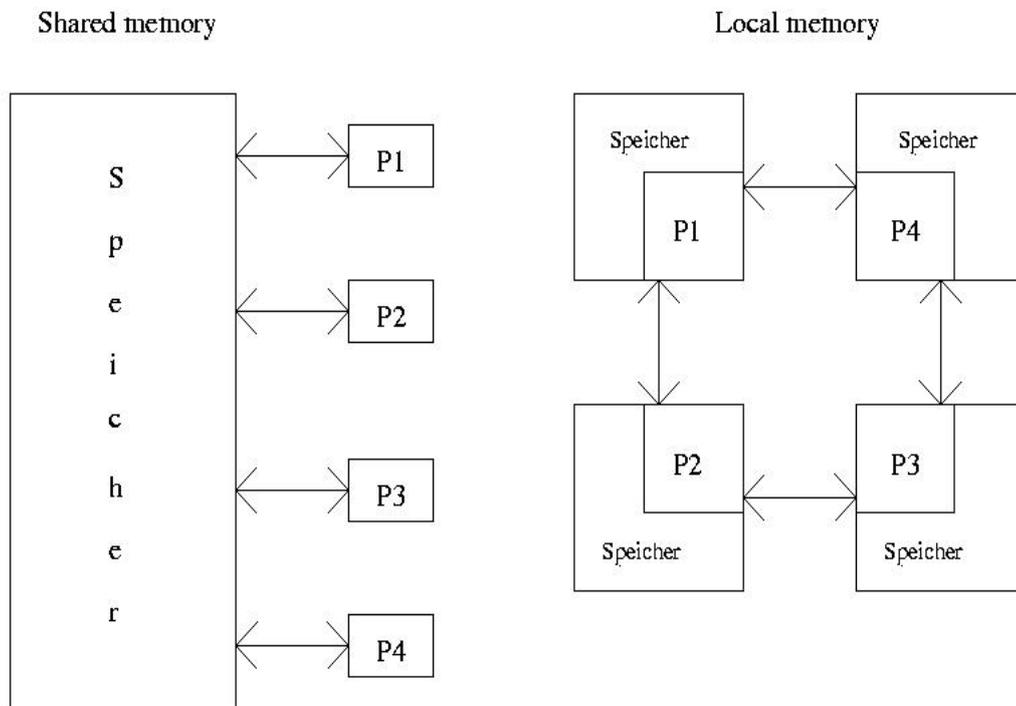
Vektoren werden in Zwischenspeichern, Vektorregistern gehalten, so lange man mit Daten aus dem aktuellen Register arbeitet sehr schnell, falls jedoch oft ein Register aus dem Hauptspeicher geladen werden muss, schlechter.

Speicher-Speicher-Maschinen

Vektorprozessor greift direkt auf Hauptspeicher. Im allgemeinen hohe Startup zeiten. Sehr gut für lange Vektoren.

c) Parallelrechner

Mehrere Prozessoren führen Operationen gleichzeitig aus. Verschiedene Modelle.

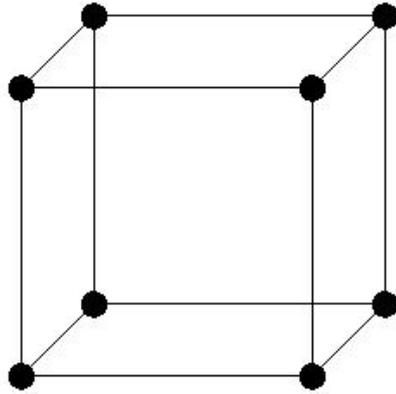


Falls ein zentraler Kontrollprozessor ein zentrales Programm abarbeitet und die Operanden dabei auf die anderen Prozessoren verteilt und alle dieselbe Operation ausführen SIMD-Rechner. Führt dagegen jeder Rechner auch sein eigenes Programm aus MIMD (multiple instruction, multiple data).

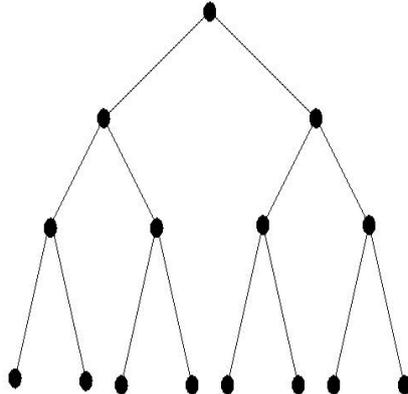
local memory bei vielen Prozessoren,
shared memory bei wenig Prozessoren.

Probleme: Verteilung der Lösung gleichmäßig
Datenkommunikation.

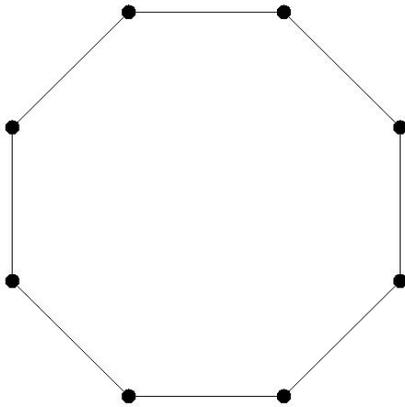
typische Architekturen: p Anzahl Prozessoren
 l Anzahl links



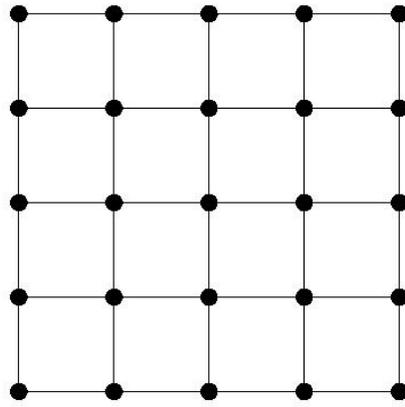
Hypercube, $p = 8, l = 3$



Binärbaum $p = 15, l \leq 3$



Ring, $p = 8, l = 2$



Gitter $p = 16, l \leq 4$

In vielen Algorithmen kann man auf Parallelität und Vektorisierbarkeit mit Rücksicht nehmen. Es ergeben sich unterschiedliche Beurteilungskriterien, je nach Rechnertyp. Es ist schwierig, allgemeine Aussagen zu machen, aber man muss diese Punkte heute immer mit berücksichtigen.

1.5 Die ':' Notation, flops und Normen

Im folgenden verwenden wir häufig die folgende Notation, die in vielen Softwareprogrammen und bei einigen Programmiersprachen verwendet wird. Für eine $m \times n$ Matrix $A = [a_{ij}]$ ($A \in \mathbb{R}^{m,n}$ oder $A \in \mathbb{C}^{m,n}$) schreiben wir

$$A(k, :) = [a_{k1}, \dots, a_{kn}],$$

dies ist also die k -te Zeile von A .

$$A(:, k) = \begin{bmatrix} a_{1k} \\ \vdots \\ a_{mk} \end{bmatrix}$$

ist die k -te Spalte von A . Allgemein steht $1 : n$ für $1, \dots, n$.

Um Algorithmen vergleichen zu können werden vielfach Operationen gezählt. Ein bewährtes Maß, das oft verwendet wird ist ein sogenannter flop, d.h. eine Gleitkommaoperation ($x, - * \div$).

Alle werden gleich gewichtet. Das war früher anderes. Früher wurden Additionen gegenüber Multiplikationen vernachlässigt oder die Kosten für die Anweisung

$$C(i, j) = C(i, j) + A(i, k) * B(k, j).$$

wurden als flop bezeichnet. Dies lässt sich bei den heutigen Rechnern nicht mehr rechtfertigen.

Normen:

Eine Vektornorm auf \mathbb{R}^n ist eine Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ welche die folgenden Bedingungen erfüllt

$$\begin{aligned} f(x) &\geq 0 && \forall x \in \mathbb{R}^n (f(x) = 0 \iff x = 0) \\ f(x+y) &\leq f(x) + f(y) && \forall x, y \in \mathbb{R}^n \\ f(\alpha x) &= |\alpha|f(x) && \forall \alpha \in \mathbb{R}, x \in \mathbb{R}^n \end{aligned} \tag{1.31}$$

Notation $f(x) = \|x\|$

p-Normen:

$$\begin{aligned} \|x\|_p &= (|x_1|^p + \dots + |x_n|^p)^{\frac{1}{p}}, p \geq 1 \\ \|x\|_1 &= |x_1| + |x_2| + \dots + |x_n| \\ \|x\|_2 &= (|x_1|^2 + |x_2|^2 + \dots + |x_n|^2)^{\frac{1}{2}} = (x^T x)^{\frac{1}{2}} \\ \|x\|_\infty &= \max_{1 \leq k \leq n} |x_k| \end{aligned}$$

Wichtige Ungleichungen:

Höldersche Ungleichung

$$|x^T y| \leq \|x\|_p \|y\|_q, \text{ für } \frac{1}{p} + \frac{1}{q} = 1 \tag{1.32}$$

$$|x^T y| \leq \|x\|_2 \|y\|_2 \quad \text{Cauchy-Schwarz} \tag{1.33}$$

$$\left. \begin{aligned} \|x\|_2 &\leq \|x\|_1 &\leq \sqrt{n} \|x\|_2 \\ \|x\|_\infty &\leq \|x\|_2 &\leq \sqrt{n} \|x\|_\infty \\ \|x\|_\infty &\leq \|x\|_1 &\leq n \|x\|_\infty \end{aligned} \right\} \text{ Normäquivalenz} \tag{1.34}$$

$f : \mathbb{R}^{m,n} \rightarrow \mathbb{R}$ ist eine Matrix Norm, falls

$$\begin{aligned} f(A) &\geq 0 && \forall A \in \mathbb{R}^{m,n} (f(A) = 0 \iff A = 0) \\ f(A+B) &\leq f(A) + f(B) && \forall A, B \in \mathbb{R}^{m,n} \\ f(\alpha A) &= |\alpha|f(A) && \forall \alpha \in \mathbb{R}, A \in \mathbb{R}^{m,n} \end{aligned} \tag{1.35}$$

Frobenius Norm:

$$\begin{aligned} \|A\|_F &= \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right)^{\frac{1}{2}} \\ (\|\bullet\|_2 - \text{Norm auf } \mathbb{R}^{mn}) \end{aligned} \tag{1.36}$$

→ Übung 1.47.

p-Normen

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \sup_{x \neq 0} \left\| A \frac{x}{\|x\|_p} \right\|_p = \max_{\|x\|_p=1} \|Ax\|_p \tag{1.37}$$

→ Übung 1.50.

Die Matrix Norm ist abhängig von der Dimension. $\|\bullet\|_2$ auf $\mathbb{R}^{3 \times 2}$ ist etwas anderes als $\|\bullet\|_2$ auf $\mathbb{R}^{5,6}$.

Konsistente Matrixnorm

$$\begin{aligned} \|AB\|_p &\leq \|A\|_p \|B\|_p \\ \|AB\|_F &\leq \|A\|_2 \|B\|_F \leq \|A\|_F \|B\|_F \end{aligned} \quad \forall A \in \mathbb{R}^{m,n}, B \in \mathbb{R}^{n,q} \quad (1.38)$$

→ Übung 1.47.

Nicht alle Normen erfüllen die Konsistenzbedingung.

Ungleichungen und Gleichungen

$$\left. \begin{aligned} \|A\|_2 &\leq \|A\|_F \leq \sqrt{n} \|A\|_2 \\ \max_{i,j} |a_{ij}| &\leq \|A\|_2 \leq \sqrt{mn} \max_{i,j} |a_{ij}| \\ \|A\|_1 &= \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \\ \|A\|_\infty &= \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \\ \frac{1}{\sqrt{n}} \|A\|_\infty &\leq \|A\|_2 \leq \sqrt{m} \|A\|_\infty \\ \frac{1}{\sqrt{m}} \|A\|_1 &\leq \|A\|_2 \leq \sqrt{n} \|A\|_1 \end{aligned} \right\} \quad (1.39)$$

Die Matrix 2-Norm ist nicht so leicht zu charakterisieren wie $\|\bullet\|_1, \|\bullet\|_\infty$.

→ Übung 1.48

Satz 1.40 Sei $A \in \mathbb{R}^{m,n}$. Dann gibt es einen Vektor $z \in \mathbb{R}^n$ mit $\|z\|_2 = 1$, so dass $A^T A z = \mu^2 z$, wobei $\mu = \|A\|_2$, d.h. $\|A\|_2^2$ ist ein Eigenwert von $A^T A$.

Beweis: Sei $z \in \mathbb{R}^n, \|z\|_2 = 1$ und $\|Az\|_2 = \|A\|_2$, dann maximiert z die Funktion

$$g(x) = \frac{1}{2} \frac{\|Ax\|_2^2}{\|x\|_2^2} = \frac{1}{2} \frac{x^T A^T A x}{x^T x}$$

also gilt $\nabla g(z) = 0$ (∇g Gradient von g)

$$\begin{aligned} \left. \frac{\partial g(x)}{\partial x_i} \right|_{x=z} &= \left[(z^T z) \sum_{j=1}^n (A^T A)_{ij} z_j - (z^T A^T A z) z \right] / (z^T z)^2 \\ \implies A^T A z &= \underbrace{(z^T A^T A z)}_{\mu^2} z = \mu^2 z. \end{aligned}$$

□

Insbesondere gilt, dass $\|A\|_2^2$ der größte Eigenwert von $A^T A$ ist.

Norminvarianz

$$\begin{aligned} \|UAV\|_2 &= \|A\|_2 \\ \|UAV\|_F &= \|A\|_F \end{aligned}, \quad U \in \mathbb{R}^{m,m}, V \in \mathbb{R}^{n,n} \text{ orthogonal}, A \in \mathbb{R}^{m,n}. \quad (1.41)$$

Absolutbeträge für Matrizen

$A \in \mathbb{R}^{m,n}$, dann ist $|A| = B$ mit $b_{ij} = |a_{ij}|$. Wir setzen

$$B \leq A, \text{ falls } b_{ij} \leq a_{ij}, \forall i = 1 : m, j = 1 : n.$$

Für das Rechnen mit Beträgen gelten folgende Regeln

$$\left. \begin{aligned} |A+B| &\leq |A| + |B| \\ |AB| &\leq |A| |B| \\ A &\leq B, C, D \geq O \implies CAD \leq CBD \\ \|A\|_p &\leq \| |A| \|_p \\ \|A\| &= \| |A| \|, \text{ für } \|\bullet\| = \|\bullet\|_1, \|\bullet\|_\infty, \|\bullet\|_F \\ |A| &\leq |B| \implies \|A\|_1 \leq \|B\|_1, \|A\|_\infty \leq \|B\|_\infty, \|A\|_F \leq \|B\|_F \end{aligned} \right\} \quad (1.42)$$

→ Übung 1.49.

Wir haben gesehen, dass sich beim Speichern oder Runden ergibt

$$[gl(A)]_{ij} = [gl(a_{ij})] = [a_{ij}(1 + \varepsilon_{ij})] \text{ mit } |\varepsilon_{ij}| \leq \text{eps}$$

Dann folgt

$$|gl(A) - A| \leq \text{eps} |A|.$$

Dieses kann auch als Normungleichung umgeschrieben werden, $\|gl(A) - A\|_1 \leq \text{eps} \|A\|_1$. Heute werden vielfach für Fehleranalysen nicht Normabschätzungen sondern elementweise Abschätzungen (LAPACK) verwendet.

1.6 Übungen zu Kapitel 1

Übung 1.43 Stelle folgende Dezimalzahlen im Dual-, Hexadezimalsystem sowie im Zwölfersystem dar (d.h. $p = 2, 16, 12$, für $p = 12$ verwende man A, B als Ziffern für 10, 11)

$$275, 0.1, 1/3, 12.125$$

Übung 1.44 a_1, \dots, a_n seien Maschinenzahlen und es sei $f(a_1, \dots, a_n) = \sum_{i=1}^n a_i$. Der Computer habe die Maschinengenauigkeit eps . Auf dem Computer wird f in der Form $\tilde{f}(a_1, \dots, a_n) = (\dots((a_1 \oplus a_2) \oplus a_3) \oplus \dots) \oplus a_n$ berechnet.

- Mache eine Rückwärtsanalyse für f . D.h. $\tilde{f}(a_1, \dots, a_n) = f(a_1(1 + \delta_1), \dots, a_n(1 + \delta_n))$ und $|\delta_1|, \dots, |\delta_n|$ sind geeignet abzuschätzen.
- Mache eine Vorwärtsanalyse für f . D.h. $\tilde{f}(a_1, \dots, a_n) = f(a_1, \dots, a_n) + \varepsilon$ und $|\varepsilon|$ ist geeignet abzuschätzen.
- Zeige, daß für den Fehler zwischen \tilde{f} und f gilt:

$$|\tilde{f}(a_1, \dots, a_n) - f(a_1, \dots, a_n)| \leq \text{eps} \sum_{i=1}^n |n - i + 1| |a_i|.$$

Terme in der Größenordnung $\text{eps}^2, \text{eps}^3, \dots$ dürfen als näherungsweise 0 vernachlässigt werden.

Übung 1.45 Ermittle in MATLAB die Realisierung von underflow/overflow sowie $x_{\min}, x_{\max}, \text{eps}$.

Übung 1.46 Gib für das Beispiel aus der Vorlesung zur Berechnung von $x^* = f(a_1, a_2) = a_1 - \sqrt{a_1^2 - a_2}$ einen besseren Algorithmus an. Begründung!

Übung 1.47 Sei $A \in \mathbb{R}^{m,n}$. Zeige:

1. $\|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2$
2. $\|AB\|_F \leq \|A\|_2 \|B\|_F \leq \|A\|_F \|B\|_F$.

Übung 1.48 Sei $A \in \mathbb{R}^{m,n}$. Zeige:

1. Ist $B \in \mathbb{R}^{n,l}$ so gilt $\|AB\|_p \leq \|A\|_p \|B\|_p$
Tipp: Zeige die Behauptung zuerst für $l = 1$ und nutze die Def.

$$2. \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|.$$

Tipp: Zeige zunächst $\|Ax\|_1 \leq (\max \Sigma \dots) \|x\|_1$ und konstruiere dann ein x für welches Gleichheit gilt.

$$3. \|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|.$$

4. Ist $\rho(A) = \max\{|\lambda| : \lambda \text{ Eigenwert von } A\}$ eine Matrixnorm?

Übung 1.49 Sei $A \in \mathbb{R}^{m,n}$, $B \in \mathbb{R}^{k,l}$. Zeige:

1. Falls $k = m, l = n$, dann gilt $|A + B| \leq |A| + |B|$.
2. Falls $k = n$, dann gilt $|AB| \leq |A| |B|$.
3. Falls $k = m, l = n$ und $A \leq B$, dann gilt für alle $C \in \mathbb{R}^{p,m}$ mit $C \geq O$: $CA \leq CB$.
4. Falls $k = m, l = n$ und $|A| \leq |B|$, dann gilt $\|A\|_\infty \leq \|B\|_\infty$.

Übung 1.50 Für eine Matrix $A \in \mathbb{R}^{m,n}$ ist die Matrix p -Norm definiert durch

$$\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} = \max_{\|x\|_p=1} \|Ax\|_p$$

1. $\|A\|_p$ ist eine Norm
2. $\|Ax\|_p \leq \|A\|_p \|x\|_p$
3. Ist c eine Konstante, so dass $\|Ax\|_p \leq c \|x\|_p$ für alle x gilt, so ist $c \geq \|A\|_p$, d.h. $c = \|A\|_p$ ist die kleinst mögliche Konstante, so dass $\|Ax\|_p \leq c \|x\|_p$ gilt.

Kapitel 2

Lösung linearer Gleichungssysteme

Problem: Berechne Lösung von $AX = B$ wobei $A \in \text{GL}_n(\mathbb{C}) = \mathbb{C}^{n,n}$ (oder $\text{GL}_n(\mathbb{R}) = \mathbb{R}^{n,n}$) und $B \in \mathbb{C}^{n,m}$ (oder $\mathbb{R}^{m,n}$) ist. Im folgenden $\mathbb{K} = \mathbb{C}$ oder $\mathbb{K} = \mathbb{R}$. Wir betrachten zuerst direkte Verfahren, d.h. Verfahren die bei exakter Rechnung in endlich vielen Schritten enden.

2.1 Gauß-Elimination

Der immer noch in vielen Variationen am weitesten bekannte Algorithmus ist die Gauß-Elimination.

2.1.1 Lösung von Dreieckssystemen

Zur Vorbereitung betrachten wir zuerst einmal spezielle Dreiecksmatrizen und nur eine rechte Seite. Sei nun $L = [l_{ij}] \in \mathbb{K}^{n,n}$ untere Δ -Matrix. Wir betrachten die Lösung von $Lx = b = [b_i]$ mit $b \in \mathbb{K}^n$. Wir erhalten sofort als Lösung

$$x_i = \left(b_i - \sum_{j=1}^{i-1} l_{ij}x_j \right) / l_{ii} \quad i = 1, \dots, n \quad (2.1)$$

Dieses Verfahren heißt Vorwärts-Einsetzen und wird durch folgenden Algorithmus realisiert.

Algorithmus 2.2

Input: $L \in \mathbb{K}^{n,n}$ untere Δ -Matrix, nichtsingulär, $b \in \mathbb{K}^n$.

Output: Lösung von $Lx = b$, überschrieben auf b .

$b(1) = b(1)/L(1, 1);$

FOR $i = 2 : n$

$b(i) = (b(i) - L(i, 1 : i - 1) * b(1 : i - 1))/L(i, i);$

END

! Achtung die Multiplikation $L(i, 1 : i - 1) * b(1 : i - 1)$ ist natürlich eine Schleife.
→ Übung 2.58

Kosten für Algorithmus 2.2: n^2 flops

Fehleranalyse für Algorithmus 2.2: Die berechnete Lösung von \tilde{x} erfüllt

$$(L + F)\tilde{x} = b, \text{ wobei } |F| \leq n \cdot \text{eps} \cdot |L| + \mathcal{O}(\text{eps}^2). \quad (2.3)$$

Der analoge Algorithmus für obere Dreiecksmatrizen heißt Rückwärts-Einsetzen. Sei $U \in [u_{ij}] \in \mathbb{K}^{n,n}$ obere Dreiecksmatrix. Für die Lösung von $Ux = b = [b_i] \in \mathbb{K}^{n,n}$ erhalten wir

$$x_i = \left(b_i - \sum_{j=i+1}^n u_{ij}x_j \right) / u_{jj} \quad i = n, n-1, \dots, 1 \quad (2.4)$$

durch den folgenden Algorithmus.

Algorithmus 2.5

Input: $U \in \mathbb{K}^{n,n}$ nichtsinguläre obere Δ -Matrix $b \in \mathbb{K}^n$.

Output: Lösung von $Ux = b$, überschrieben auf b .

```

b(n) = b(n)/U(n, n);
FOR   i = n - 1 : -1 : 1
      b(i) = (b(i) - U(i, i + 1 : n) * b(i + 1 : n))/U(i, i);
END

```

Es gilt wieder

Kosten Algorithmus 2.5: n^2 flops

Fehleranalyse Algorithmus 2.5: \tilde{x} erfüllt

$$(U + F)\tilde{x} = b, \text{ wobei } |F| \leq n \cdot \text{eps} \cdot |U| + \mathcal{O}(\text{eps}^2). \quad (2.6)$$

→ Übung 2.58

Es gibt Varianten dieser Algorithmen für Parallel- und Vektorrechner (spaltenorientierte Versionen) und entsprechende Block-Versionen für mehrfache rechte Seiten. Die Algorithmen zum Vorwärts- und Rückwärts-einsetzen bilden die Basis für die Lösung von Gleichungssystemen mittels LR -Zerlegung.

2.1.2 LR -Zerlegung

Wir kommen nun zur LR -Zerlegung einer Matrix A als $A = LR$, mit L untere und R obere Δ -Matrix. Falls man so eine Zerlegung hat, so löst man $Ax = b$ mittels Algorithmus 2.5 und 2.2 indem man $Ly = b$, $Rx = y$ nacheinander löst. Die LR -Zerlegung wird mittels des Gaußschen Eliminationsverfahrens erzeugt. Dazu verwenden wir sogenannte Gauß-Transformationen.

Sei $x \in \mathbb{K}^n$ mit $x_k \neq 0$. Sei

$$t \equiv t^{(k)} = \left[\begin{array}{c} 0 \\ \vdots \\ 0 \\ t_{k+1} \\ \vdots \\ t_n \end{array} \right] \left. \vphantom{\begin{array}{c} 0 \\ \vdots \\ 0 \\ t_{k+1} \\ \vdots \\ t_n \end{array}} \right\} k, \quad t_i = \frac{x_i}{x_k}, i = k + 1 : n.$$

→ Übung 2.62

Mehrfache Anwendung von *GAUSSAPP* führt dann auf Dreiecksgestalt.
Im k -ten Schritt erhalten wir (falls alles glatt geht)

$$A^{(k-1)} := M_{k-1} \cdots M_2 M_1 A = \begin{bmatrix} a_{11}^{(k-1)} & \cdots & a_{1,k-1}^{(k-1)} & \cdots & \cdots & a_{1,n}^{(k-1)} \\ 0 & \ddots & & & & \vdots \\ \vdots & \ddots & a_{k-1,k-1}^{(k-1)} & \cdots & \cdots & a_{k-1,n}^{(k-1)} \\ \vdots & & 0 & a_{kk}^{(k-1)} & & a_{k,n}^{(k-1)} \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & a_{n,k}^{(k-1)} & \cdots & a_{nn}^{(k-1)} \end{bmatrix}. \quad (2.15)$$

Wir fahren dann fort auf dem noch nicht reduzierten unteren Block. Die vollständige Dreiecksreduktion wird dann durch die folgende Schleife erzielt.

Algorithmus 2.16

```

n = size(A, 1);
k = 1;
WHILE A(k, k) ≠ 0 AND k < n
    t = GAUSS(A(k : n, k));
    A(k : n, :) = GAUSSAPP(A(k : n, :), t);
    k = k + 1;
END

```

! Die Elemente $A(k, k)$, die während des Algorithmus auf '0' überprüft werden müssen, heißen Pivots. Ihre relative Größe ist von größter Bedeutung für die Fehleranalyse.

Matrizentheoretisch formuliert gilt, wenn Algorithmus 2.16 mit $k = n$ endet, daß

$$M_{n-1} \cdots M_1 A =: R$$

mit R obere Δ -Matrix ist.

Für jedes $M_k = I - t^{(k)} e_k^T$ gilt $M_k^{-1} = I + t^{(k)} e_k^T$, also folgt:

$A = M_1^{-1} \cdots M_{n-1}^{-1} R =: L \cdot R$. Alle M_k sind untere Δ -Matrizen mit 1-Diagonale also auch L .

Eine Zerlegung $A = L \cdot R$ mit R obere Δ -Matrix und L untere Δ -Matrix mit 1-Diagonale heißt *LR-Zerlegung* (engl. *LU-Decomposition*).

Satz 2.17 (Existenz und Eindeutigkeit der LR-Zerlegung)

$A \in \mathbb{K}^{n,n}$ hat eine LR-Zerlegung genau dann, wenn

$$\det(A(1 : k, 1 : k)) \neq 0, \quad \text{für } k = 1 : n - 1. \quad (2.18)$$

Falls die LR-Zerlegung existiert und A nichtsingulär ist, so ist die LR-Zerlegung eindeutig und $\det A = r_{11} \cdots r_{nn}$.

Beweis: Falls A eine LR-Zerlegung

$$A = \begin{bmatrix} 1 & & & \\ l_{21} & \ddots & & \\ \vdots & \ddots & \ddots & \\ l_{n1} & \cdots & l_{n,n+1} & 1 \end{bmatrix} \begin{bmatrix} r_{11} & \cdots & \cdots & r_{1n} \\ & \ddots & & \vdots \\ & & \ddots & \vdots \\ & & & r_{n,n} \end{bmatrix}$$

besitzt, so gilt: alle r_{ii} , $i = 1 : n - 1$ sind $\neq 0$, da sie die Pivots in Algorithmus 2.16 sind.

$$\det(A(1:k, 1:k)) = \det(L(1:k, 1:k)) \cdot \det(R(1:k, 1:k)) = 1 \cdot \prod_{i=1}^k r_{ii} \neq 0.$$

Die Rückrichtung machen wir mit Induktion über k .

$k = 1$ ist klar. Angenommen $k - 1$ Schritte sind ausgeführt. $A^{(k-1)} = M_{k-1} \cdots M_1 A$ und $a_{k,k}^{(k-1)}$ ist das k -te Pivot.

$$\underbrace{M_{k-1} \cdots M_1 A}_{M^{(k-1)}} = \left[\begin{array}{ccc|ccc} 1 & & & & & \\ * & \ddots & & & & \\ \vdots & \ddots & 1 & & & \\ \vdots & & * & 1 & & \\ \vdots & & \vdots & & \ddots & \\ * & \cdots & * & & & 1 \end{array} \right] \quad A = A^{(k-1)} = \left[\begin{array}{ccc|ccc} a_{11}^{(k-1)} & \cdots & \cdots & \cdots & \cdots & * \\ & \ddots & & & & \vdots \\ & & a_{k-1,k-1}^{(k-1)} & \cdots & \cdots & * \\ \hline & & & a_{kk}^{(k-1)} & \cdots & * \\ & & & \vdots & & \vdots \\ & & & * & \cdots & * \end{array} \right].$$

Es folgt, daß

$$\begin{aligned} \det(A^{(k-1)}(1:k, 1:k)) &= \prod_{i=1}^k a_{ii}^{(k-1)} \\ &= \underbrace{\det(M^{(k-1)}(1:k, 1:k))}_{=1} \cdot \det(A(1:k, 1:k)). \end{aligned}$$

Also folgt aus $\det(A(1:k, 1:k)) \neq 0$, daß $a_{k,k}^{(k-1)} \neq 0$ ist.

Zur Eindeutigkeit: Seien $A = L_1 R_1 = L_2 R_2$ zwei LR -Zerlegungen von A , A nichtsingulär, dann sind L_i, R_i , $i = 1, 2$ auch nichtsingulär. Also folgt

$$\underbrace{L_2^{-1} L_1}_{=I} = R_2 R_1^{-1} = \left[\begin{array}{c|c} \left[\begin{array}{c|c} 1 & \\ \hline & 1 \end{array} \right] & \\ \hline & \left[\begin{array}{c|c} & \\ \hline & 1 \end{array} \right] \end{array} \right]$$

$$\implies L_1 = L_2, R_1 = R_2.$$

$$\det A = \det L \cdot \det R = 1 \cdot \det R = r_{11} \cdots r_{nn}.$$

□

! Die LR -Zerlegung läßt sich auch über anderen Körpern (Ringern) durchführen.

Einige praktische Details.

- Gauß-Transformation braucht nur auf Spalten $k : n$ angewendet werden, und selbst bei Spalte k kennen wir das Ergebnis. Also Änderung von Algorithmus 2.16

$$A(k:n, k+1:n) = \text{GAUSSAPP}(A(k:n, k+1:n), t);$$

- Die Multiplikatoren, d.h. die wichtigen Elemente von t_k können auf den entstandenen Nullen von A gespeichert werden.

Damit erhält folgenden Algorithmus für die LR -Zerlegung.

Algorithmus 2.19 (Äußere-Produkt-Form der Gauß-Elimination)

Input: $A \in \mathbb{K}^{n,n}$ mit $A(1:k, 1:k)$ nichtsingulär $k = 1 : n - 1$.

Output: Faktorisierung $M_{n-1} \cdots M_1 A = R$, mit R obere Δ -Matrix und M_i Gauß-Transformationen, R wird im oberen Δ von A und die Multiplikatoren aus M_k in $A(k+1:n, k)$, d.h. $A(k+1:n, k) = -M_k(k+1:n, k)$.

```
FOR  k = 1 : n - 1
      t = GAUSS(A(k : n, k));
      A(k + 1 : n, k) = t;
      A(k : n, k + 1 : n) = GAUSSAPP(A(k : n, k + 1 : n), t);
END
```

Kosten $\frac{2n^3}{3}$ flops, jeder Durchgang durch die k -Schleife ist ein äußeres Produkt. Wir haben gesehen, daß wir die Multiplikatoren aus $t^{(k)}$ in A speichern können. Wie erhalten wir nun L (falls wir es benötigen)?

→ Übung 2.68.

$$\begin{aligned} L &= (M_{n-1} \cdots M_1)^{-1} = M_1^{-1} \cdots M_{n-1}^{-1} \\ &= (I + t^{(1)} e_1^T) \cdots (I + t^{(n-1)} e_{n-1}^T) = I + \sum_{k=1}^{n-1} t^{(k)} e_k^T. \end{aligned}$$

Also folgt $L(k+1:n, k) = t^{(k)}$. Die Lösung eines linearen Gleichungssystems folgt nun basierend auf der LR -Zerlegung wie oben beschrieben.

! Im wesentlichen enthält Algorithmus 2.19 3 Schleifen die ineinandergeschachtelt sind. Je nach Rechnerarchitektur ist es eine andere als die beschriebene Anordnung. Dies kann zu erheblichen Einsparungen führen. Als Beispiel sei hier die sogenannte *GAXPY* LR -Zerlegung betrachtet:

$$\begin{array}{l} \text{GAXPY} \quad (y = Ax + y \quad y \in \mathbb{R}^m, A \in \mathbb{R}^{m,n}, x \in \mathbb{R}^n) \\ \text{(general } Ax \text{ plus } y) \end{array}$$

(Dies ist eine Routine von den Basic Linear Algebra Subroutines(BLAS), die auf allen Standardrechnern in sehr guter Weise implementiert sind, vektor/parallel und auf die man seine Codes aufbauen sollte.)

Betrachte die Schleifen von Algorithmus 2.19 ausgeschrieben.

```
FOR  k = 1 : n - 1
      A(k + 1 : n, k) = A(k + 1 : n, k) / A(k, k);
      FOR  j = k + 1 : n
            FOR  i = k + 1 : n
                  A(i, j) = A(i, j) - A(i, k) * A(k, j);
            END
      END
END
```

In der veränderten Variante wird die Transformation mit M_k nicht direkt ausgeführt wie in Algorithmus 2.19 sondern $A(:, j)$ wird erst im Schritt j mit $M_{j-1} \cdots M_1 A(:, j)$ überschrieben und dann M_j berechnet. Genauer, sei $1 \leq j \leq n - 1$ und $L(:, 1:j-1), R(1:j-1, 1:j-1)$

seien bekannt. Um die j -ten Spalten von L, R zu erhalten, setzen wir die j -ten Spalten von $A = LR$ gleich

$$\begin{aligned} A(1 : j - 1, j) &= L(1 : j - 1, 1 : j - 1)R(1 : j - 1, j), \\ A(j : n, j) &= \sum_{k=1}^j L(j : n, k)R(k, j). \end{aligned}$$

Die erste Gleichung ist ein Gleichungssystem mit unterer Δ -Matrix, daß gelöst werden muß. Die zweite Gleichung schreiben wir wie folgt.

$$\begin{aligned} v(j : n) &= A(j : n, j) - \sum_{k=1}^{j-1} L(j : n, k)R(k, j) \\ &= A(j : n, j) - L(j : n, 1 : j - 1)R(1 : j - 1, j). \end{aligned}$$

Dann gilt $R(j, j) = v(j)$ und $L(j + 1 : n, j) = v(j + 1 : n)/R(j, j)$.

Also ist $L(j + 1 : n, j)$ nichts anderes als eine skalierte GAXPY Operation und wir erhalten:

```
FOR   j = 1 : n
      Solve L(1 : j - 1, 1 : j - 1)R(1 : j - 1, j) = A(1 : j - 1, j);
      v(j : n) = A(j : n, j) - L(j : n, 1 : j - 1)R(1 : j - 1, j)
      R(j, j) = v(j);
      L(j + 1 : n, j) = v(j + 1 : n)/R(j, j);
END
```

Zusammengefaßt erhält man dann

Algorithmus 2.20 (GAXPY Form der Gauß-Elimination)

Input: $A \in \mathbb{K}^{n,n}$ mit $A(1 : k, 1 : k)$ nichtsingulär $k = 1 : n - 1$

Output: Faktorisierung $A = LR$ mit L untere Δ -Matrix mit 1-Diagonale, R obere Δ -Matrix. Falls $i > j$, so ist $L(i, j)$ in $A(i, j)$ gespeichert. Falls $i \leq j$, so ist $R(i, j)$ in $A(i, j)$ gespeichert.

```
FOR   j = 1 : n
      % Solve L(1 : j - 1, 1 : j - 1)R(1 : j - 1, j) = A(1 : j - 1, j)
      FOR   k = 1 : j - 1
            FOR   i = k + 1 : j - 1
                  A(i, j) = A(i, j) - A(i, k) * A(k, j);
            END
      END
      % v(j : n) = A(j : n, j) - L(j : n, 1 : j - 1) * R(1 : j - 1, j)
      FOR   k = 1 : j - 1
            FOR   i = j : n
                  A(i, j) = A(i, j) - A(i, k) * A(k, j);
            END
      END
      % L(j + 1 : n, j) = v(j + 1 : n)/v(j)
      A(j + 1 : n, j) = A(j + 1 : n, j)/A(j, j);
END
```

Kosten $\frac{2n^3}{3}$ flops.

Weitere Varianten (andere Schleifenorientierung) Crout, Prolittle, Blockvarianten, siehe Buch von Golub/Van Loan: 'Matrix Computations'.

2.1.3 Fehleranalyse der Gauß-Elimination

Anwendung der dargestellten Algorithmen zur Lösung von $Ax = b$. Fehleranalyse i.a. sehr schwierig, bei Gauß sehr viel bekannt, für andere Algorithmen teils sehr wenig. Zuerst müssen wir nun eine Störungsanalyse für Gleichungssysteme machen. Betrachte das parametrisierte System

$$(A + \varepsilon F)x(\varepsilon) = b + \varepsilon f, \quad x(0) = x, \quad F \in \mathbb{R}^{n,n}, f \in \mathbb{R}^n.$$

Falls A nichtsingulär ist, so ist $x(\varepsilon)$ differenzierbar in einer Umgebung von 0 und es gilt:

$$\dot{x}(0) = A^{-1}(f - Fx).$$

Taylorentwicklung liefert

$$x(\varepsilon) = x + \varepsilon \dot{x}(0) + \mathcal{O}(\varepsilon^2).$$

Also folgt für jede Vektornorm und konsistente Matrixnorm

$$\frac{\|x(\varepsilon) - x\|}{\|x\|} \leq |\varepsilon| \|A^{-1}\| \left\{ \frac{\|f\|}{\|x\|} + \|F\| \right\} + \mathcal{O}(\varepsilon^2). \quad (2.21)$$

Für quadratische Matrizen definiere die Kondition $\kappa_{\|\bullet\|}(A) := \|A\| \|A^{-1}\|$, diese ist abhängig von der Norm, $\kappa(A) = \infty$ für A singulär. Mit der Konsistenzungleichung $\|b\| \leq \|A\| \|x\|$ folgt

$$\frac{\|x(\varepsilon) - x\|}{\|x\|} \leq \kappa(A)(\rho_A + \rho_b) + \mathcal{O}(\varepsilon^2), \quad (2.22)$$

wobei

$$\rho_A = \varepsilon \frac{\|F\|}{\|A\|}, \quad \rho_b = \varepsilon \frac{\|f\|}{\|b\|}$$

die relativen Fehler in A, b sind.

Die Abschätzung (2.22) ist noch unbefriedigend, da sie auf „ ε klein“ beruht, eine genaue Analyse erhält man mit den folgenden Resultaten.

Lemma 2.23 Sei $F \in \mathbb{R}^{n,n}$ und $\|\bullet\|$ eine konsistente Norm. Falls $\|F\| < 1$, so ist $I - F$ nicht singulär und es gilt

$$(I - F)^{-1} = \sum_{k=0}^{\infty} F^k \quad \text{Neumannsche Reihe} \quad (2.24)$$

und

$$\|(I - F)^{-1}\| \leq \frac{1}{1 - \|F\|}. \quad (2.25)$$

Beweis: Angenommen $I - F$ singulär $\implies \exists x \neq 0$ mit

$$(I - F)x = 0 \implies \|x\| = \|Fx\| \leq \|F\| \|x\| \implies \|F\| \geq 1.$$

Widerspruch!

$$\left(\sum_{k=0}^N F^k \right) (I - F) = I - F^{N+1}.$$

Da $\|F\| < 1$ und $\|F^k\| \leq \|F\|^k$, so gilt $\lim_{k \rightarrow \infty} F^k = 0$.

$$\begin{aligned} \implies & \left(\lim_{N \rightarrow \infty} \sum_{k=0}^N F^k \right) (I - F) = I \\ \implies & \lim_{N \rightarrow \infty} \sum_{k=0}^N F^k = (I - F)^{-1} \end{aligned}$$

und

$$\|(I - F)^{-1}\| \leq \sum_{k=0}^{\infty} \|F\|^k = \frac{1}{1 - \|F\|}.$$

□

Mit diesem Lemma erhält man

Lemma 2.26 *Sei*

$$\begin{aligned} Ax &= b, & A \in \mathbb{R}^{n,n}, A \text{ nichtsingulär}, 0 \neq b \in \mathbb{R}^n \\ (A + \Delta A)y &= b + \Delta b, & \Delta A \in \mathbb{R}^{n,n}, \Delta b \in \mathbb{R}^n, \end{aligned}$$

mit $\|\Delta A\| \leq \delta \|A\|$, $\|\Delta b\| \leq \delta \|b\|$. Falls $\delta \kappa(A) = r < 1$, so ist $A + \Delta A$ nichtsingulär und

$$\frac{\|y\|}{\|x\|} \leq \frac{1+r}{1-r}.$$

Beweis: $\|A^{-1} \Delta A\| \leq \delta \|A^{-1}\| \|A\| = r < 1$. $\implies A + \Delta A$ nichtsingulär nach Lemma 2.23. $(I + A^{-1} \Delta A)y = x + A^{-1} \Delta b$ ergibt

$$\begin{aligned} \|y\| &\leq \|(I + A^{-1} \Delta A)^{-1}\| (\|x\| + \delta \|A^{-1}\| \|b\|) \\ &\leq \frac{1}{1-r} (\|x\| + \delta \|A^{-1}\| \|Ax\|) \\ &\leq \frac{1}{1-r} \left(\|x\| + \underbrace{\delta \|A^{-1}\| \|A\|}_{r} \|x\| \right). \end{aligned}$$

□

Damit haben wir jetzt ein genaues Störungsresultat:

Satz 2.27 *Unter den Voraussetzungen von Lemma 2.26 gilt*

$$\frac{\|x - y\|}{\|x\|} \leq \frac{2r}{1-r}. \quad (2.28)$$

Beweis: $y - x = A^{-1} \Delta b - A^{-1} \Delta Ay$

$$\begin{aligned} \implies \|y - x\| &\leq \delta \|A^{-1}\| \|b\| + \delta \|A^{-1}\| \|A\| \|y\| \\ &\leq \delta \|A^{-1}\| \|A\| \|x\| + \delta \|A^{-1}\| \|A\| \|y\| = r(\|x\| + \|y\|) \end{aligned}$$

und damit ergibt sich mit Lemma 2.26

$$\frac{\|y - x\|}{\|x\|} \leq r \left(1 + \frac{1+r}{1-r} \right).$$

□

Es folgt, daß der relative Fehler abhängig von der Kondition ist. Es kann allerdings sein, daß die Schranke in (2.28) eine grobe Überschätzung ist.

Beispiel 2.29

$$\begin{bmatrix} 1 & 0 \\ 0 & 10^{-6} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 10^{-6} \end{bmatrix}, \implies x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

$\kappa_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty = 10^6$. Sei $\Delta b = \begin{bmatrix} 10^{-7} \\ 0 \end{bmatrix}$, $\Delta A = 0$, so ist $\delta = 10^{-7}$, $r = 0.1$

und (2.28) ergibt $\frac{10^{-7}}{10^0} \leq \frac{2}{0.9} 0.1 = \frac{2}{9}$. Jedoch für $\Delta b = \begin{bmatrix} 0 \\ 10^{-7} \end{bmatrix}$, $\Delta A = 0$ ist weiterhin $\delta = 10^{-7}$, $r = 0.1$ und (2.28) lautet $\frac{10^{-1}}{10^0} \leq \frac{2}{9}$, also fast Gleichheit.

Eine verbesserte Fehleranalyse ist möglich, wenn man komponentenweise Störungsresultate verwendet.

Lemma 2.30 Sei

$$\begin{aligned} Ax &= b, & A \in \mathbb{R}^{n,n}, A \text{ nichtsingulär}, 0 \neq b \in \mathbb{R}^n \\ (A + \Delta A)y &= b + \Delta b, & \Delta A \in \mathbb{R}^{n,n}, \Delta b \in \mathbb{R}^n, \end{aligned}$$

mit $|\Delta A| \leq \delta|A|$, $|\Delta b| \leq \delta|b|$. Falls $\delta \| |A^{-1}| |A| \| = r < 1$, so ist $A + \Delta A$ nichtsingulär und

$$\frac{\|y\|}{\|x\|} \leq \frac{1+r}{1-r},$$

wobei $\|\bullet\| = \|\bullet\|_\infty, \|\bullet\|_1, \|\bullet\|_F$.

→ Übung 2.65.

Satz 2.31 Unter den Voraussetzungen von Lemma 2.30 gilt

$$\frac{\|y - x\|}{\|x\|} \leq \frac{2r}{1-r}, \tag{2.32}$$

wobei $\|\bullet\| = \|\bullet\|_\infty, \|\bullet\|_1, \|\bullet\|_F$.

→ Übung 2.66.

Alle Resultate lassen sich auch für \mathbb{C} zeigen. Komponentenweise Abschätzungen sind i.a. schwieriger aber oft aussagefähiger.

Beispiel 2.33 Wir betrachten Beispiel 2.29. Hier ist $\| |A^{-1}| |A| \|_\infty = 1$. Sei $\Delta b = \begin{bmatrix} 10^{-7} \\ 0 \end{bmatrix}$, $\Delta A = 0$, so ist $\delta = r = 10^{-7}$ und (2.32) ergibt $\frac{10^{-7}}{10^0} \leq \frac{2}{1-10^{-7}} 10^{-7} \approx 2 \cdot 10^{-7}$.

Für $\Delta b = \begin{bmatrix} 0 \\ 10^{-7} \end{bmatrix}$, $\Delta A = 0$ ist jetzt $\delta = r = 10^{-1}$ und (2.32) ergibt $\frac{10^{-1}}{10^0} \leq \frac{2}{0.9} 10^{-1} = \frac{2}{9}$.

Es gibt Bestrebungen, die Analyse von Algorithmen ganz auf komponentenweise Abschätzungen zu basieren, dies ist noch im Fluß. Diese Störungssätze wenden wir nun an. Zu Anfang betrachten wir nur gestörte Daten und exakte Rechnung:

$$gl(A) = A + E, gl(b) = b + e,$$

und angenommen

$$(A + E)\tilde{x} = b + e, \quad \text{wobei } \|E\|_\infty \leq \text{eps} \|A\|_\infty, \|e\|_\infty \leq \text{eps} \|b\|_\infty. \quad (2.34)$$

Störungssätze ergeben mit $\text{eps} \kappa_\infty(A) \leq \frac{1}{2}$

$$\frac{\|x - \tilde{x}\|_\infty}{\|x\|_\infty} \leq 4 \text{eps} \kappa_\infty(A). \quad (2.35)$$

Beides sind „bestmögliche“ Norm-schranken, d.h. keine allgemeine $\|\bullet\|_\infty$ -Analyse kann bessere Schranken liefern.

! Wir können konsequenterweise keinen Algorithmus als schlecht abtun, der ein ungenaues Ergebnis für eine Matrix mit $\text{eps} \kappa_\infty(A) \approx 1$ liefert. Für schlecht konditionierte Probleme kann also auch ein guter Algorithmus schlechte Ergebnisse liefern.

Nun zum Gauß-Verfahren. (Äußere Produktversion, Analyse gilt auch für GAXPY.)

Satz 2.36 Sei A eine $n \times n$ Matrix von Maschinenzahlen. Falls kein 0-Pivot während der Ausführung von Algorithmus 2.19 auftritt, dann erfüllen die berechneten Faktoren \tilde{L}, \tilde{R}

$$\tilde{L}\tilde{R} = A + H \quad (2.37)$$

mit

$$|H| \leq 3(n-1) \text{eps} (|A| + |\tilde{L}||\tilde{R}|) + \mathcal{O}(\text{eps}^2) \quad (2.38)$$

Beweis: Induktion über n . $n = 1 : \checkmark$.

Angenommen Resultat sei richtig für alle $(n-1) \times (n-1)$ Maschinenzahl-Matrizen. Sei

$$A = \left[\begin{array}{c|c} \alpha & w^T \\ \hline v & B \end{array} \right] \begin{array}{l} \} 1 \\ \} n-1 \end{array},$$

so wird im ersten Schritt von Algorithmus 2.19 $\tilde{z} = gl(v/\alpha)$ und $\tilde{A}_1 = gl(B - \tilde{z}w^T)$ berechnet. Also folgt

$$\tilde{z} = \frac{1}{\alpha}v + f, \quad \text{mit } |f| \leq \text{eps} \frac{|v|}{\alpha} \quad (2.39)$$

und

$$\tilde{A}_1 = B - \tilde{z}w^T + F, \quad \text{mit } |F| \leq 2 \text{eps} (|B| + |\tilde{z}||w^T|) + \mathcal{O}(\text{eps}^2). \quad (2.40)$$

Dann wird die LR -Zerlegung von \tilde{A}_1 berechnet. Mit Induktion folgt $\tilde{A}_1 = \tilde{L}_1\tilde{R}_1$ mit

$$\tilde{L}_1\tilde{R}_1 = \tilde{A}_1 + H_1, \quad \text{wobei } |H_1| \leq 3(n-2) \text{eps} (|\tilde{A}_1| + |\tilde{L}_1||\tilde{R}_1|) + \mathcal{O}(\text{eps}^2). \quad (2.41)$$

Also gilt

$$\begin{aligned} \tilde{L}\tilde{R} &= \begin{bmatrix} 1 & 0 \\ \tilde{z} & \tilde{L}_1 \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ 0 & \tilde{R}_1 \end{bmatrix} \\ &= A + \begin{bmatrix} 0 & 0 \\ \alpha f & H_1 + F \end{bmatrix} = A + H. \end{aligned}$$

Aus (2.40) folgt

$$|\tilde{A}_1| \leq (1 + 2 \text{eps})(|B| + |\tilde{z}| |w|^T) + \mathcal{O}(\text{eps}^2).$$

Mit (2.40), (2.41) folgt

$$|H_1 + F| \leq 3(n-1) \text{eps} (|B| + |\tilde{z}| |w|^T + |\tilde{L}_1| |\tilde{R}_1|) + \mathcal{O}(\text{eps}^2)$$

und da $|\alpha f| \leq \text{eps} |v|$, folgt

$$|H| \leq 3(n-1) \text{eps} \left\{ \begin{bmatrix} |\alpha| & |w|^T \\ |v| & |B| \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ |\tilde{z}| & |\tilde{L}_1| \end{bmatrix} \begin{bmatrix} |\alpha| & |w|^T \\ 0 & |\tilde{R}_1| \end{bmatrix} \right\} + \mathcal{O}(\text{eps}^2).$$

□

Nun haben wir eine Analyse über die LR -Zerlegung, jetzt müssen wir noch analysieren, was die Δ -Löser machen.

Satz 2.42 *Seien \tilde{L}, \tilde{R} die berechneten LR -Faktoren aus Algorithmus 2.16 oder 2.19. Bei Verwendung von Algorithmus 2.2 bzw. 2.5 zur Lösung von $\tilde{L}y = b$ und $\tilde{R}x = \tilde{y}$ ergibt sich $(A + E)\tilde{x} = b$ mit*

$$|E| \leq n \text{eps} (3|A| + 5|\tilde{L}| |\tilde{R}|) + \mathcal{O}(\text{eps}^2). \quad (2.43)$$

Beweis: Aus (2.3) und (2.6) ergibt sich

$$\begin{aligned} (\tilde{L} + F)\tilde{y} = b & \quad |F| \leq n \text{eps} |\tilde{L}| + \mathcal{O}(\text{eps}^2), \\ (\tilde{R} + G)\tilde{x} = \tilde{y} & \quad |G| \leq n \text{eps} |\tilde{R}| + \mathcal{O}(\text{eps}^2). \end{aligned}$$

Also

$$(\tilde{L} + F)(\tilde{R} + G)\tilde{x} = (\tilde{L}\tilde{R} + F\tilde{R} + \tilde{L}G + FG)\tilde{x} = b$$

$\tilde{L}\tilde{R} = A + H$ aus Satz 2.36 mit

$$|H| \leq 3(n-1) \text{eps} (|A| + |\tilde{L}| |\tilde{R}|) + \mathcal{O}(\text{eps}^2).$$

Setze $E := H + F\tilde{R} + \tilde{L}G + FG$, so gilt

$$\begin{aligned} |E| & \leq |H| + |F| |\tilde{R}| + |\tilde{L}| |G| + \mathcal{O}(\text{eps}^2) \\ & \leq 3n \text{eps} (|A| + |\tilde{L}| |\tilde{R}|) + 2n \text{eps} (|\tilde{L}| |\tilde{R}|) + \mathcal{O}(\text{eps}^2). \end{aligned}$$

□

Wäre nicht der Term $|\tilde{L}| |\tilde{R}|$ in der Abschätzung, welcher groß sein kann, so wäre der Algorithmus rückwärts stabil. Da wir jedoch nichts gegen kleine Pivots machen können, falls diese auftauchen, kann $|\tilde{L}|, |\tilde{R}|$ sehr groß werden, und der Algorithmus ist NICHT rückwärts stabil.
→ Übung

Ausweg ist Pivotisierung.

2.1.4 Partielle Pivotisierung (Spaltenpivotisierung)

Um zu vermeiden, daß 0-Pivots oder sehr kleine Pivots auftauchen, wird jeweils in der momentanen Spalte, das betragsmäßig maximale Element gesucht und durch eine Zeilenvertauschung (Permutation) in die Diagonalposition gebracht.

Grundidee:

```

FOR   k = 1 : n - 1
      Bestimme Permutationsmatrix  $P_k$ , so daß für  $z = P_k A^{(k-1)} e_k$  gilt:
       $|z(k)| = \|z(k:n)\|_\infty$ .
      Setze  $\hat{A}^{(k-1)} = P_k A^{(k-1)}$ .
      Bestimme Gauß Transformation  $M_k$ , so daß für  $v = M_k \hat{A}^{(k-1)} e_k$  gilt:
       $v(k+1:n) = 0$ .
      Setze  $A^{(k)} = M_k P_k A^{(k-1)}$ 
END

```

Diese Vorgehensweise heißt Spaltenpivotisierung oder partielle Pivotisierung und garantiert, daß alle Multiplikatoren vom Betrag ≤ 1 sind.

$$|(P_k M_{k-1} \cdots M_1 P_1 A)_{kk}| = \max_{k \leq i \leq n} |(P_k M_{k-1} \cdots M_1 P_1 A)_{ik}|, \quad k = 1 : n - 1.$$

Falls das Pivot = 0 ist, so kann man einfach den Schritt der Elimination weglassen, da dann alle Elemente in dieser Spalte = 0 sind.

Man erhält dann den folgenden Algorithmus.

Algorithmus 2.44 (Gauß–Elim. mit partieller Pivotisierung, Äußere Produktversion)

Input: $A \in \mathbb{K}^{n,n}$.

Output: Gauß–Transformationen M_1, \dots, M_{k-1} und Permutationen P_1, \dots, P_{n-1} , so daß $M_{n-1} P_{n-1} \cdots M_1 P_1 A = R$ obere Δ -Matrix.

Keiner der Multiplikatoren hat Betrag > 1 . $A(k+1:n, k)$ wird durch $-M_k(k+1:n, k)$, $k = 1 : n - 1$ überschrieben, $A(1:k, k)$ wird durch $R(1:k, k)$, $k = 1 : n$ überschrieben. Im Pivotvektor $piv(1:n-1)$ werden die Vertauschungen gespeichert. P_k vertauscht die Zeilen k und $piv(k)$, $k = 1 : n - 1$.

```

FOR   k = 1 : n - 1
      Bestimme  $\mu$  ( $k \leq \mu \leq n$ ) mit  $|A(\mu, k)| = \|A(k:n, k)\|_\infty$ 
       $A(k, k:n) \leftrightarrow A(\mu, k:n)$ 
       $piv(k) = \mu$ ;
      IF    $A(k, k) \neq 0$ 
           $t = GAUSS(A(k:n, k))$ ;
           $A(k+1:n, k) = t$ ;
           $A(k:n, k+1:n) = GAUSSAPP(A(k:n, k+1:n), t)$ ;
      END
END

```

Kosten Für die Bestimmung aller μ 's: $\mathcal{O}(n^2)$ Vergleiche, $\frac{2n^3}{3}$ flops.

Für das lineare Gleichungssystem $Ax = b$ macht man also

- $y = M_{n-1} P_{n-1} \cdots M_1 P_1 b$
- Löse $Rx = y$.

Sämtliche Informationen in A und piv .

Beispiel 2.45

$$\begin{aligned}
 A &= \begin{bmatrix} 3 & 17 & 10 \\ 2 & 4 & -2 \\ 6 & 18 & -12 \end{bmatrix} \\
 P_1 &= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} \\
 P_1 A &= \begin{bmatrix} 6 & 18 & -12 \\ 2 & 4 & -2 \\ 3 & 17 & 10 \end{bmatrix} \quad \text{piv}(1) = 3 \\
 M_1 &= \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{3} & 1 & 0 \\ -\frac{1}{2} & 0 & 1 \end{bmatrix} \quad M_1 P_1 A = \begin{bmatrix} 6 & 18 & -12 \\ 0 & -2 & 2 \\ 0 & 8 & 16 \end{bmatrix} \\
 P_2 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad \text{piv}(2) = 3 \\
 M_2 &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & \frac{1}{4} & 1 \end{bmatrix} \quad R = \begin{bmatrix} 6 & 18 & -12 \\ 0 & 8 & 16 \\ 0 & 0 & 6 \end{bmatrix}
 \end{aligned}$$

Was passiert mit L ?

Satz 2.46 Man verwende Gauß-Elimination mit partieller Pivotisierung zur Berechnung von

$$M_{n-1} P_{n-1} \cdots M_1 P_1 A = R \quad (2.47)$$

mit Algorithmus 2.44, so gilt

$$PA = LR,$$

mit $P = P_{n-1} \cdots P_1$ (Permutationsmatrix), L untere Δ -Matrix mit 1-Diagonale, $|l_{ij}| \leq 1$. Die k -te Spalte von L ist unterhalb der Diagonalen eine permutierte Version des k -ten Gauß-Vektors.

Speziell für $M_k = I - t^{(k)} e_k^T$ gilt

$$L(k+1:n, k) = g(k+1:n), \quad g = P_{n-1} \cdots P_{k+1} t^{(k)}.$$

Beweis: Aus (2.47) folgt

$$\tilde{M}_{n-1} \cdots \tilde{M}_1 P A = R \quad \text{mit}$$

$$\begin{aligned}
 \tilde{M}_{n-1} &= M_{n-1}, \\
 \tilde{M}_k &= P_{n-1} \cdots P_{k+1} M_k P_{k+1} \cdots P_{n-1}, \quad k \leq n-2.
 \end{aligned}$$

Da P_j nur Zeilen j und $\mu \geq j$ vertauscht, so folgt $P_j(1:j-1, 1:j-1) = I_{j-1}$. Also ist \tilde{M}_k eine Gauß-Transformation mit Gauß-Vektor $\tilde{t}^{(k)} = P_{n-1} \cdots P_{k+1} t^{(k)}$. \square

Durch die einfache Ersetzung von

$$\begin{aligned}
 A(k:k:n) &\longleftrightarrow A(\mu, k:n) \text{ durch} \\
 A(k, 1:n) &\longleftrightarrow A(\mu, 1:n) \text{ in Algorithmus 2.44}
 \end{aligned}$$

gilt dann wieder $A(i, j)$ enthält $L(i, j)$ für $i > j$ nach Ende des Algorithmus.

→ Übung 2.76: Analoge GAXPY Version.

Fehleranalyse: Da Permutationen rundungsfehlerfrei durchgeführt werden können, kann man analog zeigen, daß die berechnete Lösung \tilde{x} $(A + E)\tilde{x} = b$ erfüllt mit

$$\|E\| \leq n \text{ eps} \left(3\|A\| + 5\tilde{P}^T|\tilde{L}|\tilde{R} \right) + \mathcal{O}(\text{eps}^2), \quad (2.48)$$

wobei $\tilde{P}, \tilde{L}, \tilde{R}$ die berechneten P, L, R sind.

Durch die Pivotisierung folgt

$$\begin{aligned} \|\tilde{L}\|_\infty &\leq n, \text{ also folgt} \\ \|E\|_\infty &\leq n \text{ eps} \left(3\|A\|_\infty + 5n\|\tilde{R}\|_\infty \right) + \mathcal{O}(\text{eps}^2). \end{aligned}$$

Problem: Schranke für $\|\tilde{R}\|_\infty$.

Definiere den Wachstumsfaktor ρ durch

$$\rho = \max_{i,j,k} \frac{|\tilde{a}_{ij}^{(k)}|}{\|A\|_\infty}, \quad (2.49)$$

wobei $\tilde{A}^{(k)}$ die berechnete Version von $A^{(k)}$ ist. Dann folgt

$$\|E\|_\infty \leq 8n^3 \rho \|A\|_\infty \text{ eps} + \mathcal{O}(\text{eps}^2). \quad (2.50)$$

Die Schranke hängt also wesentlich von ρ ab (n^3 ist i.a. vernachlässigbar in der Praxis, da nicht auftretend).

ρ ist typisch von der Ordnung 10, kann aber 2^{n-1} sein.

→ Übung 2.71.

Im allgemeinen ist Gauß mit partieller Pivotisierung zuverlässig und kann relativ sorglos verwendet werden.

2.1.5 Vollständige Pivotisierung

Um den Wachstumsfaktor zu verkleinern gibt es eine Variante, in der das Pivot aus der gesamten Matrix $A^{(k-1)}(k : n, k : n)$ ausgewählt wird, d.h. wir bestimmen die Zerlegung

$$M_{n-1}P_{n-1} \cdots M_1P_1AQ_1 \cdots Q_{n-1} = R$$

mit Permutationsmatrizen P_i, Q_i , die so bestimmt werden, daß

$$\left| (P_k A^{(k-1)} Q_k)_{kk} \right| = \max_{k \leq i, j \leq n} \left| (P_k A^{(k-1)} Q_k)_{ij} \right|.$$

Satz 2.51 *Bei Verwendung von Gauß-Elimination mit vollständiger Pivotisierung zur Berechnung von*

$$M_{n-1}P_{n-1} \cdots M_1P_1AQ_1 \cdots Q_{n-1} = R \quad (2.52)$$

gilt

$$PAQ = LR$$

mit $P = P_{n-1} \cdots P_1$, $Q = Q_1 \cdots Q_{n-1}$ und L ist untere Δ -Matrix mit 1-Diagonale und $|l_{ij}| \leq 1$. Für $M_k = I - t^{(k)} e_k^T$ gilt

$$L(k+1 : n, k) = g(k+1 : n) \text{ mit } g = P_{n-1} \cdots P_{k+1} t^{(k)}.$$

Beweis: Analog zu Satz 2.46. □

Algorithmus 2.53 (Gauß–Elim. mit vollständiger Pivotisierung)

Input: $A \in \mathbb{K}^{n,n}$.

Output: LR -Zerlegung von PAQ mit P, Q Produkte von elementaren Permutationsmatrizen, $A(1 : k, k)$ wird durch $R(1 : k, k)$, $k = 1 : n$ und $A(k + 1 : n, k)$ durch $L(k + 1 : n, k)$, $k = 1 : n - 1$ überschrieben. P_k vertauscht Zeilen k und $p(k)$, Q_k vertauscht Spalten k und $q(k)$.

FOR $k = 1 : n - 1$

Bestimme μ, λ , so daß $|A(\mu, \lambda)| = \max \{|A(i, j)| : i, j = k : n\}$

$A(k, 1 : n) \longleftrightarrow A(\mu, 1 : n)$

$A(1 : n, k) \longleftrightarrow A(1 : n, \lambda)$

$p(k) = \mu;$

$q(k) = \lambda;$

IF $A(k, k) \neq 0$

$t = GAUSS(A(k : n, k));$

$A(k + 1 : n, k) = t;$

$A(k : n, k + 1 : n) = GAUSSAPP(A(k : n, k + 1 : n), t);$

END

END

Kosten $\frac{2n^3}{3}$ flops und Vergleiche. Vergleiche nicht vernachlässigbar.

- Vollständige Pivotisierung erlaubt LR -Zerlegung für $\text{rank } A = r < n$.
- In exakter Arithmetik gilt für Algorithmus 2.53

$$\left| a_{ij}^{(k)} \right| \leq k^{\frac{1}{2}} (2 \cdot 3^{\frac{1}{2}} \dots k^{\frac{1}{k-1}})^{\frac{1}{2}} \max |a_{ij}|.$$

Die Schranke ist sehr langsam wachsend mit k . Mit der empirischen Tatsache, daß $\rho \approx 10$, kann man aussagen, daß Gauß–Elimination mit vollständiger Pivotisierung rückwärts stabil ist, d.h. die Methode löst exakt

$$(A + E)\hat{x} = b$$

für kleines E . Wilkinsons Vermutung Schranke = n ist nicht richtig.

2.1.6 Abschätzung der Genauigkeit

Das Residuum der berechneten Lösung von $Ax = b$ ist der Vektor $b - A\tilde{x}$. Ein kleines Residuum bedeutet daß $A\tilde{x}$ eine gute Näherung von b ist. Wenn man annimmt, daß $(A + E)\tilde{x} = b$, $\|E\|_\infty \approx \text{eps } \|A\|_\infty$, so folgt $\|b - A\tilde{x}\|_\infty \approx \text{eps } \|A\|_\infty \|\tilde{x}\|_\infty$.

Heuristik I Gauß–Elimination erzeugt eine Lösung \tilde{x} mit relativ kleinem Residuum.

! Jedoch kleine Residuen implizieren nicht hohe Genauigkeit, (nur bei kleiner Kondition).

$$\frac{\|\tilde{x} - x\|_\infty}{\|x\|_\infty} \approx \text{eps } \kappa_\infty(A).$$

Heuristik II Falls $\text{eps} \approx 10^{-d}$ und $\kappa_\infty(A) \approx 10^q$, dann erzeugt Gauß–Elimination eine Lösung mit ungefähr $d - q$ korrekten Dezimalstellen.

Beispiel 2.54

$$\begin{bmatrix} .986 & .579 \\ .409 & .237 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} .235 \\ .107 \end{bmatrix}, \quad \kappa(A)_\infty \approx 700, \quad x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}.$$

eps	\tilde{x}_1	\tilde{x}_2	$\frac{\ \tilde{x}-x\ }{\ x\ _\infty}$	$\frac{\ b-A\tilde{x}\ _\infty}{\ A\ _\infty\ \tilde{x}\ _\infty}$
10^{-3}	2.11	-3.17	$5 \cdot 10^{-2}$	$2.0 \cdot 10^{-3}$
10^{-4}	1.986	-2.975	$8 \cdot 10^{-3}$	$1.5 \cdot 10^{-4}$
10^{-5}	2.0019	-3.0032	$1 \cdot 10^{-3}$	$2.1 \cdot 10^{-6}$
10^{-6}	2.00025	-3.00094	$3 \cdot 10^{-4}$	$4.2 \cdot 10^{-7}$

2.1.7 Skalierung

Eine weitere Verbesserung ist durch Skalierung der Daten zu erreichen.

$$Ax = b \iff D_1^{-1}AD_2y = D_1^{-1}b, \quad y = D_2^{-1}x$$

Falls man für D_1, D_2 Diagonalmatrizen aus Maschinenzahlen der Form

$$\text{diag}(P^{r_1}, P^{r_2}, \dots, P^{r_n})$$

nimmt, so kann die Skalierung ohne Rundungsfehler in $\mathcal{O}(n^2)$ flops durchgeführt werden. Dann gilt

$$\frac{\|D_2^{-1}(\tilde{x} - x)\|_\infty}{\|D_2^{-1}x\|_\infty} = \frac{\|\tilde{y} - y\|_\infty}{\|y\|_\infty} \approx \text{eps} \kappa_\infty(D_1^{-1}AD_2). \quad (2.55)$$

Wenn man also $\kappa_\infty(D_1^{-1}AD_2)$ gegen $\kappa_\infty(A)$ verkleinert, erwartet man eine Verbesserung des Resultats.

Varianten:

- Zeilenskalierung: $D_2 = I$, D_1 so, daß alle Zeilen \approx gleiche ∞ -Norm haben.
- Zeilen-Spalten-Gleichgewichtung: Wähle D_1, D_2 so, daß alle Zeilen und Spalten ∞ -Norm im Intervall $\left[\frac{1}{p}, 1\right]$ haben (p Basis der Maschine).

Beispiel 2.56

$$\begin{bmatrix} 10 & 100000 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 100000 \\ 2 \end{bmatrix} \iff \begin{bmatrix} 0.0001 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Bei dreistelliger Arithmetik, $p = 10$ ergibt das erste System $\tilde{x} = \begin{bmatrix} 0.00 \\ 1.00 \end{bmatrix}$ und das zweite $\begin{bmatrix} 1.00 \\ 1.00 \end{bmatrix}$. Exakte Lösung $x = \begin{bmatrix} 1.0001\dots \\ 0.9999\dots \end{bmatrix}$.

2.1.8 Iterative Verbesserung

Angenommen, wir haben $Ax = b$ mittels Gauß-Elimination mit partieller Pivotisierung gelöst, $PA = LR$, und wollen jetzt die Lösungsgenauigkeit verbessern. Wir könnten folgendes ausführen:

$$\begin{cases} r = b - A\tilde{x} & \text{(doppelt genau)}, \\ \text{Löse } Ly = Pr, \\ \text{Löse } Rz = y, \\ \text{Setze } x = \tilde{x} + z, \end{cases} \quad (2.57)$$

so folgt bei exakter Rechnung

$$Ax = A\tilde{x} + Az = (b - r) + r = b.$$

Wenn man (2.57) allerdings einfach so ausführt, ist das Ergebnis nicht genauer als \tilde{x} . Dies ist zu erwarten, denn $\tilde{r} = gl(b - A\tilde{x})$ hat im allgemeinen kaum richtige signifikante Stellen (siehe Heuristik I). Also

$$\tilde{z} = gl(A^{-1}r) = A^{-1} \cdot \text{Rauschen} = \text{Rauschen}.$$

! Um eine Verbesserung zu erhalten, sollte man daher $r = b - Ax$ in doppelter Genauigkeit rechnen (doppelt so viele Stellen wie sonst). Dieser Prozeß kann iteriert werden.

Heuristik III Bei Maschinengenauigkeit $\text{eps} = 10^{-d}$ und $\kappa_{\infty}(A) \approx 10^q$, hat nach k Schritten von (2.57) (mit doppelter Genauigkeit für das Residuum) das berechnete x ungefähr $\min\{d, k(d - q)\}$ korrekte Stellen.

Kosten $\mathcal{O}(n^2)$ pro Schritt.

Bei Rechnung von $PA = LR$ in einfacher Genauigkeit und $\text{eps} \cdot \kappa_{\infty}(A) \leq 1$, ist die Lösung korrekt in einfacher Genauigkeit nach einer Iteration. Also kaum Verteuerung und Verbesserung der Lösung. Problem: maschinenabhängig, da an Genauigkeit orientiert.

→ Übung 2.75: Konditionsschätzer

2.1.9 Übungen zu Kapitel 2.1

Übung 2.58 Schreibe einen **MATLAB**-Algorithmus zur Lösung eines Gleichungssystems $Lx = b$ mit einer unteren Dreiecksmatrix L (Algorithmus 2.2 der Vorlesung). Verfahre analog bei $Ux = b$ wobei U obere Dreiecksmatrix ist (Algorithmus 2.5 der Vorlesung).

Übung 2.59 Schreibe einen **MATLAB**-Algorithmus zur Berechnung der Determinante einer $n \times n$ -Matrix mit Hilfe des Laplaceschen Entwicklungssatzes. Der Laplacesche Entwicklungssatz zur Determinantenentwicklung (nach der ersten Zeile) lautet:

$\det A = \sum_{j=1}^n (-1)^{j+1} a_{1j} \det A(2 : n, [1 : j - 1, j + 1 : n])$. Diese Formel soll natürlich rekursiv ohne 'det' aus **MATLAB** angewendet werden.

Übung 2.60 Vergleiche folgende vier Verfahren in **MATLAB** zur Lösung von $Ax = b$, $A \in \mathbb{R}^{n,n}$, $b \in \mathbb{R}^n$. Messe die Rechenzeit für $n = 4, 6, 8, 10$ und zufällig erzeugte A, b .

1. den `\` aus **MATLAB**. (Hinter diesem steckt der Gauß-Alg.)
2. Einen selbstgeschriebenen Gauß-Algorithmus entsprechend der Vorlesung
3. Die Cramersche Regel mit Hilfe der eingebauten 'det'-Funktion
4. Die Cramersche Regel mit Hilfe des Laplaceschen Entwicklungssatzes

Die Cramersche Regel für $Ax = b$ ist wie folgt definiert:

$$x_i = \det [A(:, 1 : i - 1), b, A(:, i + 1 : n)] / \det A, \quad i = 1, \dots, n$$

Der Laplacesche Entwicklungssatz zur Determinantenentwicklung (nach der ersten Zeile) lautet:

$\det A = \sum_{j=1}^n (-1)^{j+1} a_{1j} \det A(2 : n, [1 : j - 1, j + 1 : n])$. Diese Formel soll natürlich rekursiv ohne 'det' aus **MATLAB** angewendet werden.

Übung 2.61 Die exakte Lösung des Gleichungssystems

$$\begin{pmatrix} 1 & 0.99 \\ 0.99 & 0.98 \end{pmatrix} x = \begin{pmatrix} 1.99 \\ 1.97 \end{pmatrix}$$

ist $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$. Mit $y = \begin{pmatrix} 1.94 \\ 0.050459 \end{pmatrix}$ gilt aber auch

$$\begin{pmatrix} 1 & 0.99 \\ 0.99 & 0.98 \end{pmatrix} y = \begin{pmatrix} 1.98995441 \\ 1.97004982 \end{pmatrix}.$$

Erkläre dieses Verhalten anhand der Störungstheorie der Vorlesung.

Übung 2.62 Schreibe ein **MATLAB**-Programm, welches den Gauß-Algorithmus der Vorlesung ohne Pivotisierung auf folgendes Gleichungssystem (Hilbertmatrix, in **MATLAB** 'hilb(n)') anwendet.

$$Hx = b, \quad H = \left(\frac{1}{i+j-1} \right)_{i,j=1,\dots,n}, \quad b = H(1, \dots, 1)^T.$$

Vergleiche die Lösung des Algorithmus mit der exakten Lösung $x = (1, \dots, 1)^T$ sowie der von **MATLAB** mittels \ berechneten Lösung für $n = 2, 3, 4, \dots, 30$. Was stellst Du fest? Erklärung!

Übung 2.63 Eine Matrix $A = (a_{ij})_{i,j=1,\dots,n}$ heißt strikt diagonaldominant, wenn gilt

$$|a_{jj}| - \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| > 0, \quad \text{für alle } j = 1, \dots, n.$$

Zeige:

1. Ist A strikt diagonaldominant, dann ist A nicht singulär.
2. Ist A strikt diagonaldominant, dann gilt das auch für alle Hauptabschnittsmatrizen.
3. Ist A strikt diagonaldominant, dann besitzt A eine LR-Zerlegung.

Übung 2.64 Sei

$$A = \begin{pmatrix} a_{11} & a_{12} & & & \\ a_{21} & \ddots & \ddots & & \\ & \ddots & \ddots & & a_{n-1,n} \\ & & a_{n,n-1} & a_{n,n} & \end{pmatrix}$$

tridiagonal. Wieviel flops (+, -, *, /) braucht der Gauß-Algorithmus, falls nicht permutiert werden muß und der Algorithmus durchgeht? Wie teuer sind Vorwärts- und Rückwärtseinsetzen?

Übung 2.65 Zeige Lemma 2.30 aus der Vorlesung: Sei

$$\begin{aligned} Ax &= b, & A &\in \mathbb{R}^{n,n}, 0 \neq b \in \mathbb{R}^n \\ (A + \Delta A)y &= b + \Delta b, & \Delta A &\in \mathbb{R}^{n,m}, \Delta b \in \mathbb{R}^n, \end{aligned}$$

mit $|\Delta A| \leq \delta |A|, |\Delta b| \leq \delta |b|$. Falls $\delta \| |A^{-1}| |A| \|_{\infty} = r < 1$, so ist $A + \Delta A$ nichtsingulär und

$$\frac{\|y\|_{\infty}}{\|x\|_{\infty}} \leq \frac{1+r}{1-r}.$$

Tip: Verfahre analog zum Beweis des entsprechenden Lemmas 2.26 über Normen.

Übung 2.66 Zeige Satz 2.31 aus der Vorlesung: Unter den Voraussetzungen von Lemma 2.30 gilt

$$\frac{\|y - x\|}{\|x\|} \leq \frac{2r}{1 - r},$$

wobei $\|\bullet\| = \|\bullet\|_\infty, \|\bullet\|_1, \|\bullet\|_F$. Tip: Verfahre analog zum Beweis des entsprechenden Satzes 2.27 über Normen.

Übung 2.67 Sei $1 > \varepsilon > 0$,

$$A = \begin{pmatrix} \varepsilon & -1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & \ddots & \\ & & & \ddots & -1 \\ & & & & \varepsilon \end{pmatrix} \in \mathbb{R}^{n,n}.$$

1. Zeige, es existiert eine singuläre Matrix B mit $\|B - A\|_\infty \leq \varepsilon^n$. D.h. der Abstand von A zu den singulären Matrizen ist höchstens ε^n . Tip: man füge in der Position $(n, 1)$ einen geeigneten Eintrag ein, so dass bei anschließender Determinantenbildung eine Null entsteht.
2. Berechne A^{-1}
3. Bestimme für $\|\bullet\|_\infty$ die Normen von A, A^{-1} , die Konditionszahl sowie die Skeel-Konditionszahl $\| |A| |A^{-1}| \|$.
4. Wende für $\varepsilon \leq \frac{1}{4}$ und $\delta = \varepsilon^n/4, b = (1, \dots, 1)^T, \Delta b = \delta \cdot b$ und $\Delta A = 0$ die Sätze 2.27, 2.31 der Vorlesung an. Wie scharf sind die Ergebnisse?

Übung 2.68 Sei

$$L = \begin{pmatrix} 1 & & & 0 \\ l_{21} & \ddots & & \\ \vdots & \ddots & \ddots & \\ l_{n1} & \cdots & l_{n,n-1} & 1 \end{pmatrix}, l_k = \left. \begin{pmatrix} 0 \\ \vdots \\ 0 \\ l_{k+1,k} \\ \vdots \\ l_{n,k} \end{pmatrix} \right\} k \quad k = 1, \dots, n-1.$$

Zeige:

1. $(I + l_k e_k^T)(I + l_m e_m^T) = I + l_k e_k^T + l_m e_m^T$ für alle $k \leq m$. Hier bezeichnen e_1, \dots, e_n die Einheitsvektoren.
2. $L = (I + l_1 e_1^T) \cdots (I + l_{n-1} e_{n-1}^T)$.
3. $L^{-1} = (I - l_{n-1} e_{n-1}^T) \cdots (I - l_1 e_1^T)$.

Übung 2.69 Zähle die genaue Anzahl flops in der GAXPY-Version des Gauß-Algorithmus (Algorithmus 2.20 der Vorlesung). Zeige dass $\frac{2}{3}n^3 + \mathcal{O}(n^2)$ flops benötigt werden.

Übung 2.70 Schreibe **MATLAB**-Programme, welche den Gauß-Algorithmus der Vorlesung mit partieller Pivotisierung sowie vollständiger Pivotisierung realisieren. Verwende als Testbeispiel $Hx = b$ wieder die Hilbertmatrix $H, x = (1, \dots, 1)^T$ für $n = 2, 3, 4, \dots, 30$.

Übung 2.71 Bestimme den Wachstumsfaktor $\rho = \max_{i,j,k} \frac{|a_{ij}^{(k)}|}{\|A\|_\infty}$ aus der Vorlesung bei exakter Arithmetik und partieller Pivotisierung für die Matrix

$$\begin{pmatrix} 1 & 0 & \cdots & 0 & 1 \\ -1 & 1 & \ddots & \vdots & 1 \\ -1 & -1 & \ddots & 0 & \vdots \\ \vdots & & \ddots & \ddots & \vdots \\ -1 & \cdots & \cdots & -1 & 1 \end{pmatrix}.$$

Übung 2.72 Sei $A \in \mathbb{R}^{n,n}$, $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ und nehme an, daß A und $A_{11} \in \mathbb{R}^{k,k}$ nicht singulär seien. Das Schur-Komplement ist gegeben durch $S = A_{22} - A_{21}A_{11}^{-1}A_{12}$. Zeige:

1. A läßt sich schreiben als

$$A = \begin{pmatrix} I & O \\ A_{21}A_{11}^{-1} & I \end{pmatrix} \begin{pmatrix} A_{11} & O \\ O & S \end{pmatrix} \begin{pmatrix} I & A_{11}^{-1}A_{12} \\ O & I \end{pmatrix}.$$

2. S ist nicht singulär.

3. A^{-1} läßt sich schreiben als $A^{-1} = \begin{pmatrix} * & * \\ * & S^{-1} \end{pmatrix}$ (* heißt, die Einträge interessieren uns hier nicht).

4. Ist A symmetrisch, so ist S auch symmetrisch.

5. Ist A positiv definit d.h. $x^T Ax > 0$, $\forall x \neq 0$, so ist S auch positiv definit, d.h. $y^T Sy > 0$, $\forall y \neq 0$.

6. Ist A invers nicht negativ, dann ist auch S invers nicht negativ (A heißt invers nicht negativ, falls A^{-1} nur nicht negative Einträge hat).

Übung 2.73 Sei $A \in \mathbb{R}^{n,n}$, $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ Falls $A_{11} \in \mathbb{R}^{k,k}$ nicht singulär ist, dann ist das Schur-Komplement gegeben durch $S = A_{22} - A_{21}A_{11}^{-1}A_{12}$. Zeige:

1. Ist A strikt diagonaldominant, dann (ist A_{11} nicht singulär und) S ist strikt diagonaldominant. Was bedeutet das für den Gauß-Algorithmus mit partieller Pivotisierung?

2. Ist A eine M -Matrix, dann ist auch S eine M -Matrix. (A heißt M -Matrix, falls A invers nicht negativ ist und alle Einträge von A außerhalb der Diagonalen kleiner oder gleich 0 sind).

Tip: Aus der Definition M -Matrix folgt bereits, daß die Diagonaleinträge von A positiv sind. Warum?

Übung 2.74 1. Schreibe ein **MATLAB**-Programm, welches den Gauß-Algorithmus der Vorlesung mit vollständiger Pivotisierung sowie Zeilenskalierung und iterativer Verbesserung realisiert. Verwende als Testbeispiel $Hx = b$ wieder die Hilbertmatrix H , $x = (1, \dots, 1)^T$ für $n = 2, 3, 4, \dots, 30$. (Siehe Übung 2.62).

2. Untersuche die Genauigkeitsverbesserung mittels iterativer Verbesserung dadurch, dass Daten und Rechenoperationen mit zufälligen Werten der Form $(1+\varepsilon)$ gestört werden, wobei $|\varepsilon| \leq \sqrt{\text{eps}}$ ist (einfache statt doppelte Genauigkeit).

Vergleiche die iterative Verbesserung mit und ohne gestörte Berechnung des Residuums.

Übung 2.75 Gegeben sei eine nichtsinguläre untere Dreiecksmatrix $L \in \mathbb{R}^{n,n}$.

1. Nehme an, dass $y \in \mathbb{R}^n$ gegeben ist und x die Gleichung $Lx = y$ erfüllt. Zeige: falls $c\|y\|_p \leq \|x\|_p$ gilt, so ist $\|L^{-1}\|_p \geq c$.
2. Nutze (1) zur Entwicklung eines Konditionsschätzers. Wähle ein y der Form $y^\top = (\pm 1 \ \cdots \ \pm 1)$ und betrachte folgende Variante des Vorwärtseinsetzens

```

p = 0 ∈ ℝn
for k = 1 : n
    Wähle das Vorzeichen von y(k)
    x(k) = (y(k) - p(k))/L(k, k)
    p(k + 1 : n) = p(k + 1 : n) + L(k + 1 : n, k)x(k)
end
κ = \|x\|∞ \|L\|∞

```

- (a) Schreibe ein MATLAB-Programm, in dem y so gewählt wird, dass für $k = 1, \dots, n$, $|x(k)|$ sukzessive maximiert wird.
- (b) Schreibe ein MATLAB-Programm, in dem y so gewählt wird, dass für $k = 1, \dots, n$, $|x(k)| + \|p(k + 1 : n) + L(k + 1 : n, k)x(k)\|_1$ sukzessive maximiert wird.
- (c) Vergleiche beide Varianten in MATLAB mit der eingebauten Funktion **condest** sowie der exakten Kondition **cond**. Was stellst du fest?

Übung 2.76 Schreibe ein MATLAB-Programm für die GAXPY-Version des Gauß-Algorithmus der Vorlesung mit partieller Pivotisierung. Verwende als Testbeispiel $Hx = b$ die Hilbertmatrix H , $x = (1, \dots, 1)^\top$ für $n = 2, 3, 4, \dots, 30$.

Übung 2.77 Untersuche in MATLAB die Genauigkeitsverbesserung mittels iterativer Verbesserung bei der GAXPY-Version des Gauß-Algorithmus mit partieller Pivotisierung dadurch, dass Daten und Rechenoperationen mit zufälligen Werten der Form $(1 + \varepsilon)$ gestört werden, wobei $|\varepsilon| \leq \sqrt{\text{eps}}$ ist (einfache statt doppelte Genauigkeit). Vergleiche die iterative Verbesserung mit und ohne gestörte Berechnung des Residuums. Verwende als Testbeispiel $Hx = b$ wieder die Hilbertmatrix H , $x = (1, \dots, 1)^\top$ für $n = 2, 3, 4, \dots, 30$.

2.2 Der Cholesky-Algorithmus, Bandmatrizen

Ein guter numerischer Algorithmus sollte spezielle Eigenschaften des mathematischen Problems berücksichtigen. Der Gauß-Algorithmus ist auf allgemeine lineare Gleichungssysteme zugeschnitten. Nun sind aber viele der Anwendungsproblem aus Ingenieurwissenschaft oder Physik, Chemie so, daß die Gleichungssysteme eine spezielle Struktur haben, wie zum Beispiel symmetrische, Hermitesche Matrizen, Bandstruktur etc. Wir betrachten nun zuerst Hermitesch, positiv definite Systeme

$$Ax = b \text{ mit } A = A^*, \text{ d.h. } a_{ij} = \bar{a}_{ij}, \quad (2.78)$$

$$\text{und } x^* Ax > 0, \forall x \in \mathbb{C}^n \setminus \{0\}$$

und untersuchen, wie wir den Gauß-Algorithmus verbessern können.

Satz 2.79 (Cholesky Zerlegung)

Sei $A \in \mathbb{K}^{n,n}$ ($\mathbb{K} = \mathbb{C}$ oder $\mathbb{K} = \mathbb{R}$), Hermitesch positiv definit. Dann existiert eine untere Δ -Matrix $G \in \mathbb{K}^{n,n}$ mit positiven Diagonalelementen, so daß

$$A = GG^*.$$

Beweis: Da A positiv definit ist, sind alle Hauptabschnittsmatrizen

$$A(1:k, 1:k), \quad k = 1, \dots, n$$

positiv definit, haben also Determinante $\neq 0$. Also existiert eine eindeutige LR -Zerlegung $A = LR$ mit L untere Δ -Matrix mit 1-Diagonale, und R hat Diagonalelemente > 0 (gleich den Pivots). Also kann man R schreiben als $D\tilde{R}$, wobei \tilde{R} 1-Diagonale hat. Sei $D = \text{diag}(d_1, \dots, d_n)$. Die Pivots sind gerade die Determinanten der Hauptabschnittsmatrizen, also reell positiv (alle Eigenwerte sind reell, positiv), also ist $d_i > 0, \forall i = 1, \dots, n$. Wir haben $A = LD\tilde{R} = LD^{\frac{1}{2}}D^{\frac{1}{2}}\tilde{R}$. $D^{\frac{1}{2}} = \text{diag}(d_1^{\frac{1}{2}}, \dots, d_n^{\frac{1}{2}})$.

Also folgt

$$D^{-\frac{1}{2}}L^{-1}AL^{-*}D^{-\frac{1}{2}} = D^{\frac{1}{2}}\tilde{R}L^{-*}D^{-\frac{1}{2}}.$$

Links steht eine Hermitesche Matrix, rechts eine obere Δ -Matrix mit 1-Diagonale. Also muß rechts die Einheitsmatrix stehen. Also gilt

$$D^{\frac{1}{2}}\tilde{R} = D^{\frac{1}{2}}L^* \iff \tilde{R} = L^*.$$

Setze $G := LD^{\frac{1}{2}}$. □

Damit können wir den Gauß-Algorithmus verbessern und erhalten den Cholesky-Algorithmus.

Algorithmus 2.80 (Cholesky Zerlegung GAXPY-Version)

Input: Hermitesch positiv definite Matrix $A \in \mathbb{K}^{n,n}$.

Output: Untere Δ -Matrix $G \in \mathbb{K}^{n,n}$ mit positiver Diagonale, so daß $A = GG^*$. Für $i \leq j$ überschreibt $G(i, j)$ den Wert von $A(i, j)$.

```

FOR   j = 1 : n
  IF   j > 1
      A(j : n, j) = A(j : n, j) - A(j : n, 1 : j - 1) * A(j, 1 : j - 1)*;
  END
  A(j : n, j) = A(j : n, j) / sqrt(A(j, j));
END

```

Kosten: $n^3/3$ flops.

Fehleranalyse: Analog wie bei Gauß-Zerlegung.

Wie im Fall der LR -Zerlegung kann man auch die Cholesky-Zerlegung anders aufbauen:

$$A = \begin{bmatrix} \alpha & v^* \\ v & B \end{bmatrix} = \begin{bmatrix} \beta & 0 \\ v/\beta & I_{n-1} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & B - vv^*/\alpha \end{bmatrix} \begin{bmatrix} \beta & v^*/\beta \\ 0 & I_{n-1} \end{bmatrix}$$

$\beta = \sqrt{\alpha}$. Wenn die Cholesky Zerlegung von $B - vv^*/\alpha = G_1G_1^*$ gegeben ist, dann erhält man die Zerlegung, $A = GG^*$ mit

$$G = \begin{bmatrix} \beta & 0 \\ v/\beta & G_1 \end{bmatrix}.$$

Algorithmus 2.81 (Cholesky-Zerlegung: Äußere Produkte Version)

Input: Hermitesch positiv definite Matrix $A \in \mathbb{K}^{n,n}$.

Output: Untere Δ -Matrix $G \in \mathbb{K}^{n,n}$ mit positiver Diagonale, so daß $A = GG^*$. Für $i \geq j$ überschreibt $G(i, j)$ jeweils $A(i, j)$.

```
FOR   k = 1 : n
      A(k, k) = sqrt(A(k, k));
      A(k + 1 : n, k) = A(k + 1 : n, k) / A(k, k);
      FOR   j = k + 1 : n
            A(j : n, j) = A(j : n, j) - A(j : n, k) * A(j, k);
      END
END
```

Kosten: $n^3/3$ flops.

Fehleranalyse Analog wie bei Gauß-Zerlegung.

Beachte: Die j -Schleife berechnet den unteren Δ -Anteil von

$$A(k + 1 : n, k + 1 : n) = A(k + 1 : n, k + 1 : n) - A(k + 1 : n, k)A(k + 1 : n, k)^*$$

Pivotisierung: Um Symmetrie zu erhalten, sollte man nur auf der Diagonalen pivotisieren. Also Permutationen PAP^* verwenden, die die Diagonalelemente vertauschen.

2.2.1 Verallgemeinerung für nicht positiv definite Systeme.

LDL^* -Zerlegung D blockdiagonal mit 1×1 oder 2×2 Blöcken. Hier muß man i.a. permutieren, d.h. man erzeugt $P^*AP = LDL^*$.

→ Übung.

Lösung des Gleichungssystems dann wie bei Gauß mit Vorwärts- und Rückwärtseinsetzen.

2.2.2 Bandsysteme

Eine weitere häufig vorkommende spezielle Struktur ist die Bandstruktur. Eine Matrix heißt $p - q$ Bandmatrix mit unterer Bandbreite p und oberer Bandbreite q falls $a_{ij} = 0$ für $i > j + p$ und $j > i + q$.

Beispiel 2.82

$$A = \begin{bmatrix} * & * & * & 0 & 0 \\ * & * & * & * & 0 \\ 0 & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & 0 & * & * \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad \begin{array}{l} \text{untere Bandbreite } 1 \\ \text{oberer Bandbreite } 2 \\ \text{1-2 Bandmatrix.} \end{array}$$

Satz 2.83 $A \in \mathbb{K}^{n,n}$ habe eine LR -Zerlegung $A = LR$. Falls A eine $p - q$ -Bandmatrix ist, so hat L untere Bandbreite p und R obere Bandbreite q .

Beweis: Mit vollständiger Induktion:

$$A = \begin{bmatrix} \alpha & w^* \\ v & B \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ v/\alpha & I_{n-1} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & B - vw^*/\alpha \end{bmatrix} \begin{bmatrix} \alpha & w^* \\ 0 & I_{n-1} \end{bmatrix}$$

$B - vw^*/\alpha$ ist $p - q$ Bandmatrix da nur die ersten q Komponenten von w und die ersten p Komponenten von v ungleich 0 sind. Sei $L_1 R_1$ die LR -Zerlegung von $B - vw^*/\alpha$. Nach I.V. und der Sparsität von v, w folgt daß

$$L = \begin{bmatrix} 1 & 0 \\ v/\alpha & L_1 \end{bmatrix} \quad \text{und} \quad R = \begin{bmatrix} \alpha & w^* \\ 0 & R_1 \end{bmatrix}$$

die gewünschten Bandbreiten haben und es gilt $A = LR$. □

! Es folgt also, daß man die Bandstruktur optimal ausnutzen kann vorausgesetzt, die LR -Zerlegung existiert!

Algorithmus 2.84 (Band Gauß-Elimination, äußere Produkte Version)

Input: $A \in \mathbb{K}^{n,n}$ $p - q$ -Bandmatrix, welches eine LR -Zerlegung besitzt.

Output: Zerlegung $A = LR$, $A(i, j)$ überschrieben durch $L(i, j)$ für $i > j$ und $R(i, j)$ sonst.

```

FOR   k = 1 : n - 1
      FOR   i = k + 1 : min(k + p, n)
            A(i, k) = A(i, k)/A(k, k);
      END
      FOR   j = k + 1 : min(k + q, n)
            FOR   i = k + 1 : min(k + p, n)
                  A(i, j) = A(i, j) - A(i, k) * A(k, j);
            END
      END
END
END

```

Kosten: für $n \gg p, q$: $\approx 2npq$ flops

Fehleranalyse: wie normaler Gauß.

Cholesky für Band analog.

→ Übung .

Spezielle Speichertechniken möglich. Dann kann man auch den Δ -Löser verbessern.

Algorithmus 2.85 (Band-Vorwärts-Auflösen)

Input: $L \in \mathbb{K}^{n,n}$ untere Δ -Matrix mit 1-Diagonale und unterer Bandbreite $p, b \in \mathbb{K}^n$.

Output: b überschrieben mit der Lösung von $Lx = b$.

```

FOR   j = 1 : n
      FOR   i = j + 1 : min(j + p, n)
            b(i) = b(i) - L(i, j) * b(j);
      END
END

```

Kosten: für $n \gg p$: $\approx 2np$ flops .

Algorithmus 2.86 (Band-Rückwärts-Auflösen)

Input: $R \in \mathbb{K}^{n,n}$ obere Δ -Matrix mit oberer Bandbreite $q, b \in \mathbb{K}^n$.

Output: b überschrieben mit der Lösung von $Rx = b$.

```
FOR  j = n : -1 : 1
    b(j) = b(j)/u(j, j);
    FOR  i = max(1, j - q) : j - 1
        b(i) = b(i) - L(i, j) * b(j);
    END
END
```

Tridiagonalmatrizen werden als 3 Vektoren gespeichert, oder bei Hermitesch positiv definiten Matrizen als 2 Vektoren.

Algorithmus 2.87 (Symmetrisch, positiv definit, tridiagonal Löser)

Input: $A \in \mathbb{R}^{n,n}$ symmetrisch, positiv definit, tridiagonal, $b \in \mathbb{R}^n$. A gespeichert als $d(1:n), e(1:n-1)$.

Output: Lösung von $Ax = b$ überschrieben auf b .

```
FOR  k = 2 : n
    t = e(k - 1);
    e(k - 1) = t/d(k - 1);
    d(k) = d(k) - t * e(k - 1);
END
FOR  k = 2 : n
    b(k) = b(k) - e(k - 1) * b(k - 1);
END
b(n) = b(n)/d(n);
FOR  k = n - 1 : -1 : 1
    b(k) = b(k)/d(k) - e(k) * b(k + 1);
END
```

Kosten: $8n$ flops.

→ Übung 2.94

2.2.3 Übungen zu Kapitel 2.2

Übung 2.88 Sei

$$A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{n,n}.$$

Zeige:

1. A besitzt eine LU-Zerlegung
2. partielle Pivotisierung führt zu keiner Änderung.
3. Es besteht die Beziehung $U = DL^T$. Was bedeutet das für die LU-Zerlegung?

Übung 2.89 Sei $A \in \mathbb{R}^{n,n}$, symmetrisch positiv definit, $A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$ und nehme an, dass $A_{11} \in \mathbb{R}^{k,k}$ ist. Das Schur-Komplement ist gegeben durch $S = A_{22} - A_{21}A_{11}^{-1}A_{12}$. Zeige:

1. A lässt sich schreiben als

$$A = L \begin{pmatrix} A_{11} & O \\ O & S \end{pmatrix} L^\top, \text{ wobei } L = \begin{pmatrix} I & O \\ A_{21}A_{11}^{-1} & I \end{pmatrix}.$$

2. Ist $M \in \mathbb{R}^{n,m}$ mit $\text{rang } M = m$, dann ist $M^\top A M$ symmetrisch positiv definit.

3. S ist symmetrisch positiv definit

Übung 2.90 Sei $A \in \mathbb{R}^{n,n}$ symmetrisch positiv definit und $A = LL^\top$ die Cholesky-Zerlegung von A . Zeige:

1. Für $i = 1, \dots, n$ gilt

$$l_{ii}^2 \geq \min_{x \neq 0} \frac{x^\top A x}{x^\top x} = \frac{1}{\|L^{-\top}\|_2^2}.$$

2. Für $i = 1, \dots, n$ gilt

$$\|L^\top\|_2^2 = \max_{x \neq 0} \frac{x^\top A x}{x^\top x} \geq l_{ii}^2.$$

3. Für die Konditionszahl in der 2-Norm gilt

$$\text{cond}_2(L^\top) \geq \max_{1 \leq i, k \leq n} \frac{|l_{ii}|}{|l_{kk}|}.$$

Übung 2.91 Schreibe einen **MATLAB**-Algorithmus zur Berechnung der Cholesky-Zerlegung. Vergleiche den Algorithmus bezüglich Rechengenauigkeit und Rechenzeit mit der Funktion 'chol' und dem '^' aus **MATLAB** für $n = 10, 20, 30, 40, \dots, 100$ und zufällig erzeugten positiv definiten Matrizen $A = MDM^\top$, $M = \text{rand}(n)$, $D = \text{diag}(1, 2, 4, 9, \dots, n^2)$.

Übung 2.92 Sei $A \in \mathbb{R}^{n,n}$, $A = (a_{ij})_{i,j=1,\dots,n}$, $A = A^\top$. Zeige folgende Äquivalenzen:

1. A ist positiv definit, d.h. $x^\top A x > 0$, $\forall x \neq 0$.

2. Alle Eigenwerte von A sind positiv.

3. Alle Hauptabschnittsmatrizen $(a_{ij})_{i,j=1,\dots,k}$, $k \leq n$ sind positiv definit.

4. Alle Hauptabschnittsdeterminanten $\det \left((a_{ij})_{i,j=1,\dots,k} \right)$, $k \leq n$ sind positiv.

Übung 2.93 Sei $A \in \mathbb{R}^{n,n}$ und sei $F \in \mathbb{R}^{n,r}$ mit $r \leq n$. Zeige:

1. Ist A positiv definit (d.h. $x^\top A x > 0$ für alle $x \neq 0$), dann ist $F^\top A F$ positiv definit genau dann, wenn F vollen Rang hat.

2. Welchen Rang haben Matrizen der Form $\begin{pmatrix} X \\ I_r \\ Y \end{pmatrix}$, $\begin{pmatrix} I_r \\ Y \end{pmatrix}$?

3. Ist A positiv definit, dann ist auch A^{-1} positiv definit.

Übung 2.94 1. Schreibe einen **MATLAB**-Algorithmus zur Berechnung der Cholesky-Zerlegung. Vergleiche den Algorithmus bezüglich Rechengenauigkeit und Rechenzeit mit der Funktion 'chol' und dem '^' aus **MATLAB** für $n = 10, 20, 30, 40, \dots, 100$ und zufällig erzeugten positiv definiten Matrizen $A = MDM^\top$, $M = \text{rand}(n)$, $D = \text{diag}(1, 2, 4, 9, \dots, n^2)$.

2. Schreibe einen **MATLAB**-Algorithmus zur Berechnung der Cholesky-Zerlegung für tridiagonale Matrizen. Vergleiche den Algorithmus bezüglich Rechengenauigkeit und Rechenzeit mit der Funktion 'chol' und dem '\ ' aus **MATLAB** für $n = 10, 20, 30, 40, \dots, 100$ anhand der Matrix

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

2.3 Householder und Gram-Schmidt-Orthogonalisierung

Wir haben bei der Fehleranalyse von Gauß gesehen, daß die Größe $|\tilde{L}|, |\tilde{R}|$ der berechneten Matrizen \tilde{L}, \tilde{R} in die Abschätzung des Fehlers eingeht (2.38). Nun kann natürlich $|\tilde{L}|$ sehr groß sein, wenn wir ohne Pivotisierung arbeiten. Außerdem haben wir gesehen, daß die Lösung von Gleichungssystemen mit Gauß auch bei Pivotisierung nicht numerisch rückwärts stabil ist, wegen der Wachstumsfaktoren. In diesem Kapitel betrachten wir nun Methoden, die auf Zerlegungen mit orthogonalen (unitären) Matrizen beruhen und bei denen wir numerische Stabilität zeigen können.

2.3.1 Householder Matrizen(Spiegelungen)

Sei $v \in \mathbb{K}^n \setminus \{0\}$. Eine $n \times n$ Matrix der Form

$$P = I - \frac{2vv^*}{v^*v} \quad (2.95)$$

heißt Householder Matrix.

Eine Multiplikation mit P spiegelt einen Vektor x an der Hyperebene $\text{span}(v)^\perp$. Householder Matrizen sind Hermitesch und unitär (im reellen symmetrisch und orthogonal). Sie sind Rang 1 Modifikationen der Identität. Sie werden verwendet, um bestimmte Komponenten eines Vektors zu eliminieren. (Analog wie bei Gauß-Transformationen).

Angenommen $x \in \mathbb{K}^n \setminus \{0\}$ und wir wollen v bestimmen, so daß $Px = ce_1$ (P wie in (2.95))

$$Px = \left(I - \frac{2vv^*}{v^*v} \right) x = x - \left(\frac{2v^*x}{v^*v} \right) v$$

$Px \in \text{span}\{e_1\} \implies v \in \text{span}\{x, e_1\}$. Setze $v = x + \alpha e_1$, dann gilt

$$v^*x = x^*x + \alpha x_1, v^*v = x^*x + 2 \text{Real}(\alpha x_1) + |\alpha|^2.$$

Also

$$Px = \left(1 - 2 \frac{x^*x + \alpha x_1}{x^*x + 2 \text{Real}(\alpha x_1) + |\alpha|^2} \right) x - 2\alpha \frac{v^*x}{v^*v} e_1$$

Damit der Koeffizient 0 ist, brauchen wir also nur $\alpha = e^{i\varphi} \|x\|_2 = e^{i\varphi} (x^*x)^{\frac{1}{2}}$. Dann gilt

$$v = x + e^{i\varphi} \|x\|_2 e_1 \implies Px = \left(I - 2 \frac{vv^*}{v^*v} \right) x = -e^{i\varphi} \|x\|_2 e_1.$$

2.3.2 Berechnung der Householder–Matrix

Wichtige Details:

1. Wahl des Vorzeichens. Falls $x \approx \alpha e_1$, dann hat $v = x - \text{sign}(x_1)\|x\|_2 e_1$ sehr kleine Norm \implies Auslöschung großer relativer Fehler in $\beta = 2/v^*v$. Wähle daher immer $v = x + \text{sign}(x_1)\|x\|_2 e_1$, dann gilt $\|v\|_2 \geq \|x\|_2$ und P ist fast perfekt unitär.
2. Normalisierung von v , auf $v(1) = 1$. Dies vereinfacht eine ganze Menge von Algorithmen, die auf Householder–Vektoren beruhen.

Algorithmus 2.96 (Erzeugung des Householder Vektors)

Input: $x \in \mathbb{K}^n$.

Output: $v \in \mathbb{K}^n$ mit $v(1) = 1$, so daß $(I - \frac{2vv^*}{v^*v})x = \alpha e_1$.

```

FUNCTION   v = HOUSE(x)
           n = LENGTH(x);
            $\mu = \|x\|_2$ ;
           v = x;
           IF    $\mu \neq 0$ 
                $\beta = x(1) + \text{sign}(x(1)) * \mu$ ;
                $v(2 : n) = v(2 : n) / \beta$ ;
           END
           v(1) = 1;
END HOUSE

```

Kosten: $\approx 3n$ flops.

Fehler: Berechnetes \tilde{P} erfüllt $\|\tilde{P} - P\|_2 = \mathcal{O}(\text{eps})$.

2.3.3 Anwendung von Householder–Matrizen

Bei der Multiplikation mit Householder–Matrizen kann an einigen Stellen die Struktur ausgenutzt werden.

$$PA = \left(I - \frac{2vv^*}{v^*v} \right) A = A + vw^*$$

mit $w = \beta A^*v$ und $\beta = \frac{-2}{v^*v}$.

Matrix–Vektor–Multiplikation und äußere–Produkt Aufdatierung.

Algorithmus 2.97 (Vormultiplikation mit Householder–Matrix)

Input: $A \in \mathbb{K}^{m,n}$, $v \in \mathbb{K}^m$, $v(1) = 1$.

Output: A überschrieben mit $PA = \left(I - \frac{2vv^*}{v^*v} \right) A$.

```

FUNCTION   A = ROWHOUSE(A, v)
            $\beta = -2/v^* * v$ ;
            $w = \beta * A^* * v$ ;
            $A = A + v * w^*$ ;
END ROWHOUSE

```

Kosten: $4mn$ flops.

Fehler: $gl(\tilde{P}A) = P(A + E)$, $\|E\|_2 = \mathcal{O}(\text{eps} \|A\|_2)$

Algorithmus 2.98 (Nachmultiplikation mit Householder-Matrix)

Input: $A \in \mathbb{K}^{m,n}, v \in \mathbb{K}^n, v(1) = 1.$

Output: A überschrieben mit $AP = A \left(I - \frac{2vv^*}{v^*v} \right).$

```
FUNCTION      A = COLHOUSE(A, v)
               $\beta = -2/v^* * v;$ 
               $w = \beta * A * v;$ 
               $A = A + w * v^*;$ 
END COLHOUSE
```

Kosten: $4mn$ flops.

Fehler: $gl(A\tilde{P}) = (A + E)P, \|E\|_2 = \mathcal{O}(\text{eps} \|A\|_2).$

Diese 3 Algorithmen werden verwendet um bestimmte Teile einer Matrix zu eliminieren.

Beispiel 2.99 Überschreibe $A \in \mathbb{R}^{m,n} (m \geq n)$ mit $B = Q^T A$ wobei Q orthogonal, so daß $B(j+1 : m, j) = 0$ für ein j mit $1 \leq j \leq n$. Sei weiter angenommen, daß $A(j : m, 1 : j-1) = 0$ und wir wollen den nichttrivialen Teil von v in $A(j+1 : m, j)$ speichern.

$$\begin{aligned} v(j : m) &= \text{HOUSE}(A(j : m, j)); \\ A(j : m, j : n) &= \text{ROWHOUSE}(A(j : m, j : n), v(j : m)); \\ A(j+1 : m, j) &= v(j+1 : m); \end{aligned}$$

Mathematisch gesehen haben wir damit mit der $m \times m$ Householder-Matrix

$$P = \begin{bmatrix} I_{j-1} & 0 \\ 0 & \tilde{P} \end{bmatrix} = I - \frac{zvv^T}{v^T v}$$

multipliziert. Fehleranalyse (Wilkinson) sagt, daß Berechnung und Anwendung von Householder Matrizen numerisch rückwärts stabil sind.

Faktorisierete Speicherung. Analog zur Gauß-Transformationen speichert man meistens die Produkte von P_i 's als v_i 's ab, wie im Beispiel.

2.3.4 Die QR-Zerlegung.

Die QR-Zerlegung einer Matrix $A \in \mathbb{K}^{m,n} (m \geq n)$ ist gegeben durch

$$A = QR$$

mit $Q \in \mathbb{K}^{m,m}$ unitär, $R \in \mathbb{K}^{m,n}$ obere Δ -Matrix $R = \begin{bmatrix} \square & \\ & \square \\ & & \square \\ & & & 0 \end{bmatrix}$.

Falls A vollen Rang hat, dann bilden die ersten n Spalten von Q eine Orthonormalbasis für Bild (A). Mit Hilfe der QR-Zerlegung kann man also Orthonormalbasen für Mengen von Vektoren bestimmen.

Lösung von Gleichungssystemen.

$$\begin{aligned} Ax &= b & A \in \mathbb{K}^{n,n} \\ QRx &= b \iff c = Q^H b, Rx = c \end{aligned}$$

Anwendung von Q^H auf b nur Matrix–Vektor Multiplikation oder in faktorisierter Householder Form, mehrere innere Produkte. $Rx = c$ durch rückwärts Auflösen.

Lösung von Ausgleichproblemen als nächstes:

Wie erhält man eine QR –Zerlegung?

Nacheinander Anwendung von Householder–Matrizen.

$$\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}, P_1 A = \begin{bmatrix} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \\ 0 & * & * & * & * \end{bmatrix}, P_2 P_1 A = \begin{bmatrix} * & * & * & * & * \\ 0 & * & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \\ 0 & 0 & * & * & * \end{bmatrix}, \text{ usw.}$$

Algorithmus 2.100 (Householder QR –Zerlegung)

Input: $A \in \mathbb{K}^{m,n}$ mit $m \geq n$.

Output: Householder–Matrizen P_1, \dots, P_n , so daß für $Q = P_1 \cdots P_n, Q^H A = R$ obere Δ –Matrix. Der obere Δ –Teil von A wird durch R überschrieben und die Komponenten $j+1 : m$ des j –ten Householdervektors werden auf $A(j+1 : m, j)$, $j < m$ überschrieben.

FOR $j = 1 : n$

$v(j : m) = \text{HOUSE}(A(j : m, j));$

$A(j : m, j : n) = \text{ROWHOUSE}(A(j : m, j : n), v(j : m));$

IF $j < m$

$A(j+1 : m, j) = v(j+1 : m);$

END

END

Kosten: $2n^2(m - \frac{n}{3})$ flops.

Fehleranalyse: $\tilde{Q}^*(A + E) = \tilde{R}$, wobei berechnetes \tilde{Q} exakt unitär und $Q^* \tilde{Q} = I + F$, $\|F\|_2 \approx \text{eps}$, $\|E\|_2 \approx \text{eps} \|A\|_2$.

A sieht am Ende des Algorithmus wie folgt aus

$$\begin{bmatrix} r_{11} & r_{12} & \cdots & \cdots & r_{1,n} \\ v_2^{(1)} & r_{22} & & & \vdots \\ v_3^{(1)} & v_3^{(2)} & \ddots & & \vdots \\ \vdots & \vdots & \ddots & r_{n-1,n-1} & \vdots \\ \vdots & \vdots & & v_n^{(n-1)} & r_{n,n} \\ \vdots & \vdots & & \vdots & v_{n+1}^{(n)} \\ \vdots & \vdots & & \vdots & \vdots \\ v_m^{(1)} & v_m^{(2)} & \cdots & v_m^{(n-1)} & v_m^{(n)} \end{bmatrix}$$

Die j –te Komponente von $v^{(j)}$ ist jeweils 1. Falls Q benötigt wird, kostet das

$$4(m^2 n - mn^2 + \frac{n^3}{3}) \text{ flops.}$$

2.3.5 Eigenschaften der QR-Zerlegung

Der obige Algorithmus beweist die Existenz.

Satz 2.101 Ist $A = QR$ eine QR-Zerlegung einer Matrix $A \in \mathbb{K}^{m,n}$ von vollem Rang und $Q = [q_1, \dots, q_n]$, $A = [a_1, \dots, a_n]$, dann ist

$$\text{span} \{a_1, \dots, a_k\} = \text{span} \{q_1, \dots, q_k\}, \quad k = 1 : n.$$

Für $Q_1 = Q(1 : m, 1 : n)$, $Q_2 = Q(1 : m, n + 1 : m)$ gilt:

$$\begin{aligned} \text{Bild}(A) &= \text{Bild}(Q_1), \\ \text{Bild}(A)^\perp &= \text{Bild}(Q_2) \end{aligned}$$

und $A = Q_1 R_1$ mit $R_1 = R(1 : n, 1 : n)$.

Beweis:

$$a_k = \sum_{i=1}^k r_{ik} q_i \in \text{span} \{q_1, \dots, q_k\}$$

$$\implies \text{span} \{a_1, \dots, a_k\} \subseteq \text{span} \{q_1, \dots, q_k\}$$

$$\text{rank}(A) = n \implies \dim(\text{span} \{a_1, \dots, a_k\}) = k, \forall k \leq n$$

$$\implies \text{span} \{a_1, \dots, a_k\} = \text{span} \{q_1, \dots, q_k\}.$$

Rest trivial. □

Satz 2.102 $A \in \mathbb{K}^{m,n}$ habe vollen Rang. Die „dünne“ QR-Zerlegung

$$A = Q_1 R_1 \text{ mit } Q_1 \in \mathbb{K}^{m,n}, R \in \mathbb{K}^{n,n}$$

obere Δ -Matrix $r_{ii} \in \mathbb{R}, r_{ii} > 0$, ist eindeutig. $G = R_1^*$ ist der Cholesky Faktor von A^*A .

Beweis: $A^*A = (Q_1 R_1)^*(Q_1 R_1) = R_1^* R_1$.

R_1 ist eindeutig und da A vollen Rang hat auch $Q_1 = AR_1^{-1}$. □

2.3.6 Gram-Schmidt Orthogonalisierung

Klassischer Gram-Schmidt (CGS): Die $Q_1 R_1$ Zerlegung wie in Satz 2.102 kann man natürlich auch über Gram-Schmidt erhalten. Sei $\text{rank}(A) = n$. Der klassische Gram-Schmidt ist dann gegeben durch

$$\begin{aligned} r_{ik} &= q_i^* a_k, \quad i = 1 : k - 1 \\ q_k &= (a_k - \sum_{i=1}^{k-1} r_{ik} q_i) / r_{kk}. \end{aligned}$$

Numerisch gesehen ist dieser Algorithmus sehr schlecht, da die Orthogonalität der q_k durch Rundungsfehler zerstört wird. Eine leichte Umsortierung ergibt den sogenannten modifizierten

Gram–Schmidt (MGS), der numerisch wesentlich besser ist.
 Definiere $A^{(k)} \in \mathbb{K}^{m, n-k+1}$ durch

$$A - \sum_{i=1}^{k-1} q_i r_i^* = \sum_{i=k}^n q_i r_i^* = [0 \ A^{(k)}].$$

Mit $A^{(k)} = \left[\underbrace{z}_1 \ \underbrace{B}_{n-k} \right]$ gilt: $r_{kk} = \|z\|_2$, $q_k = z/r_{kk}$ und $[r_{k,k+1}, \dots, r_{k,n}] = q_k^* B$.

Damit können wir

$$A^{(k+1)} = B - q_k [r_{k,k+1}, \dots, r_{k,n}]$$

als äußeres Produkt bestimmen und weitermachen.

Algorithmus 2.103 (Modifizierter Gram–Schmidt)

Input: $A \in \mathbb{K}^{m,n}$ Rang $(A) = n$.

Output: Zerlegung $A = Q_1 R_1$, $Q_1 \in \mathbb{K}^{m,n}$ mit orthonormalen Spalten, $R_1 \in \mathbb{K}^{n,n}$ obere Δ -Matrix.

FOR $k = 1 : n$

$$R(k, k) = \|A(1 : m, k)\|_2;$$

$$Q(1 : m, k) = A(1 : m, k)/R(k, k);$$

FOR $j = k + 1 : n$

$$R(k, j) = Q(1 : m, k)^* * A(1 : m, j);$$

$$A(1 : m, j) = A(1 : m, j) - Q(1 : m, k) * R(k, j);$$

END

END

Kosten: $2mn^2$ flops.

Fehler: $\tilde{Q}_1^* \tilde{Q}_1 = I + E$, $\|E\|_2 \approx \text{eps } \kappa_2(A)$

Zum Vergleich bei Householder $\|E\|_2 = \text{eps}$.

Kostenvergleich: Zur Berechnung einer Orthonormalbasis für $\text{Bild}(A)$ ist der modifizierte Gram–Schmidt schneller als QR.

2.3.7 Übungen zu Kapitel 2.3

Übung 2.104 Zeige:

1. Für jedes $V \in \mathbb{R}^{n,m}$, Rang $V = m$ ist $S = I - 2V(V^T V)^{-1} V^T$ symmetrisch und orthogonal und es ist $S^2 = I$.
2. S ist eine Spiegelung am orthogonalen Komplement von $\text{Bild } V$.
3. S hat nur Eigenwerte -1 m -fach, 1 $n - m$ -fach.
4. $P = I - V(V^T V)^{-1} V^T$ ist symmetrisch und erfüllt $P^2 = P$.
5. P ist eine Projektion auf das orthogonalen Komplement von $\text{Bild } V$.
6. P hat nur Eigenwerte 0 m -fach, 1 $n - m$ -fach.
7. Jede orthogonale Matrix obere Dreiecksmatrix ist bereits eine Diagonalmatrix.
8. Jede orthogonale Matrix läßt sich als Produkt von Elementarspiegelungen S (d.h. $m = 1$ in (a)) schreiben.

Übung 2.105 Sei $G = \begin{pmatrix} c & -s \\ s & c \end{pmatrix}$ orthogonal und nehme an, dass wir für vorgegebenes $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$, c, s so bestimmen können, dass $Gv = \begin{pmatrix} * \\ 0 \end{pmatrix}$ ist. Zeige:

- Ist A eine obere Hessenbergmatrix (d.h. $a_{i,j} = 0$ für alle $i > j + 1$), dann besitzt A eine QR-Zerlegung, die in $3n^2 + \mathcal{O}(n)$ flops berechnet werden kann (ohne explizite Berechnung von Q). Welches Belegungsmuster hat Q ?
- Ist A eine beliebige $n \times n$ -Matrix, dann braucht die QR-Zerlegung mit Hilfe von G , $2n^3 + \mathcal{O}(n^2)$ flops (ohne explizite Berechnung von Q). Ist das besser oder schlechter als die QR-Zerlegung mit Hilfe von Householder-Transformationen?

Übung 2.106 Zeige:

- Ist eine orthogonale Matrix gleichzeitig eine obere Dreiecksmatrix, dann ist sie bereits diagonal. Wie sehen die Diagonaleinträge aus? Gilt ein analoges Ergebnis auch für unitäre Matrizen?
- Jede orthogonale Matrix lässt sich als Produkt von Elementarspiegelungen (Householdertransformationen) schreiben. Gilt das auch, wenn man Elementarspiegelungen durch Rotationen ersetzt?

Übung 2.107 Betrachte für $c > 0, s = \sqrt{1 - c^2}$ die Matrix

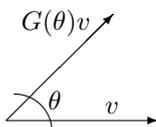
$$T_n = \text{diag}(1, s, \dots, s^{n-1}) \begin{pmatrix} 1 & -c & \cdots & -c \\ 0 & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & -c \\ 0 & \cdots & 0 & 1 \end{pmatrix} \in \mathbb{R}^{n,n}.$$

Zeige:

- $\max_{1 \leq k \leq n} \|T_n(:, k)\|_2 = \|T_n(:, 1)\|_2$.
- $\sigma_n(T_n) \equiv \min_{\|x\|_2=1} \|T_n x\|_2 \leq \frac{s^{n-1}}{c(c+1)^{n-2}}$.
- Was bedeutet das für den QR-Algorithmus mit Spaltenpivotisierung (d.h. in der Restmatrix wird die Spalte mit maximaler 2-Norm mit der führenden Spalte vertauscht) und für die Rangbestimmung von T_n ? (Betrachte z.B. $n = 100, c = 0.2$)

Übung 2.108 Sei $G(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \equiv \begin{pmatrix} c & -s \\ s & c \end{pmatrix}$. Zeige:

- die Abbildung G beschreibt eine Drehung um den Winkel θ gegen den Uhrzeigersinn in der Ebene.



- Bestimme für vorgegebenes $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$, c, s so, dass $Gv = \begin{pmatrix} * \\ 0 \end{pmatrix}$ ist. Wie müssen c, s bestimmt werden, damit ein Höchstmaß an numerischer Stabilität gesichert ist? Begründung!

Übung 2.109 Schreibe einen MATLAB-Algorithmus zur Berechnung der QR-Zerlegung mittels Householder-Transformationen.

Vergleiche den Algorithmus bezüglich Rechengenauigkeit und Rechenzeit mit der Funktion 'qr' aus MATLAB für $n = 10, 20, 30, 40, \dots, 100$ und zufällig erzeugten Matrizen

Übung 2.110 Seien $A \in \mathbb{R}^{n,m}$, $v, w \in \mathbb{R}^n$. Vergleiche den Rechenaufwand folgender Algorithmen zur Berechnung von $B = (I - vw^T)A$.

1. $P := I - vw^T, B := PA$
2. $z := w^T A, B := A - vz$

Unterscheidet sich der Rechenaufwand wesentlich?

Übung 2.111 Sei $Q = [q_1, \dots, q_n] \in \mathbb{R}^{m,n}$ mit $Q^T Q = I$.

1. Was bedeutet das für q_1, \dots, q_n ?
2. Zeige:

$$I - QQ^T = \prod_{i=1}^n (I - q_i q_i^T),$$

dabei ist die Reihenfolge, in der das Produkt gebildet wird, beliebig.

3. Sei $a \in \mathbb{R}^n$. Zeige $(I - QQ^T)a$ steht senkrecht auf q_1, \dots, q_n
4. Berechne $b = (I - QQ^T)a$ auf zwei Weisen

1.


```

      for j = 1 : n
          R_j := q_j^T a
      end
      b := a
      for j = 1 : n
          b := b - q_j * R_j
      end
      
```
2.


```

      b := a
      for j = 1 : n
          R_j := q_j^T b
          b := b - q_j * R_j
      end
      
```

Nach einer etwaigen Normierung von b am Ende, welcher der beiden Algorithmen entspricht dem klassischen und welcher dem modifizierten Gram-Schmidt-Verfahren. Wie kann man das hinsichtlich Teil (b) interpretieren?

Übung 2.112 Seien A und $B = QA$ gegeben, wobei Q orthogonal. Nehme an, dass $\nu_k = \|(a_{i,k})_{i=1,\dots,m}\|_2^2$, für $k = 1, \dots, n$ gegeben ist. Wie lassen sich $\nu'_k = \|(b_{i,k})_{i=2,\dots,m}\|_2^2$, für $k = 2, \dots, n$ möglichst schnell aus ν_2, \dots, ν_n berechnen? Welche Auswirkungen hat das auf Anzahl flops bei der QRP-Zerlegung?

Übung 2.113 Schreibe einen **MATLAB**-Algorithmus zur Berechnung

1. der QRP-Zerlegung (Spaltenpivotisierung) mittels Householder-Transformationen (hierzu wird in jedem Schritt vor der Elimination die Spalte mit der größten 2-Norm in die führende Position gebracht. D.h., ist $A(j : m, j : n)$ die Restmatrix im j -ten Schritt, so berechne μ mit $\|A(:, \mu)\|_2 = \max_{j \leq l \leq n} \|A(:, l)\|_2$ und vertausche $A(:, j)$ mit $A(:, \mu)$),
2. der QR-Zerlegung mittels klassischem Gram-Schmidt Verfahren.
3. der QR-Zerlegung mittels modifiziertem Gram-Schmidt Verfahren.

Untersuche die Algorithmen für das lineare Ausgleichsproblem $\|Ax - b\|_2 \stackrel{!}{=} \min$ und vergleiche sie bezüglich Rechengenauigkeit und Rechenzeit mit dem Cholesky-Verfahren für die Normalgleichungen $A^T Ax = A^T b$. Zur Generierung der Beispielaufgaben kann auf ein entsprechendes **MATLAB**-Scriptfile im WWW zurückgegriffen werden.

Übung 2.114 Sei

$$A = \begin{pmatrix} * & & & & & & * \\ * & & * & & & & \\ * & & & & & & * \\ * & & * & & & & \\ * & & & & * & & \\ * & * & & & & & \\ * & & & & * & & \\ * & * & * & * & * & * & * \end{pmatrix}$$

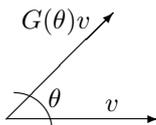
Bestimme Permutationen U, V , so dass bei einer QR-Zerlegung von UAV möglichst wenig Einträge in R entstehen.

Übung 2.115 Zeige:

1. Für jedes $V \in \mathbb{C}^{n,m}$, $\text{Rang } V = m$ ist $S = I - 2V(V^H V)^{-1}V^H$ Hermitesch und unitär und es ist $S^2 = I$.
2. S ist eine Spiegelung am orthogonalen Komplement von Bild V .
3. S hat nur die Eigenwerte -1 m -fach, 1 $n - m$ -fach.
4. Sei $x \in \mathbb{C}^n \setminus \{0\}$ und schreibe die erste Komponente x_1 von x in Polarkoordinaten als $x_1 = |x_1| e^{i\theta}$. Setze $v = x + e^{i\theta} \|x\|_2 e_1$. Dann ist $S = I - \frac{2}{v^H v} v v^H$ eine Spiegelung mit $Sx = -e^{i\theta} \|x\|_2 e_1$.

Übung 2.116 Sei $G(\theta) = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \equiv \begin{pmatrix} c & -s \\ s & c \end{pmatrix}$. Zeige:

1. die Abbildung G beschreibt eine Drehung um den Winkel θ gegen den Uhrzeigersinn in der Ebene.



2. Für vorgegebenes $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$ bewirkt die Wahl

$$c = \frac{v_1}{\sqrt{v_1^2 + v_2^2}}, \quad s = \frac{-v_2}{\sqrt{v_1^2 + v_2^2}},$$

dass $Gv = \begin{pmatrix} * \\ 0 \end{pmatrix}$ ist (Givens-Rotation). Wie sollten c, s berechnet werden, um ein Höchstmaß an numerischer Stabilität zu sichern? Begründung!

3. Sei $v \in \mathbb{R}^n$. Dann lässt sich v durch Anwendung von $n - 1$ eingebetteten Givensrotationen aus Teil 2 auf ein Vielfaches von e_1 transformieren.

Übung 2.117 Schreibe einen **MATLAB**-Algorithmus zur Berechnung der QR-Zerlegung mittels Givens-Rotationen aus Aufgabe 2.

Vergleiche den Algorithmus bezüglich Rechengenauigkeit und Rechenzeit mit der Funktion 'qr' aus **MATLAB** für $n = 10, 20, 30, 40, \dots, 100$ und zufällig erzeugten Matrizen

Übung 2.118 Sei $G = \begin{pmatrix} c & -s \\ s & c \end{pmatrix}$ orthogonal und nehme an, dass wir für vorgegebenes $v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$, c, s so bestimmen können, dass $Gv = \begin{pmatrix} * \\ 0 \end{pmatrix}$ ist. Zeige:
 Ist A eine obere Hessenbergmatrix (d.h. $a_{i,j} = 0$ für alle $i > j + 1$), dann besitzt A eine QR-Zerlegung, die in $3n^2 + \mathcal{O}(n)$ flops berechnet werden kann (ohne explizite Berechnung von Q).
 Welches Belegungsmuster hat Q ?

2.4 Ausgleichsprobleme

Bei vielen wissenschaftlichen Versuchen geht es darum, aus Messungen die Werte von Konstanten x_1, \dots, x_n zu bestimmen. Oft kann man x_i nicht direkt messen sondern nur eine leichter zugängliche Größe y , die als Funktion von x und anderen Parametern z abhängt, $y = f(z, x_1, \dots, x_n)$.

Beispiel 2.119 Bestimmung von g .

Fallversuche $h = g \frac{t^2}{2}$, h – Fallhöhe, t – Fallzeit, g – Gravitation.

Um x_i zu bestimmen, führe $m > n$ verschiedene Experimente durch und erhalte

$$y_k = f(z_k, x_1, \dots, x_n), \quad k = 1, \dots, m.$$

Dies ergibt überbestimmtes Gleichungssystem. Dies ist i.a. nicht lösbar. Daher nur beste Approximation möglich. Zum Beispiel minimiere:

$$\sum_{k=1}^m (y_k - f(z_k, x_1, \dots, x_n))^2, \quad \text{kleinste Quadrate (least square), } (\|\bullet\|_2) \quad (2.120)$$

oder

$$\max_{1 \leq k \leq m} |y_k - f(z_k, x_1, \dots, x_n)|, \quad (\|\bullet\|_\infty).$$

für (2.120) ergibt sich die notwendige Bedingung für Minimierung durch

$$\frac{\partial}{\partial x_i} \left(\sum_{k=1}^m (y_k - f_k(x_1, \dots, x_n))^2 \right) = 0, \quad i = 1 : n, \quad (2.121)$$

wobei $f_k(x_1, \dots, x_n) := f(z_k, x_1, \dots, x_n)$.

Dies sind die sogenannten Normalgleichungen. Ein wichtiger Spezialfall, das lineare Ausgleichsproblem, liegt vor falls $f_k(x_1, \dots, x_n)$ linear in x_1, \dots, x_n , d.h.

$$\begin{bmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_m(x_1, \dots, x_n) \end{bmatrix} = Ax, \quad A \in \mathbb{K}^{m,n}. \quad (2.122)$$

Dann sind die Normalgleichungen

$$\text{grad}_x ((y - Ax)^*(y - Ax)) = 2A^*Ax - 2A^*y = 0. \quad (2.123)$$

Also man hat das lineare Gleichungssystem

$$A^*Ax = A^*y \quad (2.124)$$

zur Lösung von

$$\min_{x \in \mathbb{K}^n} \|Ax - y\|_2. \quad (2.125)$$

Satz 2.126 Das lineare Ausgleichsproblem (2.125) hat mindestens eine Lösung x_0 . Ist x_1 eine weitere Lösung, so gilt $Ax_0 = Ax_1$. Das Residuum $r = y - Ax_0$ ist eindeutig bestimmt und erfüllt $A^*r = 0$. Jede Lösung x_0 erfüllt auch (2.124) und umgekehrt.

Beweis: $\mathbb{K}^m = \underbrace{\text{Bild}(A)}_L \oplus \underbrace{\text{Bild}(A)^\perp}_{L^\perp}$.

Daher gibt es eine eindeutige Zerlegung

$$y = s + r, \text{ mit } s \in L, r \in L^\perp \tag{2.127}$$

und es gibt ein x_0 mit $A^*x_0 = s$. Da $A^*r = 0$, so folgt $A^*y = A^*s = A^*Ax_0$. Also löst x_0 (2.124) die Normalgleichungen.

Umgekehrt sei x_1 Lösung von (2.124). Dann hat y die Zerlegung $y = s + r$ mit $s = Ax_1, r = y - Ax_1, s \in L, r \in L^\perp$. Wegen der Eindeutigkeit der Zerlegung folgt $Ax_1 = Ax_0$.

Weiter erfüllt jede Lösung x_0 von (2.124) das Ausgleichsproblem (2.125). Denn für beliebiges x setze $z = Ax - Ax_0, r = y - Ax_0$, dann gilt wegen $r^*z = 0$:

$$\|y - Ax\|_2^2 = \|r - z\|_2^2 = \|r\|_2^2 + \|z\|_2^2 \geq \|r\|_2^2 = \|y - Ax_0\|_2^2.$$

□

Falls A vollen Rang hat, so ist A^*A positiv definit. Nun könnte man einfach folgenden Algorithmus verwenden:

Algorithmus 2.128 (Normalgleichung)

Input: $A \in \mathbb{K}^{m,n}$ mit $\text{rank}(A) = n$, und $b \in \mathbb{K}^m$.

Output: Lösung x des Minimierungsproblems (2.125).

Berechne das untere Dreieck von $C = A^*A$.

$d = A^*b$.

Berechne Cholesky Zerlegung von $C = GG^*$.

Löse $Gy = d$ und $G^*x = y$.

Kosten: $(m + n/3)n^2$ flops.

Fehleranalyse: Angenommen, keine Fehler bei $C = A^*A$ und $d = A^*b$. Dann gilt wegen der Analyse der Cholesky Zerlegung:

$$(A^*A + E)\tilde{x} = A^*b,$$

mit $\|E\|_2 \approx \text{eps} \|A^*\|_2 \|A\|_2 \approx \text{eps} \|A^*A\|_2$, d.h. wir erwarten

$$\frac{\|x - \tilde{x}\|_2}{\|x\|_2} \approx \text{eps} \kappa_2(A^*A) = \text{eps} \kappa_2(A)^2,$$

! d.h. das Quadrat der Kondition geht ein. **Schlecht!**

Beispiel 2.129

$$A = \begin{bmatrix} 1 & 1 \\ 10^{-3} & 0 \\ 0 & 10^{-3} \end{bmatrix}, b = \begin{bmatrix} 2 \\ 10^{-3} \\ 10^{-3} \end{bmatrix}, x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \kappa_2(A) = 1.4 \cdot 10^3.$$

Mit 6-stelliger Arithmetik ist $A^*A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ singular. Mit 7-stelliger Arithmetik folgt

$$\tilde{x} = \begin{bmatrix} 2.00001 \\ 0 \end{bmatrix} \implies \frac{\|\tilde{x} - x\|_2}{\|x\|_2} \approx \underbrace{\text{eps}}_{10^{-6}} \underbrace{\kappa_2(A)^2}_{10^6}.$$

Normalengleichung kann also zu sehr schlechten Ergebnissen führen. Verwende daher QR -Zerlegung. Bilde

$$Q^*A = R = \begin{bmatrix} R_1 & \\ & 0 \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix}, \quad R_1 \text{ obere } \Delta\text{-Matrix}$$

$$Q^*b = \begin{bmatrix} c \\ d \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix}$$

Dann gilt

$$\|Ax - b\|_2^2 = \|Q^*Ax - Q^*b\|_2^2 = \|R_1x - c\|_2^2 + \|d\|_2^2,$$

da Q unitär, für alle $x \in \mathbb{K}^n$. Falls $\text{rank}(A) = \text{rank}(R_1) = n$, so gilt für die Lösung $R_1x = c$. Also hat man:

Algorithmus 2.130 (QR Lösung des Ausgleichsproblem)

Input: $A \in \mathbb{K}^{m,n}$ mit $\text{rank}(A) = n$, und $b \in \mathbb{K}^m$.

Output: $x \in \mathbb{K}^n$, so daß $\|Ax - b\|_2 \stackrel{!}{=} \min$.

Verwende Algorithmus 2.100, um A mit seiner QR -Zerlegung zu überschreiben.

FOR $j = 1 : n$

$v(j) = 1;$

$v(j+1 : m) = A(j+1 : m, j);$

$b(j : m) = \text{ROWHOUSE}(b(j : m), v(j : m));$

END

Löse $R(1 : n, 1 : n)x = b(1 : n)$ mit Rückwärts-Einsetzen.

Kosten: $2n^2(m - \frac{n}{3})$ flops.

Fehleranalyse: Das berechnete \tilde{x} löst $\|(A + \delta A)x - (b + \delta b)\|_2$ wobei

$$\|\delta A\|_F \leq (6m - 3n + 41)n \text{ eps} \|A\|_F + \mathcal{O}(\text{eps}^2)$$

$$\|\delta b\|_2 \leq (6m - 3n + 40)n \text{ eps} \|b\|_2 + \mathcal{O}(\text{eps}^2)$$

! **Vorsicht!** Probleme tauchen auf, wenn $\kappa_2(A) = \kappa_2(R) \approx \frac{1}{\text{eps}}$, oder natürlich wenn $\text{Rang}(A) < n$.

Abhilfe schafft hier die sogenannte Singulärwertzerlegung $A = U\Sigma V^*$ mit U, V unitär, $\Sigma =$

$$\begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \\ 0 & & \sigma_n & \\ \hline & & & 0 \end{bmatrix}. \quad \text{Dies kann zu diesem Zeitpunkt hier nicht gemacht werden.}$$

2.5 Iterative Verfahren

Für viele der Gleichungssysteme, die in der Praxis auftauchen, gilt $n \gg 10000$ und typischerweise $a_{ij} = 0$ für $> 90\%$ der Einträge. Wenn man derartige Probleme mit Gauß oder QR löst, kann es passieren, das in den Faktoren L, Q, R alle Einträge $\neq 0$ sind. Extrembeispiel Gauß für die Matrix

$$\begin{bmatrix} * & * & \cdots & * \\ * & * & & 0 \\ \vdots & & \ddots & \\ * & 0 & & * \end{bmatrix}.$$

Die Matrix läuft total voll. Es gibt sparse Varianten von Gauß/Cholesky/QR. Dies ist Thema einer Spezialvorlesung. Eine wichtige Idee, die insbesondere im Zusammenhang mit Vektorrechnern und Parallelrechnern von größter Bedeutung ist, ist die Konstruktion von Verfahren, die die Matrix an sich nicht verändern, und wenn möglich auf Matrix * Matrix oder Matrix * Vektor-Operationen aufgebaut sind. Das führt auf iterative Verfahren. Im wesentlichen betrachten wir hier 2 Ansätze.

2.5.1 Splitting-Verfahren

Zerlege A additiv

$$A = M - N, \tag{2.131}$$

wobei M^{-1} eine Approximation an A^{-1} ist, die leicht invertierbar ist, z.B. diagonal, block-diagonal, Δ -Matrix, ... Ersetze $Ax = b$ durch

$$Mx = Nx + b \iff x = M^{-1}Nx + M^{-1}b. \tag{2.132}$$

Dies ist eine Fixpunktgleichung.

Kanonisches Fixpunktverfahren:

$$x^{(k+1)} = M^{-1}Nx^{(k)} + M^{-1}b, \text{ mit gegebenen } x^{(0)}. \tag{2.133}$$

Satz 2.134 *Das Fixpunktverfahren (2.133) konvergiert für alle Startwerte $x^{(0)}$ gegen die Lösung von $Ax = b$ genau dann, wenn der Spektralradius $\rho(M^{-1}N) < 1$ ist, wobei $\rho(T) = \max\{|\lambda| : \lambda \text{ Eigenwert von } T\}$.*

Beweis: Banachscher Fixpunktsatz oder:

Sei $e^{(k)} = x - x^{(k)}$ der Fehler in der k -ten Iterierten, so gilt:

$$\begin{aligned} e^{(k+1)} &= x - x^{(k+1)} \\ &= (M^{-1}Nx + M^{-1}b) - (M^{-1}Nx^{(k)} + M^{-1}b) \\ &= M^{-1}N(x - x^{(k)}) \\ &= M^{-1}Ne^{(k)} \\ &= (M^{-1}N)^k e^{(0)}. \end{aligned}$$

Wenn Konvergenz für alle $x^{(0)}$ (also auch für alle $e^{(0)}$) gelten soll, dann ist dies äquivalent zu $\lim_{k \rightarrow \infty} (M^{-1}N)^k = 0$.

Sei $M^{-1}N = PJP^{-1}$ die Transformation auf Jordan-Normalform, d.h.

$$J = \begin{bmatrix} J_1 & & \\ & \ddots & \\ & & J_k \end{bmatrix}, \quad J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix},$$

so gilt $\lim_{k \rightarrow \infty} (M^{-1}N)^k = 0 \iff \lim_{k \rightarrow \infty} J^k = 0$.

J ist obere Δ -Matrix mit Eigenwerten auf der Diagonale.

Also muß $(\lambda_i)^k \rightarrow 0 \implies |\lambda_i| < 1 \implies \rho(J) < 1 \implies \rho(M^{-1}N) < 1$.

Umgekehrt, falls $\rho(J) < 1$, so folgt $\lim_{k \rightarrow \infty} J^k = 0$. Man schaue sich dann die Form von J^k an. □

Mögliche Wahlen für M, N . Setze

$$\begin{aligned} A &= \begin{matrix} D & & \\ & - & L \\ & & & - & U \end{matrix} \\ &= \begin{bmatrix} & & & & \\ & \diagdown & & & \\ & & \triangle & & \\ & & & \triangle & \\ & & & & \diagup \end{bmatrix} \end{aligned}$$

D Diagonale, L unteres \triangle , U oberes \triangle .

$M = D$, $N = L + U$ ergibt das Verfahren

$$x^{(k+1)} = D^{-1}(L + U)x^{(k)} + D^{-1}b \quad (2.135)$$

Algorithmus 2.136 (Jacobi oder Gesamtschrittverfahren)

Input: $A \in \mathbb{K}^{n,n}$, $b \in \mathbb{K}^n$, $a_{ii} \neq 0, i = 1, \dots, n$ und Startvektor $x^{(0)} = [x_1^{(0)}, \dots, x_n^{(0)}]^T$ sowie Abbruchgrenze δ .

Output: Näherung \tilde{x} für Lösung $Ax = b$

FOR $k = 0, 1, 2, \dots$

FOR $i = 1 : n$

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

END

IF $\|x^{(k+1)} - x^{(k)}\| < \delta$, STOP

END

Kosten: im wesentlichen ein Matrix-Vektor Produkt (unter Ausnutzung der Sparsität), 1 Vektoraddition pro Iterationsschritt.

Fehleranalyse: Wie bei Matrix * Vektor, Vektor+Vektor. Verfahrensfehler abhängig von $\delta, \rho(D^{-1}(L + U))$.

$M = D - L$, $N = U$ ergibt das Verfahren

$$x^{(k+1)} = (D - L)^{-1}Ux^{(k)} + (D - L)^{-1}b \quad (2.137)$$

Algorithmus 2.138 (Gauß-Seidel oder Einzelschrittverfahren)

Input: $A \in \mathbb{K}^{n,n}$, $b \in \mathbb{K}^n$, $a_{ii} \neq 0, i = 1, \dots, n$ und Startvektor $x^{(0)} = [x_1^{(0)}, \dots, x_n^{(0)}]^T$ sowie Abbruchgrenze δ .

Output: Näherung \tilde{x} für Lösung $Ax = b$

FOR $k = 0, 1, 2, \dots$

FOR $i = 1 : n$

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

END

IF $\|x^{(k+1)} - x^{(k)}\| < \delta$, STOP

END

Kosten: im wesentlichen wie Jacobi, nur schwerer als Vektoroperation zu implementieren.

Fehleranalyse: Fehler abhängig von Fehler in Lösung von $(D - L)x = c$ sowie Verfahrensfehler abhängig von $\rho((D - L)^{-1}U)$, δ .

Beide Verfahren sind einfach, aber konvergieren i.a. sehr langsam, insbesondere bei den Problemen die in der Praxis auftauchen. In Spezialfällen gilt: Falls Jacobi konvergiert, so konvergiert Gauß-Seidel schneller, falls Jacobi divergiert, so divergiert Gauß-Seidel schneller (Satz von Stein/Rosenberg).

Satz 2.139 (Ostrowski/Reich) Sei $A \in \mathbb{R}^{n,n}$ symmetrisch positiv definit, so konvergiert das Gauß-Seidel Verfahren für alle $x^{(0)}$.

Beweis: Siehe Stoer 2, Golub/Van Loan. □

2.5.2 Konvergenzbeschleunigung

Eine Möglichkeit, die Konvergenz zu beschleunigen ist die sukzessive Überrelaxation (SOR).

$M = D - \omega L, N = (1 - \omega)D + \omega U, \omega \in (0, 2), (M - N) = \omega(D - L - U)$ ergibt Verfahren

$$x^{(k+1)} = (D - \omega L)^{-1}((1 - \omega)D + \omega U)x^{(k)} + (D - \omega L)^{-1}\omega b. \quad (2.140)$$

Algorithmus 2.141 (SOR)

Input: $A \in \mathbb{K}^{n,n}, b \in \mathbb{K}^n, a_{ii} \neq 0, i = 1, \dots, n$ und Startvektor $x^{(0)} = [x_1^{(0)}, \dots, x_n^{(0)}]^T$ sowie Abbruchgrenze δ .

Output: Näherung \tilde{x} für Lösung $Ax = b$

FOR $k = 0, 1, 2 \dots$

FOR $i = 1 : n$

$$x_i^{(k+1)} = \omega \left(b_i - \sum_{j=1}^{i-1} a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)} \right) / a_{ii} + (1 - \omega)x_i^{(k)};$$

END

IF $\|x^{(k+1)} - x^{(k)}\| < \delta$, STOP

END

Kosten und Fehler wie oben ($\omega = 1$: Gauß-Seidel).

In einigen Spezialfällen von strukturierten Matrizen aus partiellen Differentialgleichungen gibt es Konvergenzbeweise, optimales ω , etc. Weitere Beschleunigungsmethoden: Tschebyschow-Beschleunigung, semiiterative Verfahren, unvollständige Dreieckszerlegungen.

2.5.3 Konjugierte Gradienten

Grundidee: Betrachte quadratisches Funktional

$$\Phi(x) = \frac{1}{2}x^T Ax - x^T b, \quad A \in \mathbb{R}^{n,n} \text{ symmetrisch, positiv definit, } b \in \mathbb{R}^n.$$

Minimum von Φ ist der Wert $-b^T A^{-1}b/2$. Dieser wird erreicht bei $x = A^{-1}b$. Also ist Minimierung von $\Phi(x)$ äquivalent zur Lösung von $Ax = b$.

Methode des Steilsten Abstiegs: In jedem Schritt lege die Richtung des steilsten Abstiegs fest und gehe in diese Richtung so weit wie möglich. Beim Punkt x_c angelangt, ist die Richtung des steilsten Abstiegs die Richtung des negativen Gradienten.

$$-\nabla\Phi|_{x_c} = b - Ax_c = r_c, \quad \text{Residuum von } x_c. \quad (2.142)$$

Falls $r_c \neq 0$, dann gibt es ein α , so daß $\Phi(x_c + \alpha r_c) < \Phi(x_c)$. Wähle

$$\alpha = \frac{r_c^T r_c}{r_c^T A r_c}. \quad (2.143)$$

Das ergibt Minimum von $\Phi(x_c + \alpha r_c)$ über α . Man hätte damit folgende Methode des steilsten Abstiegs.

```

k = 0; x_0 = 0; r_0 = b;
WHILE r_k ≠ 0
    k = k + 1;
    α_k = r_{k-1}^T r_{k-1} / r_{k-1}^T A r_{k-1};
    x_k = x_{k-1} + α_k r_{k-1};
    r_k = b - A x_k;
END

```

Sehr langsame Konvergenz wenn $\kappa_2(A)$ groß. Besser ist es, die Suchrichtungen so zu wählen, daß die Residuen senkrecht aufeinanderstehen, so daß man in exakter Rechnung in $\leq n$ Schritten fertig ist.

Algorithmus 2.144 (Konjugierte Gradienten–Verfahren)

Input: $A \in \mathbb{R}^{n,n}$, symmetrisch, positiv definit, $b \in \mathbb{R}^n$, k_{\max} (Maximum an Iterationsschritten), $\varepsilon > 0$ Toleranz.

Output: Lösung von $Ax = b$.

```

k = 0; x_0 = 0; r_0 = b; ρ_0 = ||r_0||_2^2;
WHILE √ρ_k > ε√ρ_0 AND k < k_max
    k = k + 1;
    IF k = 1
        p = r;
    ELSE
        β = ρ_{k-1} / ρ_{k-2}; % (= r_{k-1}^T r_{k-1} / r_{k-2}^T r_{k-2})
        p = r + βp; % (= p_k)
    END
    w = Ap;
    α = ρ_{k-1} / p^T w; % (= r_{k-1}^T r_{k-1} / p_k^T A p_k)
    x = x + αp; % (= x_k)
    r = r - αw; % (= r_k)
    ρ_k = ||r||_2^2; % (= r_k^T r_k)
END

```

Kosten: 1 Matrix * Vektor–Multiplikation + 2 Skalarprodukte + 2 GAXPY pro Schritt.
Fehleranalyse: Rundungsfehler zerstören Konvergenz in n Schritten und Orthogonalität der Residuen. Bemerkung:

- Für die Vektoren $p \equiv p_k$, $r \equiv r_k$ aus Algorithmus 2.144 gilt: $p_k^T A p_l = 0$, $r_k^T r_l = 0$ für alle $k > l \geq 0$.

Übung 2.150 Sei

$$A = \begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{n,n}.$$

Zeige

1. Die Eigenwerte von A sind $\lambda_j = 2 - 2 \cos \frac{j\pi}{n+1}$, und die zugehörigen Eigenvektoren lauten $v_j = (v_{ij})_{i=1,\dots,n} = \left(\sin \frac{ij\pi}{n+1} \right)_{i=1,\dots,n}$ für alle $j = 1, \dots, n$.
Hinweis: $\sin(x+y) = \sin x \cos y + \cos x \sin y$, $\cos(-x) = \cos x$, $\sin(-x) = -\sin x$.
2. Für jede mit der Matrixmultiplikation verträgliche Norm $\|\bullet\|$ gilt: $\rho(A) \leq \|A\|$.
3. In der Standard-Zerlegung $A = D - L - U$ ist das Einzelschrittverfahren definiert durch $M = D - L, N = U$ und die Iterationsmatrix $T_E = M^{-1}N$. Zeige: $T_E = (I + D^{-1}L + \dots + (D^{-1}L)^{n-1})D^{-1}U$.
4. Bestimme/schätze den Spektralradius von Gesamt-/Einzelschrittverfahren. Was kann man über die Konvergenzgeschwindigkeit des cg-Verfahrens sagen?

Übung 2.151 Sei $A \in \mathbb{R}^{n,n}$, $A = D - L - U$ Standardzerlegung, d.h. D ist Diagonale von A , $-L$ striktes unteres Dreieck und $-U$ striktes oberes Dreieck von A . Sei $\omega \in \mathbb{R}$ und betrachte die Zerlegung $\omega A = M - N$ mit $M = D - \omega L$, $N = (1 - \omega)D + \omega U$.
Zeige: Ist $\rho(M^{-1}N) < 1$, so ist $0 < \omega < 2$.

Übung 2.152 Schreibe **MATLAB**-Programme für das Jacobi-Verfahren, Gauß-Seidel-Verfahren, sowie das cg-Verfahren ohne Vorkonditionierung. Vergleiche die Verfahren bezüglich Rechenzeit und Iterationszahl anhand des Beispiels $Ax = b$,

$$A = \begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix}, \quad b = \begin{pmatrix} 1 \\ \vdots \\ \vdots \\ 1 \end{pmatrix},$$

$x^{(0)} = 0$, $n = 10, 100, 1000$. Als Abbruchkriterium nehme man $\|r_{k+1}\|_2 = \|b - Ax_{k+1}\|_2 < 10^{-8}$, $k_{\max} = 600$.

Übung 2.153 Sei $A \in \mathbb{R}^{n,n}$ symmetrisch, positiv definit und $b, x_0 \in \mathbb{R}^n$. Betrachte das cg-Verfahren ohne Vorkonditionierung aus der Vorlesung (zur besseren Analyse mit Indizes).

```

k = 0, r_0 = b - Ax_0, rho_0 = r_0^T r_0
while sqrt(rho_k) > epsilon * sqrt(rho_0) and k < k_max
    k = k + 1
    if k = 1
        p_1 = r_0
    else
        beta_k = rho_{k-1} / rho_{k-2}
        p_k = r_{k-1} + beta_k p_{k-1}
    end
    w_k = Ap_k, alpha_k = rho_{k-1} / w_k^T p_k
    x_k = x_{k-1} + alpha_k p_k
    r_k = r_{k-1} - alpha_k w_k
    rho_k = r_k^T r_k
end

```

Zeige:

1. $p_k^\top A p_l = 0$, für alle $k > l \geq 0$.

Tip: Zeige durch Induktion über k , $r_{k-1}^\top r_{l-1} = 0, p_k^\top A p_l = 0$, für alle $k > l > 0$.

2. Nehme an, es sei $p_k^\top A p_l = 0$, für alle $k > l > 0$, dann ist

$$\min_{v_1, \dots, v_k \in \mathbb{R}} \Phi(x_0 + \sum_{l=1}^k p_l v_l) = \min_{v_k \in \mathbb{R}} \left(\min_{v_1, \dots, v_{k-1} \in \mathbb{R}} \Phi(x_0 + \sum_{l=1}^k p_l v_l) \right),$$

wobei $\Phi(x) = \frac{1}{2} x^\top A x - x^\top b$. Was passiert mit dem cg-Verfahren nach maximal n Schritten?

Übung 2.154 Sei $A \in \mathbb{R}^{n,n}$ symmetrisch, positiv definit und $b \in \mathbb{R}^n$. Definiere $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}$ durch $\Phi(x) = \frac{1}{2} x^\top A x - x^\top b$. Zeige:

1. Φ besitzt ein eindeutiges Minimum in $x = A^{-1}b$

2. Sei $x_0 \in \mathbb{R}^n$ und $P = [p_1, \dots, p_k] \in \mathbb{R}^{n,k}$ so, dass $P^\top A P = D = \text{diag}(d_1, \dots, d_k)$ mit von Null verschiedenen d_1, \dots, d_k . Dann besitzt die Funktion $\Psi : \mathbb{R}^k \rightarrow \mathbb{R}$, $\Psi(v) = \Phi(x_0 + P v)$ ein eindeutiges Minimum in $v = D^{-1} P^\top (b - A x_0) = \left(\frac{p_i^\top (b - A x_0)}{d_i} \right)_{i=1, \dots, k}$

Übung 2.155 Sei $A \in \mathbb{R}^{n,n}$ symmetrisch, positiv definit und $b, x_0 \in \mathbb{R}^n$. Betrachte das cg-Verfahren ohne Vorkonditionierung aus der Vorlesung (zur besseren Analyse mit Indizes).

$k = 0, r_0 = b - A x_0, \rho_0 = r_0^\top r_0$

while $\sqrt{\rho_k} > \varepsilon \sqrt{\rho_0}$ **and** $k < k_{\max}$

$k = k + 1$

if $k = 1$

$p_1 = r_0$

else

$\beta_k = \rho_{k-1} / \rho_{k-2}$

$p_k = r_{k-1} + \beta_k p_{k-1}$

end

$w_k = A p_k$

$\alpha_k = \rho_{k-1} / w_k^\top p_k$

$x_k = x_{k-1} + \alpha_k p_k$

$r_k = r_{k-1} - \alpha_k w_k$

$\rho_k = r_k^\top r_k$

end

Zeige:

1. $\text{span}\{r_0, \dots, r_{k-1}\} = \text{span}\{p_1, \dots, p_k\}$

2. $\text{span}\{r_0, \dots, r_{k-1}\} = \text{span}\{r_0, A r_0, \dots, A^{k-1} r_0\}$

für alle $k = 1, \dots, n$, sofern der Algorithmus nicht vorher abbricht.

Kapitel 3

Lösung nichtlinearer Gleichungssysteme

Neben den linearen Gleichungssystemen führen viele Anwendungsprobleme auf die Lösung von nichtlinearen Gleichungssystemen. Gegeben

$$f : E \longrightarrow F \quad (\text{z.B. hier } E = F = \mathbb{R}^n)$$
$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \longmapsto \begin{bmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_n(x_1, \dots, x_n) \end{bmatrix}$$

Gesucht: Vektor $[x_1, \dots, x_n]^T$ so daß $f_i(x_1, \dots, x_n) = 0$, für alle $i = 1, \dots, n$.

E, F können auch unendlich dimensionale lineare Räume sein, etwa Räume von Funktionen und f kann auch ein Differentialoperator sein. Hier betrachten wir nur den Fall $E = F = \mathbb{R}^n$. Im allgemeinen ist die Lösung nur durch Näherungsverfahren zu bestimmen. Nullstellenprobleme sind eng verwandt mit Minimierungsproblemen wie

$$\text{Bestimme für } h : \mathbb{R}^n \longrightarrow \mathbb{R} \quad \min_{x \in \mathbb{R}^n} h(x). \quad (3.1)$$

In üblicher Weise muß zur Bestimmung des Minimums

$$\text{grad } (h(x)) = \begin{bmatrix} \frac{\partial h}{\partial x_1} \\ \vdots \\ \frac{\partial h}{\partial x_n} \end{bmatrix} = 0$$

berechnet werden, welches wiederum ein Nullstellenproblem ist. Neben Optimierungsproblemen treten Nullstellenbestimmungen in fast allen Anwendungen auf. Gewöhnliche oder partielle Differentialgleichungen, Ausgleichsrechnung, ...

3.1 Fixpunktverfahren

Wandle das Nullstellenproblem in ein Fixpunktproblem um.

$$x = \Phi(x). \quad (3.2)$$

Dies ist im allgemeinen auf viele Arten möglich, die äquivalent oder nicht äquivalent sein können.

Beispiel 3.3 $e^x - \sin x = 0$

Mögliche Fixpunktgleichungen (nicht alle äquivalent)

- a) $x = e^x - \sin x + x$
- b) $x = \sin x - e^x + x$
- c) $x = \arcsin(e^x)$, für $x < 0$
- d) $x = \ln(\sin x)$, für $(-2n\pi, 2n\pi + \pi), n = 1, 2, 3, \dots$

Alle haben unterschiedliche numerische Eigenschaften.

→ Übung 3.50

Die Idee des Fixpunktverfahrens ist es, die Folge

$$x^{(i+1)} = \Phi(x^{(i)}), \quad i = 0, 1, 2, \dots \quad (3.4)$$

mit gegebenem Startwert x_1 zu erzeugen.**Satz 3.5 (Banachscher Fixpunktsatz)** Sei $D \subseteq D_\Phi \subseteq \mathbb{R}^m$ ein abgeschlossenes Gebiet und Φ eine kontrahierende Abbildung von D in sich, d.h.

- a) $\Phi : D \rightarrow D$
- b) Φ kontrahierend, d.h. $\exists \alpha < 1$ und eine Norm $\|\bullet\|$, so daß

$$\|\Phi(x) - \Phi(y)\| \leq \alpha \|x - y\|, \quad \forall x, y \in D. \quad (3.6)$$

Dann gilt:

- i) Es gibt genau einen Fixpunkt \hat{x} von (3.2) in D .
- ii) Die Fixpunktiteration (3.4) konvergiert für jeden Startwert $x^{(0)} \in D$ gegen \hat{x} .
- iii) Es gelten die Fehlerabschätzungen

$$\|\hat{x} - x^{(n)}\| \leq \frac{\alpha^n}{1 - \alpha} \|x^{(1)} - x^{(0)}\| \quad \text{a priori} \quad (3.7)$$

$$\|\hat{x} - x^{(n)}\| \leq \frac{\alpha}{1 - \alpha} \|x^{(n)} - x^{(n-1)}\| \quad \text{a posteriori} \quad (3.8)$$

Beweis: $x^{(0)} \in D \implies x^{(n)} = \Phi(x^{(n-1)}) \in D$, für alle $n = 1, 2, \dots$

$$\begin{aligned} \|x^{(n+1)} - x^{(n)}\| &= \|\Phi(x^{(n)}) - \Phi(x^{(n-1)})\| \\ &\leq \alpha \|x^{(n)} - x^{(n-1)}\| \leq \alpha^2 \|x^{(n-1)} - x^{(n-2)}\| \\ &\leq \alpha^n \|x^{(1)} - x^{(0)}\| \end{aligned}$$

$$\begin{aligned} \|x^{(n+k)} - x^{(n)}\| &\leq \sum_{i=1}^k \|x^{(n+i)} - x^{(n+i-1)}\| \\ &\leq \sum_{i=1}^k \alpha^{n+i-1} \|x^{(1)} - x^{(0)}\| \\ &\leq \alpha^n \|x^{(1)} - x^{(0)}\| \sum_{i=1}^k \alpha^{i-1} \\ &\leq \frac{\alpha^n}{1 - \alpha} \|x^{(1)} - x^{(0)}\| \end{aligned}$$

\implies Folge konvergiert, da Cauchy-Folge und \mathbb{R}^m vollständig.

Φ ist stetig. Dies folgt bereits aus der Kontraktionseigenschaft (3.6). Somit ist der Grenzwert Fixpunkt von Φ , wegen

$$\hat{x} \equiv \lim_{n \rightarrow \infty} x^{(n+1)} = \lim_{n \rightarrow \infty} \Phi(x^{(n)}) = \Phi(\lim_{n \rightarrow \infty} x^{(n)}) = \Phi(\hat{x}).$$

Eindeutigkeit. Seien \hat{x}, \bar{x} Fixpunkte von Φ .

$$\|\hat{x} - \bar{x}\| = \|\Phi(\hat{x}) - \Phi(\bar{x})\| \leq \alpha \|\hat{x} - \bar{x}\|.$$

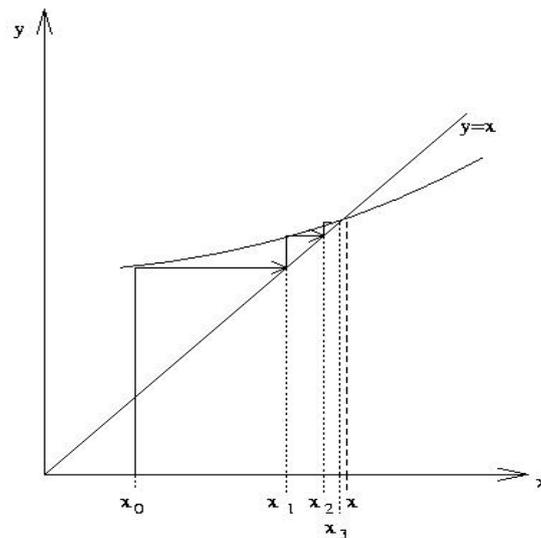
Wegen $\alpha < 1$ kann nur $\hat{x} = \bar{x}$ gelten.

$$\begin{aligned} \|x^{(n+k)} - x^{(n)}\| &\leq \frac{\alpha^n}{1-\alpha} \|x^{(1)} - x^{(0)}\| \\ \xrightarrow{k \rightarrow \infty} \|\hat{x} - x^{(n)}\| &\leq \frac{\alpha^n}{1-\alpha} \|x^{(1)} - x^{(0)}\| \iff (3.7) \end{aligned}$$

$$n = 1 \implies \|\hat{x} - x^{(1)}\| \leq \frac{\alpha}{1-\alpha} \|x^{(1)} - x^{(0)}\|$$

ersetze $x^{(1)}$ durch $x^{(n)}$, $x^{(0)}$ durch $x^{(n-1)}$ so folgt (3.8). □

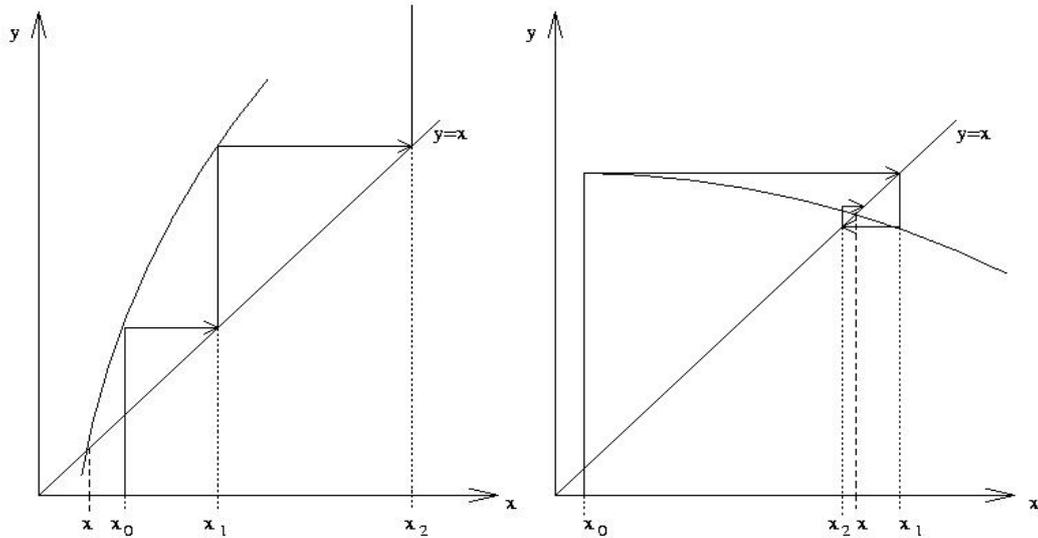
Graphik 3.9 Graphische Veranschaulichung der Konvergenz $m = 1$.



Definition 3.10 Ein Fixpunkt \hat{x} heißt anziehend, wenn die Iteration (3.4) für alle genügend nahe bei \hat{x} gelegenen Startwerte $x^{(0)} \neq \hat{x}$ gegen \hat{x} konvergiert.

Ein Fixpunkt \hat{x} heißt abstoßend, wenn es eine offene Kugel U_ε um \hat{x} gibt, so daß für alle $x^{(0)} \in U_\varepsilon \setminus \{\hat{x}\}$ ein m existiert, so daß $x^{(m)} \notin U_\varepsilon$.

Graphik 3.11 *Abstoßender Fixpunkt / Anziehender Fixpunkt.*



Satz 3.12 Sei $D \subseteq D_\Phi \subseteq \mathbb{R}^m$ abgeschlossen und konvex und nehme an, daß

$$D\Phi = \begin{bmatrix} \frac{\partial \Phi_1(x)}{\partial x_1} & \dots & \frac{\partial \Phi_1(x)}{\partial x_m} \\ \vdots & & \vdots \\ \frac{\partial \Phi_m(x)}{\partial x_1} & \dots & \frac{\partial \Phi_m(x)}{\partial x_m} \end{bmatrix} \text{ existiert und stetig ist. Falls es eine Norm gibt, so daß}$$

$$\|D\Phi(x)\| \leq \alpha < 1, \text{ so ist } \Phi \text{ in } D \text{ kontrahierend mit Lipschitzkonstante } \alpha.$$

Beweis: \implies Übung 3.49. □

Ist $D\Phi(\hat{x})$ stetig und $\|D\Phi(\hat{x})\| < 1$ so ist \hat{x} anziehend. Ein Fixpunkt ist im allgemeinen nicht anziehend, falls einer der Eigenwerte von $D\Phi(x)$ vom Betrag größer 1 ist.
 \implies Übung.

Definition 3.13 Eine Iterationsfolge $\{x^{(i)}\}_{i=0}^\infty$, die für $p \geq 1$

$$\|x^{(i+1)} - \xi\| \leq c \|x^{(i)} - \xi\|^p \quad i = 0, 1, 2, \dots \quad (3.14)$$

und $c < 1$ für $p = 1$ erfüllt, heißt Folge von mindestens p -ter Ordnung.

Korollar 3.15 Das Fixpunktverfahren (3.4) sofern es konvergiert, ist konvergent von mindestens 1. Ordnung (linear konvergent).

Kosten für Iteration (3.4): eine Auswertung von $\Phi(x)$ pro Iteration, Anzahl der Iterationen aus Fehlerschätzung ablesbar.

Versuche höhere Konvergenzordnung zu bekommen.

3.2 Das Newtonverfahren

Idee: Taylorentwicklung von Nullstelle \hat{x} von $f(x)$, für $x^{(0)}$ aus Umgebung.

$$0 = f(\hat{x}) = f(x^{(0)}) + Df(x^{(0)})(\hat{x} - x^{(0)}) + \underbrace{\mathcal{O}(\|\hat{x} - x^{(0)}\|^2)}_{\text{weglassen}}. \quad (3.16)$$

Ersetze durch

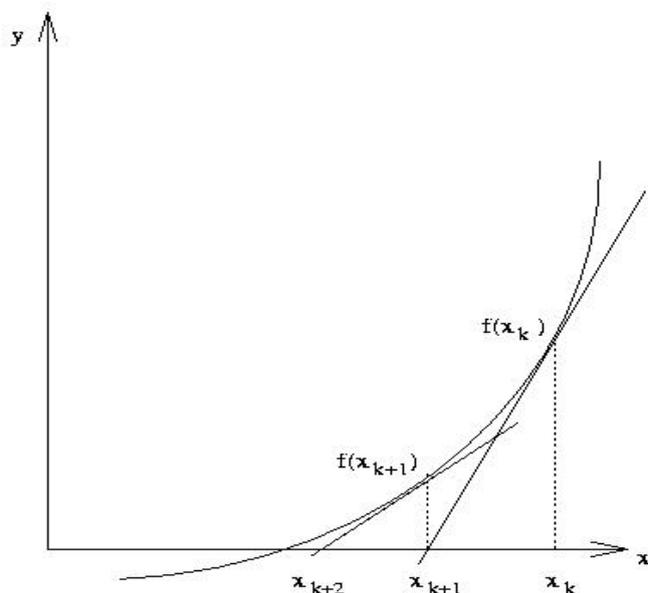
$$0 = f(x^{(0)}) + Df(x^{(0)})(x^{(1)} - x^{(0)}) \quad (3.17)$$

Angenommen $Df(x^{(0)})$ nichtsingulär, so folgt

$$Df(x^{(0)})(x^{(1)} - x^{(0)}) = -f(x^{(0)})$$

hat eindeutige Lösung $\delta^{(1)} = x^{(1)} - x^{(0)}$ und wir erhalten $x^{(1)} = x^{(0)} + \delta^{(1)}$.
 → Übung.

Graphik 3.18 *Newtonschritt*



Satz 3.19 Sei $D \subseteq \mathbb{R}^n$ offen und sei D_0 konvex mit $\bar{D}_0 \subseteq D$. $f : D \rightarrow \mathbb{R}^n$ sei für alle $x \in D_0$ differenzierbar und für alle $x \in D$ stetig. Sei $x^{(0)} \in D_0$, so dass es Konstanten $r, \alpha, \beta, \gamma, h$ gibt mit

$$\begin{aligned} S_r(x^{(0)}) &:= \{x \mid \|x - x^{(0)}\| < r\} \subseteq D_0 \\ h &= \alpha\beta\gamma/2 < 1 \\ r &= \alpha/(1 - h), \end{aligned}$$

und $f(x)$ erfülle

$$\|Df(x) - Df(y)\| \leq \gamma\|x - y\|, \quad \forall x, y \in D_0. \quad (3.20)$$

$$Df(x)^{-1} \text{ existiere } \forall x \in D_0 \text{ und sei } \|Df(x)^{-1}\| \leq \beta, \quad \forall x \in D_0, \quad (3.21)$$

sowie

$$\|Df(x^{(0)})^{-1}f(x^{(0)})\| \leq \alpha. \quad (3.22)$$

Dann gilt

a) Ausgehend von $x^{(0)}$ ist jedes $x^{(i)}$, welches durch

$$x^{(i+1)} = x^{(i)} - Df(x^{(i)})^{-1}f(x^{(i)}), \quad i = 0, 1, 2, \dots$$

gegeben ist, wohldefiniert und es ist $x^{(i)} \in S_r(x^{(0)})$, $\forall i = 0, 1, 2, \dots$

b) $\lim_{i \rightarrow \infty} x^{(i)} = \hat{x}$ existiert und es ist $f(\hat{x}) = 0$, $\hat{x} \in \overline{S_r(x^{(0)})}$.

c) Für alle $i \geq 0$ gilt $\|x^{(i)} - \hat{x}\| \leq \alpha \frac{h^{2^i-1}}{1-h^{2^i}}$

Das Newton-Verfahren ist also mindestens quadratisch konvergent.

Beweis: a) Da $Df(x)^{-1}$ für alle $x \in D_0$ existiert, so ist $x^{(i+1)}$ für alle i dann wohldefiniert, wenn $x^{(i)} \in S_r(x^{(0)})$ für alle $i \geq 0$. Wegen Voraussetzung (3.22) gilt dies für $x^{(0)}, x^{(1)}$.

I.V. Sei nun $x^{(j)} \in S_r(x^{(0)})$ für $j = 0, 1, \dots, k, k \geq 1$, so folgt aus (3.21)

$$\begin{aligned} \|x^{(k+1)} - x^{(k)}\| &= \|-Df(x^{(k)})^{-1}f(x^{(k)})\| \leq \beta \|f(x^{(k)})\| \\ &= \beta \|f(x^{(k)}) - f(x^{(k-1)}) - Df(x^{(k-1)})(x^{(k)} - x^{(k-1)})\|, \end{aligned}$$

da $f(x^{(k-1)}) + Df(x^{(k-1)})(x^{(k)} - x^{(k-1)}) = 0$.

Wende das folgende Lemma an: □

Lemma 3.23 $Df(x)$ existiere für alle $x \in D_0, D_0 \subseteq \mathbb{R}^n, D_0$ konvex und es gebe ein γ , so dass

$$\|Df(x) - Df(y)\| \leq \gamma \|x - y\|, \quad \forall x, y \in D_0.$$

Dann gilt $\forall x, y \in D_0$ die Abschätzung

$$\|f(x) - f(y) - Df(y)(x - y)\| \leq \frac{\gamma}{2} \|x - y\|^2$$

Beweis: Sei $\varphi : [0, 1] \rightarrow \mathbb{R}^n$ gegeben durch $\varphi(t) := f(y + t(x - y))$. φ ist differenzierbar $\forall t \in [0, 1], \forall x, y \in D_0$ und $\varphi'(t) = Df(y + t(x - y))(x - y)$

$$\implies \|\varphi'(t) - \varphi'(0)\| \leq \|Df(y + t(x - y)) - Df(y)\| \|x - y\| \leq \gamma t \|x - y\|^2.$$

Nun ist

$$P \equiv f(x) - f(y) - Df(y)(x - y) = \varphi(1) - \varphi(0) - \varphi'(0) = \int_0^1 (\varphi'(t) - \varphi'(0)) dt,$$

also folgt

$$\begin{aligned} \|P\| &\leq \int_0^1 \|\varphi'(t) - \varphi'(0)\| dt \\ &\leq \gamma \|x - y\|^2 \int_0^1 t dt \\ &\leq \frac{\gamma}{2} \|x - y\|^2. \end{aligned}$$

□

Weiter im Beweis von Satz 3.19. Anwendung des Lemma liefert

$$\|x^{(k+1)} - x^{(k)}\| \leq \frac{\beta\gamma}{2} \|x^{(k)} - x^{(k-1)}\|^2 \quad (3.24)$$

$$\implies \|x^{(k+1)} - x^{(k)}\| \leq \alpha h^{2^k-1}. \quad (3.25)$$

Diese letzte Abschätzung folgt aus (3.22) für $k = 0$ und mit Induktion für $k + 1$ aus (3.24). Aus (3.25) folgt

$$\begin{aligned} \|x^{(k+1)} - x^{(0)}\| &\leq \|x^{(k+1)} - x^{(k)}\| + \dots + \|x^{(1)} - x^{(0)}\| \\ &\leq \alpha(1 + h + h^3 + h^7 + \dots + h^{2^k-1}) < \frac{\alpha}{1-h} = r, \end{aligned}$$

also $x^{(k+1)} \in S_r(x_0)$.

b) Aus (3.25) folgt, dass $x^{(k)}$ Cauchy-Folge ist, denn

$$\begin{aligned} \|x^{(n+k)} - x^{(n)}\| &\leq \sum_{i=1}^k \|x^{(n+i)} - x^{(n+i-1)}\| \\ &\leq \alpha \sum_{i=1}^k h^{2^{n+i-1}-1} \\ &\leq \alpha h^{2^n-1} (1 + h^{2^n} + (h^{2^n})^2 + \dots) \\ &\leq \frac{\alpha h^{2^n-1}}{1-h^{2^n}} < \varepsilon, \end{aligned} \quad (3.26)$$

für genügend großes $n \geq N(\varepsilon)$, weil $0 < h < 1$.

$$\implies \lim_{k \rightarrow \infty} x^{(k)} = \hat{x} \in \overline{S_r(x_0)}.$$

Aus (3.26) folgt:

$$\|\hat{x} - x^{(n)}\| = \lim_{m \rightarrow \infty} \|x^{(m)} - x^{(n)}\| \leq \frac{\alpha h^{2^n-1}}{1-h^{2^n}}.$$

Noch zu zeigen ist, dass \hat{x} Nullstelle von f in $\overline{S_r(x_0)}$ ist. Da $x^{(k)} \in S_r(x^{(0)}) \forall k \geq 0$, folgt

$$\|Df(x^{(k)}) - Df(x^{(0)})\| \leq \gamma \|x^{(k)} - x^{(0)}\| < \gamma r.$$

$$\implies \|Df(x^{(k)})\| \leq \gamma r + \|Df(x^{(0)})\| =: \kappa.$$

Aus der Iterationsgleichung folgt:

$$\|f(x^{(k)})\| \leq \kappa \|x^{(k+1)} - x^{(k)}\|.$$

$$\implies \lim_{k \rightarrow \infty} \|f(x^{(k)})\| = 0$$

und wegen f stetig gilt $\lim_{k \rightarrow \infty} f(x^{(k)}) = f(\hat{x})$. □

Man kann unter noch stärkeren Voraussetzungen mit dem Satz von Newton–Kantorovich zeigen, dass \hat{x} die einzige Nullstelle in $S_r(x^{(0)})$ ist.

Voraussetzung dafür, dass das Newton-Verfahren konvergiert, ist, dass der Startwert

genügend nahe bei der gesuchten Lösung liegt. Um globale Konvergenz zu erzielen, wird das Newton-Verfahren durch Einführung eines Parameters modifiziert. Setze

$$x^{(k+1)} = x^{(k)} - \lambda_k d^{(k)}, d^{(k)} := Df(x^{(k)})^{-1} f(x^{(k)}) \quad (3.27)$$

Die λ_k werden dabei so gewählt, dass die Folge $\{h(x^{(k)})\}$ mit $h(x) = f(x)^T f(x)$ streng monoton fällt und die $x^{(k)}$ gegen ein Minimum von $h(x)$ konvergieren.

→ Übung.

Was hat das mit Nullstelle zu tun? Da

$$\begin{aligned} h(x) &\geq 0 && \forall x, \\ h(\hat{x}) = 0 &\iff f(\hat{x}) = 0 && (f(x)^T f(x) = \|f(x)\|_2^2), \end{aligned}$$

so ist jedes lokale Minimum \hat{x} von h mit $h(\hat{x}) = 0$ auch globales Minimum von h und Nullstelle von f .

Definiere Menge

$$D(\gamma, x) = \{s \in \mathbb{R}^n \mid \|s\| = 1 \text{ und } Dh(x)^T s \geq \gamma \|Dh(x)\|\}, \gamma > 0. \quad (3.28)$$

Lemma 3.29 Sei $h : \mathbb{R}^n \rightarrow \mathbb{R}$ eine Funktion, deren Gradient $Dh(x)$ für alle $x \in V(\hat{x})$ ($V(\hat{x})$ Umgebung von \hat{x}) definiert und stetig ist. Sei ferner $Dh(\hat{x}) \neq 0$ und $\gamma > 0$. Dann gibt es eine Umgebung $U(\hat{x}) \subseteq V(\hat{x})$ von \hat{x} und ein $\lambda > 0$, so dass

$$h(x - \mu s) \leq h(x) - \frac{\mu\gamma}{4} \|Dh(\hat{x})\|,$$

! für alle $x \in U(\hat{x})$, $s \in D(\gamma, x)$ und $0 \leq \mu \leq \lambda$.
(zeigt wann man $h(y) < h(x)$ finden kann.)

Beweis: Sei

$$U^1(\hat{x}) = \{x \in V(\hat{x}) \mid \|Dh(x) - Dh(\hat{x})\| < \frac{\gamma}{4} \|Dh(\hat{x})\|\}.$$

Da $Dh(\hat{x}) \neq 0$ und $Dh(x)$ stetig auf $V(\hat{x})$, so folgt $U^1(\hat{x}) \neq \emptyset$ und U^1 ist Umgebung von \hat{x} . Aus demselben Grund folgt auch, dass

$$U^2(\hat{x}) = \{x \in V(\hat{x}) \mid D(\gamma, x) \subseteq D(\frac{\gamma}{2}, \hat{x})\}$$

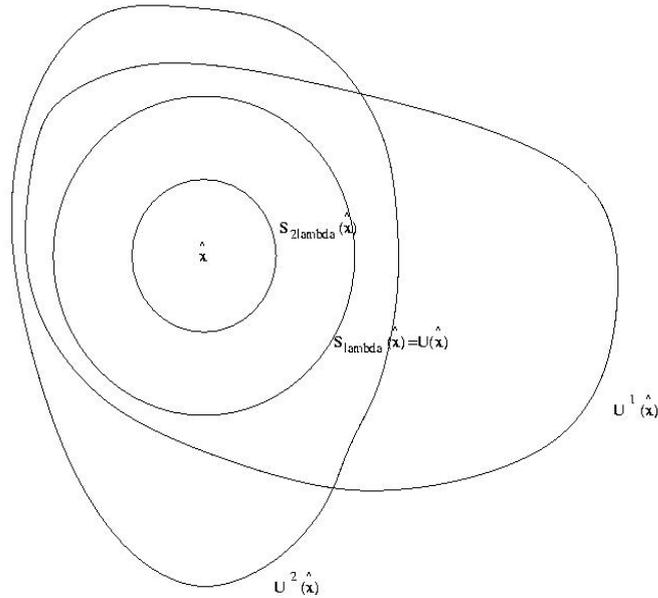
eine Umgebung von \hat{x} ist. Wähle nun $\lambda > 0$, so dass

$$\overline{S_{2\lambda}(\hat{x})} := \{x \mid \|x - \hat{x}\| \leq 2\lambda\} \subseteq U^1(\hat{x}) \cap U^2(\hat{x})$$

und setze

$$U(\hat{x}) = \overline{S_\lambda(\hat{x})} = \{x \mid \|x - \hat{x}\| \leq \lambda\}.$$

Graphik 3.30



Dann gilt für alle $x \in U(\hat{x})$ mit $0 \leq \mu \leq \lambda$, $s \in D(\gamma, x)$, dass

$$\begin{aligned} h(x) - h(x - \mu s) &= \mu Dh(x - \theta \mu s)^T s \\ &= \mu \left[(Dh(x - \theta \mu s) - Dh(\hat{x}))^T s + Dh(\hat{x})^T s \right], \end{aligned}$$

für ein θ mit $0 < \theta < 1$. Mit $x \in U(\hat{x})$ gilt auch

$$x, x - \mu s, x - \theta \mu s \in U^1 \cap U^2,$$

also folgt

$$\begin{aligned} h(x) - h(x - \mu s) &\geq -\frac{\mu\gamma}{4} \|Dh(\hat{x})\| + \mu Dh(\hat{x})^T s \\ &\geq \frac{-\mu\gamma}{4} \|Dh(\hat{x})\| + \mu \frac{\gamma}{2} \|Dh(\hat{x})\| \\ &= \frac{\mu\gamma}{4} \|Dh(\hat{x})\|. \end{aligned}$$

□

Algorithmus 3.31

Input: Differenzierbare Funktion $h(x) : \mathbb{R}^n \rightarrow \mathbb{R}$, und Zahlen $\sigma_k, \gamma_k, k = 0, 1, 2, \dots$ mit

$$\inf_k (\gamma_k) > 0, \inf_k \sigma_k > 0,$$

sowie einen Startwert $x^{(0)} \in \mathbb{R}^n$.

Output: Folge $x^{(k)}, k = 0, 1, 2, \dots$, die gegen Minimum von $h(x)$ konvergiert.

FOR $k = 0, 1, 2, \dots$ UNTIL SATISFIED

Sei $s^{(k)} \in D(\gamma_k, x^{(k)})$,

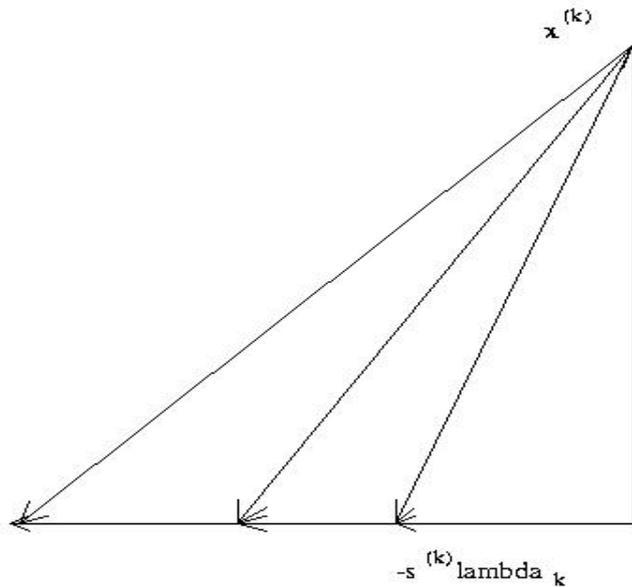
$$x^{(k+1)} = x^{(k)} - \lambda_k s^{(k)}, \quad (3.32)$$

wobei $\lambda_k \in [0, \sigma_k \|Dh(x^{(k)})\|]$, so dass

$$h(x^{(k+1)}) = \min_{\mu} \{h(x^{(k)} - \mu s^{(k)}) \mid 0 \leq \mu \leq \sigma_k \|Dh(x^{(k)})\|\}.$$

END

Graphik 3.33 I



D.h. gesucht wird $\min \varphi(\mu) = \min h(x^{(k)} + \mu s^{(k)})$, $\mu \in [0, \sigma_k \|Dh(x^{(k)})\|]$. Suche Minimum über λ_k . Das ist jetzt eine eindimensionale Minimierung, und damit zwar einfacher als $\min h(x)$ aber immer noch sehr aufwendig. Dies wird später noch vereinfacht. Konvergenz von (3.32).

Satz 3.34 Seien $h : \mathbb{R}^n \rightarrow \mathbb{R}$ und $x^{(0)} \in \mathbb{R}^n$, so dass

- $K := \{x \mid h(x) \leq h(x^{(0)})\}$ kompakt.
- h stetig differenzierbar in einer Umgebung von K .

Dann gilt für jede Folge $\{x^{(k)}\}$ in (3.32)

- $x^{(k)} \in K$ für alle $k = 0, 1, 2, \dots$

2) $\{x^{(k)}\}$ besitzt mindestens einen Häufungspunkt \hat{x} in K .

3) Jeder Häufungspunkt \hat{x} von $\{x^{(k)}\}$ erfüllt $Dh(\hat{x}) = 0$. (Stationärer Punkt)

Beweis: Aus der Definition von $x^{(k)}$ folgt $h(x^{(0)}) \geq h(x^{(1)}) \geq \dots \implies x^{(k)} \in K$ für alle k und wegen K kompakt gibt es eine Häufungspunkt \hat{x} in K . Also folgt 1), 2).

Nun zu 3). Mit Widerspruch. Angenommen $Dh(\hat{x}) \neq 0$, und sei o.B.d.A. $\lim_{k \rightarrow \infty} x^{(k)} = \hat{x}$ (sonst nehme Teilfolge). Setze $\gamma = \inf_k \gamma_k > 0, \sigma = \inf_k \sigma_k > 0$. Nach Lemma 3.29 gibt es $U(\hat{x})$ und $\lambda > 0$, so dass

$$h(x - \mu s) \leq h(x) - \mu \frac{\gamma}{4} \|Dh(\hat{x})\|, \quad \forall x \in U(\hat{x}), s \in D(\gamma, x), 0 \leq \mu \leq \lambda. \quad (3.35)$$

Da $\lim_k x^{(k)} = \hat{x}$ und $Dh(x)$ stetig ist, gibt es k_0 , so dass für alle $k \geq k_0$ gilt

$$x^{(k)} \in U(\hat{x}), \quad \|Dh(x^{(k)})\| \geq \frac{1}{2} \|Dh(\hat{x})\|.$$

Sei nun

$$\begin{aligned} \Lambda &= \min\{\lambda, \frac{1}{2}\sigma\|Dh(\hat{x})\|\}, \\ \varepsilon &= \Lambda \frac{\gamma}{4} \|Dh(\hat{x})\| > 0 \end{aligned}$$

Wegen $\sigma_k \geq \sigma$ folgt

$$[0, \Lambda] \subseteq [0, \sigma_k \|Dh(\hat{x})\|], \quad \forall k \geq k_0,$$

also $h(x^{(k+1)}) \leq \min_{\mu} \{h(x^{(k)} - \mu s^{(k)}) \mid 0 \leq \mu \leq \Lambda\}$.

$$\stackrel{(3.35)}{\implies} h(x^{(k+1)}) \leq h(x^{(k)}) - \frac{\Lambda\gamma}{4} \|Dh(\hat{x})\| = h(x^{(k)}) - \varepsilon, \quad \forall k \leq k_0,$$

da $\Lambda \leq \lambda, x^{(k)} \in U(\hat{x}), s^{(k)} \in D(\gamma_k, x^{(k)}) \subseteq D(\gamma, x^{(k)})$.

$$\implies \lim_{k \rightarrow \infty} h(x^{(k)}) = -\infty, \text{ Widerspruch!}$$

Denn $h(x^{(k)}) \geq h(x^{(k+1)}) \geq \dots \geq h(\hat{x})$. □

Dies ist ein allgemeines Minimierungsverfahren, was sehr viele Anwendungen hat, benötigt wird noch die eindimensionale Minimierung. Diese wird in Algorithmus 3.31 durch einen endlichen Suchprozess ersetzt.

a) Wähle $s^{(k)} \in D(\gamma_k, x^{(k)})$ und definiere

$$\begin{aligned} \rho_k &:= \sigma_k \|Dh(x^{(k)})\| \\ h_k(\mu) &= h(x^{(k)} - \mu s^{(k)}) \end{aligned}$$

und bestimme die kleinste Zahl $j \geq 0$ mit

$$h_k(\rho_k 2^{-j}) \leq h_k(0) - \rho_k 2^{-j} \frac{\gamma_k}{4} \|Dh(x^{(k)})\|.$$

b) Bestimme λ_k und damit

$$\begin{aligned} x^{(k+1)} &:= x^{(k)} - \lambda_k s^{(k)}, \quad \text{so dass gilt} \\ h(x^{(k+1)}) &= \min_{0 \leq i \leq j} h_k(\rho_k 2^{-i}). \end{aligned}$$

Das gesuchte j existiert. Ist $x^{(k)}$ stationär, so ist $j = 0$, andernfalls wende Lemma 3.29 auf $\hat{x} = x^{(k)}$ an. In jedem Fall kann j und damit auch λ_k in einem endlichen Suchprozess bestimmt werden und es gilt Satz 3.34 auch für die so bestimmte Folge.

→ Übung.

3.2.1 Anwendung auf modifiziertes Newtonverfahren

Wir wenden nun zur Lösung von $f(x) = 0$ das Minimierungsverfahren 3.31 auf $h(x) = f(x)^T f(x)$ an.

Als Suchrichtung $s^{(k)}$ in Punkt $x^{(k)}$ wählen wir den normierten Newtonvektor

$$s^{(k)} = \frac{d^{(k)}}{\|d^{(k)}\|} \quad \text{mit } d^{(k)} = Df(x^{(k)})^{-1} f(x^{(k)}).$$

(Immer definiert falls $Df(x^{(k)})^{-1}$ existiert.)

Lemma 3.36 Sei x so, daß $d(x) = Df(x)^{-1} f(x)$ existiert und nicht verschwindet. Dann gilt

$$s(x) = \frac{d(x)}{\|d(x)\|} \in D(\gamma, x), \quad \text{für alle } 0 < \gamma < \bar{\gamma}(x),$$

wobei

$$\bar{\gamma}(x) = \frac{1}{\text{cond}_2(Df(x))} = \frac{1}{\|Df(x)\|_2 \|Df(x)^{-1}\|_2}.$$

Beweis: $Dh(x) = 2f^T(x)Df(x)$. Submultiplikativität von $\|\bullet\|_2$ liefert

$$\begin{aligned} \|f^T(x)Df(x)\|_2 &\leq \|Df(x)\|_2 \|f(x)\|_2 \\ \|Df(x)^{-1}f(x)\|_2 &\leq \|Df(x)^{-1}\|_2 \|f(x)\|_2 \end{aligned}$$

und daher

$$\begin{aligned} \frac{Dh(x)s}{\|Dh(x)\|} &= \frac{f(x)^T Df(x) Df(x)^{-1} f(x)}{\|Df(x)^{-1} f(x)\| \|f^T(x) Df(x)\|} \\ &\geq \frac{1}{\text{cond}_2(Df(x))} \\ &> 0. \end{aligned}$$

Für alle γ mit $0 < \gamma < \frac{1}{\text{cond}_2(Df(x))}$ gilt also $s \in D(\gamma, x)$. □

Falls $Df(x)^{-1}$ existiert folgt also, daß $Dh(x) = 0 \iff f(x) = 0$. Damit erhalten wir das modifizierte Newtonverfahren.

Algorithmus 3.37

Input: Differenzierbare Funktion $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ sowie Startwert $x^{(0)} \in \mathbb{R}^n$.

Output: Folge $x^{(k)}$, $k = 0, 1, 2, \dots$ die gegen Nullstelle von f konvergiert.

FOR $k = 0, 1, 2, \dots$ UNTIL SATISFIED

Berechne $d^{(k)} = Df(x^{(k)})^{-1}f(x^{(k)})$, (lineares Gleichungssystem mit QR)

$$\gamma_k = \frac{1}{\text{cond}_2(Df(x^{(k)}))}, \text{ fast frei aus QR}$$

Setze $h_k(\tau) = h(x^{(k)} - \tau d^{(k)}) = f(x^{(k)} - \tau d^{(k)})^T f(x^{(k)} - \tau d^{(k)})$.

Bestimme $j \geq 0$, so daß $h_k(2^{-j}) \leq h_k(0) - 2^{-j} \frac{\gamma_k}{4} \|d^{(k)}\| \|Dh(x^{(k)})\|$.

Bestimme λ_k und $x^{(k+1)} = x^{(k)} - \lambda_k d^{(k)}$, so daß gilt

$$h(x^{(k+1)}) = \min_{0 \leq i \leq j} h_k(2^{-i}).$$

END

Kosten: Pro Schritt

- Berechnung von $Df(x^{(k)})$
- Lösung von $Df(x^{(k)})d^{(k)} = f(x^{(k)})$
- eindimensionale Minimierung.

Fehleranalyse: b) mit QR klar. a),c) hängt von f ab.

Über dieses Verfahren hat man die folgende Konvergenzaussage:

Satz 3.38 Gegeben sei $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $x^{(0)} \in \mathbb{R}^n$ mit

- $K = \{x \mid h(x) \leq h(x_0)\}$, $h(x) = f^T(x)f(x)$ ist kompakt.
- f ist auf Umgebung von K stetig differenzierbar.
- $Df(x)^{-1}$ existiert für alle $x \in K$.

Dann ist die Folge $\{x^{(k)}\}$, die in Algorithmus 3.37 erzeugt wird, wohldefiniert und es gilt

- $x^{(k)} \in K$, $\forall k = 0, 1, 2, \dots$, $\{x^{(k)}\}$ besitzt mindestens Häufungspunkt $\hat{x} \in K$.
- Jeder Häufungspunkt \hat{x} von $\{x^{(k)}\}$ ist Nullstelle von f .

Beweis: hier nicht, siehe z.B. Stoer. □

→ Übung.

Die Voraussetzungen geben die Klassen von Funktionen an, auf die Alg. 3.37 anwendbar ist.

3.2.2 Praktische Realisierung des modifizierten Newton-Verfahrens

Problem bei der Berechnung ist, daß in jedem Schritt $Df(x^{(k)})$ ausgerechnet werden muß. Idee: approximiere $Df(x^{(k)})$ durch $\Delta f(x^{(k)}) = (\Delta_1 f, \dots, \Delta_n f)(x^{(k)})$ (Differenzenquotientenmatrix), wobei

$$\Delta_i f(x) = \frac{f(x_1, \dots, x_{i-1}, x_i + h_i, x_{i+1}, \dots, x_n) - f(x_1, \dots, x_n)}{h_i} = \frac{f(x + h_i e_i) - f(x)}{h_i}.$$

Schwierigkeiten: h_i zu groß \rightarrow schlechte Approximation an $Df(x)$, h_i zu klein, dann $f(x + h e_i) \approx f(x)$.

Wähle h_i so, daß $f(x)$ und $f(x + h e_i)$ ungefähr die $\frac{t}{2}$ ersten Stellen gemeinsam haben (t -stellige Rechnung)

$$|h_i| \approx \sqrt{\epsilon ps} \|f(x)\| / \|\Delta_i f(x)\|$$

! Auslöschung ist dann noch erträglich.

Im allgemeinen ist jedoch auch die Berechnung von Δf zu teuer daher verwende folgende Vereinfachung mit Satz von Broyden:

Satz 3.39 Seien $A, B \in \mathbb{R}^{n,n}$, $b \in \mathbb{R}^n$, $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$ gegeben durch $F(u) = Au + b$, $x, x' \in \mathbb{R}^n$ und $p = x - x'$, $q = F(x') - F(x) = Ap$. Dann gilt für die Rang-1-Modifikation $B' = B + \frac{1}{p^T p}(q - Bp)p^T$ die Abschätzung

$$\|B' - A\|_2 \leq \|B - A\|_2$$

und $B'p = Ap = q$.

Beweis: $(B' - A)p = Bp + (q - Bp) - Ap = 0$.

Für jedes $u \in \mathbb{R}^n$ mit $\|u\|_2 = 1$ existiert ein v mit

$$u = \alpha p + v, \quad v^T p = 0, \quad \|v\|_2 \leq 1.$$

Also folgt mit $\|u\|_2 = 1$

$$\|(B' - A)u\|_2 = \|(B' - A)v\|_2 = \|(B - A)v\|_2 \leq \|B - A\|_2 \|v\|_2 \leq \|B - A\|_2.$$

Also auch für $\|B' - A\|_2$, welches das Supremum über alle u ist. □

Also folgt, daß $DF(x) = A$ von B' genau so gut approximiert wird wie von B . B' und $DF(x)$ transformieren p in denselben Vektor. Die Idee ist nun, $f(x)$ durch eine affine Abbildung zu approximieren (in der Nähe einer Nullstelle) und diese Idee anzuwenden.

Ersetze in Algorithmus 3.37 den Teil in der Schleife.

Algorithmus 3.40 (Broyden-Rang 1-Verfahren)

Input: $f(x): \mathbb{R}^n \rightarrow \mathbb{R}^n$ differenzierbar und $x^{(0)} \in \mathbb{R}^n$

Output: Folge $x^{(k)}$, $k = 0, 1, 2, \dots$

FOR $k = 0, 1, 2, \dots$ UNTIL SATISFIED

$$d^{(k)} = B_k^{-1} f(x^{(k)})$$

$$x^{(k+1)} = x^{(k)} - \lambda_k d^{(k)}$$

$$p^{(k)} = x^{(k+1)} - x^{(k)}$$

$$q^{(k)} = f(x^{(k+1)}) - f(x^{(k)})$$

$$B^{(k+1)} = B_k + \frac{1}{p^{(k)T} p^{(k)}} (q^{(k)} - B_k p^{(k)}) p^{(k)T}$$

END

λ_k wird durch näherungsweise Minimierung von $\|f(x^{(k+1)})\|^2 \approx \min_{\lambda \geq 0} \|f(x^{(k)} - \lambda d^{(k)})\|^2$ bestimmt wie vorher.

$B_0 \approx \Delta f(x^{(0)})$. Praktische Erfahrung zeigt, daß man diese Wahl von B_{k+1} nur für $0.5 \leq \lambda_k \leq 1$ machen sollte, sonst besser $B_{k+1} = \Delta f(x_{k+1})$.

Die Lösung des linearen Gleichungssystems

$$B_k d^{(k)} = f(x^{(k)})$$

kann über Rang 1 Aufdatierungsformel sehr leicht gemacht werden.

→ Übung 3.58: Sherman–Morrison Formel

→ Zusammenhang modifiziertes Newton – Sekantenverfahren – andere Verfahren.

3.3 Nullstellenbestimmung für Polynome

Vieles vereinfacht sich, falls $f(x)$ skalares Polynom ist. Dies ist praktisch fast nicht mehr relevant außer in speziellen Fällen.

Gegeben: Polynom

$$p(x) = a_0 x^n + a_1 x^{n-1} + \dots + a_{n-1} x + a_n.$$

Gesucht: reelle (alle) Nullstellen.

Newtonverfahren:

$$x^{(k+1)} = x^{(k)} - \frac{p(x^{(k)})}{p'(x^{(k)})}. \quad (3.41)$$

Probleme:

1. Preiswerte Auswertung von p, p' .
2. Wahl der Startwerte.
3. Deflation von Nullstellen.
4. Komplexe Nullstellen.

Zu Punkt 1: Horner Schema

Auswertung von $p(z)$ durch Algorithmus

$$\begin{aligned} b_0 &= a_0 \\ b_k &= b_{k-1}z + a_k, \quad k = 1, 2, \dots, n \end{aligned} \quad (3.42)$$

$$p(z) = ((\dots((a_0 z + a_1)z + a_2) \dots + a_{n-2})z + a_{n-1})z + a_n$$

kostet nur $2n$ flops und es gilt der folgende Satz:

Satz 3.43 Sei $p(x) = a_0 x^n + \dots + a_n$ und $p_1(x) = b_0 x^{n-1} + \dots + b_{n-1}$, b_i wie in (3.42) für $p(z)$. Dann gilt:

$$p(x) = (x - z)p_1(x) + p(z), \quad (3.44)$$

$$p'(z) = p_1(z). \quad (3.45)$$

Beweis:

$$\begin{aligned}
 b_n &= p(z) \\
 \implies (x-z)p_1(x) + p(z) &= xp_1(x) - zp_1(x) + p(z) \\
 &= b_0x^n + b_1x^{n-1} + \dots + b_{n-1}x + b_n - (zb_0x^{n-1} + \dots + zb_{n-1}) \\
 &= a_0x^n + \dots + a_n \\
 &= p(x). \\
 \implies p'(x) &= p_1(x) + (x-z)p'(z)
 \end{aligned}$$

Einsetzen von z gibt (3.45). □

Durch 2-faches Anwenden des Hornerchemas kann man also $p(z), p'(z)$ mit $4n - 2$ flops bestimmen.

Zu Punkt 2: Wahl der Startwerte

Satz 3.46 Für alle Nullstellen ξ_i eines Polynoms $p(x) = a_0x^n \dots + a_n$ gilt

$$\begin{aligned}
 a) \quad |\xi_i| &\leq \max \left\{ \left| \frac{a_n}{a_0} \right|, 1 + \left| \frac{a_{n-1}}{a_0} \right|, \dots, 1 + \left| \frac{a_1}{a_0} \right| \right\} \\
 b) \quad |\xi_i| &\leq \max \left\{ 1, \sum_{i=1}^n \left| \frac{a_i}{a_0} \right| \right\} \\
 c) \quad |\xi_i| &\leq \max \left\{ \left| \frac{a_n}{a_{n-1}} \right|, 2 \left| \frac{a_{n-1}}{a_{n-2}} \right|, \dots, 2 \left| \frac{a_1}{a_0} \right| \right\} \\
 d) \quad |\xi_i| &\leq \sum_{i=0}^{n-1} \left| \frac{a_{i+1}}{a_i} \right| \\
 e) \quad |\xi_i| &\leq 2 \max \left\{ \left| \frac{a_1}{a_0} \right|, \sqrt{\left| \frac{a_2}{a_0} \right|}, \sqrt[3]{\left| \frac{a_3}{a_0} \right|}, \dots, \sqrt[n]{\left| \frac{a_n}{a_0} \right|} \right\}
 \end{aligned}$$

Weitere Sätze in *Geometry of Polynomials* von M. Marden. Damit können Kreise bestimmt werden, in denen die Nullstellen liegen.

→ Übung 6.50,3.51: Gerschgorin-Satz angewendet auf Begleitmatrix.

Zu Punkt 3: Deflation

Falls man Nullstelle \hat{x} bestimmt hat, könnte man Polynom p_1 bestimmen, $p_1(x - \hat{x}) = p(x)$ bilden und weiterrechnen.

Im allgemeinen numerisch sehr schlecht, besser mit rationaler Funktion $p_1(x) = \frac{p(x)}{x - \hat{x}}$ weiterrechnen.

→ Übung 3.56

Zu Punkt 4: Komplexe Nullstellen, Bairstow-Verfahren

Für reelle Polynome will man komplexe Nullstellen mit reeller Rechnung bestimmen. Mit $z_j \in \mathbb{C} \setminus \mathbb{R}$ ist dann auch \bar{z}_j Nullstelle und $(x - z_j)(x - \bar{z}_j) = x^2 - rx + s$ reeller Faktor von $p(x)$.

$$\operatorname{Re}(z_j) = -\frac{r}{2}, \quad \operatorname{Im}(z_j) = \sqrt{s - \frac{r^2}{4}}$$

Wende Polynomdivision mit Newtonverfahren auf $p(x)/(x^2 + ux + v)$ an.

$$\frac{p(x)}{x^2 + ux + v} = Q(x) + \frac{Ax + B}{x^2 + ux + v} \tag{3.47}$$

Die so definierten A, B sind eindeutig. A, B können mit dem Horner-Schema berechnet werden.
 → Übung 3.57

Wende nun 2-dimensionales Newton-Verfahren an, um r, s so zu bestimmen, daß

$$\begin{bmatrix} A(r, s) \\ B(r, s) \end{bmatrix} = 0 \quad (3.48)$$

Dies ist das Bairstow-Verfahren.

→ Übung 3.57

→ Sekantenverfahren für Polynome.

Weitere Verfahren wie Bisektion und Sturmsche Ketten später bei Eigenwertproblemen.

3.4 Übungen zu Kapitel 3

Übung 3.49 Sei $\Phi : K \rightarrow \mathbb{R}^n$, $K \subseteq \mathbb{R}^n$ offen, konvex, Φ differenzierbar und nehme an, daß das totale Differential $D\Phi$ auf K beschränkt ist. Zeige:

$$\|\Phi(x) - \Phi(y)\|_\infty \leq \|D\Phi\|_\infty|_K \|x - y\|_\infty, \quad \forall x, y \in K.$$

Übung 3.50 Für welche der folgenden Gleichungen und welche Startwerte konvergiert die zugehörige Fixpunktiteration?

1. $x = e^x - \sin x + x$
2. $x = \sin x - e^x + x$
3. $x = \arcsin e^x, x < 0$
4. $x = \ln \sin x, x \in (0, \pi)$

Übung 3.51 Sei $p(x) = a_n x^n + \dots + a_1 x + a_0$, $a_n \neq 0$, Polynom n -ten Grades. Zeige:

1. Die Nullstellen von p sind die Eigenwerte der Begleitmatrix

$$A = \begin{pmatrix} 0 & 1 & & \\ & \ddots & \ddots & \\ & & 0 & 1 \\ -\frac{a_0}{a_n} & -\frac{a_1}{a_n} & \dots & -\frac{a_{n-1}}{a_n} \end{pmatrix}$$

2. Die Nullstellen ξ_1, \dots, ξ_n von p erfüllen für alle $i = 1, \dots, n$ die Ungleichungen

$$|\xi_i| \leq \max\left\{1, \sum_{k=0}^{n-1} \left|\frac{a_k}{a_n}\right|\right\}, \quad |\xi_i| \leq \max\left\{\left|\frac{a_0}{a_n}\right|, 1 + \left|\frac{a_1}{a_n}\right|, \dots, 1 + \left|\frac{a_{n-1}}{a_n}\right|\right\},$$

(Hinweis: Für jeden Eigenwert λ von A und jede mit der Matrixmultiplikation verträgliche Norm gilt: $|\lambda| \leq \|A\|$)

Übung 3.52 Sei $A \in \mathbb{R}^{n,n}$ und

$$f(x_1, \dots, x_n, x_{n+1}) = \begin{pmatrix} (A - x_{n+1}I) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \\ \frac{1}{2}(\sum_{i=1}^n |x_i|^2 - 1) \end{pmatrix}.$$

Zeige:

1. (x_1, \dots, x_{n+1}) ist Nullstelle von f , genau dann wenn x_{n+1} Eigenwert und $(x_1, \dots, x_n)^T$ normierter Eigenvektor von A sind.
2. Sei (x_1, \dots, x_{n+1}) Nullstelle von f . Dann ist die Ableitungsmatrix $Df(x_1, \dots, x_{n+1})$ nicht singulär, genau dann wenn x_{n+1} ein einfacher Eigenwert von A ist (Tip: Jordan-Form). Was bedeutet das für das Newton-Verfahren?

Übung 3.53 Schreibe ein **MATLAB**-Programm für das Newtonverfahren zur Bestimmung einer Nullstelle einer Funktion. Wende das Verfahren auf die Funktion $e^x - \sin x$ an und vergleiche das Programm bezüglich Rechenaufwand und Genauigkeit mit der Fixpunktiteration angewandt auf die Funktionen aus Aufgabe 3.50

Übung 3.54 Sei $p(x) = a_0x^n + a_1x^{n-1} + \dots + a_n$. Bestimme den Rechenaufwand für die Berechnung von $p(x)$, falls man gerade durchrechnet, d.h. ohne Horner-Schema.

Übung 3.55 Sei p Polynom n -ten Grades. Zeige:

1. Zu gegebenem $x \in \mathbb{R}$ lässt sich $p(x)$ mit dem Hornerschema in $2n$ flops berechnen.
2. Zu gegebenem $x \in \mathbb{R}$ lassen sich $p(x), p'(x)$ mit einem angepassten Hornerschema in $4n - 2$ flops berechnen, d.h. es werden gegenüber der Berechnung von $p(x)$ zusätzlich $2n - 2$ flops für die Berechnung von $p'(x)$ benötigt (Tip: Satz 3.43 der Vorlesung).

Übung 3.56 Schreibe ein **MATLAB**-Programm für das Newtonverfahren zur Bestimmung der reellen Nullstellen einer rationalen Funktion $f(x) = p(x)/q(x)$, wobei p, q Polynome seien. Vergleiche das Programm bezüglich Rechenaufwand und Genauigkeit mit der Funktion 'roots' aus **MATLAB** anhand des Beispiels $p_n(x) = \prod_{i=1}^n (x - i)$, $n = 2, 4, 6, \dots, 20$, $q(x) = 1$.

Übung 3.57 Sei $p(x) = a_0x^n + a_1x^{n-1} + \dots + a_n$ Polynom n -ten Grades mit $n \geq 3$. Zeige:

1. Bestimme in Abhängigkeit von gegebenen $r, q \in \mathbb{R}$ nacheinander die Koeffizienten b_0, \dots, b_n , so dass

$$p(x) = q(x)(x^2 - rx - q) + b_{n-1}(x - r) + b_n.$$

Dabei ist $q(x) = b_0x^{n-2} + b_1x^{n-3} + \dots + b_{n-2}$. Was stellst Du fest?

2. Setze $c_{j-1} = \frac{\partial b_j}{\partial r}$, $d_{j-2} = \frac{\partial b_j}{\partial q}$, $j = 0, \dots, n$. Dann gilt für die partiellen Ableitungen der Koeffizienten b_0, \dots, b_n :

$$c_{-1} = 0, \quad c_0 = b_0, \quad c_{j-1} = b_{j-1} + rc_{j-2} + qc_{j-3}, \quad j = 2, \dots, n$$

sowie

$$d_j = c_j, \quad j = -1, \dots, n.$$

Was kostet die Berechnung der partiellen Ableitungen von b_{n-1}, b_n ?

3. Das 2-dimensionale Newton-Verfahren zur Bestimmung der Nullstellen von $\Phi(r, q) = (b_{n-1}, b_n)$ lautet

$$\begin{pmatrix} r_{k+1} \\ q_{k+1} \end{pmatrix} = \begin{pmatrix} r_k \\ q_k \end{pmatrix} - \begin{pmatrix} c_{n-2} & c_{n-3} \\ c_{n-1} & c_{n-2} \end{pmatrix}^{-1} \begin{pmatrix} b_{n-1} \\ b_n \end{pmatrix}.$$

Übung 3.58 1. Sei $A \in \mathbb{R}^{n,n}$, $A = S - fg^T$ mit $f, g \in \mathbb{R}^n$ und $S \in \mathbb{R}^{n,n}$ nicht singulär. Setze $S_c = 1 - g^T S^{-1} f$. Dann ist A nicht singulär genau dann, wenn $S_c \neq 0$ ist und es gilt (Sherman-Morrison-Formel):

$$A^{-1} = S^{-1} + S^{-1} f \frac{1}{S_c} g^T S^{-1}$$

Tip: multipliziere die rechte Seite an A heran.

2. Was kostet das Lösen von m linearen Gleichungssystemen mit A verglichen mit dem Lösen m linearer Gleichungssysteme mit S ?
3. Lässt sich die Formel auf Blöcke verallgemeinern? D.h. $F, G \in \mathbb{R}^{n,r}$.

Übung 3.59 Ist $A(x)$ stetige nichtsinguläre Matrixwertige Funktion, dann gilt das auch für ihre Inverse. Tip: Cramersche Regel.

Übung 3.60 Schreibe ein **MATLAB**-Algorithmus zur Lösung einer nichtlinearen Gleichung $\Phi(x) = 0$ mit Hilfe des Newton-Verfahrens bei gegebenen Startwert $x^{(0)}$ und Ableitungsmatrix $D\Phi(x)$. Teste das Verfahren für die Funktion $\Phi(x) = \phi(x) - x$ mit $\phi(x)$ aus Aufgabe 3.60 sowie an folgender Funktion Ψ für $n = 10$ mit Startwert $(1, \dots, 1, 0)^\top$.

$$\Psi(x_1, \dots, x_n, x_{n+1}) = \begin{pmatrix} (A - x_{n+1}I) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \\ \frac{1}{2}(1 - \sum_{i=1}^n |x_i|^2) \end{pmatrix}, \quad A = \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 2 \end{pmatrix}.$$

Was stellst du fest?

Übung 3.61 Schreibe ein **MATLAB**-Algorithmus zur Lösung einer nichtlinearen Gleichung $\Phi(x) = 0$ mit Hilfe des Newton-Verfahrens bei gegebenen Startwert $x^{(0)}$ und Ableitungsmatrix $D\Phi(x)$. Teste das Verfahren für die Funktion $\Phi(x) = \phi(x) - x$ mit $\phi(x)$ aus Aufgabe 3.60 sowie an folgender Funktion Ψ für $n = 10$ mit Startwert $(1, \dots, 1, 0)^\top$.

$$\Psi(x_1, \dots, x_n, x_{n+1}) = \begin{pmatrix} (A - x_{n+1}I) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \\ \frac{1}{2}(1 - \sum_{i=1}^n |x_i|^2) \end{pmatrix}, \quad A = \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 2 \end{pmatrix}.$$

Was stellst du fest?

Übung 3.62 Schreibe ein **MATLAB**-Algorithmus zur Lösung einer nichtlinearen Gleichung $\Phi(x) = 0$ mit Hilfe des MODIFIZIERTEN Newton-Verfahrens bei gegebenen Startwert $x^{(0)}$ und Ableitungsmatrix $D\Phi(x)$. Teste das Verfahren für die Funktion $\Phi(x) = \phi(x) - x$ mit $\phi(x)$ aus Aufgabe 3.60 sowie an folgender Funktion Ψ für $n = 10$ mit Startwert $(1, \dots, 1, 0)^\top$.

$$\Psi(x_1, \dots, x_n, x_{n+1}) = \begin{pmatrix} (A - x_{n+1}I) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \\ \frac{1}{2}(1 - \sum_{i=1}^n |x_i|^2) \end{pmatrix}, \quad A = \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & \ddots & \\ & \ddots & \ddots & -1 \\ & & -1 & 2 \end{pmatrix}.$$

Was stellst du fest?

Übung 3.63 Es sei

$$f(x, y) = \begin{pmatrix} \frac{3x}{1+x} + \frac{1}{y} \\ \frac{3y}{1+y} + \frac{2}{x} \end{pmatrix}.$$

1. Man finde eine Teilmenge $I \subseteq \mathbb{R}^2$, in der f bzgl. der $\|\bullet\|_\infty$ -Norm den Voraussetzungen des Banachschen Fixpunktsatzes genügt.

2. Man führe mit $(x_0, y_0) = (2, 2)$ zwei Iterationen durch und gebe eine a-posteriori Fehlerschranke für (x_2, y_2) an.
3. Wieviele Iterationsschritte sind hinreichend, um ausgehend von $(x_0, y_0) = (2, 2)$ den Fixpunkt mit einer Genauigkeit von 10^{-4} zu berechnen

Kapitel 4

Interpolation

Ein Interpolationsproblem ist eine Aufgabe: Finde a_0, \dots, a_n so daß

$$P(x_i, a_0, \dots, a_n) = f_i \quad i = 0, \dots, n \quad (4.1)$$

für gegebene $(x_i, f_i) = (x_i, f(x_i))$, $i = 0, \dots, n$ und eine vorgebene Funktion P . (x_i, f_i) heißen Stützpunkte/Knoten.

Wichtigste Beispiele:

Polynominterpolation

$$P = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

Trigonometrische Interpolation

$$P = a_0 + a_1e^{xi} + a_2e^{2xi} + \dots + a_n e^{nxi}$$

Rationale Interpolation

$$P = \frac{a_0 + a_1x + \dots + a_px^p}{\beta_0 + \beta_1x + \dots + \beta_mx^m} \quad \text{hier nicht .}$$

Splines

stückweise Polynome.

Anwendungen:

Approximation von Funktionen, Quadratur, Lösung von gewöhnlichen oder Partiellen Differentialgleichungen, CAD, numerische Differentiation, Bildverarbeitung, Signalverarbeitung, ...

Beispiel aus alten Zeiten:

Logarithmentafel:

In Tafel:	$\log_{10} 1.0$
	$\log_{10} 1.1$
gesucht:	$\log 1.05$

Lösung lineare Interpolation: Berechne Gerade und werte diese bei 1.05 aus.

4.1 Polynominterpolation

Notation $\Pi_n = \{ \text{reelle oder komplexe Polynome vom Grad } \leq n \}$.

Satz 4.2 (Lagrange Interpolation) Seien $n + 1$ Stützpunkte (x_i, f_i) , $i = 0, \dots, n$ beliebig mit $x_i \neq x_k$ für $i \neq k$. Dann gibt es ein eindeutiges Polynom $P \in \Pi_n$, mit

$$P(x_i) = f_i \quad i = 0, \dots, n. \quad (4.3)$$

Beweis: Wir zeigen zuerst die Eindeutigkeit. Angenommen es gäbe $P_1, P_2 \in \Pi_n$

$$P_1(x_i) = P_2(x_i) = f_i \quad i = 0, \dots, n.$$

Sei $P = P_1 - P_2 \in \Pi_n$. P hat grad $\leq n$ und mindestens $n + 1$ verschiedene Nullstellen $x_i, i = 0, \dots, n$. $\implies P \equiv 0 \implies P_1 = P_2$.

Die Existenz folgt per Konstruktion. Sei

$$L_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x - x_j)}{(x_i - x_j)}. \quad (4.4)$$

Es gilt wegen a), daß L_i die eindeutigen Lösungen der Aufgabe

$$L_i(x_k) = \delta_{ik} = \begin{cases} 1 & i = k \\ 0 & i \neq k \end{cases}, k = 0, \dots, n \quad (4.5)$$

sind. Dann folgt

$$P(x) = \sum_{i=0}^n f_i L_i(x) = \sum_{i=0}^n f_i \prod_{\substack{j=0 \\ j \neq i}}^n \left(\frac{x - x_j}{x_i - x_j} \right) \quad (4.6)$$

löst die Interpolationsaufgabe. □

$L_i(x)$ in 4.4 heißen Lagrange-Interpolations-Polynome. Sie bilden Basis für Π_n . Andere Basis x^0, x^1, \dots, x^n .

Für die praktische Berechnung sind die $L_i(x)$ nicht geeignet, genauso wie die Standardbasis. Für numerische Zwecke besser Neville-Aitken-Methode oder Newton'sche Interpolationsformel. Dazu führe zuerst dividierte Differenzen ein.

Definition 4.7 Seien (x_i, f_i) , $i = 0, \dots, n$ gegeben. Der Ausdruck

$$f[x_i, x_{i+1}, \dots, x_{i+k}],$$

welcher durch die Rekursionen

$$f[x_i] = f_i \quad (4.8)$$

$$f[x_i, x_{i+1}, \dots, x_{i+k}] := \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+k}] - f[x_i, x_{i+1}, \dots, x_{i+k-1}]}{x_{i+k} - x_i}, \quad (4.9)$$

$i = 0, \dots, n$, $k = 0, \dots, n - i$ gegeben ist, heißt k -te dividierte Differenz von $(x_i, f_i), \dots, (x_{i+k}, f_{i+k})$. Es gibt ein einfaches Schema zur Berechnung der dividierten Differenzen.

Tabelle 4.10 (Schema der dividierten Differenzen)

	$k = 0$	$k = 1$	$k = 2$	$k = 3$
x_0	$f_0 = f[x_0]$			
x_1	$f_1 = f[x_1]$	$f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0}$		
x_2	$f_2 = f[x_2]$	$f[x_1, x_2] = \frac{f[x_2] - f[x_1]}{x_2 - x_1}$	$f[x_0, x_1, x_2]$	
x_3	$f_3 = f[x_3]$	$f[x_2, x_3] = \frac{f[x_3] - f[x_2]}{x_3 - x_2}$	$f[x_1, x_2, x_3]$	$f[x_0, \dots, x_3]$
x_4	\vdots	\vdots	\vdots	
\vdots				

Einfach rekursiv zu berechnen.

→ Übung 4.76

Bei der Newtoninterpolation verwendet man eine andere Basis: $N_0(x) = 1$

$$N_i(x) = \prod_{j=0}^{i-1} (x - x_j) = (x - x_0)(x - x_1) \cdots (x - x_{i-1}) \quad i = 1, \dots, n \quad (4.11)$$

Das Interpolationspolynom ist dann

$$P(x) = \sum_{i=0}^n a_i N_i(x). \quad (4.12)$$

Die Koeffizienten a_i können aus den Formeln

$$\begin{aligned} f_0 &= P(x_0) = a_0 \\ f_1 &= P(x_1) = a_0 + a_1(x_1 - x_0) \\ f_2 &= P(x_2) = a_0 + a_1(x_2 - x_0) + a_2(x_2 - x_0)(x_2 - x_1) \\ &\vdots \end{aligned} \quad (4.13)$$

bestimmt werden.

Es gilt sofort, daß

$$P_{0\dots k}(x) = \sum_{j=0}^k a_j N_j(x) \quad (4.14)$$

die Interpolationsaufgabe für $(x_i, f_i) \quad i = 0, \dots, k$ erfüllt, für $k = 0, \dots, n$.

Sei $P_{i,i+1,\dots,i+k}$ das Polynom das die Interpolationsaufgabe für $(x_j, f_j), j = i, i+1, \dots, i+k$ erfüllt.

Dann gilt

Satz 4.15

$$\begin{aligned} P_{i,i+1,\dots,i+k}(x) &= f[x_i] + f[x_i, x_{i+1}](x - x_i) + \cdots \\ &\quad + f[x_i, x_{i+1}, \dots, x_{i+k}](x - x_i) \cdots (x - x_{i+k-1}) \end{aligned}$$

Beweis: [Vollständige Induktion] Richtig für $k = 0$. Induktionsvor.: Richtig für $k - 1 \geq 0$.

Nun gilt

$$P_{i,i+1,\dots,i+k}(x) = P_{i,i+1,\dots,i+k-1}(x) + a(x - x_i) \cdots (x - x_{i+k-1}) \quad (4.16)$$

per Konstruktion, wobei a gerade der Koeffizient von x^k in $P_{i,i+1,\dots,i+k}(x)$ ist. Zu zeigen ist, daß

$$a = f[x_i, \dots, x_{i+k}] \quad (4.17)$$

Nach IV sind die höchsten Koeffizienten von $P_{i,\dots,i+k-1}$ und $P_{i+1,\dots,i+k}$ gerade $f[x_i, \dots, x_{i+k-1}]$ und $f[x_{i+1}, \dots, x_{i+k}]$, also

$$\begin{aligned} P_{i,\dots,i+k-1}(x) &= \cdots + f[x_i, \dots, x_{i+k-1}]x^{k-1} \\ P_{i+1,\dots,i+k}(x) &= \cdots + f[x_{i+1}, \dots, x_{i+k}]x^{k-1}. \end{aligned}$$

Nun erfüllen die $P_{i,\dots,i+k}$ gerade die Rekursionsformel (Neville-Formel)

$$P_{i,\dots,i+k}(x) = \frac{(x - x_{i+k})P_{i,\dots,i+k-1}(x) - (x - x_i)P_{i+1,\dots,i+k}(x)}{x_{i+k} - x_i}, \quad (4.18)$$

wie man sofort durch Einsetzen und Satz 4.2 erhält. Also folgt, daß der Koeffizient von x^k

$$\frac{f[x_{i+1}, \dots, x_{i+k}] - f[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i} = f[x_i, \dots, x_{i+k}] \quad (4.19)$$

nach (4.9) □

Also sind die Werte in der ersten Zeile des Dividierten Differenzen Schemas gerade die Koeffizienten des Newton-Polynoms.

Beispiel 4.20 $\frac{x_i = \begin{array}{|c|c|c|} \hline 0 & 1 & 3 \\ \hline \end{array}}{f_i = \begin{array}{|c|c|c|} \hline 1 & 3 & 2 \\ \hline \end{array}}$

Dividierte Differenzen Schema

$$\begin{array}{l|l} x_0 = 0 & f[x_0] = 1 \\ x_1 = 1 & f[x_1] = 3 \\ x_2 = 3 & f[x_2] = 2 \end{array} \begin{array}{l} \searrow \\ \swarrow \\ \nearrow \end{array} \begin{array}{l} f[x_0, x_1] = 2 \\ f[x_1, x_2] = -\frac{1}{2} \end{array} \begin{array}{l} \searrow \\ \nearrow \end{array} f[x_0, x_1, x_2] = -\frac{5}{6}$$

$$P_{0,1,2}(x) = 1 + 2(x - 0) - \frac{5}{6}(x - 0)(x - 1)$$

Auswertung von $P_{0,\dots,n}(x)$ mit Hornerartigem Schema.

Vorteile der dividierten Differenzen und des Newton-Schemas.

- Weitere Stützstellen möglich.
- Man erhält Interpolationspolynom, welches mit Hornerartigem Schema leicht zu berechnen ist.
- $f[x_0, \dots, x_k]$ symmetrisch in x_i , d.h. auf Anordnung kommt es nicht an. Siehe *Numerische Mathematik*, Stoer.

- Rundungsfehlereinfluß läßt sich durch Umordnung reduzieren.

- $f[x_0, \dots, x_n] = \frac{f^{(n)}(\xi)}{n!}$, für ein $\xi \in \mathcal{C}[x_0, \dots, x_n]$.

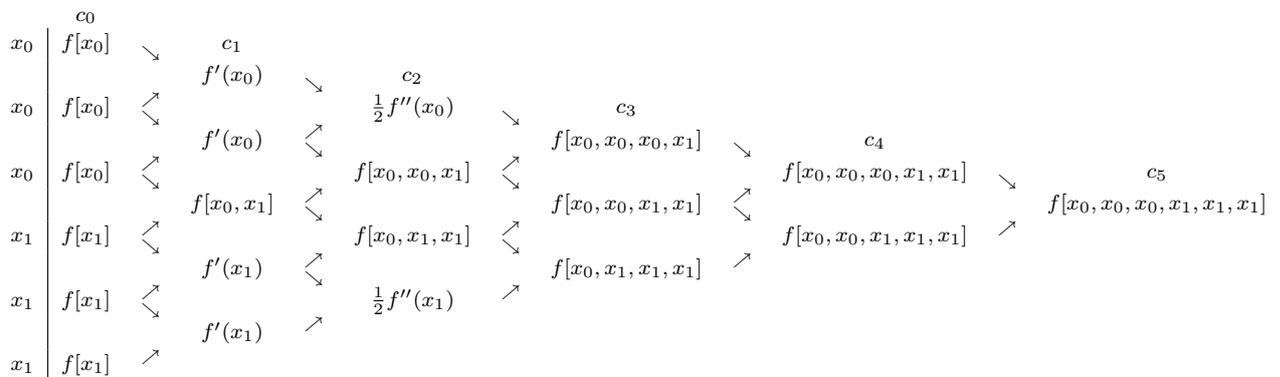
→ Übung 4.76, 4.77: einfache Interpolationsaufgabe, mehr zu Dividierten Differenzen.

Dividierte Differenzen lassen sich auf zusammenfallende Stützstellen erweitern falls $f(x_j) = f(x_i)$.

$$\begin{aligned}
 f[x_k, x_k] &:= \lim_{h \rightarrow 0} \frac{f[x_k + h] - f[x_k]}{x_k + h - x_k} = f'(x_k) \\
 f[x_k, x_k, x_k] &:= \lim_{h \rightarrow 0} \frac{f[x_k + h, x_k + 2h] - f[x_k, x_k + h]}{x_k - 2h - x_k} \\
 &= \lim_{h \rightarrow 0} \frac{f[x_k + 2h] - 2f[x_k + h] + f[x_k]}{2h^2} \\
 &= \frac{1}{2} f''(x_k) \\
 \underbrace{f[x_k, \dots, x_k]}_{m+1 \text{ mal}} &:= \frac{1}{m!} f^{(m)}(x_k).
 \end{aligned}$$

Damit kann man dann allgemeine Hermite Interpolation durchführen, bei der für mehrfache Stützstellen jeweils die entsprechenden Ableitungen mit interpoliert werden.

Beispiel 4.21



Das Interpolationspolynom ist damit

$$\begin{aligned}
 P_{000111}(x) &= c_0 + c_1(x - x_0) + c_2(x - x_0)^2 + c_3(x - x_0)^3 \\
 &\quad + c_4(x - x_0)^3(x - x_1) + c_5(x - x_0)^3(x - x_0)^2.
 \end{aligned}$$

→ Übung

Beispiel 4.22 $x_0, x_0, x_1, x_1, \dots$ Schema ergibt

$$\begin{aligned}
 P_{001122\dots kk}(x) &= [c_0 + c_1(x - x_0)] + [c_2 + c_3(x - x_1)](x - x_0)^2 \\
 &\quad + [c_4 + c_5(x - x_2)](x - x_0)^2(x - x_1)^2 + \dots \\
 &\quad + [c_{2k} + c_{2k-1}(x - x_k)](x - x_0)^2 \dots (x - x_{k-1})^2
 \end{aligned}$$

Es gilt die Bedingung

$$\left. \begin{aligned} P_{00\dots kk}(x_i) &= f(x_i) \\ P'_{00\dots kk}(x_i) &= f'(x_i) \end{aligned} \right\} i = 0, \dots, k$$

→ Übung

Fehlerabschätzung: Angenommen $f_i = f(x_i) \quad i = 0, \dots, n$.

Wie gut ist die Approximation?

Satz 4.23 Es sei f $n + 1$ -mal differenzierbar, so gibt es zu jedem \hat{x} eine Zahl ξ aus dem kleinsten Intervall I , daß alle x_i und \hat{x} enthält (konvexe Hülle $\mathcal{C}[x_i, \hat{x}]$), so daß

$$f(\hat{x}) - P_{01\dots n}(\hat{x}) = w(\hat{x}) \frac{f^{(n+1)}(\xi)}{(n+1)!} \quad (4.24)$$

wobei $w(x) = \prod_{i=0}^n (x - x_i)$.

Beweis: Sei für $\hat{x} \neq x_i$

$$F(x) := f(x) - P(x) - kw(x).$$

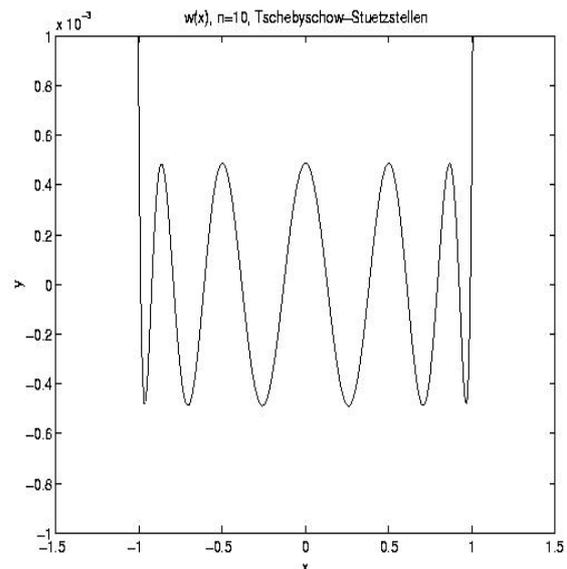
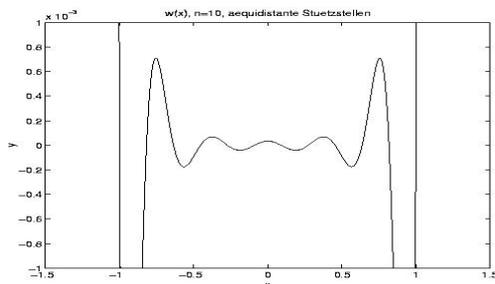
Bestimme k derart, daß $F(\hat{x}) = 0$. Dann hat $F(x)$ in $\mathcal{C}[x_0, \dots, x_n, \hat{x}]$ $n + 2$ Nullstellen x_0, \dots, x_n, \hat{x} .

Nach dem Satz von Rolle hat dann $F'(x)$ mindestens $n + 1$ Nullstellen, ..., $F^{(n+1)}(x)$ mindestens eine Nullstelle $\xi \in \mathcal{C}[x_0, \dots, x_n, \hat{x}]$.

$$\begin{aligned} F^{(n+1)}(x) \equiv 0 &\implies F^{(n+1)}(\xi) = f^{(n+1)}(\xi) - k(n+1)! = 0 \\ &\implies k = \frac{f^{(n+1)}(\xi)}{(n+1)!}. \end{aligned}$$

□

Wie sieht $w(x)$ aus?



Wächst außerhalb von $\mathcal{C}[x_0, \dots, x_n]$ stark an \implies Verwendung des Interpolationspolynoms außerhalb von $\mathcal{C}[x_0, \dots, x_n]$ vermeiden.

Sei nun $f : [a, b] \rightarrow \mathbb{R}$ gegeben, f genügend glatt. Dann kann ich f beliebig genau durch Polynome interpolieren. Man könnte nun vermuten, daß die Interpolationspolynome gegen f konvergieren, falls ich die Intervallunterteilung feiner und feiner mache.

Betrachte Folge von Unterteilungen

$$\begin{aligned}\Delta_m &= \{a = x_0^{(m)} < x_1^{(m)} < \dots < x_{n_m}^{(m)} = b\} \\ \|\Delta_m\| &= \max_i |x_{i+1}^{(m)} - x_i^{(m)}|\end{aligned}$$

Satz 4.25 [Satz von Faber] Zu jeder Folge von Intervalleinteilungen Δ_m von $[a, b]$ kann man eine auf $[a, b]$ stetige Funktion f finden, so daß die Polynome P_{Δ_m} die auf $(x_i^{(m)}, f(x_i^{(m)}))$, $i = 0, \dots, n_m$ interpolieren für $m \rightarrow \infty$ NICHT gleichmäßig gegen $f(x)$ konvergieren.

Nur für ganze Funktionen gilt für alle $\|\Delta_m\| \rightarrow 0$, $P_m \rightarrow f$ gleichmäßig.

Für beliebige stetige Funktionen gibt es spezielle Δ_m , i.a. nicht äquidistant.

→ Übung 4.82: Tschebyschowknoten für $f(x) = \frac{1}{1+25x^2}$.

→ Übung: Orthogonale Polynome für Approximation.

Rationale Interpolation, hier nicht, Übung?

4.2 Trigonometrische Interpolation

Interpolation durch trigonometrische Polynome

Betrachte Intervall $[0, 2\pi]$, Stützpunkte (x_k, f_k) , $k = 0, \dots, n-1$.

$$\begin{aligned}x_k &= k \frac{2\pi}{n} \quad f_k \in \mathbb{C}, n \text{ fest} . \\ P(x) &= \beta_0 + \beta_1 e^{ix} + \beta_2 e^{2ix} + \dots + \beta_{n-1} e^{(n-1)ix}\end{aligned}$$

soll die Interpolationsbedingungen

$$p(x_k) = f(x_k), \quad k = 0, \dots, n-1 \tag{4.26}$$

erfüllen. Periodizität \implies

$$p(x_k) = f(x_k), \quad k = \dots, -2, -1, \dots \tag{4.27}$$

Setze

$$w = e^{ix}, \quad w_k = e^{ix_k} \tag{4.28}$$

Dann ist die Interpolationsaufgabe identisch mit der komplexen Interpolationsaufgabe

$$p(w_k) = f_k \tag{4.29}$$

mit

$$p = \beta_0 + \beta_1 w + \dots + \beta_{n-1} w^{n-1}$$

Also folgt aus Satz 4.2 sofort

Satz 4.30 Sei $x_k = \frac{2\pi k}{n}$, f_k beliebig für $k = 0, \dots, n-1$, $n \in \mathbb{N}$. Dann gibt es ein eindeutiges trigonometrisches Polynom $p(w) = \beta_0 + \beta_1 w + \dots + \beta_{n-1} w^{n-1}$ mit $p(w_k) = f_k$, für $w_k = e^{ix_k}$, $k = 0, \dots, n-1$.

Zusätzlich gilt

$$\begin{aligned} i) \quad & w_k^j = w_j^k \quad \forall j, k \in \mathbb{Z} \\ ii) \quad & \sum_{k=0}^{n-1} w_k^j w_k^{-l} = \begin{cases} n & \text{für } j = l \\ 0 & \text{für } j \neq l, 0 \leq j, l \leq n-1. \end{cases} \end{aligned}$$

Beweis: i) trivial, ii) für alle $k \in \mathbb{Z}$ ist w_k eine Wurzel von

$$w^n - 1 = (w - 1)(w^{n-1} + w^{n-2} + \dots + 1) = 0$$

für $k \neq 0, \pm n, \pm 2n$, ist $w_k \neq 1$, also folgt

$$\sum_{l=0}^{n-1} w_k^l = \begin{cases} n & k = 0, \pm n, \pm 2n, \dots \\ 0 & \text{sonst.} \end{cases}$$

Die n -Vektoren $\Phi_j = (w_0^j, \dots, w_{n-1}^j)$, $j = 0, \dots, n-1$ bilden eine Orthonormalbasis von \mathbb{C}^n unter dem Skalarprodukt

$$[f, g] := \frac{1}{n} \sum_{j=0}^{n-1} f_j \bar{g}_j \quad (4.31)$$

(diskrete Version von $(f, g) = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \bar{g}(x) dx$). Denn es gilt wegen (4.27)

$$[\Phi_j, \Phi_l] = \begin{cases} 1 & j = l \\ 0 & j \neq l \quad 0 \leq j, l \leq n-1. \end{cases}$$

Damit erhält man eine explizite Formel für die Koeffizienten β_i . □

Satz 4.32 Sei $p(x) = \sum_{k=0}^{n-1} \beta_k e^{ikx}$ das interpolierende trigonometrische Polynom für $p(x_k) = f_k$, $k = 0, \dots, n-1$, so gilt

$$\beta_j = \frac{1}{n} \sum_{k=0}^{n-1} f_k w_k^{-j} = \frac{1}{n} \sum_{k=0}^{n-1} f_k \exp \left\{ \frac{-2\pi i k j}{n} \right\}, \quad j = 0, \dots, n-1. \quad (4.33)$$

Beweis: Aus dem Skalarprodukt $[\bullet, \bullet]$ folgt

$$\frac{1}{n} \sum_{k=0}^{n-1} f_k w_k^{-j} = [f, \Phi_j] = [\beta_0 \Phi_0 + \dots + \beta_{n-1} \Phi_{n-1}, \Phi_j] = \beta_j.$$

Es gilt die folgende Minimaleigenschaft der trigonometrischen Interpolation. □

Satz 4.34 Unter allen trigonometrischen Polynomen q_s der Form $q_s(x) = \gamma_0 + \gamma_1 e^{ix} + \dots + \gamma_s e^{isx}$, $s \leq n-1$, minimiert das s -te Abschnittspolynom $p_s(x) = \beta_0 + \beta_1 e^{ix} + \dots + \beta_s e^{isx}$ für $s = 0, \dots, n-1$ die Fehlerquadratsumme

$$S(q_s) = \sum_{k=0}^{n-1} |f_k - q_s(x_k)|^2 \quad (4.35)$$

Beweis: Ordne p_s, q_s die n -Tupel

$$\begin{aligned}\underline{p}_s &= (p_s(x_0), \dots, p_s(x_{n-1})) \\ \underline{q}_s &= (q_s(x_0), \dots, q_s(x_{n-1}))\end{aligned}$$

zu, die p_s, q_s eindeutig festlegen. Dann gilt

$\frac{1}{n}S(q_s) = [f - \underline{q}_s, f - \underline{q}_s]$. Da aber

$\beta_k = [f, \Phi_k]$ für $k = 0, \dots, n-1$ folgt für $f \in \mathbb{C}^n$

$$[f - \underline{p}_s, \Phi_j] = [f - \sum_{k=0}^s \beta_k \Phi_k, \Phi_j] = \beta_j - \beta_j = 0, j = 0, \dots, s \text{ und}$$

$$\begin{aligned}[f - \underline{p}_s, \underline{p}_s - \underline{q}_s] &= \sum_{j=0}^s [f - \underline{p}_s, (\beta_j - \gamma_j) \Phi_j] = 0, \\ \frac{1}{n}S(q_s) &= [f - \underline{q}_s, f - \underline{q}_s] = [(f - \underline{p}_s + \underline{p}_s - \underline{q}_s), (f - \underline{p}_s + \underline{p}_s - \underline{q}_s)] \\ &= [f - \underline{p}_s, f - \underline{p}_s] + [\underline{p}_s - \underline{q}_s, \underline{p}_s - \underline{q}_s] \\ &\geq [f - \underline{p}_s, f - \underline{p}_s] = \frac{1}{n}S(p_s)\end{aligned}$$

und Gleichheit nur für $\underline{p}_s = \underline{q}_s$ da $[\bullet, \bullet]$ eine Norm definiert. \square

Reelle Version: Beachte $e^{ix} = \cos x + i \sin x$. Setze

$$\begin{aligned}A_j &= \frac{2}{n} \sum_{k=0}^{n-1} f_k \cos \frac{2\pi jk}{n} \\ B_j &= \frac{2}{n} \sum_{k=0}^{n-1} f_k \sin \frac{2\pi jk}{n}\end{aligned} \tag{4.36}$$

für reelle f_k sind A_j, B_j reell und es gilt

→ Übung

$$\beta_{n-j} = \frac{1}{n} \sum_{k=0}^{n-1} f_k w^{j-n} w_k = \frac{1}{n} \sum_{k=0}^{n-1} f_k w_k^j \tag{4.37}$$

$$\beta_j = \frac{1}{2}[A_j - iB_j], \beta_{n-j} = \frac{1}{2}[A_j + iB_j] \tag{4.38}$$

$$\beta_j w_k^j + \beta_{n-j} w_k^{n-j} = A_j \cos jx_k + B_j \sin jx_k \tag{4.39}$$

Damit erhält man die reelle diskrete Fouriertransformation. A_j, B_j sind diskrete Approximationen an die Fourierkoeffizienten

$$\begin{aligned}a_k &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(kx) dx, \quad k = 0, \dots, n-1 \\ b_k &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(kx) dx, \quad k = 1, \dots, n-1.\end{aligned}$$

→ Übung bei Quadratur

Satz 4.40 Gegeben (x_k, f_k) , $k = 0, \dots, N-1$ mit $x_k = \frac{2\pi k}{N}$. Man setze für $j \in \mathbb{Z}$

$$\left. \begin{aligned}A_j &:= \frac{2}{N} \sum_{k=0}^{N-1} f_k \cos(kx_j) \\ B_j &:= \frac{2}{N} \sum_{k=0}^{N-1} f_k \sin(kx_j)\end{aligned} \right\} \tag{4.41}$$

Definiere für ungerades $N = 2M + 1$

$$\Psi(x) = \frac{A_0}{2} + \sum_{k=1}^M (A_k \cos kx + B_k \sin kx) \quad (4.42)$$

und für gerades $N = 2M$

$$\Psi(x) = \frac{A_0}{2} + \sum_{k=1}^{M-1} (A_k \cos kx + B_k \sin kx) + \frac{A_M}{2} \cos Mx \quad (4.43)$$

so gilt

$$\Psi(x_k) = f_k, \quad k = 0, \dots, N-1.$$

Beweis: $N = 2M \implies$

$$f_k = \sum_{j=0}^{N-1} \beta_j w_k^j = \beta_0 + \sum_{j=1}^{M-1} (\beta_j w_k^j + \beta_{N-j} w_k^{(N-j)}) + \beta_M w_k^M$$

mit (4.38), (4.39) folgt

$$B_0 = \frac{2}{N} \sum_{k=0}^{N-1} f_k \sin \frac{2\pi k \cdot 0}{N} = 0,$$

$$B_M = \frac{2}{N} \sum_{k=0}^{N-1} \sin \frac{2\pi k \cdot 1}{N} = 0.$$

Analog für ungerades N . □

Berechnung von Fouriertransformation

Fast Fourier Transformation (FFT) für $N = 2^n$ $n > 0$. Was ist zu berechnen?

$$\begin{aligned} \beta_j &= \frac{1}{N} \sum_{k=0}^{N-1} f_k e^{-\frac{2\pi i k j}{N}} \\ &= \frac{1}{N} \sum_{k=0}^{N-1} f_k \epsilon_N^{kj} \end{aligned}$$

mit $\epsilon_N = e^{-\frac{2\pi i}{N}}$.

Der Trick ist die Tasche auszunutzen, daß die $e^{-\frac{2\pi i}{N}}$ eine N -te Einheitswurzel ist, d.h. alle Potenzen liegen auch auf dem Einheitskreis und die Exponenten lassen sich nach N reduzieren.

Betrachte als Beispiel $N = 4$

$$\begin{aligned} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} &= \frac{1}{N} \begin{bmatrix} \epsilon_N^0 & \epsilon_N^0 & \epsilon_N^0 & \epsilon_N^0 \\ \epsilon_N^0 & \epsilon_N^1 & \epsilon_N^2 & \epsilon_N^3 \\ \epsilon_N^0 & \epsilon_N^2 & \epsilon_N^4 & \epsilon_N^6 \\ \epsilon_N^0 & \epsilon_N^3 & \epsilon_N^6 & \epsilon_N^9 \end{bmatrix} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \end{bmatrix} = \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \epsilon_N^1 & \epsilon_N^2 & \epsilon_N^3 \\ 1 & \epsilon_N^2 & 1 & \epsilon_N^2 \\ 1 & \epsilon_N^3 & \epsilon_N^2 & \epsilon_N^1 \end{bmatrix} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \end{bmatrix} \\ \begin{bmatrix} \beta_0 \\ \beta_2 \\ \beta_1 \\ \beta_3 \end{bmatrix} &= \frac{1}{N} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \epsilon_N^2 & 1 & \epsilon_N^2 \\ 1 & \epsilon_N & \epsilon_N^2 & \epsilon_N^3 \\ 1 & \epsilon_N^3 & \epsilon_N^2 & \epsilon_N \end{bmatrix} = \frac{1}{N} \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & \epsilon_N^2 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & \epsilon_N^2 \end{bmatrix} \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & \epsilon_N^2 & 0 \\ 0 & \epsilon_N & 0 & \epsilon_N^3 \end{bmatrix} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \end{bmatrix} \end{aligned}$$

Bilde Teilschritte: erst

$$\begin{bmatrix} z_0 \\ z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 1 & 0 & \epsilon_N^2 & 0 \\ 0 & \epsilon_N & 0 & \epsilon_N^3 \end{bmatrix} \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ f_3 \end{bmatrix}$$

Verwende $\epsilon_N^2 = -1, \epsilon_N^3 = -\epsilon_N$

$$\begin{aligned} z_0 &= f_0 + f_2 & z_2 &= f_0 - f_2 \\ z_1 &= f_1 + f_3 & z_3 &= \epsilon_N(f_1 - f_3) \end{aligned}$$

$$\begin{bmatrix} \beta_0 \\ \beta_2 \\ \beta_1 \\ \beta_3 \end{bmatrix} = \frac{1}{N} \underbrace{\begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & \epsilon_N^2 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & \epsilon_N^2 \end{bmatrix}}_{2 \text{ komplexe FT der Ordnung 2 mit } \epsilon_N^2} \begin{bmatrix} z_0 \\ z_1 \\ z_2 \\ z_3 \end{bmatrix} = \frac{1}{N} \begin{bmatrix} 1 & 1 & & \\ 1 & -1 & & \\ & & 1 & 1 \\ & & 1 & -1 \end{bmatrix} \begin{bmatrix} z_0 \\ z_1 \\ z_2 \\ z_3 \end{bmatrix}$$

2 komplexe FT der Ordnung 2 mit ϵ_N^2

$$\begin{aligned} \implies N\beta_0 &= z_0 + 1z_1, & N\beta_1 &= z_2 + z_3 \\ N\beta_2 &= z_0 - z_1, & N\beta_3 &= z_2 - z_3 \end{aligned}$$

Es findet jeweils eine Rückführung einer FFT der Ordnung $N = 2^2$ auf 2 FFT der Ordnung $\frac{N}{2} = 2^{n-1}$ statt $+O(N)$ Arbeit. Anstatt N^2 für Matrixmultiplikation erhält man dann $\frac{1}{2}N \log_2 N = 2^{n-1} \times (n-1)$ Operationen.

Da die FFT halber Dimension unabhängig laufen können lassen sie sich gut auf Parallelrechner implementieren.

→ Übung

Allgemeine Formeln:

$$\text{Sei } M = \frac{N}{2}, \text{ Setze } j = 2l$$

$$\begin{aligned} \beta_{2l} &= \sum_{k=0}^{2M-1} f_k \epsilon_N^{2lk} = \sum_{k=0}^{M-1} (f_k + f_{M+k}) \epsilon_N^{2lk} \\ &= \sum_{k=0}^{M-1} (f_k + f_{M+k}) (\epsilon_N^2)^{lk}, \end{aligned} \tag{4.44}$$

wobei wir benutzen

$$\epsilon_N^{2l(M+k)} = \epsilon_N^{2lk} \epsilon_N^{2lM} = \epsilon_N^{2lk}.$$

Mit Hilfwerten

$$z_k = f_k + f_{M+k}, \quad k = 0, \dots, M-1 \tag{4.45}$$

und mit $\epsilon_N^2 = \epsilon_M$ folgt

$$\beta_{2l} = \sum_{k=0}^{M-1} z_k \epsilon_M^{kl} \quad (l = 0, \dots, M-1). \tag{4.46}$$

Für die ungeraden Indizes gilt analog

$$\beta_{2l+1} = \sum_{k=0}^{M-1} z_{M+k} \epsilon_M^{kl} \quad l = 0, \dots, m-1 \tag{4.47}$$

mit Hilfwerten

$$z_{M+k} = (f_k - f_{M+k} \epsilon_N^k), \quad k = 0, \dots, M-1 \tag{4.48}$$

Algorithmus von Cooley/Tukey (1965)

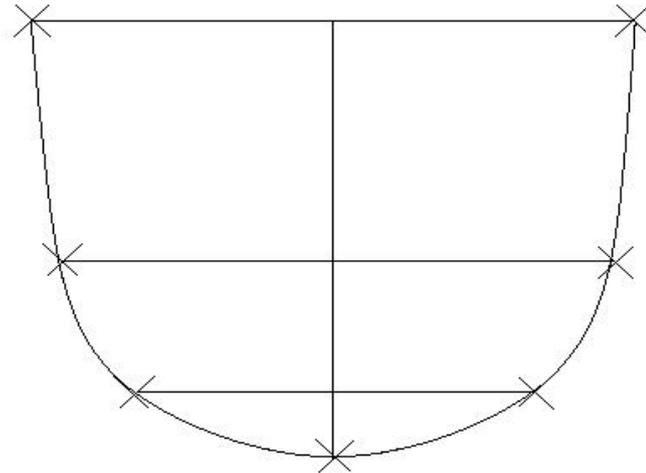
→ Übung 4.85, 4.86

Weitere Tricks möglich.

4.3 Spline Interpolation

Idee: Anstatt Polynome zu verwenden zur Interpolation, verwende stückweise Polynome. Spline kommt aus Schiffbau.

Latte um Schiffsrumpf



Anwendungen: CAD, numerische Lösung partieller + gewöhnlicher Differentialgleichungen. Mehr in Numerik III, Numerik IV; hier nur Grundlagen.

4.3.1 Allgemeine Eigenschaften

Durch $\Delta := \{a = x_0 < x_1 < \dots < x_n = b\}$ sei Unterteilung von $[a, b]$ gegeben.

Definition 4.49 Unter einer zu Δ gehörenden Spline-Funktion S_Δ versteht man eine reelle Funktion $S_\Delta : [a, b] \rightarrow \mathbb{R}$ mit

- a) $S_\Delta \in \mathcal{C}^{2k}[a, b]$ ($2k$ mal stetig differenzierbar).
- b) Auf jedem Teilintervall $[x_i, x_{i+1}]$ stimmt S_Δ mit Polynom $2k + 1$ -ten Grades überein.

Hier im folgenden $k = 1$, kubische Splines. Fast alles läßt sich übertragen.

Definition 4.50

$$\kappa^n(a, b) := \left\{ f : [a, b] \rightarrow \mathbb{R}; \begin{array}{l} f^{(i)}, i = 1, \dots, n-1, \text{ existiert und ist absolut stetig,} \\ f^{(n)} \text{ existiert fast überall auf } [a, b] \text{ und } f^{(n)} \in L^2[a, b] \end{array} \right\}$$

(f absolut stetig, falls zu jedem $\epsilon > 0$ ein $\delta > 0$ so daß für jedes endliche System von Intervallen $[a_i, b_i]$ mit $a \leq a_1 < b_1 \leq a_2 < b_2 \leq \dots \leq a_m < b_m \leq b$ mit $\sum_i |b_i - a_i| < \delta$ gilt:

$$\sum_i |f(b_i) - f(a_i)| < \epsilon.)$$

$$\begin{aligned} \kappa_p^n(a, b) &= \{f \in \kappa^n(a, b), f, f', \dots, f^{(n-1)}, b - a \text{ periodisch, d.h.} \\ &\quad f^{(i)}(a) = f^{(i)}(b) \quad i = 0, \dots, n-1\}. \end{aligned}$$

Es folgt, daß $S_\Delta \in \kappa^{2k+1}(a, b)$.

Eine Spline-Funktion heißt periodisch falls

$$S_\Delta \in \kappa_p^{2k+1}(a, b).$$

Wir wollen mit Splines interpolieren, d.h. wir wollen, daß $S_\Delta(x_i) = f_i = f(x_i)$, $i = 0, \dots, n$.

Dann S_Δ^f . Führe Maß $w(f) := \int_a^b |f''(x)|^2 dx$ ein.

$w(f) \geq 0$ jedoch $w(f) = 0 \quad \forall f = cx + 1$.

Satz 4.51 Sei $f \in \kappa^2(a, b)$, $\Delta = \{a = x_0 < \dots < x_n = b\}$ und S_Δ kubischer Spline zu Δ . Dann gilt

$$w(f - S_\Delta) = w(f) - w(S_\Delta) - 2 \left[(f' - S'_\Delta) S''_\Delta \Big|_a^b - \sum_{i=1}^n (f - S_\Delta) S'''_\Delta \Big|_{x_{i-1}^+}^{x_i^-} \right].$$

Beweis:

$$\begin{aligned} w(f - S_\Delta) &= \int_a^b |f'' - S''_\Delta|^2 dx \\ &= \int_a^b |f''|^2 dx - 2 \int_a^b f'' S''_\Delta dx + \int_a^b |S''_\Delta|^2 dx \\ &= w(f) - 2 \int_a^b (f'' - S''_\Delta) S''_\Delta dx - w(S_\Delta). \end{aligned}$$

Partielle Integration in Teilintervallen.

$$\begin{aligned} \int_{x_{i-1}}^{x_i} (f'' - S''_\Delta) S''_\Delta dx &= (f' - S'_\Delta) S''_\Delta \Big|_{x_{i-1}}^{x_i} - \int_{x_{i-1}}^{x_i} (f' - S'_\Delta) S'''_\Delta dx \\ &= (f' - S'_\Delta) S''_\Delta \Big|_{x_{i-1}}^{x_i} - (f - S_\Delta) S'''_\Delta \Big|_{x_{i-1}^+}^{x_i^-} + \int_{x_{i-1}}^{x_i} (f - S_\Delta) S^{(4)}_\Delta dx. \end{aligned}$$

$S^{(4)} \equiv 0$ auf allen $[x_{i-1}, x_i]$ und

$$\sum_{i=1}^n (f' - S'_\Delta) S''_\Delta \Big|_{x_{i-1}}^{x_i} = (f' - S'_\Delta) S''_\Delta \Big|_a^b$$

impliziert Beh. □

Dies ist Grundlage für Minimaleigenschaft der Splines.

Satz 4.52 Sei $\Delta = \{a = x_0 < x_1 < \dots < x_n = b\}$.

$F = \{f_0, \dots, f_n\}$, $f \in \kappa^2(a, b)$ mit $f(x_i) = f_i$, $i = 0, \dots, n$. Dann gilt

$$w(f - S_\Delta^f) = w(f) - w(S_\Delta^f) \geq 0 \tag{4.53}$$

für alle interpolierenden S_Δ^f die einer der 3 folgenden Bedingungen erfüllen.

$$a) \quad S_\Delta''(a) = S_\Delta''(b) = 0. \quad (4.54)$$

$$\text{oder } b) \quad f \in \kappa_p^2(a, b) \quad S_\Delta \text{ periodisch.} \quad (4.55)$$

$$\text{oder } c) \quad f'(a) = S_\Delta'(a), f'(b) = S_\Delta'(b). \quad (4.56)$$

Für jeden dieser Fälle ist S_Δ eindeutig bestimmt.

Beweis: In jedem dieser Fälle verschwindet

$$((f' - S_\Delta')S_\Delta'')|_b^a - \sum_{i=1}^n (f - S_\Delta)S_\Delta'''|_{x_{i-1}^+}^{x_i^-} = 0.$$

Eindeutigkeit folgt so. Angenommen es gäbe \tilde{S}_Δ^f mit den gleichen Eigenschaften. Dann kann ich $f(x) = \tilde{S}_\Delta^f(x)$ setzen. \implies

$$w(\tilde{S}_\Delta^f(x) - S_\Delta^f(x)) = w(\tilde{S}_\Delta^f(x)) - w(S_\Delta^f(x)) \geq 0 \text{ und genauso } w(S_\Delta^f(x)) - w(\tilde{S}_\Delta^f(x)) \geq 0$$

also $\int_a^b (\tilde{S}_\Delta'' - S_\Delta'')^2 dx = 0$ beide $\tilde{S}_\Delta'', S_\Delta''$ stetig $\implies \tilde{S}_\Delta'' = S_\Delta'' \implies \tilde{S}_\Delta^f = S_\Delta^f + cx + d$, wegen $\tilde{S}_\Delta^f(x) = S_\Delta^f(x)$ für a, b folgt $c = d = 0$. \square

Dieser Satz besagt, daß die Splines unter allen Funktionen aus $\kappa^2(a, b)$ die interpolierende Funktion mit der geringsten Krümmung sind.

$$\kappa(x) = \frac{f''(x)}{\sqrt{1+f'(x)^2}} \approx f''(x) \text{ für } f'(x) \text{ klein.}$$

Analoge Eigenschaften für gelten andere Splines.

4.3.2 Berechnung der Splines

Setze $h_{j+1} = x_{j+1} - x_j$, $j = 0, \dots, n-1$ und

$$M_j = S_\Delta''(x_j), \quad j = 0, \dots, n. \quad (4.57)$$

M_j heißen Momente (aus Physik 2. Ableitungen). Berechne Splines die (4.54), (4.55) oder (4.32) erfüllen mit Hilfe der Momente.

Da S_Δ kubisch in $[x_i, x_{i+1}]$ ist, folgt daß S_Δ'' linear in $[x_i, x_{i+1}]$ und es gilt

$$S_\Delta''(x) = M_j \frac{x_{j+1} - x}{h_{j+1}} + M_{j+1} \frac{x - x_j}{h_{j+1}}, \quad x \in [x_j, x_{j+1}] \quad (2 \text{ Punkte Formel}) \quad (4.58)$$

Integration

$$S_\Delta'(x) = -M_j \frac{(x_{j+1} - x)^2}{2h_{j+1}} + M_{j+1} \frac{(x - x_j)^2}{2h_{j+1}} + \alpha_j \quad (4.59)$$

$$S_\Delta(x) = M_j \frac{(x_{j+1} - x)^3}{6h_{j+1}} + M_{j+1} \frac{(x - x_j)^3}{6h_{j+1}} + \alpha_j(x - x_j) + \beta_j. \quad (4.60)$$

Interpolationsbedingungen

$$\begin{aligned}
S_{\Delta}(x_j) &= f_j \\
\implies M_j \frac{h_{j+1}^2}{6} + \beta_j &= f_j \implies \beta_j = f_j - M_j \frac{h_{j+1}^2}{6}, \\
S_{\Delta}(x_{j+1}) &= f_{j+1} \\
\implies M_{j+1} \frac{h_{j+1}^2}{6} + \alpha_j h_{j+1} + \beta_j &= f_{j+1} \implies \alpha_j = \frac{f_{j+1} - f_j}{h_{j+1}} - \frac{h_{j+1}}{6} (M_{j+1} - M_j),
\end{aligned}$$

und damit

$$S_{\Delta}(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3 \quad (4.61)$$

mit

$$\begin{aligned}
a_j &= f_j, c_j = \frac{M_j}{2}, b_j = -\frac{M_j h_{j+1}}{2} + \alpha_j = S'_{\Delta}(x_j) \\
&= \frac{f_{j+1} - f_j}{h_{j+1}} - \frac{2M_j + M_{j+1} h_{j+1}}{6}, d_j = \frac{M_{j+1} - M_j}{6h_{j+1}} = \frac{S'''_{\Delta}(x_j^+)}{6}
\end{aligned} \quad (4.62)$$

Wir brauchen also nur die M_j zu berechnen. Wir haben gefordert, daß $S'_{\Delta}(x)$ an den Knoten $x = x_j$, $j = 1, \dots, n-1$ stetig sein soll. $\implies n-1$ Gleichungen. Es gilt

$$\begin{aligned}
S'_{\Delta}(x) &= -M_j \frac{(x_{j+1} - x)^2}{2h_{j+1}} + M_{j+1} \frac{(x - x_j)^2}{2h_{j+1}} + \frac{f_{j+1} - f_j}{h_{j+1}} \\
&\quad - \frac{h_{j+1}}{6} (M_{j+1} - M_j)
\end{aligned} \quad (4.63)$$

wenn man die Werte einsetzt. Damit ergibt sich für $j = 1, \dots, n-1$

$$\begin{aligned}
S'_{\Delta}(x_j^-) &= \frac{f_j - f_{j-1}}{h_j} + \frac{h_j}{3} M_j + \frac{h_j}{6} M_{j-1} \\
S'_{\Delta}(x_j^+) &= \frac{f_{j+1} - f_j}{h_{j+1}} - \frac{h_{j+1}}{3} M_j - \frac{h_{j+1}}{6} M_{j+1}
\end{aligned} \quad (4.64)$$

Gleichsetzen wegen Stetigkeit \implies

$$\begin{aligned}
\frac{h_j}{6} M_{j-1} + \frac{h_j + h_{j+1}}{3} M_j + \frac{h_{j+1}}{6} M_{j+1} &= \frac{f_{j+1} - f_j}{h_{j+1}} - \frac{f_j - f_{j-1}}{h_j} \\
&= f[x_j, x_{j+1}] - f[x_{j-1}, x_j], \quad j = 1, \dots, n-1
\end{aligned} \quad (4.65)$$

Das sind $n-1$ Gleichungen für M_0, \dots, M_n .

Die zwei weiteren notwendigen Bedingungen erhält man aus (4.54), (4.55), (4.56).

Aus (4.54)

$$S''_{\Delta}(a) = M_0 = M_n = S''_{\Delta}(b) = 0.$$

Aus (4.55)

$$\begin{aligned}
S''_{\Delta}(a) = S''_{\Delta}(b) &\implies M_0 = M_n. \\
S'_{\Delta}(a) = S'_{\Delta}(b) &\implies \frac{h_n}{6} M_{n-1} + \frac{h_n + h_1}{3} M_n + \frac{h_1}{6} M_1 \\
&= \frac{f_1 - f_n}{h_1} - \frac{f_n - f_{n-1}}{h_n}
\end{aligned}$$

Beweis: Sei

$$A = \begin{bmatrix} 2 & \lambda_0 & & & \\ \mu_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \lambda_{n-1} & \\ & & \mu_n & 2 & \end{bmatrix}$$

Sei $x, y \in \mathbb{R}^{n+1}$ mit $Ax = y$ und setze $\lambda_n = \mu_0 = 0$. Dann gilt

$$\max_i |x_i| \leq \max_i |y_i| \quad (A \text{ ist monoton}).$$

Denn sei r der Index für den $\max_i |x_i|$ angenommen wird.

Dann gilt

$$\begin{aligned} \mu_r x_{r-1} + 2x_r + \lambda_r x_{r+1} &= y_r \\ \implies \max_i |y_i| \geq |y_r| &\geq 2|x_r| - \mu_r |x_{r-1}| - \lambda_r |x_{r+1}| \\ &\geq 2|x_r| - (\mu_r + \lambda_r) |x_r| \\ &\geq |x_r| = \max_i |x_i|. \end{aligned}$$

Angenommen A singular $\implies x \neq 0$ mit $Ax = 0$.

Monotonie führt sofort zum Widerspruch. Anderes System analog. \square

Lösung der Gleichungssysteme:

\longrightarrow Übung: Gauß ohne Pivot (Diagonaldominanz).

Zyklische Reduktion, Sherman–Morrison, iterative Verfahren; alles geht sehr gut.

\longrightarrow Übung!

Typischerweise werden jedoch Bandlöser benutzt, da die Matrizen diagonalähnlich zu einer positiv definiten Matrix sind.

4.3.3 Fehlerabschätzung und Konvergenz bei Interpolation mit Splines

Sei $\Delta_m = \{a = x_0^{(m)} < x_1^{(m)} < \dots < x_{n_m}^{(m)} = b\}$ eine Folge von Unterteilungen von $[a, b]$ und sei $\sigma_m = \max_i (x_{i+1}^{(m)} - x_i^{(m)})$.

Satz 4.70 Sei $f \in C^4[a, b]$ und $|f^{(4)}(x)| \leq L \forall x \in [a, b], \Delta_m$ wie oben mit $\sup_{m,j} \frac{\sigma_m}{x_{j+1}^{(m)} - x_j^{(m)}} \leq \kappa < \infty$. Seien S_{Δ_m} die zu f gehörigen Splines, die die Interpolationsbedingungen

$$\begin{aligned} S_{\Delta_m}(\zeta) &= f(\zeta) & \text{für } \zeta \in \Delta_m \\ S'_{\Delta_m}(x) &= f'(x) & \text{für } x = a, b \end{aligned} \quad (4.71)$$

erfüllen. Dann gibt es von Δ_m unabhängige Konstanten $C_i \leq 2$, so daß für alle $x \in [a, b]$

$$\left| f^{(i)}(x) - S_{\Delta_m}^{(i)}(x) \right| \leq C_i L \kappa \sigma_m^{4-i}, \quad i = 0, 1, 2, 3. \quad (4.72)$$

Beweis: Betrachte feste Zerlegung Δ von $[a, b]$. Die Momente M_j erfüllen für den Fall den wir hier haben, das Gleichungssystem

$$\begin{bmatrix} 2 & \lambda_0 & & & \\ \mu_1 & 2 & \ddots & & \\ & \ddots & \ddots & \lambda_{n-1} & \\ & & \mu_n & 2 & \end{bmatrix} \begin{bmatrix} M_0 \\ \vdots \\ \vdots \\ M_n \end{bmatrix} = \begin{bmatrix} d_0 \\ \vdots \\ \vdots \\ d_n \end{bmatrix}$$

$$\lambda_0 = \mu_n = 1, \lambda_j = \frac{h_{j+1}}{h_j + h_{j+1}}, \mu_j = 1 - \lambda_j, j = 1, \dots, n-1$$

$$d_j = 6f[x_{j-1}, x_j, x_{j+1}], j = 1, \dots, n-1$$

$$d_0 = \frac{6}{h_1} (f[x_0, x_1] - f'_0),$$

$$d_n = \frac{6}{h_n} (f'_n - f[x_{n-1}, x_n]).$$

Definiere

$$F = \begin{bmatrix} f''(x_0) \\ \vdots \\ f''(x_n) \end{bmatrix}, r = d - AF, M = \begin{bmatrix} M_0 \\ \vdots \\ M_n \end{bmatrix}.$$

Dann gilt $r_0 = d_0 - 2f''(x_0) - f''(x_1)$. Schätze zuerst $\|r_0\|_\infty$ ab. Taylorentwicklung um x_0 liefert

$$\begin{aligned} r_0 &= \frac{6}{h_1} (f[x_0, x_1] - f'_0) - 2f''(x_0) - f''(x_1) \\ &= \frac{6}{h_1} \left(f'(x_0) + \frac{h_1}{2} f''(x_0) + \frac{h_1^2}{6} f'''(x_0) + \frac{h_1^3}{24} f^{(4)}(\tau_1) - f'(x_0) \right) \\ &\quad - 2f''(x_0) - \left(f''(x_0) + h_1 f'''(x_0) + \frac{h_1^2}{2} f^{(4)}(\tau_2) \right) \\ &= \frac{h_1^2}{4} f^{(4)}(\tau_1) - \frac{h_1^2}{2} f^{(4)}(\tau_2), \text{ für } \tau_1, \tau_2 \in [x_0, x_1] \end{aligned}$$

$$\implies |r_0| \leq \frac{3}{4} L \sigma_m^2. \quad (4.73)$$

Analog gilt für $r_n = d_n - f''(x_{n-1}) - 2f''(x_n)$ die Abschätzung

$$|r_n| \leq \frac{3}{4} L \sigma_m^2. \quad (4.74)$$

und für $j = 1, \dots, n-1$ mit Taylorentwicklung um x_j .

$$\begin{aligned} r_j &= d_j - \mu_j f''(x_{j-1}) - 2f''(x_j) - \lambda_j f''(x_{j+1}) \\ &= 6f[x_{j-1}, x_j, x_{j+1}] - \frac{h_j}{h_j + h_{j+1}} f''(x_{j-1}) - 2f''(x_j) - \frac{h_{j+1}}{h_j + h_{j+1}} f''(x_{j+1}) \\ &= \frac{1}{h_j + h_{j+1}} \left\{ 6 \left(f'(x_j) + \frac{h_{j+1}}{2} f''(x_j) + \frac{h_{j+1}^2}{6} f'''(x_j) + \frac{h_{j+1}^3}{24} f^{(4)}(\tau_{j1}) \right. \right. \\ &\quad \left. \left. - f'(x_j) + \frac{h_j}{2} f''(x_j) - \frac{h_j^2}{6} f'''(x_j) + \frac{h_j^3}{24} f^{(4)}(\tau_{j2}) \right) \right. \\ &\quad \left. - h_j \left(f''(x_j) - h_j f'''(x_j) + \frac{h_j^2}{2} f^{(4)}(\tau_{j3}) \right) \right. \\ &\quad \left. - 2f''(x_j)(h_j + h_{j+1}) \right. \\ &\quad \left. - h_{j+1} \left(f''(x_j) + h_{j+1} f'''(x_j) + \frac{h_{j+1}^2}{2} f^{(4)}(\tau_{j4}) \right) \right\} \\ &= \frac{1}{h_j + h_{j+1}} \left(\frac{h_{j+1}^3}{4} f^{(4)}(\tau_{j1}) + \frac{h_j^3}{4} f^{(4)}(\tau_{j2}) - \frac{h_j^3}{2} f^{(4)}(\tau_{j3}) - \frac{h_{j+1}^3}{2} f^{(4)}(\tau_{j4}) \right), \end{aligned}$$

mit $\tau_{j_i} \in [x_{j-1}, x_{j+1}]$, $i = 1, 2, 3, 4$.

$$\implies |r_j| \leq \frac{3}{4}L \frac{h_{j+1}^3 + h_j^3}{h_j + h_{j+1}} \leq \frac{3}{4}L\sigma_m^2.$$

Insgesamt gilt

$$\begin{aligned} \|r\|_\infty &\leq \frac{3}{4}L\sigma_m^2 \implies \\ \|A(M - F)\|_\infty &\leq \frac{3}{4}L\sigma_m^2 \implies \|M - F\|_\infty \leq \frac{3}{4}L\sigma_m^2, \end{aligned} \quad (4.75)$$

da A monoton.

Sei $x \in [x_{j-1}, x_j]$. Dann folgt

$$\begin{aligned} S_\Delta''' - f'''(x) &= \frac{M_j - M_{j-1}}{h_j} - f'''(x) \\ &= \frac{M_j - f''(x_j)}{h_j} - \frac{M_{j-1} - f''(x_{j-1})}{h_j} \\ &\quad + \frac{[f''(x_j) - f''(x)] - [f''(x_{j-1}) - f''(x)]}{h_j} - f'''(x). \end{aligned}$$

Taylorentwicklung um x und (4.75) \implies

$$\begin{aligned} |S_\Delta'''(x) - f'''(x)| &\leq \frac{3}{2}L \frac{\sigma_m^2}{h_j} + \frac{1}{h_j} \left| (x_j - x)f'''(x) + \frac{(x_j - x)^2}{2}f^{(4)}(\eta_1) - (x_{j-1} - x)f'''(x) \right. \\ &\quad \left. - \frac{(x_{j-1} - x)^2}{2}f^{(4)}(\eta_2) - h_j f'''(x) \right| \\ &\leq \frac{3}{2}L \frac{\sigma_m^2}{h_j} + \frac{1}{2} \frac{\sigma_m^2}{h_j} L, \quad \tau_1, \tau_2 \in [x_{j-1}, x_j] \\ &\leq 2L\kappa\sigma_m, \quad \text{wegen Vor. } \frac{\sigma_m}{h_j} \leq \kappa. \end{aligned}$$

Hierbei wurde noch $\max_{x \in [x_{j-1}, x_j]} \{(x_j - x)^2 + (x - x_{j-1})^2\} = h_j^2$ benutzt. Nun gibt es für jedes $x \in (a, b)$ ein x_j mit

$$\begin{aligned} |x_j - x| &\leq \frac{1}{2}\sigma_m \\ \implies f''(x) - S_\Delta''(x) &= f''(x_j) - S_\Delta''(x_j) + \int_{x_j}^x (f'''(t) - S_\Delta'''(t)) dt \end{aligned}$$

\implies mit der Abschätzung für $f''' - S_\Delta''' \implies$

$$\begin{aligned} |f''(x) - S_\Delta''(x)| &\leq \frac{3}{4}L\sigma_m^2 + \frac{1}{2}\sigma_m \cdot 2L\kappa\sigma_m \\ &\leq \frac{7}{4}L\kappa\sigma_m^2, \quad \text{da } \kappa \geq 1. \end{aligned}$$

Damit haben wir die Abschätzung für $f'' - S_\Delta''$.

Es gilt $f(x_{j-1}) = S_\Delta(x_{j-1})$, $f(x_j) = S_\Delta(x_j)$. Außer den Randpunkten $\alpha_0 = a$, $\alpha_{n+1} = b$ gibt es also noch den Satz von Rolle n Stellen $\alpha_j \in (x_{j-1}, x_j)$ mit

$$f'(\alpha_j) = S'_\Delta(\alpha_j) \quad j = 0, \dots, n+1.$$

Es gilt also für jedes $x \in [a, b]$ ein α_j mit $|\alpha_j - x| < \sigma_m$. Daher gilt für alle $x \in [a, b]$:

$$f'(x) - S'_\Delta(x) = \int_{\alpha_j(x)}^x (f''(t) - S''_\Delta(t)) dt$$

$$\implies |f'(x) - S'_\Delta(x)| \leq \frac{7}{4} L\kappa\sigma_m^2 \cdot \sigma_m = \frac{7}{4} L\kappa\sigma_m^3.$$

Analog folgt für $x \in [a, b]$ (Hier wieder $|x - x_j| \leq \frac{1}{2}\sigma_m$)

$$f(x) - S_\Delta(x) = \int_{x_j}^x f'(t) - S'_\Delta(t) dt$$

und damit

$$|f(x) - S_\Delta(x)| \leq \frac{7}{4} L\kappa\sigma_m^3 \frac{1}{2}\sigma_m = \frac{7}{8} L\kappa\sigma_m^4.$$

□

Es gäbe noch viel mehr zu Splines zu sagen.

4.4 Übungen zu Kapitel 4

Übung 4.76 Seien $x_0, \dots, x_n \in \mathbb{R}$ paarweise verschieden und $f_0, \dots, f_n \in \mathbb{R}$ zugehörige Funktionswerte. Definiere die dividierten Differenzen $f[x_i, x_{i+1}, \dots, x_{i+k}]$ rekursiv durch

$$f[x_i] = f_i, \quad i = 0, \dots, n,$$

sowie

$$f[x_i, x_{i+1}, \dots, x_{i+k}] := \frac{f[x_{i+1}, x_{i+2}, \dots, x_{i+k}] - f[x_i, x_{i+1}, \dots, x_{i+k-1}]}{x_{i+k} - x_i},$$

$i = 0, \dots, n, k = 0, \dots, n - i$. Nehme an, daß

$$f[x_k], f[x_{k-1}, x_k], f[x_{k-2}, x_{k-1}, x_k], \dots, f[x_0, \dots, x_k]$$

bereits vorhanden sind. Wie kann man daraus $f[x_0, \dots, x_{k+1}]$ berechnen? Was folgt daraus insgesamt für die Berechnung von $f[x_n], f[x_{n-1}, x_n], f[x_{n-2}, x_{n-1}, x_n], \dots, f[x_0, \dots, x_n]$?

Übung 4.77 Seien $x_0, \dots, x_n \in \mathbb{R}$ paarweise verschieden und $f_0, \dots, f_n \in \mathbb{R}$ zugehörige Funktionswerte und $p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0$ das eindeutige Polynom n -ten Grades mit $p(x_i) = f_i, i = 0, \dots, n$.

1. Bestimme den Rechenaufwand (flops) zur Bestimmung von p für die Newtondarstellung von p (d.h. die dividierten Differenzen) sowie für die Berechnung der Koeffizienten von p durch Lösen des linearen Gleichungssystems $Va = f$. Dabei ist V die Matrix aus Aufgabe 3 und $a = (a_0, \dots, a_n)^T, f = (f_0, \dots, f_n)^T$.
2. Wie teuer ist die Auswertung von p in der Newtondarstellung bei gegebenen dividierten Differenzen, bzw. bei gegebenen Koeffizienten a_0, \dots, a_n verglichen mit der Lagrange-Darstellung von p ?
3. Welche der drei Varianten ist die billigste, wenn man $p(x)$ ein einziges Mal auswerten will. Wie sieht es für eine große Anzahl von Werten x (Größenordnung n) aus?

Übung 4.78 Zeige, daß

$$\frac{(b-x)^2}{(b-a)^2} \left[f_1 \frac{2x+b-3a}{b-a} + s_1(x-a) \right] + \frac{(a-x)^2}{(a-b)^2} \left[f_2 \frac{2x+a-3b}{a-b} + s_2(x-b) \right]$$

das eindeutige Polynom dritten Grades ist, welches für $a < b$ und vorgegebene Werte $f_1, s_1, f_2, s_2 \in \mathbb{R}$ die Interpolationsaufgabe

$$p(a) = f_1, p'(a) = s_1, p(b) = f_2, p'(b) = s_2$$

erfüllt ('Hermite-Interpolierende').

Übung 4.79 Seien $x_0, \dots, x_n \in \mathbb{R}$ und

$$V = \begin{pmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^n \end{pmatrix}.$$

Zeige, dass V (Vandermonde-Matrix) invertierbar ist genau dann, wenn x_0, \dots, x_n paarweise verschieden sind.

Übung 4.80 Zu vorgegebenen Stützstellen $t_0 < t_1 < \dots < t_n$ mit Funktionswerten f_0, \dots, f_n ist das Neville-Schema zur Berechnung des zugehörigen Interpolationspolynoms p an der Stelle x wie folgt definiert: $P_{i,0} := f_i, n \geq i \geq 0$.

Für $k = 1, \dots, n$ setze $P_{i,k} := \frac{(x - t_{i-k})P_{i,k-1} - (x - t_i)P_{i-1,k-1}}{t_i - t_{i-k}}, n \geq i \geq k$.

Dann ist $p(x) = P_{n,n}$.

- Lässt sich das Neville-Schema in einem ähnlichen Schema wie die dividierten Differenzen aufschreiben?
- Wie teuer ist die Berechnung von $P_{n,n}$ im günstigsten Fall und ist dies günstiger oder teurer als die Newton-Darstellung mittels dividierter Differenzen?
- Nehme an, dass $x = 0, t_i = h^2 4^{-i}, i = 0, \dots, n$. Schreibe das Neville-Schema für diesen Fall um.

Übung 4.81 Schreibe ein **MATLAB**-Programm, welches für beliebige Vorgabe $t_0 < t_1 < \dots < t_n$ von Stützstellen mit Funktionswerten f_0, \dots, f_n das zugehörige Interpolationspolynom p an der Stelle x mit Hilfe des Neville-Schemas auswertet, welches wie folgt definiert ist:

$P_{i,0} := f_i, n \geq i \geq 0$.

Für $k = 1, \dots, n$ setze $P_{i,k} := \frac{(x - t_{i-k})P_{i,k-1} - (x - t_i)P_{i-1,k-1}}{t_i - t_{i-k}}, n \geq i \geq k$.

Dann ist $p(x) = P_{n,n}$. Wende das Programm auf $f(t) = \log_{10}(t) - \frac{t-1}{t}, x = 5.25$ und folgende Stützstellenverteilung an und interpretiere die Ergebnisse.

- 1.0, 2.0, 4.0, 8.0, 10.0
- 2.0, 4.0, 8.0, 10.0
- 4.0, 8.0, 10.0
- 2.0, 4.0, 8.0

Übung 4.82 Bestimme mit dem Newtonschema die Interpolationspolynome p zu den Funktionen $f(t) = \frac{1}{1+25t^2}$ und $g(t) = \sqrt{|t|}$. Als Stützstellen verwende jeweils

- i) $t_i = -1 + ih, i = 0, \dots, n; \quad h = \frac{2}{n},$
 ii) $t_i = \cos \frac{(2i+1)\pi}{2(n+1)}, i = 0, \dots, n,$

für $n = 2, 4, 6, \dots, 20$. Werte das Polynom an den Stellen $y_j^i = t_i + j \frac{t_{i+1} - t_i}{21}, j = 1, \dots, 20$ aus. Vergleiche als Schätzung für den maximalen Fehler $\|p - f\|_{\infty [-1,1]}$

$$\max_{i \in \{0, \dots, n-1\}, j \in \{1, \dots, 20\}} |p(y_j^i) - f(y_j^i)|.$$

Plotte die Interpolationspolynome zu den Stützstellen in i), ii) für $n = 20$. Erläutere die Ergebnisse.

Übung 4.83 Gegeben sei die Funktion $f : \mathbb{R} \mapsto \mathbb{R}$ durch

$$f|_{[0,2\pi)}(t) := \begin{cases} 1 & \text{für } 0 \leq t < \pi \\ -1 & \text{für } \pi \leq t < 2\pi \end{cases}$$

und $f(t + 2k\pi) = f(t)$ für alle $k \in \mathbb{Z}$. (Rechteckimpuls erster Art)

i) Bestimme die zu f gehörige Fourierreihe

$$\frac{1}{2}a_0 + \sum_{k=1}^{\infty} (a_k \cos kx + b_k \sin kx), \quad a_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \cos kx \, dx, \quad b_k = \frac{1}{\pi} \int_0^{2\pi} f(x) \sin kx \, dx.$$

ii) Zeichne (oder Plotte) die ersten fünf Partialsummen und vergleiche mit f .

Übung 4.84 Seien $x_0, x_1, x_2 \in \mathbb{R}$ mit $x_0 \neq x_2$. Zeige:

1. es gibt eindeutige Polynome p_1, p_2, p_3, p_4 dritten Grades, welche den Bedingungen

$$\begin{array}{cccc} p_1(x_0) = 1 & p_2(x_0) = 0 & p_3(x_0) = 0 & p_4(x_0) = 0 \\ p_1(x_2) = 0 & p_2(x_2) = 1 & p_3(x_2) = 0 & p_4(x_2) = 0 \\ p_1'(x_1) = 0 & p_2'(x_1) = 0 & p_3'(x_1) = 1 & p_4'(x_1) = 0 \\ p_1''(x_1) = 0 & p_2''(x_1) = 0 & p_3''(x_1) = 0 & p_4''(x_1) = 1 \end{array}$$

genügen.

2. Es existiert ein eindeutiges Polynom p dritten Grades gibt, daß für die Werte f_0, f_2, f_1', f_1'' die Interpolationsaufgabe

$$p(x_0) = f_0, \quad p(x_2) = f_2, \quad p'(x_1) = f_1', \quad p''(x_1) = f_1''$$

erfüllt.

Übung 4.85 Sei $c \in \mathbb{C}^n, J \in \mathbb{C}^{n \times n}$ definiert durch

$$c := \begin{pmatrix} c_0 \\ c_1 \\ \vdots \\ c_{n-1} \end{pmatrix}, \quad J := \begin{pmatrix} 0 & & & 1 \\ 1 & \ddots & & \\ & \ddots & \ddots & \\ & & & 1 & 0 \end{pmatrix}$$

Definiere die Matrix $C \in \mathbb{C}^{n \times n}$ durch

$$C := (c, Jc, J^2c, \dots, J^{n-1}c)$$

(C bezeichnet man auch als Zirkulante).

Sei $\omega_n := e^{-i \frac{2\pi}{n}}$ und $W \in \mathbb{C}^{n \times n}$ definiert als $W := \frac{1}{\sqrt{n}} (\omega_n^{kl})_{k,l=0, \dots, n-1}$. Zeige:

- i) $W^*W = I$, d.h. W ist unitär.
- ii) W^*CW ist eine Diagonalmatrix. (Tip: Zeige zuerst, daß W^*JW diagonal ist und nutze dann die Def. von C).
- iii) Was sind die Eigenwerte von C und wie lassen sie sich berechnen?
- iv) Sei $n = 2^k$. Zeige, daß sich das Gleichungssystem $Cx = b$, $b \in \mathbb{C}^n$ gegeben, durch drei FFT berechnen läßt. Bestimme die wesentliche Rechenzeit.

Übung 4.86 Gegeben die Werte

$$\begin{array}{cccc} y_0 = 0.580 & y_1 = 0.951 & y_2 = 0.786 & y_3 = 0.298 \\ y_4 = 0.454 & y_5 = 0.006 & y_6 = 0.276 & y_7 = 0.306 \end{array}$$

Führe eine FFT durch, d.h. berechne

$$c_k = \sum_{m=0}^{n-1} \omega_n^{km} y_m, \quad \text{für } k = 0, \dots, n-1$$

mittels FFT, wobei $\omega_n := e^{-i\frac{2\pi}{n}}$.

Übung 4.87 Seien $t_{-3} < t_{-2} < \dots < t_{n+2} < t_{n+3}$ und $t_{i+1} - t_i = h$ für alle $i = -3, \dots, n+2$. Definiere für $i = -1, \dots, n+1$ die Funktion

$$B_i(t) = \frac{1}{h^3} \begin{cases} (t - t_{i-2})^3 & t \in [t_{i-2}, t_{i-1}) \\ h^3 + 3h^2(t - t_{i-1}) + 3h(t - t_{i-1})^2 - 3(t - t_{i-1})^3 & t \in [t_{i-1}, t_i) \\ h^3 + 3h^2(t_{i+1} - t) + 3h(t_{i+1} - t)^2 - 3(t_{i+1} - t)^3 & t \in [t_i, t_{i+1}) \\ (t_{i+2} - t)^3 & t \in [t_{i+1}, t_{i+2}) \\ 0 & \text{sonst} \end{cases}$$

Zeige:

- $B_i(t)$ ist zweimal stetig differenzierbar.
- Skizziere $B_i(t)$.
- Betrachte für $p(t) = \sum_{i=-1}^{n+1} a_i B_i(t)$ und Funktionswerte f_0, \dots, f_n die Interpolationsaufgabe $p(t_i) = f_i$, $i = 0, \dots, n$, $p''(t_0) = 0 = p''(t_n)$. Lässt sich die Aufgabe eindeutig lösen und wenn ja, wie teuer ist die Berechnung von a_0, \dots, a_n ? Wie teuer ist die Auswertung von $p(x)$ bei gegebenen a_0, \dots, a_n ?
- Was ändert sich bei der abgewandelten Aufgabe $p(t_i) = f_i$, $i = 0, \dots, n$, $p'(t_0) = f'_0$, $p'(t_n) = f'_n$ bei zusätzlichen Funktionswerten f'_0, f'_n ?

Übung 4.88 Sei $A \in \mathbb{R}^{n,n}$.

$$1. \text{ Sei } B \in \mathbb{R}^{n,n}, B = \begin{pmatrix} 0 & 0 & b_{1n} \\ 0 & 0 & 0 \\ b_{n1} & 0 & 0 \end{pmatrix} \text{ Finde } U, V \in \mathbb{R}^{n,2}, \text{ so dass } B = UV^T.$$

2. Zeige die Sherman-Morrison-Woodbury Formel, d.h. sind A und $A + UV^T$ nicht singulär, dann ist auch $C = I + V^T A^{-1} U$ nichtsingulär und es gilt:

$$(A + UV^T)^{-1} = A^{-1} - A^{-1}UC^{-1}V^T A^{-1}.$$

3. Betrachte das lineare Gleichungssystem $(A + B)x = b$, wobei A und $A + B$ nichtsingulär seien mit B aus Teil 1. Nehme an, dass A tridiagonal ist und A und $A + B$ eine LU-Zerlegung besitzen (ohne Pivotisierung). Wie teuer ist die Lösung von $(A + B)x = b$ wenn man einmal eine LU-Zerlegung von $A + B$ macht verglichen mit der Variante, bei der man zunächst eine LU-Zerlegung von A macht und dann die Sherman-Morrison-Woodbury Formel verwendet.

Übung 4.89 Bestimme die kubischen Splines S_Δ zu den Funktionen $f(t) = \frac{1}{1+25t^2}$ und $g(t) = \sqrt{|t|}$, wobei $S''_\Delta(-1) = S''_\Delta(1) = 0$ (natürliche Splineinterpolation). Als Stützstellen verwende $t_i = -1 + ih$, $i = 0, \dots, n$; $h = \frac{2}{n}$, für $n = 2, 4, 6, \dots, 20$. Werte den Spline an den Stellen $y_j^i = t_i + j \frac{t_{i+1} - t_i}{21}$, $j = 1, \dots, 20$ aus. Vergleiche als Schätzung für den maximalen Fehler $\|S_\Delta - f\|_\infty$ auf $[-1, 1]$

$$\max_{i \in \{0, \dots, n-1\}, j \in \{1, \dots, 20\}} |S_\Delta(y_j^i) - f(y_j^i)|.$$

Plotte die Splines zu den Stützstellen für $n = 20$. Erläutere die Ergebnisse.

Übung 4.90 Seien $\alpha, \beta \in \mathbb{R}$ und die Funktion $s_{\alpha, \beta} : [-1, 2] \rightarrow \mathbb{R}$ definiert durch

$$s_{\alpha, \beta}(x) = \begin{cases} (x+1)^4 + \alpha(x-1)^4 + 1 & x \in [-1, 0] \\ -x^3 - 8\alpha x + 1 & x \in (0, 1] \\ \beta x^3 + 8x^2 + \frac{11}{3} & x \in (1, 2] \end{cases}$$

Bestimme α, β so, dass $s_{\alpha, \beta}$ ein kubischer Spline bezüglich der Knotenverteilung $\Delta = \{-1, 0, 1, 2\}$ ist.

Kapitel 5

Numerische Integration

Aufgabe der numerischen Integration (Quadratur) ist die Berechnung von

$$\int_a^b f(x)dx, \quad a, b < \infty$$

Typischerweise ist f nicht direkt integrierbar.

Idee: Approximiere $f(x)$ durch Funktion, die einfach zu integrieren ist; Polynome, trigonometrische Polynome, Splines, rationale Funktionen, Exponentialsummen, ...

5.1 Newton–Cotes Formeln

Die Idee der Newton–Cotes Formeln ist $f(x)$ durch ein Polynom zu interpolieren und dieses zu integrieren. Bilde äquidistante Unterteilung von $[a, b]$ $x_i = a + ih$ $i = 0, \dots, n$, mit $h = \frac{b-a}{n}$ $n > 0$ $n \in \mathbb{N}$ und bilde Interpolationspolynom in Π_n , das

$$P_n(x_i) = f(x_i) \quad i = 0, \dots, n \tag{5.1}$$

erfüllt. Satz (4.2) \implies

$$P_n(x) = \sum_{i=0}^n f_i L_i(x). \tag{5.2}$$

\implies

$$\begin{aligned} \int_a^b P_n(x)dx &= \sum_{i=0}^n f_i \int_a^b L_i(x)dx \quad \text{mit } x = a + hs \\ &= h \sum_{i=0}^n f_i \int_0^n \prod_{\substack{k=0 \\ k \neq i}}^n \frac{s-k}{i-k} ds \\ &= h \sum_{i=0}^n f_i \alpha_i = \frac{b-a}{n} \sum_{i=0}^n \sigma_i f_i \quad \sigma_i \in \mathbb{Z}. \end{aligned} \tag{5.3}$$

Gewichte α_i sind unabhängig von f_i, a, b jedoch abhängig von n .

Die Gewichte liegen tabelliert vor. Diese Formeln heißen Newton–Cotes–Formeln.

Fehler:

$$\int_a^b (P_n(x) - f(x)) dx = h^{p+1} \kappa \cdot f^{(p)}(\xi), \xi \in (a, b). \quad (5.4)$$

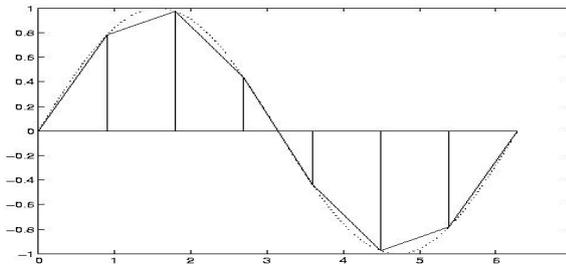
Tabelle 5.5

n	σ_i	ns	Fehler	Name
1	1 1	2	$\frac{h^3}{12} f^{(2)}(\xi)$	Trapezregel
2	1 4 1	6	$\frac{h^5}{90} f^{(4)}(\xi)$	Simpsonregel
3	1 3 3 1	8	$\frac{h^5}{80} f^{(4)}(\xi)$	3/8 Regel
4	7 32 12 32 7	90	$\frac{h^7}{945} f^{(6)}(\xi)$	Milne-Regel
5	19 75 50 50 75 19	288	$\frac{h^7}{12096} f^{(6)}(\xi)$	–
6	41 216 27 272 27 216 41	840	$\frac{h^9}{1400} f^{(8)}(\xi)$	Weddle-Regel

größere $n \implies$ negative Gewichte.

Wegen der schlechten Approximationseigenschaften (global) von Polynomen werden die Regeln stückweise auf Teilintervalle angewendet und dann aufsummiert.

Beispiel 5.6 (Summierte Trapezregel: stückweise lineare Interpolation)



Teilintervalle $[x_i, x_{i+1}]$, $x_i = a + ih$ $i = 0, \dots, N$ $h = \frac{b-a}{N}$. Dann ist die summierte Trapezregel

$$\begin{aligned} T(h) &:= \sum_{i=0}^{N-1} \frac{h}{2} [f(x_i) + f(x_{i+1})] \\ &= h \left[\frac{f(a)}{2} + f(a+h) + \dots + f(b-h) + \frac{f(b)}{2} \right] \end{aligned} \quad (5.7)$$

Für jedes Teilintervall ist der Fehler

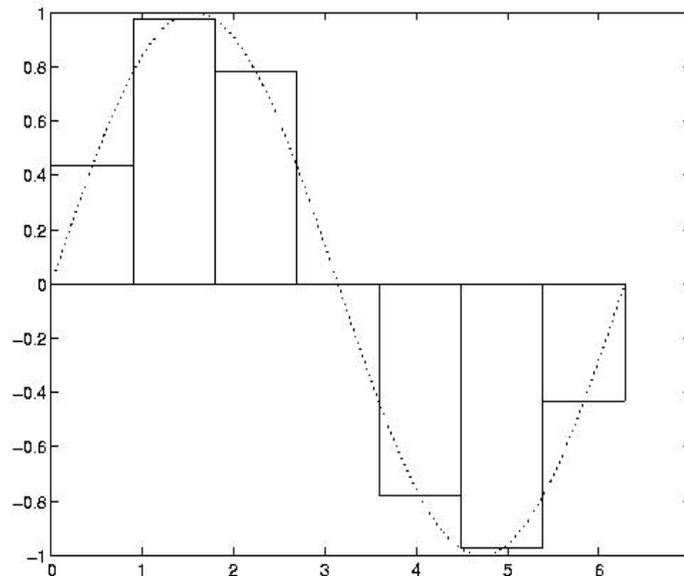
$$\frac{h}{2} [f(x_i) - f(x_{i+1})] - \int_{x_i}^{x_{i+1}} f(x) dx = \frac{h^3}{12} f^{(2)}(\xi_i) \text{ für } \xi_i \in (x_i, x_{i+1}). \quad (5.8)$$

Insgesamt ergibt sich also

$$\begin{aligned} \left| T(h) - \int_a^b f(x) dx \right| &= \left| \sum_{i=0}^{n-1} \frac{h^3}{12} f^{(2)}(\xi_i) \right| \\ &\leq \frac{h^3}{12} N \cdot \underbrace{\sup_{\xi \in [a,b]} f''(\xi)}_{=M} = \frac{b-a}{12} h^2 M. \end{aligned} \quad (5.9)$$

Fehler geht mit h^2 gegen 0, Verfahren 2. Ordnung.

Beispiel 5.10 (Mittelpunktregel)



Beispiel 5.11 Sei N gerade und verwende Simpsonregel auf $[x_{2i}, x_{2i+1}, x_{2i+2}]$, $i = 0, \dots, \frac{N}{2} - 1$. Der Näherungswert auf jedem der Teilintervalle ist

$$\frac{h}{3} [f(x_{2i}) + 4f(x_{2i+1}) + f(x_{2i+2})].$$

Summation ergibt

$$\begin{aligned} S(h) &= \frac{h}{3} [(f(a) + 4f(a+h) + f(a+2h)) + \\ &\quad + 4f(a+3h) + 2f(a+4h) + 4f(a+5h) + f(a+6h)) + \\ &\quad + 2f(b-2h) + 4f(b-h) + f(b)] \end{aligned} \quad (5.12)$$

und als Fehler

$$\left| S(h) - \int_a^b f(x) dx \right| \leq \frac{h^5}{90} \cdot \frac{N}{2} \cdot \underbrace{\sup_{\xi \in [a,b]} f^{(4)}(\xi)}_{=M} \quad (5.13)$$

$$= \frac{b-a}{180} h^4 \cdot M.$$

Verfahren 4. Ordnung.

Weitere Quadraturformeln erhält man durch andere Approximationsaufgaben für $f(x)$.
Zum Beispiel Gauss-Quadratur für Integrale der Form

$$\int_a^b w(x)f(x)dx \approx \sum_{i=1}^n w_i f(x_i).$$

wobei w_i, f_i so gewählt werden, daß der Fehler

$$\int_a^b w(x)f(x)dx - \sum_{i=1}^n w_i f(x_i)$$

für Polynome möglichst hohen Grades verschwindet \rightarrow Übung.
Andere Ideen später bei der Lösung von Dgl.

5.2 Extrapolation

Untersuche das asymptotische Verhalten der Trapezsumme $T(h)$ für $h \rightarrow 0$. Dazu brauchen wir zuerst einige Eigenschaften der Bernoulli-Polynome, die wie folgt rekursiv definiert sind.

Definition 5.14 Die Polynome $B_k(x), k = 0, 1, 2, \dots$, die durch die rekursive Folge

$$\begin{aligned} B_0(x) &= 1, B'_k(x) = k \cdot B_{k-1}(x) \quad k \geq 1, \\ \int_0^1 B_k(x) dx &= 0, \quad k \geq 1 \end{aligned} \tag{5.15}$$

definiert sind, heißen Bernoulli-Polynome.

Eigenschaften der Bernoulli-Polynome:

$$B_k(x) = a_k + k \int_0^x B_{k-1}(t) dt \quad k \geq 1 \tag{5.16}$$

wobei a_k so zu bestimmen ist, daß $\int_0^1 B_k(x) dx = 0$. Man erhält sofort

$$\begin{aligned} B_1(x) &= x - \frac{1}{2}, B_2 = x^2 - x + \frac{1}{6}, B_3(x) = x^3 - \frac{3}{2}x^2 + \frac{1}{2}x \\ B_4(x) &= x^4 - 2x^3 + x^2 - \frac{1}{30}, \dots \end{aligned} \tag{5.17}$$

und es gilt das folgende Lemma.

Lemma 5.18

i) $B_k(0) = B_k(1)$ für $k \geq 2$.

ii) Die Polynome $P_k(x) = B_k(\frac{1}{2} + x)$ sind für gerades k gerade, d.h. $P(x) = P(-x)$ und für ungerades k ungerade Polynome, d.h. $P(x) = -P(-x)$ in x .

iii) $B_{2k+1}(0) = 0$ für alle $k \geq 1$.

Beweis:

i) Aus (5.16)

$$B_k(1) - B_k(0) = k \int_0^1 B_{k-1}(t) dt = 0 \text{ für } k \geq 2.$$

ii) $P_0(x) \equiv 1$ ist gerade. Sei $P_{2k}(x)$ für $k \geq 0$ gerades Polynom. Zeige, dass $P_{2k+2}(x)$ wieder gerade ist. Es gilt

$$P'_{2k+1}(x) = B'_{2k+1}(\frac{1}{2} + x) = (2k+1)B_{2k}(\frac{1}{2} + x) = (2k+1)P_{2k}(x)$$

Also

$$\int_{-\frac{1}{2}}^{\frac{1}{2}} P'_{2k+1}(x) dx = 0.$$

Aus $P_{2k}(x)$ gerade folgt, daß

$$q(x) := (2k+1) \int_0^x P_{2k}(t) dt$$

ein ungerades Polynom ist und es gilt

$$\int_{-\frac{1}{2}}^{\frac{1}{2}} q(x) dx = - \int_{+\frac{1}{2}}^{-\frac{1}{2}} q(-x) dx = \int_{\frac{1}{2}}^{-\frac{1}{2}} q(x) dx = - \int_{-\frac{1}{2}}^{\frac{1}{2}} q(x) dx = 0.$$

Weiterhin ist $q'(x) = (2k+1)P_{2k}(x)$ und $\int_{-\frac{1}{2}}^{\frac{1}{2}} q(x) dx = 0 \implies q(x) = P_{2k+1}(x)$ ungerade.

$P_{2k+2}(x) = (2k+2) \int_0^x P_{2k+1}(t) dt + a \implies P_{2k+2}$ ist wieder gerade. \implies ii).

iii) Für $k \geq 1$ gilt

$$\begin{aligned} B_{2k+1}(1) &= B_{2k+1}(0) = P_{2k+1}(-\frac{1}{2}) = -P_{2k+1}(\frac{1}{2}) = -B_{2k+1}(1) \\ \implies B_{2k+1}(0) &= B_{2k+1}(1) = 0. \end{aligned}$$

□

$\beta_k = (-1)^{k+1} B_{2k}(0)$ heißen Bernoullische Zahlen.

$\beta_1 = \frac{1}{6}, \beta_2 = \frac{1}{30}, \beta_3 = \frac{1}{42}, \beta_4 = \frac{1}{30}$.

Wir definieren nun neue Funktionen

$$S_k(x) := B_k(x-i) \text{ für } x \in E_i := \{x | i < x < i+1\}, i \leq k. \quad (5.19)$$

Diese Funktionen sind periodisch mit Periode 1 und durch die folgende Rekursion bestimmt.

$$S_k(x) = S_k(0) + k \int_0^x S_{k-1}(t) dt, \quad k \geq 2 \quad (5.20)$$

S_1 ist stückweise linear mit Sprungstellen bei $i \in \mathbb{N}$ und Nullstellen bei $i + \frac{1}{2}$ und es gilt wegen Lemma 5.18

$$\begin{aligned} B_k(0) &= S_k(0) = S_k(1) = S_k(2) = \dots \\ S_{2k+1}(0) &= 0, S_{2k}(0) = (-1)^{k+1} \beta_k. \end{aligned} \quad (5.21)$$

Damit können wir nun die Euler–MacLaurin'sche Summenformel beweisen.

Satz 5.22 Für $f \in \mathcal{C}^{2m+2}[a, b]$ hat $T(h)$ die Entwicklung

$$T(h) = \tau_0 + \tau_1 h^2 + \tau_2 h^4 + \dots + \tau_m h^{2m} + \alpha_{m+1}(h) h^{2m+2} \quad (5.23)$$

mit $\tau_0 := \int_a^b f(x) dx$. Dabei sind die τ_i von h unabhängige Konstanten und $|\alpha_{m+1}(h)| \leq M$ für alle $h = \frac{b-a}{n}, n \in \mathbb{N}$.

Beweis: Betrachte $g(t) := f(a + th)$. Damit können wir nun das Integral

$$\int_0^n S_1(t) g'(t) dt$$

durch partielle Integration umformen. $S_1(x)$ stückweise linear aus $B_1(x) = x - \frac{1}{2}$.

$$\int_0^1 S_1(t) g'(t) dt = \int_0^1 B_1(t) g'(t) dt = B_1(t) g(t) \Big|_0^1 - \int_0^1 B_0(t) g(t) dt = \frac{1}{2} [g(1) + g(0)] - \int_0^1 g(t) dt.$$

Analog

$$\begin{aligned} \int_i^{i+1} S_1(t) g'(t) dt &= \int_0^1 B_1(t) g'(t+i) dt \\ &= \frac{1}{2} [g(i+1) + g(i)] - \int_i^{i+1} g(t) dt, \quad i = 0, \dots, n-1 \\ \implies \frac{g(0)}{2} + g(1) + \dots + g(n-1) + \frac{g(n)}{2} - \int_0^n g(t) dt &= \int_0^n S_1(t) g'(t) dt \end{aligned} \quad (5.24)$$

Für die rechte Seite gilt

$$\begin{aligned} \int_0^n S_1(t) g'(t) dt &= \frac{1}{2!} \left(S_2(t) g'(t) \Big|_0^n - \int_0^n S_2(t) g''(t) dt \right) \\ &= \frac{\beta_1}{2!} [g'(n) - g'(0)] - \frac{1}{2!} \int_0^n S_2(t) g''(t) dt \end{aligned}$$

$$\begin{aligned}
&= \frac{\beta_1}{2!} [g'(n) - g'(0)] - \frac{1}{3!} S_3(t) g''(t) \Big|_0^n + \frac{1}{3!} \int_0^n S_3(t) g^{(3)}(t) dt \\
&= \frac{\beta_1}{2!} [g'(n) - g'(0)] + \frac{1}{3!} \int_0^n S_3(t) g^{(3)}(t) dt
\end{aligned}$$

$2m$ mal partielle Integration ergibt

$$\begin{aligned}
&\frac{g(0)}{2} + g(1) + \dots + g(n-1) + \frac{g(n)}{2} - \int_0^n g(t) dt \\
&= \sum_{k=1}^m (-1)^{k+1} \frac{\beta_k}{(2k)!} [g^{(2k-1)}(n) - g^{(2k-1)}(0)] + R_{m+1}
\end{aligned} \tag{5.25}$$

mit

$$\begin{aligned}
R_{m+1} &= \frac{1}{(2m+1)!} \int_0^n S_{2m+1}(t) g^{(2m+1)}(t) dt \\
&= \frac{1}{(2m+2)!} (S_{2m+2}(t) - S_{2m+2}(0)) g^{(2m+1)}(t) \Big|_0^n \\
&\quad - \frac{1}{(2m+2)!} \int_0^n [S_{2m+2}(t) - S_{2m+2}(0)] g^{(2m+2)}(t) dt \\
&= -\frac{1}{(2m+2)!} \int_0^n [S_{2m+2}(t) - S_{2m+2}(0)] g^{(2m+2)}(t) dt
\end{aligned} \tag{5.26}$$

Es ist $\int_0^n g(t) dt = \frac{1}{h} \int_a^b f(x) dx$ und $g^{(k)}(t) = h^k f^{(k)}(a+th)$ $k = 0, 1, 2, \dots$

$$\begin{aligned}
\Rightarrow \frac{g(0)}{2} + g(1) + \dots + g(n-1) + \frac{g(n)}{2} &= \frac{f(a)}{2} + f(a+h) + \dots + f(b-h) + \frac{f(b)}{2} \\
&= \frac{1}{h} T(h).
\end{aligned}$$

Also folgt

$$T(h) = \tau_0 + \tau_1 h^2 + \dots + \tau_m h^{2m} + \alpha_{m+1}(h) h^{2m+2} \tag{5.27}$$

mit $\tau_0 = \int_a^b f(x) dx$.

$$\tau_k = \frac{(-1)^{k+1}}{(2k)!} \beta_k [f^{(2k-1)}(b) - f^{(2k-1)}(a)] \tag{5.28}$$

$$\alpha_{m+1} = \frac{-1}{(2m+2)!} \int_a^b f^{(2m+2)}(x) \left[S_{2m+2} \left(\frac{x-a}{n} \right) - S_{2m+2}(0) \right] dx \tag{5.29}$$

Da S_{2m+2} stetig und periodisch, folgt, daß es eine von h unabhängige Schranke für $|\alpha_{m+1}^{(n)}|$ gibt. \square

Die Euler–MacLaurinformel läßt sich nun zur Extrapolation verwenden (Romberg/Bulisch).

Wenn man das Restglied vernachlässigt, ist $T(h)$ ein Polynom in h^2 , das für $h = 0$ den Wert $\tau_0 = \int_a^b f(x) dx$ liefert.

Idee: Bestimme für verschiedene Schrittweiten

$h_0 = b - a, h_1 = \frac{h_0}{n_1}, \dots, h_m = \frac{h_0}{n_m}, n_i \in \mathbb{N}$, bestimme die zugehörige Trapezsumme

$$T_{i,0} = T(h_i), \quad i = 0, \dots, m \quad (5.30)$$

und dann durch Interpolation das Polynom

$$\tilde{T}_{m,m}(h) = a_0 + a_1 h^2 + \dots + a_m h^{2m}, \quad (5.31)$$

daß die Interpolationsbedingungen

$$\tilde{T}_{m,m}(h_i) = T(h_i), \quad i = 0, \dots, m \quad (5.32)$$

erfüllt. Werte dann dies Polynom bei $h = 0$ mit dem Neville-Schema aus, dies liefert i.a. sehr gute Näherung für Integral.

Sei $\tilde{T}_{i,k}(h)$ das Polynom vom Grad k in h^2 , daß

$$\tilde{T}_{i,k}(h_j) = T_{j,0}, \quad j = i - k, \dots, i \quad (5.33)$$

erfüllt. Damit gilt für die extrapolierten Werte $T_{i,k} := \tilde{T}_{i,k}(0)$.

$$T_{i,k} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\left(\frac{h_{i-k}}{h_i}\right)^2 - 1}, \quad 1 \leq k \leq i \leq m \quad (5.34)$$

→ Übung 5.42 $\left\{ \begin{array}{l} h_0 = b - a, h_1 = \frac{h_0}{2}, \dots \implies T_{11} \text{ Simpsonregel.} \\ h_0 = b - a, h_1 = \frac{h_0}{3}, \dots \implies T_{11} \frac{3}{8} \text{ - Regel.} \end{array} \right.$

Extrapolationsfolgen

a) $h_0 = b - a, h_1 = \frac{h_0}{2}, h_2 = \frac{h_1}{2}, h_3 = \frac{h_2}{2}, \dots$ (Romberg Folge)

b) $h_0 = b - a, h_1 = \frac{h_0}{2}, h_2 = \frac{h_0}{3}, h_3 = \frac{h_1}{2}, h_4 = \frac{h_2}{2}, h_5 = \frac{h_3}{2}, \dots$ (Bulirschfolge).

b) besser, da Rechenaufwand weniger schnell ansteigt. Vorherige Werte können verwendet werden. → Übung 5.42.

Fehlerbetrachtung: Für ein $\xi \in [a, b]$

$$T_{mm} - \int_a^b f(x) dx = -\frac{1}{(2m+2)!} f^{(2m+2)}(\xi) \int_a^b K(x) dx \quad (5.35)$$

mit

$$K(x) = \sum_{i=0}^m c_{m_i} h_i^{2m+2} \left[S_{2m+2} \left(\frac{x-a}{h_i} \right) - S_{2m+2}(0) \right] \quad (5.36)$$

$$\implies T_{mm} - \int_a^b f(x) dx = (b-a) h_0^2 \dots h_m^2 \frac{\beta_{m+1}}{(2m+1)!} f^{(2m+2)}(\xi) \quad (5.37)$$

Details siehe z.B. Stoer, *Numerische Mathematik*.

5.3 Übungen zu Kapitel 5

Übung 5.38 Bestimme für eine auf dem Intervall $[a, b]$ stetig differenzierbare Funktion f eine Quadraturformel für $\int_a^b f(t) dt$, indem Du Stützstellen $a = x_0 < x_1 < \dots < x_n = b$ einführst und auf dem Intervall $[x_i, x_{i+1}]$ die Funktion f durch die (stückweise) Hermite-Interpolierende aus Übung 12, Nr.2 ersetzt (neben $f(x_0), \dots, f(x_n)$ sollen auch $f'(x_0), \dots, f'(x_n)$ gegeben sein).

Übung 5.39 Sei $f : [a, b] \rightarrow \mathbb{R}$ eine zweimal stetig differenzierbare Funktion.

a) Sei F Stammfunktion zu f , d.h. $F' = f$. Seien $x, h \in \mathbb{R}$ mit $[x - h, x + h] \subseteq [a, b]$. Zeige:

$$F(x + \frac{h}{2}) - F(x - \frac{h}{2}) = hf(x) + \frac{h^3}{24} \frac{f''(\theta_+) + f''(\theta_-)}{2}$$

für geeignete zu wählende $\theta_+, \theta_- \in [a, b]$.

b) Zeige: Sei $l \in \mathbb{N} \setminus \{0\}$, $\theta_1, \dots, \theta_l \in [a, b]$. Dann existiert ein $\theta \in [a, b]$, so daß

$$f''(\theta) = \frac{1}{l} \sum_{k=1}^l f''(\theta_k)$$

ist.

c) Sei $z_i = a + ih$, $i = 0, \dots, N$, $h = \frac{b-a}{N}$. Zeige:

$$\int_{z_i}^{z_{i+1}} f(t) dt - hf\left(\frac{z_i + z_{i+1}}{2}\right) = \frac{h^3}{24} f''(\theta_i)$$

für ein geeignetes $\theta_i \in [a, b]$, $i = 0, \dots, N - 1$.

d) Sei $M(h) = h \sum_{i=0}^{N-1} f\left(\frac{z_i + z_{i+1}}{2}\right)$. (Summierte Mittelpunkregel) Zeige:

$$\int_a^b f(t) dt - M(h) = (b-a) \frac{h^2}{24} f''(\theta),$$

wobei $\theta \in [a, b]$ geeignet zu wählen ist.

Übung 5.40 Bestimme für eine auf dem Intervall $[a, b]$ stetig differenzierbare Funktion f eine Quadraturformel für $\int_a^b f(t) dt$, indem Du Stützstellen $a = x_0 < x_1 < \dots < x_n = b$ einführst und f durch den kubischen Spline S ersetzt (neben $f(x_0), \dots, f(x_n)$ sollen auch $f'(x_0)$ und $f'(x_n)$ gegeben sein). Tipp: Interpretiere S als Hermite-Interpolierende, bei der $s_i = S'(x_i)$, $0 < i < n$ als Unbekannte auftritt.

Übung 5.41 Verwende zur Berechnung der Integrale

$$\int_0^{\frac{\pi}{2}} \sin x dx, \int_0^1 \sqrt{x} dx, \int_0^1 x\sqrt{x} dx, \int_0^1 x^2\sqrt{x} dx$$

a) die summierte Trapezregel für $N = 1, 2, 4, 8, 16, \dots, 2^9$ und folgende Unterteilung:

$$x_i = i \cdot h, \quad i = 0, 1, \dots, N, \quad h = \frac{\pi}{2N}.$$

b) die summierte Mittelpunkregel für $N = 1, 2, 4, 8, 16, \dots, 2^9$ und der Unterteilung aus a).

c) die summierte Simpsonregel für $N = 2, 4, 8, 16, \dots, 2^9$ und der Unterteilung aus a).

Vergleiche die Verfahren bezüglich Genauigkeit und Rechenaufwand.

Übung 5.42 Verwende zur Berechnung des Integrals aus Aufgabe 5.41 das Extrapolationsverfahren zur summierten Trapezregel für $k = 0, 1, 2, \dots, 9$ mit der Rombergfolge

$$h_l = \frac{\pi}{2^{l+1}}, \quad l = 0, \dots, k.$$

Untersuche das Verfahren bezüglich Genauigkeit und Rechenzeit.

Übung 5.43 Zu einer gegebenen stetigen Funktion $p : [-1, 1] \rightarrow \mathbb{R}$ betrachte folgende Quadraturformel

$$\int_{-1}^1 p(x) \, dx \approx p\left(-\frac{1}{\sqrt{3}}\right) + p\left(\frac{1}{\sqrt{3}}\right).$$

Zeige: Die Formel ist für Polynome p vom Grade höchstens 3 exakt.

Tipp: Zerlege p als $p(x) = L(x)q(x) + r(x)$, wobei $L(x)$ ein quadratisches Polynom ist, welches $\pm \frac{1}{\sqrt{3}}$ als Nullstelle hat und r ein höchstens lineares Polynom ist.

Übung 5.44 Gegeben sei eine Quadraturformel $\int_{-1}^1 f(x) \, dx = Q(f) + R_Q(f)$, wobei $Q(f) = \sum_{i=1}^s \omega_i f(\xi_i)$, $-1 \leq \xi_1 < \dots < \xi_s \leq 1$ und $\omega_i \geq 0$, $\sum_{i=1}^s \omega_i = 2$. Sei ausserdem eine stetige $F : [-1, 1]^2 \rightarrow \mathbb{R}$ gegeben.

1. Konstruiere mit Hilfe der Quadraturformel Q eine Kubaturformel K , d.h. eine Vorschrift zur Berechnung von

$$\int_{-1}^1 \int_{-1}^1 F(x, y) \, dx \, dy = K(F) + R_K(F),$$

wobei $R_K(F)$ zugehöriger Approximationsfehler.

2. Zeige, dass für den Fehler $R_K(F)$ der Kubaturformel die folgende Abschätzung gilt.

$$R_K(F) \leq 2 \left(\max_{x \in [-1, 1]} |R_Q(F(x, \bullet))| + \max_{y \in [-1, 1]} |R_Q(F(\bullet, y))| \right).$$

Kapitel 6

Eigenwertprobleme

6.1 Grundlagen (Wiederholung)

Eigenwerte einer Matrix $A = [a_{ij}]_{i,j=1,\dots,n} \in \mathbb{C}^{n,n}$ sind die n Nullstellen $\lambda_1, \dots, \lambda_n$ des charakteristischen Polynoms. $p(z) = \det(zI - A)$.

$\sigma(A) = \{\lambda_1, \dots, \lambda_n\} = \underline{\text{Spektrum}}$.

$\det(A) = \lambda_1 \cdots \lambda_n$, $\underline{\text{Spur}}(A) = \sum_{i=1}^n \lambda_i = \sum_{i=1}^n a_{i,i}$.

Sei $\lambda \in \sigma(A)$, dann heißen Vektoren $x \in \mathbb{C}^n \setminus \{0\}$ die $Ax = \lambda x$ erfüllen, Eigenvektoren.

Ein Unterraum $S \in \mathbb{C}^n$ mit der Eigenschaft $x \in S \implies Ax \in S$ heißt invariant. Falls $AX = XB$ $B \in \mathbb{C}^{k,k}$, $X \in \mathbb{C}^{n,k}$, $\text{rank } X = k$, so spannen die Spalten von X einen k dimensionalen invarianten Unterraum auf und $By = \lambda y \implies AXy = \lambda Xy$.

Falls X vollen Spaltenrang hat gilt daher $\sigma(B) \subset \sigma(A)$. Falls X quadratisch, so folgt $\sigma(B) = \sigma(A)$, und $A, B = X^{-1}AX$ heißen ähnlich. X ist Ähnlichkeitstransformation

Lemma 6.1 Sei $A \in \mathbb{C}^{n,n}$, $B \in \mathbb{C}^{k,k}$, $X \in \mathbb{C}^{n,k}$, $AX = XB$ und $\text{rang}(X) = k$, dann gibt es unitäre Matrix $Q \in \mathbb{C}^{n,n}$, so daß

$$Q^*AQ = T = \begin{bmatrix} T_{11} & T_{12} \\ 0 & T_{22} \end{bmatrix} \begin{matrix} k \\ n-k \end{matrix} \quad (6.2)$$

und $\sigma(T_{11}) = \sigma(A) \cap \sigma(B)$.

Beweis: Sei $X = Q \begin{bmatrix} R \\ 0 \end{bmatrix}$ die QR-Zerlegung von X . Es folgt:

$$AQ \begin{bmatrix} R \\ 0 \end{bmatrix} = Q \begin{bmatrix} R \\ 0 \end{bmatrix} B \iff (Q^*AQ) \begin{bmatrix} R \\ 0 \end{bmatrix} = \begin{bmatrix} R \\ 0 \end{bmatrix} B =: \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix} \begin{bmatrix} R \\ 0 \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}.$$

R ist nichtsingulär da $\text{Rang}(X) = k$. Aus $T_{21}R = 0 \cdot B = 0$ folgt $T_{21} = 0$ und aus $T_{11}R = RB$ folgt $\sigma(T_{11}) = \sigma(B)$. $\sigma(A) = \sigma(T) = \sigma(T_{11}) \cup \sigma(T_{22}) \implies \text{Beh.}$ \square

Damit erhält man den folgenden Satz

Satz 6.3 (Satz von Schur (1909)) Sei $A \in \mathbb{C}^{n,n}$. Dann gibt es $Q \in \mathbb{C}^{n,n}$ unitär so daß

$$Q^*AQ = T = D + N, \quad (6.4)$$

mit $D = \text{diag}(\lambda_1, \dots, \lambda_n)$, N strikte obere Δ -Matrix. Q kann so gewählt werden, daß Eigenwerte in beliebiger Folge auf der Diagonale stehen.

Beweis: Mit Induktion $n = 1$ trivial. Induktionsvor.: Sei Beh. richtig für $n - 1$. Sei $Ax = \lambda x$, für $x \neq 0$. Mit Lemma 6.1 $B = \lambda$ gibt es unitäres U so daß

$$U^*AU = \begin{bmatrix} \lambda & w^* \\ 0 & c \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}$$

Mit Induktion gibt es \tilde{U} so daß $\tilde{U}^*C\tilde{U}$ obere Δ -Matrix.

Sei $Q = U \begin{bmatrix} T & \tilde{w}^* \\ 0 & \tilde{U} \end{bmatrix}$ dann ist Q^*AQ obere Δ -Matrix. □

Sei $Q = [q_1, \dots, q_n]$, q_i heißen Schur-Vektoren, und spannen $S_k = \text{span} \{q_1, \dots, q_k\}$ invariant auf. Eine Matrix heißt normal falls $A^*A = AA^*$.

Korollar 6.5 $A \in \mathbb{C}^{n,n}$ ist normal $\iff \exists$ unitäres $Q \in \mathbb{C}^{n,n}$, so daß $Q^*AQ = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$.

Beweis: A normal $\iff Q^*AQ$ normal $\forall Q$ unitär, denn

$$(Q^*AQ)^*(Q^*AQ) = Q^*A^*AQ = Q^*AA^*Q = (Q^*AQ)(Q^*AQ)^*.$$

$Q^*AQ = T$ obere Δ -matrix, normal. $\iff T$ diagonal. □

Es gibt eine reelle Variante des Satzes von Schur:

Satz 6.6 Sei $A \in \mathbb{R}^{n,n}$. Es gibt orthogonale Matrix $Q \in \mathbb{R}^{n,n}$ so daß

$$Q^T A Q = \begin{bmatrix} R_{11} & \cdots & R_{1m} \\ & \ddots & \vdots \\ & & R_{mm} \end{bmatrix}$$

wobei jedes R_{ii} entweder 1×1 oder 2×2 mit komplex-konjugierten Eigenwerten ist.

Beweis \implies Übung.

Satz 6.7 (Jordan'sche Normalform) Sei $A \in \mathbb{C}^{n,n}$. Dann gibt es $X \in \mathbb{C}^{n,n}$ so daß $X^{-1}AX = \text{diag}(J_1, \dots, J_t)$, wobei

$$J_i = \begin{bmatrix} \lambda_i & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{bmatrix} \in \mathbb{C}^{m_i, m_i}$$

und $m_1 + \dots + m_t = n$.

Beweis: bekannt ekelhaft. □

Die Summe der Dimensionen m_i zu einem Eigenwert λ heißt algebraische Vielfachheit, falls diese 1 ist, so heißt λ einfacher Eigenwert. Die Anzahl der Jordanblöcke zu λ heißt geometrische Vielfachheit. Es gilt immer algebraische Vielfachheit \geq geometrische Vielfachheit. Falls algebraische Vielfachheit $>$ geometrische Vielfachheit, so heißt λ defektiver Eigenwert und A defektiv.

Korollar 6.8 $A \in \mathbb{C}^{m,n}$ ist nicht defektiv genau, wenn es $X \in \mathbb{C}^{n,n}$ nicht singular gibt, so daß

$$X^{-1}AX = \text{diag}(\lambda_1, \dots, \lambda_n). \quad (6.9)$$

Beweis: Klar am Satz 6.7. Jordanform numerisch schwer zu berechnen. Kleine Störungen in der Matrix ändern Jordanform stark, numerisch ist jede Matrix diagonalisierbar, diagonalisierbare Matrizen liegen dicht in der Menge der Matrizen. Matrizen die fast defektiv sind, können schlecht konditionierte Eigenvektormatrizen haben. Und dies kann drastische Folgen haben, falls man nichtunitäre Transformationen zur Eigenwertberechnung verwendet. Denn es gilt

$$fl(X^{-1}AX) = X^{-1}AX + E \quad (6.10)$$

wobei

$$\|E\|_2 \cong \text{eps } \kappa_2(X) \|A\|_2 \quad (6.11)$$

Falls X nicht unitär ist, kann E sehr groß werden.

Beispiel 6.12

$$\begin{aligned} A &= \begin{bmatrix} 5 & -1 \\ 2 & 2 \end{bmatrix} & \sigma(A) &= \{3, 4\}, \\ Q &= \begin{bmatrix} 0.6 & 0.8 \\ -0.8 & 0.6 \end{bmatrix} & \kappa_2(Q) &= 1, \\ X &= \begin{bmatrix} 1.0 & 1.01 \\ 2. & 2.00 \end{bmatrix} & \kappa_2(X) &= 167. \end{aligned}$$

Sei $\beta = 10, t = 3$ dann ergibt sich mit Rundung

$$\begin{aligned} \|fl(X^{-1}AX) - X^{-1}AX\|_2 &\cong 1.9 \\ \|fl(Q^*AQ - Q^*AQ)\|_2 &\cong 10^{-3}. \end{aligned}$$

Satz 6.13 (Satz von Gerschgorin) Sei

$$A = D + F \text{ mit } D = \text{diag}(d_1, \dots, d_n)$$

und F habe Nulldiagonale, dann gilt

$$\sigma(A) \subseteq \bigcup_{i=1}^n D_i$$

wobei

$$D_i = \left\{ z \in \mathbb{C} : |z - d_i| \leq \sum_{j=1}^n |f_{ij}| \right\}.$$

Beweis: → Übung 6.50. □

Satz 6.14 (Satz von Bauer–Fike.) Sei μ ein Eigenwert von $A + E \in \mathbb{C}^{n,n}$ und $X^{-1}AX = D = \text{diag}(\lambda_1, \dots, \lambda_n)$ diagonal, dann gilt

$$\min_{\lambda \in \sigma(A)} |\lambda - \mu| \leq \kappa_p(X) \|E\|_p$$

Beweis: → Übung 6.51 □

Der Satz von Bauer–Fike besagt, daß $\kappa(X) \cdot \|E\|$ große Störungen in den Eigenwerten auftreten können. X orthogonal $\implies \kappa_2(X) = 1$. Resultate für invariante Unterräume im allgemeinen sehr schwer.

Satz 6.15 (Satz von Stewart) Sei Q unitär

$$Q^*AQ = \begin{bmatrix} T_{11} & T_{12} \\ 0 & T_{22} \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}$$

eine Schur–Zerlegung mit $Q = [Q_1 \quad Q_2]$. Sei $\text{sep}(T_{11}, T_{22}) = \min_{X \neq 0} \frac{\|T_{11}X - XT_{22}\|_F}{\|X\|_F}$ und sei $E \in \mathbb{C}^{n,n}$ eine Störung. Sei

$$Q^*EQ = \begin{bmatrix} E_{11} & E_{12} \\ E_{21} & E_{22} \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}.$$

Falls $\delta = \text{sep}(T_{11}, T_{22}) - \|E_{11}\|_2 - \|E_{22}\|_2 > 0$ und $\|E_{21}\|_2(\|T_{12}\|_2 + \|E_{21}\|_2) \leq \frac{\delta^2}{4}$, so gibt es $P \in \mathbb{C}^{n-k,k}$ mit $\|P\|_2 \leq 2\|E_{21}\|_2/\delta$ und die Spalten von $\hat{Q} = (Q_1 + Q_2P)(I + P^*P)^{-\frac{1}{2}}$ bilden eine Orthonormalbasis für invarianten Unterraum von $A + E$.

6.2 Der QR–Algorithmus für unsymmetrische Matrizen

Wir betrachten hier nur den reellen Fall, der Berechnung der reellen Schur–Form. Der beste Algorithmus zur Berechnung der Schur–Form ist der QR–Algorithmus von Francis, der im Prinzip die folgende Iterationsschrift ausführt.

$$H_0 = Q_0^T A Q_0$$

FOR $k = 1, 2, \dots$

$$\begin{aligned} H_{k-1} &= Q_k R_k && \text{QR–Zerlegung} \\ H_k &= R_k Q_k \end{aligned}$$

END

Diese Iteration würde pro Schritt $\mathcal{O}(n^3)$ flops brauchen (ungefähr $2n$ Iterationen) also $\mathcal{O}(n^4)$ insgesamt. Um das zu verbessern benutze Q_0 , um A auf sogenannte Hessenbergform zu transformieren.

$$Q_0^T A Q_0 = H \equiv \begin{bmatrix} \square & & \\ & \square & \\ & & \square \end{bmatrix}$$

Dies geschieht mit Hilfe von Householder-Transformation (HOUSE, COLHOUSE, ROWHOUSE)

$$\begin{aligned}
 A &= \begin{bmatrix} x & x & x & x & x & x \\ x & x & x & x & x & x \\ x & x & x & x & x & x \\ x & x & x & x & x & x \\ x & x & x & x & x & x \\ x & x & x & x & x & x \end{bmatrix} \rightarrow P_1^T A P_1 = \begin{bmatrix} x & x & x & x & x & x \\ x & x & x & x & x & x \\ 0 & x & x & x & x & x \\ 0 & x & x & x & x & x \\ 0 & x & x & x & x & x \\ 0 & x & x & x & x & x \end{bmatrix} \rightarrow \\
 P_2^T P_1^T A P_1 P_2 &= \begin{bmatrix} x & x & x & x & x & x \\ x & x & x & x & x & x \\ 0 & x & x & x & x & x \\ 0 & 0 & x & x & x & x \\ 0 & 0 & x & x & x & x \\ 0 & 0 & x & x & x & x \end{bmatrix} \rightarrow \\
 P_3^T P_2^T P_1^T A P_1 P_2 P_3 &= \begin{bmatrix} x & x & x & x & x & x \\ x & x & x & x & x & x \\ 0 & x & x & x & x & x \\ 0 & 0 & x & x & x & x \\ 0 & 0 & 0 & x & x & x \\ 0 & 0 & 0 & x & x & x \end{bmatrix} \rightarrow P_4^T P_3^T P_2^T P_1^T A P_1 P_2 P_3 P_4 = \dots
 \end{aligned}$$

Algorithmus 6.16 (Householder-Reduktion auf Hessenberg-Form)

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

Output: $H \in \mathbb{R}^{n,n}$ $H = Q^T A Q$, Q Produkt von Householder-Transformation gespeichert auf unterem Teil von H .

FOR $k = 1 : n - 2$

$v(k+1:n) = \text{HOUSE}(A(k+1:n, k));$

$A(k+1:n, k:n) = \text{ROWHOUSE}(A(k+1:n, k:n), v(k+1:n));$

$A(1:n, k+1:n) = \text{COLHOUSE}(A(1:n, k+1:n), v(k+1:n));$

$A(k+2:n, k) = v(k+2:n);$

END

Kosten: $20n^3/3$ flops + $\frac{4n^3}{3}$ flops falls Q gebraucht wird.

Fehleranalyse: $\tilde{H} = Q^T(A + E)Q$ mit $Q^T Q = I$ und $\|E\|_F \leq cn^2 \text{eps} \|A\|_F$.

→ Übung 6.55.

Die Hessenberg-Reduktion ist nicht eindeutig. Falls Hessenberg-Matrix unreduziert ist, d.h. alle Elemente in der ersten unteren Nebendiagonalen $\neq 0$ sind, so ist die Hessenreduktion fast eindeutig durch die erste Spalte von Q bestimmt, d.h. wird diese festgelegt, so ist der Rest eindeutig bis auf Vorzeichenmuster.

Satz 6.17 (Implizites Q -Theorem) Sei $A \in \mathbb{R}^{n,n}$, seien $Q = [q_1, \dots, q_n], V = [v_1, \dots, v_n]$ orthogonale Matrizen, so daß $Q^T A Q = H, V^T A V = G$ beides Hessenbergmatrizen. Sei k der kleinste Index so daß $h_{k+1,k} = 0$ ($k = n$, falls H unreduziert). Falls $v_1 = q_1$ so gilt $v_i = \pm q_i$ und $|h_{i,i-1}| = |q_{i,i-1}|$ für $i = 2, k$. Falls $k < n$ so gilt $g_{k+1,k} = 0$.

Beweis: Sei $W = V^T Q$. Dann gilt $GW = WH$. Also folgt für $i = 2 : k$

$$h_{i,i-1}w_i = Gw_{i-1} - \sum_{j=1}^{i-1} h_{j,i-1}w_j$$

Da $w_1 = e_1$ folgt, daß $[w_1, \dots, w_i]$ obere Δ -Matrix ist. $w_i = v_i^T q_i$ und $h_{i,i-1} = w_i^T G w_{i-1}$ implizieren

$v_i = \pm q_i, |h_{i,i-1}| = |g_{i,i-1}|$ für $i = 2 : k$. Falls $h_{k-1,k} = 0$ so folgt (bis auf Vorzeichen)

$$\begin{aligned} g_{k-1,k} &= e_{k+1}^T G e_k = e_{k+1}^T G W e_k = (e_{k+1}^T W)(H e_k) \\ &= e_{k+1}^T \sum_{i=1}^k h_{ik} W e_i = \sum_{i=1}^k h_{ik} e_{k+1}^T e_i = 0. \end{aligned}$$

□

⇒ Erste Spalte festlegen ergibt im wesentlichen gleiche Hessenbergreduktion (bis auf Vorzeichen).

Damit kann man nun eine praktische Variante des QR-Algorithmus machen.

```

H = Q_0^T A Q_0   Hessenberg-Reduktion
FOR   k = 1, 2
    H = QR       QR-Zerlegung
    H = RQ
END

```

H bleibt Hessenberg, denn da $H = QR$ und R obere Δ -Matrix folgt, daß $Q = HR^{-1}$ obere Hessenberg-Matrix ⇒ das neue $H = RQ$ ebendfalls obere Hessenbergmatrix.

Im folgenden kann man annehmen, daß $h_{i+1,i} \neq 0$ für die Hessenberg-Matrix. Sonst kann man das Problem zerlegen.

Falls

$$H = \left[\begin{array}{c|c} H_{11} & H_{12} \\ \hline 0 & H_{22} \end{array} \right] \begin{array}{l} p \\ n-p \end{array}$$

mit z. B.

$$|h_{p+1,p}| \leq c \text{ eps} (|h_{pp}| + |h_{p+1,p+1}|) \quad (6.18)$$

so spaltet man H in zwei Teilprobleme auf: H_{11}, H_{22} .

Beschleunigung durch Shifts

Satz 6.19 Sei μ ein Eigenwert einer unreduzierten oberen Hessenbergmatrix H . Sei $H - \mu I = QR$ QR-Zerlegung und $\tilde{H} = RQ + \mu I$ so gilt $\tilde{h}_{n,n-1} = 0, \tilde{h}_{nn} = \mu$.

Beweis: H unreduziert ⇒ $H - \mu I$ unreduziert $U^T(H - \mu I) = R$ ist singular und da man zeigen kann, daß $|r_{ii}| \geq |h_{i+1,i}| \Rightarrow r_m = 0. \Rightarrow$ letzte Zeile von \tilde{H} ist 0. □

Idee ist also mit Näherungen an Eigenwerte zu shiften z.B. mit einem Eigenwert von

$$\left[\begin{array}{cc} h_{n-1,n-1} & h_{n-1,n} \\ h_{n,n-1} & h_{n,n} \end{array} \right]$$

falls der reell ist oder mit zwei aufeinanderfolgenden komplex konjugierten Shifts $\lambda, \bar{\lambda}$. $(A - \lambda I)(A - \bar{\lambda} I)$, um komplexe Arithmetik zu vermeiden.

$$\begin{array}{l} A_i - \lambda I = Q_1 R_1 \\ A_{i+1} = R_1 Q_1 + \lambda I \\ A_{i+1} - \bar{\lambda} I = Q_2 R_2 \\ A_{i+2} = R_2 Q_2 + \bar{\lambda} I \end{array} \left| \begin{array}{l} (A_i - \lambda I)(A_i - \bar{\lambda} I) = (Q_1 Q_2)(R_2 R_1) \end{array} \right.$$

→ Übung $(A_i - \lambda I)(A_i - \bar{\lambda} I) = A_i^2 + 2\operatorname{Re}(\lambda)A_i + |\lambda|^2$ ist reell auch wenn λ nicht reell.
 Man könnte $A_i^2 - 2\operatorname{Re}(\lambda)A_i + |\lambda|^2$ berechnen und davon QR -Zerlegungen $Q_i R_i$ bilden und dann $A_{i+2} = Q_i^T A_i Q_i$. Dies ist zu teuer. Besser nutze implizites Q -Theorem.

- i) Berechne erste Spalte $z = (A_i - \lambda_i I)(A_i - \bar{\lambda}_i I)e_1$.
- ii) Bestimme Householder-Matrix P_0 , so daß $P_0 z = \alpha e_1$.
- iii) Bestimme Householder-Matrizen P_1, P_2, \dots, P_{n-2} so daß $P_{n-2}^T \dots P_0^T A_i P_0 \dots P_{n-2}$ wieder obere Hessenbergmatrix.

Algorithmus 6.20 (Francis QR -Schritt)

Input: Unreduzierte obere Hessenbergmatrix $H \in \mathbb{R}^{n,n}$, dessen untere rechte 2×2 Matrix die Eigenwerte a_1, a_2 hat.

Output: $Z^T H Z$ wobei $Z = P_1 \dots P_{n-2}$ Produkt von Householder-Matrizen ist und $Z^T (H - a_1 I)(H - a_2 I)$ obere Δ -Matrix.

$m = n - 1$

$s = H(m, m) + H(n, n)$

$t = H(m, m) * H(n, n) - H(m, n) * H(n, m)$

$x = H(1, 1) * H(1, 1) + H(1, 2) * H(2, 1) - s * H(1, 1) + t$

$y = H(2, 1) * H(1, 1) + H(2, 2) + s$

$z = H(2, 1) * H(3, 2)$

FOR $k = 0 : n - 3$

$v = \text{HOUSE}((x, y, z)^T)$

$H(k + 1 : k + 3, k + 1 : n) = \text{ROWHOUSE}(H(k + 1 : k + 3, k + 1 : n), v)$

$r = \min\{k + 4, n\}$

$H(1 : r, k + 1 : k + 3) = \text{COLHOUSE}(H(1 : r, k + 1 : k + 3), v)$

$x = H(k + 2, k + 1)$

$y = H(k + 3, k + 1)$

IF $k < n - 3$

$z = H(k + 4, k + 1)$

END

END

$v = \text{HOUSE}([x, y]^T)$

$H(n - 1 : n, n - 2 : n) = \text{ROWHOUSE}(H(n - 1 : n, n - 2 : n), v)$

$H(1 : n, n - 1 : n) = \text{COLHOUSE}(H(1 : n, n - 1 : n), v)$

Kosten $10n^2$ flops + $20n^2$ flops für Akkumulation von Q .

→ Übung 6.58

Algorithmus 6.21 (QR -Algorithmus)

Input: $A \in \mathbb{R}^{n,n}$, Toleranz $tol > \epsilon$

Output: Reelle Schurform $Q^T A T = T$. A wird mit Hessenbergform von A überschrieben. Falls Q, T gebraucht werden, so wird T in A gespeichert. Falls nur Eigenwerte gebraucht werden, so werden die Diagonal Blöcke in den entsprechenden Positionen gespeichert.

(1) Verwende Algorithmus 6.18 zur Transformation auf Hessenbergform.

UNTIL $q == n$

Setze alle Subdiagonalelemente zu Null, die $|h_{i,i-1}| \leq tol (|h_{ii}| + |h_{i-1,i-1}|)$ erfüllen.

Finde größtes $q \geq 0$ und kleinstes $p \geq 0$ so daß

$$H = \begin{bmatrix} H_{11} & H_{12} & H_{13} \\ 0 & H_{22} & H_{23} \\ 0 & 0 & H_{33} \end{bmatrix} \begin{matrix} p \\ n - p - q \\ q \end{matrix},$$

wobei H_{33} quasi obere Δ -Matrix und H_{22} unreduziert.

IF $q < n$

Mache einen Francis QR-Schnitt mit H_{22} , $H_{22} = Z^T H_{22} Z$

IF Q benötigt wird

$$Q = Q \cdot \text{diag}(I_p, Z, I_q)$$

$$H_{12} = H_{12} Z$$

$$H_{23} = Z^T H_{23}$$

END

END

END

Bringe alle 2×2 Blöcke mit reellen Eigenwerten auf obere Δ -Form und akkumulierte ggf. Transformationen

Kosten: $10n^3$ für Eigenwerte sowie $25n^3$ für Q und T

Beispiel 6.22

$$A = \begin{bmatrix} 2 & 3 & 4 & 5 & 6 \\ 4 & 4 & 5 & 6 & 7 \\ 0 & 3 & 6 & 7 & 8 \\ 0 & 0 & 2 & 8 & 9 \\ 0 & 0 & 0 & 1 & 10 \end{bmatrix}$$

Iteration	$\iota(h_{21})$	$\iota(h_{32})$	$\iota(h_{43})$	$\iota(h_{34})$
1	10^0	10^0	10^0	10^0
2	10^0	10^0	10^0	10^0
3	10^0	10^0	10^{-1}	10^0
4	10^0	10^0	10^{-3}	10^{-3}
5	10^0	10^0	10^{-6}	10^{-5}
6	10^{-1}	10^0	10^{-13}	10^{-13}
7	10^{-1}	10^0	10^{-28}	10^{-13}
8	10^{-4}	10^0	<i>conv.</i>	<i>conv.</i>
9	10^{-8}	10^0		
10	10^{-8}	10^0		
11	10^{-16}	10^0		
12	10^{-32}	10^0		
13	<i>conv.</i>	<i>conv.</i>		

Fehleranalyse: Algorithmus berechnet $\tilde{T} = Q^T(A + E)Q$

$$\tilde{Q}^T \tilde{Q} = I \quad \|E\|_2 \approx \text{eps} \|A\|_2$$

$$\tilde{Q} \text{ fast orthogonal} \quad \tilde{Q}^T \tilde{Q} = I + F \quad \|F\|_2 \approx \text{eps}$$

Was macht man nun wenn man nur bestimmte Eigenvektoren zu bestimmten Eigenwerten will, nicht den invarianten Unterraum.

Es gibt verschiedene Möglichkeiten; entweder man sortiert um oder man verwendet die sogenannte Inverse Iteration. Sei μ eine Eigenwertnäherung und sei $(A - \mu I)$ nicht exakt singulär, so ist der Eigenvektor zum größten Eigenwert von $(A - \mu I)^{-1}$ gerade der gesuchte.

Algorithmus 6.23 (Inverse Iteration) Gegeben $A \in \mathbb{C}^{n,n}$, $\mu \in \mathbb{C}$ Eigenwertnäherung und Startvektor $q^{(0)}$

FOR $k = 1, 2, \dots$

$$\text{Löse } (A - \mu I)t^{(k)} = q^{(k-1)} \quad q^{(k)} = z^{(k)} / \|z^{(k)}\|_2$$

$$\lambda^{(k)} = (q^{(k)})^T A q^{(k)}$$

END

Kann das überhaupt klappen? $A - \mu I$ ist doch fast singulär.

Angenommen A hat Basis aus Eigenvektoren $\{x_1, \dots, x_n\}$ und $Ax_i = \lambda_i x_i$, $i = 1, \dots, n$. Falls

$$q^{(0)} = \sum_{i=1}^n \beta_i x_i$$

so ist $q^{(k)}$ ein Vektor in der Richtung

$$(A - \mu I)^{-k} q^{(0)} = \sum_{i=1}^n \frac{\beta_i}{(\lambda_i - \mu)^k} x_i. \quad (6.24)$$

Das Hauptgewicht liegt also in der Richtung von x_i , der zu $\lambda_i \approx \mu$ gehört. In diese Richtung ist das Gleichungssystem auch lösbar, da x_i im Bild von $A - \mu I$ liegt.

Abbruchkriterium Sei $r^{(k)} = (A - \mu I)q^{(k)}$ das Residuum. Man bricht die Iteration ab wenn

$$\|r^{(k)}\|_\infty \leq c \cdot \text{eps} \|A\|_\infty. \quad (6.25)$$

Falls man Transformation auf Hessenbergform schon hat, so braucht jeder Schritt nur $0(n^2)$ flops.

Für symmetrische Matrizen kann man man QR genau so anwenden. Man braucht weniger Speicher und flops, symmetrische Speicherung und man erhält schnellere Konvergenz. Doppelshifts sind nicht notwendig. Viele Tricks und Varianten. Für ausführliche Inhalte siehe B. Parlett, *The Symmetric Eigenvalue Problem*.

6.3 Der QR -Algorithmus für symmetrische Matrizen

→ Übung 6.56. Reduktion symmetrischer Matrizen auf Tridiagonalgestalt.

→ Übung 6.52, 6.57. Symmetrischer QR -Algorithmus mit Wilkinson-Shift.

6.4 Die Singulärwertzerlegung

Eines der wichtigsten Hilfsmittel in allen Bereichen der numerischen Mathematik ist heutzutage die Singulärwertzerlegung einer Matrix.

Satz 6.26 Sei A eine reelle $m \times n$ Matrix. Dann gibt es orthogonale Matrizen $U = [u_1, \dots, u_m] \in \mathbb{R}^{m,m}$ und $V = [v_1, \dots, v_n] \in \mathbb{R}^{n,n}$, so daß

$$U^T A V = \text{diag}(\sigma_1, \dots, \sigma_p) \in \mathbb{R}^{m,n}, \quad p = \min(m, n) \quad (6.27)$$

wobei $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_p \geq 0$.

Beweis: Seien $x \in \mathbb{R}^n, y \in \mathbb{R}^m$ Vektoren mit $\|x\|_2 = \|y\|_2 = 1$, so daß $Ax = \sigma y, \sigma = \|A\|_2$. Dies ist möglich, da $\|A\|_2 = \sup_{\|x\|_2=1} \|Ax\|_2$ angenommen wird.

Da jede Menge von orthogonalen Vektoren zu einer Basis des ganzen Raumes ergänzt werden kann, gibt es $V_1 \in \mathbb{R}^{n,(n-1)}$ und $U_1 \in \mathbb{R}^{m,(m-1)}$, so daß

$V = [x, V_1] \in \mathbb{R}^{n,n}, U = [y, U_1] \in \mathbb{R}^{m,m}$ orthogonal sind. Dann gilt

$$\begin{aligned} U^T A V &= \begin{bmatrix} y^T \\ U_1^T \end{bmatrix} A \begin{bmatrix} x & V_1 \end{bmatrix} = \begin{bmatrix} y^T \\ U_1^T \end{bmatrix} \begin{bmatrix} \sigma y & A V_1 \end{bmatrix} \\ &= \begin{bmatrix} \sigma & w^T \\ 0 & B \end{bmatrix} =: A_1. \end{aligned} \quad (6.28)$$

Da

$$\left\| A_1 \begin{bmatrix} \sigma \\ w \end{bmatrix} \right\|_2^2 \geq (\sigma^2 + w^T w)^2, \quad (6.29)$$

so folgt, daß

$$\|A_1\|_2^2 \geq \sigma^2 + w^T w. \quad (6.30)$$

Aber $\sigma^2 = \|A\|_2^2 = \|A_1\|_2^2$. Also folgt $w = 0$. Per Induktion folgt das Resultat. \square

Die Werte σ_i heißen Singularwerte, die u_i, v_i linke und rechte Singulärvektoren. Es gilt

$$A v_i = \sigma_i u_i, \quad A^T u_i = \sigma_i v_i, \quad i = 1, \dots, \min(m, n) \quad (6.31)$$

→ Übung

Die Singulärwerte sind die Längen der Halbachsen der Hyperellipsoide $E = \{Ax : \|x\| = 1\}$.

→ Übung

Sei $\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$, so gilt

$$\begin{aligned} \text{rang}(A) &= r \\ \text{Kern}(A) &= \text{span}\{v_{r+1}, \dots, v_n\} \\ \text{Bild}(A) &= \text{span}\{u_1, \dots, u_r\} \\ \|A\|_p^2 &= \sigma_1^2 + \dots + \sigma_p^2 \quad p = \min\{m, n\} \\ \|A\|_2 &= \sigma_1 \\ \text{cond}_2(A) &= \frac{\sigma_1}{\sigma_p} \end{aligned}$$

Die Singulärwertzerlegung wird daher zur Rangbestimmung verwendet. Weitere Anwendungen, Bildverarbeitung, Ausgleichsprobleme, Abspaltung von Kern und Co-Kern.

Satz 6.32 Sei $A = U\Sigma V^T$ die SWZ von $A \in \mathbb{R}^{m,n}$. Sei $k < r = \text{rang}(A)$ und

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T. \quad (6.33)$$

Dann gilt

$$\min_{\text{rang}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1}. \quad (6.34)$$

Beweis: $U^T A_k V = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0) \implies \text{rang}(A_k) = k$ und

$U^T(A - A_k)V = \text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_p) \implies \|A - A_k\|_2 = \sigma_{k+1}$.

Angenommen $\text{rang}(B) = k$ für $B \in \mathbb{R}^{m,n}$. Dann gibt es orthonormale Vektoren x_1, \dots, x_{n-k} die Kern (B) aufspannen. Kern $(B) = \text{span}\{x_1, \dots, x_{n-k}\}$. Aus Dimensionsgründen gilt

$$w = \text{span}\{x_1, \dots, x_{n-k}\} \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \{0\} \quad (6.35)$$

($z = \sum \alpha_i v_i \quad \alpha_i = v_i^T z$).

Sei $z \in W, \|z\|_2 = 1$. Da $Bz = 0$ und

$Az = \sum_{i=1}^{k+1} \sigma_i (v_i^T z) u_i$, so folgt

$$\|A - B\|_2^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 = \sum_{i=1}^{k+1} \sigma_i^2 (v_i^T z)^2 \geq \sigma_{k+1}^2. \quad (6.36)$$

6.4.1 Die Berechnung der Singulärwertzerlegung

Man könnte die Berechnung der Singulärwertzerlegung direkt auf ein symmetrisches Eigenwertproblem zurückführen, denn falls $A = U\Sigma V^T \in \mathbb{R}^{m,n}$, $m \geq n$, so gilt

$AA^T = U\Sigma V^T V \Sigma^T U^T = U\Sigma^2 U^T = U \text{diag}(\sigma_1^2, \dots, \sigma_n^2) U^T$ und da AA^T symmetrisch ist dies eine Schurform mit spezieller Anordnung. Analog

$A^T A = V \Sigma^2 V^T = V \text{diag}(\sigma_1^2, \dots, \sigma_n^2, 0, \dots, 0) V^T$. Man könnte also $A^T A$ und AA^T bilden, aber wie wir schon beim Ausgleichsproblem gesehen haben, würde dieses zur Quadrierung der Kondition führen.

Besser ist daher der Algorithmus von Golub/Kahan der auf den symmetrischen QR-Algorithmus aufbaut.

(1) Schritt transformiere A auf Bidiagonalform mit Householder-Transformationen.

Algorithmus 6.37 (Householder-Bidiagonalisierung)

Input: $A \in \mathbb{R}^{m,n}, m \geq n$

Output: Der obere Bidiagonalteil von A wird dem oberen Bidiagonalteil von $B = U^T A V$ überschrieben, B obere Bidiagonalmatrix.

$$U = U_1 \cdots U_{n-1} (\cdot U_n), \quad V = V_1 \cdots V_{n-2}$$

Produkt von Householder-Matrizen. Der wesentliche Teil von U_j ist in $A(j+1 : m, j)$ und der wesentliche Teil von V_j in $A(j, j+2 : m)$ gespeichert.

FOR $j = 1 : n - 1$

$v(j : m) = \text{HOUSE}(A(j : m, j));$

$A(j : m, j : n) = \text{ROWHOUSE}(A(j : m, j : n), v(j : m));$

$A(j+1 : m, j) = v(j+1 : m);$

IF $j \leq n - 2$

$v(j+1 : n) = \text{HOUSE}(A(j, j+1 : n)^T);$

$A(j : m, j+1 : n) = \text{COLHOUSE}(A(j : m, j+1 : n), v(j+1 : n));$

$A(j, j+2 : n) = v(j+2 : n)^T;$

END

END

IF $m > n$

$v(n : m) = \text{HOUSE}(A(n : m, n));$

$A(n : m, n) = \text{ROWHOUSE}(A(n : m, n), v(n : m));$

$A(n+1 : m, n) = v(n+1 : m);$

END

Kosten: $4mn^2 - 4n^3/3$ flops

Falls U, V explizit benötigt werden, diese erhält man mit $4m^2n - 4n^3/3$ für U und $4n^3/3$ für V .

Fehleranalyse: $B = U^T(A + E)V$ mit $U^T U = I, V^T V = I, \|E\|_F \leq cn^2 \text{ eps } \|A\|_F$.

Vorgehen:

$$\begin{aligned} & \begin{bmatrix} x & x & x & x \\ x & x & x & x \end{bmatrix} \xrightarrow{U_1} \begin{bmatrix} x & x & x & x \\ 0 & x & x & x \end{bmatrix} \xrightarrow{V_1} \begin{bmatrix} x & x & 0 & 0 \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \xrightarrow{V_2} \begin{bmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{bmatrix} \\ & \xrightarrow{U_3} \begin{bmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \\ 0 & 0 & 0 & x \end{bmatrix} \xrightarrow{U_4} \begin{bmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \\ 0 & 0 & 0 & 0 \end{bmatrix}. \end{aligned}$$

(2) Wende symmetrischen QR-Algorithmus implizit auf $T = B^T B$ (tridiagonal) an,

$$B = \begin{bmatrix} d_1 & f_1 & & & \\ & \ddots & \ddots & & \\ & & \ddots & f_{n-1} & \\ \hline & & & & d_n \\ & & & & 0 \end{bmatrix}, \quad B^T B = \begin{bmatrix} d_1^2 & f_1 d_1 & & & \\ f_1 d_1 & d_2^2 + f_1^2 & f_2 d_2 & & \\ \ddots & \ddots & \ddots & & \\ & & & f_{n-1} d_{n-1} & f_{n-1} d_{n-1} \\ & & & & d_n^2 + f_{n-1}^2 \end{bmatrix}$$

Berechne Eigenwert λ von

$$T(n-1:n, n-1:n) = \begin{bmatrix} d_{n-1}^2 + f_{n-2}^2 & d_{n-1} f_{n-1} \\ d_{n-1} f_{n-1} & d_n^2 + f_{n-1}^2 \end{bmatrix}$$

der näher an $d_n^2 + f_{n-1}^2$ liegt.

Berechne die beiden Elemente von

$$(T - \lambda I)e_1 = \begin{bmatrix} d_1^2 - \lambda \\ d_1 f_1 \end{bmatrix}$$

und eine Givensrotation $G = \begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}$

$$\begin{bmatrix} c_1 & s_1 \\ -s_1 & c_1 \end{bmatrix}^T \begin{bmatrix} d_1^2 - \lambda \\ d_1 f_1 \end{bmatrix} = \begin{bmatrix} * \\ 0 \end{bmatrix}$$

und transformiere

$$BG_1 = \begin{bmatrix} * & * & & & \\ \oplus & * & * & & \\ & & \ddots & \ddots & \\ & & & \ddots & * \\ & & & & * \end{bmatrix}$$

wieder auf Bidiagonalform.

→ Übung

Berechnung und Anwendung von Givensrotationen

Bestimme $c, s \in \mathbb{R}$ mit $c^2 + s^2 = 1$, so daß

$$\begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} * \\ 0 \end{bmatrix}.$$

Wichtige Details:

Stabilität der Berechnung von c, s . Im Prinzip ist $c = \pm b/\sqrt{a^2 + b^2}, s = \mp a/\sqrt{a^2 + b^2}$. Man berechnet aber den Quotienten $t = a/b$ bzw. $t = b/a$ und macht Fallunterscheidung. Aus Stabilitätgründen sollte $|t| \leq 1$ gelten.

Algorithmus 6.38 (Erzeugung der Givens-Parameter)

Input: $a, b \in \mathbb{R}$.

Output: $c, s \in \mathbb{R}$ mit $c^2 + s^2 = 1$, so daß $\begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} * \\ 0 \end{bmatrix}$.

```

FUNCTION   [c, s] = GIVENS(a, b)
IF         |b| > |a|
            t = -a/b;
            s = 1/√(1 + t²);
            c = s * t;
ELSE
            t = -b/a;
            c = 1/√(1 + t²);
            s = c * t;
END
END GIVENS

```

Kosten: 5 flops, 1 $\sqrt{\quad}$.

Algorithmus 6.39 (Vormultiplikation mit Givens-Rotation)

Input: $A \in \mathbb{K}^{2,n}$, $c, s \in \mathbb{R}$, $c^2 + s^2 = 1$.

Output: A überschrieben mit $\begin{bmatrix} c & -s \\ s & c \end{bmatrix} A$.

```

FUNCTION   A = ROWROT(A, c, s)
            A = [c, -s; s, c] * A;
END ROWROT

```

Kosten: $6n$ flops.

Fehler: $gl(\tilde{G}A) = G(A + E)$, $\|E\|_2 = \mathcal{O}(\text{eps} \|A\|_2)$

Algorithmus 6.40 (Nachmultiplikation mit Givens-Rotation)

Input: $A \in \mathbb{K}^{m,2}$, $c, s \in \mathbb{R}$, $c^2 + s^2 = 1$.

Output: A überschrieben mit $A \begin{bmatrix} c & s \\ -s & c \end{bmatrix}$.

```

FUNCTION   A = COLROT(A, c, s)
            A = A * [c, s; -s, c];
END COLROT

```

Kosten: $6m$ flops.

Fehler: $gl(A\tilde{G}) = (A + E)G$, $\|E\|_2 = \mathcal{O}(\text{eps} \|A\|_2)$

Diese 3 Algorithmen werden verwendet um bestimmte einzelne Einträge einer Matrix zu eliminieren.

Dieses ist ideal für den SWZ-Schritt.

Algorithmus 6.41 (Golub/Kahan SWZ–Schritt)

Input: Bidiagonalmatrix $B \in \mathbb{R}^{m,n}$, die keine Nullen auf Diagonale und Superdiagonale hat.

Output: B überschrieben mit der Bidiagonalmatrix $\bar{B} = \bar{U}^T B \bar{V}$, \bar{U}, \bar{V} orthogonal und \bar{V} ist im wesentlichen die orthogonale Matrix die man erhalten würde, wenn man die symmetrische Version des QR–Algorithmus auf $T = B^T B$ anwenden würde.

Sei μ der Eigenwert der rechten unteren 2×2 Matrix von $B^T B$ der näher am letzten Diagonalelement von $B^T B$ ist.

$$y = B(1, 1)^2 - \mu$$

$$z = B(1, 1) * B(1, 2)$$

FOR $k = 1 : n - 1$

$$[c, s] = \text{GIVENS}(y, z);$$

$$l = \max\{1, k - 1\}$$

$$B(l : k + 1, k : k + 1) = \text{COLROT}(B(l : k + 1, k : k + 1), c, s);$$

$$y = B(k, k);$$

$$z = B(k + 1, k);$$

$$[c, s] = \text{GIVENS}(y, z);$$

$$l = \min\{k + 2, n\}$$

$$B(k : k + 1, k : l) = \text{ROWROT}(B(k : k + 1, k : l), c, s);$$

IF $k < n - 1$

$$y = B(k, k + 1);$$

$$z = B(k, k + 2);$$

END

END

Kosten: $30n$ flops, $2n \sqrt{\quad}$

Akkumulation von U $6mn$ flops

Akkumulation von V $6n^2$ flops

Fehler: Wie bei QR.

Deflation wie bei QR–Algorithmus

Die Anfangsreduktion auf Bidiagonalform und die Iterationsfolge zusammen mit Deflation wird zusammengefaßt im SWZ–Algorithmus.

Algorithmus 6.42 (SWZ Berechnung–Golub/Reinsch)

Input: $A \in \mathbb{R}^{m,n}$, $m \geq n$, $\varepsilon = c \cdot \text{eps}$, c klein.

Output: A überschrieben mit $U^T AV = D + E$ U, V orthogonal, D diagonal, $\|E\|_2 \approx \text{eps} \|A\|_2$.

Bidiagonalisiere A mit Algorithmus 6.37. $q = 0$;

WHILE $q < n$

Setze $A(i, i+1)$ zu Null, falls

$|A(i, i+1)| \leq \varepsilon(|A(i, i)| + |A(i+1, i+1)|)$, $i = 1 : n-1$.

Bestimme größtes q und kleinstes p , so daß in

$$B = \begin{bmatrix} B_{11} & 0 & 0 \\ 0 & B_{22} & 0 \\ 0 & 0 & B_{33} \end{bmatrix} \begin{matrix} p \\ n-p-q \\ q+m-n \end{matrix}$$

$p \quad n-p-q \quad q$

B_{33} diagonal ist, und B_{22} nur Nebendiagonalelemente $\neq 0$ hat.

IF $q < n$

Setze in B_{22} alle $A(i, i)$ zu Null, falls $|A(i, i)| \leq \varepsilon \|B_{22}\|$.

IF Diagonalelement von B_{22} Null

 mache Spezialschritt (unten)

END

 Wende Algorithmus 6.41 auf B_{22} an

END

END

Spezialschritt:

Sei k der größte Index, so daß für die Einträge in B_{22} $A(k, k) = 0$ ist und $A(l, l) \neq 0$, $l > k$.

Schiebe Eintrag an Position $(k, k+1)$ nach rechts heraus

FOR $j = k+1 : n-p-q$

 Eliminiere Eintrag in Position (k, j)

$[c, s] = \text{GIVENS}(A(j, j), A(k, j))$;

$A([j, k], j) = \text{ROWROT}(A([j, k], j), c, s)$;

 Neuer Eintrag in Position $(k, j+1)$ kommt herein

 IF $j < n-p-q$

$A([j, k], j+1) = \text{ROWROT}(A([j, k], j+1), c, s)$;

 END

END

Schiebe Eintrag an Position $(k-1, k)$ nach oben heraus

FOR $j = k-1 : -1 : p+1$

 Eliminiere Eintrag in Position (j, k)

$[c, s] = \text{GIVENS}(A(j, j), A(j, k))$;

$A(j, [j, k]) = \text{COLROT}(A(j, [j, k]), c, s)$;

 Neuer Eintrag in Position $(j-1, k)$ kommt herein

 IF $j > p+1$

$A(j-1, [j, k]) = \text{COLROT}(A(j-1, [j, k]), c, s)$;

 END

END

	flops
Σ	$4mn^2 - 4n^3/3$
<u>Koten:</u> Σ, V	$4mn^2 + 8n^3$
Σ, U	$4m^2n - 8mn^2$
Σ, U, V	$4m^2n + 8mn^2 + 9n^3$

Den Spezielschritt verdeutlichen wir an folgendem Beispiel
Zuerst schieben wir den Eintrag $(k, k + 1)$ nach rechts heraus

$$\left[\begin{array}{cc|cc} * & * & & \\ & * & \otimes & \\ \hline & & 0 & \otimes \\ \hline & & & * & * \\ & & & & * \end{array} \right] \mapsto \left[\begin{array}{cc|cc} * & * & & \\ & * & \otimes & \\ \hline & & 0 & 0 & \otimes \\ \hline & & & \otimes & \otimes \\ & & & & * \end{array} \right] \mapsto \left[\begin{array}{cc|cc} * & * & & \\ & * & * & \\ \hline & & 0 & 0 & 0 \\ \hline & & & \otimes & \otimes \\ & & & & \otimes \end{array} \right]$$

Dann schieben wir den Eintrag $(k - 1, k)$ nach oben heraus

$$\left[\begin{array}{cc|cc} * & * & & \\ & * & \otimes & \\ \hline & & 0 & 0 & 0 \\ \hline & & & \otimes & \otimes \\ & & & & \otimes \end{array} \right] \mapsto \left[\begin{array}{cc|cc} * & \otimes & \otimes & \\ & \otimes & 0 & \\ \hline & & 0 & 0 & 0 \\ \hline & & & \otimes & \otimes \\ & & & & \otimes \end{array} \right] \mapsto \left[\begin{array}{cc|cc} \otimes & \otimes & 0 & \\ & \otimes & 0 & \\ \hline & & 0 & 0 & 0 \\ \hline & & & \otimes & \otimes \\ & & & & \otimes \end{array} \right]$$

Beispiel 6.43

	$\mathcal{O}(a_{21})$	$\mathcal{O}(a_{23})$	$\mathcal{O}(a_{34})$
1	10^0	10^0	10^0
2	10^0	10^0	10^0
3	10^0	10^0	10^0
4	10^0	10^{-1}	10^{-2}
5	10^0	10^{-3}	10^{-8}
6	10^0	10^{-1}	10^{-27}
7	10^0	10^{-1}	–
8	10^0	10^{-4}	–
9	10^{-1}	10^{-14}	–
10	10^{-1}	–	–
11	10^{-4}	–	–
12	10^{-12}	–	–
13	–	–	–

Also kubische Konvergenz.
Weiteres siehe Golub / van Loan.

6.5 Der Bisektionsverfahren zur Berechnung der Eigenwerte symmetrischer tridiagonaler Matrizen

Sei T_r die führende Hauptabschnittsmatrix von

$$T = \begin{bmatrix} a_1 & b_1 & & & \\ \bar{b}_1 & a_2 & b_2 & & \\ & \bar{b}_2 & \ddots & \ddots & \\ & & \ddots & \ddots & b_{n-1} \\ & & & \bar{b}_{n-1} & a_n \end{bmatrix} = T^* \quad (6.44)$$

und sei $p_r(x) = \det(T_r - xI)$, $r = 1 : n$ dann gilt

$$p_r(x) = (a_r - x)p_{r-1}(x) - (b_{r-1})^2 p_{r-2}(x) \quad r = 2 : n \quad (6.45)$$

→ Übung 6.59 Man kann $p_r(x)$ in $O(r)$ flops auswerten.

Algorithmus 6.46 *Bisektion zur Berechnung der Eigenwerte von T .*

Input: T wie in (6.44) symmetrisch, tridiagonal, Toleranzgrenze tol , $y, z \in \mathbb{R}, y < z, p_n(y)p_n(z) < 0$.

WHILE $|y - z| > tol (|y| + |z|)$

$x = \frac{y+z}{2}$
 IF $p_n(x)p_n(x) < 0$
 $z = x$
 ELSE
 $y = x$

END

END

Kosten: $O(h)$ flops. Fehler halbiert sich pro Schritt, lineare Konvergenz.

Parallel auswertbar.

Satz 6.47 (Sturm'sche Kette) Sei T in (6.44) unreduziert dann separieren die Eigenwerte von T_{r-1} die Eigenwerte von T_r .

$$\lambda_r(T_r) < \lambda_{r-1}(T_{r-1}) < \dots < \lambda_2(T_r) < \lambda_1(T_{r-1}) < \lambda_1(T_r) \quad (6.48)$$

Falls $a(\lambda)$ die Anzahl der Zeichenwechsel in der Folge $\{p_0(\lambda), \dots, p_n(\lambda)\}$ ist, so ist $a(x)$ die Anzahl der Eigenwerte von $T < als x .$

Beweis: z.B. Stoer □

Damit können wir einen weiteren Algorithmus zur Bisektion bekommen.

Algorithmus 6.49 (Sturm'sche Ketten, Bisektion) *Input:* T wie in (6.44), Toleranz tol und Polynomfolge wie in 6.45, $[a, b]$, so daß alle Eigenwerte von T in $[a, b]$. Sei $x = b, y = a$.

WHILE $|y - z| > tol (|y| + |z|)$

$$x = (y + z)/2$$

$$IF a(x) \geq k$$

$$z = x$$

ELSE

$$y = x$$

END

END

Man kann $y(\lambda)$ auch mittels Cholesky-Zerlegung bestimmen.

6.6 Übungen zu Kapitel 6

Übung 6.50 a) Zeige den Satz von Gerschgorin: Sei $A = (a_{ij})_{i,j=1,\dots,n} \in \mathbb{C}^{n,n}$. Dann sind die Eigenwerte von A enthalten in der Vereinigung der Kreisscheiben

$$\bigcup_{i=1}^n K_i, \quad K_i = \{z : |z - a_{ii}| \leq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|\}.$$

b) Zeige folgende verschärfte Variante des Satzes von Gerschgorin: Sind i_1, \dots, i_r (paarweise verschieden) so gewählt, daß

$$\left(\bigcup_{k=1}^r K_{i_k} \right) \cap \left(\bigcup_{\substack{i=1 \\ i \neq i_1, \dots, i_r}}^n K_i \right),$$

dann enthält $\bigcup_{k=1}^r K_{i_k}$ bereits r (nicht zwingend verschiedene) Eigenwerte von A

Hinweis: Die Eigenwerte einer Matrix hängen stetig von den Koeffizienten der Matrix ab.

Übung 6.51 Zeige den Satz von Bauer und Fike: Sei $A = (a_{ij})_{i,j=1,\dots,n} \in \mathbb{C}^{n,n}$ diagonalisierbar, d.h., nehme an, daß $A = X\Lambda X^{-1}$ ist, wobei $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ ist und X nicht singular ist. Sei $E \in \mathbb{C}^{n,n}$ irgend eine Matrix. Dann existiert zu jedem Eigenwert μ von $A + E$ ein Eigenwert λ_j von A derart, daß

$$|\lambda_j - \mu| \leq \|X\|_p \|X^{-1}\|_p \|E\|_p.$$

Dabei ist $\|\bullet\|_p$ die übliche induzierte Matrixnorm zur p -Norm für Vektoren.

Hinweis: Betrachte o.B.d.A. $\mu \neq \lambda_j$ und betrachte $(\mu I - \Lambda)X^{-1}v$ für ein geeignetes v .

Übung 6.52 Schreibe ein **MATLAB**-Programm für den QR-Algorithmus für symmetrische Matrizen einmal mit Wilkinson-Shift und einmal ohne Shift (Der Wilkinson-Shift einer symmetrischen Matrix A ist der Eigenwert von $\begin{pmatrix} a_{n-1,n-1} & a_{n-1,n} \\ a_{n-1,n} & a_{nn} \end{pmatrix}$, welcher dichter bei a_{nn} gelegen ist).

Für die Reduktion auf symmetrische Tridiagonalgestalt darf die entsprechende Routine 'hess' aus **MATLAB** verwendet werden. Für die Iteration reicht die explizite Variante.

Vergleiche beide Varianten bezüglich Rechenzeit und Genauigkeit anhand der Matrix $A =$

$$\begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & & \\ & & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix} \text{ für Dimension } n = 10, 20, \dots, 100.$$

Übung 6.53 Sei $A \in \mathbb{C}^{n,n}$ eine nichtzerfallende obere Hessenbergmatrix und $\mu \in \mathbb{C}$. Zeige:

1. Falls $A - \mu I = QR$ eine QR-Zerlegung von $A - \mu I$ ist, dann ist Q eine obere Hessenbergmatrix (d.h. $q_{ij} = 0, i > j + 1$). Ist μ ein Eigenwert von A , dann ist $r_{nn} = 0$.
2. $H = Q^H A Q$ ist eine obere Hessenbergmatrix. Ist μ ein Eigenwert von A , dann ist $h_{n,n-1} = 0$. Was bedeutet das für die letzte Spalte von Q ?
3. Nehme an A sei reell. Dann ist auch $(A - \mu I)(A - \bar{\mu} I)$ reell.
4. Sei A reell und $\mu = \alpha + i\beta$ ein Eigenwert von A . Falls $A^2 - 2\alpha A + (\alpha^2 + \beta^2)I = QR$ eine QR-Zerlegung ist, dann ist $\begin{pmatrix} r_{n-1,n-1} & r_{n-1,n} \\ 0 & r_{nn} \end{pmatrix}$ singular.

Übung 6.54 Sei $A \in \mathbb{C}^{n,n}$ obere Hessenbergmatrix und p ein Polynom. Nehme an, dass $p(A)$ nicht-singular ist. Zeige: Ist $p(A) = QR$ eine QR-Zerlegung von $p(A)$, dann ist $B = Q^H A Q$ eine obere Hessenbergmatrix. Gibt es eine Möglichkeit B zu berechnen, ohne das ganze Polynom $p(A)$ auszurechnen?

Übung 6.55 Schreibe ein **MATLAB**-Programm für die Reduktion auf obere Hessenberggestalt mittels Householder-Transformationen. Vergleiche das Programm bezüglich Genauigkeit und Rechenaufwand mit der **MATLAB**-Routine 'hess' anhand zufällig erzeugter Beispiele der Dimensionen $n = 10, 20, \dots, 100$.

Übung 6.56 Schreibe ein **MATLAB**-Programm für die Reduktion symmetrischer Matrizen auf Tridiagonalgestalt mittels Householder-Transformationen. Vergleiche das Programm bezüglich Genauigkeit und Rechenaufwand mit der **MATLAB**-Routine 'hess' anhand zufällig erzeugter Beispiele der Dimensionen $n = 10, 20, \dots, 100$.

Übung 6.57 Schreibe ein **MATLAB**-Programm für den QR-Algorithmus für symmetrische Matrizen einmal mit Wilkinson-Shift und einmal ohne Shift (Der Wilkinson-Shift einer symmetrischen Matrix A ist der Eigenwert von $\begin{pmatrix} a_{n-1,n-1} & a_{n-1,n} \\ a_{n-1,n} & a_{nn} \end{pmatrix}$, welcher dichter bei a_{nn} gelegen ist). Für die Reduktion auf symmetrische Tridiagonalgestalt darf die entsprechende Routine 'hess' aus **MATLAB** verwendet werden. Für die Iteration soll die implizite Variante verwendet werden.

Vergleiche beide Varianten bezüglich Rechenzeit und Genauigkeit anhand der Matrix $A = \begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & & \\ & & & -1 & \\ & & & -1 & 2 \end{pmatrix}$ für Dimension $n = 10, 20, \dots, 100$.

Übung 6.58 Schreibe ein **MATLAB**-Programm für den QR-Algorithmus für allgemeine reelle Matrizen mit implizitem Francis-Shift (Eigenwerte der Matrix $\begin{pmatrix} a_{n-1,n-1} & a_{n-1,n} \\ a_{n,n-1} & a_{nn} \end{pmatrix}$). Für die Reduktion auf Hessenberggestalt darf die entsprechende Routine 'hess' aus **MATLAB** verwendet werden. Vergleiche das Verfahren bezüglich Rechenzeit und Genauigkeit mit der Funktion 'schur' aus **MATLAB** anhand zufällig erzeugter Matrizen der Dimensionen $n = 10, 20, \dots, 100$.

Übung 6.59 Gegeben sei eine symmetrische Tridiagonalmatrix

$$A = \begin{pmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & & \\ & & & \beta_n & \\ & & & \beta_n & \alpha_n \end{pmatrix} \in \mathbb{R}^{n,n}.$$

Zeige: Das charakteristische Polynom $p_n(t) = \det(A - tI)$ von A lässt sich rekursiv berechnen durch

$$p_0(t) = 1$$

$$p_1(t) = \alpha_1 - t$$

$$p_j(t) = (\alpha_j - t)p_{j-1}(t) - \beta_j^2 p_{j-2}(t), \quad j = 2, 3, \dots, n.$$