

# VORLESUNGSSKRIPT EINFÜHRUNG IN DIE NUMERIK

SOMMERSEMESTER 2022

Roland Herzog\*

2022-07-18

\*Interdisciplinary Center for Scientific Computing, Heidelberg University, 69120 Heidelberg, Germany  
([roland.herzog@iwr.uni-heidelberg.de](mailto:roland.herzog@iwr.uni-heidelberg.de), <https://scoop.iwr.uni-heidelberg.de/team/roland-herzog>).

Material für 13 Wochen.

# Inhaltsverzeichnis

o	Einführung	7
§ 1	Was ist Numerische Mathematik?	7
§ 2	Wiederholung und Notation	11
§ 2.1	Vektor- und Matrix-Normen	12
§ 2.2	Differenzierbarkeit und der Satz von Taylor	17
§ 2.3	Landau-Notation	21
1	Sensitivitätsanalyse und Kondition von Aufgaben	24
§ 3	Sensitivitätsanalyse und Kondition	24
§ 3.1	Absolute Betrachtungsweise	25
§ 3.2	Relative Betrachtungsweise	29
§ 3.3	Kondition linearer Gleichungssysteme	34
§ 3.4	Kondition zusammengesetzter Funktionen	36
§ 4	Singulärwertzerlegung von Matrizen	38
§ 4.1	Spektralzerlegung	38
§ 4.2	Singulärwertzerlegung	39
2	Zahldarstellung und Rechnerarithmetik	47
§ 5	Zahldarstellung	47
§ 6	Rechnerarithmetik	52
3	Stabilität von Algorithmen	54
§ 7	Stabilitätsanalyse	54
§ 7.1	Vorwärtsanalyse	54
§ 7.2	Rückwärtsanalyse	69
§ 7.3	Allgemeine Stabilitätsanalyse	71
4	Direkte Lösungsverfahren für lineare Gleichungssysteme	72
§ 8	LR-Zerlegung und das Gaußsche Eliminationsverfahren	72
§ 9	Fehleranalyse bei der Lösung linearer Gleichungssysteme	84
§ 9.1	Fall ohne Pivotsuche	88

§ 9.2	Fall mit Spaltenpivotsuche	94
§ 10	LR-ähnliche Zerlegungen für spezielle Matrizen	96
§ 10.1	Matrizen, für die keine Pivotsuche erforderlich ist	96
§ 10.2	Cholesky-Zerlegung für symmetrische positiv definite Matrizen	96
5	Ausgleichsprobleme	100
§ 11	Einführung	100
§ 12	QR-Zerlegung	104
§ 12.1	Gram-Schmidt-Verfahren	105
§ 12.2	Givens-Rotationen	106
§ 12.3	Householder-Transformationen	109
§ 12.4	Nutzung der QR-Zerlegung	111
§ 13	Nutzung der Singulärwertzerlegung	112
6	Interpolation und Approximation	114
§ 14	Einführung	114
§ 15	Polynominterpolation	115
§ 15.1	Monombasis	116
§ 15.2	Lagrange-Polynombasis	118
§ 15.3	Newton-Polynombasis	121
§ 15.4	Approximationsfehler und Konvergenz bei der Interpolation	123
§ 15.5	Kondition der Interpolationsaufgabe	127
§ 15.6	Anwendungen der Interpolation	129
§ 16	Bernstein-Polynome und Bézier-Kurven	130
§ 17	Approximation	138
§ 17.1	Tschebyschow-Approximation	138
§ 17.2	Anwendung der Tschebyschow-Approximation zur Optimierung der Stützstellen	141
7	Numerische Integration	145
§ 18	Interpolatorische Quadraturformeln	145
§ 18.1	Abgeschlossene Newton-Cotes-Formeln	146
§ 18.2	Offene Newton-Cotes-Formeln	150
§ 19	Summierte Quadraturformeln	151
§ 20	Gauß-Quadraturformeln	153
§ 21	Quadratur basierend auf Richardson-Extrapolation	155
8	Iterative Lösungsverfahren für nichtlineare Gleichungssysteme	158

---

§ 22	Bisektionsverfahren	158
§ 23	Fixpunktverfahren	159
§ 24	Newton-Verfahren	161
§ 24.1	Einige Hilfsresultate	163
§ 24.2	Das lokale Newton-Verfahren	165
9	Iterative Lösungsverfahren für lineare Gleichungssysteme	168
§ 25	Stationäre Verfahren	169
§ 26	Instationäre Verfahren	175



# Kapitel 0 Einführung

## § 1 WAS IST NUMERISCHE MATHEMATIK?

Die **numerische Mathematik** (kurz: **Numerik**, englisch: *numerical analysis*) beschäftigt sich mit der Konstruktion und Analyse von Rechenverfahren (Algorithmen) zur Lösung von Aufgaben der «kontinuierlichen Mathematik» (im Gegensatz zur diskreten Mathematik) auf Computern. Das englische Originalzitat «Numerical analysis is the study of algorithms for the problems of continuous mathematics.» hinter dieser Definition stammt von [Trefethen, 1992](#). Ein typisches Beispiel einer Grundaufgabe der Numerik ist etwa das Lösen linearer Gleichungssysteme  $Ax = b$  mit Hilfe des Gauß-Algorithmus.

Die Numerik ist damit ein Teilgebiet des **wissenschaftlichen Rechnens** (englisch: *scientific computing, computational science*), bei dem es darum geht, eine konkrete Fragestellung (etwa aus der Physik) in ein Modell und anschließend in ein numerisches Verfahren zu übersetzen, dieses auf geeigneter Hardware effizient zu implementieren und die Simulationsergebnisse mit realen Beobachtungen oder Experimenten abzugleichen, um damit die ursprüngliche Fragestellung beantworten zu können.

Nicholas Higham hat 2016 eine Liste der «Top 10»-Algorithmen in der angewandten Mathematik veröffentlicht<sup>1</sup>, indem er die Verweise im Index des über 1000-seitigen Buches *The Princeton Companion to Applied Mathematics* ([Higham u. a., 2016](#)) gezählt hat. Herausgekommen ist dabei folgende Liste:

- (1) Newton and quasi-Newton methods
- (2) Matrix factorizations (LU, Cholesky, QR)
- (3) Singular value decomposition, QR and QZ algorithms
- (4) Monte-Carlo methods
- (5) Fast Fourier transform
- (6) Krylov subspace methods (conjugate gradients, Lanczos, GMRES, MINRES)
- (7) JPEG
- (8) PAGERANK
- (9) Simplex algorithm
- (10) Kalman filter

Wir behandeln im Verlauf dieser Lehrveranstaltung von dieser Liste voraussichtlich zumindest teilweise die Punkte (1) bis (3) und (6) ~~Punkte (1) bis (3), (6) und (8)~~. Für das Simplex-Verfahren (Punkt (9)) verweisen wir auf die Lehrveranstaltung *Grundlagen der Optimierung*.

Zur Einführung gehen wir hier schon einmal auf die PAGERANK-Aufgabe ein. Dabei handelt es sich um

---

<sup>1</sup><https://press.princeton.edu/ideas/nicholas-higham-on-the-top-10-algorithms-in-applied-mathematics>

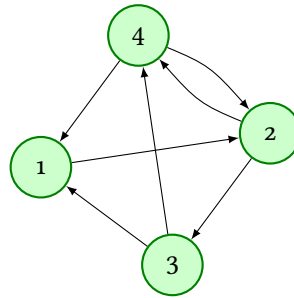


Abbildung 1.1: Ein Miniatur-Web zur Demonstration des PAGERANK-Verfahrens.

die von Sergey Brin und Lawrence Page in [Brin, Page, 1998](#) vorgeschlagenen Ranking-Methode für Webseiten, die seit dem Gründungsjahr 1998 von Google bestimmt, in welcher Reihenfolge Suchergebnisse der Google-Suche erscheinen. Unsere Beschreibung folgt [Higham u. a., 2016](#), Chapter VI.9.

Die in das Ranking einzubeziehenden Webseiten werden als sogenannter einfacher gerichteter Graph dargestellt, siehe [Abbildung 1.1](#). Dabei sind die Webseiten die **Knoten** (englisch: *vertex*) des Graphen, und ein Pfeil (eine **Kante**, englisch: *edge, arc*) bedeutet, dass eine Webseite per Link auf eine andere verweist.

Der Graph wird mittels einer **Adjazenzmatrix** (englisch: *adjacency matrix*) codiert:

$$A = \begin{matrix} & \textcircled{1} & \textcircled{2} & \textcircled{3} & \textcircled{4} \\ \begin{matrix} \textcircled{1} \\ \textcircled{2} \\ \textcircled{3} \\ \textcircled{4} \end{matrix} & \begin{bmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \end{matrix} \quad (1.1)$$

Dabei sind in der ersten Spalte diejenigen Kanten aufgeführt, die in Knoten  $\textcircled{1}$  starten usw.<sup>2</sup> In der ersten Zeile finden wir diejenigen Kanten, die im Knoten  $\textcircled{1}$  enden.

**Quizfrage:** Welche Bedeutung haben die Spaltensummen von  $A$ ? Welche Bedeutung haben die Zeilensummen von  $A$ ?

Wir stellen uns nun eine Person vor, die sich aktuell auf einer dieser Webseiten  $\textcircled{7}$  aufhält. Ihr Verhalten wird wie folgt modelliert. Mit einer Wahrscheinlichkeit  $\alpha$  wird die Person einem (von der aktuellen Webseite wegführenden) Link folgen. Jedem zur Verfügung stehenden Link wird dabei mit derselben Wahrscheinlichkeit gefolgt. Andernfalls, also mit Wahrscheinlichkeit  $1 - \alpha$ , begibt sich die Person zu einer zufälligen Webseite aus  $\textcircled{1} - \textcircled{4}$ .

Befindet sich die Person beispielsweise zu einem Zeitpunkt auf Webseite  $\textcircled{2}$ , so wird sie sich zum

<sup>2</sup>**Beachte:** Oft wird die Adjazenzmatrix gerichteter Graphen gerade als die Transponierte unserer Matrix  $A$  definiert. Dann stehen die Zeilen für die Anfangsknoten und die Spalten für die Endknoten. Für uns ist aber die Konvention wie in (1.1) günstig.



nächsten Zeitpunkt

auf Webseite ① mit Wahrscheinlichkeit  $\alpha \cdot \frac{0}{2} + (1 - \alpha) \cdot \frac{1}{4}$   
 auf Webseite ② mit Wahrscheinlichkeit  $\alpha \cdot \frac{0}{2} + (1 - \alpha) \cdot \frac{1}{4}$   
 auf Webseite ③ mit Wahrscheinlichkeit  $\alpha \cdot \frac{1}{2} + (1 - \alpha) \cdot \frac{1}{4}$   
 auf Webseite ④ mit Wahrscheinlichkeit  $\alpha \cdot \frac{1}{2} + (1 - \alpha) \cdot \frac{1}{4}$

befinden. (**Quizfrage:** Klar?) Der Faktor  $\frac{1}{2}$  kommt daher, dass die Webseite ② gerade zwei wegführende Links hat.

Insgesamt kann man diese Überlegungen in der Matrix  $\Pi$  der Übergangswahrscheinlichkeiten zusammenfassen, die in unserem Beispiel wie folgt aussieht:

$$\Pi = \alpha \begin{bmatrix} 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 & & & \\ & \frac{1}{2} & & \\ & & \frac{1}{2} & \\ & & & \frac{1}{2} \end{bmatrix} + (1 - \alpha) \frac{1}{4} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}. \quad (1.2)$$

Die Einträge  $\pi_{ij}$  dieser Matrix beschreiben, mit welcher Wahrscheinlichkeit die Person von Webseite ① (Spalte) im nächsten Zeitpunkt auf die Webseite ① (Zeile) gelangt.

Besteht das Netzwerk nun allgemeiner aus  $n$  Webseiten, ist weiterhin  $A$  die Adjazenzmatrix und  $D$  die Diagonalmatrix mit den Kehrwerten der Anzahlen der wegführenden Links<sup>3</sup> sowie  $e = (1, 1, \dots, 1)^T \in \mathbb{R}^n$ , dann hat die Matrix der Übergangswahrscheinlichkeiten die Gestalt

$$\Pi = \alpha A D + (1 - \alpha) \frac{1}{n} e e^T. \quad (1.3)$$

**Quizfrage:** Welche Eigenschaften fallen Ihnen zu dieser Matrix ein? Was sind insbesondere die Spaltensummen  $\sum_{i=1}^n \pi_{ij}$ ?

Welcher Zusammenhang besteht nun zwischen der Matrix  $\Pi$  und dem Ranking der Webseiten in unserem Netzwerk? Dazu betrachten wir folgendes Gedankenexperiment. Nehmen wir an, die Person startet zum Zeitpunkt 0 auf Webseite ②, also mit einer Wahrscheinlichkeitsverteilung von  $x^{(0)} = (0, 1, 0, 0)^T$ . Zum Zeitpunkt 1 ist sie dann mit den Wahrscheinlichkeiten auf den Seiten ③–④ angekommen, die sich aus den Einträgen von

$$x^{(1)} = \Pi x^{(0)}$$

ergeben. Zum Zeitpunkt 2 haben wir dann die Verteilung  $x^{(2)} = \Pi x^{(1)} = \Pi^2 x^{(0)}$  usw. Man kann im Fall von  $\alpha \in (0, 1)$  zeigen, dass diese Folge  $(x^{(n)})$  der Wahrscheinlichkeiten gegen eine stationäre Verteilung konvergiert, und zwar unabhängig davon, mit welcher Anfangsverteilung  $x^{(0)}$  wir beginnen.<sup>4</sup> Die Eigenschaft, dass  $x^{(n)}$  eine Wahrscheinlichkeitsverteilung beschreibt (also  $x^{(n)} \geq 0$  und  $e^T x^{(n)} = 1$ ) überträgt sich von der Anfangsverteilung  $x^{(0)}$  auf alle Nachfolger. **Quizfrage:** Warum?

<sup>3</sup>Wir gehen hier der Einfachheit davon aus, dass jede Webseite im betrachteten Netzwerk wenigstens einen von der Seite wegführenden Link hat.

<sup>4</sup>Wir haben es hier mit einer stationären Markow-Kette mit endlichem Zustandsraum zu tun.

Die Stationarität bedeutet, dass sich die Wahrscheinlichkeiten nicht mehr ändern, dass also gilt:

$$\Pi x = x. \quad (1.4)$$

Diese Gleichung (1.4) kann man auf verschiedene Arten interpretieren. Eine davon besagt, dass  $x$  ein Eigenvektor der Matrix  $\Pi$  zum Eigenwert  $\lambda = 1$  des Eigenwertproblems

$$\Pi x = \lambda x \quad (1.5)$$

ist. Man kann unter geeigneten Annahmen zeigen (Stichwort: **Perron-Frobenius-Theorem**),

- dass tatsächlich  $\lambda = 1$  immer ein Eigenwert der Übergangsmatrix  $\Pi$  ist und daher eine Lösung  $x \neq 0$  von (1.5) existiert,
- dass alle anderen Eigenwerte von  $\Pi$  vom Betrag her  $\leq 1$  sind,
- dass  $\lambda = 1$  der einzige strikt positive, reelle Eigenwert von  $\Pi$  ist – die anderen Eigenwerte sind entweder negativ reell oder komplexwertig,
- dass die algebraische, also auch die geometrische Vielfachheit des Eigenwerts  $\lambda = 1$  gleich eins ist,
- dass der zugehörige Eigenvektor  $x$  so gewählt werden kann, dass er nur Einträge  $\geq 0$  hat und für die Gesamtsumme gilt:  $e^\top x = 1$ .

Die Größe der Einträge des gesuchten Eigenvektors  $x$  gibt gerade die Wichtigkeit der zugehörigen Webseite ①–④ im Sinne der PAGERANK-Idee an. Die PAGERANK-Aufgabe ist also letztendlich eine **Eigenwertaufgabe** (englisch: *eigenvalue problem*)!

Angesichts der zur Zeit (Stand: 2022) geschätzten Anzahl von etwa  $n \approx 2 \cdot 10^9$  Webseiten und der entsprechenden Dimensionen der Matrizen  $A$  und  $\Pi$  ist neben der Frage, wie/ob man diese Matrizen speichert, auch die Frage relevant, wie man den gesuchten Eigenvektor effizient berechnet.

**Quizfrage:** Die Adjazenzmatrix  $A$  des heutigen Internet ist extrem dünn besetzt. Trifft das auch für die Matrix  $\Pi$  zu? Was hat das für Konsequenzen?

Wir nehmen hier vorweg, dass glücklicherweise für die Bestimmung von Eigenpaaren für betragsgrößte Eigenvektoren einer Matrix mit der **Potenzmethode** (auch **Vektoriteration**, **Von-Mises-Iteration**, englisch: *power iteration*) eine matrix-freie Methode zur Verfügung steht. Wir geben sie im Folgenden an und kommen darauf später in der Lehrveranstaltung zurück. Für die Matrix  $\Pi$  aus (1.2) kann man die Potenzmethode wie in **Algorithmus 1.1** gezeigt aufschreiben.

**Algorithmus 1.1** (PAGERANK-Algorithmus mittels Potenzmethode).

**Eingabe:** Adjazenzmatrix  $A \in \{0, 1\}^{n \times n}$

**Eingabe:** Link-Folge-Wahrscheinlichkeit  $\alpha \in (0, 1)$

**Ausgabe:** ungefähre Eigenvektor  $x \geq 0$  der Matrix  $\Pi$  aus (1.3) zum Eigenwert  $\lambda = 1$

- 1: Initialisiere  $x^{(0)} := e/n$   $\| e = (1, 1, \dots, 1)^\top \in \mathbb{R}^n$
- 2: Setze  $k := 0$
- 3: **while** noch nicht konvergiert **do**
- 4:     Setze  $x^{(k+1)} := \alpha A D x^{(k)} + (1 - \alpha) \frac{1}{n} e$   $\| x^{(k+1)} = \Pi x^{(k)}$
- 5:     Setze  $x^{(k+1)} := x^{(k+1)} / e^\top x^{(k+1)}$   $\| \text{Normalisierung}$

```
6:   Setze  $k := k + 1$ 
7: end while
```

Als Konvergenzkriterium in [Zeile 3](#) kann man etwa heranziehen, ob sich die Einträge von  $x^{(k+1)}$  nur noch sehr wenig von denen von  $x^{(k)}$  unterscheiden. Die Initialisierung in [Zeile 1](#) kann auch anders erfolgen, etwa indem man  $x^{(0)}$  mit Werten in  $[0, 1]$  zufällig wählt und dann normalisiert, sodass  $e^T x^{(0)} = 1$  gilt, oder durch die Wahl eines Einheitsvektors (ein Eintrag eins, Rest Nullen).

**Quizfrage:** Wie kommt man auf die Anweisung in [Zeile 4](#)? Was ist mit der Matrix in (1.3) passiert, die als lauter Einsen besteht?

**Quizfrage:** Wozu könnte der Schritt in [Zeile 5](#) dienen? Sollte dieser nicht eigentlich überflüssig sein?

Für unser obiges Beispiel aus [Abbildung 1.1](#) ergibt sich bei  $\alpha = 0.85$  und Startvektor  $x^{(0)} = (0, 1, 0, 0)^T$  auf vier Stellen gerundet die Folge

$$\begin{aligned} x^{(0)} &= (0.0000, 1.0000, 0.0000, 0.0000)^T \\ x^{(1)} &= (0.0375, 0.0375, 0.4625, 0.4625)^T \\ x^{(2)} &= (0.4306, 0.2659, 0.0534, 0.2500)^T \\ &\vdots \\ x^{(23)} &= (0.2240, 0.3374, 0.1809, 0.2578)^T \\ x^{(24)} &= (0.2239, 0.3374, 0.1809, 0.2578)^T \\ x^{(25)} &= (0.2239, 0.3374, 0.1809, 0.2578)^T \\ &\vdots \end{aligned}$$

Die Abhängigkeit der PAGERANK:s von der Link-Folge-Wahrscheinlichkeit  $\alpha$  für das Beispiel-Netzwerk aus [Abbildung 1.1](#) wird in [Abbildung 1.2](#) gezeigt. (**Quizfrage:** Wie kann man [Abbildung 1.2](#), vor allem für  $\alpha$  in der Nähe von 0 und 1 erklären?) Im Original-Paper [Brin, Page, 1998](#) wird eine Link-Folge-Wahrscheinlichkeit von  $\alpha = 0.85$  angenommen.

## § 2 WIEDERHOLUNG UND NOTATION

Wir gehen davon aus, dass der Leser mit Grundkonzepten der mehrdimensionalen Analysis und der linearen Algebra vertraut ist, darunter insbesondere Stetigkeit und Differenzierbarkeit von Funktionen sowie Matrix- und Vektorrechnung im Vektorraum  $\mathbb{R}^n$ .<sup>5</sup>

Wir bezeichnen die nicht-negativen bzw. die positiven reellen Zahlen mit

$$\mathbb{R}_{\geq 0} := \{x \in \mathbb{R} \mid x \geq 0\} \quad \text{und} \quad \mathbb{R}_{> 0} := \{x \in \mathbb{R} \mid x > 0\}.$$

Die natürlichen Zahlen sind  $\mathbb{N} = \{1, 2, \dots\}$ , und wir setzen  $\mathbb{N}_0 := \mathbb{N} \cup \{0\}$ .

<sup>5</sup>Vieles kann direkt auch auf den komplexen Vektorraum  $\mathbb{C}^n$  übertragen werden.

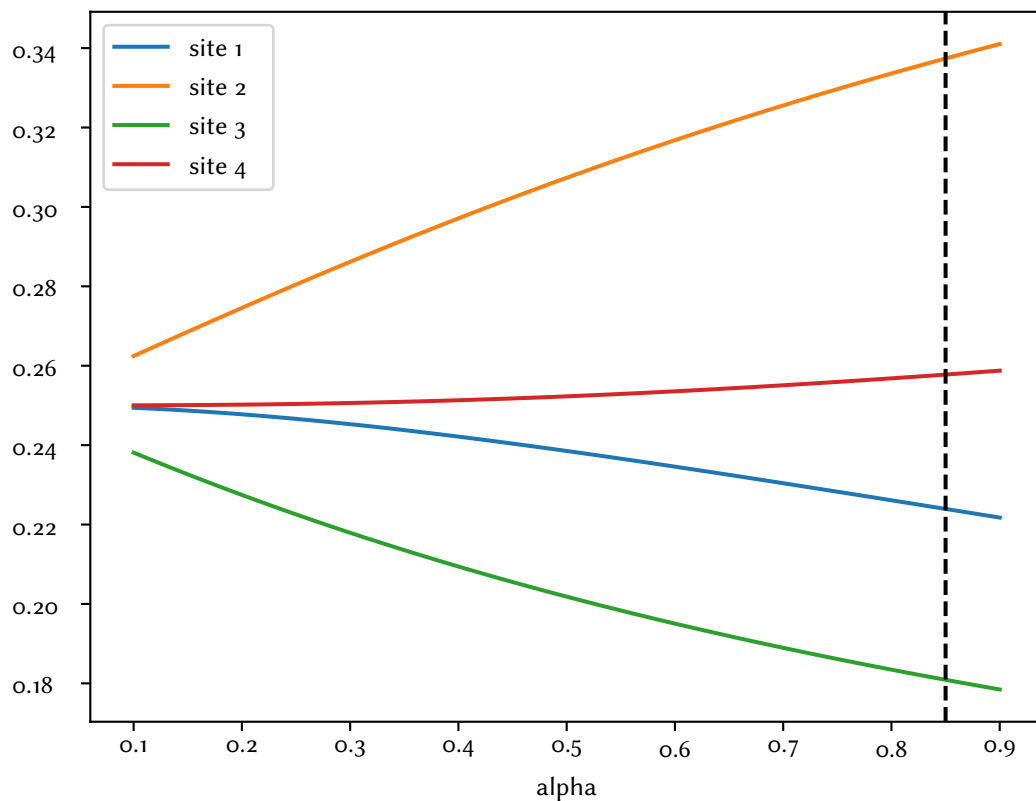


Abbildung 1.2: Gewichte (Einträge im Eigenvektor zum Eigenwert  $\lambda = 1$ ) der Webseiten für das Beispiel-Netzwerk aus [Abbildung 1.1](#) in Abhängigkeit der Link-Folge-Wahrscheinlichkeit  $\alpha$ .

## § 2.1 VEKTOR- UND MATRIX-NORMEN

**Definition 2.1** (Norm). Es sei  $V$  irgendein Vektorraum über  $\mathbb{R}$ . Eine Abbildung  $\|\cdot\|: V \rightarrow \mathbb{R}$  heißt eine **Norm** (englisch: **norm**) auf  $V$ , wenn für beliebige  $x, y \in V$  und  $\alpha \in \mathbb{R}$  folgende Eigenschaften gelten:

(i)  $\|x\| \geq 0$  und  $\|x\| = 0 \Rightarrow x = 0$ .  
(**Definitheit**)

(ii)  $\|\alpha x\| = |\alpha| \|x\|$ .  
(**positive Homogenität**)

(iii)  $\|x + y\| \leq \|x\| + \|y\|$ .  
(**Subadditivität / Dreiecksungleichung**)

## VEKTORNORMEN

Wir arbeiten in dieser Lehrveranstaltung im Vektorraum  $\mathbb{R}^n$  der Spaltenvektoren, soweit nichts anderes gesagt wird, mit dem **Euklidischen Innenprodukt** (englisch: *Euclidean inner product*)

$$x^T y = \sum_{i=1}^n x_i y_i \quad \text{für } x, y, \in \mathbb{R}^n,$$

der daraus erzeugten **Euklidischen Norm** (auch: **2-Norm**)

$$\|x\|_2 := \sqrt{x^T x} \quad \text{für } x \in \mathbb{R}^n$$

und dem wiederum daraus erzeugten Abstand (der **Metrik**, (englisch: *metric*))

$$\|x - y\|_2 \quad \text{für } x, y, \in \mathbb{R}^n.$$

Wie für jedes Innenprodukt gilt die **Cauchy-Schwarz-Ungleichung** (englisch: *Cauchy-Schwarz inequality*)

$$x^T y \leq |x^T y| \leq \|x\|_2 \|y\|_2 \quad \text{für alle } x, y, \in \mathbb{R}^n,$$

wobei die Gleichheit  $|x^T y| = \|x\|_2 \|y\|_2$  genau dann gilt, wenn  $x$  und  $y$  linear abhängig sind.

**Quizfrage:** Wie sehen das Euklidische Innenprodukt und die daraus erzeugte Euklidische Norm für **Zeilenvektoren** (englisch: *row vectors*) im  $\mathbb{R}_n$  aus?

Im Folgenden gelten alle Aussagen für  $\mathbb{R}^n$  sinngemäß auch für  $\mathbb{R}_n$ .

Wir benötigen außerdem die **1-Norm** und die  **$\infty$ -Norm** von Vektoren im  $\mathbb{R}^n$ , definiert durch

$$\|x\|_1 := \sum_{i=1}^n |x_i| \quad \text{und} \quad \|x\|_\infty := \max_{i=1, \dots, n} |x_i|.$$

**Beachte:** In den endlich-dimensionalen Vektorräumen  $\mathbb{R}^n$  sind alle Normen äquivalent. Es ist also unerheblich, in welcher Norm wie z. B. Konvergenz oder Fehler messen. Die Wahl erfolgt häufig danach, wie es am einfachsten geht. Für die 1-, 2- und  $\infty$ -Normen im  $\mathbb{R}^n$  gelten die folgenden Abschätzungen, die die Äquivalenz dieser Normen untereinander bestätigen:

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2, \tag{2.1a}$$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty, \tag{2.1b}$$

$$\|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty. \tag{2.1c}$$

**Quizfrage:** Für welche Vektoren im  $\mathbb{R}^n$  sind die Abschätzungen im (2.1) jeweils scharf?

Für  $\varepsilon > 0$  und  $x^* \in \mathbb{R}^n$  ist

$$B_\varepsilon(x^*) := \{x \in \mathbb{R}^n \mid \|x - x^*\|_2 < \varepsilon\}$$

die **offene  $\varepsilon$ -Umgebung** von  $x^*$  oder auch die **offene  $\varepsilon$ -Kugel** um  $x^*$ .

## MATRIXNORMEN

Matrizen kennzeichnen wir i. d. R. mit Großbuchstaben und ihre Einträge mit kleinen Buchstaben, beispielsweise hat eine Matrix  $A \in \mathbb{R}^{m \times n}$  die Einträge  $a_{ij}$  mit  $i = 1, \dots, m$  und  $j = 1, \dots, n$ . Manchmal schreiben wir statt  $A$  dann auch  $(a_{ij})$ , um die gesamte Matrix zu adressieren. Die **Einheitsmatrix** bezeichnen wir mit  $\text{Id}$ . Wenn wir die Dimension betonen wollen, schreiben wir auch  $\text{Id}_n$  für die Einheitsmatrix der Größe  $n \times n$ .

Auch im Vektorraum der Matrizen  $A \in \mathbb{R}^{m \times n}$  bzw. der durch sie repräsentierten linearen Abbildung  $\mathbb{R}^n \ni x \mapsto Ax \in \mathbb{R}^m$  lassen sich Normen definieren. Wir betrachten hier vorwiegend die durch die Vektornormen  $\|\cdot\|_2$ ,  $\|\cdot\|_1$  und  $\|\cdot\|_\infty$  induzierten Operatornormen:

**Definition 2.2** (Operatornorm einer Matrix). *Es sei  $A \in \mathbb{R}^{m \times n}$ . Weiter seien  $\|\cdot\|_X$  bzw.  $\|\cdot\|_Y$  die in den Vektorräumen  $\mathbb{R}^n$  und  $\mathbb{R}^m$  verwendeten Vektornormen. Die Größe*

$$\|A\|_{X \rightarrow Y} := \sup_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_Y}{\|x\|_X} \quad (2.2)$$

heißt die (durch die Normen  $\|\cdot\|_X$  im  $\mathbb{R}^n$  und  $\|\cdot\|_Y$  im  $\mathbb{R}^m$  induzierte) **Operatornorm** (englisch: **operator norm**) von  $A$ .

**Lemma 2.3** (Eigenschaften der Operatornorm). *Es seien  $A \in \mathbb{R}^{m \times n}$  und  $B \in \mathbb{R}^{\ell \times m}$ . In den Räumen  $\mathbb{R}^n$ ,  $\mathbb{R}^m$  und  $\mathbb{R}^\ell$  verwenden wir die Normen  $\|\cdot\|_X$ ,  $\|\cdot\|_Y$  bzw.  $\|\cdot\|_Z$ . Für die Operatornorm (2.2) gelten folgende Eigenschaften:*

(i) Die Operatornorm ist eine Norm auf  $\mathbb{R}^{m \times n}$ .

(ii)  $\|Ax\|_Y \leq \|A\|_{X \rightarrow Y} \|x\|_X$  für alle  $x \in \mathbb{R}^n$ .

(iii) Es gilt

$$\|A\|_{X \rightarrow Y} = \sup_{\substack{x \in \mathbb{R}^n \\ \|x\|_X = 1}} \|Ax\|_Y = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_X = 1}} \|Ax\|_Y.$$

(iv) Es gilt

$$\begin{aligned} \|A\|_{X \rightarrow Y} &= \inf \{c \geq 0 \mid \|Ax\|_Y \leq c \|x\|_X \text{ für alle } x \in \mathbb{R}^n\} \\ &= \min \{c \geq 0 \mid \|Ax\|_Y \leq c \|x\|_X \text{ für alle } x \in \mathbb{R}^n\}. \end{aligned}$$

(v)  $\|BA\|_{X \rightarrow Z} \leq \|B\|_{Y \rightarrow Z} \|A\|_{X \rightarrow Y}$  (**Submultiplikativität**).

**Beachte:** Aussage (iv) bedeutet: Kennt man eine Zahl  $c \geq 0$ , mit der  $\|Ax\|_Y \leq c \|x\|_X$  für alle  $x \in \mathbb{R}^n$  gilt, so besteht die Abschätzung  $\|A\|_{X \rightarrow Y} \leq c$ .

*Beweis.* Findet sich als Übungsaufgabe 1 auf Übungsblatt 1. □

**Lemma 2.4** (Operatornormen). Für Matrizen  $A \in \mathbb{R}^{m \times n}$  und ihre Operatornormen, induziert durch die Verwendung jeweils der 1-Norm, der 2-Norm bzw. der  $\infty$ -Norm im Definitionsraum  $\mathbb{R}^n$  wie auch im Zielraum  $\mathbb{R}^m$  gelten die folgenden Aussagen:

$$\|A\|_{1 \rightarrow 1} = \max_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_1}{\|x\|_1} = \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}| \quad (\text{Spaltensummennorm}), \quad (2.3a)$$

$$\|A\|_{2 \rightarrow 2} = \max_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_2}{\|x\|_2} = \text{Wurzel des maximalen Eigenwerts von } A^T A \quad (\text{Spektralnrm}), \quad (2.3b)$$

$$\|A\|_{\infty \rightarrow \infty} = \max_{x \in \mathbb{R}^n \setminus \{0\}} \frac{\|Ax\|_\infty}{\|x\|_\infty} = \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}| \quad (\text{Zeilensummennorm}). \quad (2.3c)$$

*Beweis.* Zum Beweis von (2.3a) schätzen wir ab:

$$\begin{aligned} \|Ax\|_1 &= \sum_{i=1}^m \left| \sum_{j=1}^n a_{ij} x_j \right| \\ &\leq \sum_{i=1}^m \sum_{j=1}^n |a_{ij}| |x_j| && (\text{Dreiecksungleichung}) \\ &= \sum_{j=1}^n |x_j| \sum_{i=1}^m |a_{ij}| \\ &\leq \sum_{j=1}^n |x_j| \max_{k=1, \dots, n} \sum_{i=1}^m |a_{ik}| \\ &= \|x\|_1 \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}| \quad \text{für alle } x \in \mathbb{R}^n. \end{aligned}$$

Es gilt also (siehe Text nach Lemma 2.3)  $\|A\|_{1 \rightarrow 1} \leq \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|$ . Um die Gleichheit zu zeigen, sei  $j^*$  einer derjenigen Indizes, für die

$$\sum_{i=1}^m |a_{ij^*}| = \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|$$

gilt. Wir wählen  $x^* := e_{j^*}$  als  $j^*$ -ten Einheitsvektor in  $\mathbb{R}^n$ . Dann gilt  $\|x^*\|_1 = 1$  und

$$\|Ax^*\|_1 = \|Ae_{j^*}\|_1 = \left\| \begin{pmatrix} a_{1j^*} \\ \vdots \\ a_{mj^*} \end{pmatrix} \right\|_1 = \sum_{i=1}^m |a_{ij^*}| = \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}| = \|x^*\|_1 \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|,$$

also  $\|A\|_{1 \rightarrow 1} \geq \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|$ . Zusammen mit der obigen umgekehrten Abschätzung ist (2.3a) gezeigt.

Zum Beweis von (2.3c) schätzen wir wie folgt ab:

$$\begin{aligned}
 \|Ax\|_\infty &= \max_{i=1,\dots,m} \left| \sum_{j=1}^n a_{ij} x_j \right| \\
 &\leq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}| |x_j| \quad (\text{Dreiecksungleichung}) \\
 &\leq \max_{j=1,\dots,n} |x_j| \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}| \\
 &= \|x\|_\infty \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}| \quad \text{für alle } x \in \mathbb{R}^n.
 \end{aligned}$$

Es gilt also  $\|A\|_{\infty \rightarrow \infty} \leq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|$ . Um die Gleichheit zu zeigen, sei  $i^*$  einer derjenigen Indizes, für die

$$\sum_{j=1}^n |a_{i^*j}| = \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|$$

gilt. Wir wählen  $x^* := (\text{sgn } a_{i^*1}, \dots, \text{sgn } a_{i^*n})^\top$ .<sup>6</sup> Dann gilt  $\|x^*\|_\infty = 1$  und

$$\|Ax^*\|_\infty = \max_{i=1,\dots,m} \left| \sum_{j=1}^n a_{ij} x_j^* \right| \geq \left| \sum_{j=1}^n a_{i^*j} x_j^* \right| = \left| \sum_{j=1}^n |a_{i^*j}| \right| = \sum_{j=1}^n |a_{i^*j}| = \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|,$$

also  $\|A\|_{\infty \rightarrow \infty} \geq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|$ . Zusammen mit der obigen umgekehrten Abschätzung ist (2.3c) gezeigt.

Den Beweis von (2.3b) liefern wir in Abschnitt 4 nach. □

Für die Operatornormen von Spaltenvektoren  $y \in \mathbb{R}^m \cong \mathbb{R}^{m \times 1}$  und Zeilenvektoren  $y \in \mathbb{R}_n \cong \mathbb{R}^{1 \times n}$  gilt Folgendes:

**Folgerung 2.5** (Operatornormen von Zeilen- und Spaltenvektoren).

(i) Für **Spaltenvektoren**  $y \in \mathbb{R}^m$ , die wir auch als Matrix  $y \in \mathbb{R}^{m \times 1}$  interpretieren, gilt:

$$\|y\|_{1 \rightarrow 1} = \|y\|_1, \quad \|y\|_{2 \rightarrow 2} = \|y\|_2, \quad \|y\|_{\infty \rightarrow \infty} = \|y\|_\infty \quad (2.4)$$

(ii) Für **Zeilenvektoren**  $y \in \mathbb{R}_n$ , die wir auch als Matrix  $y \in \mathbb{R}^{1 \times n}$  interpretieren, gilt:

$$\|y\|_{1 \rightarrow 1} = \|y\|_\infty, \quad \|y\|_{2 \rightarrow 2} = \|y\|_2, \quad \|y\|_{\infty \rightarrow \infty} = \|y\|_1 \quad (2.5)$$

**Beweis.** Aussage (i): Für  $y \in \mathbb{R}^m \cong \mathbb{R}^{m \times 1}$  gilt nach (2.3a) bzw. (2.3c)

$$\|y\|_{1 \rightarrow 1} = \sum_{i=1}^m |y_{i1}| = \|y\|_1 \quad \text{und} \quad \|y\|_{\infty \rightarrow \infty} = \max_{i=1,\dots,m} |y_{i1}| = \|y\|_\infty.$$

<sup>6</sup>Das **Signum** einer reellen Zahl  $x$  ist definiert als  $\text{sgn } x = 1$  für  $x > 0$ ,  $\text{sgn } x = -1$  für  $x < 0$  und  $\text{sgn } x = 0$  für  $x = 0$ .



Da wir (2.3b) noch nicht bewiesen haben, untersuchen wir die 2-Norm direkt:

$$\|y\|_{2 \rightarrow 2} = \max_{\substack{\alpha \in \mathbb{R} \\ |\alpha|=1}} \|\alpha y\|_2 = \|y\|_2.$$

**Aussage (ii):** Für  $y \in \mathbb{R}^n \cong \mathbb{R}^{1 \times n}$  gilt nach (2.3a) bzw. (2.3c)

$$\|y\|_{1 \rightarrow 1} = \max_{j=1, \dots, n} |y_{1j}| = \|y\|_\infty \quad \text{und} \quad \|y\|_{\infty \rightarrow \infty} = \sum_{j=1}^n |y_{1j}| = \|y\|_1.$$

Die 2-Norm untersuchen wir wieder direkt:

$$\|y\|_{2 \rightarrow 2} = \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} |y x| \leq \max_{\substack{x \in \mathbb{R}^n \\ \|x\|_2=1}} \|y\|_2 \|x\|_2.$$

Im Fall  $y = 0$  sind ohnehin beide Seiten gleich Null. Andernfalls erreichen wir durch die Wahl von  $x = y/\|y\|_2$  die Gleichheit (Cauchy-Schwarz-Ungleichung).  $\square$

## § 2.2 DIFFERENZIERBARKEIT UND DER SATZ VON TAYLOR

Zwei wesentliche Hilfsmittel zur Analyse von Aufgaben und Algorithmen sind die Kettenregel und der Satz von Taylor, die wir hier kurz wiederholen:

**Definition 2.6** (Differenzierbarkeit, siehe z. B. Heuser, 2002, Definition 164). Es sei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  eine Funktion.  $F$  heißt **differenzierbar** (englisch: **differentiable**) an der Stelle  $x_0 \in \mathbb{R}^n$ , wenn es eine Matrix  $A \in \mathbb{R}^{m \times n}$  gibt, sodass für die durch

$$r(x_0; \Delta x) := F(x_0 + \Delta x) - F(x_0) - A \Delta x$$

definierte **Restgliedfunktion** (englisch: **remainder function**) gilt:<sup>7</sup>

$$\lim_{\Delta x \rightarrow 0} \frac{r(x_0; \Delta x)}{\|\Delta x\|} = 0. \quad (2.6)$$

In diesem Fall heißt die durch die Definition eindeutig bestimmte Matrix  $A$  die **Ableitung** (englisch: **derivative**) von  $F$  an der Stelle  $x_0$  und wird mit  $F'(x_0)$  bezeichnet. Sie stimmt überein mit der **Jacobi-matrix** (der Matrix aus allen partiellen Ableitungen erster Ordnung, englisch: **Jacobian**) von  $F$  an der Stelle  $x_0$ :

$$J_F(x_0) := \begin{bmatrix} \frac{\partial F_1(x_0)}{\partial x_1} & \dots & \frac{\partial F_1(x_0)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial F_m(x_0)}{\partial x_1} & \dots & \frac{\partial F_m(x_0)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^{m \times n}. \quad (2.7)$$

<sup>7</sup>Da diese Aussage wegen der Äquivalenz der Normen auf  $\mathbb{R}^n$  für jede Vektornorm  $\|\cdot\|$  gilt, wird die Norm hier nicht näher spezifiziert.

Partielle Ableitungen zweiter Ordnung werden, sofern sie existieren, nach dem Muster

$$\frac{\partial^2 F_j(x_0)}{\partial x_2 \partial x_1} := \frac{\partial}{\partial x_2} \frac{\partial F_j}{\partial x_1}(x_0)$$

gebildet. Analog gibt es auch Ableitungen höherer Ordnung. **Quizfrage:** Wieviele verschiedene partielle Ableitungen der Ordnung  $k \in \mathbb{N}_0$  gibt es?

**Definition 2.7** ( $C^k$ -Funktion). Es sei  $k \in \mathbb{N}_0$  und  $U \subseteq \mathbb{R}^n$  offen. Eine Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  heißt  **$k$ -mal stetig partiell differenzierbar** (englisch:  *$k$  times continuously partially differentiable*) auf  $U$  oder eine  **$C^k$ -Funktion** auf  $U$  (kurz:  $C^k(U, \mathbb{R}^m)$ ), wenn  $F$  sowie alle partiellen Ableitungen bis einschließlich zur Ordnung  $k$  existieren und auf  $U$  stetig sind.

**Beachte:** Eine Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , die auf der offenen Menge  $U \subseteq \mathbb{R}^n$  eine  $C^1$ -Funktion ist, ist in jedem Punkt von  $U$  auch differenzierbar, erfüllt also die Restgliedbedingung (2.6) für jedes feste  $x_0 \in U$ . Kurz: Eine einmal stetig partiell differenzierbare Funktion ist differenzierbar.

**Satz 2.8** (Kettenregel).

Es seien  $F: \mathbb{R}^m \rightarrow \mathbb{R}^p$  und  $G: \mathbb{R}^n \rightarrow \mathbb{R}^m$  zwei Funktionen, und es sei  $G$  an der Stelle  $x_0 \in \mathbb{R}^n$  differenzierbar sowie  $F$  an der Stelle  $G(x_0)$  differenzierbar. Dann ist auch die Komposition  $F \circ G$  an der Stelle  $x_0$  differenzierbar, und es gilt

$$(F \circ G)'(x_0) = F'(G(x_0)) G'(x_0)$$

oder ausgeschrieben

$$\begin{bmatrix} \frac{\partial (F \circ G)_1(x_0)}{\partial x_1} & \dots & \frac{\partial (F \circ G)_1(x_0)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial (F \circ G)_p(x_0)}{\partial x_1} & \dots & \frac{\partial (F \circ G)_p(x_0)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \frac{\partial F_1(G(x_0))}{\partial x_1} & \dots & \frac{\partial F_1(G(x_0))}{\partial x_m} \\ \vdots & & \vdots \\ \frac{\partial F_p(G(x_0))}{\partial x_1} & \dots & \frac{\partial F_p(G(x_0))}{\partial x_m} \end{bmatrix} \begin{bmatrix} \frac{\partial G_1(x_0)}{\partial x_1} & \dots & \frac{\partial G_1(x_0)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial G_m(x_0)}{\partial x_1} & \dots & \frac{\partial G_m(x_0)}{\partial x_n} \end{bmatrix}.$$

Die Ableitung (Jacobimatrix) der zusammengesetzten Funktion  $F \circ G$  ergibt sich also aus dem Produkt der Jacobimatrizen  $F'(G(x_0))$  und  $G'(x_0)$ .

Wir geben nun zwei Versionen des Satzes von Taylor für *reellwertige* Funktionen an, die sich in der Darstellung des Restglieds unterscheiden. Analog zur Schreibweise für Intervalle

$$[a, b] := \{x \in \mathbb{R} \mid a \leq x \text{ und } x \leq b\} \quad (2.8)$$

ist es günstig, eine ähnliche Schreibweise auch für die **Verbindungsstrecke** (englisch: *segment*) zwischen Punkten  $a, b$  im  $\mathbb{R}^n$  einzuführen, also

$$\langle a, b \rangle := \{(1-t)a + tb \mid t \in [0, 1]\} \quad (2.9)$$

für  $a, b \in \mathbb{R}^n$ . **Quizfrage:** Gilt für reelle Zahlen  $a$  und  $b$  immer  $[a, b] = \langle a, b \rangle$ ?

**Definition 2.9** (Konvexe Menge). Eine Menge  $A \subseteq \mathbb{R}^n$  heißt **konvex** (englisch: *convex*), wenn für  $x, y \in A$  immer auch die gesamte Verbindungsstrecke in  $A$  liegt, also  $\langle x, y \rangle \subseteq A$ .

**Satz 2.10** (Taylor mit Restglied in Zwischenwertform, siehe z. B. Heuser, 2002, Satz 168.1).

Es sei  $U \subseteq \mathbb{R}^n$  offen,  $k \in \{1, 2\}$  und  $f: U \rightarrow \mathbb{R}$  eine  $C^k$ -Funktion auf  $U$ , kurz:  $C^k(U, \mathbb{R})$ . Die Punkte  $x_0$  und  $x_0 + d$  seien so gewählt, dass  $\langle x_0, x_0 + d \rangle$  in  $U$  liegt. Dann existiert eine Zahl  $\xi \in [0, 1]$ , sodass gilt:

$$\text{im Fall } k = 1: \quad f(x_0 + d) = f(x_0) + f'(x_0 + \xi d) d,$$

$$\text{im Fall } k = 2: \quad f(x_0 + d) = f(x_0) + f'(x_0) d + \frac{1}{2} d^T f''(x_0 + \xi d) d.$$

Die erste Aussage ist gerade der **Mittelwertsatz** (englisch: *mean value theorem*). Weiter sind  $f'(x_0)$  und  $f''(x_0)$  die Ableitungen erster bzw. zweiter Ordnung von  $f$  an der Stelle  $x_0$ , die folgende Gestalt besitzen:

$$f'(x_0) = \begin{bmatrix} \frac{\partial f(x_0)}{\partial x_1} & \dots & \frac{\partial f(x_0)}{\partial x_n} \end{bmatrix} \quad (\text{Jacobimatrix})$$

$$f''(x_0) = \begin{bmatrix} \frac{\partial^2 f(x_0)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x_0)}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f(x_0)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(x_0)}{\partial x_n \partial x_n} \end{bmatrix} \quad (\text{Hessematrix, englisch: Hessian}).$$

Als Resultat des Satzes von Schwarz (Heuser, 2002, Satz 162.2) ist die Hessematrix  $f''(x_0)$  dabei symmetrisch, also die Reihenfolge der partiellen Differentiationen unerheblich.

**Quizfrage:** Gilt der Satz von Taylor 2.10 auch für vektorwertige  $C^k$ -Funktionen  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ?

**Satz 2.11** (Taylor mit integralem Restglied, siehe z. B. Heuser, 2002, Satz 168.1).

Unter denselben Voraussetzungen wie in Satz 2.10 gilt

$$\text{im Fall } k = 1: \quad f(x_0 + d) = f(x_0) + \left[ \int_0^1 f'(x_0 + t d) dt \right] d$$

$$= f(x_0) + \int_0^1 f'(x_0 + t d) d dt,$$

$$\text{im Fall } k = 2: \quad f(x_0 + d) = f(x_0) + f'(x_0) d + \frac{1}{2} d^T \left[ \int_0^1 (1-t) f''(x_0 + t d) dt \right] d$$

$$= f(x_0) + f'(x_0) d + \frac{1}{2} \int_0^1 (1-t) d^T f''(x_0 + t d) d dt.$$

Die erste Aussage folgt direkt aus dem **Hauptsatz der Differential- und Integralrechnung**.<sup>8</sup>

**Beachte:** Unter dem jeweils ersten Integral steht eine vektorwertige bzw. matrixwertige Funktion, die einfach komponentenweise integriert wird. Im jeweils zweiten Integral wurden die konstanten (von  $t$  unabhängigen) Inkremente  $d$  in das Integral hineingezogen, wodurch der Integrand reellwertig wird.

**Quizfrage:** In dieser Form gilt Satz 2.11 auch für vektorwertige  $C^k$ -Funktionen  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ . Wie genau sollte der Satz dann formuliert werden?

<sup>8</sup>«Eine differenzierbare Funktion ist das Integral ihrer Ableitung.»

**Definition 2.12** (Lipschitz-Stetigkeit). Eine Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  heißt **Lipschitz-stetig** (englisch: **Lipschitz continuous**) auf der Menge  $A \subseteq \mathbb{R}^n$ , falls es eine Zahl  $L \geq 0$  (genannt: **Lipschitz-Konstante**, englisch: **Lipschitz constant**) gibt mit der Eigenschaft<sup>9</sup>

$$\|F(x) - F(y)\|_2 \leq L \|x - y\|_2 \quad \text{für alle } x, y \in A. \quad (2.10)$$

In diesem Fall heißt  $F$  dann auch  **$L$ -Lipschitz-stetig** auf  $A$ .

Eine einfache aber wichtige Konsequenz aus dem **Satz von Taylor 2.10** ist die Tatsache, dass  $C^1$ -Funktionen genau dann Lipschitz-stetig sind, wenn ihre Ableitung beschränkt ist, genauer:

**Folgerung 2.13** (Lipschitz-Stetigkeit von  $C^1$ -Funktionen). Es sei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  eine  $C^1$ -Funktion auf der offenen, konvexen Menge  $U \subseteq \mathbb{R}^n$ . Dann sind folgende Aussagen äquivalent:

- (i)  $F$  ist Lipschitz-stetig auf  $U$ .
- (ii) Die Ableitung  $F'$  ist auf  $U$  beschränkt, also  $\sup\{\|F'(x)\|_{2 \rightarrow 2} \mid x \in U\} < \infty$ .

Falls die Aussagen wahr sind, dann ist  $L^* := \sup\{\|F'(x)\|_{2 \rightarrow 2} \mid x \in U\} < \infty$  die kleinstmögliche Lipschitz-Konstante in (2.10).

**Beachte:**  $\|F'(x)\|_{2 \rightarrow 2}$  ist die durch die 2-Norm induzierte Operatornorm der Jacobimatrix von  $F$  an der Stelle  $x$ , vgl. auch (2.3b).

**Beweis.** Wir nehmen zunächst **Aussage (i)** an. Es seien  $x \in U$  und  $d \in \mathbb{R}^n$  fest gewählt sowie  $\alpha > 0$ . Aus der Dreiecksungleichung folgt

$$\begin{aligned} \|F'(x)(\alpha d)\|_2 &= \|F'(x)(\alpha d) \pm [F(x + \alpha d) - F(x)]\|_2 \\ &\leq \|F(x + \alpha d) - F(x) - F'(x)(\alpha d)\|_2 + \|F(x + \alpha d) - F(x)\|_2. \end{aligned}$$

Es gilt also auch

$$\|F'(x)d\|_2 = \frac{\|F'(x)(\alpha d)\|_2}{\alpha} \leq \frac{\|F(x + \alpha d) - F(x) - F'(x)(\alpha d)\|_2}{\alpha} + \frac{\|F(x + \alpha d) - F(x)\|_2}{\alpha}.$$

Für hinreichend kleine  $\alpha > 0$  liegt  $x + \alpha d \in U$ , und wir können für den ersten Term auf der rechten Seite die Voraussetzung der Differenzierbarkeit und für den zweiten Term die Voraussetzung der  $L$ -Lipschitz-Stetigkeit von  $F$  auf  $U$  mit einer Lipschitz-Konstanten  $L$  nutzen. Im Grenzübergang  $\alpha \searrow 0$  folgt:

$$\|F'(x)d\|_2 \leq 0 + L \|d\|_2.$$

Da  $d \in \mathbb{R}^n$  beliebig war, folgt aus **Lemma 2.3** jetzt  $\|F'(x)\|_2 \leq L$ .

<sup>9</sup>Auf der linken und rechten Seite von (2.10) dürften wegen der Äquivalenz der Normen auch andere Vektornormen im  $\mathbb{R}^n$  bzw.  $\mathbb{R}^m$  stehen. Die konkrete Wahl beeinflusst nicht die Frage, ob eine Funktion Lipschitz-stetig ist, sondern nur, welche Lipschitz-Konstanten gültig sind. Die Wahl der 2-Norm ist hier günstig wegen der Formulierung der **Folgerung 2.13**.

Wir nehmen nun **Aussage (ii)** an. Es seien  $x, y \in U$ . Aufgrund der Konvexität von  $U$  liegt auch  $\langle x, y \rangle$  in  $U$ . Aus dem **Satz von Taylor 2.11** folgt

$$F(y) - F(x) = \int_0^1 F'(x + t(y - x)) dt (y - x),$$

also (**Quizfrage:** Warum gelten jeweils die Abschätzungen?)

$$\begin{aligned} \|F(y) - F(x)\|_2 &= \left\| \int_0^1 F'(x + t(y - x)) (y - x) dt \right\| \\ &\leq \int_0^1 \|F'(x + t(y - x)) (y - x)\| dt \\ &\leq \int_0^1 \|F'(x + t(y - x))\|_2 \|y - x\|_2 dt. \end{aligned}$$

Der Integrand  $\|F'(x + t(y - x))\|_2$  ist aber durch  $\sup\{\|F'(x)\|_2 \mid x \in U\}$  beschränkt, da  $x + t(y - x)$  auf der Verbindungsstrecke  $\langle x, y \rangle$  und damit in  $U$  liegt. Damit ist  $L^* := \sup\{\|F'(x)\|_2 \mid x \in U\} < \infty$  eine gültige Lipschitz-Konstante für  $F$  auf  $U$ , also gilt **Aussage (i)**.

Wir müssen noch zeigen, dass  $L^*$  die kleinstmögliche Lipschitz-Konstante ist. Aus dem Beweisteil **Aussage (i)  $\Rightarrow$  Aussage (ii)** erhalten wir aber

$$F \text{ ist } L\text{-Lipschitz} \quad \Rightarrow \quad \|F'(x)\|_2 \leq L \text{ für alle } x \in U \quad \Rightarrow \quad L^* := \sup\{\|F'(x)\|_2 \mid x \in U\} \leq L.$$

□

## § 2.3 LANDAU-NOTATION

Wir geben jetzt noch die (Teile der) Definition der **Landau-Notation** (englisch: *Landau notation*) an, mit deren Hilfe man asymptotisches Verhalten, z. B. von Fehlern, charakterisieren kann.

**Definition 2.14** (Landau-Notation für Funktionen).

Es sei  $g: \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}_{\geq 0}$  eine Funktion.

(i) Es sei  $m \in \mathbb{N}$ . Wir bezeichnen die Menge von Funktionen

$$O(g) := \left\{ F: \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}^m \left| \begin{array}{l} \text{es existiert ein } C \geq 0 \text{ und ein } \varepsilon > 0, \text{ sodass gilt:} \\ \|F(x)\| \leq C g(x) \text{ für alle } x \text{ mit } 0 < \|x\| < \varepsilon \end{array} \right. \right\}$$

als «**Groß-O von g bei Null**» (englisch: *big-O of g near zero*).

Für  $F \in O(g)$  sagt man auch,  $g$  sei eine **asymptotisch obere Schranke** für  $F$  bei Null.

(ii) Es sei  $m \in \mathbb{N}$ . Wir bezeichnen die Menge von Funktionen

$$o(g) := \left\{ F: \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}^m \left| \begin{array}{l} \text{für alle } c > 0 \text{ existiert ein } \varepsilon > 0, \text{ sodass gilt:} \\ \|F(x)\| \leq c g(x) \text{ für alle } x \text{ mit } 0 < \|x\| < \varepsilon \end{array} \right. \right\}$$

als «**Klein-o von  $g$  bei Null**» (englisch: little-o of  $g$  near zero).

Für  $F \in o(g)$  sagt man auch,  $F$  sei gegenüber  $g$  **asymptotisch vernachlässigbar** bei Null.

**Beachte:** Wir lassen in der obigen Definition offen, welche Norm im  $\mathbb{R}^m$  wir verwenden, um  $\|F(x)\|$  zu messen, weil es wegen der Äquivalenz der Normen unerheblich ist.

**Quizfrage:** Wie kann man die Restgliedbedingung (2.6) in der Landau-Notation ausdrücken?

**Quizfrage:** Welche Beziehung besteht zwischen den Funktionenklassen  $O(\|\Delta x\|)$  und  $O(\|\Delta x\|^2)$ ?

**Beispiel 2.15** (Beispiele für Landau-Notation für Funktionen).

(i) Die Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ , definiert durch  $f(x) = \sin(x)$ , gehört zu  $O(g)$  für  $g(x) = |x|$ .  
Kurz schreibt man auch:  $\sin(x) \in O(|x|)$ .

(ii) Die Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ , definiert durch  $f(x) = x^2$ , gehört zu  $o(g)$  für  $g(x) = |x|$ .  
Kurz schreibt man auch:  $x^2 \in o(|x|)$ .

(iii) Die Funktion  $F: \mathbb{R}^2 \rightarrow \mathbb{R}^2$ , definiert durch  $F(x) = \begin{pmatrix} x_1 - x_2 \\ x_1^2 + x_2 \end{pmatrix}$ , gehört zu  $O(g)$  für  $g(x) = \|x\|$ .  
Kurz schreibt man auch:  $F(x) \in O(\|x\|)$ .

Auch für die Beschreibung der Asymptotik von Folgen ist die Landau-Notation hilfreich. Dann ist nicht der Fall «bei Null», sondern der Fall «bei Unendlich» interessant.

**Definition 2.16** (Landau-Notation für Folgen).

Es sei  $(y^{(n)})$  eine  $\mathbb{R}_{\geq 0}$ -wertige Folge.

(i) Es sei  $m \in \mathbb{N}$ . Wir bezeichnen die Menge von Folgen

$$O(y^{(n)}) := \left\{ (x^{(n)}) : \mathbb{N} \rightarrow \mathbb{R}^m \left| \begin{array}{l} \text{es existiert ein } C \geq 0 \text{ und ein } n_0 \in \mathbb{N}, \text{ sodass gilt:} \\ \|x^{(n)}\| \leq C y^{(n)} \text{ für alle } n \geq n_0 \end{array} \right. \right\}$$

als «**Groß-O von  $(y^{(n)})$  bei Unendlich**».

(ii) Es sei  $m \in \mathbb{N}$ . Wir bezeichnen die Menge von Folgen

$$o(y^{(n)}) := \left\{ (x^{(n)}) : \mathbb{N} \rightarrow \mathbb{R}^m \left| \begin{array}{l} \text{für alle } c > 0 \text{ existiert ein } n_0 \in \mathbb{N}, \text{ sodass gilt:} \\ \|x^{(n)}\| \leq c y^{(n)} \text{ für alle } n \geq n_0 \end{array} \right. \right\}$$

als «**Klein-o von  $(y^{(n)})$  bei Unendlich**».

**Beachte:** Auch in dieser Definition ist wieder irrelevant, welche Norm im  $\mathbb{R}^m$  wir verwenden, um  $\|x^{(n)}\|$  zu messen.

**Beispiel 2.17** (Beispiele für Landau-Notation für Folgen).

- (i) Die reellwertige Folge  $x^{(n)} = 1/n$  gehört zu  $O(y^{(n)})$  für  $y^{(n)} = 1$ .  
Kurz schreibt man auch:  $1/n \in O(1)$ .
- (ii) Die reellwertige Folge  $x^{(n)} = 1/n$  gehört zu  $o(y^{(n)})$  für  $y^{(n)} = 1$ .  
Kurz schreibt man auch:  $1/n \in o(1)$ .
- (iii) Die  $\mathbb{R}^2$ -wertige Folge  $x^{(n)} = \begin{pmatrix} 1/n \\ 1/n^2 \end{pmatrix}$  gehört zu  $O(y^{(n)})$  für  $y^{(n)} = 1/n$ .  
Kurz schreibt man auch:  $x^{(n)} \in O(1/n)$ .

**Quizfrage:** Was bedeuten diese Aussagen?

Ende der Woche 1

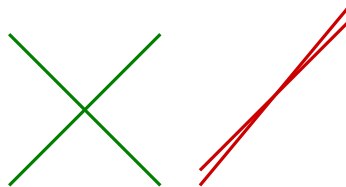
# Kapitel 1    Sensitivitätsanalyse und Kondition von Aufgaben

In diesem Kapitel beschäftigen wir uns zunächst mit Eigenschaften von Aufgaben, insbesondere mit dem Konzept der **Kondition** eines Problems. Der Grund dafür ist, dass jedes numerische Verfahren höchstens so gut sein kann, wie die Kondition der Aufgabe es zulässt.

## § 3    SENSITIVITÄTSANALYSE UND KONDITION

Zur Einstimmung zunächst eine anschauliche Situation.

**Beispiel 3.1** (Schnitt zweier Geraden). Wir bestimmen grafisch den Schnittpunkt zweier Geraden im  $\mathbb{R}^2$ . Im ersten Fall (links im Bild) ist das Ergebnis, also die Lage des Schnittpunkts, vergleichsweise unempfindlich gegenüber Störungen in den Daten (der Eingabe) der Aufgabe.<sup>1</sup> Wir sprechen von einer **gut konditionierten Aufgabe**. Im zweiten Fall (rechts im Bild) ist dagegen das Ergebnis viel empfindlicher gegenüber Veränderungen in der Eingabe. Wir sprechen dann von einer **schlecht konditionierten Aufgabe**.



**Quizfrage:** Welche weiteren Beispiele für Aufgaben, für die Sie eine gute oder schlechte Kondition vermuten, fallen Ihnen ein?

Die **Sensitivitätsanalyse** (englisch: *sensitivity analysis*) beschäftigt sich mit der Frage, wie stark das Ergebnis  $y \in \mathbb{R}^m$  einer mathematischen Funktion

$$y = F(x)$$

von Änderungen bzw. Störungen in der Eingabe (den Daten)  $x \in \mathbb{R}^n$  abhängt. Die Darstellung  $y = F(x)$  bedeutet nicht, dass  $y$  notwendigerweise durch eine endliche Abfolge von Operationen explizit

<sup>1</sup>Die Störungen können etwa durch die begrenzte Zeichengenauigkeit auftreten.



berechenbar ist. Vielmehr kann  $y = F(x)$  zum Beispiel auch dafür stehen, dass  $y$  die Lösung eines nichtlinearen Gleichungssystems ist, in das  $x$  als Datum eingeht. Wie man aber bereits an einfachen Beispielen wie  $y = x^2$  sieht, wird die Antwort *lokal* sein, also abhängig davon, welchen Wert die Eingabe  $x$  selbst hat; vgl. [Abbildung 3.1](#).

Wir setzen voraus, dass  $F$  am interessierenden Punkt  $x$  differenzierbar ist. Dann erhalten wir mit Hilfe der Definition der Differenzierbarkeit, dass

$$F(x + \Delta x) \in F(x) + F'(x) \Delta x + o(\|\Delta x\|) \quad (3.1)$$

gilt.<sup>2</sup> Bezeichnen wir mit

$$\Delta y := F(x + \Delta x) - F(x)$$

die durch die Änderung  $\Delta x$  der Eingangsdaten hervorgerufenen Änderungen im Resultat an der Stelle  $x$ , so können wir (3.1) schreiben als

$$\Delta y \in F'(x) \Delta x + o(\|\Delta x\|). \quad (3.2)$$

Das heißt also, dass die durch die Ableitung  $F'(x)$  definierte lineare Abbildung  $\Delta x \mapsto F'(x) \Delta x$  eine derart gute Näherung für  $\Delta y$  ist, dass der dabei gemachte Fehler asymptotisch vernachlässigbar gegenüber  $\|\Delta x\|$  ist. Man schreibt statt (3.2) auch:

$$\Delta y \doteq F'(x) \Delta x \quad (3.3)$$

und sagt,  $\Delta y$  sei gleich  $F'(x) \Delta x$  «bis auf Terme, die gegenüber  $\|\Delta x\|$  vernachlässigbar sind» oder «bis auf Terme höherer Ordnung in  $\Delta x$ ». Letztere Sprechweise ist dadurch begründet, dass unter der Voraussetzung, dass  $F$  eine  $C^2$ -Funktion in einer Umgebung von  $x$  ist, aus dem [Satz von Taylor 2.11](#) folgt, dass statt (3.2) sogar gilt:  $\Delta y \in F'(x) \Delta x + O(\|\Delta x\|^2)$ . Der letzte Term ist dabei «von höherer Ordnung» (nämlich quadratisch) in  $\|\Delta x\|$ .

In der **differentiellen Sensitivitätsanalyse** (englisch: *differential sensitivity analysis*) vernachlässigt man tatsächlich die Restgliedterme und stützt sich allein auf die Ableitung  $F'(x)$ . Daraus wollen wir nun im Folgenden verschiedene Informationen gewinnen.

### § 3.1 ABSOLUTE BETRACHTUNGSWEISE

Zunächst lautet (3.3) in Komponenten ausgeschrieben für eine Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

$$\Delta y = \begin{pmatrix} \Delta y_1 \\ \vdots \\ \Delta y_m \end{pmatrix} \doteq F'(x) \Delta x = \begin{bmatrix} \frac{\partial F_1(x)}{\partial x_1} & \cdots & \frac{\partial F_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial F_m(x)}{\partial x_1} & \cdots & \frac{\partial F_m(x)}{\partial x_n} \end{bmatrix} \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{pmatrix}. \quad (3.4)$$

<sup>2</sup>**Beachte:** Auf der rechten Seite von (3.1) steht die Summe aus einer Funktion, nämlich  $F(x) + F'(x) \Delta x$ , und einer Menge von Funktionen  $O(\|\Delta x\|)$ . Das ist im Sinne der Minkowski-Summe zu verstehen, d. h., auf der rechten Seite steht die Menge aller Funktionen der Bauart  $\Delta x \mapsto F(x) + F'(x) \Delta x + G(\Delta x)$ , wobei  $G(\Delta x)$  zu  $O(\|\Delta x\|)$  gehört. Alternativ können wir (3.1) auch lesen als:  $F(x + \Delta x) - F(x) - F'(x) \Delta x \in O(\|\Delta x\|)$ .

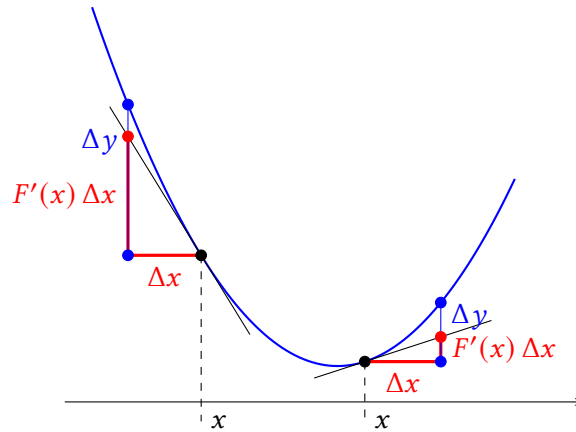


Abbildung 3.1: Wie stark das Ergebnis  $y = F(x) := x^2$  von Änderungen  $\Delta x$  in der Eingabe  $x$  abhängt, ist i. d. R. abhängig vom Wert von  $x$  abhängig.

Der Eintrag

$$K_{ij}(x) := \frac{\partial F_i(x)}{\partial x_j} \quad (3.5)$$

in Zeile  $i$  und Spalte  $j$  der Jacobimatrix gibt also den Faktor an, mit dem sich Änderungen  $\Delta x_j$  in der  $j$ -ten Komponente der Eingabe  $x$  auswirken auf die  $i$ -te Komponente des Ergebnisses.

**Definition 3.2** (Absolute partielle Konditionszahlen). Die Zahlen  $K_{ij}(x)$  heißen die **absoluten partiellen Konditionszahlen** (englisch: **absolute partial condition numbers**) der Funktion  $F$  an der Stelle  $x$ .

Wollen wir Störungen in **allen Komponenten der Eingabe** gemeinsam betrachten, aber uns im **Ergebnis** auf die  **$i$ -te Komponente** beschränken, also in (3.4) nur die  $i$ -Zeile

$$\Delta y_i \doteq \underbrace{\begin{bmatrix} \frac{\partial F_i(x)}{\partial x_1} & \dots & \frac{\partial F_i(x)}{\partial x_n} \end{bmatrix}}_{F'_i(x)} \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{pmatrix} = \sum_{j=1}^n \frac{\partial F_i(x)}{\partial x_j} \Delta x_j \quad (3.6)$$

herausgreifen, so können wir (3.6) mit Hilfe der Cauchy-Schwarz-Ungleichung abschätzen durch<sup>3</sup>

$$|\Delta y_i| \leq \|F'_i(x)\|_2 \|\Delta x\|_2. \quad (3.7)$$

Dabei ist  $\|F'_i(x)\|_2$  die 2-Norm des Zeilenvektors bestehend aus den partiellen Ableitungen der  $i$ -ten Komponentenfunktion  $F_i$  von  $F$ , also die Norm der  $i$ -ten Zeile der Jacobimatrix  $F'(x)$ . Diese Abschätzung gilt gleichmäßig, unabhängig von der Richtung des Störungsvektors  $\Delta x$ , es handelt sich also um eine *Worst-Case*-Abschätzung. Aus der Cauchy-Schwarz-Ungleichung wissen wir auch, dass die Abschätzung genau dann mit Gleichheit angenommen wird, also  $|\Delta y_i| \doteq \|F'_i(x)\|_2 \|\Delta x\|_2$  gilt, wenn die Störungsrichtung  $\Delta x$  und die Ableitung  $F'_i(x)^T$  linear abhängig sind, also einer dieser Vektoren ein Vielfaches des anderen ist.

<sup>3</sup>Auch hier bedeutet der Punkt auf dem Symbol  $\leq$ , dass die Abschätzung gilt bis auf Terme, die gegenüber  $\|\Delta x\|$  vernachlässigbar sind.

Möchte man umgekehrt die Auswirkung der Störung in nur **einer Komponente  $x_j$  der Eingabe**, aber gleichzeitig auf **alle Komponenten des Ergebnisses** betrachten, so setzen wir in (3.4) alle  $\Delta x_k = 0$  bis auf den Index  $k = j$  und erhalten

$$\Delta y \doteq \underbrace{\begin{bmatrix} \frac{\partial F_1(x)}{\partial x_j} \\ \vdots \\ \frac{\partial F_m(x)}{\partial x_j} \end{bmatrix}}_{\frac{\partial F(x)}{\partial x_j}} \Delta x_j \quad (3.8)$$

und daraus die Abschätzung

$$\|\Delta y\|_2 \leq \left\| \frac{\partial F(x)}{\partial x_j} \right\|_2 |\Delta x_j|. \quad (3.9)$$

Wollen wir nun abschließend gleichzeitig Störungen **aller Komponenten in der Eingabe** und die Auswirkung auf **alle Komponenten im Ergebnis** betrachten, so können wir (3.4) mit Hilfe von Lemma 2.3 abschätzen als

$$\|\Delta y\|_2 \leq \|F'(x)\|_{2 \rightarrow 2} \|\Delta x\|_2. \quad (3.10)$$

Dabei ist  $\|F'(x)\|_{2 \rightarrow 2}$  die durch die 2-Norm induzierte Operatornorm (Definition 2.2) der Jacobimatrix  $F'(x)$ . Die Abschätzung (3.10) besagt also, dass die Größe (Norm)  $\|\Delta y\|_2$  der Störung  $\Delta y$  im Ergebnis nicht größer sein kann als die Operatornorm der Jacobimatrix, multipliziert mit der Größe (Norm)  $\|\Delta x\|_2$  der Störung  $\Delta x$  in der Eingabe. Diese Abschätzung gilt wieder gleichmäßig, unabhängig von der Richtung des Störungsvektors  $\Delta x$ , es handelt sich also um eine *Worst-Case*-Abschätzung. In Abschnitt 4 werden wir sehen, für welche Richtungen  $\Delta x$  diese Abschätzung scharf ist.

**Definition 3.3** (Absolute Konditionszahl). *Es sei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  an der Stelle  $x$  differenzierbar.*

(i) *Dann heißt die Operatornorm*

$$K(x) := \|F'(x)\|_{2 \rightarrow 2} = \|(K_{ij}(x))\|_{2 \rightarrow 2} \quad (3.11)$$

*die **absolute Konditionszahl** (englisch: **absolute condition number**) von  $F$  an der Stelle  $x$ . Man spricht auch von der absoluten Konditionszahl der Aufgabe,  $y = F(x)$  auszuwerten.*

(ii) *Die Funktion  $F$  heißt an der Stelle  $x$  **schlecht konditioniert im absoluten Sinne**, wenn  $\|F'(x)\|_{2 \rightarrow 2}$  «sehr groß» ist, ansonsten **gut konditioniert im absoluten Sinne**.*

**Beachte:** Was «sehr groß» im Einzelfall bedeutet, wird von den Anforderungen der jeweiligen Aufgabe bestimmt und kann nicht allgemeingültig festgelegt werden. Man beachte dabei auch, dass sich beim Umskalieren der Eingabe oder des Resultats — z. B. durch Wechsel der physikalischen Einheit von m zu mm — die absolute Konditionszahl ändert.

**Beispiel 3.4** (Absolute Konditionierung). *Bei einer gleichmäßig beschleunigten Bewegung mit der Beschleunigung  $a$  und Anfangsgeschwindigkeit  $v_0$  legt ein Körper innerhalb der Zeit  $t$  die Strecke*

$$s(t, v_0, a) = v_0 t + \frac{1}{2} a t^2$$

zurück. Wir betrachten die Strecke  $s$  als Funktion der Eingangsgrößen  $x = (t, v_0, a)$ . Bei einer Beschleunigung von  $a = 5 \text{ m s}^{-2}$ , einer Anfangsgeschwindigkeit von  $v_0 = 0 \text{ m s}^{-1}$  und einer Zeit von  $t = 10 \text{ s}$  beispielsweise ist  $s = 250 \text{ m}$ . Die Jacobimatrix an dieser Stelle ist

$$\left[ \frac{\partial s}{\partial t}(t, v_0, a) \quad \frac{\partial s}{\partial v_0}(t, v_0, a) \quad \frac{\partial s}{\partial a}(t, v_0, a) \right] = \left[ v_0 + a t \quad t \quad \frac{1}{2} t^2 \right] = \left[ 50 \text{ m s}^{-1} \quad 10 \text{ s} \quad 50 \text{ s}^2 \right].$$

Die Daten  $x = (t, v_0, a)$  stammen vielleicht aus einer Messung und sind Unsicherheiten unterworfen. Änderungen einer einzelnen der drei Größen der Größenordnung  $\Delta t = 0.1 \text{ s}$ ,  $\Delta v_0 = 0.5 \text{ m s}^{-1}$  bzw.  $\Delta a = 0.2 \text{ m s}^{-2}$  rufen beispielsweise folgende Änderungen im Ergebnis hervor:

$$\begin{aligned} \Delta s &\doteq \frac{\partial s}{\partial t}(t, v_0, a) \Delta t = 50 \text{ m s}^{-1} \cdot 0.1 \text{ s} = 5 \text{ m}, \\ \Delta s &\doteq \frac{\partial s}{\partial v_0}(t, v_0, a) \Delta v_0 = 10 \text{ s} \cdot 0.5 \text{ m s}^{-1} = 5 \text{ m}, \\ \Delta s &\doteq \frac{\partial s}{\partial a}(t, v_0, a) \Delta a = 50 \text{ s}^2 \cdot 0.2 \text{ m s}^{-2} = 10 \text{ m}. \end{aligned}$$

**Quizfrage:** Wie ändern sich die Zahlenwerte der absoluten partiellen Konditionszahlen, wenn man die zurückgelegte Strecke  $s$  in mm misst statt wie oben in m? **Quizfrage:** Und was ändert sich, wenn man die Anfangsgeschwindigkeit  $v_0$  in  $\text{km h}^{-1}$  misst statt wie oben in  $\text{m s}^{-1}$ ? **Quizfrage:** Ist die Frage nach der absoluten Konditionszahl (3.11) überhaupt physikalisch sinnvoll?

**Beachte:** Im Fall von  $F: \mathbb{R} \rightarrow \mathbb{R}$  gibt es nur eine einzige absolute partielle Konditionszahl  $K_{11} = F'(x)$ , siehe (3.5). Diese ist vorzeichenbehaftet, während die absolute Konditionszahl nach Definition 3.3 immer nicht-negativ ist. Im betrachteten Fall gilt  $\|F'(x)\|_{2 \rightarrow 2} = |F'(x)| = |K_{11}|$ .

**Bemerkung 3.5** (Alternative Definition der absoluten Konditionszahl). In der Literatur findet man auch die folgende Definition<sup>4</sup> der absoluten Konditionszahl einer Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

$$\limsup_{\Delta x \rightarrow 0} \frac{\|F(x + \Delta x) - F(x)\|_2}{\|\Delta x\|_2} = \lim_{\substack{\varepsilon \searrow 0 \\ \Delta x \neq 0}} \sup_{\|\Delta x\|_2 < \varepsilon} \frac{\|F(x + \Delta x) - F(x)\|_2}{\|\Delta x\|_2}. \quad (3.12)$$

Bzgl. welcher Norm die Größe  $\Delta x$  im Supremum kleiner als  $\varepsilon$  genommen wird, ist hier unerheblich. Die Gleichheit der beiden Ausdrücke in (3.12) gilt nach Definition des Limes superior. Man betrachtet also beliebige Störungen  $\Delta x$  in immer kleiner werdenden Kugeln und bestimmt die größtmögliche Änderung, die solche Störungen im Ergebnis – bezogen auf ihre Größe  $\|\Delta x\|_2$  – hervorrufen können.

Man kann zeigen, dass (3.12) genau dann endlich ist, wenn  $F$  an der Stelle  $x$  differenzierbar ist. In diesem Fall stimmt (3.12) mit (3.11) aus Definition 3.3, also mit  $\|F'(x)\|_{2 \rightarrow 2}$ , überein. Wenn  $F$  an der Stelle  $x$  nicht differenzierbar ist, dann ist (3.12) gleich  $\infty$ .

**Beachte:** Die Wahl der  $2 \rightarrow 2$ -Norm (also der Wahl der 2-Norm in Zähler und Nenner von (3.12)) ist willkürlich. Stattdessen hätten wir z. B. auch die  $\infty \rightarrow \infty$ -Norm nehmen können, wie es im Anschluss bei der relativen Betrachtungsweise gleich der Fall sein wird. Aufgrund der Äquivalenz der Normen ändert sich dadurch die absolute Konditionszahl nur um einen (allerdings dimensionsabhängigen) konstanten Faktor.

<sup>4</sup>vgl. etwa Bornemann, 2018, Kapitel 11.2, S.41

### § 3.2 RELATIVE BETRACHTUNGSWEISE

Neben den bisherigen absoluten Betrachtungsweise kann man die Sensitivitäten auch relativ betrachten. Das bedeutet, dass wir die Störung in der Eingabe  $\Delta x_j$  auf  $x_j$  beziehen und die Änderung in der Ausgabe  $\Delta y_i$  auf  $y_i$ . Aus (3.4) wird also

$$\begin{pmatrix} \frac{\Delta y_1}{y_1} \\ \vdots \\ \frac{\Delta y_m}{y_m} \end{pmatrix} \doteq \begin{bmatrix} \frac{\partial F_1(x)}{\partial x_1} \cdot \frac{x_1}{y_1} & \dots & \frac{\partial F_1(x)}{\partial x_n} \cdot \frac{x_n}{y_1} \\ \vdots & & \vdots \\ \frac{\partial F_m(x)}{\partial x_1} \cdot \frac{x_1}{y_m} & \dots & \frac{\partial F_m(x)}{\partial x_n} \cdot \frac{x_n}{y_m} \end{bmatrix} \begin{pmatrix} \frac{\Delta x_1}{x_1} \\ \vdots \\ \frac{\Delta x_n}{x_n} \end{pmatrix}. \quad (3.13)$$

Die Größe

$$k_{ij}(x) := \frac{\partial F_i(x)}{\partial x_j} \cdot \frac{x_j}{y_i} = \frac{\partial F_i(x)}{\partial x_j} \cdot \frac{x_j}{F_i(x)} = K_{ij}(x) \cdot \frac{x_j}{F_i(x)} \quad (3.14)$$

gibt den Faktor an, mit dem sich relative Änderungen  $\Delta x_j/x_j$  in der  $j$ -ten Komponente der Eingabe relativ auswirken auf die  $i$ -te Komponente des Ergebnisses. Die relative Betrachtungsweise setzt natürlich voraus, dass die Komponenten  $y_i$  des Ergebnisses  $y = F(x)$  ungleich null sind und außerdem die Komponenten  $x_j$  der Eingabe ungleich null sind.

**Definition 3.6** (Relative partielle Konditionszahlen). Es sei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  an der Stelle  $x$  differenzierbar. Die Zahlen  $k_{ij}(x)$  heißen die **relativen partiellen Konditionszahlen**<sup>5</sup> (englisch: **relative partial condition numbers**) der Funktion  $F$  an der Stelle  $x$ .

In Analogie zu Definition 3.3 und Bemerkung 3.5 wollen wir noch eine **relative Konditionszahl**  $k(x)$  für die Auswertung von  $F$  an einer Stelle  $x$  angeben. Es stellt sich dabei – um die relative Konditionszahl aus den relativen partiellen Konditionszahlen  $k_{ij}(x)$  berechenbar zu machen – als praktisch heraus, bzgl.  $x$  komponentenweise vorzugehen.

Wir betrachten dazu zunächst eine reellwertige Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  und untersuchen die größtmögliche relative Änderung, die Störungen  $\Delta x$  im Ergebnis  $f(x)$  hervorrufen können, wenn sie klein werden. Die Größe der Störung  $\Delta x$  messen wir dabei wie gesagt erstens komponentenweise und zweitens relativ zur Größe von  $x$ . Es sei dazu  $D$  die Diagonalmatrix  $D = \text{diag}(x_1, \dots, x_n)$ . Dann hat  $D^{-1}\Delta x$  die Gestalt

$$D^{-1}\Delta x = \begin{pmatrix} \frac{\Delta x_1}{x_1} \\ \vdots \\ \frac{\Delta x_n}{x_n} \end{pmatrix}.$$

Aufgrund der komponentenweise Betrachtung von  $\Delta x$  messen wir die Größe dieses Ausdrucks als

$$\|D^{-1}\Delta x\|_\infty = \max \left\{ \left| \frac{\Delta x_i}{x_i} \right| \mid i = 1, \dots, n \right\}.$$

Die obige Frage nach der größtmöglichen relativen Änderung in  $f(x)$  durch Störungen  $\Delta x$  führt uns nun zu dem Ausdruck

$$\limsup_{\Delta x \rightarrow 0} \frac{|f(x + \Delta x) - f(x)|}{|f(x)| \|D^{-1}\Delta x\|_\infty},$$

<sup>5</sup>In den Wirtschaftswissenschaften werden die relativen partiellen Konditionszahlen  $k_{ij}(x)$  auch als **Elastizitäten** bezeichnet.

den wir jetzt untersuchen. Es gilt

$$\begin{aligned} \limsup_{\Delta x \rightarrow 0} \frac{|f(x + \Delta x) - f(x)|}{|f(x)| \|D^{-1} \Delta x\|_\infty} \\ = \lim_{\varepsilon \searrow 0} \sup_{\substack{\|\Delta x\| < \varepsilon \\ \Delta x \neq 0}} \frac{|f(x + \Delta x) - f(x)|}{|f(x)| \|D^{-1} \Delta x\|_\infty} \end{aligned}$$

Bzgl. welcher Norm  $\Delta x$  im Supremum kleiner als  $\varepsilon$  genommen wird, ist hier unerheblich. Weiter:

$$\begin{aligned} &= \lim_{\varepsilon \searrow 0} \sup_{\substack{\|\Delta x\| < \varepsilon \\ \Delta x \neq 0}} \frac{|f'(x) \Delta x|}{|f(x)| \|D^{-1} \Delta x\|_\infty} \\ &= \sup_{\Delta x \neq 0} \frac{|f'(x) \Delta x|}{|f(x)| \|D^{-1} \Delta x\|_\infty} \\ &= \sup_{\Delta y \neq 0} \frac{|f'(x) D \Delta y|}{|f(x)| \|\Delta y\|_\infty} \quad (\text{setze } \Delta y = D^{-1} \Delta x) \\ &= \frac{\|D f'(x)^\top\|_1}{|f(x)|} \quad (\text{wegen (2.5) mit der } \|\cdot\|_{\infty \rightarrow \infty}\text{-Norm des Zeilenvektors } f'(x) D) \\ &= \sum_{j=1}^n \left| \frac{\partial f(x)}{\partial x_j} \right| \frac{|x_j|}{|f(x)|} = \sum_{j=1}^n |k_{1j}|. \end{aligned}$$

Ist nun  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  vektorwertig, so betrachten wir die obige Rechnung komponentenweise in  $F_i$ . Wir erhalten ganz analog wie oben

$$\max_{i=1, \dots, m} \limsup_{\Delta x \rightarrow 0} \frac{|F_i(x + \Delta x) - F_i(x)|}{|F_i(x)| \|D^{-1} \Delta x\|_\infty} = \max_{i=1, \dots, m} \sum_{j=1}^n \left| \frac{\partial f(x)}{\partial x_j} \right| \frac{|x_j|}{|f(x)|} = \max_{i=1, \dots, m} \sum_{j=1}^n |k_{ij}(x)| \quad (3.15)$$

als Maß für die relative Änderung im Ergebnis. Das ist gerade die Zeilensummennorm der Matrix, die aus den Einträgen  $k_{ij}(x)$  besteht.

**Quizfrage:** Wie würde dieser Ausdruck aussehen, wenn wir statt der  $\infty$ -Normen überall die 2-Normen verwenden würden?

**Definition 3.7** (Relative Konditionszahl). Es sei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  an der Stelle  $x$  differenzierbar. Weiter sei  $x \neq 0$  und  $F(x) \neq 0$ .

(i) Dann heißt der Ausdruck

$$k(x) := \|(k_{ij}(x))\|_{\infty \rightarrow \infty} = \max_{i=1, \dots, m} \sum_{j=1}^n |k_{ij}(x)| \quad (3.16)$$

die **relative Konditionszahl** (englisch: **relative condition number**) von  $F$  an der Stelle  $x$ .

(ii) Die Funktion  $F$  heißt an der Stelle  $x$  **schlecht konditioniert im relativen Sinne**, wenn (3.16)  $\gg 1$  ist, ansonsten **gut konditioniert im relativen Sinne**.

**Beispiel 3.8** (Relative Konditionierung). Wir greifen das [Beispiel 3.4](#) nochmals auf, dieses Mal aber in relativer Betrachtungsweise. Die Matrix der relativen partiellen Konditionszahlen der Funktion  $s(t, v_0, a)$  an der Stelle  $x = (t, v_0, a) = (10 \text{ s}, 0 \text{ m s}^{-1}, 5 \text{ m s}^{-2})$  ist gleich

$$\begin{bmatrix} 50 \text{ m s}^{-1} \cdot \frac{10 \text{ s}}{250 \text{ m}} & 10 \text{ s} \cdot \frac{0 \text{ m s}^{-1}}{250 \text{ m}} & 50 \text{ s}^2 \cdot \frac{5 \text{ m s}^{-2}}{250 \text{ m}} \end{bmatrix} = \begin{bmatrix} 2 & 0 & 1 \end{bmatrix}.$$

Diese Zahlen sind wie erwartet einheitenlos.

Eine Änderung von  $\Delta t = 0.1 \text{ s}$  wie in [Beispiel 3.4](#) entspricht in relativer Betrachtungsweise gegenüber dem ungestörten Wert von  $t = 10 \text{ s}$  einer Änderung von  $\frac{0.1 \text{ s}}{10 \text{ s}} = 0.01$ . Diese erzeugt (bis auf Terme höherer Ordnung) eine relative Änderung von  $2 \cdot 0.01 = 0.02$  im Ergebnis, also absolut gesehen  $0.02 \cdot 250 \text{ m} = 5 \text{ m}$ . Das stimmt (wie erwartet) mit dem aus [Beispiel 3.4](#) bekannten Ergebnis überein. Analog ergibt die Änderung von  $\Delta a = 0.2 \text{ m s}^{-2}$ , also in relativer Betrachtung von  $\frac{0.2 \text{ m s}^{-2}}{5 \text{ m s}^{-2}} = 0.04$ , eine relative Änderung von  $1 \cdot 0.04 = 0.04$  im Ergebnis, also absolut gesehen  $0.04 \cdot 250 \text{ m} = 10 \text{ m}$ , was wiederum mit dem aus [Beispiel 3.4](#) bekannten Ergebnis übereinstimmt.

Die relative Konditionszahl bzgl. der Anfangsgeschwindigkeit  $v_0$  ist an der betrachteten Stelle  $x$  ohne Aussagekraft, da  $v_0 = 0 \text{ m s}^{-1}$  ist. Änderungen  $\Delta v_0$  in  $v_0$  lassen sich also nicht relativ ausdrücken, vgl. (3.13).

In relativer Betrachtungsweise lassen sich in jedem Fall die partiellen relativen Konditionszahlen  $k_{ij}(x)$  in physikalisch sinnvoller Weise zur relativen Konditionszahl gemäß (3.16) zusammenfassen, da die  $k_{ij}(x)$  immer einheitenlos sind. **Quizfrage:** Wie ändern sich die Zahlenwerte der relativen partiellen Konditionszahlen, wenn man die zurückgelegte Strecke  $s$  in mm misst statt wie oben in m? **Quizfrage:** Und was ändert sich, wenn man die Anfangsgeschwindigkeit  $v_0$  in  $\text{km h}^{-1}$  misst statt wie oben in  $\text{m s}^{-1}$ ? **Quizfrage:** Wie groß ist die relative Konditionszahl (3.16) in diesem Beispiel?

Im weiteren Verlauf der Vorlesung werden vor allem relative Konditionszahlen betrachtet. Das liegt (neben der Tatsache, dass sie vor allem bei Einbeziehung von Einheiten sinnvoller ist) auch daran, dass später wesentliche Fehlereinflüsse durch die Zahldarstellung im Computer eingebracht werden, die gut im relativen Sinn gemessen werden können, siehe [Kapitel 2](#).

Wir geben nun die absoluten und relativen partiellen Konditionszahlen der Grundrechenarten und einiger weiterer elementarer Operationen an. Es handelt sich dabei jeweils um Funktionen  $F: \mathbb{R}^2 \rightarrow \mathbb{R}$  bzw.  $F: \mathbb{R} \rightarrow \mathbb{R}$ .

**Lemma 3.9** (Absolute und relative partielle Konditionszahlen für einige Grundoperationen). Die absoluten und relativen Konditionszahlen (im Fall von  $x \neq 0$ ) der folgenden Grundoperationen sind (auf ihren jeweiligen Definitionsbereichen) gegeben durch

Operation	$K_{11}(x)$	$K_{12}(x)$	$k_{11}(x)$	$k_{12}(x)$	$k(x)$
$x_1 + x_2$	1	1	$\frac{x_1}{x_1+x_2}$	$\frac{x_2}{x_1+x_2}$	$\frac{ x_1 + x_2 }{ x_1+x_2 }$
$x_1 - x_2$	1	-1	$\frac{x_1}{x_1-x_2}$	$\frac{-x_2}{x_1-x_2}$	$\frac{ x_1 + x_2 }{ x_1-x_2 }$
$x_1 \cdot x_2$	$x_2$	$x_1$	1	1	2
$\frac{x_1}{x_2}$	$\frac{1}{x_2}$	$-\frac{x_1}{x_2^2}$	1	-1	2
$\frac{1}{x_1}$	$-\frac{1}{x_1^2}$		-1		1
$\sqrt{x_1}$	$\frac{1}{2\sqrt{x_1}}$		$\frac{1}{2}$		$\frac{1}{2}$

Beweis. Nachrechnen

□

**Bemerkung 3.10** (Bedeutung von Lemma 3.9).

- (i) Die Ergebnisse von Lemma 3.9 besagen, dass die Multiplikation, Division, Kehrwerte und Quadratwurzeln auf ihrem gesamten Definitionsbereich im relativen Sinne überall gut konditioniert sind.
- (ii) Für die Addition ist die relative Konditionszahl  $k(x) = 1$ , sofern  $x_1$  und  $x_2$  dasselbe Vorzeichen haben. Man bezeichnet diesen Fall auch als **echte Addition**. Dort wo  $x_1$  und  $x_2$  jedoch verschiedene Vorzeichen haben (also  $x_1 + x_2$  einer **echten Subtraktion** zweier positiver Zahlen entspricht), ist die relative Konditionszahl  $k(x) > 1$ . Für  $x_1 \approx -x_2$  kann die relative Konditionszahl beliebig groß werden, siehe Abbildung 3.2. Beispielsweise für festes  $x_1 > 0$  und  $x(\varepsilon) = (x_1, \varepsilon - x_2)^T$  mit  $\varepsilon > 0$  ist  $k(x(\varepsilon)) \in O(\varepsilon^{-1})$ . Ein analoges Resultat gilt natürlich auch für die Subtraktion zweier Zahlen, wobei dann der Fall kritisch ist, bei dem beide dasselbe Vorzeichen haben.

**Beachte:** Die echte Subtraktion zweier fast gleicher Zahlen ist schlecht konditioniert!

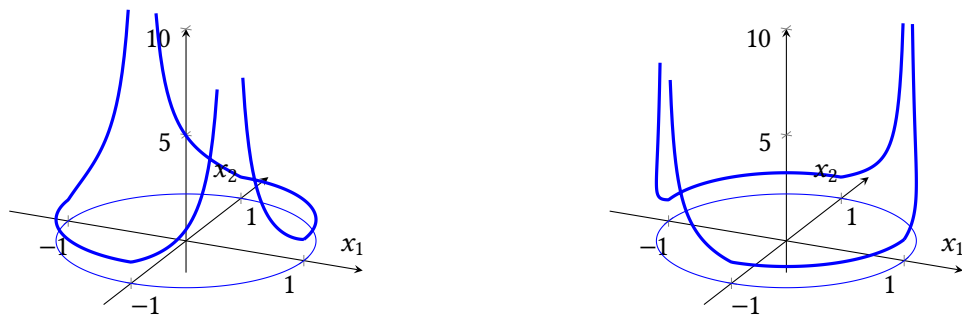


Abbildung 3.2: Relative Konditionszahl  $k(x)$  der Addition (links) und der Subtraktion (rechts) zweier Zahlen  $x_1, x_2$  auf dem Einheitskreis.

Das folgende Beispiel wurde nachträglich eingefügt.



**Beispiel 3.11** (Auslöschung). Die Subtraktion von

$$x_1 = 1.000\,001$$

$$x_2 = 1.000\,000$$

ergibt das Ergebnis  $y = x_1 - x_2 = 10^{-6}$ . Die relative partielle Konditionszahl  $k_{11}(x)$  ist

$$k_{11}(x) = \frac{x_1}{x_1 - x_2} \approx 1.0 \cdot 10^6.$$

Änderungen etwa von  $\Delta x_1 = 10^{-6}$  ergeben eine Änderung von  $\Delta y = 10^{-6}$  im Ergebnis. Relativ gesehen sind diese Änderungen  $\Delta x_1/x_1 \approx 10^{-6}$  und  $\Delta y/y \approx 10^{-6}/10^{-6} = 1$ . Wie erwartet verstärkt sich die relative Störung in der Eingabe um den Faktor  $k_{11}(x) \approx 10^6$ .

Nochmal konkret: Statt  $y = 1 \cdot 10^{-6}$  ergibt sich bei Änderung von  $x_1$  zu  $x_1 + \Delta x_1$  nun das Ergebnis  $y + \Delta y = 2 \cdot 10^{-6}$ . Während sich also bei  $x_1$  die siebte Dezimalstelle ändert, führt dies im Ergebnis  $y$  zu einer Änderung bereits der ersten Dezimalstelle. Wir verlieren damit sechs Dezimalstellen an Genauigkeit! Diesen Effekt nennt man **Auslöschung** (englisch: **cancellation**).

**Beachte:** Gegenüber den Daten verliert man (im schlechtesten Fall) etwa  $\log_{10} k(x)$  Dezimalziffern an Genauigkeit, wobei  $k(x)$  die relative Konditionszahl der betrachteten Operation darstellt.

Der Klarheit halber liefern wir noch eine genaue Definition der Begriffe **echte Addition** und **echte Subtraktion** nach:

**Definition 3.12** (Echte Addition, echte Subtraktion).

- (i) Wir bezeichnen die Operation  $x_1 + x_2$  als **echte Addition** der Zahlen  $x_1$  und  $x_2$ , wenn  $|x_1 + x_2| = |x_1| + |x_2|$  gilt. Andernfalls heißt sie **echte Subtraktion**.
- (ii) Wir bezeichnen die Operation  $x_1 - x_2$  als **echte Addition** der Zahlen  $x_1$  und  $x_2$ , wenn  $|x_1 - x_2| = |x_1| + |x_2|$  gilt. Andernfalls heißt sie **echte Subtraktion**.

Das folgende Beispiel wurde nachträglich eingefügt.

**Beispiel 3.13** (Echte Addition, echte Subtraktion).

- (i) Die folgenden Operationen sind echte Additionen:

$$3 + 5, \quad -3 - 5$$

- (ii) Die folgenden Operationen sind echte Subtraktionen:

$$-3 + 5, \quad 5 - 3$$

### § 3.3 KONDITION LINEARER GLEICHUNGSSYSTEME

Lineare Gleichungssysteme sind eine Grundaufgabe in der angewandten und insbesondere in der numerischen Mathematik. Wir betrachten hier Systeme

$$Ax = b \quad \Leftrightarrow \quad x = A^{-1}b$$

mit invertierbaren Matrizen  $A \in \mathbb{R}^{n \times n}$  und rechten Seiten  $b \in \mathbb{R}^n$ . Zugunsten der für Gleichungssysteme üblichen Notation müssen wir in unserem allgemeinen Setting  $y = F(x)$  also entsprechende Umbenennungen vornehmen.

Wir betrachten zunächst Störungen  $\Delta b$  der rechten Seite  $b$ . Dann übernimmt  $b$  die Rolle der Eingabe  $x$ , und die Lösung  $x$  übernimmt die Rolle des Ergebnisses  $y$ . Statt  $y = F(x)$  haben wir also  $x = A^{-1}b$ . Diese Funktion ist linear in der Eingabe  $b$ , sodass wir an Stelle der Abschätzung (3.3) sogar die Gleichheit

$$\Delta x = A^{-1} \Delta b \tag{3.17}$$

bekommen. Hier könnten wir jetzt komponentenweise wie in Abschnitt 3.2 vorgehen. Es zeigt sich aber, dass wir das einfachere Resultat erhalten, wenn wir an Stelle der Norm der Vektoren der relativen Änderungen  $(\Delta b_i/b_i)$  und  $(\Delta x_i/x_i)$  nur die relativen Änderungen in den Normen  $\|\Delta b\|/\|b\|$  und  $\|\Delta x\|/\|x\|$  betrachten und dabei jeweils die 2-Norm verwenden. Die zugehörigen Abschätzungen

$$\begin{aligned} \|\Delta x\|_2 &= \|A^{-1} \Delta b\|_2 \leq \|A^{-1}\|_{2 \rightarrow 2} \|\Delta b\|_2, \\ \|b\|_2 &= \|Ax\|_2 \leq \|A\|_{2 \rightarrow 2} \|x\|_2 \end{aligned}$$

lassen sich kombinieren zu

$$\frac{\|\Delta x\|_2}{\|x\|_2} \leq \|A\|_{2 \rightarrow 2} \|A^{-1}\|_{2 \rightarrow 2} \frac{\|\Delta b\|_2}{\|b\|_2}. \tag{3.18}$$

Wir betrachten nun Störungen  $\Delta A$  der Matrix  $A$ . Jetzt übernimmt also  $A$  die Rolle der Eingabe  $x$ . Die differentielle Sensitivitätsanalyse, siehe (3.3), verlangt nun die Bestimmung der Ableitung der Funktion  $A \mapsto F(A) := A^{-1}b$ , also insbesondere die Bestimmung der Ableitung der Matrixinversion.

**Lemma 3.14** (Ableitung der Matrixinversion). *Es sei  $\Phi(A) := A^{-1}$  die Matrixinversions-Funktion. Diese ist für jede reguläre (englisch: **non-singular**) Matrix  $A \in \mathbb{R}^{n \times n}$  in einer ganzen Umgebung von  $A$  definiert und dort stetig differenzierbar.<sup>6</sup> Es gilt die Ableitungsformel*

$$\Phi'(A) \Delta A = -A^{-1} \Delta A A^{-1} \quad \text{für alle regulären } A \in \mathbb{R}^{n \times n}. \tag{3.19}$$

*Beweis.* Wir wollen den Satz über implizite Funktionen anwenden auf die implizite Gleichung

$$A \Phi(A) = \text{Id}$$

für  $\Phi(A)$ . Die Differentiation dieser Beziehung in Richtung einer beliebigen Matrix  $\Delta A \in \mathbb{R}^{n \times n}$  ergibt nach Produktregel

$$\Delta A \Phi(A) + A (\Phi'(A) \Delta A) = 0_{n \times n}. \tag{*}$$

<sup>6</sup> $\Phi$  ist sogar eine  $C^\infty$ -Funktion, was wir hier aber nicht benötigen.

Durch Umstellen erhalten wir zunächst

$$A (\Phi'(A) \Delta A) = -\Delta A \Phi(A)$$

und dann durch Multiplikation von links mit  $A^{-1}$

$$\Phi'(A) \Delta A = -A^{-1} \Delta A \Phi(A) = -A^{-1} \Delta A A^{-1}.$$

Da also die Auflösung der Gleichung (\*) nach der gesuchten Richtungsableitung  $\Phi'(A) \Delta A$  eindeutig möglich war, können wir aus dem Satz über implizite Funktionen schließen, dass die Funktion  $\Phi$  in einer Umgebung von  $A$  tatsächlich definiert und stetig differenzierbar ist und dass die Richtungsableitungen  $\Phi'(A) \Delta A$  durch (3.19) gegeben sind.  $\square$

Mit dem Wissen aus Lemma 3.14 können wir nun die Gleichung (3.3) für unsere Aufgabe mit  $F(A) = A^{-1}b$  auswerten, also

$$\Delta x \doteq F'(A) \Delta A = [\Phi'(A) \Delta A] b = -A^{-1} \Delta A A^{-1} b = -A^{-1} \Delta A x. \quad (3.20)$$

Wir schätzen ähnlich wie oben ab:

$$\|\Delta x\|_2 = \|A^{-1} \Delta A x\|_2 \leq \|A^{-1}\|_{2 \rightarrow 2} \|\Delta A\|_{2 \rightarrow 2} \|x\|_2$$

und erhalten schließlich

$$\frac{\|\Delta x\|_2}{\|x\|_2} \leq \|A^{-1}\|_{2 \rightarrow 2} \|A\|_{2 \rightarrow 2} \frac{\|\Delta A\|_{2 \rightarrow 2}}{\|A\|_{2 \rightarrow 2}}. \quad (3.21)$$

**Definition 3.15** (Konditionszahl einer Matrix). *Es sei  $A \in \mathbb{R}^{n \times n}$ . Wenn  $A$  regulär ist, dann bezeichnen wir die Zahl*

$$\kappa(A) := \|A\|_{2 \rightarrow 2} \|A^{-1}\|_{2 \rightarrow 2} \quad (3.22)$$

*als die **Konditionszahl der Matrix**  $A$ . Wenn  $A$  singular ist, dann setzen wir die Konditionszahl auf  $\kappa(A) := \infty$ .*

Nach (3.18) und (3.21) gibt die Konditionszahl eine obere Schranke für den Sensitivitätsfaktor an, mit dem sich eine relative Änderung der 2-Norm der rechten Seite bzw. der  $2 \rightarrow 2$ -Norm der Matrix überträgt auf die relative Änderung der 2-Norm der Lösung. In Abschnitt 4 kommen wir nochmal darauf und auch auf Möglichkeiten zur Berechnung von  $\kappa(A)$  zurück.

Die Konditionszahl einer Matrix hat aber noch eine andere Interpretation:

**Satz 3.16** (Kahan, 1966, vgl. Bornemann, 2018, Kapitel 11.9). *Es sei  $A \in \mathbb{R}^{n \times n}$  regulär. Dann gilt*

$$\frac{1}{\kappa(A)} = \min \left\{ \frac{\|E\|_{2 \rightarrow 2}}{\|A\|_{2 \rightarrow 2}} \mid E \in \mathbb{R}^{n \times n}, A + E \text{ ist singular} \right\}. \quad (3.23)$$

Der Kehrwert der Konditionszahl einer Matrix ist also ihr relativer Abstand zur Menge der singulären Matrizen in der  $2 \rightarrow 2$ -Norm.

*Beweis.* Wenn  $A + E$  eine singuläre Matrix ist, dann gibt es ein  $x \neq 0$  mit der Eigenschaft  $(A + E)x = 0$ , also auch  $A^{-1}(A + E)x = x + A^{-1}Ex = 0$ . Es gilt also

$$0 < \|x\|_2 = \|A^{-1}Ex\|_2 \leq \|A^{-1}\|_{2 \rightarrow 2} \|E\|_{2 \rightarrow 2} \|x\|_2 = \kappa(A) \frac{\|E\|_{2 \rightarrow 2}}{\|A\|_{2 \rightarrow 2}} \|x\|_2.$$

Daraus folgt

$$\frac{1}{\kappa(A)} \leq \frac{\|E\|_{2 \rightarrow 2}}{\|A\|_{2 \rightarrow 2}} \frac{\|x\|_2}{\|x\|_2}.$$

Um den Beweis abzuschließen, geben wir jetzt eine spezielle Matrix  $E$  an, für die hier die Gleichheit gilt, also  $\kappa(A) = \frac{\|A\|_{2 \rightarrow 2}}{\|E\|_{2 \rightarrow 2}}$  oder äquivalent  $\|E\|_{2 \rightarrow 2} = 1/\|A^{-1}\|_{2 \rightarrow 2}$ . Dazu sei  $y \in \mathbb{R}^n$  ein Vektor mit  $\|y\|_2 = 1$ , der das Maximum in der Definition der  $2 \rightarrow 2$ -Operatornorm für  $A^{-1}$  realisiert, also

$$\|A^{-1}\|_{2 \rightarrow 2} = \max_{\substack{z \in \mathbb{R}^n \\ \|z\|_2=1}} \|A^{-1}z\|_2 = \|A^{-1}y\|_2,$$

vgl. Lemma 2.3 (iii). Außerdem sei  $x := A^{-1}y \neq 0$ . Wir wählen  $E := -\frac{yx^T}{\|x\|_2^2}$ . Es gilt

$$\|E\|_{2 \rightarrow 2} = \frac{\|x\|_2 \|y\|_2}{\|x\|_2^2} = \frac{1}{\|x\|_2} = \frac{1}{\|A^{-1}y\|_2} = \frac{1}{\|A^{-1}\|_{2 \rightarrow 2}}.$$

Für den Beweis der ersten Gleichheit (Spektralnorm für Rang-1-Matrizen) betrachten wir

$$\begin{aligned} \|E\|_{2 \rightarrow 2} &= \max_{\substack{z \in \mathbb{R}^n \\ \|z\|_2=1}} \|Ez\|_2 \quad (\text{Lemma 2.3}) \\ &= \max_{\substack{z \in \mathbb{R}^n \\ \|z\|_2=1}} \frac{\|yx^Tz\|_2}{\|x\|_2^2} \quad (\text{Definition von } E) \\ &= \max_{\substack{z \in \mathbb{R}^n \\ \|z\|_2=1}} \frac{\|y\|_2}{\|x\|_2^2} |x^Tz| \quad (\text{positive Homogenität der Norm}) \\ &= \frac{\|y\|_2}{\|x\|_2^2} \frac{x^Tx}{\|x\|_2} \quad (z = x/\|x\|_2 \text{ maximiert den Ausdruck, Cauchy-Schwarz-Ungleichung}) \\ &= \frac{\|x\|_2 \|y\|_2}{\|x\|_2^2}. \end{aligned}$$

□

### § 3.4 KONDITION ZUSAMMENGESETZTER FUNKTIONEN

In diesem Abschnitt geht es darum, dass wir eine Funktionsauswertung  $y = F(x)$  für  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  in kleinere Einheiten zerlegen, im einfachsten Fall

$$y = F(x) = (H \circ G)(x) = H(G(x))$$

mit Funktionen

$$G: \mathbb{R}^n \rightarrow \mathbb{R}^\ell \quad \text{und} \quad H: \mathbb{R}^\ell \rightarrow \mathbb{R}^m.$$

Die Funktionen  $G$  und  $H$  seien an der Stelle  $x$  bzw. an der Stelle  $G(x)$  differenzierbar. Nach Kettenregel gilt für die differentielle Störungsanalyse, vgl. (3.3),

$$\Delta y \doteq (H \circ G)'(x) \Delta x = H'(G(x)) G'(x) \Delta x. \quad (3.24)$$

Wir können also die absolute Konditionszahl von  $F$  an der Stelle  $x$  mit Hilfe von Lemma 2.3 (v) (Submultiplikativität der Matrixnormen) abschätzen durch

$$K(x) = \|F'(x)\|_{2 \rightarrow 2} \leq \|H'(G(x))\|_{2 \rightarrow 2} \|G'(x)\|_{2 \rightarrow 2}. \quad (3.25)$$

Eine solche Abschätzung wird i. A. nicht scharf sein (**Quizfrage:** Warum?) Dennoch ist die Betrachtung später hilfreich für die Stabilitätsanalyse von Algorithmen.

Wir können eine zu (3.24) analoge Beziehung auch für die relative Sensitivitätsanalyse finden. Wir schreiben dazu (3.24) einmal komponentenweise auf:

$$\begin{pmatrix} \Delta y_1 \\ \vdots \\ \Delta y_m \end{pmatrix} \doteq \begin{bmatrix} \frac{\partial H_1(G(x))}{\partial z_1} & \dots & \frac{\partial H_1(G(x))}{\partial z_\ell} \\ \vdots & & \vdots \\ \frac{\partial H_m(G(x))}{\partial z_1} & \dots & \frac{\partial H_m(G(x))}{\partial z_\ell} \end{bmatrix} \begin{bmatrix} \frac{\partial G_1(x)}{\partial x_1} & \dots & \frac{\partial G_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial G_\ell(x)}{\partial x_1} & \dots & \frac{\partial G_\ell(x)}{\partial x_n} \end{bmatrix} \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{pmatrix}$$

und ergänzen wie in (3.13) die Terme, die zu einer relativen Betrachtungsweise führen:

$$\begin{pmatrix} \frac{\Delta y_1}{y_1} \\ \vdots \\ \frac{\Delta y_m}{y_m} \end{pmatrix} \doteq \underbrace{\begin{bmatrix} \frac{\partial H_1(G(x))}{\partial z_1} \cdot \frac{z_1}{y_1} & \dots & \frac{\partial H_1(G(x))}{\partial z_\ell} \cdot \frac{z_\ell}{y_1} \\ \vdots & & \vdots \\ \frac{\partial H_m(G(x))}{\partial z_1} \cdot \frac{z_1}{y_m} & \dots & \frac{\partial H_m(G(x))}{\partial z_\ell} \cdot \frac{z_\ell}{y_m} \end{bmatrix}}_{=:(k_{ik}^H(G(x)))} \underbrace{\begin{bmatrix} \frac{\partial G_1(x)}{\partial x_1} \cdot \frac{x_1}{z_1} & \dots & \frac{\partial G_1(x)}{\partial x_n} \cdot \frac{x_n}{z_\ell} \\ \vdots & & \vdots \\ \frac{\partial G_\ell(x)}{\partial x_1} \cdot \frac{x_1}{z_1} & \dots & \frac{\partial G_\ell(x)}{\partial x_n} \cdot \frac{x_n}{z_\ell} \end{bmatrix}}_{=:(k_{kj}^G(x))} \begin{pmatrix} \frac{\Delta x_1}{x_1} \\ \vdots \\ \frac{\Delta x_n}{x_n} \end{pmatrix} \quad (3.26)$$

Auch hier können wir also wieder die Submultiplikativität der Matrixnormen (Lemma 2.3 (v)) verwenden, um abzuschätzen:

$$\left\| \begin{pmatrix} \frac{\Delta y_1}{y_1} \\ \vdots \\ \frac{\Delta y_m}{y_m} \end{pmatrix} \right\|_\infty \leq \| (k_{ik}^H(G(x))) \|_{\infty \rightarrow \infty} \| (k_{kj}^G(x)) \|_{\infty \rightarrow \infty} \left\| \begin{pmatrix} \frac{\Delta x_1}{x_1} \\ \vdots \\ \frac{\Delta x_n}{x_n} \end{pmatrix} \right\|_\infty. \quad (3.27)$$

Auch diese Abschätzung ist i. A. nicht scharf, aber eine ähnliche Betrachtung wird später bei der Stabilitätsanalyse von Algorithmen noch nützlich sein.

**Bemerkung 3.17.** Wir betonen nochmal zur Verdeutlichung, dass die Kondition eine Eigenschaft der betrachteten Aufgabe ist. Sie hat nichts damit zu tun, mit welchem Verfahren wir einmal diese Aufgabe versuchen werden zu lösen.

Bei einer schlecht konditionierten Aufgabe muss man sich die Frage stellen, ob die Aufgabe überhaupt sachgerecht formuliert ist.

## § 4 SINGULÄRWERTZERLEGUNG VON MATRIZEN

### § 4.1 SPEKTRALZERLEGUNG

Eine Zahl  $\lambda \in \mathbb{C}$  heißt **Eigenwert** (englisch: *eigenvalue*) und ein Vektor  $v \in \mathbb{C}^n \setminus \{0\}$  zugehöriger **Eigenvektor** (englisch: *eigenvector*) einer Matrix  $A \in \mathbb{R}^{n \times n}$  (oder  $A \in \mathbb{C}^{n \times n}$ ), wenn gilt:

$$Av = \lambda v, \quad (4.1)$$

d. h., dass  $v$  unter der Abbildung  $v \mapsto Av$  auf ein Vielfaches von sich selbst abgebildet wird. Zusammen heißt  $(\lambda, v)$  ein **Eigenpaar** von  $A$ . Im Allgemeinen muss man damit rechnen, dass Eigenwerte und Eigenvektoren komplex sind, selbst wenn  $A$  eine reelle Matrix ist.

In der linearen Algebra untersucht man u. a. die Frage, unter welchen Umständen eine Basis aus lauter Eigenvektoren von  $A$  existiert und wann diese als Orthonormalbasis gewählt werden kann.

**Definition 4.1** (Orthogonale Matrix). Eine Matrix  $V \in \mathbb{R}^{n \times n}$  heißt **orthogonal**, wenn  $V^T V = \text{Id}$  ist.

Diese Bedingung bedeutet, dass die Spaltenvektoren  $v_i$  von  $V$  paarweise senkrecht (bzgl. des Euklidischen Innenprodukts) aufeinander stehen und außerdem  $\|v_i\|_2 = 1$  ist, denn es gilt

$$(V^T V)_{ij} = v_i^T v_j = \begin{cases} \|v_i\|_2^2 = 1 & \text{im Fall } i = j \\ 0 & \text{im Fall } i \neq j. \end{cases}$$

Eine Matrix  $V \in \mathbb{R}^{n \times n}$  ist also genau dann orthogonal, wenn ihre Spaltenvektoren eine (reelle) **Orthonormalbasis** von  $\mathbb{R}^n$  bilden, also eine Menge linear unabhängiger Vektoren, die paarweise senkrecht stehen und deren 2-Norm gleich eins ist. Für orthogonale Matrizen gilt offenbar  $V^{-1} = V^T$  und daher auch  $VV^T = \text{Id}$ .

Außerdem ist das Produkt orthogonaler Matrizen wieder eine orthogonale Matrix.<sup>7</sup>

**Beachte:** Der Begriff «orthogonale Matrix» ist für reelle Matrizen reserviert. Das komplexe Analogon heißt **unitäre Matrix**.

Wir interessieren uns hier nur für die Frage, wann eine Matrix  $A \in \mathbb{R}^{n \times n}$  die Eigenschaft besitzt, eine Orthonormalbasis von  $\mathbb{R}^n$  aus lauter reellen Eigenvektoren zusammenstellen zu können. Die Antwort gibt folgender Satz.

**Satz 4.2** (Spektralsatz). Es sei  $A \in \mathbb{R}^{n \times n}$ . Der Raum  $\mathbb{R}^n$  besitzt genau dann eine reelle Orthonormalbasis aus lauter Eigenvektoren von  $A$ , wenn  $A = A^T$  ist. In diesem Fall besitzt  $A$  die folgende **Spektralzerlegung** (englisch: *spectral decomposition*):

$$A = V \Lambda V^T, \quad (4.2)$$

<sup>7</sup>D. h., die orthogonalen Matrizen bilden die algebraische Struktur einer Gruppe, die sog. **orthogonale Gruppe** (englisch: *orthogonal group*). Diese ist eine Untergruppe der **allgemeinen linearen Gruppe** (englisch: *general linear group*), die aus den invertierbaren  $n \times n$ -Matrizen besteht.

wobei  $\Lambda \in \mathbb{R}^{n \times n}$  die Diagonalmatrix der Eigenwerte ist und die orthogonale Matrix  $V \in \mathbb{R}^{n \times n}$  aus Eigenvektoren von  $A$  besteht.

Die Gleichung (4.2) können wir auch lesen als

$$A \begin{bmatrix} | & & | \\ v_1 & \cdots & v_n \\ | & & | \end{bmatrix} = \begin{bmatrix} | & & | \\ v_1 & \cdots & v_n \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix},$$

also spaltenweise:  $A v_i = \lambda_i v_i$ . (**Quizfrage:** Warum?)

## § 4.2 SINGULÄRWERTZERLEGUNG

Die Spektralzerlegung mit Hilfe der Eigenpaare ist nicht immer die geeignete Darstellung einer Matrix. Insbesondere für rechteckige Matrizen  $A \in \mathbb{R}^{m \times n}$  existieren überhaupt keine Eigenpaare. Daher führen wir jetzt die Singulärwertzerlegung ein, die immer existiert.

**Definition 4.3** (Singulärwertzerlegung). Unter einer **Singulärwertzerlegung** (englisch: **singular value decomposition**, kurz: **SVD**) einer Matrix  $A \in \mathbb{R}^{m \times n}$  versteht man eine Darstellung der Gestalt

$$A = U \Sigma V^T, \quad (4.3)$$

wobei

- $U \in \mathbb{R}^{m \times m}$  und  $V \in \mathbb{R}^{n \times n}$  orthogonale Matrizen sind und
- $\Sigma \in \mathbb{R}^{m \times n}$  eine (rechteckige) Diagonalmatrix ist, deren Diagonaleinträge  $\sigma_1, \dots, \sigma_{\min\{m,n\}} \geq 0$  sind.

Die Spaltenvektoren  $u_i$  von  $U$  heißen **Linkssingulärvektoren** (englisch: **left singular vectors**) von  $A$ . Die Spaltenvektoren  $v_i$  von  $V$  heißen **Rechtssingulärvektoren** (englisch: **right singular vectors**) von  $A$ . Die Zahlen  $\sigma_i$  heißen **Singulärwerte** (englisch: **singular values**) von  $A$ .

Im Falle von  $n > m$  hat die Singulärwertzerlegung beispielsweise die folgende Gestalt:

$$A = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_m \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} 0 & \cdots & 0 \\ \vdots & & \vdots \\ 0 & \cdots & 0 \end{bmatrix} \begin{bmatrix} \text{---} & v_1^T & \text{---} \\ & \vdots & \\ \text{---} & v_n^T & \text{---} \end{bmatrix} \quad (4.4a)$$

und im Falle von  $m > n$

$$A = \begin{bmatrix} | & & & | \\ u_1 & \cdots & \cdots & u_m \\ | & & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ \hline 0 & \cdots & 0 & \\ \vdots & & \vdots & \\ 0 & \cdots & 0 & \end{bmatrix} \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_n^\top & - \end{bmatrix}. \quad (4.4b)$$

Jede Matrix  $A \in \mathbb{R}^{m \times n}$  besitzt eine Singulärwertzerlegung. Um das zu begründen (**Quizfrage:** Können Sie die Details ausarbeiten?), kann man Spektralzerlegungen

$$A^\top A = V \Lambda V^\top \quad \text{und} \quad A A^\top = U \tilde{\Lambda} U^\top \quad (4.5)$$

der reellen, symmetrischen, positiv semi-definiten Matrizen  $A^\top A \in \mathbb{R}^{n \times n}$  und  $A A^\top \in \mathbb{R}^{m \times m}$  gemäß [Satz 4.2](#) betrachten. Die Singulärwerte sind gerade die Wurzeln der Diagonaleinträge der kleineren der beiden Diagonalmatrizen  $\Lambda \in \mathbb{R}^{n \times n}$  und  $\tilde{\Lambda} \in \mathbb{R}^{m \times m}$ .

**Bemerkung 4.4** (Eindeutigkeit und Anordnung der Singulärwerte). *Die Singulärwerte sind eindeutig und werden i. d. R. absteigend angeordnet:*

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_{\min\{m,n\}} \geq 0. \quad (4.6)$$

*Wir gehen ab jetzt immer von absteigend sortierten Singulärwerten aus.*

Die Singulärwertzerlegung enthält sehr viel Information über die Matrix  $A$  und die zu ihr gehörende lineare Abbildung  $x \mapsto Ax$ . Beispielsweise ist die Anzahl  $0 \leq r \leq \min\{m, n\}$  der echt positiven Singulärwerte gerade gleich  $\text{Rang}(A)$ , also gleich der Dimension des Bildraumes  $\text{Bild } A$ . Die «späten» Rechtssingulärvektoren  $v_{r+1}, \dots, v_n$  bilden eine orthonormale Basis von  $\ker A$ , und die «frühen» Rechtssingulärvektoren  $v_1, \dots, v_r$  bilden eine orthonormale Basis des orthogonalen Komplements  $(\ker A)^\perp$ . Weiterhin bilden die «frühen» Linkssingulärvektoren  $u_1, \dots, u_r$  eine Basis von  $\text{Bild } A$ , und die «späten» Linkssingulärvektoren  $u_{r+1}, \dots, u_m$  bilden eine Basis von  $(\text{Bild } A)^\perp$ . **Quizfrage:** Klar?

**Quizfrage:** Was vermuten Sie: Ist die Singulärwertzerlegung einer Matrix eindeutig?

**Bemerkung 4.5** (Dünne Singulärwertzerlegung). *Im Falle von  $n > m$  sehen wir aus der Darstellung (4.4a), dass die «späten» Zeilen von  $V^\top$  (Spalten von  $V$ ) mit den Indizes  $m+1, \dots, n$  zur Darstellung von  $A$  nicht benötigt werden. Es gilt also auch*

$$A = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_m \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \\ & & & \end{bmatrix} \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_m^\top & - \end{bmatrix}. \quad (4.7a)$$



Im Falle von  $m > n$  werden hingegen die «späten» Spalten von  $U$  mit den Indizes  $n + 1, \dots, m$  nicht benötigt. Es gilt also auch

$$A = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_n \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{bmatrix} \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_n^\top & - \end{bmatrix}. \quad (4.7b)$$

Man spricht dann von einer **dünnen Singulärwertzerlegung** (englisch: **thin SVD**, **economy-size SVD**).

**Beispiel 4.6** (Singulärwertzerlegung). Die Matrix

$$A = \begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 4 & 6 & 8 \\ 1 & 1 & 0 & 0 \end{bmatrix}$$

besitzt die Singulärwertzerlegung (aus Darstellungsgründen auf zwei Nachkommastellen gerundet)

$$A = \begin{bmatrix} -0.45 & 0.02 & -0.89 \\ -0.89 & 0.04 & 0.45 \\ -0.05 & -1.00 & 0.00 \\ \text{Bild } A & & (\text{Bild } A)^\perp \end{bmatrix} \begin{bmatrix} 12.26 & 0 & 0 & 0 \\ 0 & 1.30 & 0 & 0 \\ 0 & 0 & 0.00 & 0 \end{bmatrix} \begin{bmatrix} -0.19 & -0.37 & -0.55 & -0.73 \\ -0.69 & -0.61 & 0.23 & 0.31 \\ -0.70 & 0.70 & -0.08 & -0.11 \\ 0.00 & 0.00 & -0.80 & 0.60 \\ & & (\ker A)^\perp & \ker A \end{bmatrix}^\top.$$

Der Rang der Matrix ist also gleich 2. Mit Hilfe von `NUMPY` kann diese Singulärwertzerlegung wie folgt berechnet werden:

```
import numpy as np
A = np.array([[1, 2, 3, 4], [2, 4, 6, 8], [1, 1, 0, 0]])
U, Sigma, VT = np.linalg.svd(A)
```

Die Darstellung (4.3) erlaubt es uns, die Matrix  $A$  als eine Summe von  $r$  Rang-1-Matrizen darzustellen:

$$A = U \Sigma V^\top = \sum_{i=1}^r \sigma_i \underbrace{u_i v_i^\top}_{\text{Rang-1-Matrix}}. \quad (4.8)$$

Damit können wir auch Matrix-Vektor-Produkte genau analysieren:

$$Ax = U \Sigma V^\top x = \sum_{i=1}^r \sigma_i u_i (v_i^\top x). \quad (4.9)$$

**(Quizfrage:** Wie können wir diese Darstellung interpretieren? Was also «macht»  $v_i^\top x$ ? Was bedeutet die Multiplikation mit  $\sigma_i$ ? Und was die Multiplikation mit  $u_i$ ?) Insbesondere folgt aus (4.9) sofort

$$Av_i = \begin{cases} \sigma_i u_i & \text{für } i = 1, \dots, r, \\ 0 & \text{für } i = r + 1, \dots, n. \end{cases}$$

Mit Hilfe der Darstellung (4.9) können wir nun die Frage auflösen, warum die Matrixnorm  $\|A\|_{2 \rightarrow 2}$  durch (2.3b) gegeben ist. Die (quadrierte) Norm des Bildes  $Ax$  lässt sich nun unter Ausnutzung der Orthonormalität der Vektoren  $u_i$  wie folgt bestimmen:

$$\begin{aligned}
 \|Ax\|_2^2 &= \left\| \sum_{i=1}^r \sigma_i u_i (v_i^\top x) \right\|_2^2 = \left( \sum_{i=1}^r \sigma_i u_i (v_i^\top x), \sum_{j=1}^r u_j \sigma_j (v_j^\top x) \right) \quad (\text{wegen (4.9)}) \\
 &= \sum_{i=1}^r \sum_{j=1}^r \sigma_i \sigma_j u_i^\top u_j (v_i^\top x) (v_j^\top x) \\
 &= \sum_{i=1}^r \sigma_i^2 \|u_i\|_2^2 (v_i^\top x)^2 \quad (\text{wegen der Orthogonalität von } u_i \text{ und } u_j \text{ für } i \neq j) \\
 &= \sum_{i=1}^r \sigma_i^2 (v_i^\top x)^2 \quad (\text{wegen } \|u_i\|_2 = 1) \\
 &\leq \sigma_1^2 \sum_{i=1}^r (v_i^\top x)^2 \\
 &= \sigma_1^2 \|x\|_2^2.
 \end{aligned}$$

Die obige Abschätzung ist mit Gleichheit erfüllt, wenn  $x$  ein Vielfaches von  $v_1$  ist. In dem Fall gilt  $\|Ax\|_2 = \sigma_1 \|x\|_2$ . Da der größte Singulärwert  $\sigma_1$  gleichzeitig die Wurzel des größten Eigenwertes von  $A^\top A$  ist, haben wir damit (2.3b) gezeigt und außerdem:

**Satz 4.7** (Die  $\|\cdot\|_{2 \rightarrow 2}$ -Norm einer Matrix ist ihr größter Singulärwert).

(i) Für jede Matrix  $A \in \mathbb{R}^{m \times n}$  gilt:

$$\|A\|_{2 \rightarrow 2} = \sigma_1, \quad (4.10)$$

also ist die Spektralnorm gerade der größte Singulärwert der Matrix.

(ii) Für *symmetrische* Matrizen  $A \in \mathbb{R}^{n \times n}$  gilt gleichzeitig:

$$\|A\|_{2 \rightarrow 2} = \max\{|\lambda| \mid \lambda \text{ ist ein Eigenwert von } A\}, \quad (4.11)$$

also ist die Spektralnorm gerade der Betrag des betragsgrößten Eigenwerts der Matrix.

*Beweis.* Aussage (i) haben wir oben bereits gezeigt. Es gilt nach Konstruktion der Singulärwertzerlegung

$$\sigma_1^2 = \max\{|\lambda| \mid \lambda \text{ ist ein Eigenwert von } A^\top A\}.$$

Wenn  $A$  symmetrisch ist, dann existiert die Spektralzerlegung  $A = V\Lambda V^\top$ , siehe (4.2) mit einer orthogonalen Matrix  $V$ . Es gilt also  $A^\top A = V\Lambda V^\top V\Lambda V^\top = V\Lambda^2 V^\top$ . Die Eigenwerte von  $A^\top A$  sind demnach – wenn  $A$  symmetrisch ist – gerade die Quadrate der Eigenwerte von  $A$ . Wir haben also

$$\begin{aligned}
 \sigma_1^2 &= \max\{|\lambda| \mid \lambda \text{ ist ein Eigenwert von } A^\top A\} \\
 &= \max\{\lambda^2 \mid \lambda \text{ ist ein Eigenwert von } A\}.
 \end{aligned}$$

Durch Wurzelziehen folgt Aussage (ii). □

Im Beweis von Satz 3.16 haben wir die Norm  $\|y x^T\|_{2 \rightarrow 2}$  der Rang-1-Matrix  $y x^T$  benötigt. **Quizfrage:** Können Sie diese Norm mit dem jetzigen Wissen über die Singulärwertzerlegung bestimmen?

Falls  $A \in \mathbb{R}^{n \times n}$  invertierbar ist (**Quizfrage:** Wie drückt man das mit Hilfe der Singulärwerte von  $A$  aus?), dann erlaubt es uns die Singulärwertzerlegung  $A = U \Sigma V^T$ , auch direkt eine Singulärwertzerlegung von  $A^{-1}$  anzugeben:

$$A^{-1} = V \Sigma^{-1} U^T. \quad (4.12)$$

Hieraus sieht man sofort, dass  $1/\sigma_n$  der maximale Singulärwert von  $A^{-1}$  ist.

Weiter ergibt sich aus diesen Überlegungen, dass die Konditionszahl (3.22) einer regulären Matrix  $A \in \mathbb{R}^{n \times n}$  die Beziehung

$$\kappa(A) = \|A\|_{2 \rightarrow 2} \|A^{-1}\|_{2 \rightarrow 2} = \frac{\sigma_1}{\sigma_n} \quad (4.13)$$

erfüllt.

**Quizfrage:** Was ist die Konditionszahl einer orthogonalen Matrix? **Quizfrage:** Welche quadratischen (notwendigerweise invertierbaren) Matrizen haben Konditionszahl gleich eins? **Quizfrage:** Wie sieht die Singulärwertzerlegung einer Diagonalmatrix aus?

Auf Basis der SVD kann man sehr schön sehen, warum die Abschätzung

$$\|BA\|_{2 \rightarrow 2} \leq \|B\|_{2 \rightarrow 2} \|A\|_{2 \rightarrow 2}$$

aus Lemma 2.3 (v) (Submultiplikativität) insbesondere für die Spektralnorm i. A. nicht scharf ist. Wir betrachten dazu Singulärwertzerlegungen beider Matrizen:

$$B = U^{(2)} \Sigma^{(2)} V^{(2)T} \quad \text{und} \quad A = U^{(1)} \Sigma^{(1)} V^{(1)T}.$$

Es gilt also  $\|A\|_{2 \rightarrow 2} = \sigma_1^{(1)}$  und  $\|B\|_{2 \rightarrow 2} = \sigma_1^{(2)}$ . Es sei nun  $x$  ein Rechtssingulärvektor von  $A$  zum größten Singulärwert  $\sigma_1^{(1)}$  mit  $\|x\|_2 = 1$ . Dann gilt  $Ax = \sigma_1^{(1)} u_1^{(1)}$  mit  $\|Ax\|_2 = \sigma_1^{(1)}$ . Falls nun  $u^{(1)}$  ein Rechtssingulärvektor von  $B$  zum größten Singulärwert  $\sigma_1^{(2)}$  ist, so gilt weiter

$$BAx = \sigma_1^{(1)} B u^{(1)} = \sigma_1^{(1)} \sigma_1^{(2)} u^{(2)} \quad \text{mit} \quad \|BAx\|_2 = \sigma_1^{(1)} \sigma_1^{(2)}.$$

In diesem Fall haben wir also tatsächlich einen Vektor gefunden, für den  $\|BAx\|_2 = \|B\|_{2 \rightarrow 2} \|A\|_{2 \rightarrow 2} \|x\|_2$  gilt und damit  $\|BA\|_{2 \rightarrow 2} = \|B\|_{2 \rightarrow 2} \|A\|_{2 \rightarrow 2}$ . Im anderen Fall besitzt  $u^{(1)}$  auch Anteile von Rechtssingulärvektoren der Matrix  $B$ , die nicht zum größten Singulärwert von  $B$  gehören. Das impliziert  $\|B u^{(1)}\|_2 < \sigma_1^{(2)} \|u^{(1)}\|_2$ .

Mit diesen Überlegungen kann man folgenden Satz beweisen:

**Satz 4.8** (Submultiplikativität der  $2 \rightarrow 2$ -Norm). *Es seien  $A \in \mathbb{R}^{m \times n}$  und  $B \in \mathbb{R}^{\ell \times m}$ . Dann sind die folgenden Aussagen äquivalent:*

- (i) *Es gilt  $\|BA\|_{2 \rightarrow 2} = \|B\|_{2 \rightarrow 2} \|A\|_{2 \rightarrow 2}$ .*

- (ii) Es gibt einen Vektor  $x \in \mathbb{R}^n$ ,  $x \neq 0$ , mit der Eigenschaft  $\|Ax\|_2 = \sigma_1^{(1)} \|x\|_2$  und  $\|BAx\|_2 = \sigma_1^{(2)} \|Ax\|_2$ .
- (iii) Der Unterraum von  $\mathbb{R}^m$ , der durch die Linkssingulärvektoren von  $A$  aufgespannt wird, die zum größten Singulärwert  $\sigma_1^{(1)}$  gehören, schneidet den Unterraum von  $\mathbb{R}^m$ , der durch die Rechtssingulärvektoren von  $B$  aufgespannt wird, die zum größten Singulärwert  $\sigma_1^{(2)}$  gehören, in mehr als nur dem Nullvektor.

**Folgerung 4.9** (Submultiplikativität der  $2 \rightarrow 2$ -Norm). Insbesondere in folgenden Fällen sind die Aussagen (i) bis (iii) wahr:

- Die «mittlere» Dimension ist  $m = 1$ .
- *Mindestens* eine der Matrizen  $A$  oder  $B$  ist die Nullmatrix.
- *Mindestens* eine der Matrizen  $A$  oder  $B$  ist eine orthogonale Matrix.

**Quizfrage:** Beweis?

Eine weitere wichtige Folgerung daraus ist:

**Folgerung 4.10** (Multiplikation mit einer orthogonalen Matrix ändert die Konditionszahl nicht). Es sei  $A \in \mathbb{R}^{n \times n}$  beliebig und  $V \in \mathbb{R}^{n \times n}$  orthogonal. Dann gilt  $\|AV\|_{2 \rightarrow 2} = \|VA\|_{2 \rightarrow 2} = \|A\|_{2 \rightarrow 2}$ .

Beweis.

□

Schließlich geben wir noch eine Eigenschaft der Singulärwertzerlegung an.

**Satz 4.11** (Die SVD löst die Aufgabe der Best-Approximation mit konstantem Rang, vgl. Golub, van Loan, 1983, Theorem 2.5.3). Es sei  $A \in \mathbb{R}^{m \times n}$  eine Matrix mit Singulärwertzerlegung  $A = U\Sigma V^T$ . Weiter sei  $r = \text{Rang}(A)$  und  $0 \leq k \leq r$ . Dann löst die Matrix

$$A_k = \sum_{i=1}^k \sigma_i u_i v_i^T = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{bmatrix} \begin{bmatrix} - & v_1^T & - \\ & \vdots & \\ - & v_k^T & - \end{bmatrix} \in \mathbb{R}^{m \times n} \quad (4.14)$$

die folgende Optimierungsaufgabe:

$$\begin{aligned} &\text{Minimiere} \quad \|A - B\|_{2 \rightarrow 2}, \quad B \in \mathbb{R}^{m \times n} \\ &\text{unter der Bedingung} \quad \text{Rang}(B) = k. \end{aligned} \quad (4.15)$$

Der Optimalwert ist  $\|A - A_k\|_2 = \sigma_{k+1}$ .

*Beweis.* Wir überzeugen uns zunächst davon, dass die Matrix  $A_k$  den Rang  $k$  hat und dadurch als Kandidat in der Optimierungsaufgabe (4.15) zur Verfügung steht. Unter Ausnutzung der Orthogonalität von  $U$  und  $V$  ergibt sich aus (4.14)

$$U^T A_k V = \left[ \begin{array}{ccc|ccc} \sigma_1 & & & 0 & \cdots & 0 \\ & \ddots & & \vdots & & \vdots \\ & & \sigma_k & 0 & \cdots & 0 \\ \hline 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{array} \right] \in \mathbb{R}^{m \times n}.$$

Da die Multiplikation von  $A_k$  mit invertierbaren Matrizen den Rang nicht verändert und  $\sigma_1 \geq \cdots \geq \sigma_k > 0$  gilt, gilt tatsächlich  $\text{Rang}(A_k) = k$ . **Quizfrage:** Warum lässt die Multiplikation einer Matrix mit einer invertierbaren Matrix den Rang unverändert? **Quizfrage:** Warum ist  $\sigma_k > 0$ ?

Weiter gilt (Illustration für den Fall  $m > n$ )

$$U^T (A - A_k) V = \left[ \begin{array}{ccc|ccc} 0 & & & 0 & \cdots & 0 \\ & \ddots & & \vdots & & \vdots \\ & & 0 & 0 & \cdots & 0 \\ \hline 0 & \cdots & 0 & \sigma_{k+1} & & \\ \vdots & & \vdots & & \ddots & \\ 0 & \cdots & 0 & & & \sigma_p \\ \hline 0 & \cdots & 0 & 0 & \cdots & 0 \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & 0 & \cdots & 0 \end{array} \right] \in \mathbb{R}^{m \times n}$$

mit  $p = \min\{m, n\}$ . Aus Folgerung 4.10 folgt  $\|A - A_k\|_{2 \rightarrow 2} = \|U^T (A - A_k) V\|_{2 \rightarrow 2}$ , und letztere Matrix hat offenbar  $\sigma_{k+1}$  als größten Singulärwert, also gilt  $\|A - A_k\|_{2 \rightarrow 2} = \sigma_{k+1}$ .

Wir zeigen jetzt noch, dass für jede beliebige Matrix  $B \in \mathbb{R}^{m \times n}$  mit  $\text{Rang}(B) = k$  gilt:  $\|A - B\|_{2 \rightarrow 2} \geq \sigma_{k+1}$ . Damit ist dann die behauptete Optimalität von  $A_k$  gezeigt und der Beweis abgeschlossen. Die Dimensionsformel  $k = \dim \text{Bild } B = n - \dim \ker B$  zeigt, dass der Kern von  $B$  die Dimension  $n - k$  besitzt. Wir können also eine Menge orthonormaler Vektoren  $x_1, \dots, x_{n-k}$  finden, sodass

$$\ker B = \text{span}\{x_1, \dots, x_{n-k}\}$$

gilt. Aus Dimensionsgründen gilt

$$\text{span}\{x_1, \dots, x_{n-k}\} \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \{0\}.$$

Es sei nun  $z$  ein Vektor in diesem Durchschnitt mit  $\|z\|_2 = 1$ . Wegen  $Bz = 0$  und

$$Az = \sum_{i=1}^r \sigma_i u_i (v_i^T z) = \sum_{i=1}^{k+1} \sigma_i u_i (v_i^T z)$$

erhalten wir

$$\begin{aligned}
 \|A - B\|_{2 \rightarrow 2}^2 &\geq \|(A - B)z\|_{2 \rightarrow 2}^2 \quad \text{siehe Lemma 2.3} \\
 &= \|Az\|_{2 \rightarrow 2}^2 \quad \text{wegen } Bz = 0 \\
 &= \left\| \sum_{i=1}^{k+1} \sigma_i u_i (v_i^\top z) \right\|_2^2 \\
 &= \sum_{i=1}^{k+1} \sigma_i^2 (v_i^\top z)^2 \quad \text{wegen der Orthonormalität der Vektoren } u_i \\
 &\geq \sigma_{k+1}^2 \sum_{i=1}^{k+1} (v_i^\top z)^2 \quad \text{wegen } \sigma_1 \geq \dots \geq \sigma_{k+1} \\
 &= \sigma_{k+1}^2 \quad \text{wegen } 1 = \|z\|_2^2 = \left\| \sum_{i=1}^{k+1} (v_i^\top z) v_i \right\|_2^2 = \sum_{i=1}^{k+1} (v_i^\top z)^2.
 \end{aligned}$$

□

Dieses Resultat zeigt, dass die führenden Terme in einer Singulärwertzerlegung die beste Approximation der gegebenen Matrix im Sinne der  $2 \rightarrow 2$ -Norm unter allen Matrizen mit gegebenem **Rang** <sup>GM</sup> darstellen. Der dabei gemachte Approximationsfehler ist gerade der erste Singulärwert, der nicht berücksichtigt wird.

**Folgerung 4.12** (Abstand zu Matrizen mit Rangdefekt). *Es sei  $A \in \mathbb{R}^{m \times n}$  eine Matrix und  $p = \min\{m, n\}$ . Dann gilt*

$$\min\{\|A - B\|_{2 \rightarrow 2} \mid B \in \mathbb{R}^{m \times n}, \text{Rang}(B) < p\} = \sigma_p > 0. \quad (4.16)$$

*Beweis.*

□

Der kleinste Singulärwert gibt also den Abstand einer Matrix  $A$  zur Menge der Matrizen gleicher Dimensionen an, die nicht vollen Rang haben.

Ende der Woche 3

## Kapitel 2 Zahldarstellung und Rechnerarithmetik

### § 5 ZAHLDARSTELLUNG

Computer rechnen mit bestimmten Zahldarstellungen. Man unterscheidet zwischen **Ganzzahlen** (englisch: *integer numbers*), **Festkommazahlen** (englisch: *fixed-point numbers*) und **Fließkommazahlen** (englisch: *floating-point numbers*). Da sich die numerische Mathematik mit Aufgaben der «kontinuierlichen Mathematik» beschäftigt, sind Ganzzahlen i. d. R. nur als Zähler usw. interessant. Wir behandeln deshalb jetzt hier vorrangig Fließkommazahlen.

**Definition 5.1** (Fließkommazahl, vgl. Bornemann, 2018, Kapitel 12.1). Eine **Fließkommazahl** zur **Basis** (englisch: *base, radix*)  $\beta \in \{2, 3, \dots\} \subset \mathbb{N}$  ist eine Zahl der Form

$$m \beta^e. \quad (5.1)$$

Dabei heißt

$$m = \pm((m_0.m_{-1} \cdots m_{1-r}))_\beta = \pm(m_0 \beta^0 + m_{-1} \beta^{-1} + \cdots + m_{1-r} \beta^{1-r}) \quad (5.2)$$

mit Ziffern  $m_0, \dots, m_{1-r} \in \{0, 1, \dots, \beta - 1\}$  die (rationale) **Mantisse** (englisch: *significand, mantissa*) der Zahl und

$$e = \pm(e_{s-1}e_{s-2} \cdots e_0)_\beta = \pm(e_{s-1} \beta^{s-1} + e_{s-2} \beta^{s-2} + \cdots + e_1 \beta^1 + e_0 \beta^0) \quad (5.3)$$

mit Ziffern  $e_{s-1}, \dots, e_0 \in \{0, 1, \dots, \beta - 1\}$  ihr (ganzzahliger) **Exponent** (englisch: *exponent*). Es gilt die Konvention, dass die führende Stelle der Mantisse  $m_0 \neq 0$  ist.<sup>1</sup> Die einzige Ausnahme ist die Zahl 0, die durch die Mantisse  $m = 0$  und beliebigen Exponenten  $e$  dargestellt wird.

Wir bezeichnen die Menge der darstellbaren Fließkommazahlen (zu gegebener Basis, Mantissenlänge und Wertebereich für den Exponenten) mit  $\mathbb{F}$ . Man spricht manchmal auch von einem **Fließkommagitter** (englisch: *floating point grid*).

Die Notation  $(\cdot)_\beta$  in (5.2) und (5.3) bedeutet, dass die Zahl zur Basis  $\beta$  notiert ist. Im Fall der Basis  $\beta = 10$  lässt man den tiefgestellten Index  $(\cdot)_{10}$  häufig weg.

**Beachte:** Die Fließkommazahlen  $\mathbb{F}$  sind eine Teilmenge der rationalen Zahlen  $\mathbb{Q}$ . (**Quizfrage:** Klar?)

Mit Hilfe von Fließkommazahlen lassen sich Zahlen sehr unterschiedlicher Größenordnung darstellen, etwa (zur Basis  $\beta = 10$ )

<sup>1</sup>Zur Verdeutlichung kann man dann auch von einer **normalisierten Fließkommazahl** sprechen.

- die Ruhemasse eines Elektrons, etwa  $9.109\,383 \cdot 10^{-31}$  kg,
- die Lichtgeschwindigkeit im Vakuum, etwa  $2.997\,924\,58 \cdot 10^8$  m s<sup>-1</sup>,
- die Avogadro-Konstante, etwa  $6.022\,140\,76 \cdot 10^{23}$  mol.

Wie oben bereits bemerkt: Wird eine Zahl nicht zur Basis  $\beta = 10$  notiert, so deutet man die Basis häufig durch einen tiefgestellten Index an. Beispielsweise ist

$$1001_2 = 9 \quad \text{und} \quad 1001.11_2 = 9 + \frac{1}{2} + \frac{1}{4} = 9.75.$$

Die Anzahl  $r$  der Stellen der Mantisse sowie die Anzahl  $s$  der Stellen des Exponenten in einem System von Fließkommazahlen sind in der Regel fest. Die Menge der damit darstellbaren Zahlen ist somit endlich.

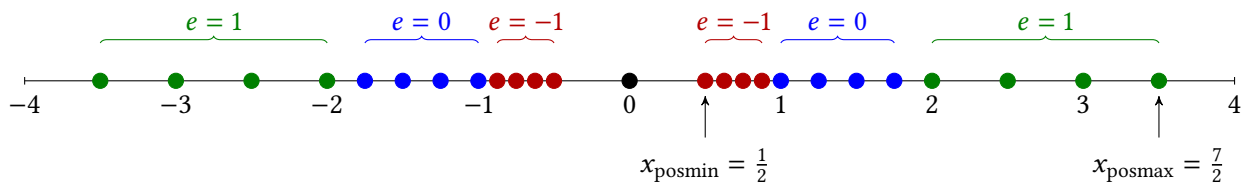
**Beispiel 5.2** (Darstellbare Fließkommazahlen). *Im Fall der Basis  $\beta = 2$ , Mantissenlänge  $r = 3$  und Exponentenlänge  $s = 1$  sind die folgenden positiven Mantissen möglich:*

$$\begin{aligned} m = 1.00_2 &= \frac{4}{4}, \\ m = 1.01_2 &= \frac{5}{4}, \\ m = 1.10_2 &= \frac{6}{4}, \\ m = 1.11_2 &= \frac{7}{4}. \end{aligned}$$

Durch Kombination mit den zur Verfügung stehenden Exponenten  $e \in \{-1, 0, 1\}$ , also den Faktoren  $2^{-1}$ ,  $2^0$  und  $2^1$ , sind also die positiven Zahlen

$$\left\{ \frac{4}{8}, \frac{5}{8}, \frac{6}{8}, \frac{7}{8}, \frac{4}{4}, \frac{5}{4}, \frac{6}{4}, \frac{7}{4}, \frac{4}{2}, \frac{5}{2}, \frac{6}{2}, \frac{7}{2} \right\}$$

darstellbar, dazu noch ihre negativen Gegenstücke sowie schließlich die Zahl 0. (**Quizfrage:** Ist die Darstellung einer Fließkommazahl immer eindeutig?)



Wie das **Beispiel 5.2** zeigt, sind Fließkommazahlen nicht gleichmäßig verteilt. Weiterhin gibt es eine kleinste und eine größte positive Zahl in  $\mathbb{F}$ . Die kleinste positive Zahl  $x_{\text{posmin}}$  wird dargestellt durch die Mantisse  $m = (1.0 \cdots 0)_\beta = 1$  und den Exponenten

$$e = -((\beta - 1)(\beta - 1) \cdots (\beta - 1))_\beta = -(\beta - 1) \sum_{k=0}^{s-1} \beta^k = -(\beta - 1) \frac{\beta^s - 1}{\beta - 1} = 1 - \beta^s.$$



Es ergibt sich also

$$x_{\text{posmin}} = \beta^{1-\beta^s} \quad (5.4)$$

Die größte positive Zahl  $x_{\text{posmax}}$  dagegen wird dargestellt durch die Mantisse

$$m = ((\beta-1)(\beta-1) \cdots (\beta-1))_\beta = (\beta-1) \sum_{k=1-r}^0 \beta^k = (\beta-1) \beta^{1-r} \sum_{k=0}^{r-1} \beta^k = (\beta-1) \beta^{1-r} \frac{\beta^r - 1}{\beta - 1} = \beta^{1-r} (\beta^r - 1)$$

und den Exponenten

$$e = ((\beta-1)(\beta-1) \cdots (\beta-1))_\beta = (\beta-1) \sum_{k=0}^{s-1} \beta^k = (\beta-1) \frac{\beta^s - 1}{\beta - 1} = \beta^s - 1.$$

Daraus ergibt sich

$$x_{\text{posmax}} = \beta^{1-r} (\beta^r - 1) \beta^{\beta^s - 1} = (\beta - \beta^{1-r}) \beta^{\beta^s - 1}. \quad (5.5)$$

Im [Beispiel 5.2](#) etwa mit  $\beta = 2$ ,  $r = 3$  und  $s = 1$  gilt

$$\begin{aligned} x_{\text{posmin}} &= \beta^{1-\beta^s} = 2^{1-2} = \frac{1}{2}, \\ x_{\text{posmax}} &= (\beta - \beta^{1-r}) \beta^{\beta^s - 1} = (2 - 2^{1-3}) 2^{2^1 - 1} = \frac{7}{4} \cdot 2^1 = \frac{7}{2}. \end{aligned}$$

Das stimmt mit unseren Beobachtungen in [Beispiel 5.2](#) überein.

**Definition 5.3** (Zulässiger Bereich). Die Menge

$$\mathbb{D} := [-x_{\text{posmax}}, -x_{\text{posmin}}] \cup \{0\} \cup [x_{\text{posmin}}, x_{\text{posmax}}] \subset \mathbb{R} \quad (5.6)$$

heißt der **zulässige Bereich** eines Fließkommagitters  $\mathbb{F}$ .

**Quizfrage:** Was ist die Bedeutung des zulässigen Bereichs?

**Definition 5.4** (IEEE<sup>2</sup>-Standard 754). Der 1985 eingeführte **IEEE-Standard 754** ist auch heute noch die Grundlage der Darstellung von Fließkommazahlen auf sehr vielen Computerarchitekturen. In diesem Standard spielt das Format **double precision** (auch: **binary64**) die wichtigste Rolle. Es verwendet die Basis  $\beta = 2$ , Mantissenlänge  $r = 53$  und ganzzahlige Exponenten  $e \in [-1022, 1023]$ . Daraus können wir ähnlich wie oben die kleinste und größte positive darstellbare (normalisierte) Fließkommazahl herleiten:

$$\begin{aligned} x_{\text{posmin}} &= 1 \cdot 2^{-1022} = 2^{-1022} \approx 2.2 \cdot 10^{-308}, \\ x_{\text{posmax}} &= 2^{1-53} (2^{53} - 1) \cdot 2^{1023} = (2 - 2^{-52}) \cdot 2^{1023} \approx 1.8 \cdot 10^{308}. \end{aligned}$$

An Stelle eines vorzeichenbehafteten Exponenten  $e$  wird dieser als  $e = c - 1023$  interpretiert, wobei die Zahl  $c \in [0, 2047]$  mit 11 Bits gespeichert wird. Da die Extremfälle  $c = 0$  und  $c = 2047$  besondere Bedeutung haben, ergeben sich die nutzbaren Werte  $e \in [-1022, 1023]$  für den Exponenten. Da die Mantisse normalisiert ist, ist die führende Ziffer für Fließkommazahlen  $x \neq 0$  immer  $m_0 = 1$  und muss nicht gespeichert werden. Zusammen mit dem Vorzeichen der Mantisse benötigt eine Zahl nach dem double precision-Standard daher  $1 + 52 + 11 = 64$  Bits oder 8 Bytes an Speicher.

Zusätzlich definiert der IEEE-Standard die erweiterten Fließkommazahlen  $+\infty$ ,  $-\infty$  und NaN («not a number»), die mit Hilfe von  $c = 2047$  kodiert werden.

<sup>2</sup>Institute of Electrical and Electronics Engineers

**Definition 5.5** (Rundung). Es sei  $\mathbb{F}$  eine Menge von Fließkommazahlen. Eine Funktion  $\text{rd}: \mathbb{R} \cup \{\pm\infty\} \rightarrow \mathbb{F} \cup \{\pm\infty\}$  heißt **Rundung**, wenn sie folgende Eigenschaften erfüllt:

(i) Sie ist monoton, d. h.,

$$\text{rd}(x) \leq \text{rd}(y) \quad \text{für } x, y \in \mathbb{R} \cup \{\pm\infty\} \text{ mit } x \leq y.$$

(ii) Sie ist idempotent, d. h.,

$$\text{rd}(x) = x \quad \text{für } x \in \mathbb{F} \cup \{\pm\infty\}.$$

**Quizfrage:** Wofür könnten die beiden Eigenschaften der Rundung nützlich sein?

Wir betrachten hier nur die gängigste Strategie des IEEE-Standards, die **Rundung zur nächstgelegenen Fließkommazahl** (englisch: *round to nearest*). Diese Funktion bildet Zahlen im Bereich  $[-x_{\text{posmax}}, x_{\text{posmax}}]$  auf die jeweils nächstgelegene Zahl in  $\mathbb{F}$  ab, die bis auf endlich viele Ausnahmewerte für  $x$  eindeutig ist. Die Ausnahmewerte sind gerade diejenigen Zahlen  $x$ , die genau in der Mitte zwischen zwei benachbarten Zahlen in  $\mathbb{F}$  liegen. Diese Grenzfälle werden bei der Unterstrategie *round to nearest, ties to even* so gerundet, dass das niederwertigste Bit der Mantisse = 0 wird. Weiterhin werden Zahlen  $|x| > x_{\text{posmax}} + C$  auf  $+\infty$  bzw.  $-\infty$  gerundet, wobei  $C = \frac{1}{2} 0.0 \dots 01_{\beta} \cdot \beta^E$  mit dem maximalen Exponenten  $E$  ist. **Quizfrage:** Was könnte der Grund für diese Regel sein? (Diese Frage kann erst nach Lemma 5.6 beantwortet werden.)

Wenn wir in Zukunft von Rundung sprechen, so ist immer diese Strategie gemeint.

**Lemma 5.6** (Relativer Fehler bei der Rundung). Der relative Fehler bei der Rundung in einem System von Fließkommazahlen  $\mathbb{F}$  mit Mantissenlänge  $r$  erfüllt die Abschätzung

$$\frac{|x - \text{rd}(x)|}{|x|} \leq \epsilon_{\text{mach}} := \frac{1}{2} \beta^{1-r} \quad (5.7)$$

für alle Zahlen  $x \neq 0$ , die im zulässigen Bereich  $\mathbb{D}$  liegen. Dabei ist  $\epsilon_{\text{mach}}$  die sogenannte **Maschinengenauigkeit** (englisch: *machine precision, unit roundoff*).

**Beweis.** Wir müssen nur den Fall  $x > 0$  betrachten, da der Beweis im Fall  $x < 0$  analog läuft. Aufgrund der Annahme an  $x$  existieren zwei benachbarte Zahlen  $\underline{x}$  und  $\bar{x}$  in  $\mathbb{F}$  mit der Eigenschaft  $0 < \underline{x} \leq x < \bar{x}$ . Nach Definition der Rundung gilt

$$|x - \text{rd}(x)| \leq \frac{1}{2} (\bar{x} - \underline{x}).$$

Aus der Darstellung (4.2)–(5.3) sieht man, dass sich  $\underline{x}$  und  $\bar{x}$  um die Zahl

$$(0.00 \dots 01)_{\beta} \cdot \beta^e = \beta^{1-r} \cdot \beta^e$$

unterscheiden, wobei  $e$  der Exponent von  $\underline{x}$  ist. (**Quizfrage:** Warum stimmt das auch dann, wenn  $\underline{x} = ((\beta - 1) \cdot (\beta - 1) \cdot (\beta - 1) \dots (\beta - 1))_{\beta} \cdot \beta^e$  ist?)

Es gilt also

$$|x - \text{rd}(x)| \leq \frac{1}{2} \beta^{1-r} \cdot \beta^e$$

und

$$x \geq \underline{x} \geq (1.00 \cdots 00)_\beta \cdot \beta^e,$$

also

$$\frac{|x - \text{rd}(x)|}{|x|} \leq \frac{1}{2} \beta^{1-r} \cdot \beta^e \cdot \frac{1}{\beta^e} = \frac{1}{2} \beta^{1-r}$$

wie behauptet. □

**Quizfrage:** Gilt die Abschätzung (5.7) auch für Zahlen  $x$  mit  $|x| < x_{\text{posmin}}$ ? Falls nicht, wie groß kann der relative Fehler bei der Rundung für diese Zahlen höchstens werden?

**Folgerung 5.7** (zur Maschinengenauigkeit). Für alle  $x \in \mathbb{D}$  (einschließlich  $x = 0$ ) gilt

$$\text{rd}(x) = x(1 + \varepsilon) \tag{5.8}$$

mit einer Zahl  $\varepsilon \in \mathbb{R}$ , die  $|\varepsilon| \leq \varepsilon_{\text{mach}}$  erfüllt.

**Beispiel 5.8** (Maschinengenauigkeit).

(i) Für das Fließkommagitter aus [Beispiel 5.2](#) gilt

$$\varepsilon_{\text{mach}} = \frac{1}{2} 2^{1-3} = 2^{-3} = \frac{1}{8}.$$

(ii) Im Falle des double precision-Standards haben wir

$$\varepsilon_{\text{mach}} = \frac{1}{2} 2^{1-53} = 2^{-53} \approx 1.11 \cdot 10^{-16}.$$

(iii) Im Zehnersystem ( $\beta = 10$ ) mit zwei Nachkommastellen, also für Zahlen der Bauart 1.23 mit Mantissenlänge  $r = 3$ , ergibt sich die «Maschinen»genauigkeit  $\varepsilon_{\text{mach}} = \frac{1}{2} 10^{1-3} = 0.005$ .

**Bemerkung 5.9** (Maschinengenauigkeit).

(a) Die Maschinengenauigkeit ist gerade die Hälfte des Stellenwertes der letzten Ziffer der Mantisse; vgl. (5.2). Auf die darstellbaren Exponenten kommt es dabei nicht an.

(b) Die doppelte Fließkommagenauigkeit  $2\varepsilon_{\text{mach}}$  ist gerade der Abstand zwischen der Fließkommazahl 1 und der nächst größeren Fließkommazahl.

**Bemerkung 5.10** (Alternative Definition der Maschinengenauigkeit). Manche Autoren definieren die Maschinengenauigkeit als das Doppelte von unserer Definition. Das liegt darin begründet, dass (5.7) dann auch für andere Rundungsarten als round to nearest gilt, etwa round toward 0. Siehe auch die Diskussion auf [https://en.wikipedia.org/wiki/Machine\\_epsilon](https://en.wikipedia.org/wiki/Machine_epsilon).

## § 6 RECHNERARITHMETIK

Die Grundoperationen  $+$ ,  $-$ ,  $\cdot$  und  $/$  liefern, selbst für Argumente in  $\mathbb{F}$ , im Allgemeinen Ergebnisse, die keine Fließkommazahlen in  $\mathbb{F}$  sind. (**Quizfrage:** Können Sie Beispiele dafür angeben, wann das der Fall ist, etwa für das Fließkommasystem  $\mathbb{F}$  aus [Beispiel 5.2](#)?) Die Grundoperationen lassen sich also nicht fehlerfrei in  $\mathbb{F}$  ausführen. Stattdessen werden Ersatzoperationen, genannt **Maschinenoperationen**  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$  implementiert (in Hardware), die Ergebnisse in  $\mathbb{F}$  liefern. Der IEEE-Standard sieht vor, dass diese Grundoperationen mit korrekter Rundung ausgeführt werden. Das heißt konkret, dass folgende Beziehungen gelten:

$$x \oplus y = \text{rd}(x + y), \quad (6.1a)$$

$$x \ominus y = \text{rd}(x - y), \quad (6.1b)$$

$$x \odot y = \text{rd}(x \cdot y), \quad (6.1c)$$

$$x \oslash y = \text{rd}(x/y). \quad (6.1d)$$

Dabei sind jeweils  $x, y \in \mathbb{F}$ , und das Ergebnis liegt in  $\mathbb{F} \cup \{\pm\infty\}$  bzw. ergibt NaN im Falle von  $0 \oslash 0$ . Aus [Lemma 5.6](#) bzw. [Folgerung 5.7](#) ergibt sich direkt folgendes Resultat:

**Satz 6.1** (Relativer Fehler der Maschinenoperationen  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$ ). *Es seien  $x, y \in \mathbb{F}$ .*

(i) *Falls  $x + y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \oplus y = (x + y)(1 + \varepsilon).$$

(ii) *Falls  $x - y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \ominus y = (x - y)(1 + \varepsilon).$$

(iii) *Falls  $x \cdot y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \odot y = (x \cdot y)(1 + \varepsilon).$$

(iv) *Falls  $y \neq 0$  ist und  $x/y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \oslash y = (x/y)(1 + \varepsilon).$$

**Beachte:** Die Voraussetzung  $x + y \in \mathbb{D}$  besagt, dass das korrekte Resultat im Bereich der von den zur Verfügung stehenden Fließkommazahlen überdeckten Intervalle [\(5.6\)](#) liegt.

Wie passt die Aussage von [Satz 6.1](#), dass  $x \ominus y$  nur um einen kleinen relativen Faktor vom echten Ergebnis  $x - y$  abweicht, mit der in [Abschnitt 3](#) gemachten Beobachtung zusammen, dass die Aufgabe,  $x - y$  zu berechnen, eine beliebig große relative Konditionszahl haben kann?

Das sind zwei verschiedene Paar Schuhe. Die Aussage von [Satz 6.1](#) bezieht sich auf die Berechnung von  $x \ominus y$ , wobei  $x$  und  $y$  bereits Fließkommazahlen sind und keiner Rundung unterliegen. (Eine Rundung erfolgt erst **nach** der Berechnung von  $x - y$ .) Die Aussage zur Kondition bezieht sich darauf, dass das Ergebnis von  $x - y$  stark abweichen kann, wenn die Daten  $x$  und  $y$  nicht exakt sind, etwa durch Rundung **vor** der Berechnung, um  $x$  und  $y$  in Fließkommazahlen zu verwandeln. Allerdings bestehen Algorithmen in der Regel aus einer Kette von Operationen, bei denen zwischendurch immer wieder Rundungen auf Fließkommazahlen stattfinden. Dadurch werden «früh» im Verfahren gemachte Rundungsfehler in allen späteren Schritten sichtbar, wobei die Konditionszahlen als (Verstärkungs-) Faktoren eingehen. Das untersuchen wir genauer in [Kapitel 3](#).

**Bemerkung 6.2** (Assoziativ- und Distributivgesetz).

- (a) Das Assoziativgesetz und das Distributivgesetz gelten für die Maschinenoperationen nicht mehr! Es ist also i. A.

$$\begin{aligned}(x \oplus y) \oplus z &\neq x \oplus (y \oplus z), \\ (x \oplus y) \odot z &\neq (x \odot z) \oplus (y \odot z)\end{aligned}$$

für  $x, y, z \in \mathbb{F}$ .

- (b) Das Kommutativgesetz der Addition und der Multiplikation gelten für  $\oplus$  und  $\odot$  aber weiter, d. h.,

$$\begin{aligned}x \oplus y &= y \oplus x, \\ x \odot y &= y \odot x\end{aligned}$$

für  $x, y \in \mathbb{F}$ .

- (c) Diese Beobachtung führt dazu, dass man Algorithmen, bevor man den Einfluss von Rundungsfehlern untersucht, zweifelsfrei notieren muss, weil bereits die Reihenfolge einen Einfluss auf das Ergebnis haben kann. Dabei helfen Klammersetzung und Einführung von Zwischenresultaten. Beispielsweise sind zur Auswertung von  $y = x_1^2 - x_2^2 = (x_1 + x_2)(x_1 - x_2)$  folgende Realisierungen in Algorithmen denkbar:

$u_1 := x_1 \odot x_1$	$u_1 := x_1 \oplus x_2$	$u_1 := x_1 \ominus x_2$
$u_2 := x_2 \odot x_2$	$u_2 := x_1 \ominus x_2$	$u_2 := x_1 \oplus x_2$
$y := u_1 \ominus u_2$	$y := u_1 \odot u_2$	$y := u_1 \odot u_2$

Mit den Operationen  $+$ ,  $-$ ,  $\cdot$  in  $\mathbb{R}$  führen diese Algorithmen zu identischen Ergebnissen, in Maschinearithmetik mit  $\oplus$ ,  $\ominus$ ,  $\odot$  sind die zweite und dritte Variante identisch (**Quizfrage:** Warum?), aber verschieden von der ersten.

Ende der Woche 4

# Kapitel 3 Stabilität von Algorithmen

## § 7 STABILITÄTSANALYSE

Bei der Stabilitätsanalyse von Algorithmen geht es darum, zu beurteilen, welchen Einfluss die in einer konkreten Implementierung mit Maschinenoperationen unvermeidbar auftretenden Rundungsfehler auf das Endergebnis haben. Diese Beurteilung kann man auf unterschiedliche Arten vornehmen, was zu verschiedenen Stabilitätsbegriffen führt. Wir beginnen mit dem Konzept im Sinne der Vorwärtsanalyse.

### § 7.1 VORWÄRTSANALYSE

Wir bezeichnen mit

$F(x)$  die auszuwertende Funktion,

$\widehat{F}(x)$  die durch den implementierten Algorithmus realisierte Funktion.

Wir beginnen mit einem einfachen Beispiel.

**Beispiel 7.1** (Fehlerfortpflanzung bei  $(a+b)+c$ ). Wir wollen  $y = F(x) = a+b+c$  für  $x = (a, b, c)^T \in \mathbb{F}^3$  auswerten. Wir gehen hier zunächst der Einfachheit halber davon aus, dass die Summanden bereits als Fließkommazahlen vorliegen. Wir betrachten den Algorithmus  $\widehat{F}(x) = ((a \oplus b) \oplus c)$ . Um Fehlerabschätzungen zu erhalten, gehen wir davon aus, dass die Zwischenergebnisse (vor der Rundung)  $a+b$  sowie  $(a \oplus b) + c$  im zulässigen Bereich  $\mathbb{D}$  liegen. Dann erhalten wir aus [Satz 6.1](#) die Aussage, dass Zahlen  $\varepsilon^{(1)}, \varepsilon^{(2)} \in \mathbb{R}$  mit  $|\varepsilon_i| \leq \varepsilon_{\text{mach}}$  existieren, sodass gilt:

$$\begin{aligned} a \oplus b &= (a+b)(1+\varepsilon^{(1)}) \\ (a \oplus b) \oplus c &= ((a \oplus b) + c)(1+\varepsilon^{(2)}) \\ &= ((a+b)(1+\varepsilon^{(1)}) + c)(1+\varepsilon^{(2)}) \\ &= (a+b+c) \left[ 1 + \frac{a+b}{a+b+c} \varepsilon^{(1)} (1+\varepsilon^{(2)}) + \varepsilon^{(2)} \right]. \end{aligned}$$

Für den relativen Fehler gilt daher

$$\frac{\widehat{F}(x) - F(x)}{F(x)} = \frac{(a \oplus b) \oplus c - (a+b+c)}{a+b+c} = \frac{a+b}{a+b+c} \varepsilon^{(1)} (1+\varepsilon^{(2)}) + \varepsilon^{(2)}.$$

Unter Vernachlässigung des gegenüber  $\varepsilon^{(i)}$  typischerweise kleinen Terms  $\varepsilon^{(1)}\varepsilon^{(2)}$  vereinfacht sich dieser Ausdruck zu

$$\frac{(a \oplus b) \oplus c - (a + b + c)}{a + b + c} \doteq \frac{a + b}{a + b + c} \varepsilon^{(1)} + \varepsilon^{(2)}. \quad (7.1)$$

Da wir die Fehlerbeiträge als voneinander unabhängig betrachten wollen, verwenden wir die Dreiecksungleichung und erhalten schließlich

$$\frac{|(a \oplus b) \oplus c - (a + b + c)|}{|a + b + c|} \leq \frac{|a + b|}{|a + b + c|} |\varepsilon^{(1)}| + |\varepsilon^{(2)}|. \quad (7.2)$$

Wir sehen also hier, dass der erste (früher) auftretende relative Rundungsfehler  $\varepsilon^{(1)}$  mit dem Faktor  $\frac{|a+b|}{|a+b+c|}$  fortgepflanzt wird. Dieser Faktor ist der Betrag der relativen partiellen Konditionszahl  $k_{11}$  der Additionsaufgabe  $(a + b) + c$ , vgl. Lemma 3.9. (**Quizfrage:** Warum ist hier nur die eine relative partielle Konditionszahl  $k_{11}$  relevant, während die andere  $k_{12}$  nicht vorkommt?) Der zweite (später) auftretende Rundungsfehler  $\varepsilon^{(2)}$  taucht dagegen ohne einen solchen Faktor in der Abschätzung auf. **Quizfrage:** Warum ist das so?

**Quizfrage:** Wie sieht im Vergleich zu (7.1) die Analyse für den alternativen Algorithmus  $\widehat{F}(x) = (a \oplus (b \oplus c))$  aus?

Das Beispiel 7.1 zeigt bereits wesentliche Bestandteile der Stabilitätsanalyse eines Algorithmus im Vorwärtssinne. In jedem Schritt kommen neue Rundungsfehler hinzu, und bereits vorhandene Fehler werden – mit der entsprechenden relativen Konditionszahl multipliziert – weitergegeben. Wir wollen dies nun allgemein analysieren. Dazu betrachten wir eine auszuwertende Funktion  $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$  und einen Algorithmus, der aus der Eingabe  $x =: x^{(0)} \in \mathbb{R}^{n^{(0)}}$  ein erstes Zwischenergebnis  $x^{(1)} := F^{(0)}(x^{(0)}) \in \mathbb{R}^{n^{(1)}}$  erzeugt, anschließend ein weiteres  $x^{(2)} := F^{(1)}(x^{(1)}) \in \mathbb{R}^{n^{(2)}}$  usw., bis schließlich das Endergebnis

$$y := x^{(N+1)} := F^{(N)}(x^{(N)}) \in \mathbb{R}^{n^{(N+1)}}$$

vorliegt. Für die erste und letzte Dimension gilt  $n^{(0)} = n$  und  $n^{(N+1)} = m$ . Die durch den Algorithmus realisierte Funktion wird also in Teilabbildungen zerlegt:

$$F = F^{(N)} \circ \dots \circ F^{(0)}.$$

Dabei bestehen die Komponenten der vektorwertigen Teilabbildungen  $F^{(k)}$  jeweils aus elementaren Abbildungen wie  $+$ ,  $\cdot$ ,  $\cos$  etc..

Demgegenüber steht in der tatsächlichen Implementierung des Verfahrens die Realisierung der Teilabbildungen durch Maschinenoperationen  $\widehat{F}^{(k)}$ . Dabei nehmen wir an, dass (komponentenweise)

$$\widehat{F}^{(k)} := \text{rd}(F^{(k)}) \quad (7.3)$$

gilt, d. h., die implementierten Operationen sind bis auf die nachfolgende Rundung exakt.

**Beispiel 7.2.** Für den «Algorithmus  $(a + b) + c$ » aus [Beispiel 7.1](#) gilt

$$x^{(0)} = \begin{pmatrix} a \\ b \\ c \end{pmatrix} \xrightarrow{F^{(0)}} \begin{pmatrix} a+b \\ c \end{pmatrix} \xrightarrow{F^{(1)}} (a+b) + c = x^{(2)} =: y \quad \text{und}$$

$$\widehat{x}^{(0)} = \begin{pmatrix} \widehat{a} \\ \widehat{b} \\ \widehat{c} \end{pmatrix} \xrightarrow{\widehat{F}^{(0)}} \begin{pmatrix} \widehat{a} \oplus \widehat{b} \\ \widehat{c} \end{pmatrix} \xrightarrow{\widehat{F}^{(1)}} (\widehat{a} \oplus \widehat{b}) \oplus \widehat{c} = \widehat{x}^{(2)} =: \widehat{y}.$$

Wir untersuchen jetzt den Einfluss der Rundungsfehler, also den Unterschied zwischen

$$y = (F^{(N)} \circ \dots \circ F^{(0)})(x^{(0)}) \quad \text{und}$$

$$\widehat{y} = (\widehat{F}^{(N)} \circ \dots \circ \widehat{F}^{(0)})(\widehat{x}^{(0)})$$

Wir lassen im Gegensatz zu unserem anfänglichen [Beispiel 7.1](#) nun auch (Rundungs-)fehler in der Eingabe zu, die zu einem anfänglichen Fehler

$$\Delta x^{(0)} := \widehat{x}^{(0)} - x^{(0)}$$

führen. Wegen [Lemma 5.6](#) gilt dabei

$$\Delta x^{(0)} = \begin{bmatrix} \varepsilon_1^{(0)} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \varepsilon_{n^{(0)}}^{(0)} \end{bmatrix} x^{(0)}$$

mit Zahlen  $\varepsilon_1^{(0)}, \dots, \varepsilon_{n^{(0)}}^{(0)}$ , die  $|\varepsilon_j^{(0)}| \leq \varepsilon_{\text{mach}}$  erfüllen. Also gilt insbesondere  $\|\Delta x^{(0)}\|_\infty \leq \varepsilon_{\text{mach}} \|x^{(0)}\|_\infty$ .

Im nächsten Schritt erhalten wir die Zwischenergebnisse

$$x^{(1)} = F^{(0)}(x^{(0)}) \quad \text{und}$$

$$\widehat{x}^{(1)} = \widehat{F}^{(0)}(\widehat{x}^{(0)}).$$

Für den Fehler

$$\Delta x^{(1)} := \widehat{x}^{(1)} - x^{(1)}$$

gilt

$$\begin{aligned} \Delta x^{(1)} &= \widehat{F}^{(0)}(\widehat{x}^{(0)}) - F^{(0)}(x^{(0)}) \\ &= [\widehat{F}^{(0)}(\widehat{x}^{(0)}) - F^{(0)}(\widehat{x}^{(0)})] + [F^{(0)}(\widehat{x}^{(0)}) - F^{(0)}(x^{(0)})] \\ &= \underbrace{[\text{rd}(F^{(0)}(\widehat{x}^{(0)})) - F^{(0)}(\widehat{x}^{(0)})]}_{\text{neu hinzukommender Rundungsfehler}} + \underbrace{[F^{(0)}(\widehat{x}^{(0)}) - F^{(0)}(x^{(0)})]}_{\text{Fortpflanzung des bereits vorhandenen Fehlers}}. \end{aligned} \tag{7.4}$$

Für den Term in der ersten Klammer gilt wegen [Lemma 5.6](#):

$$\text{rd}(F^{(0)}(\widehat{x}^{(0)})) - F^{(0)}(\widehat{x}^{(0)}) = \underbrace{\begin{bmatrix} \varepsilon_1^{(1)} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \varepsilon_{n^{(1)}}^{(1)} \end{bmatrix}}_{=: E^{(1)}} F^{(0)}(\widehat{x}^{(0)}) = E^{(1)} F^{(0)}(\widehat{x}^{(0)}) \tag{7.5}$$



mit Zahlen  $\varepsilon_1^{(1)}, \dots, \varepsilon_n^{(1)}$ , die  $|\varepsilon_j^{(1)}| \leq \varepsilon_{\text{mach}}$  erfüllen. Nach dem [Satz von Taylor 2.10](#) haben wir außerdem

$$F^{(0)}(\widehat{x}^{(0)}) \doteq F^{(0)}(x^{(0)}) \quad (7.6)$$

bis auf einen Fehlerterm, in den die Größe  $\|\Delta x^{(0)}\|_\infty \leq \varepsilon_{\text{mach}} \|x^{(0)}\|_\infty$  eingeht sowie eine Schranke für die Ableitung von  $F$  in einer Umgebung von  $x^{(0)}$ . Wenn wir jetzt (7.6) in die rechte Seite von (7.5) einsetzen, so erhalten wir also schließlich

$$\text{rd}(F^{(0)}(\widehat{x}^{(0)})) - F^{(0)}(\widehat{x}^{(0)}) \doteq E^{(1)} F^{(0)}(x^{(0)}) = E^{(1)} x^{(1)} \quad (7.7)$$

bis auf Terme, in die die komponentenweisen relativen Rundungsfehler  $\varepsilon_j^{(1)}$  quadratisch eingehen.

Der zweite Term in (7.4) lässt sich nach dem [Satz von Taylor 2.10](#) approximieren durch<sup>1</sup>

$$F^{(0)}(\widehat{x}^{(0)}) - F^{(0)}(x^{(0)}) \doteq DF^{(0)}(x^{(0)}) (\widehat{x}^{(0)} - x^{(0)}) = DF^{(0)}(x^{(0)}) \Delta x^{(0)},$$

wobei die Restterme wiederum quadratisch in den relativen Rundungsfehlern  $\varepsilon_j^{(1)}$  sind.

Zusammenfassend können wir also (7.4) in einer Approximation erster Ordnung schreiben als

$$\begin{array}{c} \text{neu hinzukommender Rundungsfehler} \\ \Delta x^{(1)} \doteq \overbrace{E^{(1)} x^{(1)} + DF^{(0)}(x^{(0)}) \Delta x^{(0)}}^{\text{Fortpflanzung des bereits vorhandenen Fehlers}}. \end{array} \quad (7.8)$$

Diese Beziehung gilt nicht nur für den ersten Schritt  $\cdot^{(0)} \rightsquigarrow \cdot^{(1)}$  im Algorithmus, sondern auch für alle weiteren Schritte  $\cdot^{(k)} \rightsquigarrow \cdot^{(k+1)}$ , also

$$\Delta x^{(k+1)} \doteq E^{(k+1)} x^{(k+1)} + DF^{(k)}(x^{(k)}) \Delta x^{(k)} \quad (7.9)$$

für  $k = 0, \dots, N$ . Durch Einsetzen erhalten wir schließlich für den absoluten Fehler  $\Delta y = \Delta x^{(N+1)}$  in der Ausgabe die Beziehung

$$\begin{aligned} \Delta y &\doteq D(F^{(N)} \circ \dots \circ F^{(0)})(x^{(0)}) \Delta x \\ &\quad + D(F^{(N)} \circ \dots \circ F^{(1)})(x^{(1)}) E^{(1)} x^{(1)} \\ &\quad + \dots \\ &\quad + DF^{(N)}(x^{(N)}) E^{(N)} x^{(N)} \\ &\quad + E^{(N+1)} x^{(N+1)}, \end{aligned} \quad (7.10)$$

wobei wie eingangs gesagt  $\Delta x := \Delta x^{(0)}$  den Fehler in der Eingabe bezeichnet. (**Quizfrage:** Können Sie die einzelnen Zeilen in (7.10) interpretieren? Wo entstehen sie jeweils?) Wir können an der Gleichung (7.10) sehen, dass die Norm der Jacobimatrix der «Restabbildung»  $F^{(N)} \circ \dots \circ F^{(k)}$  eine obere Schranke für den Einfluss des im Schritt  $\cdot^{(k-1)} \rightsquigarrow \cdot^{(k)}$  neu hinzukommenden Rundungsfehlers  $E^{(k)} x^{(k)}$  auf den Fehler  $\Delta y$  im Endergebnis angibt.

<sup>1</sup>Wir schreiben hier jetzt  $DF^{(0)}$  für die Ableitung (Jacobimatrix) von  $F^{(0)}$ , um Ausdrücke wie  $(F^{(0)})'$  zu vermeiden.

**Beispiel 7.3** (Fehlerfortpflanzung bei  $(a + b) + c$ ). Wir führen die Rechnung aus [Beispiel 7.1](#) für den «Algorithmus  $(a + b) + c$ » mit der soeben entwickelten Technik nochmal durch. Wir definieren das Verfahren mit Hilfe der Zwischenergebnisse

$$x = x^{(0)} = \begin{pmatrix} a \\ b \\ c \end{pmatrix}, \quad x^{(1)} = F^{(0)}(x^{(0)}) = \begin{pmatrix} a + b \\ c \end{pmatrix} \quad \text{und} \quad y = x^{(2)} = F^{(1)}(x^{(1)}) = a + b + c.$$

Es gilt also  $N = 1$  (Anzahl der Zwischenschritte), und wir benötigen die Teilabbildungen

$$\begin{aligned} F^{(0)}(x) = \begin{pmatrix} x_1 + x_2 \\ x_3 \end{pmatrix} &\Rightarrow DF^{(0)}(x) = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \\ F^{(1)}(x) = x_1 + x_2 &\Rightarrow DF^{(1)}(x) = [1 \quad 1]. \end{aligned}$$

Weiter benötigen wir die Jacobimatrix der Abbildung  $F^{(1)} \circ F^{(0)}$ , also

$$D(F^{(1)} \circ F^{(0)})(x) = \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}.$$

Die Matrizen für die auf jeder Ebene entstehenden relativen Fehler sind  $E^{(1)} = \begin{bmatrix} \varepsilon_1^{(1)} & 0 \\ 0 & \varepsilon_2^{(1)} \end{bmatrix}$  und  $E^{(2)} = \varepsilon_1^{(2)}$ .

Für  $\varepsilon_2^{(1)}$  dürfen wir sogar Null annehmen. (**Quizfrage:** Warum?) Somit ergibt sich für den absoluten Fehler im Ergebnis also wie in (7.10):

$$\begin{aligned} \Delta y &\doteq D(F^{(1)} \circ F^{(0)})(x^{(0)}) \Delta x + DF^{(1)}(x^{(1)}) E^{(1)} x^{(1)} + E^{(2)} x^{(2)} \\ &= \begin{bmatrix} 1 & 1 & 1 \end{bmatrix} \Delta x + \begin{bmatrix} 1 & 1 \end{bmatrix} \begin{bmatrix} \varepsilon_1^{(1)} & 0 \\ 0 & 0 \end{bmatrix} \begin{pmatrix} a + b \\ c \end{pmatrix} + \varepsilon_1^{(2)} (a + b + c) \\ &= \Delta a + \Delta b + \Delta c + \varepsilon_1^{(1)} (a + b) + \varepsilon_1^{(2)} (a + b + c). \end{aligned}$$

In relativer Betrachtungsweise haben wir daher

$$\frac{\Delta y}{y} \doteq \frac{\Delta a + \Delta b + \Delta c}{a + b + c} + \frac{a + b}{a + b + c} \varepsilon_1^{(1)} + \varepsilon_1^{(2)}. \quad (7.11)$$

**Quizfrage:** Wodurch kommen in (7.11) die zusätzlichen Terme im Vergleich zu unserer früheren Rechnung (7.1) zustande?

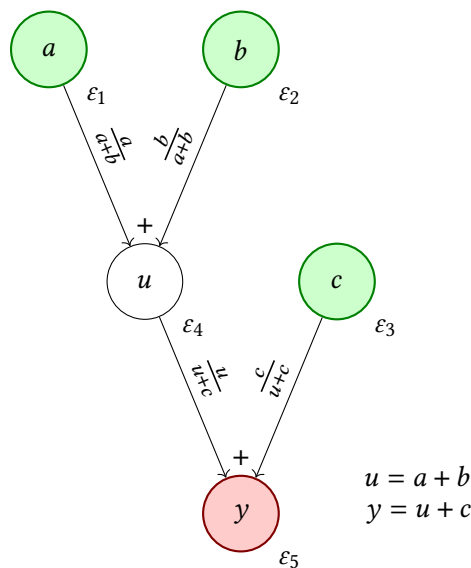
Die oben besprochene Technik zur Analyse des Einflusses von Rundungsfehlern (7.10) heißt (**differentielle**) **Vorwärtsanalyse**, weil man die Eingabefehler und zwischendurch auftretende Rundungsfehler vorwärts durch das Verfahren transportiert und weil die Faktoren, die dabei auftreten, aus der Ableitung gewonnen werden.

**Bemerkung 7.4** (zur differentiellen Vorwärtsanalyse).

- (1) Die Vorwärtsanalyse setzt natürlich voraus, dass alle Teilabbildungen  $F^{(k)}$  und deren Realisierungen durch Maschinenoperationen  $\widehat{F}^{(k)}$  überall definiert sind, wo sie benötigt werden.

- (2) Wie man bereits an einfachen Beispielen wie [Beispiel 7.3](#) sieht, führen verschiedene Realisierungen mathematisch äquivalenter Ausdrücke wie  $(a + b) + c$  und  $a + (b + c)$  zu unterschiedlichen Abschätzungen für die Fehlerfortpflanzung!
- (3) Die Faktoren in den Abschätzungen der Fehlerfortpflanzung wie beispielsweise (7.11) hängen von der Auswertungsstelle  $x$  ab.
- (4) Für die Einträge  $\varepsilon_j^{(k)}$  der Matrizen der relativen Fehler gilt die Abschätzung  $|\varepsilon_j^{(k)}| \leq \varepsilon_{\text{mach}}$ , zumindest dann, wenn die Zwischenergebnisse  $F^{(k)}(\hat{x}^{(k)})$  vor der Rundung komponentenweise im Bereich  $\mathbb{D}$  liegen; siehe [Lemma 5.6](#).

Die Notation der Vorwärtsanalyse mit Hilfe von Jacobimatrizen wird schnell unübersichtlich. Einfacher ist es, zu einer grafischen Notation überzugehen, die auf [Bauer, 1974](#) zurückgeht. Dies führen wir jetzt an der Aufgabe aus [Beispiel 7.3](#) vor. Dazu stellen wir einen **Berechnungsgraphen** (englisch: *computational graph*) auf. Die Eingaben, Zwischenergebnisse und das Endergebnis<sup>2</sup> sind die Knoten dieses Graphen. Knoten, die **Eingaben** bzw. **Ausgaben** darstellen, sind farbig markiert. Die an einem Knoten eingehenden Kanten entsprechen einer unären (wie  $\sqrt{a}$ ) oder binären (wie  $a \cdot b$ ) Operation:



Die in jedem Knoten entstehenden relativen Rundungsfehler notieren wir in der Form  $\varepsilon_1$  etc. An die Kanten schreiben wir die relativen partiellen Konditionszahlen der jeweiligen Operation. Diese geben die Faktoren an, mit denen die entstehenden Rundungsfehler jeweils weitergegeben werden. Mit  $\varepsilon_u$  usw. bezeichnen wir den relativen Fehler des jeweiligen Zwischenergebnisses.

Mit Hilfe eines solchen Graphen können wir folgende Arten von Fragen beantworten:

- (1) Welchen Einfluss hat der Fehler in einem Knoten  $i$  auf einen später im Graph folgenden Knoten  $j$ ?

<sup>2</sup>Im Falle einer vektorwertigen Ausgabe gibt es mehrere Ergebnisknoten.

- (2) Wie hängt der Fehler in einem Knoten  $j$  von den Fehlern in allen früher im Graph vorkommenden Knoten  $i$  ab?

**Frage (1)** können wir beantworten, indem wir die Konditionszahlen entlang des Pfades von  $i$  nach  $j$  multiplizieren. Falls es mehrere solcher Kanten gibt, müssen wir die Ergebnisse aller möglichen Pfade von  $i$  nach  $j$  summieren. Im Beispiel trägt etwa der relative Fehler  $\varepsilon_a$  im Knoten  $a$  (der ausschließlich aus dem relativen Eingangsfehler  $\varepsilon_1$  stammt) zum relativen Fehler  $\varepsilon_y$  im Ergebnisknoten  $y$  mit dem Faktor  $\frac{a}{a+b} \cdot \frac{u}{u+c} = \frac{a}{a+b} \cdot \frac{a+b}{a+b+c} = \frac{a}{a+b+c}$  bei, also gilt

$$\varepsilon_y \doteq \frac{a}{a+b+c} \varepsilon_a + \dots,$$

wobei  $\dots$  für die Fehlerbeiträge aller anderen Knoten steht, von denen ausgehend man ebenfalls auf mindestens einem Pfad zum Ergebnisknoten  $y$  gelangt. (**Quizfrage:** Wie groß ist der Einfluss von  $\varepsilon_c$  auf  $\varepsilon_y$ ?)

Für die Antwort auf **Frage (2)** müssen wir die Fehlerbeiträge entlang aller Pfade, die von irgendwelchen Knoten  $i$  zum betrachteten Knoten  $j$  führen, aufsummieren. Dabei ist auch der im Knoten  $j$  selbst entstehende Rundungsfehler zu berücksichtigen sowie alle Rundungsfehler, die in den Vorgängerknoten  $i$  entstehen. Beispielsweise können wir aus dem Graphen ablesen, dass

$$\varepsilon_u \doteq \varepsilon_4 + \frac{a}{a+b} \varepsilon_1 + \frac{b}{a+b} \varepsilon_2$$

gilt. Für den Ergebnisknoten  $y$  gilt entsprechend

$$\begin{aligned} \varepsilon_y &\doteq \varepsilon_5 + \frac{c}{u+c} \varepsilon_3 + \frac{u}{u+c} \varepsilon_u \\ &= \varepsilon_5 + \frac{c}{u+c} \varepsilon_3 + \frac{u}{u+c} \left[ \varepsilon_4 + \frac{a}{a+b} \varepsilon_1 + \frac{b}{a+b} \varepsilon_2 \right] \\ &= \varepsilon_5 + \frac{c}{a+b+c} \varepsilon_3 + \frac{a+b}{a+b+c} \left[ \varepsilon_4 + \frac{a}{a+b} \varepsilon_1 + \frac{b}{a+b} \varepsilon_2 \right] \\ &= \frac{a}{a+b+c} \varepsilon_a + \frac{b}{a+b+c} \varepsilon_b + \frac{c}{a+b+c} \varepsilon_c + \frac{a+b}{a+b+c} \varepsilon_4 + \varepsilon_5. \end{aligned} \quad (7.12)$$

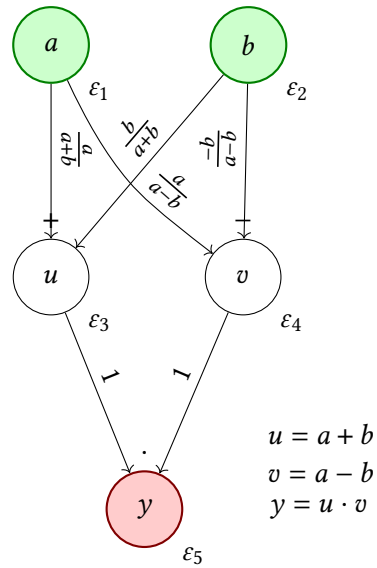
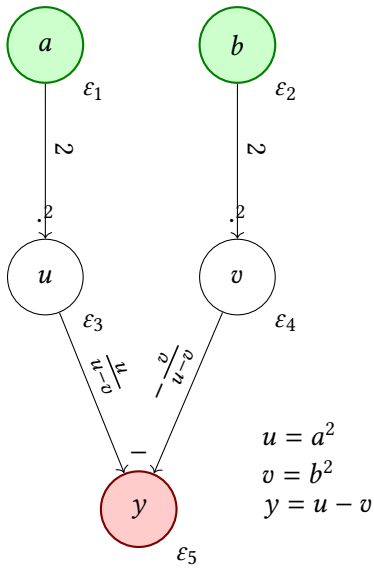
Diese Rechnung unterscheidet sich von unserem früheren Ergebnis

$$\varepsilon_y = \frac{\Delta y}{y} \doteq \frac{\Delta a + \Delta b + \Delta c}{a+b+c} + \frac{a+b}{a+b+c} \varepsilon_1^{(1)} + \varepsilon_1^{(2)} \quad (7.11)$$

lediglich in der Notation, da per Definition  $\varepsilon_a = \Delta a/a$  gilt etc.

Wir betrachten noch weitere Beispiele zur Vorwärtsanalyse.

**Beispiel 7.5** (Vorwärtsanalyse für  $a^2 - b^2$  in zwei Varianten, vgl. [Freund, Hoppe, 2007](#), Kapitel 1.3, Beispiel 5, S.18). Wir werten den Ausdruck  $a^2 - b^2 = (a+b)(a-b)$  auf die beiden offensichtlichen Weisen aus. Dazu gehören folgende Berechnungsgraphen:



Im ersten (linken) Algorithmus erhalten wir

$$\begin{aligned}
 \varepsilon_y &\doteq \varepsilon_5 + \frac{u}{u-v} \varepsilon_u - \frac{v}{u-v} \varepsilon_v \\
 &= \varepsilon_5 + \frac{u}{u-v} [\varepsilon_3 + 2 \varepsilon_1] - \frac{v}{u-v} [\varepsilon_4 + 2 \varepsilon_2] \\
 &= \varepsilon_5 + \frac{a^2}{a^2 - b^2} [\varepsilon_3 + 2 \varepsilon_1] - \frac{b^2}{a^2 - b^2} [\varepsilon_4 + 2 \varepsilon_2] \\
 &= \frac{2 a^2}{a^2 - b^2} \varepsilon_a - \frac{2 b^2}{a^2 - b^2} \varepsilon_b + \frac{a^2}{a^2 - b^2} \varepsilon_3 - \frac{b^2}{a^2 - b^2} \varepsilon_4 + \varepsilon_5,
 \end{aligned} \tag{7.13}$$

im zweiten (rechten) Algorithmus dagegen

$$\begin{aligned}
 \varepsilon_y &\doteq \varepsilon_5 + \varepsilon_u + \varepsilon_v \\
 &= \varepsilon_5 + \left[ \varepsilon_3 + \frac{a}{a+b} \varepsilon_1 + \frac{b}{a+b} \varepsilon_2 \right] + \left[ \varepsilon_4 + \frac{a}{a-b} \varepsilon_1 - \frac{b}{a-b} \varepsilon_2 \right] \\
 &= \left( \frac{a}{a+b} + \frac{a}{a-b} \right) \varepsilon_1 + \left( \frac{b}{a+b} - \frac{b}{a-b} \right) \varepsilon_2 + \varepsilon_3 + \varepsilon_4 + \varepsilon_5 \\
 &= \frac{2 a^2}{a^2 - b^2} \varepsilon_a - \frac{2 b^2}{a^2 - b^2} \varepsilon_b + \varepsilon_3 + \varepsilon_4 + \varepsilon_5.
 \end{aligned} \tag{7.14}$$

Die Terme mit  $\varepsilon_a$  und  $\varepsilon_b$  sind in beiden Verfahren gleich (**Quizfrage:** Warum muss das so sein?). Die verbleibenden Terme können wir im Fall von (7.13) abschätzen durch

$$\left| \frac{a^2}{a^2 - b^2} \varepsilon_3 - \frac{b^2}{a^2 - b^2} \varepsilon_4 + \varepsilon_5 \right| \leq \frac{a^2 + b^2 + |a^2 - b^2|}{|a^2 - b^2|} \varepsilon_{\text{mach}}, \tag{7.15}$$

wohingegen wir im Fall von (7.14)

$$|\varepsilon_3 + \varepsilon_4 + \varepsilon_5| \leq 3 \varepsilon_{\text{mach}} \tag{7.16}$$

erhalten. Wir vergleichen jetzt die beiden oberen Schranken. Die zweite ist genau dann kleiner gleich der ersten, also

$$a^2 + b^2 + |a^2 - b^2| \geq 3 |a^2 - b^2|,$$

wenn  $\left|\frac{a}{b}\right| \in \left[\frac{1}{3}, 3\right]$  gilt, wenn also die Daten  $a$  und  $b$  betragsmäßig relativ nahe beieinander liegen. In diesem Fall ist also die Auswertung von der Form  $(a+b)(a-b)$  vorteilhaft.

**Beispiel 7.6** (Vorwärtsanalyse der  $p$ - $q$ -Formel in zwei Varianten, vgl. Freund, Hoppe, 2007, Kapitel 1.4, Beispiel 1, S.23). Die quadratische Gleichung

$$y^2 + 2 p y - q = 0 \tag{7.17}$$

hat bekanntlich die beiden Lösungen

$$y_1 = -p + \sqrt{p^2 + q} \quad \text{und} \quad y_2 = -p - \sqrt{p^2 + q}.$$

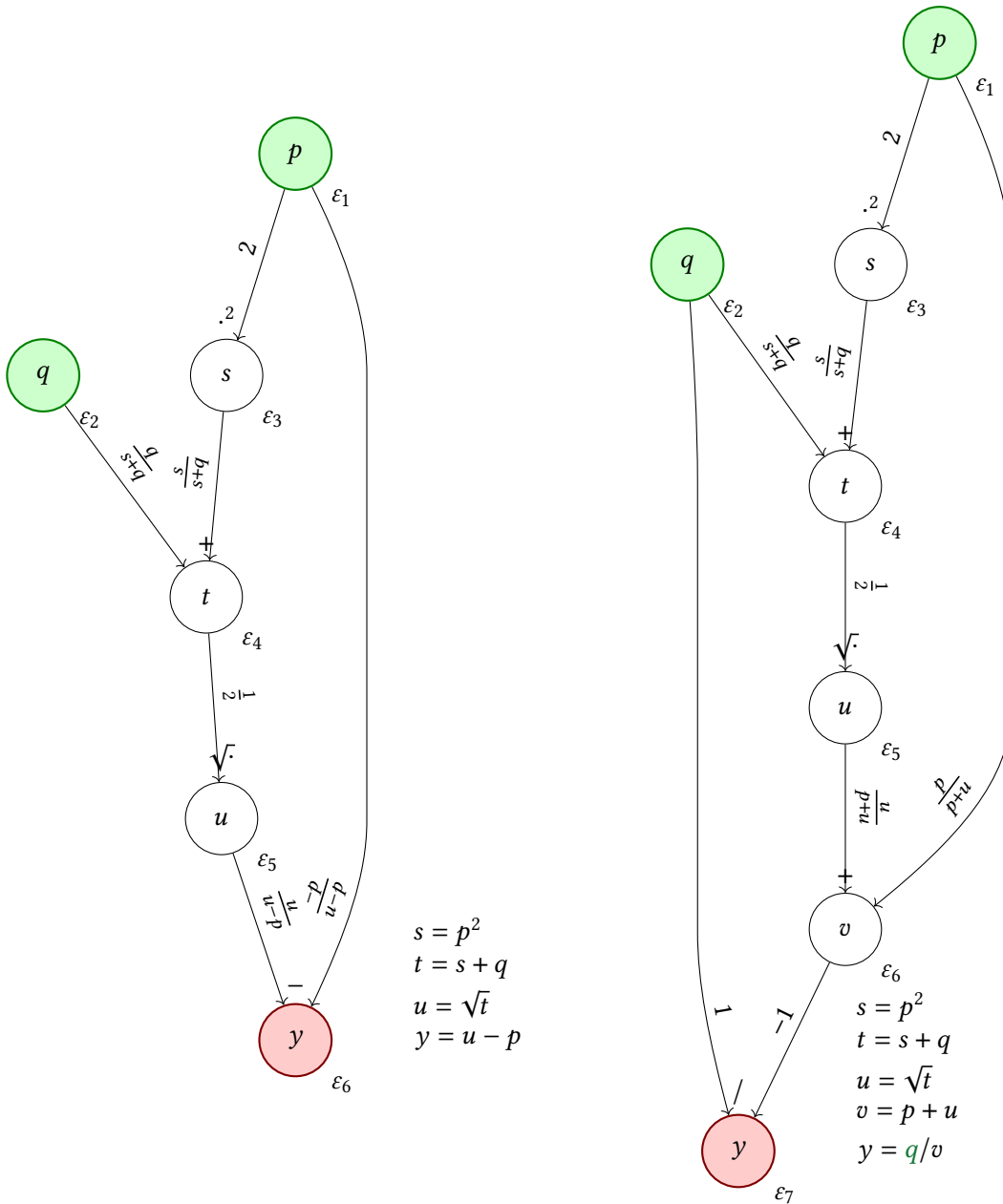
Nach dem *Wurzelsatz von Vieta*

$$y_1 y_2 = -q$$

(elementares Nachrechnen) hat man aber auch die Möglichkeit, eine der beiden Wurzeln aus der Kenntnis der anderen zu berechnen, etwa  $y_1 = -q/y_2$ . Was ist numerisch besser?

Wir betrachten hier nur den Fall  $p > 0$  und  $q > 0$ . Mit diesen Daten ist die Wurzel  $y_1$  numerisch die interessantere. (**Quizfrage:** Warum?)

Wir diskutieren jetzt beide Varianten,  $F(p, q) = -p + \sqrt{p^2 + q}$  zu berechnen: **Der rechte Berechnungsgraph war fehlerhaft und wurde korrigiert.**<sup>RH</sup>



Anstatt wie bisher die komplette Fortpflanzung aller Fehler anzusehen, suchen wir im Berechnungsgraphen nach den Stellen, an denen möglicherweise (abhängig von den Eingangsdaten) große Faktoren stehen. In der ersten Variante sind alle Schritte gut konditioniert (**Quizfrage:** Klar?) bis möglicherweise auf den letzten Schritt. An den Kanten, die zum Ergebnisknoten führen, stehen die Faktoren  $\frac{u}{u-p}$  und  $\frac{-p}{u-p}$ . Wir betrachten exemplarisch nur den ersten, der ausgeschrieben

$$\frac{u}{u-p} = \frac{\sqrt{p^2+q}}{\sqrt{p^2+q}-p} = \frac{\sqrt{p^2+q}}{\sqrt{p^2+q}-p} \frac{\sqrt{p^2+q}+p}{\sqrt{p^2+q}+p} = \frac{p^2+q+p\sqrt{p^2+q}}{q} \geq \frac{2p^2}{q}$$

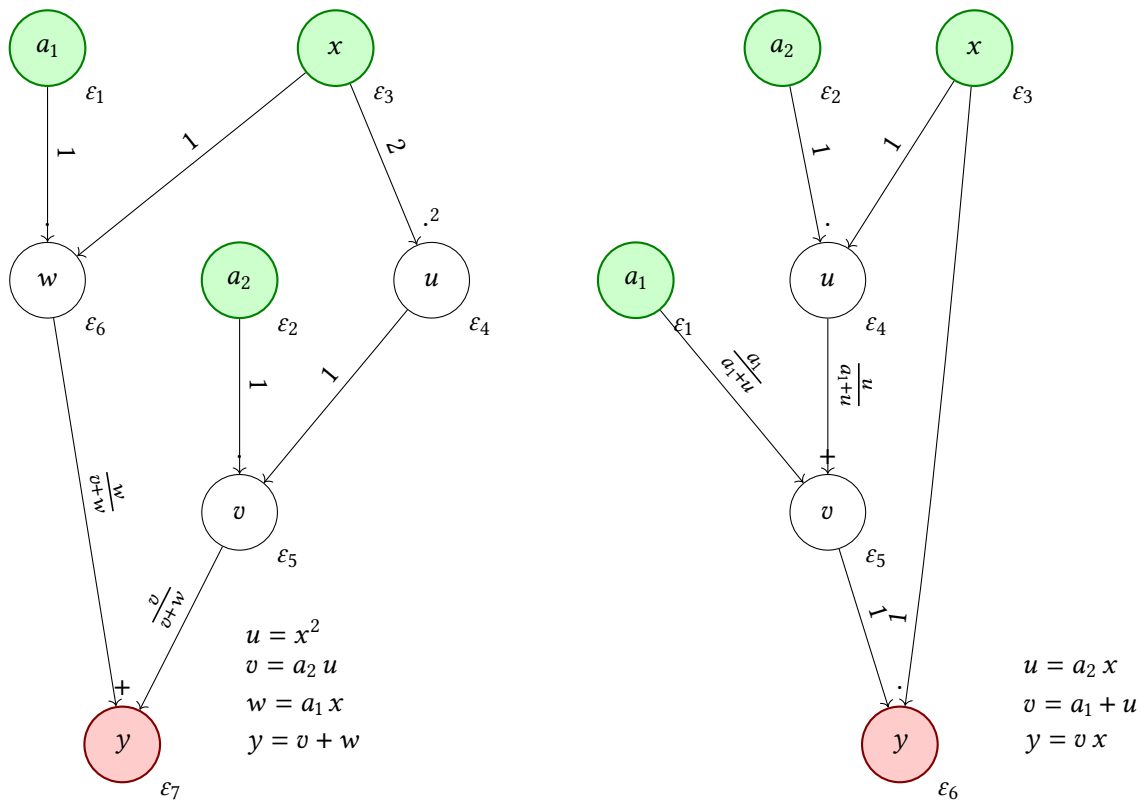
lautet. Die Abschätzung gilt wegen der oben getroffenen Annahme  $p > 0$  und  $q > 0$ . Bei konstantem  $q$  und  $p \rightarrow \infty$  geht also die relative partielle Konditionszahl  $\frac{u}{u-p}$  auf der betrachteten Kante gegen  $\infty$ . Bereits

dieser eine Beitrag zum Gesamtfehler  $\varepsilon_y$  kann also für große Werte von  $p$  inakzeptabel groß werden.

**Quizfrage:** Wie verhält sich der andere Faktor  $\frac{-p}{u-p}$ ?

Betrachten wir nun die zweite Variante. In dieser liegen alle auftretenden Kantenfaktoren betragsmäßig im Intervall  $[0, 2]$ , was kein Problem darstellt.

**Beispiel 7.7** (Horner-Schema). Das **Horner-Schema** ist ein Verfahren zur Auswertung eines durch seine Koeffizienten gegebenen Polynoms an einer gegebenen Stelle  $x \in \mathbb{R}$ . Für Polynome der Bauart  $p(x) = a_1 x + a_2 x^2 = x(a_1 + a_2 x)$  gibt die rechte Darstellung die Auswertungsvorschrift nach dem Horner-Schema an. Wir vergleichen sie mit der »naiven« (linken) Darstellung.



Aus diesen Berechnungsgraphen können wir folgende Darstellungen für den relativen Fehler ablesen: Für die naive Auswertung gilt

$$\begin{aligned} \varepsilon_y &\doteq \varepsilon_7 + \frac{a_1}{a_1 + a_2 x} (\varepsilon_{a_1} + \varepsilon_x + \varepsilon_6) + \frac{a_2 x}{a_1 + a_2 x} (\varepsilon_{a_2} + 2\varepsilon_x + \varepsilon_4 + \varepsilon_5) \\ &= \varepsilon_{a_1} + \varepsilon_x + \varepsilon_6 + \varepsilon_7 + \frac{x}{a_1/a_2 + x} (-\varepsilon_{a_1} + \varepsilon_{a_2} + \varepsilon_x + \varepsilon_4 + \varepsilon_5 - \varepsilon_6), \end{aligned}$$

und für das Horner-Schema haben wir

$$\begin{aligned} \varepsilon_y &\doteq \varepsilon_5 + \varepsilon_6 + \varepsilon_x + \frac{a_1}{a_1 + a_2 x} \varepsilon_{a_1} + \frac{a_2 x}{a_1 + a_2 x} (\varepsilon_{a_2} + \varepsilon_x + \varepsilon_4) \\ &= \varepsilon_{a_1} + \varepsilon_x + \varepsilon_5 + \varepsilon_6 + \frac{x}{a_1/a_2 + x} (-\varepsilon_{a_1} + \varepsilon_{a_2} + \varepsilon_x + \varepsilon_4). \end{aligned}$$



Im Vergleich fällt also das Horner-Schema bei der Vorwärtsanalyse günstiger aus, vor allem, wenn die Auswertungsstelle nahe der Nullstelle  $-a_1/a_2$  liegt. Außerdem benötigt es eine Multiplikation weniger.

Der Vollständigkeit halber geben wir das Horner-Schema noch für allgemeine Polynome vom Grad  $n$  an. Es ist das gängige Verfahren zur Auswertung von Polynomen.

**Algorithmus 7.8** (Horner-Schema für  $p(x) = a_0 + a_1 x + \dots + a_n x^n$ ).

**Eingabe:** Koeffizienten  $a_0, \dots, a_n$  eines Polynoms  $p(x) = a_0 + a_1 x + \dots + a_n x^n$

**Eingabe:** Auswertungsstelle  $x$

**Ausgabe:** Funktionswert  $y = p(x)$

```

1: Setze  $y := a_n$ 
2: for  $k = n - 1, n - 2, \dots, 0$  do
3:   Setze  $y := x y + a_k$ 
4: end for

```

Um die Stabilität eines konkreten Algorithmus zur Lösung einer Aufgabe  $y = F(x)$  zu beurteilen, bietet sich ein Vergleich mit einem hypothetischen Modellalgorithmus an, und zwar mit

$$\bar{F} := \text{rd} \circ F \circ \text{rd}. \quad (7.18)$$

Das heißt, die Eingabe wird auf das verwendete Fließkommagitter gerundet, dann wird ohne Fehler die exakte Funktion  $F$  ausgewertet, und anschließend wird das Ergebnis wieder auf das Fließkommagitter gerundet. Der Modellalgorithmus (7.18) stellt damit den geringstmöglichen Eingriff in die exakte Funktionsauswertung dar, der unter Verwendung von Fließkommaarithmetik möglich ist.

Mit den uns bekannten Techniken können wir den relativen Fehler  $\varepsilon_{\bar{y}}$  von  $\bar{y} := \bar{F}(x)$  gegenüber dem exakten Wert  $y = F(x)$  abschätzen. Es gilt folgendes Resultat:

$$\|\bar{F}(x) - F(x)\|_2 \leq (K(x)\|x\|_2 + \|y\|_2) \varepsilon_{\text{mach}}, \quad (7.19a)$$

$$\left\| \left[ \frac{\bar{F}_i(x) - F_i(x)}{F_i(x)} \right]_{i=1}^m \right\|_{\infty} \leq (k(x) + 1) \varepsilon_{\text{mach}} \quad (7.19b)$$

mit den absoluten und relativen Konditionszahlen  $K(x) = \|F'(x)\|_{2 \rightarrow 2}$  und  $k(x) = \|(k_{ij}(x))\|_{\infty \rightarrow \infty}$ .

**Quizfrage:** Beweis?

Die Schranke auf der rechten Seite von (7.19a) können wir als den unvermeidbaren (absoluten) Fehler ansehen, der allein durch den Übergang zu Fließkommazahlen passiert und der wesentlich durch die Kondition, also eine Eigenschaft der Aufgabe, bestimmt ist. Dieser unvermeidbare Fehler dient als Referenz, mit der wir die Abschätzungen der Rundungsfehler, die wir für konkrete Algorithmen wie etwa im Beispiel 7.5 zeigen können, vergleichen. Entsprechend gibt die Schranke auf der rechten Seite von (7.19b) diesen unvermeidbaren Fehler in relativen Termen an.

**Beachte:** Die Abschätzung (7.19b) ergibt sich gerade aus der Vorwärtsanalyse, wenn man die in den Zwischenschritten auftretenden Rundungsfehler gleich Null setzt, also bis auf den Eingangsfehler und die Rundung des Ergebnisses keine weiteren Rundungsfehler zwischendurch berücksichtigt. In den beiden Algorithmen aus Beispiel 7.5 sind dazu beispielsweise die Zahlen  $\varepsilon_3$  und  $\varepsilon_4$  gleich Null zu setzen.

**Definition 7.9** (Vorwärtsstabilität). Ein mit Maschinenoperationen realisierter Algorithmus  $\widehat{F}(x)$  zur Auswertung von  $y = F(x)$  heißt **vorwärtsstabil** (englisch: **forward stable**) an der Stelle  $x$ , wenn gilt: Der relative Fehler  $\varepsilon_y$  der Ausgabe  $y$  liegt in derselben Größenordnung wie die Schranke in (7.19b):

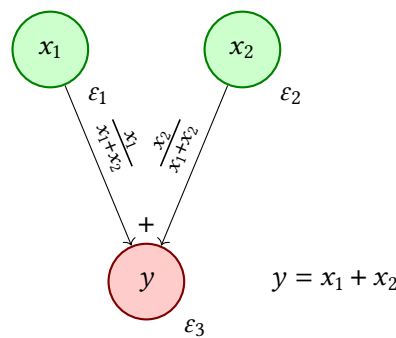
$$\|\varepsilon_y\|_\infty = \frac{\|\widehat{F}(x) - F(x)\|_\infty}{\|F(x)\|_\infty} \leq (C(x) k(x) + 1) \varepsilon_{\text{mach}} \quad (7.20)$$

mit einem «nicht zu großen» Faktor  $C(x)$ .

**Beachte:** Algorithmen mit vektorwertigem Ergebnis sind genau dann vorwärtsstabil an der Stelle  $x$ , wenn sie dort komponentenweise vorwärtsstabil sind.

**Beispiel 7.10** (Vorwärtsstabilität der Maschinenoperationen  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$ ).

*Lemma 5.6 und Satz 6.1 zeigen, dass die Maschinenoperation  $\oplus$  mit vorangegangener Rundung der Eingabe an allen Stellen  $x = (x_1, x_2)^\top$  vorwärtsstabil sind, für die  $x_1, x_2$  und  $x_1 + x_2$  in  $\mathbb{D}$  liegen. Dasselbe gilt für  $\ominus$ ,  $\odot$  und  $\oslash$ . Der Faktor  $C(x)$  kann gleich 1 gewählt werden. (Quizfrage: Klar?) Für  $\oplus$  beispielsweise haben wir:*



Damit ergibt sich  $\varepsilon_y = \varepsilon_3 + \frac{x_1}{x_1 + x_2} \varepsilon_1 + \frac{x_2}{x_1 + x_2} \varepsilon_2$  und daher  $|\varepsilon_y| \leq \varepsilon_{\text{mach}} + \frac{|x_1| + |x_2|}{|x_1 + x_2|} \varepsilon_{\text{mach}} = (k(x) + 1) \varepsilon_{\text{mach}}$ .

**Beispiel 7.11** (Vorwärtsstabilität für  $a^2 - b^2$  in zwei Varianten, vgl. [Beispiel 7.5](#)).

Die Abschätzungen (7.15) und (7.16) zeigen, dass beide Algorithmen zur Bestimmung von  $a^2 - b^2$  aus [Beispiel 7.5](#) «überall» (sofern  $a, b, a + b, a - b$  und  $a^2 - b^2$  in  $\mathbb{D}$  liegen) vorwärtsstabil sind, da die durch die Rundungsfehler «zwischen» entstehenden Terme mit  $\varepsilon_3$  und  $\varepsilon_4$  jeweils in derselben Größenordnung wie die restlichen Fehlerterme liegen.

**Beispiel 7.12** (Vorwärtsstabilität des Innenprodukts).

Das auf die offensichtliche Art und Weise berechnete Innenprodukt  $F(x, y) = x^\top y = \sum_{i=1}^n x_i y_i$  zweier Vektoren  $x, y \in \mathbb{R}^n$  ist vorwärtsstabil, sofern die Eingaben und alle Zwischenergebnisse in  $\mathbb{D}$  liegen.

*Beweis.*

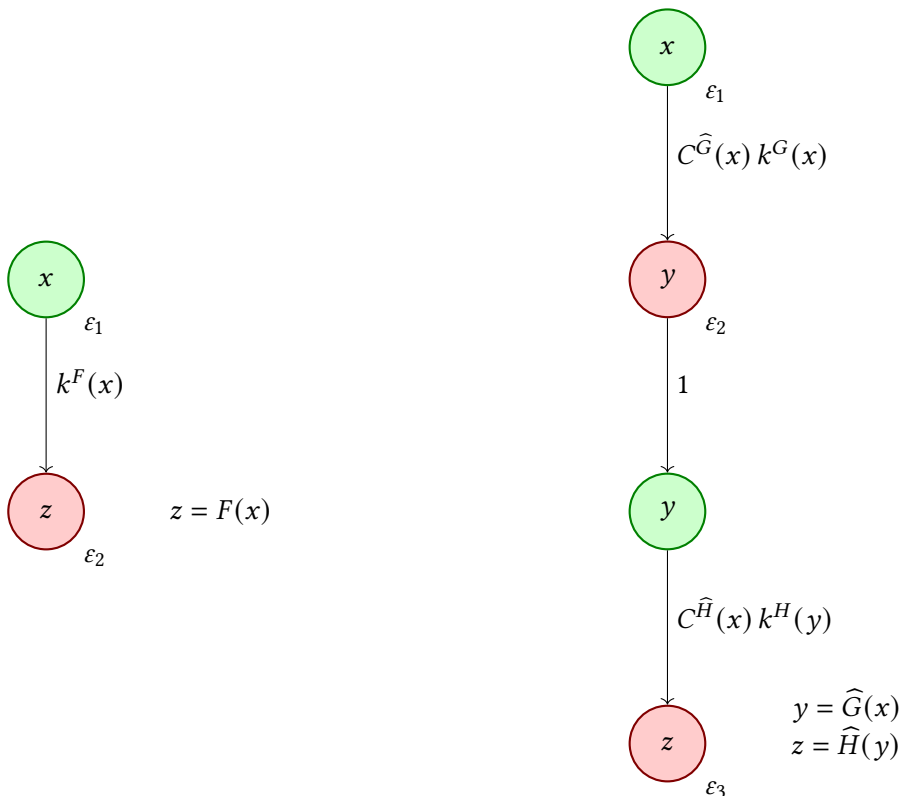
□

**Bemerkung 7.13** (zur Vorwärtsstabilität).

- (1) Der Übergang von einem vorwärtsstabilen Algorithmus zu einem nicht vorwärtsstabilen Algorithmus ist fließend, da in (7.20) verlangt wird, dass der Faktor  $C(x)$  «nicht zu groß» ist.
- (2) Wie wir schon in Beispiel 7.11 gesehen haben, hängt es i. A. von den Eingabedaten ab, welche von zwei (oder mehr) Varianten eines Algorithmus die bessere Schranke für den Vorwärtsfehler liefert. Es bietet sich also an, beide Varianten zu implementieren und mittels einer Fallunterscheidung zur Laufzeit die passendere auszuwählen!

Die Vorwärtsstabilität eines Algorithmus kann also nur dadurch gefährdet sein, dass die bei den Zwischenergebnissen auftretenden Rundungsfehler auf ihrem weiteren Weg hin zum Endergebnis eine derart große Verstärkung erfahren, dass ihr Beitrag zum Gesamtfehler die Schranke auf der rechten Seite von (7.19b) um Größenordnungen überschreitet. Leider bedeutet das auch, dass die Hintereinanderausführung zweier für sich genommen vorwärtsstabiler Algorithmen nicht notwendig ebenfalls wieder vorwärtsstabil ist!

Um das einzusehen, betrachten wir eine Funktion  $F = H \circ G$ . Die folgenden Berechnungsgraphen stellen die Situationen dar, dass wir einerseits ein Verfahren für die komplette Funktion  $F$  betrachten und andererseits zwei Verfahren für die Teilabbildungen  $H$  und  $G$ :<sup>3</sup>



<sup>3</sup>Im Unterschied zu unseren bisherigen Berechnungsgraphen sind hier die Ein- und Ausgaben möglicherweise vektorwertig.

Das erste (hypothetische) Verfahren dient nur zur Festlegung der Referenz, vgl. (7.18) und (7.19b), wobei ein unvermeidbarer relativer Fehler der Größe

$$\begin{aligned}\varepsilon_z &\doteq (k_{ij}^F(x)) \varepsilon_1 + \varepsilon_2 \\ \|\varepsilon_z\|_\infty &\leq (k^F(x) + 1) \varepsilon_{\text{mach}}\end{aligned}\quad (7.21)$$

produziert wird. Bei der Hintereinanderausführung der Algorithmen für die Teilabbildungen  $G$  und  $H$  passiert Folgendes. Das Zwischenergebnis  $y$  weist wegen der Vorwärtsstabilität des Verfahrens  $\widehat{G}$  den relativen Fehler

$$\|\varepsilon_y\|_\infty \leq [C^{\widehat{G}}(x) k^G(x) + 1] \varepsilon_{\text{mach}}$$

auf. Das Endergebnis nach Ausführung des zweiten Teilalgorithmus hat also den relativen Fehler

$$\begin{aligned}\|\varepsilon_z\|_\infty &\leq C^{\widehat{H}}(y) k^H(y) \|\varepsilon_y\|_\infty + \|\varepsilon_3\|_\infty \\ &\leq [C^{\widehat{H}}(y) k^H(y) (C^{\widehat{G}}(x) k^G(x) + 1) + 1] \varepsilon_{\text{mach}}.\end{aligned}\quad (7.22)$$

Der möglicherweise problematische Term ist in **rot** markiert. Der Faktor  $\varepsilon_{\text{mach}}$  repräsentiert den bei der Bildung des Zwischenergebnisses  $y$  unvermeidbaren Rundungsfehler (im Berechnungsgraphen oben  $\varepsilon_2$ ). Dieser wird dann mit der Konditionszahl  $k^H(y)$  (und dem wegen der Vorwärtsstabilität gutartigen Faktor  $C^{\widehat{H}}(y)$ ) durch die nachfolgende Teilabbildung verstärkt und liefert den Beitrag

$$C^{\widehat{H}}(y) k^H(y) \varepsilon_{\text{mach}}$$

zum relativen Gesamtfehler im Ergebnis  $z$ . Um die Vorwärtsstabilität der Hintereinanderausführung  $\widehat{H} \circ \widehat{G}$  nachzuweisen, müssten wir insbesondere diesen Term nach oben gegen die rechte Seite in (7.21) (mal Faktor) abschätzen.

Wegen der Submultiplikativität der Matrixnormen, vgl. (3.27), gilt zwar

$$k^F(x) \leq k^H(y) k^G(x),$$

aber wir würden die umgekehrte Abschätzung benötigen, die nicht gilt.

Der Punkt ist, dass die Teilabbildung  $H$  eine wesentlich größere Konditionszahl besitzen kann als die Gesamtabbildung  $F = H \circ G$ ! Da es hierbei um Eigenschaften der Abbildungen, also um die Aufgabe selbst geht, tritt dieses Problem unabhängig davon auf, welche konkreten Algorithmen man für die Teilaufgaben  $G$  und  $H$  verwenden würde!

**Regel:** Unterteile eine Aufgabe niemals so in Teilaufgaben, dass von diesen eine wesentlich schlechter konditioniert ist als die Gesamtaufgabe!

**Beispiel 7.14.** Die in Vorlesungen zur linearen Algebra eingeübte Technik, die Eigenwerte einer (symmetrischen) Matrix über die Nullstellen ihres charakteristischen Polynoms zu berechnen, ist numerisch kein geeignetes Verfahren! Man kann zeigen, dass die Aufgabe, die Abbildung «Matrix  $\mapsto$  Eigenwerte» gut konditioniert ist und ebenso die erste Teilabbildung «Matrix  $\mapsto$  Koeffizienten des charakteristischen Polynoms». Die Berechnung der Nullstellen eines Polynoms aus seinen Koeffizienten ist jedoch keine gut konditionierte Aufgabe (**Quizfrage:** Beispiel?). Siehe dazu auch Bornemann, 2018, Kapitel 18.

## § 7.2 RÜCKWÄRTSANALYSE

In **Abschnitt 7.1** haben wir die Qualität von Algorithmen mit Hilfe der Vorwärtsanalyse beurteilt. Ein anderer Zugang ist der der Rückwärtsanalyse. Wir betrachten wiederum

$F(x)$  die auszuwertende Funktion,

$\widehat{F}(x)$  die durch den implementierten Algorithmus realisierte Funktion.

Bei der **Rückwärtsanalyse** von  $\widehat{F}$  stellen wir die Frage, ob wir nicht das numerische Ergebnis  $\widehat{F}(x)$  interpretieren können als  $F(\bar{x})$ , also als exaktes Ergebnis zu einer veränderten Eingabe  $\bar{x}$ . Da ein solches  $\bar{x}$  i. A. nicht eindeutig ist, könnten wir die Frage nach einem  $\bar{x}$  mit minimaler Norm stellen, das also die Aufgabe

$$\begin{aligned} &\text{Minimiere} \quad \|\bar{x}\|_\infty, \quad \bar{x} \in \mathbb{R}^n \\ &\text{unter der Bedingung} \quad F(\bar{x}) = \widehat{F}(x) \end{aligned}$$

löst. Da die Bestimmung eines normminimalen  $\bar{x}$  allerdings aufwändig sein kann, begnügt man sich häufig mit einem  $\bar{x}$  mit ausreichend kleiner Norm, das die Bedingung  $F(\bar{x}) = \widehat{F}(x)$  erfüllt.

**Definition 7.15** (Rückwärtsstabilität). *Ein mit Maschinenoperationen realisierter Algorithmus  $\widehat{F}(x)$  zur Auswertung von  $y = F(x)$  heißt **rückwärtsstabil** (englisch: **backward stable**) an der Stelle  $x$ , wenn gilt: Es gibt ein  $\bar{x}$  mit der Eigenschaft  $\widehat{F}(\bar{x}) = F(x)$ , sodass gilt:*

$$\frac{\|\bar{x} - x\|_\infty}{\|x\|_\infty} \leq B(x) \varepsilon_{\text{mach}} \quad (7.23)$$

mit einem «nicht zu großen» Faktor  $B(x)$ .

Um es mit den Worten von **Trefethen, Bau, 1997**, S.104 zu sagen:

A backward stable algorithm gives exactly the right answer to nearly the right question.

**Satz 7.16** (Rückwärtsanalyse für die Maschinenoperationen  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$ ). *Es seien  $x, y \in \mathbb{F}$ .*

(i) *Falls  $x + y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \oplus y = x(1 + \varepsilon) + y(1 + \varepsilon).$$

(ii) *Falls  $x - y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \ominus y = x(1 + \varepsilon) - y(1 + \varepsilon).$$

(iii) *Falls  $x \cdot y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :*

$$x \odot y = x(1 + \varepsilon) \cdot y.$$

(iv) Falls  $y \neq 0$  ist und  $x/y \in \mathbb{D}$  liegt, so gilt mit einer Zahl  $\varepsilon \in \mathbb{R}$ ,  $|\varepsilon| \leq \varepsilon_{\text{mach}}$ :

$$x \oslash y = x(1 + \varepsilon)/y.$$

*Beweis.* Der Beweis folgt sofort aus [Satz 6.1](#). □

**Folgerung 7.17** (Rückwärtsstabilität der Maschinenoperationen  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$ ). Die Maschinenoperationen  $\oplus$ ,  $\ominus$ ,  $\odot$  und  $\oslash$  sind also rückwärtsstabil an den in [Satz 7.16](#) genannten Stellen. Der Faktor  $B(x)$  kann gleich 1 gewählt werden.

**Beispiel 7.18** (Rückwärtsstabilität des Innenprodukts).

Das auf die offensichtliche Art und Weise berechnete Innenprodukt  $F(x, y) = x^T y = \sum_{i=1}^n x_i y_i$  zweier Vektoren  $x, y \in \mathbb{R}^n$  ist rückwärtsstabil, sofern die Eingaben und alle Zwischenergebnisse in  $\mathbb{D}$  liegen.

*Beweis.* □

**Beispiel 7.19** (Rückwärtsstabilität der Polynomauswertung).

Die Auswertung eines Polynoms  $p(x) = a_1 x + a_2 x^2 = x(a_1 + a_2 x)$  nach der naiven Methode und nach dem Horner-Schema, vgl. [Beispiel 7.7](#), sind beide rückwärtsstabil.

*Beweis.* □

Was können wir über den Vorwärtsfehler an der Stelle  $x$  eines rückwärtsstabilen Algorithmus sagen? Dazu sei  $\bar{x}$  wie in [Definition 7.15](#) gewählt. Wir schätzen ab: **Das hier komponentenweise schreiben<sup>RH</sup>** Das hier gilt nur, wenn wir in der Definition der Vorwärtsstabilität mit  $\leq$  arbeiten<sup>GM</sup>

$$\begin{aligned} \frac{\|\widehat{F}(x) - F(x)\|_\infty}{\|F(x)\|_\infty} &= \frac{\|F(\bar{x}) - F(x)\|_\infty}{\|F(x)\|_\infty} \\ &\leq k(x) \frac{\|\bar{x} - x\|_\infty}{\|x\|_\infty} \quad \text{nach Definition der relativen Konditionszahl } k(x) \\ &= B(x) k(x) \varepsilon_{\text{mach}} \quad \text{nach Definition 7.15 der Rückwärtsstabilität.} \end{aligned} \quad (7.24)$$

Damit erfüllt  $\widehat{F}$  also auch die [Definition 7.9](#) der Vorwärtsstabilität, und zwar mit der gleichen Konstante  $C(x) = B(x)$ ! Wir haben damit bewiesen:

**Satz 7.20** (Rückwärtsstabilität impliziert Vorwärtsstabilität). Es sei  $\widehat{F}(x)$  ein mit Maschinenoperationen realisierter Algorithmus zur Auswertung von  $y = F(x)$ . Wenn der Algorithmus  $\widehat{F}$  an der Stelle  $x$  rückwärtsstabil ist, dann ist er dort auch vorwärtsstabil. Kurz: Ein rückwärtsstabiles Verfahren ist auch vorwärtsstabil.

Die Umkehrung von Satz 7.20 gilt aber nicht, wie folgendes Beispiel zeigt:

**Beispiel 7.21** (Äußere Produkte sind vorwärtsstabil, aber nicht rückwärtsstabil). Wir betrachten die Aufgabe, das äußere Produkt zweier Vektoren  $x, y \in \mathbb{R}^n$  zu berechnen, also

$$F(x, y) = x y^T.$$

Das Ergebnis ist eine Matrix, liegt also — entgegen unserem allgemeinen Setting — nicht im  $\mathbb{R}^m$ , sondern im  $\mathbb{R}^{n \times n}$ . Das macht aber keine Schwierigkeiten, da wir  $x y^T$  einfach vektorisieren (z. B. spaltenweise untereinander schreiben) und dann z. B. die  $\infty$ -Norm für Vektoren verwenden können, um Fehler zu messen.

Der offensichtliche Algorithmus,  $\widehat{F}_{ij}(x, y) = \text{rd}(x_i) \odot \text{rd}(y_j)$  zu bestimmen, ist komponentenweise vorwärtsstabil nach Beispiel 7.10, also insgesamt vorwärtsstabil. Dieser Algorithmus ist allerdings nicht rückwärtsstabil, denn das Ergebnis  $(\widehat{F}_{ij}(x))$  wird i. A. keine Rang-1-Matrix sein, d. h.,  $\widehat{F}(x)$  lässt sich überhaupt nicht schreiben als äußeres Produkt von zwei Vektoren.

### § 7.3 ALLGEMEINE STABILITÄTSANALYSE

Das allgemeinste Konzept zur Stabilitätsanalyse kombiniert die Vorwärts- und Rückwärtsanalyse.

**Definition 7.22** (Stabilität). Ein mit Maschinenoperationen realisierter Algorithmus  $\widehat{F}(x)$  zur Auswertung von  $y = F(x)$  heißt **(numerisch) stabil** (englisch: **(numerically) stable**) an der Stelle  $x$ , wenn gilt: Es gibt ein  $\bar{x}$ , sodass gilt:

$$\frac{\|\widehat{F}(x) - F(\bar{x})\|_\infty}{\|F(\bar{x})\|_\infty} \leq C(x) \varepsilon_{\text{mach}} \quad (7.25a)$$

$$\frac{\|\bar{x} - x\|_\infty}{\|x\|_\infty} \leq B(x) \varepsilon_{\text{mach}} \quad (7.25b)$$

mit «nicht zu großen» Faktoren  $B(x)$  und  $C(x)$ .

**Beachte:** Der hypothetische Modellalgorithmus (7.18), also  $\widehat{F} := \text{rd} \circ F \circ \text{rd}$  erfüllt diese Definition mit  $B(x) = C(x) = 1$ . (**Quizfrage:** Beweis?) Die Definition 7.22 kann also so gelesen werden, dass man vom Algorithmus  $\widehat{F}$  qualitativ dieselbe Abschätzung verlangt, die für den Modellalgorithmus gilt.

Um es nochmal mit den Worten von Trefethen, Bau, 1997, S.104 zu sagen:

A stable algorithm gives nearly the right answer to nearly the right question.

Wie man leicht sieht, gilt

$$\text{Rückwärtsstabilität} \Rightarrow \text{Stabilität}.$$

**Quizfrage:** Gilt auch Vorwärtsstabilität  $\Rightarrow$  Stabilität?

Ende der Woche 5

# Kapitel 4 Direkte Lösungsverfahren für lineare Gleichungssysteme

Das Lösen linearer Gleichungssysteme ist eine Grundaufgabe in der numerischen Mathematik. Wir betrachten in diesem Kapitel **direkte Lösungsverfahren**.

## § 8 LR-ZERLEGUNG UND DAS GAUSSSCHE ELIMINATIONSVERFAHREN

Wir betrachten in diesem Abschnitt lineare Gleichungssysteme

$$Ax = b$$

mit quadratischen Matrizen  $A \in \mathbb{R}^{n \times n}$  und rechten Seiten  $b \in \mathbb{R}^n$ . Wenn  $A$  invertierbar (regulär, englisch: *non-singular*) ist, dann ist  $Ax = b$  eindeutig lösbar. Andernfalls ist das System entweder überhaupt nicht lösbar, oder es besitzt unendlich viele Lösungen. Bei numerischen Lösungsverfahren für lineare Gleichungssysteme geht es immer auch darum, festzustellen, ob die Matrix invertierbar ist bzw. numerisch nahe an der Nicht-Invertierbarkeit.

Wir beginnen mit einem einfachen Beispiel:

$$\left. \begin{array}{rrcr} 6x_1 & + & 1x_2 & + & 1x_3 & = & 11 \\ 3x_1 & + & 3x_2 & + & 1x_3 & = & 5 \\ -6x_1 & - & 6x_2 & - & 3x_3 & = & -9 \end{array} \right\} \Leftrightarrow \begin{bmatrix} 6 & 1 & 1 \\ 3 & 3 & 1 \\ -6 & -6 & -3 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 11 \\ 5 \\ -9 \end{pmatrix}.$$

**Beachte:** Jede Zeile der Matrix steht für eine Gleichung, jede Spalte gehört zu einer der Variablen.

Folgende Operationen lassen die Lösungsmenge eines linearen Gleichungssystems unverändert:

- (1) Vertauschen von Zeilen (Gleichungen),
- (2) Multiplikation einer Zeile mit einer Zahl  $\neq 0$ ,
- (3) Addition einer Zeile zu einer anderen Zeile.

Diese Operationen wirken sich natürlich immer auch auf die rechte Seite aus.

Die oben genannten Möglichkeiten der Äquivalenztransformationen nutzt das Gaußsche Eliminationsverfahren, das Sie vermutlich kennen, aus, um die Gestalt des Systems zu vereinfachen. Dabei



wird zunächst während der sogenannten **Vorwärtselimination** (englisch: *forward elimination*) dafür gesorgt, dass die Anzahl der in den Gleichungen vorkommenden Variablen von Zeile zu Zeile um eins abnimmt. Die Matrix wird also auf **obere (rechte) Dreiecksgestalt** gebracht.

In einem ersten Schritt werden dazu die Einträge der ersten Spalte unterhalb der Diagonalen durch eine geeignete Linearkombination der betreffenden Zeile mit der ersten Zeile zu Null gemacht. Wir addieren also in unserem Beispiel zur Zeile 2 das  $-\frac{3}{6}$ -Fache der Zeile 1. Wir addieren weiterhin zur Zeile 3 das  $\frac{6}{6}$ -Fache der Zeile 1. Die erste Zeile selbst bleibt unverändert:

$$\begin{array}{c} \begin{array}{c} \text{Zeile 2} \\ \text{Zeile 3} \end{array} \begin{array}{c} -\frac{3}{6} \\ \frac{6}{6} \end{array} \end{array} \left[ \begin{array}{ccc|c} 6 & 1 & 1 & 11 \\ 3 & 3 & 1 & 5 \\ -6 & -6 & -3 & -9 \end{array} \right] \rightsquigarrow \left[ \begin{array}{ccc|c} 6 & 1 & 1 & 11 \\ 0 & \frac{5}{2} & \frac{1}{2} & -\frac{1}{2} \\ 0 & -5 & -2 & 2 \end{array} \right].$$

Im zweiten Schritt werden die Einträge der zweiten Spalte unterhalb der Diagonalen zu Null gemacht, indem geeignete Linearkombinationen mit der zweiten Zeile gebildet werden. In diesem Fall addieren wir das  $\frac{5}{5/2} = 2$ -Fache der Zeile 2 zur Zeile 3:

$$\begin{array}{c} \text{Zeile 3} \\ \frac{5}{5/2} \end{array} \left[ \begin{array}{ccc|c} 6 & 1 & 1 & 11 \\ 0 & \frac{5}{2} & \frac{1}{2} & -\frac{1}{2} \\ 0 & -5 & -2 & 2 \end{array} \right] \rightsquigarrow \left[ \begin{array}{ccc|c} 6 & 1 & 1 & 11 \\ 0 & \frac{5}{2} & \frac{1}{2} & -\frac{1}{2} \\ 0 & 0 & -1 & 1 \end{array} \right].$$

Die Phase der Vorwärtselimination ist hier beendet. Das System lautet jetzt ausgeschrieben

$$\left. \begin{array}{rrcr} 6x_1 & + & 1x_2 & + & 1x_3 & = & 11 \\ 0x_1 & + & \frac{5}{2}x_2 & + & \frac{1}{2}x_3 & = & -\frac{1}{2} \\ 0x_1 & + & 0x_2 & - & x_3 & = & 1 \end{array} \right\}.$$

Von hier aus können wir es mittels **Rückwärtssubstitution** (englisch: *back substitution*) leicht lösen:

$$\begin{aligned} x_3 &= \frac{1}{-1} (1) &&= -1, \\ x_2 &= \frac{2}{5} \left( -\frac{1}{2} - \frac{1}{2}x_3 \right) &&= 0, \\ x_1 &= \frac{1}{6} (11 - 1x_2 - 1x_3) &&= 2. \end{aligned}$$

Wir wollen das Vorgehen noch einmal rekapitulieren. Wir haben  $A$  durch Zeilenmanipulationen auf eine obere Dreiecksmatrix transformiert. Welche Zeilenmanipulationen waren dazu notwendig? Im ersten Schritt haben wir  $A$  (und die rechte Seite  $b$ ) von links mit der Matrix

$$G^{(1)} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{a_{21}}{a_{11}} & 1 & 0 \\ -\frac{a_{31}}{a_{11}} & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{3}{6} & 1 & 0 \\ -\frac{-6}{6} & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}$$

multipliziert. Es entstand die Matrix

$$A^{(2)} := G^{(1)}A = \begin{bmatrix} 6 & 1 & 1 \\ 0 & \frac{5}{2} & \frac{1}{2} \\ 0 & -5 & -2 \end{bmatrix},$$

deren erste Zeile nach Konstruktion mit der ersten Zeile der Vorgängermatrix  $A^{(1)} := A$  identisch ist. Die Matrix  $A^{(2)}$  wurde dann von links mit

$$G^{(2)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{a_{32}^{(2)}}{a_{22}^{(2)}} & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 2 & 1 \end{bmatrix}$$

multipliziert. Es entstand die Matrix

$$A^{(3)} := G^{(2)} A^{(1)} = \begin{bmatrix} 6 & 1 & 1 \\ 0 & \frac{5}{2} & \frac{1}{2} \\ 0 & 0 & -1 \end{bmatrix} =: R,$$

deren ersten zwei Zeilen mit den ersten zwei Zeilen der Vorgängermatrix  $A^{(2)}$  identisch sind. Die Transformation auf die obere Dreiecksmatrix  $R$  war damit abgeschlossen.

Die obige Rechnung zeigt den Zusammenhang

$$G^{(2)} G^{(1)} A = R \quad \Leftrightarrow \quad A = (G^{(1)})^{-1} (G^{(2)})^{-1} R.$$

Matrizen, die sich wie  $G^{(k)}$  nur darin von der Einheitsmatrix unterscheiden, dass in genau einer Spalte unterhalb der Diagonale Einträge ungleich Null stehen dürfen, heißen **Frobeniusmatrizen**. Genauer sprechen wir von einer  **$k$ -Frobeniusmatrix**, wenn dies in der  $k$ -ten Spalte der Fall ist. Frobeniusmatrizen sind immer invertierbar (**Quizfrage:** Warum?), und sie haben folgende für uns nützliche Eigenschaften:

**Lemma 8.1** (Eigenschaften von Frobeniusmatrizen). *Es sei  $G^{(k)} \in \mathbb{R}^{n \times n}$  eine  $k$ -Frobeniusmatrix,  $k = 1, \dots, n$ . Weiter sei  $G'^{(k)} := G^{(k)} - \text{Id}$ .<sup>1</sup> Dann gilt:*

$$(i) \quad (G^{(k)})^{-1} = \text{Id} - G'^{(k)}.$$

$$(ii) \quad G^{(1)} G^{(2)} \dots G^{(n)} = \text{Id} + \sum_{j=1}^n G'^{(j)}, \text{ ist also eine untere Dreiecksmatrix.}$$

$$(iii) \quad (G^{(1)})^{-1} (G^{(2)})^{-1} \dots (G^{(n)})^{-1} = \text{Id} - \sum_{j=1}^n G'^{(j)}, \text{ ist also eine untere Dreiecksmatrix.}$$

*Beweis.*

□

**Quizfrage:** Gelten die Aussagen (ii) und (iii) auch dann, wenn die Matrizen in anderer Reihenfolge stehen?

<sup>1</sup> $G'^{(k)}$  entspricht also  $G^{(k)}$  mit «genullter» Diagonale.

Die Inverse einer Frobeniusmatrix lässt sich also leicht angeben, man negiert einfach die Vorzeichen der Einträge unterhalb der Diagonale (**Aussage (i)**). Produkte von Inversen (in der richtigen Reihenfolge!) werden einfach spaltenweise zusammengesetzt (**Aussage (iii)**).

In unserem Beispiel ist  $G^{(1)}$  eine 1-Frobeniusmatrix und

$$(G^{(1)})^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{1}{2} & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix},$$

und  $G^{(2)}$  ist eine 2-Frobeniusmatrix mit

$$(G^{(2)})^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 2 & 1 \end{bmatrix}^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -2 & 1 \end{bmatrix}.$$

Wir definieren

$$L := (G^{(1)})^{-1}(G^{(2)})^{-1} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -2 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ -1 & -2 & 1 \end{bmatrix},$$

eine **untere (linke) Dreiecksmatrix**.

Wir haben damit eine Zerlegung der Form

$$A = LR$$

erhalten mit einer unteren Dreiecksmatrix  $L$  mit lauter Einsen auf der Diagonale und einer oberen Dreiecksmatrix  $R$ . Eine solche Zerlegung nennt man eine **LR-Zerlegung** (oder **LU-Zerlegung**, englisch: *LU decomposition*) der quadratischen Matrix  $A$ .

**Frage:** Existiert die LR-Zerlegung immer? Ist sie eindeutig?

Wir beantworten die zweite Frage zuerst.

**Lemma 8.2** (Eindeutigkeit der LR-Zerlegung). *Es sei  $A \in \mathbb{R}^{n \times n}$  eine reguläre Matrix. Wenn  $A$  die beiden LR-Zerlegungen  $A = L_1 R_1 = L_2 R_2$  besitzt, dann gilt  $L_1 = L_2$  und  $R_1 = R_2$ .*

*Beweis.*

□

Wie das Beispiel der regulären (und gut konditionierten<sup>2</sup>) Matrix

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$$

<sup>2</sup> $A^T A = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}$  hat die Eigenwerte  $\frac{3}{2} \pm \frac{1}{2}\sqrt{5}$ , deren Wurzeln die Singulärwerte von  $A$  sind. Der größere durch den kleineren Singulärwert ist also die Wurzel aus  $\frac{3+\sqrt{5}}{3-\sqrt{5}}$ , was gerade etwa  $\kappa(A) \approx 2.62$  ergibt.

zeigt, besitzen nicht alle regulären Matrizen eine LR-Zerlegung. **Quizfrage:** Was geht hier schief?

Wir können das Problem aber beheben, indem wir Zeilenvertauschungen in  $A$  zulassen. Solche Vertauschungen werden durch Permutationsmatrizen realisiert.

**Definition 8.3** (Permutationsmatrix). Eine Matrix  $\Pi \in \mathbb{R}^{n \times n}$  heißt **Permutationsmatrix** (englisch: **permutation matrix**), wenn in jeder Zeile und jeder Spalte genau ein Eintrag gleich 1 ist und alle anderen Einträge gleich 0 sind. Eine Permutationsmatrix  $\Pi$  heißt **Vertauschungsmatrix** oder **Transpositionsmatrix** (englisch: **transposition matrix**), wenn es genau zwei Indizes  $i \neq j$  gibt mit der Eigenschaft  $\Pi_{ij} = 1$ . Wir bezeichnen sie dann auch mit  $\Pi^{(ij)}$ .

Aus Gründen einer vereinfachten Notation lassen wir auch den Fall  $i = j$  zu und setzen  $\Pi^{(ii)} = \text{Id}$  für beliebiges  $1 \leq i \leq n$ .

Jede  $n \times n$ -Permutationsmatrix  $\Pi$  entspricht genau einer Permutation  $\pi$  der Zahlen  $\{1, \dots, n\}$ , also einer bijektiven Abbildung von  $\{1, \dots, n\}$  auf sich selbst. Dabei gilt

$$\Pi = \begin{bmatrix} - & e_{\pi(1)}^\top & - \\ & \vdots & \\ - & e_{\pi(n)}^\top & - \end{bmatrix}$$

mit den Einheitsvektoren  $e_j$  des  $\mathbb{R}^n$ . Permutationsmatrizen sind orthogonale Matrizen, d. h., es gilt insbesondere  $\Pi^{-1} = \Pi^\top$ . Weil jede Permutation als Produkt von Transpositionen (Vertauschungen genau zweier Elemente) geschrieben werden kann, kann jede Permutationsmatrix als Produkt von Transpositionsmatrizen geschrieben werden. Transpositionsmatrizen sind ihre eigenen Inversen.

**Beispiel 8.4** (Permutationsmatrizen). Die Transpositionsmatrix  $\Pi^{(24)}$  der Größe  $5 \times 5$  ist

$$\Pi^{(24)} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

Sie entspricht der Permutation  $\pi(1) = 1, \pi(2) = 4, \pi(3) = 3, \pi(4) = 2, \pi(5) = 5$ , kurz:  $\pi = (14325)$ . Die Permutationsmatrix

$$\Pi = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \end{bmatrix}$$

ist selbst keine Transpositionsmatrix, kann aber als Produkt der Transpositionsmatrizen

$$\Pi = \Pi^{(24)} \Pi^{(45)}$$

geschrieben werden. Die Inverse ist demnach  $\Pi^{-1} = \Pi^{(45)} \Pi^{(24)}$ .

Die Multiplikation einer Matrix  $A$  mit einer Transpositionsmatrix oder allgemeiner mit einer Permutationsmatrix  $\Pi$  von links (also  $\Pi A$ ) entspricht einer Permutation der Zeilen von  $A$ . **Quizfrage:** Was passiert, wenn man  $A$  von rechts mit  $\Pi$  multipliziert, also  $A\Pi$  bildet?

Wir beleuchten zunächst wieder an einem Beispiel, dass Zeilenvertauschungen i. A. erforderlich sind, um die Existenz der LR-Zerlegung zu garantieren. Wir betrachten dazu das Gleichungssystem

$$\left[ \begin{array}{ccc|c} 0 & 1 & 1 & 4 \\ 2 & 1 & 3 & 7 \\ \textcircled{3} & 1 & 6 & 2 \end{array} \right].$$

Bevor wir die Einträge der ersten Spalte unterhalb der Diagonalen zu Null machen können, benötigen wir ein Element  $\neq 0$  im Eintrag  $(1,1)$ . Dazu tauschen wir zwei Zeilen so, dass ein betragsgrößtes Element der ersten Spalte in der ersten Zeile steht. Im vorliegenden Fall tauschen wir also die Zeilen 1 und 3 durch Multiplikation von links mit der Transpositionsmatrix  $\Pi^{(13)}$ . Dann erzeugen wir wie gewohnt die Nullen unterhalb der Hauptdiagonale:

$$\begin{array}{c} \textcircled{3} \\ \textcircled{2} \\ 0 \end{array} \left[ \begin{array}{ccc|c} \textcircled{3} & 1 & 6 & 2 \\ 2 & 1 & 3 & 7 \\ 0 & 1 & 1 & 4 \end{array} \right] \rightsquigarrow \left[ \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 0 & \frac{1}{3} & -1 & \frac{17}{3} \\ 0 & \textcircled{1} & 1 & 4 \end{array} \right].$$

*(Note: The diagram shows row operations: Row 1 ← Row 3, Row 2 ← Row 2 - (2/3)Row 3, Row 3 ← Row 3 - (1/3)Row 3. The pivot element 1 in row 3, column 2 is circled.)*

Im zweiten Schritt tauschen wir zunächst wieder die Zeilen von der zweiten Zeile an abwärts derart, dass ein betragsgrößtes Element der zweiten Spalte in der zweiten Zeile steht. (Das ist hier zwar jetzt nicht unbedingt erforderlich, da das  $(2,2)$ -Element nicht Null ist, erfolgt aber bereits im Hinblick auf die numerische Stabilität des Verfahrens.) Dazu benötigen wir die Transpositionsmatrix  $\Pi^{(23)}$ . Danach erzeugen wir die letzte Null unterhalb der Hauptdiagonale:

$$\begin{array}{c} 3 \\ 0 \\ 0 \end{array} \left[ \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 0 & \textcircled{1} & 1 & 4 \\ 0 & \frac{1}{3} & -1 & \frac{17}{3} \end{array} \right] \rightsquigarrow \left[ \begin{array}{ccc|c} 3 & 1 & 6 & 2 \\ 0 & 1 & 1 & 4 \\ 0 & 0 & -\frac{4}{3} & \frac{13}{3} \end{array} \right].$$

*(Note: The diagram shows row operations: Row 3 ← Row 3 - (1/3)Row 2. The pivot element 1 in row 2, column 2 is circled.)*

Die rot eingekreisten Elemente heißen **Pivotelemente** (englisch: *pivot elements, pivots*).<sup>3</sup>

Die obige Rechnung entspricht der Abfolge

$$G^{(2)} P^{(2)} \overbrace{G^{(1)} P^{(1)} A}^{=: A^{(2)}} = R \quad (8.1)$$

*(Note: The diagram shows the sequence of operations: G^(2) P^(2) G^(1) P^(1) A = R, where G^(1) P^(1) A is grouped as A^(2).)*

mit Transpositionsmatrizen  $P^{(k)}$  (die auch Einheitsmatrizen sein könnten) für den Zeilentausch und Frobeniusmatrizen  $G^{(k)}$  für die Erzeugung der Nullen. In unserem Beispiel war

$$G^{(2)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{1}{3} & 1 \end{bmatrix}, \quad P^{(2)} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}, \quad G^{(1)} = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{2}{3} & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad P^{(1)} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix}$$

<sup>3</sup>vom Französischen *pivot*: Dreh- und Angelpunkt

Die Transpositionsmatrix  $P^{(r)} := \Pi^{(rs)}$  wurde jeweils so gewählt, dass dabei die  $r$ -te Zeile getauscht wurde mit einer späteren Zeile Nummer  $s$ , die die Bedingung

$$|a_{sr}^{(r)}| = \max\{|a_{ir}^{(r)}| \mid r \leq i \leq n\} \quad (8.2)$$

erfüllt. Diese Strategie heißt **Spaltenpivotsuche** (englisch: *column pivoting*). Der Name kommt daher, dass man in der aktuellen Spalte Nummer  $r$  ab dem Diagonalelement abwärts nach einem betragsgrößten Element sucht.

Um (8.1) als LR-Zerlegung von  $A$  mit Zeilentausch zu schreiben, müssen wir die abwechselnden Produkte von Transpositions- und Frobeniusmatrizen auflösen. Dabei hilft folgendes Resultat:

**Lemma 8.5** (Produkt von Transpositions- und Frobeniusmatrizen). *Es sei  $G^{(k)} \in \mathbb{R}^{n \times n}$  eine  $k$ -Frobeniusmatrix und  $G'^{(k)} := G^{(k)} - \text{Id}$ . Mit der Transpositionsmatrix  $\Pi^{(rs)}$  gilt unter der Voraussetzung  $1 \leq k < r, s \leq n$ :*

$$\Pi^{(rs)} G^{(k)} = H^{(k)} \Pi^{(rs)}, \quad (8.3)$$

wobei

$$H^{(k)} := \text{Id} + \Pi^{(rs)} G'^{(k)}$$

wieder eine Frobeniusmatrix ist.

**Beachte:** Das Lemma besagt, dass eine  $k$ -Frobeniusmatrix  $G^{(k)}$ , in der zwei Zeilen  $(r, s)$  unterhalb von  $k$  getauscht werden, als eine andere  $k$ -Frobeniusmatrix  $H^{(k)}$  geschrieben werden kann, in der die zwei Spalten mit denselben Nummern  $(r, s)$  rechts von  $k$  getauscht werden.

Im Fall  $r = s$  ist  $\Pi^{(rs)} = \text{Id}$ ,  $H^{(k)} = G^{(k)}$  und die Aussage (8.3) klar. An Stelle eines formalen Beweises illustrieren wir die Aussage im Fall  $n = 3$ ,  $k = 1$ ,  $r = 2$  und  $s = 3$ . Die Matrix auf der linken Seite von (8.3) ist

$$\underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}}_{=\Pi^{(23)}} \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ g_2 & 1 & 0 \\ g_3 & 0 & 1 \end{bmatrix}}_{=:G^{(k)}} = \begin{bmatrix} 1 & 0 & 0 \\ g_3 & 0 & 1 \\ g_2 & 1 & 0 \end{bmatrix}.$$

Die Matrix auf der rechten Seite von (8.3) ist tatsächlich identisch:

$$\left( \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 \\ g_2 & 0 & 0 \\ g_3 & 0 & 0 \end{bmatrix} \right) \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ g_3 & 1 & 0 \\ g_2 & 0 & 1 \end{bmatrix}}_{=:H^{(k)}} \underbrace{\begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}}_{=\Pi^{(23)}} = \begin{bmatrix} 1 & 0 & 0 \\ g_3 & 0 & 1 \\ g_2 & 1 & 0 \end{bmatrix}.$$

Man sieht hier auch, wie die neue Frobeniusmatrix  $H^{(k)}$  entsteht, nämlich durch den Austausch der Elemente  $(r, k)$  und  $(s, k)$  in der  $k$ -ten Spalte von  $G^{(k)}$ .

Das in

$$G^{(2)} P^{(2)} G^{(1)} P^{(1)} A = R \quad (8.1)$$

vorkommende Produkt  $P^{(2)}G^{(1)}$  erfüllt die Voraussetzungen von Lemma 8.5, da  $G^{(1)}$  eine 1-Frobeniusmatrix ist und  $P^{(2)}$  entweder die Einheitsmatrix oder eine Transpositionsmatrix ist, die die zweite Zeile ( $r = 2$ ) mit einer späteren Zeile ( $s > r$ ) vertauscht. Daher können wir (8.1) auch schreiben in der Form

$$G^{(2)} H^{(1)} P^{(2)} P^{(1)} A = R. \quad (8.4)$$

Wir können diese Gleichung nun von links mit der unteren Dreiecksmatrix (beachte Lemma 8.1 (iii))

$$L := (H^{(1)})^{-1} (G^{(2)})^{-1}$$

multiplizieren und erhalten schließlich eine **LR-Zerlegung mit Pivotisierung**

$$PA = LR, \quad (8.5)$$

wobei die Permutationsmatrix

$$P := P^{(2)} P^{(1)}$$

durch das Produkt der Transpositionsmatrizen entsteht, die in den einzelnen Zerlegungsschritten bestimmt werden. Die Matrix  $L$  kann unter Ausnutzung von Aussage (iii) aus Lemma 8.1 spaltenweise von links nach rechts aufgebaut werden.

In der Praxis wird die LR-Zerlegung wie folgt implementiert:

- (1) Die Matrix  $A$  wird nach und nach *in-place* überschrieben und nimmt spaltenweise die Einträge von  $L$  strikt unterhalb der Diagonale auf und zeilenweise die Einträge von  $R$  ab der Diagonale aufwärts. (**Quizfrage:** Warum ist das Überschreiben der Matrix  $A$  während der noch laufenden LR-Zerlegung möglich?)
- (2) Die Vertauschungen werden in einem separaten Vektor  $p \in \mathbb{N}^n$  gespeichert.
- (3) Statt zunächst mit den Matrizen  $G^{(k)}$  zu arbeiten, diese anschließend gemäß Lemma 8.5 umzuschreiben und dann zu invertieren, wird direkt die Matrix  $L$  Spalte für Spalte aufgebaut.

In unserem obigen Beispiel führt das zu folgenden Schritten. Zunächst initialisieren wir den Permutationsvektor als  $p := (1, 2, \dots, n)^T$ . Dann pivotisieren wir die erste Spalte, d. h., wir tauschen die Zeilen 1 und 3 und entsprechend auch die Einträge 1 und 3 im Vektor  $p$ :

$$\begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} \rightsquigarrow \begin{bmatrix} 3 & 1 & 6 \\ 2 & 1 & 3 \\ 0 & 1 & 1 \end{bmatrix}, \quad p := \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}.$$

Wir erzeugen Nullen in der ersten Spalte (die nicht gespeichert werden) und speichern die (vorläufige) erste Spalte von  $L$  — das ist die erste Spalte von  $(G^{(1)})^{-1}$  — strikt unterhalb der Diagonale an Stelle der entstehenden Nullen. Die erste Zeile von  $R$  liegt ebenfalls schon fest:

$$\begin{bmatrix} 3 & 1 & 6 \\ 2 & 1 & 3 \\ 1 & 1 & 1 \end{bmatrix} \rightsquigarrow \begin{bmatrix} 3 & 1 & 6 \\ 2 & \frac{1}{3} & -1 \\ 0 & 1 & 1 \end{bmatrix}, \quad p := \begin{pmatrix} 3 \\ 2 \\ 1 \end{pmatrix}.$$

Wir pivotisieren jetzt die zweite Spalte, d. h., tauschen die Zeilen 2 und 3 und entsprechend die Einträge 2 und 3 im Vektor  $p$ . Dabei werden auch die Einträge der bereits gespeicherten ersten Spalte von  $L$  mit getauscht:

$$\begin{bmatrix} 3 & 1 & 6 \\ 2 & \frac{1}{3} & -1 \\ 0 & 1 & 1 \end{bmatrix} \rightsquigarrow \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 2 & \frac{1}{3} & -1 \end{bmatrix}, \quad p := \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}.$$

Wir erzeugen jetzt die letzte fehlende Null in der zweiten Spalte (die wieder nicht gespeichert wird) und speichern die zweite Spalte von  $L$  strikt unterhalb der Diagonale an Stelle der entstehenden Null. Die zweite Zeile von  $R$  liegt ebenfalls schon fest:

$$\begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 2 & \frac{1}{3} & -1 \end{bmatrix} \rightsquigarrow \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 2 & \frac{1}{3} & -\frac{4}{3} \end{bmatrix}, \quad p := \begin{pmatrix} 3 \\ 1 \\ 2 \end{pmatrix}.$$

Damit ist die Faktorisierung beendet, und wir erhalten

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 2 & \frac{1}{3} & 1 \end{bmatrix} \quad \text{und} \quad R = \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 0 & 0 & -\frac{4}{3} \end{bmatrix}.$$

Die Permutationsmatrix  $P$  ergibt sich aus dem Vektor  $p$ . Dessen Einträge geben an, in welcher Reihenfolge die Zeilen der  $n \times n$ -Einheitsmatrix in  $P$  erscheinen. Im Beispiel ist

$$P = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}.$$

Wir führen noch die Probe durch:

$$PA = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 1 & 1 \\ 2 & 1 & 3 \\ 3 & 1 & 6 \end{bmatrix} = \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 2 & 1 & 3 \end{bmatrix},$$

$$LR = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 2 & \frac{1}{3} & 1 \end{bmatrix} \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 0 & 0 & -\frac{4}{3} \end{bmatrix} = \begin{bmatrix} 3 & 1 & 6 \\ 0 & 1 & 1 \\ 2 & 1 & 3 \end{bmatrix}.$$

Mit Hilfe von NUMPY und SCIPY kann eine LR-Zerlegung von  $A$  wie folgt bestimmt werden:

```
import numpy as np
import scipy.linalg
A = np.array([[0,1,1],[2,1,3],[3,1,6]])
P, L, R = scipy.linalg.lu(A)
```

**Beachte:** Durch `scipy.linalg.lu` wird eine Zerlegung bestimmt, sodass  $A = P L R$  gilt, also  $P^T A = L R$ . Im Vergleich zu der von uns benutzten Konvention ist die Permutationsmatrix also gerade transponiert!

Mit den gerade beschriebenen Techniken kann man einen Satz über die Existenz von LR-Zerlegung mit Spaltenpivotisierung zeigen. Dazu muss man im Prinzip nur zeigen, dass es immer gelingt, in jeder Spalte von der Diagonale an abwärts mindestens ein Element zu finden, das nicht Null ist. Diese Zerlegung



ist nicht mehr eindeutig (vgl. [Lemma 8.2](#)), da es i. A. verschiedene mögliche Permutationsmatrizen  $P$  gibt, die zu nicht-verschwindenden Pivot-Elementen führen.

**Satz 8.6** (Existenz einer LR-Zerlegung, vgl. [Bornemann, 2018](#), Satz 7.11). *Es sei  $A \in \mathbb{R}^{n \times n}$  regulär. Wird die Permutationsmatrix  $P \in \mathbb{R}^{n \times n}$  gemäß dem Prinzip der Spaltenpivotsuche gewählt, dann besitzt  $PA$  eine (eindeutige) LR-Zerlegung. Zu einem solchen  $P$  existiert also eine untere Dreiecksmatrix  $L$  mit lauter Einsen auf der Diagonale und eine obere Dreiecksmatrix  $R$ , sodass  $PA = LR$  ist. Weiterhin gilt, dass die Einträge von  $L$  betragsweise  $\leq 1$  sind.*

**Quizfrage:** Können Sie begründen, warum die Einträge von  $L$  betragsweise  $\leq 1$  sind, also  $|L| \leq 1$  gilt?

**Beachte:** Durch die Spaltenpivotsuche wird  $P$  nicht unbedingt eindeutig festgelegt, da durchaus mehrere Einträge der aktuellen Spalte gleichzeitig betragsmaximal sein können.

Da klar ist, dass wir selbst bei exakter Rechnung eine Spaltenpivotisierung für die LR-Zerlegung benötigen, werden wir das in Zukunft nicht immer betonen und einfach von LR-Zerlegungen von  $A$  sprechen.

**Frage:** Was sind die Vorteile, eine LR-Zerlegung einer Matrix  $A$  zu erstellen und zu speichern?

Wenn eine LR-Zerlegung von  $A$  vorliegt, können wir jedes lineare Gleichungssystem mit regulärer Koeffizientenmatrix  $A$  schnell lösen:

**Algorithmus 8.7** (Lösung linearer Gleichungssysteme bei Vorliegen einer LR-Zerlegung mit Spaltenpivotisierung).

**Eingabe:** untere Dreiecksmatrix  $L$  mit lauter Einsen auf der Diagonale

**Eingabe:** obere Dreiecksmatrix  $R$

**Eingabe:** Permutationsmatrix  $P$ , sodass gilt:  $PA = LR$

**Eingabe:** Vektor der rechten Seite  $b$

**Ausgabe:** Lösung des linearen Gleichungssystems  $Ax = b$

1: Löse  $Ly = Pb$  durch Vorwärtseinsetzen

2: Löse  $Rx = y$  durch Rückwärtseinsetzen

**Quizfrage:** Können Sie die Korrektheit von [Algorithmus 8.7](#) begründen, also dass durch ihn  $Ax = b$  gelöst wird?

Die Zwischengröße  $y$  berechnet sich dabei aus  $c := Pb$  gemäß

$$y_i := c_i - \sum_{k=1}^{i-1} \ell_{ik} y_k, \quad i = 1, 2, \dots, n. \quad (8.6)$$

Die anschließende Rückwärtssubstitution gelingt durch

$$x_i := \frac{1}{r_{ii}} \left( y_i - \sum_{k=i+1}^n r_{ik} x_k \right), \quad i = n, n-1, \dots, 1. \quad (8.7)$$

Bei bereits vorliegender LR-Zerlegung ist der numerische Aufwand für die Lösung von  $Ax = b$  durch **Algorithmus 8.7** von der Größenordnung  $2n^2$  Grundoperationen (+, −, · und /), **genauer:  $n^2 - n$  für die Lösung des Systems mit  $L$  und  $n^2$  für die Lösung des Systems mit  $R$** . Man spricht dabei auch von **FLOPs** (englisch: *floating point operations*). **Quizfrage:** Ist die Ermittlung des Aufwands klar?

**Quizfrage:** Wie können wir mit Hilfe von **Algorithmus 8.7** mehrere lineare Gleichungssysteme gleichzeitig lösen, die dieselbe Koeffizientenmatrix  $A$  besitzen? Wie groß ist dann der Aufwand?

Der Aufwand für die eigentliche Erzeugung der LR-Zerlegung wird durch folgendes Resultat bestimmt:

**Satz 8.8** (Aufwand der LR-Zerlegung). *Der numerische Aufwand für die Erzeugung der LR-Zerlegung einer regulären Matrix  $A \in \mathbb{R}^{n \times n}$  mit Spaltenpivotsuche ist von der Größenordnung  $\frac{2}{3}n^3 + O(n^2)$ .*

*Beweis.* Die Manipulation der  $k$ -ten Spalte erfordert

- einen Aufwand von  $n - k + 1$  Beträgen und Vergleichen für die Bestimmung des Pivotelements (**Quizfrage:** Warum?)
- einen Aufwand von  $n - k + 1$  Multiplikationen und ebensovielen Additionen für das Bilden der Linearkombination der  $k$ -ten Zeile mit *einer* der darunter stehenden Zeilen (**Quizfrage:** Warum sind das nicht  $n$  Multiplikation und Additionen?)
- Es gibt  $n - k$  solcher darunter stehenden Zeilen, also ist der Aufwand für die Linearkombinationen (die zur Erzeugung der Nullen unterhalb der Hauptdiagonale in Schritt  $k$  führen) insgesamt  $2(n - k)(n - k + 1)$  FLOPs.

Die Summation über die Zeilennummern  $k = 1, \dots, n$  liefert einen Gesamtaufwand proportional zu

$$\sum_{k=1}^n (n - k + 1) = \frac{1}{2}n(n + 1) \in O(n^2)$$

FLOPs für die Spaltenpivotsuche und von

$$\sum_{k=1}^n 2(n - k)(n - k + 1) = \frac{2}{3}n(n^2 - 1) \in \frac{2}{3}n^2 + O(n^2)$$

FLOPs für die Linearkombinationen. Die Spaltenpivotsuche fällt dabei also nicht weiter ins Gewicht.  $\square$

**Beachte:** Die Lösung eines linearen Gleichungssystems  $Ax = b$  mit Hilfe der Gauß-Elimination hat den Aufwand  $\frac{2}{3}n^3 + O(n^2)$ . Wenn man die dabei ohnehin notwendigen Manipulationen von  $A$  in Form einer LR-Zerlegung speichert, dann kann man nachfolgende Gleichungssysteme mit derselben Matrix  $A$  in  $2n^2$  FLOPs lösen. Es ist also sinnvoll und kein zusätzlicher Aufwand, bei der Lösung eines linearen Gleichungssystems mit Hilfe des Gauß-Verfahrens immer die LR-Zerlegung für eine eventuelle zukünftige Verwendung zu speichern.

**Quizfrage:** Kann man die LR-Zerlegung von  $A$  auch dazu verwenden, wenn man lineare Gleichungssysteme mit  $A^T$  lösen muss? Oder mit  $A^2$ ?

Die Erstellung der LR-Zerlegung einer regulären Matrix  $A$  hat noch weitere Vorteile:

- (1) Wir können mittels einer LR-Zerlegung die Determinante berechnen:

$$\det A = \det(P^T L R) = (\det P)(\det L)(\det R) = \pm \det R,$$

wobei sich das Vorzeichen danach richtet, ob die durch  $P$  dargestellte Permutation das Produkt einer geraden oder einer ungeraden Anzahl von Transpositionen ist.

- (2) Wir können mittels einer LR-Zerlegung die Inverse berechnen: Dazu lösen wir  $n$  lineare Gleichungssysteme  $AX = B$  mit rechter Seite  $B = \text{Id}$ . Die Lösung  $X$  ist dann die Inverse von  $A$ .

**Bemerkung 8.9** (Totalpivotisierung). An Stelle der bisher verwendeten Spaltenpivotsuche kann auch eine **totale Pivotsuche** (englisch: **total pivoting**) durchgeführt werden. Dabei wird an Stelle von (8.2) im  $r$ -ten Schritt die Strategie

$$|a_{st}^{(r)}| = \max\{|a_{ij}^{(r)}| \mid r \leq i, j \leq n\} \quad (8.8)$$

verfolgt. Es wird also bei der Suche nach dem betragsgrößten Pivotelement nicht nur nach unten geschaut, sondern auch nach rechts. Entsprechend werden nicht nur die Zeilen  $r$  und  $s$  getauscht, sondern zusätzlich auch die Spalten  $r$  und  $t$ . Das entspricht einer Umsortierung der Variablen im Gleichungssystem. Dieses Verfahren führt zu einer **LR-Zerlegung mit Totalpivotisierung** der Form

$$PAQ^T = LR, \quad (8.9)$$

wobei  $Q$  ebenfalls eine Permutationsmatrix ist. Der Aufwand für die Suche nach dem Pivotelement ist jetzt ebenfalls von der Größenordnung  $O(n^3)$ , sodass das Gesamtverfahren für die LR-Zerlegung noch immer von der Größenordnung  $O(n^3)$  FLOPs ist. **Quizfrage:** Wieviele Vergleiche sind bei der totalen Pivotsuche genau erforderlich?

Der Vollständigkeit halber geben wir auch für den Fall der LR-Zerlegung mit totaler Pivotisierung das Verfahren zur Lösung eines linearen Gleichungssystems  $Ax = b$  analog zu **Algorithmus 8.7** an:

**Algorithmus 8.10** (Lösung linearer Gleichungssysteme bei Vorliegen einer LR-Zerlegung mit Totalpivotisierung).

**Eingabe:** untere Dreiecksmatrix  $L$  mit lauter Einsen auf der Diagonale

**Eingabe:** obere Dreiecksmatrix  $R$

**Eingabe:** Permutationsmatrizen  $P$  und  $Q$ , sodass gilt:  $PAQ^T = LR$

**Eingabe:** Vektor der rechten Seite  $b$

**Ausgabe:** Lösung des linearen Gleichungssystems  $Ax = b$

1: Löse  $Ly = Pb$  durch Vorwärtseinsetzen

2: Löse  $Rx = y$  durch Rückwärtseinsetzen

3: Setze  $x := Q^T x$

Ende der Woche 6

## § 9 FEHLERANALYSE BEI DER LÖSUNG LINEARER GLEICHUNGSSYSTEME

Wir untersuchen jetzt die Stabilität der algorithmischen Komponenten bei der Lösung linearer Gleichungssysteme. Wir beginnen mit der Lösung von Systemen mit unteren und oberen Dreiecksmatrizen,  $Lx = b$  bzw.  $Rx = b$ . Etwas allgemeiner als eigentlich erforderlich bestehen wir dabei nicht darauf, dass die Diagonale von  $L$  aus lauter Einsen besteht.<sup>4</sup> Wir betrachten folgende konkrete Implementierungen der Vorwärtssubstitution (8.6) bzw. Rückwärtssubstitution (8.7):

**Algorithmus 9.1** (Lösung linearer Gleichungssysteme mit unterer Dreiecksmatrix durch Vorwärtssubstitution).

**Eingabe:** reguläre untere Dreiecksmatrix  $L$  (nicht notwendig mit Einsen auf der Diagonale)

**Eingabe:** Vektor der rechten Seite  $b$

**Ausgabe:** Lösung des linearen Gleichungssystems  $Ly = c$

```

1: for  $i = 1, 2, \dots, n$  do
2:   Setze  $s_i := c_i$ 
3:   for  $j = i - 1, i - 2, \dots, 1$  do
4:     Setze  $s_i := s_i - \ell_{ij} y_j$ 
5:   end for
6:   Setze  $y_i := \frac{s_i}{\ell_{ii}}$ 
7: end for
```

**Algorithmus 9.2** (Lösung linearer Gleichungssysteme mit oberer Dreiecksmatrix durch Rückwärtssubstitution, vgl. Higham, 2002, Algorithm 8.1).

**Eingabe:** reguläre obere Dreiecksmatrix  $R$

**Eingabe:** Vektor der rechten Seite  $b$

**Ausgabe:** Lösung des linearen Gleichungssystems  $Rx = b$

```

1: for  $i = n, n - 1, \dots, 1$  do
2:   Setze  $s_i := b_i$ 
3:   for  $j = i + 1, i + 2, \dots, n$  do
4:     Setze  $s_i := s_i - r_{ij} x_j$ 
5:   end for
6:   Setze  $x_i := \frac{s_i}{r_{ii}}$ 
7: end for
```

**Beachte:** Beide Algorithmen sind so formuliert, dass die Summation in  $s_i$  der Beiträge «von der Diagonale aus nach außen» erfolgt. Diese Reihenfolge erweist sich als günstig für die Fehlerabschätzung in den Sätzen 9.4 und 9.5.

Für die Fehlerabschätzung benötigen wir folgendes Lemma.

**Lemma 9.3** (vgl. Higham, 2002, Lemma 3.1). Für  $n \in \mathbb{N}$  seien  $\delta_1, \dots, \delta_n$  Zahlen mit  $|\delta_i| \leq \varepsilon_{\text{mach}}$  und

<sup>4</sup>Das ist praktisch u. a. bereits im Vorgriff auf die Cholesky-Zerlegung in Abschnitt 10.2.

$\rho_1, \dots, \rho_n \in \{\pm 1\}$ . Falls  $n \varepsilon_{\text{mach}} \leq C < 1$  gilt, dann ist

$$\prod_{i=1}^n (1 + \delta_i)^{\rho_i} = 1 + \theta,$$

und  $\theta$  erfüllt die Abschätzung

$$|\theta| \leq \gamma_n := \frac{n \varepsilon_{\text{mach}}}{1 - n \varepsilon_{\text{mach}}} \leq \frac{n \varepsilon_{\text{mach}}}{1 - C}.$$

*Beweis.* Wir führen einen Induktionsbeweis. Für  $n = 1$  gilt im Fall von  $\rho_1 = 1$

$$\theta = \delta_1$$

und daher  $|\theta| \leq \varepsilon_{\text{mach}} < \frac{1 \varepsilon_{\text{mach}}}{1 - 1 \varepsilon_{\text{mach}}}$ . Im Fall von  $\rho_1 = -1$  gilt

$$\theta = \frac{1}{1 + \delta_1} - 1 = -\text{GM} \frac{\delta_1}{1 + \delta_1},$$

und damit

$$|\theta| = \frac{|\delta_1|}{|1 + \delta_1|} \leq \frac{1 \varepsilon_{\text{mach}}}{1 - 1 \varepsilon_{\text{mach}}}.$$

~~und wegen  $|1 + \delta_1| \geq 1$  gilt ebenfalls  $|\theta| \leq \varepsilon_{\text{mach}} < \frac{1 \varepsilon_{\text{mach}}}{1 - 1 \varepsilon_{\text{mach}}}$  GM~~

Nun zum Induktionsschritt. Die Behauptung sei bereits bewiesen für gegebenes  $n$ . Im Fall von  $\rho_{n+1} = 1$  gilt

$$\prod_{i=1}^{n+1} (1 + \delta_i)^{\rho_i} = (1 + \theta)(1 + \delta_{n+1}) =: 1 + \theta',$$

also ist  $\theta' = \theta(1 + \delta_{n+1}) + \delta_{n+1}$ . Wir schätzen ab:

$$\begin{aligned} |\theta'| &\leq |\theta| |1 + \delta_{n+1}| + |\delta_{n+1}| \\ &\leq \frac{n \varepsilon_{\text{mach}}}{1 - n \varepsilon_{\text{mach}}} (1 + \varepsilon_{\text{mach}}) + \varepsilon_{\text{mach}} \quad (\text{nach Induktionsvoraussetzung}) \\ &= \frac{n \varepsilon_{\text{mach}} + n \varepsilon_{\text{mach}}^2 + \varepsilon_{\text{mach}} - n \varepsilon_{\text{mach}}^2}{1 - n \varepsilon_{\text{mach}}} \\ &= \frac{(n + 1) \varepsilon_{\text{mach}}}{1 - n \varepsilon_{\text{mach}}} \\ &\leq \frac{(n + 1) \varepsilon_{\text{mach}}}{1 - (n + 1) \varepsilon_{\text{mach}}}. \end{aligned}$$

Im Fall von  $\rho_{n+1} = -1$  haben wir hingegen

$$\prod_{i=1}^{n+1} (1 + \delta_i)^{\rho_i} = \frac{1 + \theta}{1 + \delta_{n+1}} =: 1 + \theta',$$

also ist

$$\theta' = \frac{1 + \theta}{1 + \delta_{n+1}} - 1 = \frac{\theta - \delta_{n+1}}{1 + \delta_{n+1}}.$$

Wir schätzen ab: (Die folgende Gleichung wurde ausgetauscht.)<sup>GM</sup>

$$\begin{aligned}
 |\theta'| &= \frac{|\theta - \delta_{n+1}|}{|1 + \delta_{n+1}|} \\
 &\leq \frac{1}{|1 + \delta_{n+1}|} (|\theta| + |\delta_{n+1}|) \\
 &\leq \frac{1}{1 - \varepsilon_{\text{mach}}} \left( \frac{n\varepsilon_{\text{mach}}}{1 - n\varepsilon_{\text{mach}}} + \varepsilon_{\text{mach}} \right) \\
 &= \frac{1}{1 - \varepsilon_{\text{mach}}} \left( \frac{n\varepsilon_{\text{mach}} + \varepsilon_{\text{mach}} - n\varepsilon_{\text{mach}}^2}{1 - n\varepsilon_{\text{mach}}} \right) \\
 &\leq \frac{(n+1)\varepsilon_{\text{mach}}}{(1 - \varepsilon_{\text{mach}})(1 - n\varepsilon_{\text{mach}})} \\
 &= \frac{(n+1)\varepsilon_{\text{mach}}}{1 + n\varepsilon_{\text{mach}}^2 - \varepsilon_{\text{mach}} - n\varepsilon_{\text{mach}}^2} \\
 &\leq \frac{(n+1)\varepsilon_{\text{mach}}}{1 - (n+1)\varepsilon_{\text{mach}}}
 \end{aligned}$$

Damit ist in beiden Fällen der Induktionsschritt gezeigt. Die Aussage  $\gamma_n \leq \frac{n\varepsilon_{\text{mach}}}{1-c}$  ist wegen  $1 - n\varepsilon_{\text{mach}} \geq 1 - c$  klar.  $\square$

Wir können nun ein Resultat zur Rückwärtsstabilität von [Algorithmus 9.1](#) zeigen, wobei nur die Matrix  $L$  gestört ist und die rechte Seite nicht.

**Satz 9.4** (Rückwärtsanalyse von [Algorithmus 9.1](#), vgl. [Higham, 2002](#), Theorem 8.3). *Es sei  $L \in \mathbb{F}^{n \times n}$  eine reguläre untere Dreiecksmatrix,  $c \in \mathbb{F}^n$  und  $n\varepsilon_{\text{mach}} \leq C < 1$ . Weiter sei  $\hat{y}$  die von [Algorithmus 9.1](#) in Maschinenoperationen berechnete Näherungslösung. Dann gibt es eine untere Dreiecksmatrix  $\Delta L$ , sodass gilt*

$$(L + \Delta L) \hat{y} = c \quad \text{mit} \quad |\Delta \ell_{ij}| \leq \begin{cases} \gamma_i |\ell_{ii}| & \text{für } i = j, \\ \gamma_{i-j} |\ell_{ij}| & \text{für } i > j \end{cases}$$

$$\text{mit } \gamma_i := \frac{i\varepsilon_{\text{mach}}}{1-i\varepsilon_{\text{mach}}} \leq \frac{i\varepsilon_{\text{mach}}}{1-C}.$$

**Quizfrage:** Inwiefern ist dies ein Resultat über die Rückwärtsstabilität von [Algorithmus 9.1](#)?

*Beweis.* Wir beginnen mit  $\hat{y}_1 = \hat{s}_1 \oslash \ell_{11} = \frac{c_1}{\ell_{11}}(1 + \delta_1)$ . Dabei ist  $\hat{s}_1 = s_1 = c_1$  (ohne Fehler, da  $c \in \mathbb{F}^n$ ) und  $|\delta_1| \leq \varepsilon_{\text{mach}}$ . Daraus folgt

$$\frac{1}{1 + \delta_1} \ell_{11} \hat{y}_1 = (1 + \theta) \ell_{11} \hat{y}_1 = c_1$$

mit einem  $|\theta| \leq \gamma_1$ . Es gilt also  $\Delta \ell_{11} := \theta \ell_{11}$  mit  $|\Delta \ell_{11}| = |\theta| |\ell_{11}| \leq \gamma_1 |\ell_{11}|$ .

An Stelle einer formalen Induktion zeigen wir noch die nächsten zwei Schritte, zunächst also  $i = 2$ . Dabei gilt

$$\begin{aligned}
 \hat{s}_2 &= c_2 \ominus (\ell_{21} \odot \hat{y}_1) \\
 &= c_2 \ominus (\ell_{21} \hat{y}_1 (1 + \varepsilon_1)) \\
 &= c_2 (1 + \varepsilon_2) - \ell_{21} \hat{y}_1 (1 + \varepsilon_1) (1 + \varepsilon_2)
 \end{aligned}$$

und weiter

$$\widehat{y}_2 = \widehat{s}_2 \oslash \ell_{22} = \frac{\widehat{s}_2}{\ell_{22}} (1 + \delta_2).$$

Daraus folgt

$$\frac{1}{1 + \delta_2} \ell_{22} \widehat{y}_2 = \widehat{s}_2 = c_2 (1 + \varepsilon_2) - \ell_{21} \widehat{y}_1 (1 + \varepsilon_1) (1 + \varepsilon_2).$$

Sortieren der Terme ergibt

$$\frac{1}{1 + \delta_2} \frac{1}{1 + \varepsilon_2} \ell_{22} \widehat{y}_2 + \ell_{21} \widehat{y}_1 (1 + \varepsilon_1) = c_2.$$

Wir setzen  $1 + \theta := \frac{1}{1 + \delta_2} \frac{1}{1 + \varepsilon_2}$ . Dann gilt  $|\theta| \leq \gamma_2$  und  $|\varepsilon_1| \leq \gamma_1$ , und mit  $\Delta \ell_{22} := \theta \ell_{22}$  und  $\Delta \ell_{21} := \varepsilon_1 \ell_{21}$  ist auch die zweite Zeile von  $(L + \Delta L) \widehat{y} = c$  erfüllt.

Schließlich betrachten wir noch den Index  $i = 3$ . Dabei gilt

$$\begin{aligned} \widehat{s}_3 &= (c_3 \ominus (\ell_{32} \oslash \widehat{y}_2)) \ominus (\ell_{31} \oslash \widehat{y}_1) \\ &= (c_3 \ominus (\ell_{32} \widehat{y}_2 (1 + \varepsilon_3))) \ominus (\ell_{31} \oslash \widehat{y}_1) \\ &= (c_3 (1 + \varepsilon_4) - \ell_{32} \widehat{y}_2 (1 + \varepsilon_3) (1 + \varepsilon_4)) \ominus (\ell_{31} \widehat{y}_1 (1 + \varepsilon_5)) \\ &= c_3 (1 + \varepsilon_4) (1 + \varepsilon_6) - \ell_{32} \widehat{y}_2 (1 + \varepsilon_3) (1 + \varepsilon_4) (1 + \varepsilon_6) - \ell_{31} \widehat{y}_1 (1 + \varepsilon_5) (1 + \varepsilon_6) \end{aligned}$$

und weiter

$$\widehat{y}_3 = \widehat{s}_3 \oslash \ell_{33} = \frac{\widehat{s}_3}{\ell_{33}} (1 + \delta_3).$$

Daraus folgt

$$\frac{1}{1 + \delta_3} \ell_{33} \widehat{y}_3 = \widehat{s}_3 = c_3 (1 + \varepsilon_4) (1 + \varepsilon_6) - \ell_{32} \widehat{y}_2 (1 + \varepsilon_3) (1 + \varepsilon_4) (1 + \varepsilon_6) - \ell_{31} \widehat{y}_1 (1 + \varepsilon_5) (1 + \varepsilon_6).$$

Sortieren der Terme ergibt

$$\frac{1}{1 + \delta_3} \frac{1}{1 + \varepsilon_4} \frac{1}{1 + \varepsilon_6} \ell_{33} \widehat{y}_3 + \ell_{32} \widehat{y}_2 (1 + \varepsilon_3) + \ell_{31} \widehat{y}_1 \frac{1 + \varepsilon_5}{1 + \varepsilon_4} = c_3.$$

Wir setzen  $1 + \theta := \frac{1}{1 + \delta_3} \frac{1}{1 + \varepsilon_4} \frac{1}{1 + \varepsilon_6}$  und  $1 + \theta' := \frac{1 + \varepsilon_5}{1 + \varepsilon_4}$ . Dann gilt  $|\theta| \leq \gamma_3$  und  $|\theta'| \leq \gamma_2$  sowie  $|\varepsilon_3| \leq \gamma_1$ , und mit  $\Delta \ell_{33} := \theta \ell_{33}$ ,  $\Delta \ell_{32} := \varepsilon_3 \ell_{32}$  sowie  $\Delta \ell_{31} := \theta' \ell_{31}$  ist auch die dritte Zeile von  $(L + \Delta L) \widehat{y} = c$  erfüllt.  $\square$

Ein analoges Resultat gilt auch für Gleichungssysteme mit oberer Dreiecksmatrix und [Algorithmus 9.2](#):

**Satz 9.5** (Rückwärtsanalyse von [Algorithmus 9.2](#), vgl. [Higham, 2002](#), Theorem 8.3). *Es sei  $R \in \mathbb{F}^{n \times n}$  eine reguläre obere Dreiecksmatrix,  $b \in \mathbb{F}^n$  und  $n \varepsilon_{\text{mach}} \leq C < 1$ . Weiter sei  $\widehat{x}$  die von [Algorithmus 9.2](#) in Maschinenoperationen berechnete Näherungslösung. Dann gibt es eine obere Dreiecksmatrix  $\Delta R$ , sodass gilt*

$$(R + \Delta R) \widehat{x} = b \quad \text{mit} \quad |\Delta r_{ij}| \leq \begin{cases} \gamma_{n+1-i} |r_{ii}| & \text{für } i = j, \\ \gamma_{j-i} |r_{ij}| & \text{für } i < j \end{cases}$$

$$\text{mit } \gamma_i := \frac{i \varepsilon_{\text{mach}}}{1 - i \varepsilon_{\text{mach}}} \leq \frac{i \varepsilon_{\text{mach}}}{1 - C}.$$

Die größte in den Abschätzungen von [Satz 9.4](#) und [Satz 9.5](#) vorkommende Konstante ist

$$\gamma_n = \frac{n \varepsilon_{\text{mach}}}{1 - n \varepsilon_{\text{mach}}} \leq \frac{n \varepsilon_{\text{mach}}}{1 - C}.$$

Daraus folgt sofort folgendes Resultat:

**Folgerung 9.6.** *Unter den Voraussetzungen von [Satz 9.4](#) und [Satz 9.5](#) gilt:*

$$(L + \Delta L) \hat{y} = c \quad \text{mit} \quad |\Delta L| \leq \gamma_n |L| \leq \frac{n}{1 - C} |L| \varepsilon_{\text{mach}}, \quad (9.1a)$$

$$(R + \Delta R) \hat{x} = b \quad \text{mit} \quad |\Delta R| \leq \gamma_n |R| \leq \frac{n}{1 - C} |R| \varepsilon_{\text{mach}}. \quad (9.1b)$$

Wenn

Wir verwenden hier und im Folgenden die Notation  $|A|$  für den komponentenweisen Betrag einer Matrix  $A$ , also für diejenige Matrix, deren Einträge sich aus den Absolutbeträgen der Einträge von  $A$  ergeben.

## § 9.1 FALL OHNE PIVOTSUCHE

Es bleibt noch eine Fehlerabschätzung für die Erstellung einer LR-Zerlegung für eine Matrix  $A$  herzuleiten. Wir tun dies zuerst für den Fall, dass die LR-Zerlegung ohne Pivotisierung berechnet wird. Wir betrachten dazu konkret das folgende Verfahren:

**Algorithmus 9.7** (Berechnung der LR-Zerlegung ohne Pivotisierung, vgl. [Golub, van Loan, 1983](#), Algorithm 3.2.1).

**Eingabe:** reguläre Matrix  $A$

**Ausgabe:** Faktoren  $L$  und  $R$  der LR-Zerlegung  $A = LR$

```

1: for  $k = 1, 2, \dots, n - 1$  do
2:   Setze  $I := k + 1 : n = \{k + 1, k + 2, \dots, n\}$ 
3:   Setze  $A(I, k) := A(I, k) / A(k, k)$ 
4:   Setze  $A(I, I) := A(I, I) - A(I, k) A(k, I)$ 
5: end for
6: Extrahiere  $L$  aus dem strikteren unteren Dreieck von  $A$  und ergänze Einsen auf der Diagonale
7: Extrahiere  $R$  aus dem oberen Dreieck von  $A$ 
```

**Quizfrage:** Ist klar, dass [Algorithmus 9.7](#) das oben («In der Praxis wird die LR-Zerlegung wie folgt implementiert») anhand von Beispielen erklärte Verfahren der Transformation *in-place* (im Fall ohne Pivotisierung) implementiert?

Die [Zeilen 6](#) und [7](#) erzeugen die Faktoren  $L$  und  $R$  der Zerlegung  $A = LR$ , ohne weitere Rundungsfehler hinzuzufügen. **Quizfrage:** Klar?



**Satz 9.8** (Rückwärtsanalyse der LR-Zerlegung ohne Pivotisierung, vgl. [Golub, van Loan, 1983](#), Theorem 3.3.1). Es sei  $A \in \mathbb{F}^{n \times n}$  eine reguläre Matrix. Weiter seien  $\widehat{L}$  und  $\widehat{R}$  die von [Algorithmus 9.7](#) in Maschinenoperationen berechnete Näherungslösung der LR-Zerlegung von  $A = LR$ . Dabei nehmen wir an, dass in jeder Iteration das Pivotelement  $\widehat{A}(k, k) \neq 0$  ist. Dann gibt es eine Matrix  $\Delta A$ , sodass gilt:

$$\widehat{L} \widehat{R} = A + \Delta A \quad \text{mit} \quad |\Delta A| \leq 3(n-1)(|A| + |\widehat{L}| |\widehat{R}|) \varepsilon_{\text{mach}}. \quad (9.2)$$

**Beachte:**  $|\widehat{L}| |\widehat{R}|$  bezeichnet das Produkt der Matrizen  $|\widehat{L}|$  und  $|\widehat{R}|$ , also  $(|\widehat{L}| |\widehat{R}|)_{ij} = \sum_{k=1}^n |\widehat{\ell}_{ik}| |\widehat{r}_{kj}|$ .

*Beweis.* Wir führen einen Induktionsbeweis über die Dimension  $n$ . Für  $n = 1$  führt [Algorithmus 9.7](#) keine Berechnungen durch und liefert  $\widehat{A} = A$  zurück. Also ist  $\Delta A = 0$  und (9.2) erfüllt.

Wir machen jetzt den Induktionsschritt von  $n$  auf  $n + 1$ . Dazu benennen wir die Einträge der Matrix  $A \in \mathbb{F}^{(n+1) \times (n+1)}$  wie folgt:

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

mit  $v, w \in \mathbb{F}^n$  und  $\alpha \in \mathbb{F}$ . Der [Algorithmus 9.7](#) führt in der ersten Iteration ( $k = 1$ ) folgende Schritte aus: Zunächst wird (siehe [Zeile 3](#) in [Algorithmus 9.7](#))

$$\widehat{z} = v \oslash \alpha$$

bestimmt und in den Einträgen  $(2, 1)$  bis  $(n + 1, 1)$  der Matrix gespeichert (als erste Spalte von  $L$ ). Gegenüber dem exakten Ergebnis  $z = v/\alpha$  erhalten wir nach [Satz 6.1](#) einen Fehlervektor

$$f := \widehat{z} - z \quad \text{mit} \quad |f| \leq \frac{|v|}{|\alpha|} \varepsilon_{\text{mach}}.$$

Anschließend werden die Linearkombinationen der Zeilen  $2, \dots, n + 1$  mit der ersten Zeile bestimmt ([Zeile 4](#)). In Maschinenoperationen bedeutet das die Berechnung der Untermatrix

$$(\widehat{A}_1)_{ij} := b_{ij} \ominus (\widehat{z}_i \odot w_j), \quad i, j = 2, \dots, n + 1.$$

Gemäß [Satz 6.1](#) gilt

$$(\widehat{A}_1)_{ij} = [b_{ij} - \widehat{z}_i w_j (1 + \varepsilon_1)] (1 + \varepsilon_2)$$

mit relativen Fehlern  $|\varepsilon_1| \leq \varepsilon_{\text{mach}}$  und  $|\varepsilon_2| \leq \varepsilon_{\text{mach}}$ , die natürlich vom Index  $(i, j)$  abhängen. Damit erhalten wir für den Fehler  $\Delta A_1 := \widehat{A}_1 - A_1$  gegenüber der exakten Matrix  $A_1 = B - \widehat{z} w^T$  die Abschätzung

$$|\Delta A_1| \leq (|B| + 2 |\widehat{z} w^T|) \varepsilon_{\text{mach}} \leq 2 (|B| + |\widehat{z} w^T|) \varepsilon_{\text{mach}}. \quad (9.3)$$

Am Ende der Iteration  $k = 1$  sieht die Matrix jetzt wie folgt aus:

$$\widehat{A}^{(1)} = \left[ \begin{array}{c|c} \alpha & w^T \\ \hline \widehat{z} & \widehat{A}_1 \end{array} \right]$$

Die Berechnung der weiteren Schritte in [Algorithmus 9.7](#) sind genau dieselben, als würden wir [Algorithmus 9.7](#) neu starten und auf die Untermatrix  $\widehat{A}_1 \in \mathbb{F}^{n \times n}$  anwenden. Aufgrund der Induktionsvoraussetzung gilt für die in Maschinenoperationen bestimmte LR-Zerlegung von  $\widehat{A}_1$  aber

$$\widehat{L}_1 \widehat{R}_1 = \widehat{A}_1 + \Delta \widehat{A}_1 \quad \text{mit} \quad |\Delta \widehat{A}_1| \leq 3(n-1)(|\widehat{A}_1| + |\widehat{L}_1| |\widehat{R}_1|) \varepsilon_{\text{mach}}. \quad (9.4)$$

Die in Maschinenoperationen berechnete Gesamtzerlegung von  $A$  hat also schließlich die Darstellung

$$\widehat{L}\widehat{R} = \begin{bmatrix} 1 & 0 \\ \widehat{z} & \widehat{L}_1 \end{bmatrix} \begin{bmatrix} \alpha & w^\top \\ 0 & \widehat{R}_1 \end{bmatrix} = \begin{bmatrix} \alpha & w^\top \\ \alpha \widehat{z} & \widehat{z} w^\top + \widehat{L}_1 \widehat{R}_1 \end{bmatrix}.$$

Für die Fehlermatrix  $\Delta A := \widehat{A} - A = \widehat{L}\widehat{R} - A$  gilt also

$$\Delta A = \begin{bmatrix} 0 & 0 \\ \alpha \widehat{z} - \alpha z & \widehat{z} w^\top + \widehat{L}_1 \widehat{R}_1 - B \end{bmatrix}$$

und wegen  $f = \widehat{z} - z$ ,  $\widehat{A}_1 = A_1 + \Delta A_1$  und  $A_1 = B - \widehat{z} w^\top$  weiter auch

$$\Delta A = \begin{bmatrix} 0 & 0 \\ \alpha f & \Delta A_1 + \Delta \widehat{A}_1 \end{bmatrix}. \quad (9.5)$$

Wir schätzen noch den in (9.4) vorkommenden Term  $|\widehat{A}_1|$  ab:

$$\begin{aligned} |\widehat{A}_1| &= |A_1 + \Delta A_1| \quad \text{nach Definition von } \delta A_1 \\ &\leq |A_1| + |\Delta A_1| \\ &\leq |B - \widehat{z} w^\top| + 2(|B| + |\widehat{z} w^\top|) \varepsilon_{\text{mach}} \quad \text{nach Definition von } A_1 \text{ und (9.3)} \\ &\leq (1 + 2 \varepsilon_{\text{mach}})(|B| + |\widehat{z} w^\top|) \end{aligned} \quad (9.6)$$

und erhalten damit

$$\begin{aligned} |\Delta A_1 + \Delta \widehat{A}_1| &\leq |\Delta A_1| + |\Delta \widehat{A}_1| \\ &\leq 2(|B| + |\widehat{z} w^\top|) \varepsilon_{\text{mach}} + 3(n-1)(|\widehat{A}_1| + |\widehat{L}_1| |\widehat{R}_1|) \varepsilon_{\text{mach}} \quad \text{nach (9.3) und (9.4)} \\ &\leq 2(|B| + |\widehat{z} w^\top|) \varepsilon_{\text{mach}} + 3(n-1)(|B| + |\widehat{z} w^\top| + |\widehat{L}_1| |\widehat{R}_1|) \varepsilon_{\text{mach}} \quad \text{nach (9.6)} \\ &\leq 3n(|B| + |\widehat{z} w^\top| + |\widehat{L}_1| |\widehat{R}_1|) \varepsilon_{\text{mach}}. \end{aligned}$$

Mit dieser Abschätzung und  $|\alpha f| \leq |v| \varepsilon_{\text{mach}}$  erhalten wir schließlich aus (9.5):

$$\begin{aligned} |\Delta A| &\leq \begin{bmatrix} 0 & 0 \\ |v| & 3n(|B| + |\widehat{z} w^\top| + |\widehat{L}_1| |\widehat{R}_1|) \end{bmatrix} \varepsilon_{\text{mach}} \\ &\leq 3n \left( \begin{bmatrix} |\alpha| & |w^\top| \\ |v| & |B| \end{bmatrix} + \begin{bmatrix} 1 & 0 \\ |\widehat{z}| & |\widehat{L}_1| \end{bmatrix} \begin{bmatrix} |\alpha| & |w^\top| \\ 0 & |\widehat{R}_1| \end{bmatrix} \right) \varepsilon_{\text{mach}}, \end{aligned}$$

was zu beweisen war. □

Was bedeuten die Resultate der [Sätze 9.4, 9.5](#) und [9.8](#) nun für die Lösung eines linearen Gleichungssystems  $Ax = b$ , bei dem zunächst die LR-Zerlegung von  $A$  (ohne Pivotisierung, [Algorithmus 9.7](#)) bestimmt wird und dann die Lösung  $x$  mittels Vorwärts- und Rückwärtssubstitution  $Ly = b$  ([Algorithmus 9.1](#)) und  $Rx = y$  ([Algorithmus 9.2](#))?

Wir können folgende Abschätzung des Rückwärtsfehlers zeigen:

**Satz 9.9** (Rückwärtsanalyse der Lösung eines linearen Gleichungssystems mittels LR-Zerlegung ohne Pivotisierung, vgl. Golub, van Loan, 1983, Theorem 3.3.2). Es gelten die Voraussetzungen von Satz 9.8 und  $n \varepsilon_{\text{mach}} \leq C < 1$ . Es werde die LR-Zerlegung von  $A \in \mathbb{F}^{n \times n}$  ohne Pivotisierung mittels Algorithmus 9.7 bestimmt und anschließend das lineare Gleichungssystem  $Ax = b$  mit  $b \in \mathbb{F}^n$  mit Hilfe der Vorwärts- und Rückwärtssubstitution (Algorithmen 9.1 und 9.2) gelöst. Die erhaltene Näherungslösung sei  $\hat{x}$ . Dann gilt: Es gibt eine Matrix  $E$ , sodass gilt:

$$(A + E) \hat{x} = b \quad \text{mit} \quad |E| \leq 3(n-1)(|A| + |\hat{L}| |\hat{R}|) \varepsilon_{\text{mach}} + 2 \gamma_n |\hat{L}| |\hat{R}|. \quad (9.7)$$

*Beweis.* Es seien  $\hat{L}$  und  $\hat{R}$  die mittels Algorithmus 9.7 (ohne Pivotisierung) in Maschinenoperationen bestimmten Faktoren. Wir wissen nach Folgerung 9.6, dass die numerische Lösung  $\hat{x}$  die Bedingungen

$$\begin{aligned} (\hat{L} + \Delta \hat{L}) \hat{y} &= b, \\ (\hat{R} + \Delta \hat{R}) \hat{x} &= \hat{y} \end{aligned}$$

erfüllt. Dabei gilt  $|\Delta \hat{L}| \leq \gamma_n |\hat{L}^{\text{GM}}|$  und  $|\Delta \hat{R}| \leq \gamma_n |\hat{R}^{\text{GM}}|$ . Einsetzen ergibt

$$(\hat{L} + \Delta \hat{L})(\hat{R} + \Delta \hat{R}) \hat{x} = (\hat{L} \hat{R} + \hat{L} \Delta \hat{R} + \Delta \hat{L} \hat{R} + \Delta \hat{L} \Delta \hat{R}) \hat{x} = b.$$

Gemäß Satz 9.8 gibt es eine Matrix  $\Delta A$ , sodass  $\hat{L} \hat{R} = A + \Delta A$  gilt und  $\Delta A$  die Abschätzung (9.2) erfüllt. Wir haben also  $(A + E) \hat{x} = b$  mit der Matrix

$$E = \Delta A + \hat{L} \Delta \hat{R} + \Delta \hat{L} \hat{R} + \Delta \hat{L} \Delta \hat{R}.$$

Wir schätzen die Größe von  $E$  komponentenweise ab mit

$$\begin{aligned} |E| &\leq |\Delta A| + |\Delta \hat{L}| |\hat{R}| + |\hat{L}| |\Delta \hat{R}| + |\Delta \hat{L}| |\Delta \hat{R}| \\ &\leq 3(n-1)(|A| + |\hat{L}| |\hat{R}|) \varepsilon_{\text{mach}} + \gamma_n |\hat{L}| |\hat{R}| + \gamma_n |\hat{L}| |\hat{R}|. \end{aligned}$$

Damit ist die Behauptung gezeigt. □

**Bemerkung 9.10** (Bedeutung von Satz 9.8 und Satz 9.9).

- (i) Der Satz 9.8 ist ein Resultat zur Rückwärtsanalyse der LR-Zerlegung ohne Pivotisierung mittels Algorithmus 9.7. Dabei ist die Matrix  $A$  die Eingabe des Verfahrens, und die Faktoren  $L$  und  $R$  sind die Ausgabe.

Um die Rückwärtsstabilität beurteilen zu können, schreiben wir die Abschätzung (9.2) in die Form

$$|\Delta A| \leq 3(n-1)(|A| + |\hat{L}| |\hat{R}|) \varepsilon_{\text{mach}} = 3(n-1) \left( 1 + \frac{|\hat{L}| |\hat{R}|}{|A|} \right) |A| \varepsilon_{\text{mach}}$$

um. Dabei steht  $\frac{|\hat{L}| |\hat{R}|}{|A|}$  für den komponentenweisen Quotienten aus den Einträgen von  $|\hat{L}| |\hat{R}|$  und  $|A|$ , und  $\mathbf{1}$  ist die Matrix bestehend aus lauter Einsen. Wir müssen hier voraussetzen, dass  $A$  keine Nulleinträge hat. Der Faktor  $B(x)$  in der Definition (7.23) der Rückwärtsstabilität entspricht hier also der Größe

$$3(n-1) \left( 1 + \frac{|\hat{L}| |\hat{R}|}{|A|} \right).$$

Die Größen  $\widehat{L}$  und  $\widehat{R}$  «dürfen» hier vorkommen, da sie ja aus der Eingabe  $A$  berechenbar sind.

Die Abschätzung liefert den Hinweis, dass der Fall kritisch ist, bei dem mindestens ein Eintrag von  $|\widehat{L}| |\widehat{R}|$  sehr groß ist gegenüber dem entsprechenden Eintrag von  $|A|$ .

- (ii) Der [Satz 9.9](#) betrifft die Rückwärtsanalyse der auf dieser Zerlegung basierenden Lösung eines linearen Gleichungssystems  $Ax = b$ . Dabei sind  $A$  und  $b$  die Eingaben, und die Lösung  $x$  ist die Ausgabe. Aus der Abschätzung (9.7) erhalten wir genau dieselbe Information über den kritischen Fall.

Wir betrachten jetzt eine typische Situation für diesen kritischen Fall, dass ein Eintrag von  $|\widehat{L}| |\widehat{R}|$  sehr groß ist gegenüber dem entsprechenden Eintrag von  $|A|$ .

**Beispiel 9.11** ( $|\widehat{L}| |\widehat{R}|$  groß gegenüber  $|A|$ , vgl. [Higham, 2002](#), Chapter 1.12.1). Wir betrachten folgende Matrix und ihre LR-Zerlegung ohne Pivotisierung in exakter Rechnung:

$$A = \begin{bmatrix} \varepsilon & -1 \\ 1 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} \varepsilon & -1 \\ 0 & 1 + \frac{1}{\varepsilon} \end{bmatrix}.$$

Die Matrix  $A$  ist gut konditioniert<sup>5</sup>.

Wodurch entstehen jetzt die Probleme bei der Fließkommarechnung? Nehmen wir dazu an, dass  $\varepsilon$  und  $\frac{1}{\varepsilon}$  im Fließkommagitter  $\mathbb{F}$  liegen. Dann liegt  $A \in \mathbb{F}^{2 \times 2}$ , und wir erhalten in Fließkommarechnung mit Hilfe von [Algorithmus 9.7](#) die folgende LR-Zerlegung von  $A$  mit Faktoren

$$\widehat{L} = L = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \quad \text{und} \quad \widehat{R} = \begin{bmatrix} \varepsilon & -1 \\ 0 & 1 \oplus \frac{1}{\varepsilon} \end{bmatrix}.$$

**Quizfrage:** Können Sie nachvollziehen, dass das tatsächlich die Faktoren sind, die wir in Fließkommarechnung erhalten?

Nun sei  $\varepsilon > 0$  so gewählt, dass zusätzlich zu den obigen Bedingungen  $1 \oplus \frac{1}{\varepsilon} = \frac{1}{\varepsilon}$  gilt. **Quizfrage:** Wann passiert das? Welche Auswirkungen hat nun dieser kleine relative Fehler in einem einzigen Eintrag der LR-Zerlegung bei der Lösung eines linearen Gleichungssystems  $Ax = b$  mit Hilfe der Vorwärts- und Rückwärtssubstitution ([Algorithmen 8.7, 9.1 und 9.2](#))?

Wir betrachten dazu den Fall einer ausgewählten rechten Seite und eines Parameters  $\varepsilon$ , für den sich während der Rechnung möglichst wenig Rundungsfehler ergeben, sodass wir den Fehlereinfluss gut verfolgen können. Jegliche Rundungsfehler werden in **rot** markiert.

- (i) Wir betrachten dazu jetzt die rechte Seite  $b = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \in \mathbb{F}^2$ . Die exakte Lösung ergibt sich aus

$$Ly = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \quad \Rightarrow \quad y = \begin{pmatrix} 1 \\ 1 - \frac{1}{\varepsilon} \end{pmatrix}$$

und weiter

$$Rx = \begin{bmatrix} \varepsilon & -1 \\ 0 & 1 + \frac{1}{\varepsilon} \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 - \frac{1}{\varepsilon} \end{pmatrix} \quad \Rightarrow \quad x = \frac{1}{\varepsilon + 1} \begin{pmatrix} 2 \\ \varepsilon - 1 \end{pmatrix}.$$

---

<sup>5</sup> $\kappa(A) \leq \left( \frac{3+\sqrt{5}}{3-\sqrt{5}} \right)^{1/2} \approx 2.62$  für alle  $\varepsilon \in [0, 1]$

Aus der LR-Zerlegung in Fließkommarechnung erhalten wir dagegen

$$\widehat{L}\widehat{y} = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{pmatrix} \widehat{y}_1 \\ \widehat{y}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow \widehat{y} = \begin{pmatrix} 1 \\ 1 \ominus (\frac{1}{\varepsilon} \odot 1) \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \ominus \frac{1}{\varepsilon} \end{pmatrix}.$$

Wir betrachten weiter den Fall eines Parameters  $\varepsilon$ , sodass  $1 \ominus \frac{1}{\varepsilon} = 1 - \frac{1}{\varepsilon}$  und  $\varepsilon \ominus 1 = \varepsilon - 1$  gilt.

**Quizfrage:** Wann tritt dieser Fall auf? Dann haben wir weiter

$$\widehat{R}\widehat{x} = \begin{bmatrix} \varepsilon & -1 \\ 0 & \frac{1}{\varepsilon} \end{bmatrix} \begin{pmatrix} \widehat{x}_1 \\ \widehat{x}_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 - \frac{1}{\varepsilon} \end{pmatrix}.$$

Daher gilt

$$\widehat{x}_2 = \text{rd} \left( \frac{1 - \frac{1}{\varepsilon}}{\frac{1}{\varepsilon}} \right) = \text{rd}(\varepsilon - 1) = \varepsilon - 1$$

und

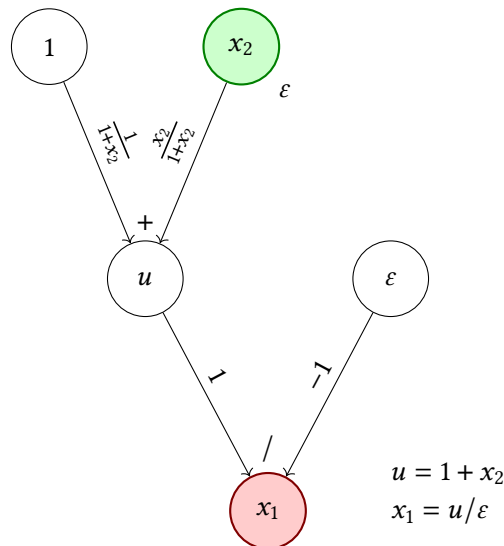
$$\widehat{x}_1 = \text{rd} \left( \frac{(1 \oplus (\widehat{x}_2 \odot 1))}{\varepsilon} \right) = \text{rd} \left( \frac{(1 \oplus (\varepsilon - 1))}{\varepsilon} \right) = 1,$$

also gilt

$$\widehat{x} = \begin{pmatrix} 1 \\ \varepsilon - 1 \end{pmatrix}.$$

Der relative Fehler von  $\widehat{x}$  gegenüber  $x$  wird dominiert durch die erste Komponente, wo er etwa  $1/2$  beträgt.

Wodurch kommt dieser große Fehler zustande? In unserem Beispiel entsteht der erste numerische Fehler erst bei der Rückwärtssubstitution (Lösen mit  $\widehat{R}$ ), und zwar durch den geringfügig falschen  $(2, 2)$ -Eintrag  $\frac{1}{\varepsilon}$  in  $\widehat{R}$ . Der dadurch verursachte relative Fehler in  $\widehat{x}_2$  ist mit  $\varepsilon$  gering. Im nächsten Schritt wird er jedoch bei der Berechnung von  $\widehat{x}_1$  wie folgt verstärkt:



Dabei ist die relative Konditionszahl mit  $\frac{x_2}{1+x_2} = \frac{\varepsilon-1}{2\varepsilon} \approx -\frac{1}{2\varepsilon}$  sehr groß (Auslöschung bei einer echten Subtraktion!), verstärkt also den relativen Fehler in  $\widehat{x}_2$  enorm zur beobachteten Größe von etwa  $1/2$ .

(ii) Wir betrachten jetzt noch die alternative rechte Seite  $b = \begin{pmatrix} \varepsilon - 1 \\ 2 \end{pmatrix} \in \mathbb{R}^2$ . Die exakte Lösung ergibt sich aus

$$Ly = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ 2 \end{pmatrix} \Rightarrow y = \begin{pmatrix} \varepsilon - 1 \\ 1 + \frac{1}{\varepsilon} \end{pmatrix}$$

und weiter

$$Rx = \begin{bmatrix} \varepsilon & -1 \\ 0 & 1 + \frac{1}{\varepsilon} \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ 1 + \frac{1}{\varepsilon} \end{pmatrix} \Rightarrow x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Aus der LR-Zerlegung in Fließkommarechnung erhalten wir dagegen

$$\widehat{L}\widehat{y} = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{pmatrix} \widehat{y}_1 \\ \widehat{y}_2 \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ 2 \end{pmatrix} \Rightarrow \widehat{y} = \begin{pmatrix} \varepsilon - 1 \\ 2 \ominus (1 - \frac{1}{\varepsilon}) \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ \text{rd}(1 + \frac{1}{\varepsilon}) \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ \frac{1}{\varepsilon} \end{pmatrix}.$$

Dann haben wir weiter

$$\widehat{R}\widehat{x} = \begin{bmatrix} \varepsilon & -1 \\ 0 & \frac{1}{\varepsilon} \end{bmatrix} \begin{pmatrix} \widehat{x}_1 \\ \widehat{x}_2 \end{pmatrix} = \begin{pmatrix} \varepsilon - 1 \\ \frac{1}{\varepsilon} \end{pmatrix} \Rightarrow \widehat{x} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

In diesem Fall ist die numerische Lösung – trotz der Rundungsfehler im Faktor  $\widehat{R}$  und des Rundungsfehlers in  $\widehat{y}$  – am Ende dennoch exakt.

Die Fehlerschranke (9.7) der Rückwärtsanalyse enthält den Term

$$|\widehat{L}| |\widehat{R}| = \begin{bmatrix} 1 & 0 \\ \frac{1}{\varepsilon} & 1 \end{bmatrix} \begin{bmatrix} \varepsilon & 1 \\ 0 & \frac{1}{\varepsilon} \end{bmatrix} = \begin{bmatrix} \varepsilon & 1 \\ 1 & \frac{2}{\varepsilon} \end{bmatrix},$$

wobei der **rot** markierte Eintrag sehr groß gegenüber dem entsprechenden Eintrag in

$$|A| = \begin{bmatrix} \varepsilon & 1 \\ 1 & 1 \end{bmatrix}$$

ist. Die relative Fehlerschranke des Rückwärtsfehlers ist also groß und damit auch die aus Satz 7.20 erhaltene relative Fehlerschranke des Vorwärtsfehlers  $\frac{\widehat{x} - x}{x}$ . Ob tatsächlich ein Fehler in der Lösung dieser Größenordnung auftritt oder nicht, hängt – wie das Beispiel zeigt – von der konkreten rechten Seite ab!

## § 9.2 FALL MIT SPALTENPIVOTSUCHE

Welchen Vorteil bringt nun die Spaltenpivotisierung? Wir geben zunächst das Verfahren für die LR-Zerlegung mit Spaltenpivotisierung an, damit klar ist, welchen Algorithmus wir untersuchen.

**Algorithmus 9.12** (Berechnung der LR-Zerlegung mit Spaltenpivotisierung, vgl. Golub, van Loan, 1983, Algorithm 3.4.1).

**Eingabe:** reguläre Matrix  $A$

**Ausgabe:** Faktoren  $L$  und  $R$  sowie Permutationsmatrix  $P$  der LR-Zerlegung  $PA = LR$

1: Setze  $p := (1, 2, \dots, n)^\top$

2: **for**  $k = 1, 2, \dots, n - 1$  **do**

3:     Setze  $I := k + 1 : n = \{k + 1, k + 2, \dots, n\}$

4: Bestimme einen Index  $s \in I \cup \{k\}$ , sodass gilt (vgl. (8.2))

$$|a_{sk}| = \max\{|a_{ik}| \mid i \in I \cup \{k\}\},$$

5: Tausche  $A(k, :)$  mit  $A(s, :)$  aus

6: Tausche  $p(k)$  mit  $p(s)$  aus

7: Setze  $A(I, k) := A(I, k)/A(k, k)$  // relevante Einträge der  $k$ -ten Spalte des Faktors  $L$

8: Setze  $A(I, I) := A(I, I) - A(I, k) A(k, I)$

9: **end for**

10: Extrahiere  $L$  aus dem strikteren unteren Dreieck von  $A$  und ergänze Einsen auf der Diagonale

11: Extrahiere  $R$  aus dem oberen Dreieck von  $A$

12: Setze

$$P := \begin{bmatrix} - & e_{p(1)}^T & - \\ & \vdots & \\ - & e_{p(n)}^T & - \end{bmatrix}$$

Der Vorteil der Spaltenpivotisierung ist gar nicht so einfach zu charakterisieren. Es gilt aber folgendes Analogon zu Satz 9.8:

**Satz 9.13** (Rückwärtsanalyse der LR-Zerlegung mit Spaltenpivotisierung, vgl. Golub, van Loan, 1983, Chapter 3.4.6). Es sei  $A \in \mathbb{F}^{n \times n}$  eine reguläre Matrix. Weiter seien  $\widehat{L}$  and  $\widehat{R}$  die von Algorithmus 9.12 in Maschinenoperationen berechnete Näherungslösung der LR-Zerlegung von  $PA = LR$ . Dabei nehmen wir an, dass in jeder Iteration das Pivotelement  $\widehat{A}(k, k) \neq 0$  ist. Dann gibt es eine Matrix  $\Delta A$ , sodass gilt:

$$\widehat{L} \widehat{R} = PA + \Delta A \quad \text{mit} \quad |\Delta A| \leq 3(n-1)(|PA| + |\widehat{L}| |\widehat{R}|) \epsilon_{\text{mach}}. \quad (9.8)$$

Alle Vorkommen von der Permutationsmatrix  $P$  sind der Übersichtlichkeit halber rot hervorgehoben.

**Beweis.** Es ist lediglich der Satz 9.8 auf  $PA$  an Stelle von  $A$  anzuwenden. Die zusätzlichen Schritte in Algorithmus 9.12 gegenüber Algorithmus 9.7 haben keinen Einfluss auf die Abschätzung der Rundungsfehler.  $\square$

**Quizfrage:** Warum wurde in Satz 9.13 wie bereits in Satz 9.8 angenommen, dass die Pivotelemente  $\widehat{A}(k, k) \neq 0$  sind? Folgt das nicht bereits aus der Annahme der Regularität von  $A$  und der Verwendung der Spaltenpivotsuche?

Mit diesem Resultat können wir nun folgendes Analogon zu Satz 9.9 zeigen:

**Satz 9.14** (Rückwärtsanalyse der Lösung eines linearen Gleichungssystems mittels LR-Zerlegung mit Spaltenpivotisierung, vgl. Golub, van Loan, 1983, Chapter 3.4.6). Es gelten die Voraussetzungen von Satz 9.13 und  $n \epsilon_{\text{mach}} \leq C < 1$ . Es werde die LR-Zerlegung von  $A \in \mathbb{F}^{n \times n}$  mit Spaltenpivotisierung mittels Algorithmus 9.12 bestimmt und anschließend das lineare Gleichungssystem  $Ax = b$  mit  $b \in \mathbb{F}^n$  mit Hilfe der Vorwärts- und Rückwärtssubstitution (Algorithmen 9.1 und 9.2) in der Form  $Ly = Pb$  und  $Rx = y$  gelöst. Die erhaltene Näherungslösung sei  $\widehat{x}$ . Dann gilt: Es gibt eine Matrix  $E$ , sodass gilt:

$$(A + E) \widehat{x} = b \quad \text{mit} \quad |E| \leq 3(n-1)(|A| + P^T |\widehat{L}| |\widehat{R}|) \epsilon_{\text{mach}} + 2 \gamma_n P^T |\widehat{L}| |\widehat{R}|. \quad (9.9)$$

Beweis.

□

In der Praxis kann man durch Spaltenpivotisierung oft die Größe des Terms  $\frac{P^T \widehat{L} \widehat{R}}{|A|}$  im Vergleich zum gleichen Term ohne  $P$  günstig beeinflussen. Damit reduziert sich dann die relative Fehlerschranke des Rückwärtsfehlers (9.9) gegenüber (9.7) und damit (Satz 7.20) auch die relative Fehlerschranke des Vorwärtsfehlers  $\frac{\widehat{x}_i - x_i}{x_i}$ . **Beachte:** Für  $P \neq \text{Id}$  sehen natürlich auch die Faktoren  $\widehat{L}$  und  $\widehat{R}$  anders aus, daher ist eine Vorhersage der Auswirkungen der Spaltenpivotisierung i. A. schwierig. Es gibt Beispiele, bei denen eine Spaltenpivotisierung tatsächlich keine Vorteile bringt. In der Praxis ist die Spaltenpivotisierung aber eine gute Wahl; siehe Übung.

## § 10 LR-ÄHNLICHE ZERLEGUNGEN FÜR SPEZIELLE MATRIZEN

### § 10.1 MATRIZEN, FÜR DIE KEINE PIVOTSUCHE ERFORDERLICH IST

Es gibt Klassen von Matrizen, für die wir auch ohne Spaltenpivotsuche eine brauchbare Fehlerschranke (9.7) erhalten. Dazu zählen

- (1) **total nichtnegative Matrizen**  $A$ . Das sind Matrizen, für die gilt, dass jede quadratische Untermatrix, die durch eine Anzahl von Zeilen und Spalten gewonnen wird, positive Determinante besitzt. Die LR-Faktorisierung dieser Matrizen erfüllt (ohne Pivotisierung) die Bedingungen  $L \geq 0$  und  $R \geq 0$ , woraus  $|L| |R| = |LR| = |A|$  folgt.
- (2) **diagonal-dominante Matrizen**  $A$ . Das sind Matrizen, für die der Diagonaleintrag in jeder Zeile betragsmäßig mindestens so groß ist wie die Summe der Beträge der restlichen Elemente derselben Zeile:

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad \text{für alle } i = 1, \dots, n.$$

Für reguläre, diagonal-dominante Matrizen  $A$  kann man zeigen (Higham, 2002, Theorem 9.9), dass  $\| |L| |R| \|_\infty \leq (2n - 1) \|A\|_\infty$  gilt.

- (3) Wenn die Matrix  $A$  zusätzlich zur Diagonaldominanz noch **tridiagonal** ist, also  $a_{ij} = 0$  gilt für alle Indizes mit  $|i - j| > 1$ , dann kann man sogar  $|L| |R| \leq 3 |A|$  zeigen (Higham, 2002, Theorem 9.13).

### § 10.2 CHOLSKY-ZERLEGUNG FÜR SYMMETRISCHE POSITIV DEFINITE MATRIZEN

Eine Matrix  $A \in \mathbb{R}^{n \times n}$  heißt **positiv definit**, wenn für alle  $x \neq 0$  gilt:  $x^T A x > 0$ . Die Cholesky-Zerlegung verallgemeinert den Begriff der Wurzel von positiven Zahlen auf symmetrische, positiv definite Matrizen. Man kürzt die Eigenschaft einer Matrix, «symmetrisch und positiv definit» zu sein, häufig mit **s. p. d.** ab.



**Definition 10.1** (Cholesky-Zerlegung). Es sei  $A \in \mathbb{R}^{n \times n}$  s. p. d.. Eine Zerlegung der Form  $A = LL^T$ , wobei  $L$  eine untere Dreiecksmatrix mit positiver Diagonale ist, heißt eine **Cholesky-Zerlegung** von  $A$ .

**Satz 10.2** (Cholesky-Zerlegung, vgl. Bornemann, 2018, Kapitel 8.3). Jede s. p. d. Matrix  $A \in \mathbb{R}^{n \times n}$  besitzt eine eindeutige Cholesky-Zerlegung.

*Beweis.* Wir führen einen konstruktiven Beweis per Induktion über die Dimension  $n$ . Für  $n = 1$  ist  $A \in \mathbb{R}$  eine positive Zahl und besitzt eine eindeutige positive Wurzel  $L$ . Für den Induktionsschritt sei die Behauptung bereits bewiesen für Dimension  $n$ . Wir betrachten den Ansatz

$$A' = \left[ \begin{array}{c|c} A & a \\ \hline a^T & \alpha \end{array} \right] = \left[ \begin{array}{c|c} L & 0 \\ \hline \ell^T & \lambda \end{array} \right] \left[ \begin{array}{c|c} L^T & \ell \\ \hline 0 & \lambda \end{array} \right] =: L'(L')^T,$$

wobei  $A' \in \mathbb{R}^{(n+1) \times (n+1)}$  s. p. d. ist. Die Dimensionen von  $A$  und  $L$  sind  $n \times n$ , die Dimensionen von  $a$  und  $\ell$  sind  $n \times 1$ , und  $\alpha$  und  $\lambda$  sind Zahlen. Am Ansatz lesen wir folgende Bedingungen ab:

$$\begin{aligned} A &= LL^T \\ a &= L\ell \\ \alpha &= \ell^T \ell + \lambda^2. \end{aligned}$$

**Quizfrage:** Was ist mit der vierten Gleichung?

Der Faktor  $L'$  soll eine untere Dreiecksmatrix mit positiver Diagonale sein. Insbesondere ist daher  $L$  ebenfalls eine untere Dreiecksmatrix mit positiver Diagonale. Aufgrund der Induktionsvoraussetzung ist  $L$  aber durch die erste Gleichung  $A = LL^T$  bereits eindeutig bestimmt. Da  $L$  auch invertierbar ist (**Quizfrage:** Klar?), können wir die zweite Gleichung auflösen und erhalten  $\ell = L^{-1}a$ . Schließlich ergibt sich aus der letzten Gleichung  $\lambda^2 = \alpha - \ell^T \ell$ . Wenn wir gleich noch zeigen können, dass dieser Ausdruck positiv ist, dann ist auch  $\lambda$  als die (positive) Wurzel eindeutig bestimmt. Wir wählen dazu  $x := -L^{-T}\ell$  und bilden mit dem Vektor  $x' := \begin{pmatrix} x \\ 1 \end{pmatrix}$  die quadratische Form:

$$\begin{aligned} \begin{pmatrix} x \\ 1 \end{pmatrix}^T A' \begin{pmatrix} x \\ 1 \end{pmatrix} &= \begin{pmatrix} x \\ 1 \end{pmatrix}^T \left[ \begin{array}{c|c} LL^T & L\ell \\ \hline \ell^T L^T & \alpha \end{array} \right] \begin{pmatrix} x \\ 1 \end{pmatrix} \\ &= \ell^T L^{-1} LL^T L^{-T} \ell - 2 \ell^T L^{-1} L \ell + \alpha \\ &= \alpha - \ell^T \ell. \end{aligned}$$

Da nach Voraussetzung  $A'$  positiv definit ist, ist die linke Seite positiv, also auch die rechte. Damit haben wir gezeigt, dass auch die s. p. d. Matrix  $A'$  eine Cholesky-Zerlegung besitzt, die auch eindeutig ist.  $\square$

Der Beweis von Satz 10.2 ist konstruktiv und kann direkt in ein numerisches Verfahren umgesetzt werden:

**Algorithmus 10.3** (Cholesky-Zerlegung einer s. p. d. Matrix, vgl. Golub, van Loan, 1983, Algorithm 3.4.1).

**Eingabe:** s. p. d. Matrix  $A$

**Ausgabe:** Faktor  $L$  der Cholesky-Zerlegung  $A = LL^T$

```

1: Setze  $L := \mathbf{0}$  (Nullmatrix)
2: for  $k = 1, 2, \dots, n$  do
3:   Setze  $I := 1 : k - 1 = \{1, 2, \dots, k - 1\}$ 
4:   Löse  $L(I, I) \ell = A(I, k)$ 
5:   Setze  $L(k, I) := \ell^T$ 
6:   Setze  $L(k, k) := \sqrt{A(k, k) - \ell^T \ell}$ 
7: end for

```

**Beachte:** In der ersten Iteration ( $k = 1$ ) ist  $I = \emptyset$ . Das bedeutet, dass in den Zeilen 4 und 5 nichts getan wird und  $\ell \in \mathbb{R}^{0 \times 1}$  als leerer Vektor zu deuten ist mit  $\ell^T \ell = 0$  in Zeile 6.

**Satz 10.4** (Aufwand der Cholesky-Zerlegung). Der numerische Aufwand für die Erzeugung der Cholesky-Zerlegung einer s. p. d. Matrix  $A \in \mathbb{R}^{n \times n}$  ist von der Größenordnung  $\frac{1}{3}n^3 + O(n^2)$ .

*Beweis.* Der wesentliche Aufwand in jeder Iteration ist die Lösung des Gleichungssystems in Zeile 4 mittels Vorwärtssubstitution (Algorithmus 9.1) mit dem Aufwand  $(k - 1)^2$ . In Zeile 6 kommen  $(k - 1)$  Multiplikationen,  $k - 1$  Additionen und eine Wurzel hinzu, zusammen  $2k - 1$  FLOPs. Die Summation über  $k = 1, 2, \dots, n$  ergibt  $\frac{1}{3}n^3 + O(n^2)$  FLOPs.  $\square$

**Bemerkung 10.5** (Zur Cholesky-Zerlegung).

- (i) Der Aufwand für die Berechnung der Cholesky-Zerlegung einer s. p. d. Matrix ist mit dem führenden Term  $\frac{1}{3}n^3$  nur halb so groß wie der Aufwand für die LR-Zerlegung einer allgemeinen regulären Matrix.
- (ii) Algorithmus 10.3 greift nur auf das obere Dreieck von  $A$  zu. Man kann das Verfahren also auch für Datenstrukturen implementieren, bei denen überhaupt nur die obere Dreiecksmatrix von  $A$  gespeichert wird, womit der Speicheraufwand i. W. halbiert wird.
- (iii) An Stelle der Zerlegung  $A = LL^T$  wird manchmal auch die Zerlegung  $A = LDL^T$  verwendet, wobei  $L$  eine untere Dreiecksmatrix mit lauter Einsen auf der Diagonale und  $D$  eine positive Diagonalmatrix ist.
- (iv) Es gibt eine alternative Definition der Cholesky-Zerlegung

$$A = R^T R,$$

bei der eine obere Dreiecksmatrix  $R$  mit positiver Diagonale verwendet wird.<sup>6</sup> Es gilt der Zusammenhang  $L = R^T$ . Natürlich gibt es hier auch wieder die Version  $A = R^T D R$ .

- (v) Die Existenz und Eindeutigkeit der Cholesky-Zerlegung (Satz 10.2) kann auch als Spezialfall der LR-Zerlegung hergeleitet werden. **Quizfrage:** Können Sie das?

<sup>6</sup>Das ist zum Beispiel der Standard in MATLABs `chol` und `scipy.linalg.cholesky`.

Wenn die Cholesky-Zerlegung  $A = LL^T$  vorliegt, so können wir jedes lineare Gleichungssystem  $Ax = b$  mit Hilfe der beiden Dreieckssysteme  $Ly = b$  und  $L^Tx = y$  mit dem Gesamtaufwand von  $n^2$  FLOPs lösen.

Ende der Woche 7

---

# Kapitel 5 Ausgleichsprobleme

## § 11 EINFÜHRUNG

Die (lineare) Ausgleichsrechnung befasst sich mit Lösungen überbestimmter linearer Gleichungssysteme  $Ax = b$  mit  $A \in \mathbb{R}^{m \times n}$  und  $b \in \mathbb{R}^m$ , wobei  $m \geq n$  gilt. Solche Aufgaben entstehen vor allem bei der Bestimmung von Parametern  $x$  in (linearen) Modellen aufgrund von Messdaten. Dabei heißt  $A$  die **Designmatrix** oder **Modellmatrix** (englisch: *design matrix, model matrix*) und  $b$  der **Vektor der Messwerte** (englisch: *vector of measurements*).

Da man i. A. nicht alle Gleichungen in  $Ax = b$  gleichzeitig erfüllen kann, betrachtet man Lösungen in einem ausgleichenden Sinne, die die **Fehlerquadratsumme** (englisch: *sum of squares*) minimieren, also die **Kleinste-Quadrate-Aufgabe** (englisch: *least-squares problem*)

$$\text{Minimiere } \frac{1}{2} \|Ax - b\|_2^2, \quad x \in \mathbb{R}^n \quad (11.1)$$

lösen. Man spricht auch von einem **Ausgleichsproblem**, **Regressionsanalyse**, **Parameterschätzung**, **Parameteridentifikation** oder **Modellkalibrierung** (englisch: *regression analysis, parameter estimation, parameter identification, model calibration*). Der Vektor

$$r := Ax - b \quad (11.2)$$

heißt der zu  $x$  gehörende **Residuenvektor**, kurz: das **Residuum** (englisch: *residual (vector)*).<sup>1</sup>

**Beispiel 11.1** (Bremsweg). Für den Bremsweg  $s$  eines Autos in Abhängigkeit der Geschwindigkeit  $v$  liegen folgende Messwerte vor:

Geschwindigkeit $v$ in $\text{km h}^{-1}$	10	20	25	30	40	50
Bremsweg $s$ in m	4	10	14	20	30	42

Der Zusammenhang wird durch eine quadratische Abhängigkeit der Form

$$s(v) = a v^2 + b v$$

mit noch unbekannten Modellparametern  $(a, b) \in \mathbb{R}^2$  modelliert. Diese Parameter sollen so bestimmt werden, dass das Modell und die Messdaten möglichst gut übereinstimmen (im Sinne der kleinsten

<sup>1</sup>Fälschlicherweise wird in der Literatur leider oft auch die Norm des Residuums als «Residuum» bezeichnet.

Fehlerquadratsumme). Die Designmatrix und der Vektor der Messwerte (ohne Einheiten) sind in diesem Beispiel gegeben durch

$$A = \begin{bmatrix} 10^2 & 10 \\ 20^2 & 20 \\ 25^2 & 25 \\ 30^2 & 30 \\ 40^2 & 40 \\ 50^2 & 50 \end{bmatrix}, \quad b = \begin{pmatrix} 4 \\ 10 \\ 14 \\ 20 \\ 30 \\ 42 \end{pmatrix}$$

mit Dimensionen  $n = 2$  and  $m = 6$ .

**Satz 11.2** (Lösung von (11.1)). Es sei  $A \in \mathbb{R}^{m \times n}$  und  $b \in \mathbb{R}^m$ .

- (i) Die Aufgabe (11.1) besitzt immer wenigstens eine (global optimale) Lösung.
- (ii) Die optimalen Lösungen von (11.1) sind genau die Lösungen der **Normalengleichung(en)** (auch: **Normalgleichung(en)**, englisch: **normal equations**)

$$A^T A x = A^T b. \quad (11.3)$$

*Beweis.* Wir benennen die Zielfunktion in (11.1) mit

$$f(x) := \frac{1}{2} \|Ax - b\|_2^2.$$

Wir verwenden die Singulärwertzerlegung  $A = U \Sigma V^T$ , die die Gestalt (4.4b) hat. Wie in (4.6) nehmen wir an, dass

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 = \sigma_{r+1} = \dots = \sigma_n$$

gilt, dass also  $r = \text{Rang}(A)$  ist für ein  $0 \leq r \leq n$ .

Zwei Punkte  $x, y \in \mathbb{R}^n$ , sodass  $x - y \in \ker A$  liegt, besitzen denselben Funktionswert  $f(x) = f(y)$ . (**Quizfrage:** Klar?) Wir vergeben uns also nichts, wenn wir uns bei der Suche nach Minimierern auf  $(\ker A)^\perp$  beschränken. Für solche  $x \in (\ker A)^\perp$  folgt aus der Singulärwertzerlegung, dass  $\|Ax\|_2 \geq \sigma_r \|x\|_2$  gilt. (**Quizfrage:** Klar?) Es gilt also für alle  $x \in (\ker A)^\perp$ :

$$f(x) = \frac{1}{2} \|Ax\|_2^2 - b^T A x + \frac{1}{2} \|b\|_2^2 \geq \sigma_r^2 \|x\|_2^2 - \|A^T b\|_2 \|x\|_2 + \frac{1}{2} \|b\|_2^2.$$

Daher gilt insbesondere, dass  $f(x) > f(0)$  wird für alle  $x \in (\ker A)^\perp$  mit  $\|x\|_2 > 2 \frac{\|A^T b\|_2}{\sigma_r^2}$ . Daher können wir uns bei der Suche nach Minimierern auf die Menge

$$K := (\ker A)^\perp \cap \left\{ x \in \mathbb{R}^n \mid \|x\|_2 \leq \frac{\|A^T b\|_2}{\sigma_r^2} \right\}$$

einschränken. Diese Menge ist aber kompakt, und die Funktion  $f$  ist stetig. Nach dem **Satz von Weierstraß** existiert also ein Minimierer von  $f$  über  $K$ , der gleichzeitig Minimierer von  $f$  über ganz  $\mathbb{R}^n$  ist. (**Quizfrage:** Klar?) Das zeigt **Aussage (i)**.

Aus der mehrdimensionalen Differenzialrechnung sollte bekannt sein, dass

$$\nabla f(x) = A^T(Ax - b) \stackrel{!}{=} 0,$$

also gerade die Normalengleichung (11.3), eine notwendige Bedingung dafür ist, dass  $x$  ein (lokaler) Minimierer von  $f$  ist. Wir zeigen jetzt noch, dass diese Bedingung sogar hinreichend ist, und zwar für globale Optimalität. Es sei dazu  $x \in \mathbb{R}^n$  ein Punkt, der die Normalengleichung (11.3) erfüllt, und  $y \in \mathbb{R}^n$  ein beliebiger anderer Punkt. Dann gilt

$$\begin{aligned} f(y) - f(x) &= \frac{1}{2} \|Ay - b\|_2^2 - \frac{1}{2} \|Ax - b\|_2^2 \\ &= \frac{1}{2} \|A(y - x) + Ax - b\|_2^2 - \frac{1}{2} \|Ax - b\|_2^2 \\ &= \frac{1}{2} \|A(y - x)\|_2^2 + (y - x)^T \underbrace{A^T(Ax - b)}_{=0} + \frac{1}{2} \|Ax - b\|_2^2 \\ &= \frac{1}{2} \|A(y - x)\|_2^2 \\ &\geq 0. \end{aligned}$$

Das zeigt, dass jedes  $x \in \mathbb{R}^n$ , das die Normalengleichung (11.3) erfüllt, tatsächlich global optimal ist. Damit ist auch Aussage (ii) gezeigt.  $\square$

**Folgerung 11.3** (Eindeutigkeit der Lösung der Kleinste-Quadrate-Aufgabe (11.1)). *Die folgenden Aussagen sind äquivalent:*

- (i) Die Lösung der Kleinste-Quadrate-Aufgabe (11.1) ist eindeutig.
- (ii) Die zu  $A$  gehörige Abbildung ist injektiv, also  $\ker A = \{0\}$ .
- (iii)  $A$  hat vollen Spaltenrang, also  $\text{Rang}(A) = n$ .

*Beweis.*

$\square$

Wir gehen im Folgenden, sofern nichts anderes gesagt wird, davon aus, dass die Bedingungen aus Folgerung 11.3 zutreffen, sodass die Normalengleichung (11.3) eine eindeutige Lösung hat. Die Matrix  $A^T A$  ist in diesem Fall positiv definit und nicht nur positiv semidefinit (und ohnehin symmetrisch). Man könnte also im Prinzip die Matrix  $A^T A$  formen, ihre Cholesky-Zerlegung bestimmen und dann die Normalengleichung durch die Lösung zweier Dreieckssysteme mit  $R$  bzw.  $R^T$  bestimmen.

Dieses Vorgehen hat aber gewisse Nachteile, die wir jetzt aufzeigen. Dazu sei  $A = U\Sigma V^T$  eine Singulärwertzerlegung von  $A$ . In Verallgemeinerung von (4.13) können wir

$$\kappa(A) := \frac{\sigma_1}{\sigma_{\min\{m,n\}}} \tag{11.4}$$

als die **Konditionszahl** einer nicht notwendig quadratischen Matrix  $A$  definieren.<sup>2</sup> Wegen  $m \geq n$  gilt hier also  $\kappa(A) = \frac{\sigma_1}{\sigma_n}$ . Die Matrix  $A^T A$  hat dann die Darstellung

$$A^T A = V \Sigma^2 V^T,$$

die gleichzeitig die Spektral- wie auch die Singulärwertzerlegung von  $A^T A$  ist<sup>GM</sup>. Die Konditionszahl von  $A^T A$  ist also gerade gleich  $(\kappa(A))^2$ !

**Quizfrage:** Wie groß ist der Aufwand für das Aufstellen von  $A^T A$ , für die Cholesky-Zerlegung und für die Lösung der Normalengleichung mit Hilfe der Cholesky-Zerlegung?

**Beispiel 11.4** (Auswirkungen der quadrierten Konditionszahl, vgl. das sogenannte Beispiel von Läuchli, Bornemann, 2018, S.70). In exakter Rechnung hat die Aufgabe mit

$$A = \begin{bmatrix} 1 & 1 \\ \varepsilon & 0 \\ 0 & \varepsilon \end{bmatrix}, \quad b = \begin{pmatrix} 2 \\ \varepsilon \\ \varepsilon \end{pmatrix}$$

und  $\varepsilon > 0$  die exakte Lösung  $x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ . Die Eigenwerte von

$$A^T A = \begin{bmatrix} 1 + \varepsilon^2 & 1 \\ 1 & 1 + \varepsilon^2 \end{bmatrix}$$

sind  $2 + \varepsilon^2$  und  $\varepsilon^2$ . Die Kondition von  $A$  erfüllt demnach  $\kappa(A) = \frac{\sqrt{2 + \varepsilon^2}}{\varepsilon} \in O(\varepsilon^{-1})$ , und  $\kappa(A^T A) = \frac{2}{\varepsilon^2} + 1 \in O(\varepsilon^{-2})$ .

Bei Rechnung in double precision etwa mit  $\varepsilon = \text{rd}(10^{-7})$  fallen bei der Berechnung von  $A^T A$  und von  $A^T b = \begin{pmatrix} 2 + \varepsilon^2 \\ 2 + \varepsilon^2 \end{pmatrix}$  Rundungsfehler der relativen Größenordnung  $\varepsilon_{\text{mach}}$  an. (**Quizfrage:** Warum?) Diese lassen selbst bei exakter Weiterrechnung einen relativen Vorwärtsfehler in der 2-Norm der Lösung von bis zu  $\kappa(A^T A) \varepsilon_{\text{mach}} \approx 1.4 \cdot 10^{14} \cdot 2.2 \cdot 10^{-16} \approx 4.4 \cdot 10^{-2}$  erwarten, siehe (3.18), (3.21). In der Tat ergibt sich die numerische Lösung mittels Cholesky-Zerlegung von  $A^T A$  und der anschließenden Vorwärts-/Rückwärtssubstitution

```
import numpy as np
import scipy.linalg
e = 1e-7;
A = np.array([[1, 1], [e, 0], [0, e]])
b = np.array([2, e, e])
x = scipy.linalg.solve(A.T @ A, A.T @ b, assume_a = 'pos') # select Cholesky solver
```

zu

$$\hat{x} \approx \begin{pmatrix} 1.011235955056179 \\ 0.988764044943822 \end{pmatrix}.$$

Die numerische Lösung hat also einen relativen Fehler von etwa  $1.1 \cdot 10^{-2}$ .

<sup>2</sup>In der Tat ist das die Konditionszahl von  $A$  wie wir sie aus (4.13) für invertierbare, quadratische Matrizen kennen, wobei aber der Zielraum  $\mathbb{R}^m$  auf den tatsächlichen Bildraum von  $A$  eingeschränkt wird.

## § 12 QR-ZERLEGUNG

Es stellt sich die Frage, ob man die Normalengleichung (11.3) nicht auch direkter lösen kann, insbesondere ohne  $A^T A$  aufzustellen. Eine Möglichkeit dazu bietet die QR-Zerlegung von  $A$ . Wir nehmen dafür wieder an, dass  $A \in \mathbb{R}^{m \times n}$  vollen Spaltenrang hat, insbesondere also  $m \geq n$  gilt.

Das Ziel ist eine Zerlegung von  $A$  in der folgenden Form.

**Definition 12.1** (QR-Zerlegung). *Es sei  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$  und vollem Spaltenrang. Eine Zerlegung der Form*

$$A = QR, \quad (12.1)$$

wobei  $Q \in \mathbb{R}^{m \times n}$  orthonormale Spaltenvektoren hat und  $R \in \mathbb{R}^{n \times n}$  eine obere Dreiecksmatrix ist, heißt eine (**reduzierte**) **QR-Zerlegung** (englisch: (reduced, economic, thin) QR decomposition) von  $A$ . Eine QR-Zerlegung heißt **normiert**, wenn die Diagonale von  $R$  positiv ist.

**Beachte:** Die Spalten von  $Q$  bilden unter den Voraussetzungen von Definition 12.1 eine Orthonormalbasis von Bild  $A$ .

**Satz 12.2** (Existenz und Eindeutigkeit der QR-Zerlegung, vgl. Bornemann, 2018, Satz 9.3). *Es sei  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$  und vollem Spaltenrang. Dann besitzt  $A$  eine eindeutige normierte QR-Zerlegung. Die Faktoren  $Q$  und  $R$  ergeben sich aus*

$$\begin{aligned} A^T A &= R^T R && \text{(Cholesky-Zerlegung von } A^T A), \\ R^T Q^T &= A^T && \text{(Lösung von } n \text{ linearen Gleichungssystemen mit Koeffizientenmatrix } R^T). \end{aligned}$$

*Beweis.* Wir beweisen zunächst die Eindeutigkeit. Dazu nehmen wir an, dass  $A = QR$  irgendeine normierte QR-Zerlegung von  $A$  ist. Dann folgt

$$A^T A = R^T Q^T Q R = R^T R.$$

Die Matrix  $A^T A$  ist s. p. d.. Nach Satz 10.2 (Existenz und Eindeutigkeit der Cholesky-Zerlegung) muss  $R$  der eindeutige Cholesky-Faktor von  $A^T A$  sein. Aus  $A = QR$  folgt durch Multiplikation mit  $R^{-1}$  notwendig

$$Q = AR^{-1},$$

also ist auch  $Q$  eindeutig festgelegt.

Kommen wir nun zur Existenz einer normierten QR-Zerlegung. Dazu definieren wir zunächst  $R$  als Cholesky-Faktor von  $A^T A = R^T R$ , was nach Satz 10.2 möglich ist. Per Definition besitzt  $R$  positive Diagonaleinträge. Weiter definieren wir  $Q := AR^{-1}$ . Wir müssen noch zeigen, dass die Spalten von  $Q$  orthonormal sind. Wir bilden dazu

$$Q^T Q = R^{-T} A^T A R^{-1} = R^{-T} R^T R R^{-1} = \text{Id}.$$

□



Wie oben bereits bemerkt, ist das Vorgehen wie im Beweis ungeeignet zur numerischen Bestimmung von  $R$ . Es kommt noch hinzu, dass die Orthonormalität der Spalten von  $Q$  bei der Definition gemäß  $Q := AR^{-1}$  nicht explizit eingefordert wird, sondern sich indirekt ergibt. Es ist also aufgrund von Rundungsfehlern damit zu rechnen, dass die Spalten von  $Q$  «nicht besonders» orthonormal sind.

Wir besprechen daher nun drei Alternativen zur Bestimmung der QR-Zerlegung.

## § 12.1 GRAM-SCHMIDT-VERFAHREN

Das **Gram-Schmidt-Verfahren** ist ein Verfahren, eine gegebene (endliche) Folge linear unabhängiger Vektoren  $a_1, \dots, a_n \in \mathbb{R}^m$  so in eine Folge von Vektoren  $q_1, \dots, q_n$  zu modifizieren, dass folgende Eigenschaften erfüllt sind:

- (1)  $\text{span}\{q_1, \dots, q_k\} = \text{span}\{a_1, \dots, a_k\}$  für alle  $k = 1, \dots, n$ .
- (2) Die Vektoren  $\text{span}^{\text{GM}}\{q_1, \dots, q_n\}$  sind paarweise orthonormal.

Zu diesem Zweck definieren wir  $q_1$  als normierte Version von  $a_1$ , definieren anschließend  $q_2$  als  $a_2$  ohne die Anteile von  $a_2$ , die bereits in dem Unterraum liegen, der von  $q_1$  aufgespannt wird, und normieren anschließend wieder usw. Man spricht auch davon, dass man den Vektor  $a_k$  «gegen die vorhergehenden Vektoren orthonormalisiert».

Dieses Verfahren wenden wir nun auf die Spaltenvektoren von  $A$  an. Das führt zu [Algorithmus 12.3](#).

**Algorithmus 12.3** (Gram-Schmidt-Verfahren zur Berechnung der QR-Zerlegung).

**Eingabe:** Matrix  $A \in \mathbb{R}^{m \times n}$  mit  $\text{Rang}(A) = n$

**Ausgabe:** Faktoren  $Q$  und  $R$  der normierten QR-Zerlegung  $A = QR$

```

1: for  $k = 1, 2, \dots, n$  do
2:   for  $\ell = 1, 2, \dots, k - 1$  do
3:     Berechne  $r_{\ell k} := a_k^T q_\ell$                                 // Bestimmung der Orthogonalisierungskoeffizienten
4:   end for
5:   Setze  $q_k := a_k$ 
6:   for  $\ell = 1, 2, \dots, k - 1$  do
7:     Setze  $q_k := q_k - r_{\ell k} q_\ell$                                 // Orthogonalisierung von  $a_k$  gegen  $q_\ell$ 
8:   end for
9:   Setze  $r_{kk} := \|q_k\|_2$ 
10:  Setze  $q_k := \frac{q_k}{r_{kk}}$                                 // Normalisierung von  $q_k$ 
11: end for
```

**Quizfrage:** In der Praxis kann [Algorithmus 12.3](#) so implementieren, dass kein zusätzlicher Speicher für  $Q$  benötigt wird, sondern dass  $A$  *in place* überschrieben wird. Was ist dafür zu ändern?

[Algorithmus 12.3](#) bestimmt zunächst die Koeffizienten  $r_{k\ell}$  zur Orthogonalisierung von  $a_k$  gegen alle früheren Vektoren  $q_\ell$ , bevor die Orthogonalisierung ausgeführt wird ([Zeile 7](#)). Stattdessen könnten wir

auch die Orthogonalisierung nach und nach ausführen und dabei die Berechnung des Koeffizienten  $r_{k\ell}$  hinauszögern, bis die Orthogonalisierung von  $a_k$  gegen ein bestimmtes  $q_\ell$  an der Reihe ist. In [Algorithmus 12.3](#) führt das dazu, dass wir [Zeilen 2 bis 8](#) ersetzen durch

```

2: Setze  $q_k := a_k$ 
3: for  $\ell = 1, 2, \dots, k-1$  do
4:   Berechne  $r_{\ell k} := q_k^\top q_\ell$  // Bestimmung der Orthogonalisierungskoeffizienten
5:   Setze  $q_k := q_k - r_{\ell k} q_\ell$  // sukzessive Orthogonalisierung von  $a_k$  gegen  $q_\ell$ 
6: end for

```

Das ist in exakter Rechnung äquivalent zu [Algorithmus 12.3](#), aber numerisch von Vorteil. Diese Variante heißt **modifiziertes Gram-Schmidt-Verfahren**. Wir können es, wie in [Algorithmus 12.4](#) gezeigt, so implementieren, dass die Anteile einer Spalte  $q_k$ , sobald diese fertig vorliegt, sofort aus allen späteren Vektoren  $a_\ell$  herausgerechnet wird.

**Algorithmus 12.4** (modifiziertes Gram-Schmidt-Verfahren zur Berechnung der QR-Zerlegung).

**Eingabe:** Matrix  $A \in \mathbb{R}^{m \times n}$  mit  $\text{Rang}(A) = n$

**Ausgabe:** Faktoren  $Q$  und  $R$  der normierten QR-Zerlegung  $A = QR$

```

1: Setze  $Q := A$ 
2: for  $k = 1, 2, \dots, n$  do
3:   Setze  $r_{kk} := \|q_k\|_2$ 
4:   Setze  $q_k := \frac{q_k}{r_{kk}}$  // Normalisierung von  $q_k$ 
5:   for  $\ell = k+1, k+2, \dots, n$  do
6:     Berechne  $r_{k\ell} := q_k^\top q_\ell$  // Bestimmung der Orthogonalisierungskoeffizienten
7:     Setze  $q_\ell := q_\ell - r_{k\ell} q_k$  // sukzessive Orthogonalisierung von  $q_\ell$  gegen  $q_k$ 
8:   end for
9: end for

```

**Quizfrage:** In welcher Reihenfolge werden in [Algorithmus 12.3](#) bzw. in [Algorithmus 12.4](#) die Einträge der Matrizen  $Q$  und  $R$  berechnet? **Quizfrage:** Woran könnte es liegen, dass die Variante [Algorithmus 12.4](#) numerisch vorteilhafter ist als [Algorithmus 12.3](#)?

**Quizfrage:** Was passiert, wenn die Matrix  $A$  nicht vollen Spaltenrang hat?

**Satz 12.5** (Aufwand des modifizierten Gram-Schmidt-Verfahrens). *Der numerische Aufwand für die Erzeugung der QR-Zerlegung einer Matrix  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$  durch [Algorithmus 12.4](#) ist von der Größenordnung  $2mn^2$ .*

*Beweis.*

□

## § 12.2 GIVENS-ROTATIONEN

Das (modifizierte) Gram-Schmidt-Verfahren in exakter Rechnung bricht ab, falls  $A$  nicht vollen Spaltenrang besitzt. In numerischer Rechnung kann es auch bis zum Ende durchlaufen, wenn der betreffende

Vektor  $q_k$  nicht exakt der Nullvektor ist. Das Verfahren produziert dann eine Matrix  $Q$ , die stark von der Orthogonalität abweicht.

Das in diesem Abschnitt entwickelte Verfahren hat diese Nachteile nicht. Es stellt keine Anforderung an den Spaltenrang von  $A$ . Es erstellt für eine Matrix  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$  außerdem eine **(volle) QR-Zerlegung** (englisch: *full QR decomposition*)

$$A = QR = [Q_1 \mid Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix}, \quad (12.2)$$

wobei  $Q \in \mathbb{R}^{m \times m}$  orthonormale Spaltenvektoren hat und  $R \in \mathbb{R}^{m \times n}$  eine obere Dreiecksmatrix ist. Die Teilmatrizen haben die Dimensionen  $Q_1 \in \mathbb{R}^{m \times n}$ ,  $Q_2 \in \mathbb{R}^{m \times (m-n)}$  und  $R_1 \in \mathbb{R}^{n \times n}$ .

Die Idee ist es,  $A$  durch Multiplikation mit einer Folge von Matrizen von links auf obere Dreiecksgestalt zu transformieren, indem in geeigneter Reihenfolge Nullen unterhalb der Diagonalen erzeugt werden. Im Unterschied zur LR-Zerlegung erfolgt dies aber nicht durch Frobeniusmatrizen, sondern mit Hilfe orthogonaler Rotationsmatrizen, sogenannter Givens-Rotationen.

**Definition 12.6** (Givens-Rotation). Eine quadratische Matrix der Form<sup>3</sup>

$$\Omega^{(k,\ell)} := \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & c & & -s \\ & & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \\ & & & s & & c & \\ & & & & & & \ddots \\ & & & & & & & 1 \\ & & & & & & & & 1 \end{bmatrix}$$

$\uparrow \quad \quad \uparrow$   
 $\ell \quad \quad k$

$\leftarrow \ell$   
 $\leftarrow k$

heißt **Givens-Rotation**, wenn  $c$  und  $s$  die Bedingung  $c^2 + s^2 = 1$  erfüllen. Es gibt daher einen eindeutigen Winkel  $\theta \in [0, 2\pi)$ , sodass  $c = \cos(\theta)$  und  $s = \sin(\theta)$  gilt.

Eine Givens-Rotation ist eine orthogonale Matrix mit Determinante 1, repräsentiert also tatsächlich eine Rotation. (**Quizfrage:** Beweis?)

**Frage:** Wie wirkt eine Givens-Rotation?

Es sei dazu  $x := \begin{pmatrix} a \\ b \end{pmatrix} \in \mathbb{R}^2 \setminus \{0\}$  und  $r := \|x\|_2 = \sqrt{a^2 + b^2}$ . Setzen wir  $c := \frac{a}{r}$  und  $s := -\frac{b}{r}$ , dann wirkt die Givens-Rotation

$$\Omega = \begin{bmatrix} c & -s \\ s & c \end{bmatrix}$$

<sup>3</sup>Nicht notierte Einträge sind Null.

wie folgt auf den Vektor  $x$ :

$$\Omega x = \begin{bmatrix} c & -s \\ s & c \end{bmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}.$$

**Quizfrage:** Was ist der zu dieser Transformation gehörende Winkel  $\theta$ ?

Diese Beobachtung können wir nutzen, um in geeigneter Reihenfolge Nullen unterhalb der Diagonale von  $A$  zu erzeugen.

**Beispiel 12.7** (Transformation einer Matrix auf obere Dreiecksgestalt durch Givens-Rotationen). Es sei  $A \in \mathbb{R}^{4 \times 2}$ . Wir führen die folgenden Givens-Rotationen aus. Dabei bezeichnet # ein im aktuellen Schritt geändertes Element. Die blau markierten Einträge sind diejenigen, die in der nächsten Givens-Rotation adressiert werden.

$$\begin{array}{c}
 \begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix} \xrightarrow{\Omega^{(3,4)}} \begin{bmatrix} * & * & * \\ * & * & * \\ \# & \# & \# \\ 0 & \# & \# \end{bmatrix} \xrightarrow{\Omega^{(2,3)}} \begin{bmatrix} * & * & * \\ \# & \# & * \\ 0 & \# & * \\ 0 & * & * \end{bmatrix} \xrightarrow{\Omega^{(1,2)}} \begin{bmatrix} \# & \# & \# \\ 0 & \# & \# \\ 0 & * & * \\ 0 & * & * \end{bmatrix} \\
 \\
 \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & * & * \\ 0 & * & * \end{bmatrix} \xrightarrow{\Omega'^{(3,4)}} \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & \# & \# \\ 0 & 0 & \# \end{bmatrix} \xrightarrow{\Omega'^{(2,3)}} \begin{bmatrix} * & * & * \\ 0 & \# & \# \\ 0 & 0 & \# \\ 0 & 0 & * \end{bmatrix} \\
 \\
 \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ 0 & 0 & * \end{bmatrix} \xrightarrow{\Omega''^{(3,4)}} \begin{bmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & \# \\ 0 & 0 & 0 \end{bmatrix}.
 \end{array}$$

Am Ende dieses Vorgangs haben wir die Darstellung

$$\Omega''^{(3,4)} \Omega'^{(2,3)} \Omega'^{(3,4)} \Omega^{(1,2)} \Omega^{(2,3)} \Omega^{(3,4)} A = R.$$

Die gesuchte orthogonale Transformationsmatrix ist also die Inverse:

$$Q = (\Omega^{(3,4)})^T (\Omega^{(2,3)})^T (\Omega^{(1,2)})^T (\Omega'^{(3,4)})^T (\Omega'^{(2,3)})^T (\Omega''^{(3,4)})^T.$$

**Quizfrage:** Warum erzeugt man die Nullen Spalte für Spalte von unten nach oben und nicht in einer anderen Reihenfolge? Zum Beispiel Spalte für Spalte von oben nach unten? Oder erst alle Nullen in der letzten Zeile, dann die Nullen in der Zeile darüber usw.? Warum bleiben einmal erzeugte Nullen erhalten?

**Bemerkung 12.8** (QR-Zerlegung mit Givens-Rotationen).

- (i) In der Praxis berechnet man  $Q$  möglichst nicht als Matrix, sondern speichert stattdessen die Faktoren  $c$  und  $s$  für die Teilmatrizen. Wenn man später das Matrix-Vektor-Produkt  $Qx$  oder  $Q^T y$  benötigt, dann kann es Schritt für Schritt berechnet werden.
- (ii) Durch Givens-Rotationen erhält man i. A. keine positive Diagonale in  $R$ , also keine normierte QR-Zerlegung.

- (iii) Givens-Rotationen sind besonders dann günstig, wenn die Matrix  $A$  bereits viele Null-Einträge besitzt.
- (iv) Das beschriebene Verfahren berechnet eine volle QR-Zerlegung (12.2), die natürlich die Faktoren der reduzierten QR-Zerlegung als Teilmatrizen enthält.

Wir betonen nochmals, dass die QR-Zerlegung mittels Givens-Rotationen auch dann möglich ist, wenn  $A$  nicht vollen Spaltenrang besitzt.

### § 12.3 HOUSEHOLDER-TRANSFORMATIONEN

Die Idee einer Householder-Transformation ist es, einen Vektor  $a \in \mathbb{R}^m$  an einer Hyperebene mit geeignetem Normalenvektor  $v$  zu reflektieren, sodass  $a$  auf  $\|a\|_2 e_1$  oder auf  $-\|a\|_2 e_1$  abgebildet wird, wobei  $e_1 = (1, 0, \dots, 0)^T \in \mathbb{R}^m$  der erste Einheitsvektor ist.

**Frage:** Wie sehen solche Reflexionsmatrizen aus?

Der Anteil eines beliebigen Vektors  $x$  in Richtung von  $v$  ist gegeben durch

$$\frac{v^T x}{\|v\|_2^2}.$$

Bei der Reflexion soll dieser Anteil das Vorzeichen wechseln. Also hat die gesuchte Reflexion die Darstellung

$$Qx = x - 2 \frac{v^T x}{\|v\|_2^2} v, \quad \text{d. h. } Q = \text{Id} - 2 \frac{v v^T}{\|v\|_2^2}. \quad (12.3)$$

**Beachte:** Die Länge von  $v$  ist dabei unerheblich.

**Definition 12.9** (Householder-Transformation). Eine Abbildung bzw. eine Matrix der Form (12.3) heißt **Householder-Transformation** oder **Householder-Reflexion** (englisch: Householder transformation, Householder reflection).

**Beachte:** Solche Matrizen sind orthogonal und symmetrisch, erfüllen also auch  $Q^2 = \text{Id}$ . **Quizfrage:** Klar?

Wie muss  $v$  gewählt werden, damit  $Q$  einen gegebenen Vektor  $a$  auf  $\|a\|_2 e_1$  abbildet? Aus (12.3) erhalten wir die Bedingung

$$Qa = a - 2 \frac{v^T a}{\|v\|_2^2} v \stackrel{!}{=} \|a\|_2 e_1.$$

Umstellen ergibt

$$2 \frac{v^T a}{\|v\|_2^2} v = a - \|a\|_2 e_1,$$

also ist  $v$  (irgendein) Vielfaches von  $a - \|a\|_2 e_1$ , und wir setzen

$$v := a - \|a\|_2 e_1. \quad (12.4)$$

Möchte man  $v$  stattdessen auf  $-\|a\|_2 e_1$  abbilden, so kann

$$v := a + \|a\|_2 e_1 \quad (12.5)$$

verwendet werden. Damit liegt die gesuchte Matrix  $Q$  in (12.3) fest. In der Praxis wählt man (12.4), wenn der erste Eintrag des Vektors  $a$  negativ ist, sonst (12.5). **Quizfrage:** Warum ist das günstig?

Im ersten Schritt der QR-Zerlegung mittels Householder-Reflexionen wendet man nun eine solche Matrix  $Q^{(1)}$  auf  $A$  an, wobei als Vektor  $a$  die erste Spalte von  $A$  verwendet wird. In diesem Schritt ergibt sich also

$$\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{bmatrix} \xrightarrow{Q^{(1)}} \begin{bmatrix} \# & \# & \# \\ 0 & \# & \# \\ 0 & \# & \# \\ 0 & \# & \# \end{bmatrix}.$$

Anschließend wird dasselbe Verfahren nochmal auf die **rot** markierte Untermatrix von  $Q^{(1)}A$  angewendet. Die dabei verwendete Matrix  $Q \in \mathbb{R}^{(m-1) \times (m-1)}$  wird derart zu  $Q^{(2)}$  ergänzt, dass die erste Zeile von  $Q^{(1)}A$  unverändert bleibt, also

$$Q^{(2)} = \left[ \begin{array}{c|c} 1 & 0 \\ \hline 0 & Q \end{array} \right].$$

Durch wiederholte Anwendung dieser Idee erhalten wir die Darstellung

$$Q^{(n)} \dots Q^{(1)} A = R.$$

Die gesuchte Transformationsmatrix  $Q$  ist also die Inverse:

$$Q = Q^{(1)} \dots Q^{(n)}.$$

**Quizfrage:** Warum werden hier die Faktoren nicht transponiert?

**Bemerkung 12.10** (QR-Zerlegung mit Householder-Transformationen).

- (i) In der Praxis berechnet man  $Q$  möglichst nicht als Matrix, sondern speichert stattdessen die Vektoren  $v$  aus den einzelnen Teilschritten. Wenn man später das Matrix-Vektor-Produkt  $Qx$  oder  $Q^T y$  benötigt, dann kann es Schritt für Schritt berechnet werden.
- (ii) Durch Householder-Transformationen erhält man i. A. keine positive Diagonale in  $R$ , also keine normierte QR-Zerlegung.
- (iii) Das beschriebene Verfahren berechnet eine volle QR-Zerlegung (12.2), die natürlich die Faktoren der reduzierten QR-Zerlegung als Teilmatrizen enthält.
- (iv) Householder-Transformationen werden bei der QR-Zerlegung durch MATLABS `qr`-Routine sowie in `numpy.linalg.qr` und `scipy.linalg.qr` verwendet.

Wir betonen auch hier, dass die QR-Zerlegung mittels Householder-Transformationen auch dann möglich ist, wenn  $A$  nicht vollen Spaltenrang besitzt.

## § 12.4 NUTZUNG DER QR-ZERLEGUNG

Wir nehmen jetzt wieder an, dass  $A$  vollen Spaltenrang  $n$  besitzt. Die Optimalitätsbedingung (Normalengleichung) der Kleinste-Quadrate-Aufgabe (11.1)

$$A^T A x = A^T b \quad (11.3)$$

lässt sich bei Vorliegen einer (reduzierten) QR-Zerlegung (12.1) von  $A$  einfach lösen, denn (11.3) ist äquivalent zu  $R^T Q^T Q R x = R^T Q^T b$ , also auch zu

$$R x = Q^T b. \quad (12.6)$$

**Beispiel 12.11** (Lösung von Beispiel 11.4 mittels QR-Zerlegung). Wir greifen nochmal das Beispiel 11.4 auf, lösen das Ausgleichsproblem dieses Mal aber mittels QR-Zerlegung (durch Householder-Transformationen):

```
import numpy as np
import scipy.linalg
e = 1e-7;
A = np.array([[1, 1], [e, 0], [0, e]])
b = np.array([2, e, e])
Q, R = np.linalg.qr(A, mode = 'reduced')
# Q, R = scipy.linalg.qr(A, mode = 'economic')
x = scipy.linalg.solve_triangular(R, Q.T @ b)
```

Das ergibt eine numerische Lösung von

$$\hat{x} \approx \begin{pmatrix} 1.0000000000000000 \\ 1.0000000000000004 \end{pmatrix}.$$

Die numerische Lösung hat also nur noch einen relativen Fehler von etwa  $4.4 \cdot 10^{-16}$ .

Wenn die rechte Seite  $b$  bereits zu dem Zeitpunkt bekannt ist, wenn man die QR-Zerlegung von  $A$  berechnet, so kann man sie an die Matrix  $A$  anhängen und die QR-Zerlegung von  $[A \ b]$  bestimmen lassen. Dann bekommt man im Ergebnis  $Q^T b$  gleich mitgeliefert und benötigt die Matrix  $Q$  später überhaupt nicht mehr. Durch Verzicht auf die Rückgabe von  $Q$  bei der QR-Zerlegung kann nochmal numerischer Aufwand gespart werden.

```
import numpy as np
import scipy.linalg
e = 1e-7;
A = np.array([[1, 1], [e, 0], [0, e]])
b = np.array([2, e, e])
R = np.linalg.qr(np.hstack((A, b[:,None])), mode = 'r') # return (full size) R only
QTb = R[:,2] # extract Q' * b
R = R[:,2:] # extract reduced R
x = scipy.linalg.solve_triangular(R, QTb)
```

Das ergibt eine numerische Lösung von

$$\hat{x} \approx \begin{pmatrix} 0.9999999999999996 \\ 1.0000000000000004 \end{pmatrix}$$

mit einem relativen Fehler von etwa  $6.3 \cdot 10^{-16}$ .

## § 13 NUTZUNG DER SINGULÄRWERTZERLEGUNG

In diesem Abschnitt betrachten wir noch, wie wir alternativ zur QR-Zerlegung auch die Singulärwertzerlegung<sup>4</sup>

$$A = U\Sigma V^T$$

der Matrix  $A \in \mathbb{R}^{m \times n}$  nutzen können, um die Kleinste-Quadrate-Aufgabe (11.1) zu lösen. Dafür ist es nicht einmal erforderlich, dass  $A$  vollen Spaltenrang besitzt. Sollte dies nicht der Fall sein, so ist die Lösung nicht eindeutig (Folgerung 11.3). In dem Fall können wir aber nach der eindeutigen Lösung mit der kleinsten Norm fragen. Anders ausgedrückt fragen wir nach derjenigen Lösung, die in  $(\ker A)^\perp$  liegt. **Quizfrage:** Warum ist die normminimale Lösung der Normalengleichung (11.3) gerade auch diejenige, die in  $(\ker A)^\perp$  liegt?

**Satz 13.1** (Lösung der Kleinste-Quadrate-Aufgabe durch Singulärwertzerlegung, vgl. Rannacher, 2017, Satz 4.13). *Es sei  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$ , und es gelte  $r = \text{Rang}(A)$ . Dann ist die eindeutig bestimmte Lösung der Kleinste-Quadrate-Aufgabe (11.1) mit minimaler Norm gegeben durch*

$$x = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i. \quad (13.1)$$

Das zugehörige Residuum hat die quadrierte Norm

$$\|r\|_2^2 = \|Ax - b\|_2^2 = \sum_{i=r+1}^m (u_i^T b)^2. \quad (13.2)$$

*Beweis.* Für jedes  $x \in \mathbb{R}^n$  gilt

$$\|Ax - b\|_2^2 = \|AVV^T x - b\|_2^2 = \|U^T AVV^T x - U^T b\|_2^2 = \|\Sigma V^T x - U^T b\|_2^2.$$

Wir setzen zur Abkürzung  $z := V^T x \in \mathbb{R}^n$ . Dann lautet der Ausdruck in der Norm ausgeschrieben:

$$\begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \\ \hline 0 & \cdots & 0 & \\ \vdots & & \vdots & \\ 0 & \cdots & 0 & \end{bmatrix} \begin{pmatrix} z_1 \\ \vdots \\ z_n \end{pmatrix} - \begin{bmatrix} - & u_1^T & - \\ & \vdots & \\ - & u_n^T & - \\ & u_{n+1}^T & - \\ & \vdots & \\ - & u_m^T & - \end{bmatrix} \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix}.$$

<sup>4</sup>deren Berechnung übrigens auch Householder-Transformationen nutzt



Durch die Wahl  $z_i = \frac{u_i^T b}{\sigma_i}$  für  $i = 1, \dots, r$  können wir die ersten  $r$  Komponenten zu Null machen, wodurch  $\|\Sigma V^T x - U^T b\|_2^2$  minimiert wird. Sollte  $r < n$  sein, so wählen wir  $z_i = 0$  für die verbleibenden Indizes  $i = r+1, \dots, n$ , um die Lösung mit minimaler Norm zu erhalten. (**Quizfrage:** Warum führt eine minimale Norm in  $z$  auch zu einer minimaler Norm in  $x = Vz$ ?) Die Aussage (13.2) folgt unmittelbar.  $\square$

```
import numpy as np
import scipy.linalg
e = 1e-7;
A = np.array([[1, 1], [e, 0], [0, e]])
b = np.array([2, e, e])
U, Sigma, VT = np.linalg.svd(A, full_matrices = False)
x = np.dot(VT.T, np.dot(np.diag(1/Sigma), np.dot(U.T, b)))
```

Das ergibt eine numerische Lösung von

$$\hat{x} \approx \begin{pmatrix} 0.9999999999999997 \\ 1.0000000000000002 \end{pmatrix}$$

mit einem relativen Fehler von etwa  $4.0 \cdot 10^{-16}$ .

Ende der Woche 8

# Kapitel 6 Interpolation und Approximation

## § 14 EINFÜHRUNG

In diesem Kapitel geht es darum, «beliebige» Funktionen  $f$  durch einfach zu handhabende Funktionen  $g$  zu approximieren (anzunähern). Dafür gibt es mindestens drei Motivationen:

- (1) Die Auswertung von  $g$  (zum Beispiel ein Polynom) kann numerisch wesentlich günstiger sein als die Auswertung der (komplizierten) Funktion  $f$ .
- (2) Die Funktion  $g$  kann durch i. d. R. wenige Parameter repräsentiert werden, also mit wenig Speicherbedarf dargestellt werden.
- (3) Von der Funktion  $f$  sind nur endlich viele Funktionswerte (etwa aus Messungen) bekannt, die etwa in einer Tabelle von Stützstellen  $x_i$  und zugehörigen Funktionswerten  $f(x_i)$  aufgelistet sind. Wir wollen für diese Funktion  $f$  aber gerne auch sinnvolle Werte zwischen den Stützstellen angeben können, sie also zu einer überall definierten Funktion  $g$  fortsetzen.

Es gibt verschiedene Möglichkeiten, eine Approximation  $g$  zu einer gegebenen Funktion  $f$  festzulegen. Eine Strategie ist es,  $g$  aus einer Menge von zur Verfügung stehenden Funktionen (meist ein endlich-dimensionaler Vektorraum) so auszuwählen, dass der Abstand (Approximationsfehler)  $\|f - g\|$  minimiert wird. Dabei ist  $\|\cdot\|$  eine geeignete Norm auf dem Funktionenraum, dem  $f$  angehört, und  $g$  wird in einem endlich-dimensionalen Unterraum gesucht. Solche Approximationsaufgaben behandeln wir in [Abschnitt 17](#).

Eine andere Strategie besteht darin,  $g$  in einer Menge zur Verfügung stehender Funktionen dadurch festzulegen, dass  $g$  mit  $f$  an gewissen Stützstellen übereinstimmt. Man spricht dann davon, dass die Funktion  $g$  die Funktion  $f$  an diesen Stützstellen interpoliert. Solche Aufgaben sind Gegenstand von [Abschnitt 15](#).

Wir behandeln im gesamten [Kapitel 6](#) nur Aufgaben für Funktionen, die auf  $\mathbb{R}$  definiert sind, also nur von einer einzigen Variablen  $x \in \mathbb{R}$  abhängen. Typische Funktionen  $g$ , die man für die Approximation

oder Interpolation verwendet, sind:

$$\begin{aligned} & a_0 + a_1x + \dots + a_nx^n && \text{(Polynome),} \\ & \frac{a_0 + a_1x + \dots + a_nx^n}{b_0 + b_1x + \dots + b_mx^m} && \text{(rationale Funktionen),} \\ & \frac{1}{2}a_0 + \sum_{k=1}^n (a_k \cos(kx) + b_k \sin(kx)) && \text{(trigonometrische Polynome),} \\ & \sum_{k=1}^n a_k \exp(b_kx) && \text{(Exponentialsummen).} \end{aligned}$$

Die **Parameter**, die die Funktionen beschreiben, sind jeweils **hervorgehoben**.

## § 15 POLYNOMINTERPOLATION

Wir behandeln in diesem Abschnitt die **Interpolation** (englisch: *interpolation*) mit Hilfe von Polynomen. Dazu bezeichnen wir mit

$$P_n := \left\{ p(x) = \sum_{k=0}^n a_k x^k \mid a_0, \dots, a_n \in \mathbb{R} \right\} \quad (15.1)$$

den Vektorraum der **Polynome** (englisch: *polynomials*) vom Höchstgrad  $n \in \mathbb{N}_0$ . Die Zahlen  $a_0, \dots, a_n$  heißen die **Koeffizienten** (englisch: *coefficients*) des Polynoms  $p$ .

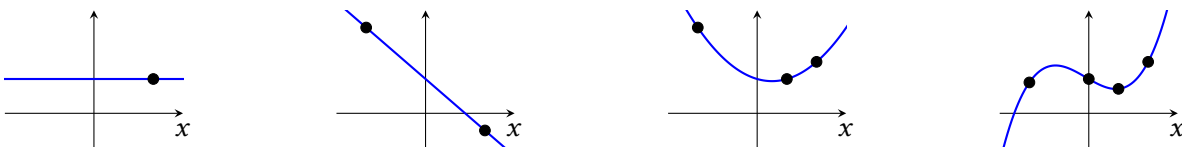
**Definition 15.1** (Interpolationsaufgabe). Zu gegebenen **Stützstellen** (englisch: *supporting points, sampling points, interpolation points*)  $x_0, \dots, x_n \in \mathbb{R}$  und zugehörigen Funktionswerten (**Stützwerte**)  $y_0, \dots, y_n \in \mathbb{R}$  lautet die **Interpolationsaufgabe** (englisch: *interpolation problem*) wie folgt:

$$\begin{aligned} & \text{Finde } p \in P_n, \text{ sodass die } \textbf{Interpolationsbedingungen} \text{ (englisch: } \textit{interpolation conditions}) \\ & p(x_i) = y_i \quad \text{für alle } i = 0, 1, \dots, n \text{ gelten.} \end{aligned} \quad (15.2)$$

**Beachte:** Wir beginnen mit der Nummerierung der Stützstellen beim Index 0, damit der Index  $n$  der letzten Stützstelle mit dem Maximalgrad des Polynoms übereinstimmt.

Für die Polynomgrade  $n = 0, 1, 2, 3$  spricht man speziell von **konstanter**, **linearer quadratischer** bzw. **kubischer Interpolation**.

Zu einer Aussage über die eindeutige Lösbarkeit der Interpolationsaufgabe kommen wir gleich. Zunächst eine Illustration verschiedener Interpolationsaufgaben mit Polynomgraden  $n = 0, 1, 2, 3$ .



**Quizfrage:** Erwarten Sie, dass die Lösung der Interpolationsaufgabe (15.2) immer ein Polynom vom Grad  $n$  ist? Oder wann ist das ggf. nicht der Fall?

## § 15.1 MONOMBASIS

Setzen wir die Stützstellen  $x_0, \dots, x_n$  nacheinander in die Darstellung (15.1) ein, so erhalten wir mit

$$\begin{bmatrix} 1 & x_0 & \cdots & x_0^n \\ 1 & x_1 & \cdots & x_1^n \\ \vdots & & & \vdots \\ 1 & x_n & \cdots & x_n^n \end{bmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix} \quad (15.3)$$

eine Bestimmungsgleichung für die unbekannten Koeffizienten  $a_0, \dots, a_n$ .

**Beispiel 15.2.** Für die Paare von Stützstellen und Stützwerten  $(-1, 5)$ ,  $(0, 3)$ ,  $(1, -7)$  und  $(2, 0)$  ergibt sich das System

$$\begin{bmatrix} 1 & -1 & 1 & -1 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \\ 1 & 2 & 4 & 8 \end{bmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{pmatrix} = \begin{pmatrix} 5 \\ 3 \\ -7 \\ 0 \end{pmatrix}$$

mit der (eindeutigen) Lösung  $\frac{1}{6}(18, -61, -24, 25)^\top$ . Das gesuchte Polynom ist also

$$p(x) = 3 - \frac{61}{6}x - 4x^2 + \frac{25}{6}x^3.$$

Die Matrix in (15.3) heißt die **Vandermonde-Matrix**  $V[x_0, \dots, x_n]$  zu den Stützstellen  $x_0, \dots, x_n$ . Sie hat folgende Eigenschaften:

**Satz 15.3** (Eigenschaften der Vandermonde-Matrix).

- (i)  $\det V[x_0, \dots, x_n] = \prod_{0 \leq i < j \leq n} (x_j - x_i)$ .
- (ii)  $V[x_0, \dots, x_n]$  ist genau dann regulär, wenn die Stützstellen paarweise verschieden sind.
- (iii) In diesem Fall erfüllt die Konditionszahl die Abschätzung

$$\kappa(V[x_0, \dots, x_n]) \geq \frac{2^{n-1}}{\sqrt{n+1}}. \quad (15.4)$$

Falls alle  $|x_j| \geq 1$  sind oder aber alle  $|x_j| \leq 1$ , dann gilt sogar

$$\kappa(V[x_0, \dots, x_n]) \geq 2^{n-1}. \quad (15.5)$$

*Beweis.* Zu **Aussage (i)**: Der Beweis verwendet Induktion über die Dimension der Matrix und folgt [Beutelspacher, 2014](#), Kapitel 7.5. Für  $n = 0$  ist die Indexmenge  $(i, j)$  in  $\prod_{0 \leq i < j \leq n}$  leer. Definitionsgemäß ist daher das Produkt gleich 1, was mit  $\det V[x_0] = \det(1) = 1$  übereinstimmt. Es sei nun die Behauptung für  $n - 1$  bereits gezeigt. Wir multiplizieren die vorletzte Spalte mit  $x_0$  und subtrahieren sie von der letzten. Dann multiplizieren wir die vorvorletzte Spalte mit  $x_0$  und subtrahieren sie von der vorletzten usw. In Matrixform geschrieben bilden wir also das Produkt

$$\begin{bmatrix} 1 & x_0 & \cdots & x_0^n \\ 1 & x_1 & \cdots & x_1^n \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^n \end{bmatrix} \begin{bmatrix} 1 & -x_0 & & \\ & 1 & -x_0 & \\ & & \ddots & \ddots \\ & & & 1 & -x_0 \\ & & & & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 1 & x_1 - x_0 & x_1(x_1 - x_0) & \cdots x_1^{n-1}(x_1 - x_0) \\ \vdots & \vdots & \vdots & \ddots \vdots \\ 1 & x_n - x_0 & x_n(x_n - x_0) & \cdots x_n^{n-1}(x_n - x_0) \end{bmatrix}.$$

Die Determinante ändert sich durch diese Manipulation nicht. (**Quizfrage:** Warum?) Die Determinante der erhaltenen Produktmatrix kann mit Hilfe des Entwicklungssatzes nach der ersten Zeile entwickelt werden. Wir erhalten

$$\det \begin{pmatrix} x_1 - x_0 & x_1(x_1 - x_0) & \cdots & x_1^{n-1}(x_1 - x_0) \\ \vdots & \vdots & \ddots & \vdots \\ x_n - x_0 & x_n(x_n - x_0) & \cdots & x_n^{n-1}(x_n - x_0) \end{pmatrix}.$$

Nun können wir aus der ersten Zeile den Faktor  $x_1 - x_0$  ausklammern, aus der zweiten den Faktor  $x_2 - x_0$  ausklammern usw. und erhalten

$$(x_1 - x_0)(x_2 - x_0) \cdots (x_n - x_0) \det \begin{pmatrix} 1 & x_1 & \cdots & x_1^{n-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \cdots & x_n^{n-1} \end{pmatrix}.$$

Die verbleibende Matrix ist dabei eine Vandermonde-Matrix der Dimension  $(n - 1) \times (n - 1)$ . Hier können wir die Induktionsannahme verwenden und erhalten schließlich

$$(x_1 - x_0)(x_2 - x_0) \cdots (x_n - x_0) \prod_{1 \leq i < j \leq n} (x_j - x_i) = \prod_{0 \leq i < j \leq n} (x_j - x_i),$$

was den Induktionsbeweis abschließt.

**Aussage (ii)** folgt sofort aus **Aussage (i)**, weil die Determinante genau dann ungleich Null ist, wenn die Stützstellen paarweise verschieden sind.

**Aussage (iii)** wurde in [Tyrtyshnikov, 1994](#), Theorem 4.1 bewiesen. □

**Folgerung 15.4** (Eindeutige Lösbarkeit von Interpolationsaufgabe (15.2)). *Die Interpolationsaufgabe (15.2) ist genau dann eindeutig lösbar, wenn die Stützstellen  $x_0, \dots, x_n$  paarweise verschieden sind. (Auf die Stützwerte  $y_0, \dots, y_n$  kommt es nicht an.)*

Neben der schlechten Kondition der Matrix in (15.3) ist auch noch der Aufwand  $\sim n^3$  als Nachteil unseres Ansatzes in diesem Abschnitt zu nennen. Dazu kommt, dass beim Hinzufügen einer weiteren

Stützstelle das Gleichungssystem (15.3) neu gelöst werden muss. Alle Nachteile kommen daher, dass offenbar die gewählte Darstellung

$$p(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

des gesuchten Polynoms nicht gut geeignet ist. Diese Darstellung kam daher, dass wir (ohne uns Gedanken darüber zu machen) als Basis des Raumes  $P_n$  die **Monome**  $\{1, x, x^2, \dots, x^n\}$  verwendet haben. Daher behandeln wird in den folgenden zwei Abschnitten zwar dieselbe Interpolationsaufgabe (15.2), jedoch mit anderen Basen für den Polynomraum  $P_n$ .

## § 15.2 LAGRANGE-POLYNOMBASIS

Zu gegebenen, paarweise verschiedenen Stützstellen  $x_0, \dots, x_n \in \mathbb{R}$  heißt das Polynom

$$L_i^{(n)}(x) := \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} \quad (15.6)$$

das  $i$ -te **Lagrange-Polynom** (englisch: *Lagrange polynomial*).

**Satz 15.5** (Eigenschaften der Lagrange-Polynome). *Es seien  $x_0, \dots, x_n$  gegeben. Dann gilt für die Lagrange-Polynome (15.6):*

(i)

$$L_i^{(n)}(x_j) = \begin{cases} 1, & \text{falls } i = j, \\ 0, & \text{falls } i \neq j. \end{cases} \quad (15.7)$$

(ii) Die Lagrange-Polynome  $L_0^{(n)}, \dots, L_n^{(n)}$  bilden eine Basis von  $P_n$ .

(iii) Die Summe aller Lagrange-Polynome ist konstant gleich 1:

$$\sum_{i=0}^n L_i^{(n)} \equiv 1. \quad (15.8)$$

*Beweis.*

□

Aufgrund der Eigenschaft (15.7) ist die Interpolationsaufgabe (15.2) bei einer Darstellung in der Lagrange-Basis besonders einfach zu lösen, denn es gilt

$$p(x) = \sum_{i=0}^n y_i L_i^{(n)}(x). \quad (15.9)$$

Es muss also kein Gleichungssystem gelöst werden. Wir haben jedoch auch Nachteile. Der Aufwand für die Auswertung des Interpolationspolynoms (15.9), wenn man die  $L_i^{(n)}$  in der Darstellung (15.6)

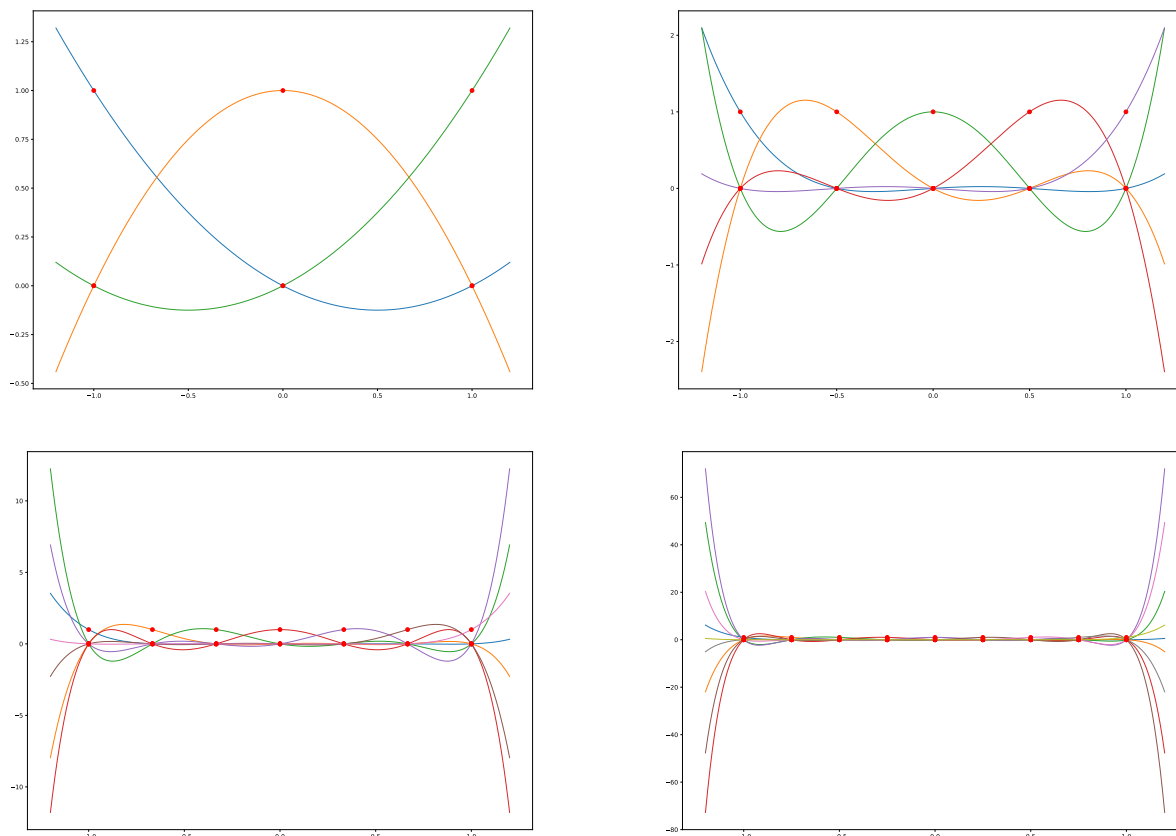


Abbildung 15.1: Lagrange-Polynome für verschiedene Anzahlen  $n + 1 \in \{3, 5, 7, 9\}$  von (äquidistant auf  $[-1, 1]$  gewählten) Stützstellen  $x_0, \dots, x_n$ .

belässt, ist von der Größenordnung  $n^2$  FLOPs. (**Quizfrage:** Klar?) Multipliziert man stattdessen das Polynom aus und wertet es mit dem Horner-Schema aus (Algorithmus 7.8), so sinkt der Aufwand auf die Größenordnung  $n$  FLOPs. Es ergibt sich dann allerdings der Nachteil, dass sich beim Hinzufügen einer neuen Stützstelle alle Polynome  $L_i^{(n)}$  ändern und das ausmultiplizierte Polynom neu berechnet werden muss.

Es gibt jedoch eine Möglichkeit, beide Probleme (Aufwand bei der Auswertung und bei der Aufdatierung) gleichzeitig zu adressieren. Dazu definieren wir das **Knotenpolynom**

$$\ell(x) := \prod_{j=0}^n (x - x_j) \quad (15.10)$$

und erhalten mit der Produktregel folgende Darstellung der Ableitung:

$$\ell'(x) = \sum_{k=0}^n \prod_{\substack{j=0 \\ j \neq k}}^n (x - x_j). \quad (15.11)$$

Insbesondere gilt

$$\ell'(x_i) = \sum_{k=0}^n \prod_{\substack{j=0 \\ j \neq k}}^n (x_i - x_j) = \prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j), \quad (15.12)$$

denn alle Summanden bis auf denjenigen für  $k = i$  sind gleich Null.

Wir erhalten also folgende alternative Darstellung des Lagrange-Polynoms  $L_i^{(n)}$ :

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} = \frac{\prod_{\substack{j=0 \\ j \neq i}}^n (x - x_j)}{\prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)} = \frac{\frac{\ell(x)}{x - x_i}}{\ell'(x_i)} = \frac{\ell(x)}{(x - x_i) \ell'(x_i)} \quad \text{wegen (15.12)} \quad (15.13)$$

für alle Punkte  $x \neq x_i$  (wo  $L_i^{(n)}$  aber bekanntermaßen gleich 1 ist).

Wir führen jetzt zur Abkürzung die **baryzentrischen Gewichte**

$$\omega_i := \frac{1}{\prod_{\substack{j=0 \\ j \neq i}}^n (x_i - x_j)} = \frac{1}{\ell'(x_i)}, \quad i = 0, \dots, n \quad (15.14)$$

ein. Gemäß (15.13) können wir dann

$$L_i^{(n)}(x) = \frac{\ell(x) \omega_i}{x - x_i} \quad (15.15)$$

schreiben, und da die Summe aller Lagrange-Polynome konstant gleich 1 ist (**Quizfrage:** Warum ist das so?), gilt weiter

$$L_i^{(n)}(x) = \frac{\ell(x) \omega_i}{x - x_i} = \frac{\ell(x) \frac{\omega_i}{x - x_i}}{\sum_{j=0}^n L_j^{(n)}(x)} = \frac{\frac{\omega_i}{x - x_i}}{\sum_{j=0}^n L_j^{(n)}(x) \frac{1}{\ell(x)}} = \frac{\frac{\omega_i}{x - x_i}}{\sum_{j=0}^n \frac{\omega_j}{x - x_j}}, \quad (15.16)$$

wobei die letzte Gleichheit aus (15.15) folgt. **Beachte:** Diese Darstellung ist verwendbar für alle Punkte  $x$ , die keine Stützstelle sind.

Damit können wir nun schließlich die Lösung der Interpolationsaufgabe (15.2) statt wie in (15.9) auch wie folgt darstellen:

$$p(x) = \frac{\sum_{i=0}^n y_i \frac{\omega_i}{x - x_i}}{\sum_{j=0}^n \frac{\omega_j}{x - x_j}}. \quad (15.17)$$

Diese Darstellung heißt die **baryzentrische Darstellung** des Lagrangeschen Interpolationspolynoms, und sie gilt für alle  $x \notin \{x_0, \dots, x_n\}$ . (Dort ist aber ja  $p(x_i) = y_i$  ohnehin bekannt.)



Da die Gewichte  $\omega_i$ , siehe (15.14), nur von den Stützstellen abhängen, können sie vorher berechnet werden. (Dafür benötigen wir Größenordnung  $n^2$  FLOPs.) Die Auswertung von (15.17) an einem Punkt  $x$  erfordert dann nur noch Größenordnung  $n$  FLOPs. (**Quizfrage:** Wieviele genau?) Zudem können die Gewichte bei Hinzukommen einer weiteren Stützstelle (und übrigens auch beim Entfernen oder Verschieben) mit Größenordnung  $n$  FLOPs aufdatiert werden.

### § 15.3 NEWTON-POLYNOMBASIS

Die Lösung der Interpolationsaufgabe (15.2) kann auch sukzessive aus der Lösung mit weniger Punkten aufgebaut werden.

**Definition 15.6** (Newtonsche Basispolynome). Zu gegebenen, paarweise verschiedenen Stützstellen  $x_0, \dots, x_n \in \mathbb{R}$  heißt das Polynom

$$N_i(x) := \prod_{j=0}^{i-1} (x - x_j) \quad (15.18)$$

das  $i$ -te **Newtonsche Basispolynom** (englisch: **Newton basis polynomial**).

**Beachte:** Insbesondere ist  $N_0(x) \equiv 1$  (leeres Produkt).

**Satz 15.7** (Eigenschaften der Newtonschen Basispolynome). Es seien  $x_0, \dots, x_n$  gegeben. Dann gilt für die Newtonschen Basispolynome (15.18):

$$(i) \quad N_i(x_j) = 0 \quad \text{für } 1 \leq j < i \leq n. \quad (15.19)$$

(ii) Die  $N_0, \dots, N_n$  bilden eine Basis von  $P_n$ .

Beweis. □

Wir schreiben also nun das gesuchte Interpolationspolynom zur Aufgabe (15.2) in der Form

$$p(x) = \sum_{i=0}^n c_i N_i(x). \quad (15.20)$$

Setzen wir die Stützstellen  $x_0, \dots, x_n$  nacheinander in die Darstellung (15.20) ein, so erhalten wir das lineare Gleichungssystem

$$\begin{bmatrix} N_0(x_0) & N_1(x_0) & \cdots & N_n(x_0) \\ N_0(x_1) & N_1(x_1) & \cdots & N_n(x_1) \\ \vdots & & & \vdots \\ N_0(x_n) & N_1(x_n) & \cdots & N_n(x_n) \end{bmatrix} \begin{pmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

Durch Ausnutzen der Eigenschaften von  $N_i(x_j)$  können wir dies äquivalent schreiben als

$$\begin{bmatrix} 1 & & & & \\ 1 & x_1 - x_0 & & & \\ 1 & x_2 - x_0 & (x_2 - x_0)(x_2 - x_1) & & \\ \vdots & \vdots & \vdots & \ddots & \\ 1 & x_n - x_0 & (x_n - x_0)(x_n - x_1) & \cdots & \prod_{j=0}^{n-1} (x_n - x_j) \end{bmatrix} \begin{pmatrix} c_0 \\ c_1 \\ c_2 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad (15.21)$$

wobei nicht bezeichnete Einträge in der Matrix für Nullen stehen. Dieses untere Dreieckssystem können wir sukzessive durch Vorwärtseinsetzen lösen. Die ersten Schritte dieses Vorgehens führen auf

$$\begin{aligned} c_0 &= y_0 \\ c_1 &= \frac{y_1 - c_0}{x_1 - x_0} = \frac{y_1 - y_0}{x_1 - x_0} \\ c_2 &= \frac{y_2 - c_0 - (x_2 - x_0)c_1}{(x_2 - x_0)(x_2 - x_1)} = \frac{y_2 - y_0 - (x_2 - x_0)\frac{y_1 - y_0}{x_1 - x_0}}{(x_2 - x_0)(x_2 - x_1)} = \frac{\frac{y_2 - y_1}{x_2 - x_1} - \frac{y_1 - y_0}{x_1 - x_0}}{x_2 - x_0} \end{aligned}$$

usw. Das motiviert die folgende Definition.

**Definition 15.8** (Dividierte Differenzen). Zu gegebenen paarweise verschiedenen Stützstellen  $x_0, \dots, x_n$  und zugehörigen Stützwerten  $y_0, \dots, y_n$  heißen die rekursiv definierten Größen

$$y[x_i] := y_i, \quad i = 0, \dots, n \quad (15.22a)$$

$$y[x_i, \dots, x_{i+k}] := \frac{\overbrace{y[x_{i+1}, \dots, x_{i+k}]}^{\text{Index startet später}} - \overbrace{y[x_i, \dots, x_{i+k-1}]}^{\text{Index endet früher}}}{x_{i+k} - x_i}, \quad i \geq 0, i+k \leq n, k \geq 0 \quad (15.22b)$$

die **(Newtonsche) dividierten Differenzen** (englisch: **(Newton) divided differences**).

Man kann die dividierten Differenzen in folgendem Schema anordnen:

$$\begin{array}{ccccccc} y_0 = y[x_0] & \xrightarrow{\text{red}} & y[x_0, x_1] & \xrightarrow{\text{blue}} & y[x_0, x_1, x_2] & \xrightarrow{\text{green}} & y[x_0, x_1, x_2, x_3] \\ & \nearrow & & \nearrow & & \nearrow & \\ y_1 = y[x_1] & \xrightarrow{\text{blue}} & y[x_1, x_2] & \xrightarrow{\text{green}} & y[x_1, x_2, x_3] & & \\ & \nearrow & & \nearrow & & & \\ y_2 = y[x_2] & \xrightarrow{\text{green}} & y[x_2, x_3] & & & & \\ & \nearrow & & & & & \\ y_3 = y[x_3] & & & & & & \end{array}$$

Die gesuchten Koeffizienten stehen in der ersten Zeile. Man berechnet die Größen im Schema üblicherweise Diagonale für Diagonale, die jeweils in einer anderen Farbe kodiert sind.<sup>1</sup> Jede Diagonale entspricht der Hinzunahme eines weiteren Stützpunkts, hier:  $x_1$ ,  $x_2$  und  $x_3$ .

<sup>1</sup>Eine Berechnung Spalte für Spalte ist ebenfalls möglich.

**Satz 15.9** (Lösung der Interpolationsaufgabe mit Hilfe dividierter Differenzen). *Die eindeutige Lösung der Interpolationsaufgabe (15.2) in Newtonscher Darstellung (15.20) ist gegeben durch*

$$p(x) = \sum_{i=0}^n y[x_0, \dots, x_i] N_i(x). \quad (15.23)$$

*Beweis.* Wir können beispielsweise durch Induktion zeigen, dass die gesuchten Koeffizienten  $c_i$  des Interpolationspolynoms in der Newtonschen Darstellung (15.20) die Beziehung  $c_i = y[x_0, \dots, x_i]$  erfüllen. Für einen vollständigen Beweis verweisen wir auf [Rannacher, 2017](#), Satz 2.2.  $\square$

Die Auswertung eines Polynoms in der Form (15.20), also

$$c_0 + c_1 (x - x_0) + c_2 (x - x_0)(x - x_1) + c_3 (x - x_0)(x - x_1)(x - x_2) + \dots, \quad (15.24)$$

erfolgt ähnlich zum Horner-Schema:

**Algorithmus 15.10** (Horner-Schema für (15.24)).

**Eingabe:** Koeffizienten  $c_0, \dots, c_n$  eines Polynoms der Form (15.24)

**Eingabe:** Auswertungsstelle  $x$

**Ausgabe:** Funktionswert  $y = p(x)$

```

1: Setze  $y := c_n$ 
2: for  $k = n - 1, n - 2, \dots, 0$  do
3:   Setze  $y := (x - x_k) y + c_k$ 
4: end for
```

Die Berechnung der dividierten Differenzen erfordert  $\sim n^2$  FLOPs. Das Hinzufügen einer neuen Stützstelle  $x_{n+1}$  zu  $x_0, \dots, x_n$  erfordert die Berechnung einer neuen Diagonalzeile mit dem Aufwand  $\sim n$  FLOPs. Die Auswertung mit [Algorithmus 15.10](#) geht wiederum in  $\sim n$  FLOPs.

## § 15.4 APPROXIMATIONSFEHLER UND KONVERGENZ BEI DER INTERPOLATION

**Frage:** Um wieviel kann das Interpolationspolynom von der Funktion, die es interpoliert, abweichen?

**Beachte:** Das ist eine neue Art von Fehler, die hier auftritt. Wir können ihn als **Approximationsfehler** oder **Verfahrensfehler** bezeichnen. Er entsteht dadurch, dass wir eine gegebene Funktion  $f$  durch eine andere Funktion (hier ein Polynom) approximieren.

Der folgende Satz liefert eine Abschätzung des Approximationsfehlers. Auf welche Art und Weise ([Abschnitte 15.1](#) bis [15.3](#)) das Interpolationspolynom dabei gewonnen wird, ist völlig unerheblich.

In Erweiterung von  $\langle a, b \rangle$ , siehe (2.9), definieren wir  $\langle a_0, a_1, \dots, a_n \rangle$  als das kleinste Intervall, das alle genannten Zahlen  $a_0, a_1, \dots, a_n \in \mathbb{R}$  enthält.<sup>2</sup>

<sup>2</sup>Man könnte auch sagen:  $\langle a_0, a_1, \dots, a_n \rangle$  ist die konvexe Hülle der Zahlen  $a_0, a_1, \dots, a_n$ .

**Satz 15.11** (Approximationsfehler vgl. [Rannacher, 2017](#), Satz 2.3). Es sei  $[a, b]$  ein Intervall und  $x_0, x_1, \dots, x_n$  paarweise verschiedene Stützstellen in  $[a, b]$ . Weiter sei<sup>3</sup>  $f \in C^{n+1}([a, b])$  und  $p$  das nach [Folgerung 15.4](#) eindeutig definierte Interpolationspolynom zu den Stützstellen  $x_0, \dots, x_n$ . Dann existiert für jedes  $x \in [a, b]$  ein  $\xi \in \langle x_0, x_1, \dots, x_n, x \rangle \subseteq [a, b]$ , sodass gilt:

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \ell(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{i=0}^n (x - x_i) \quad (15.25)$$

und folglich

$$\|f - p\|_{C([a,b])} := \sup_{x \in [a,b]} |f(x) - p(x)| \leq \sup_{x \in [a,b]} \frac{|f^{(n+1)}(x)|}{(n+1)!} \sup_{x \in [a,b]} |\ell(x)|. \quad (15.26)$$

Hierbei ist  $\ell(x) = \prod_{i=0}^n (x - x_i)$  wieder das Knotenpolynom aus (15.10). Die Norm  $\|\cdot\|_{C([a,b])}$  heißt die **Supremumsnorm** oder **Norm der gleichmäßigen Konvergenz** auf  $[a, b]$  (englisch: *supremum norm, norm of uniform convergence*). Statt dem Supremum dürften wir in (15.26) auch überall Maximum schreiben. (**Quizfrage:** Warum?)

*Beweis.* Für  $x = x_i$  ist nichts zu zeigen. (**Quizfrage:** Klar?) Es sei also  $x \in [a, b] \setminus \{x_0, \dots, x_n\}$  beliebig, aber fest gewählt. Wir setzen

$$c := \frac{f(x) - p(x)}{\ell(x)}.$$

(Der Nenner ist ungleich Null.) Dann besitzt die Funktion  $F(t) := f(t) - p(t) - c \ell(t)$  mindestens die  $n+2$  Nullstellen  $x_0, \dots, x_n$  und das gewählte  $x$  im Intervall  $\langle x_0, x_1, \dots, x_n, x \rangle$ . Nach dem [Satz von Rolle](#)<sup>4</sup> besitzt also  $F'$  in  $\langle x_0, x_1, \dots, x_n, x \rangle$  mindestens  $n+1$  Nullstellen. Man argumentiert so weiter, bis man die Aussage erhält, dass die  $(n+1)$ -te Ableitung  $F^{(n+1)}$  mindestens noch eine Nullstelle  $\xi \in \langle x_0, x_1, \dots, x_n, x \rangle$  besitzt. Es gilt also

$$\begin{aligned} 0 &= F^{(n+1)}(\xi) = f^{(n+1)}(\xi) - p^{(n+1)}(\xi) - c \ell^{(n+1)}(\xi) \\ &= f^{(n+1)}(\xi) - 0 - c (n+1)!, \end{aligned}$$

d. h.  $c = \frac{f^{(n+1)}(\xi)}{(n+1)!}$ . (**Quizfrage:** Warum ist  $p^{(n+1)}(\xi) = 0$ ?) (**Quizfrage:** Und warum ist  $\ell^{(n+1)}(\xi) = (n+1)!$ ?) Wir erhalten also

$$f(x) - p(x) = c \ell(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \ell(x),$$

was zu zeigen war. □

**Beachte:** Der Auswertungspunkt  $x$  kann durchaus außerhalb des Intervalls  $\langle x_0, x_1, \dots, x_n \rangle$  der Stützstellen liegen. Wenn man das Interpolationspolynom  $p(x)$  an einem solchen Punkt auswertet, spricht man auch von **Extrapolation** (englisch: *extrapolation*).

<sup>3</sup>Wir hatten die Räume  $C^k$  in [Definition 2.7](#) nur für offene Mengen definiert. Mit  $C^{n+1}([a, b])$  ist der Vektorraum der Funktionen gemeint, die in  $C^{n+1}((a, b))$  liegen und deren Ableitungen der Ordnungen  $0, 1, \dots, n+1$  sich stetig auf  $[a, b]$  fortsetzen lassen.

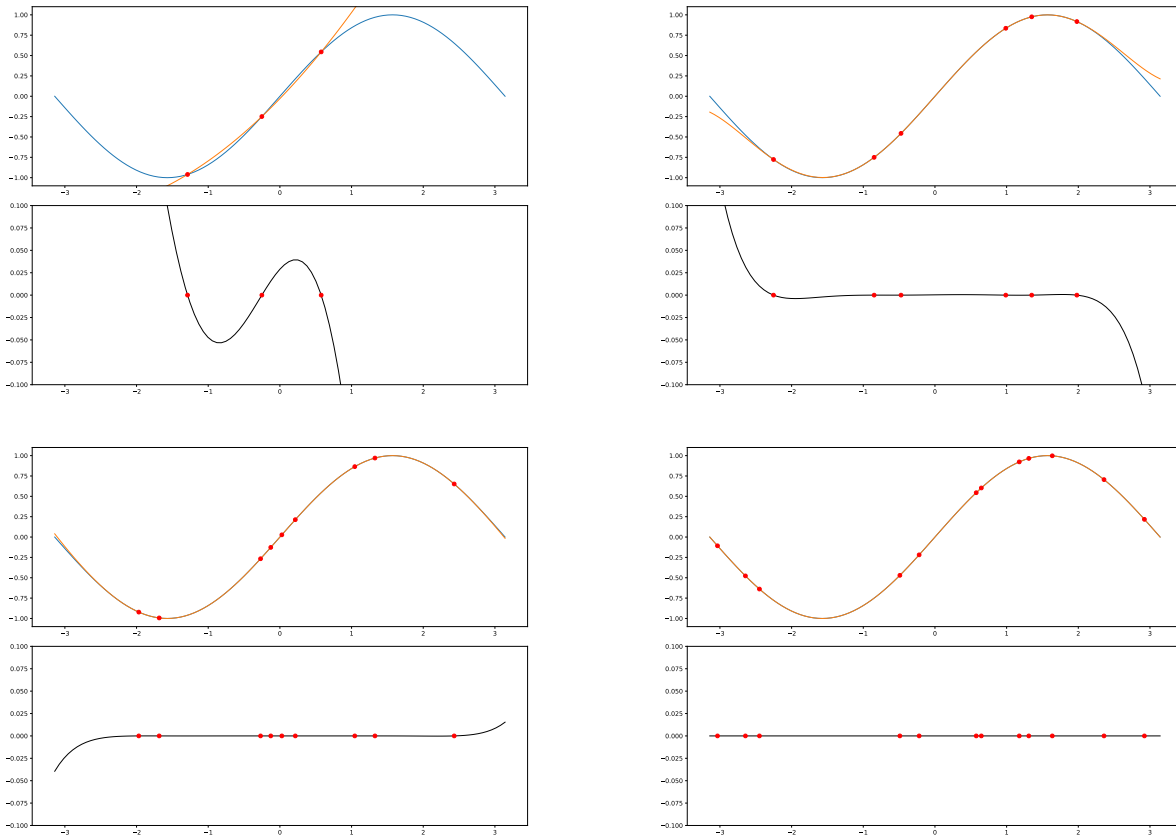
<sup>4</sup>Zwischen zwei Nullstellen einer differenzierbaren Funktion liegt (mindestens) eine Nullstelle ihrer Ableitung.

**Beispiel 15.12** (Fehlerabschätzung). Wir betrachten die Funktion  $f(x) = \sin(x)$  auf irgendeinem Intervall  $[a, b]$  mit irgendwelchen(!) Stützstellen  $x_0, x_1, \dots, x_n \in [a, b]$ . Dann folgt aus (15.26) die Fehlerabschätzung

$$\begin{aligned} \|f - p\|_{C([a,b])} &= \sup_{x \in [a,b]} |f(x) - p(x)| \leq \sup_{x \in [a,b]} \frac{|f^{(n+1)}(x)|}{(n+1)!} \sup_{x \in [a,b]} |\ell(x)| \\ &\leq \frac{1}{(n+1)!} (b-a)^n. \end{aligned}$$

**(Quizfrage:** Wieso können wir  $|f^{(n+1)}(x)| \leq 1$  abschätzen?) Für  $n \rightarrow \infty$  konvergiert die rechte Seite gegen Null. Mit anderen Worten konvergiert jede Folge von Interpolationspolynomen mit beliebigen Stützstellen auf jedem festen Intervall  $[a, b]$  gleichmäßig gegen die Funktion  $f(x) = \sin(x)$ .

Die folgenden Abbildungen illustrieren den Sachverhalt aus Beispiel 15.12 auf dem Intervall  $[a, b] = [-\pi, \pi]$  mit zufälligen Interpolationspunkten wachsender Anzahl  $n+1 \in \{3, 6, 9, 12\}$ . Gezeigt ist jeweils die Funktion  $f(x) = \sin(x)$  (blau), das Interpolationspolynom (orange) und darunter der Fehler  $f(x) - p(x)$  (schwarz).



Die Situation aus Beispiel 15.12 ist deshalb bemerkenswert, weil die Ableitungen  $f^{(n+1)}$  der Sinusfunktion alle auf dem betrachteten Intervall beschränkt sind. Für viele andere Funktionen wachsen die Ableitungen für  $n \rightarrow \infty$  dagegen stark an. In diesem Fall ist die gleichmäßige Konvergenz der

Interpolationspolynome auf kompakten Intervallen nicht garantiert, sondern hängt davon ab, wie die Stützstellen gewählt werden (siehe [Satz 15.14](#)).

Ein bekanntes Beispiel für die Nicht-Konvergenz ist das folgende.

**Beispiel 15.13** (Runge). Wir definieren die Funktion

$$f(x) := \frac{1}{1 + 25x^2}$$

und wählen die äquidistanten Stützstellen

$$x_i := \frac{2i}{n} - 1, \quad i = 0, \dots, n$$

auf  $[-1, 1]$ . Dann gilt für die Folge  $(p^{(n)})$  der durch diese Stützstellen festgelegten Interpolationspolynome

$$\lim_{n \rightarrow \infty} \|f - p^{(n)}\|_{C([-1,1])} = \infty.$$

Dabei sind die Werte der Fehlerfunktion  $f - p^{(n)}$  an den Rändern des Intervalls  $[-1, 1]$  maßgeblich, siehe [Abbildung 15.2](#). Diese Beobachtung ist als **Runge-Phänomen** bekannt.

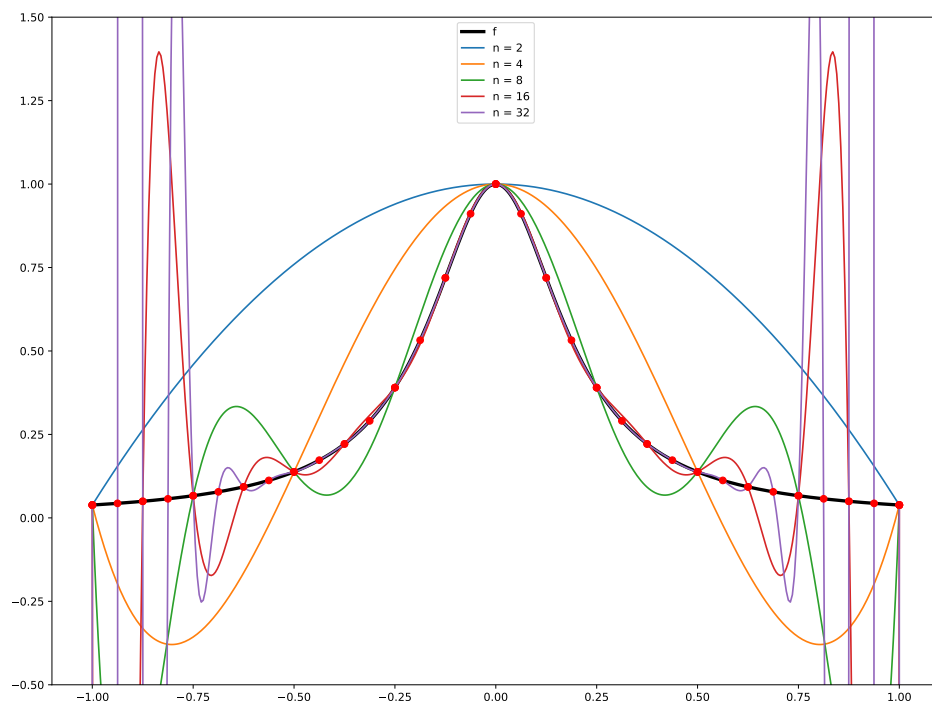


Abbildung 15.2: Illustration zu [Beispiel 15.13](#). Funktion  $f(x) = \frac{1}{1+25x^2}$  (schwarz) und ihre Interpolation durch Polynome mit verschiedenen Anzahlen äquidistanter Stützstellen auf  $[-1, 1]$ .

Das [Beispiel von Runge 15.13](#) ist kein Widerspruch zum berühmten [Approximationssatz von Weierstraß](#). Dieser besagt, dass jede stetige Funktion auf einem kompakten Intervall  $[a, b]$  dort gleichmäßig durch eine Folge von Polynomen approximiert werden kann. Diese Polynome können aber nicht als beliebige Folge von Interpolationspolynome gewählt werden. Stattdessen gilt:

**Satz 15.14** (Gleichmäßige Konvergenz der Interpolationspolynome bei geeigneter Wahl der Stützstellen). Es sei  $f \in C^0([a, b])$ . Dann existiert eine Folge  $(T^{(n)})$  von Stützstellen –  $T^{(n)}$  ist eine Teilmenge von  $[a, b]$  der Mächtigkeit  $n + 1$  –, sodass die Folge der zugehörigen Interpolationspolynome gleichmäßig auf  $[a, b]$  gegen  $f$  konvergiert.

*Beweis.* Wir skizzieren hier nur den Beweis. Wir betrachten die Aufgabe der Bestapproximation von  $f$  auf  $[a, b]$  durch ein Polynom gegebenen Höchstgrades  $n \in \mathbb{N}_0$ :

$$\text{Minimiere } \|f - p\|_{C([a,b])} := \sup_{x \in [a,b]} |f(x) - p(x)|, \quad p \in P_n. \quad (15.27)$$

Man kann zeigen, dass diese Aufgabe für jedes  $n \in \mathbb{N}_0$  eine Lösung besitzt, die sogar eindeutig ist. Zudem gilt für das optimale Polynom  $p^{(n)}$  der sogenannte **Alternantensatz 17.2** (englisch: *equioscillation theorem*), der besagt, dass es (mindestens)  $n + 2$  Extremstellen  $a \leq \xi_0 < \xi_1 < \dots < \xi_{n+1} \leq b$  von  $f - p^{(n)}$  gibt mit der Eigenschaft

$$f(\xi_i) - p^{(n)}(\xi_i) = s(-1)^i \|f - p^{(n)}\|_{C([a,b])}$$

mit  $s = 1$  oder  $s = -1$ . Das heißt, dass das für (15.27) optimale Polynom  $p^{(n)}$  die Eigenschaft hat, dass die Fehlerfunktion  $f - p^{(n)}$  abwechselnd und insgesamt mindestens  $n + 2$  Mal an der oberen bzw. unteren Grenze anschlägt. Nach dem **Zwischenwertsatz** für stetige Funktionen gibt es also zwischen zwei benachbarten  $\xi_i$  immer eine Zahl  $\xi_i < x_i < \xi_{i+1}$ , wo  $f(x_i) - p^{(n)}(x_i) = 0$  ist,  $i = 0, 1, \dots, n + 1$ . Diese  $x_i$  sind gerade die gesuchten Interpolationspunkte, die die Menge  $T^{(n)}$  bilden.

Dass die Folge  $\|f - p^{(n)}\|_{C([a,b])}$  gegen Null konvergiert – mit anderen Worten, dass die Folge  $(p^{(n)})$  gleichmäßig gegen  $f$  konvergiert, folgt aus dem Approximationssatz von Weierstraß. Dieser besagt, dass es eine Folge  $(g^{(n)})$  gibt mit  $g^{(n)} \in P_n$ , sodass  $\|f - g^{(n)}\|_{C([a,b])} \rightarrow 0$  konvergiert. Wegen der Bestapproximationseigenschaft der Lösung  $p^{(n)}$  von (15.27) gilt ja sogar

$$\|f - p^{(n)}\|_{C([a,b])} \leq \|f - g^{(n)}\|_{C([a,b])} \rightarrow 0$$

für  $n \rightarrow \infty$ . □

Der Nachteil an diesem Resultat ist natürlich, dass die «geeignete Folge» von Stützstellen von der Funktion  $f$  abhängt und man sie i. A. nicht kennt.

## § 15.5 KONDITION DER INTERPOLATIONSAUFGABE

Wir betrachten noch die Frage der Kondition der Interpolationsaufgabe (15.2), genauer: das Änderungsverhalten des Interpolationspolynoms bei festen (paarweise verschiedenen) Stützstellen  $x_0, \dots, x_n$ , aber veränderlichen Stützwerten  $y_0, \dots, y_n$ , die wir zum Vektor  $y \in \mathbb{R}^{n+1}$  zusammenfassen.

Wir betrachten dazu einmal das eindeutige Interpolationspolynom  $p(\cdot; y)$  zu den nominalen Stützwerten  $y$  und einmal das Interpolationspolynom  $p(\cdot; y + \Delta y)$  zu gestörten Stützwerten  $y + \Delta y$ . Dann gilt

wegen der linearen Abhängigkeit  $y \mapsto p(\cdot; y)$  nach (15.9):

$$\begin{aligned} p(x; y + \Delta y) - p(x; y) &= \sum_{i=0}^n (y_i + \Delta y_i) L_i^{(n)}(x) - \sum_{i=0}^n y_i L_i^{(n)}(x) \\ &= \sum_{i=0}^n \Delta y_i L_i^{(n)}(x) \end{aligned}$$

bzw. in relativer Betrachtungsweise

$$\frac{p(x; y + \Delta y) - p(x; y)}{p(x; y)} = \sum_{i=0}^n \frac{y_i L_i^{(n)}(x)}{p(x; y)} \frac{\Delta y_i}{y_i}.$$

Die punktweise relative partielle Konditionszahl an der Stelle  $x \in \mathbb{R}$

$$\frac{y_i L_i^{(n)}(x)}{p(x; y)} = \frac{y_i L_i^{(n)}(x)}{\sum_{j=0}^n y_j L_j^{(n)}(x)}$$

wird maßgeblich beeinflusst von dem Wert des  $i$ -ten Lagrange-Polynoms  $L_i^{(n)}(x)$ , der vor allem außerhalb des Intervalls  $\langle x_0, x_1, \dots, x_n \rangle$  der Stützstellen bereits für moderate Werte von  $n$  sehr groß sein kann, siehe [Abbildung 15.1](#).

Dazu noch ein Beispiel.

**Beispiel 15.15** (vgl. [Rannacher, 2017](#), Beispiel 2.2). Wir betrachten das Interpolationspolynom zu den  $n+1 = 2m+1$  äquidistanten Stützstellen auf  $[-1, 1]$  und Stützwerten  $y_i = 0$  mit Ausnahme des mittleren Wertes von  $y_m = \varepsilon$  (bei  $x_m = 0$ ). Das Interpolationspolynom hat die Darstellung

$$p(x) = \varepsilon L_m^{(n)}(x) = \varepsilon \prod_{\substack{j=0 \\ j \neq m}}^n \frac{x - x_j}{0 - x_j}.$$

Einige dieser Polynome für verschiedene Werte von  $n$  (und mit  $\varepsilon = 1$ ) werden in [Abbildung 15.3](#) gezeigt.

**Bemerkung 15.16** (Polynominterpolation). Die Interpolation mit Polynomen hat, wie wir bisher gesehen haben, gewisse Nachteile:

- (1) Die Auswirkung der Änderung von Stützwerten ist nicht lokal. Ändern wir also auch nur einen der Stützwerte  $y_i$ , so ändert sich das Interpolationspolynom an allen Punkten außer an den Stützstellen  $x_0, \dots, x_n$ . (**Quizfrage:** Wie sieht man das?)
- (2) Die Interpolationsaufgabe ist bereits für moderate Polynomgrade schlecht konditioniert.

Das ist der Grund, warum man bei einer großen Anzahl von Stützstellen bevorzugt mit stückweisen Polynomen interpoliert, die man glatt aneinanderstückelt. Solche Funktionen heißen **Splines**. Eine klassische Wahl sind **kubische Splines**, die man im Funktionenraum

$$\{s \in C^2([a, b]) \mid s|_{[x_i, x_{i+1}]} \in P_3, \ i = 0, 1, \dots, n-1\}$$

sucht.



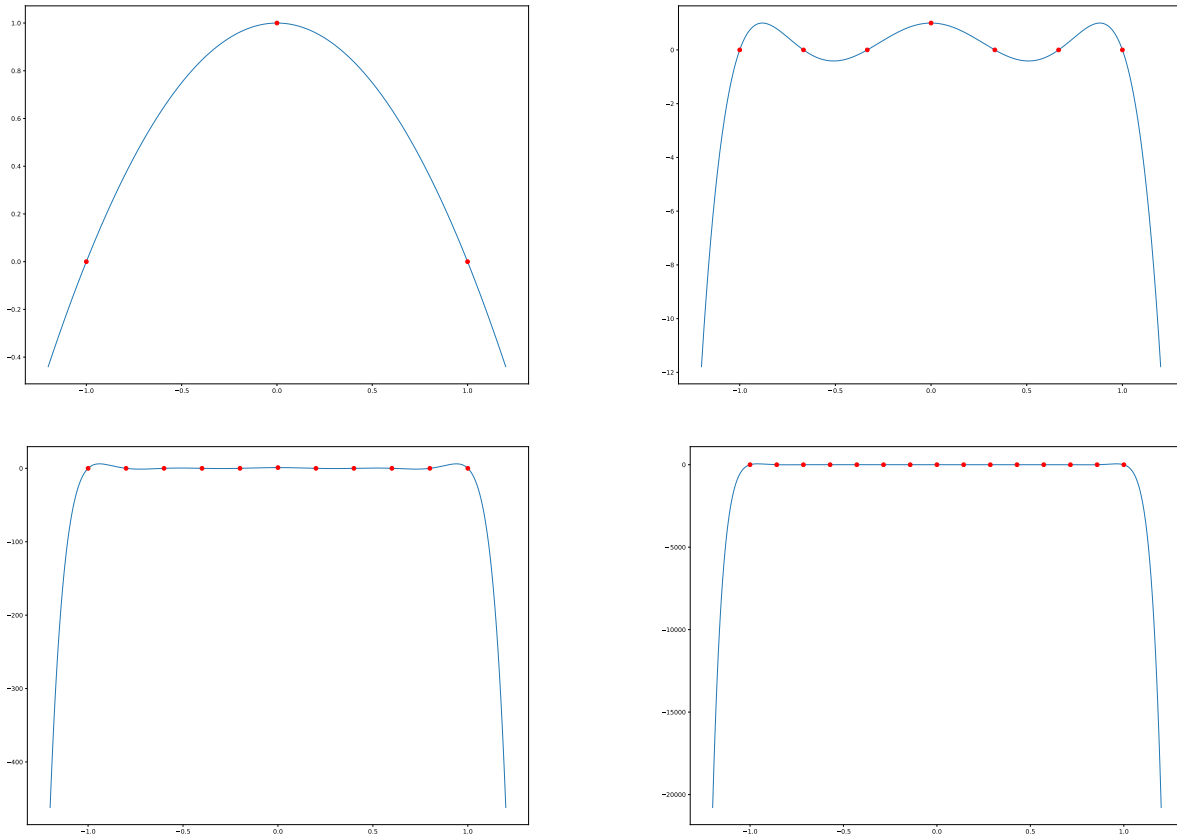


Abbildung 15.3: Illustration zu Beispiel 15.15 für verschiedene Anzahlen  $n + 1 \in \{3, 7, 11, 15\}$  von äquidistant auf  $[-1, 1]$  gewählten Stützstellen  $x_0, \dots, x_n$ .

**Bemerkung 15.17** (Interpolationsaufgaben mit Polynomen). *Wir haben uns hier nur mit Interpolationsaufgaben beschäftigt, bei denen die Funktionswerte  $y_0, \dots, y_n$  vorgegeben waren. Genauer spricht man dabei von **Lagrange-Interpolationsaufgaben** (englisch: **Lagrange interpolation**).<sup>5</sup> Allgemeiner kann man auch Aufgaben betrachten, bei denen man zusätzlich an einigen oder allen Stützstellen Werte der ersten und ggf. höherer Ableitungen des gesuchten Interpolationspolynoms vorgibt. Solche Aufgaben heißen **Hermite-Interpolationsaufgaben** (englisch: **Hermite interpolation**).*

## § 15.6 ANWENDUNGEN DER INTERPOLATION

Eine wichtige Anwendung der Polynominterpolation liegt in der Entwicklung von Formeln für die **numerische Differentiation** (englisch: *numerical differentiation*) durch **Finite-Differenzen-Approximation** (englisch: *finite difference approximation*) der Ableitung. Dabei geht es darum, die erste oder auch höhere Ableitungen einer Funktion  $f$  an einer Stelle  $x_0$  auszuwerten, wobei aber nur Funktionswerte von  $f$  verwendet werden sollen. Diese sehr häufig verwendete Technik ist zum

<sup>5</sup>Dabei ist nicht gemeint, dass man notwendigerweise die Darstellung (15.9) des Interpolationspolynoms unter Verwendung der Lagrange-Polynome benutzt, sondern eben nur, dass man nur Funktionswerte interpoliert, keine Ableitungen.

Beispiel dann hilfreich, wenn Alternativen wie Ableiten durch symbolisches oder algorithmisches Differenzieren nicht möglich ist.

Häufig verwendet man dabei äquidistante Stützstellen, die wir hier mit  $x_i = x_0 + i h$  bezeichnen. Dabei kann der Index  $i$  auch negative Werte annehmen, um Funktionsauswertungen «links» von  $x_0$  mit einzubeziehen. Die Idee ist, die Funktion  $f$  durch ihr Interpolationspolynom  $p \in P_n$  zu ersetzen und dann dessen Ableitung  $p'$  (stellvertretend für die Ableitung von  $f$ ) am Punkt  $x_0$  auszuwerten. Allgemeiner können wir statt des ersten Ableitung auch Ableitungen der Ordnung  $m$  verwenden, sofern  $1 \leq m \leq n$  gilt. (**Quizfrage:** Warum darf  $m$  nicht größer als  $n$  sein?)

Wir geben die wichtigsten Beispiele in [Tabelle 15.1](#) an.

Stützstellen	$n$	$m$	Formel	Fehler	Vor.	Name
$x_0, x_1$	1	1	$f'(x_0) \approx \frac{f(x_0+h)-f(x_0)}{h}$	$O(h)$	$f \in C^2$	Vorwärtsdifferenz
$x_{-1}, x_0$	1	1	$f'(x_0) \approx \frac{f(x_0)-f(x_0-h)}{h}$	$O(h)$	$f \in C^2$	Rückwärtsdifferenz
$x_{-1}, x_0, x_1$	2	1	$f'(x_0) \approx \frac{f(x_0+h)-f(x_0-h)}{2h}$	$O(h^2)$	$f \in C^3$	zentrale Diff.
$x_{-1}, x_0, x_1$	2	2	$f''(x_0) \approx \frac{f(x_0+h)-2f(x_0)+f(x_0-h)}{h^2}$	$O(h^2)$	$f \in C^4$	zentrale Diff. für 2. Abl.

Tabelle 15.1: Einige wichtige Finite-Differenzen-Formeln.

Die Angabe des Fehlers in [Tabelle 15.1](#) ist beispielsweise im Falle der Vorwärtsdifferenz so zu verstehen, dass

$$\frac{f(x_0+h)-f(x_0)}{h} \in f'(x_0) + O(h)$$

gilt, also es existiert eine Konstante  $C \geq 0$ , sodass

$$\left| \frac{f(x_0+h)-f(x_0)}{h} - f'(x_0) \right| \leq C h$$

für alle  $0 < h < \varepsilon$  gilt. Das kann man beispielsweise mit dem [Satz von Taylor 2.10](#) zeigen unter der Voraussetzung, dass  $f$  eine  $C^2$ -Funktion auf einer offenen Obermenge von  $\langle x_0, x_1 \rangle = [x_0, x_0+h]$  ist. **Quizfrage:** Können Sie eine Vermutung darüber anstellen, wie der Zusammenhang zwischen der benötigten Glattheitsordnung  $C^k$  der Funktion  $f$ , der Fehlerordnung und der Ordnung  $m$  der Ableitung ist?

Ende der Woche 9

## § 16 BERNSTEIN-POLYNOME UND BÉZIER-KURVEN

Bernstein-Polynome werden nicht zur Interpolation eingesetzt, sondern vorrangig zur Approximation und zur Beschreibung von Kurven und Flächen in der Computergrafik und Konstruktion (CAD). Mit ihnen kann man auch den Approximationssatz von Weierstraß konstruktiv beweisen. Bernstein-Polynome sind eine Familie reeller Polynome mit ganzzahligen Koeffizienten bzgl. der Monombasis.

**Definition 16.1** (Bernstein-Polynome). Es sei  $n \in \mathbb{N}_0$ . Die Polynome<sup>6</sup>

$$B_i^{(n)}(t) := \binom{n}{i} (1-t)^{n-i} t^i, \quad i = 0, \dots, n \quad (16.1)$$

heißen die **Bernstein-Polynome** (englisch: **Bernstein polynomials**) vom Grad  $n$  auf dem Intervall  $[0, 1]$ .

Die Bernstein-Polynome für verschiedene Grade sind in **Abbildung 16.1** gezeigt. (**Quizfrage:** Ist klar, dass die Bernstein-Polynome ganzzahlige Koeffizienten bzgl. der Monombasis haben?)

Man kann die Bernstein-Polynome auf beliebige Intervalle  $[a, b]$  transformieren. Sie haben dann die Gestalt

$$B_i^{(n)}(t) := \frac{1}{(b-a)^n} \binom{n}{i} (b-t)^{n-i} (t-a)^i$$

und heißen auch **verallgemeinerte Bernstein-Polynome** (englisch: *generalized Bernstein polynomials*). Da sich die Eigenschaften aber nicht wesentlich verändern, betrachten wir nur die Polynome (16.1).

**Satz 16.2** (Eigenschaften der Bernstein-Polynome). Die Bernstein-Polynome (16.1) haben folgende Eigenschaften:

(i) Alle Bernstein-Polynome  $B_i^{(n)}$ ,  $i = 0, \dots, n$ , sind Polynome vom Grad  $n$ .

(ii) Die Bernstein-Polynome haben auf  $[0, 1]$  Werte in  $[0, 1]$  und summieren sich zu eins:

$$0 \leq B_i^{(n)}(t) \leq 1 \quad \text{für alle } t \in [0, 1] \quad \text{und} \quad \sum_{i=0}^n B_i^{(n)}(t) \equiv 1 \quad \text{für alle } t \in \mathbb{R}. \quad (16.2)$$

Man spricht auch davon, dass die Funktionen  $B_0^{(n)}, \dots, B_n^{(n)}$  eine **Zerlegung der Eins(funktion)** (englisch: **partition of unity**) bilden.

(iii) Für  $n \geq 1$  gilt sogar

$$0 < B_i^{(n)}(t) < 1 \quad \text{für alle } t \in (0, 1). \quad (16.3)$$

(iv) Die Bernstein-Polynom erfüllen folgende Symmetriebedingung:

$$B_i^{(n)}(t) = B_{n-i}^{(n)}(1-t). \quad (16.4)$$

(v)  $B_i^{(n)}(t)$  hat  $t = 0$  als  $i$ -fache Nullstellen und  $t = 1$  als  $(n-i)$ -fache Nullstelle.

(vi) Für  $n \geq 1$  hat  $B_i^{(n)}(t)$  hat in  $[0, 1]$  genau eine Maximalstelle, und zwar bei  $i/n$ .

(vii) Die Polynome  $\{B_0^{(n)}, \dots, B_n^{(n)}\}$  bilden eine Basis von  $P_n$ .

<sup>6</sup>Dabei ist  $\binom{n}{i} = \frac{n!}{i!(n-i)!}$  der Binomialkoeffizient (« $n$  über  $i$ »).

(viii) Die Bernstein-Polynome haben eine rekursive Darstellung:

$$B_i^{(n)}(t) = \begin{cases} (1-t) B_0^{(n-1)}(t) & \text{im Fall } i = 0, \\ t B_{i-1}^{(n-1)}(t) + (1-t) B_i^{(n-1)}(t) & \text{im Fall } 0 < i < n, \\ t B_{n-1}^{(n-1)}(t) & \text{im Fall } i = n. \end{cases} \quad (16.5)$$

(ix) Die erste Ableitung der Bernstein-Polynome hat ebenfalls eine rekursive Darstellung:

$$\frac{d}{dt} B_i^{(n)}(t) = \begin{cases} -n B_0^{(n-1)}(t) & \text{im Fall } i = 0, \\ n B_{i-1}^{(n-1)}(t) - n B_i^{(n-1)}(t) & \text{im Fall } 0 < i < n, \\ n B_{n-1}^{(n-1)}(t) & \text{im Fall } i = n. \end{cases} \quad (16.6)$$

Beweis.

□

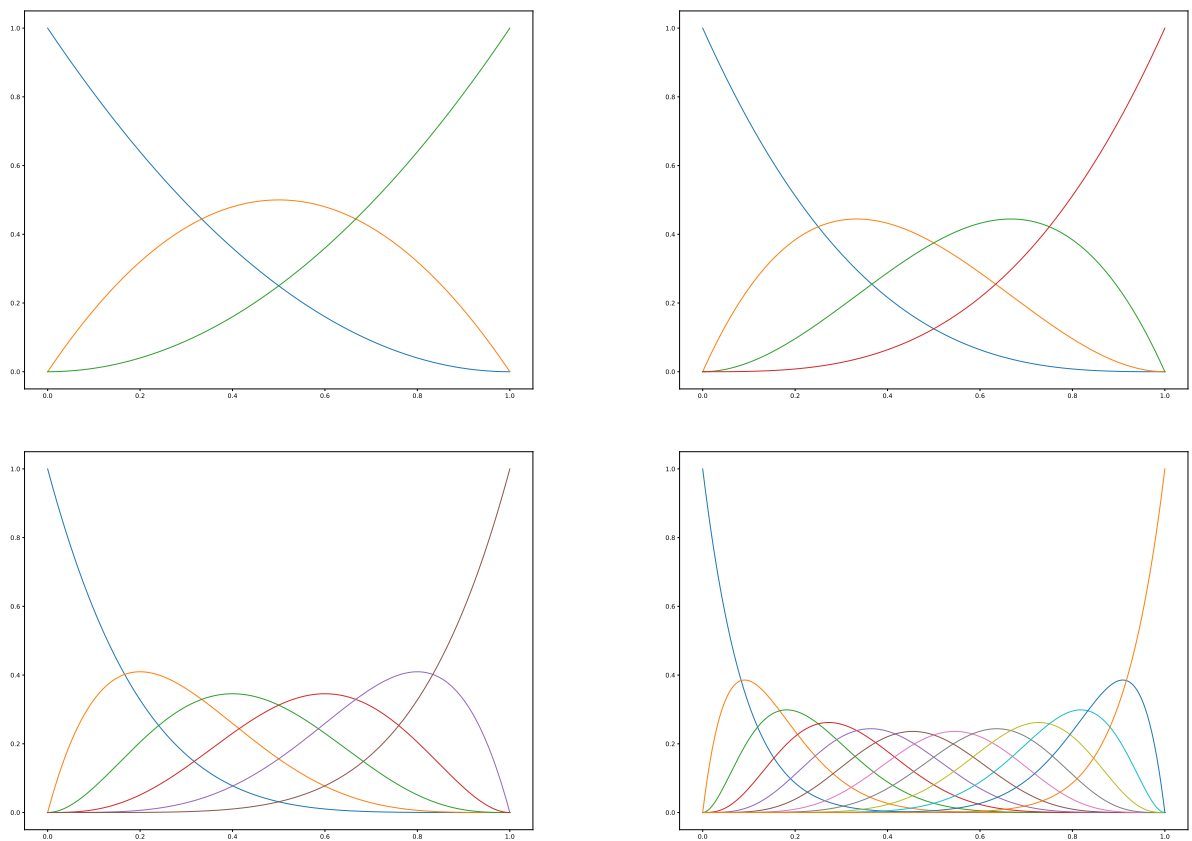


Abbildung 16.1: Illustration zu Bernstein-Polynome für verschiedene Grade  $n \in \{2, 3, 5, 11\}$  auf  $[0, 1]$ .

Wir sehen jetzt, wie Bernstein-Polynome verwendet werden, um glatte Kurven in der Ebene oder im Raum zu definieren.

**Definition 16.3** (Bézier-Kurve). Für gegebene Punkte  $b_0, \dots, b_n \in \mathbb{R}^d$  heißt das vektorwertige Polynom

$$B(t) := \sum_{i=0}^n B_i^{(n)}(t) b_i, \quad t \in [0, 1] \quad (16.7)$$

die zugehörige **Bézier-Kurve** (englisch: **Bézier curve**). Die Punkte  $b_0, \dots, b_n$  heißen die **Kontrollpunkte** (englisch: **control points**) der Bézier-Kurve.

**Quizfrage:** Kommt es auf die Reihenfolge der Kontrollpunkte an?

**Beispiel 16.4** (Bézier-Kurve). Gegeben sind die Kontrollpunkte

$$\begin{bmatrix} b_0 & b_1 & b_2 & b_3 \end{bmatrix} = \begin{bmatrix} 0 & \frac{1}{3} & \frac{2}{3} & 1 \\ 0 & -1 & -1 & 1 \end{bmatrix}$$

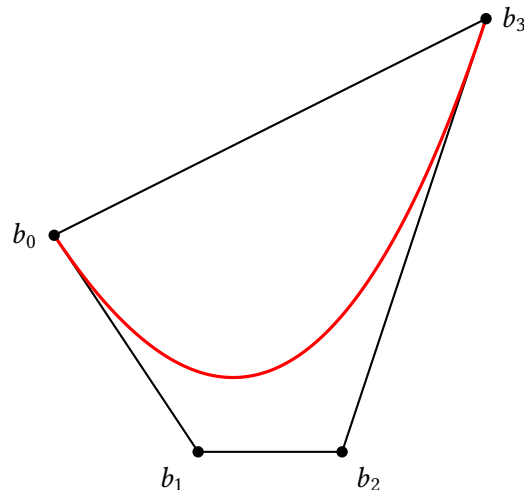
in  $\mathbb{R}^2$ . Die Bernstein-Polynome vom Grad 3 sind

$$\begin{aligned} B_0^{(3)}(t) &= 1 - 3t + 3t^2 - t^3, \\ B_1^{(3)}(t) &= 0 + 3t - 6t^2 + 3t^3, \\ B_2^{(3)}(t) &= 0 + 0t + 3t^2 - 3t^3, \\ B_3^{(3)}(t) &= 0 + 0t + 0t^2 + t^3. \end{aligned}$$

Die zu den Kontrollpunkten  $b_0, \dots, b_3$  gehörige Bézier-Kurve ist daher gegeben durch

$$\begin{aligned} B(t) &= (1 - 3t + 3t^2 - t^3) \begin{pmatrix} 0 \\ 0 \end{pmatrix} + (3t - 6t^2 + 3t^3) \begin{pmatrix} \frac{1}{3} \\ -1 \end{pmatrix} + (3t^2 - 3t^3) \begin{pmatrix} \frac{2}{3} \\ -1 \end{pmatrix} + (t^3) \begin{pmatrix} 1 \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} t \\ t^3 + 3t^2 - 3t \end{pmatrix} \end{aligned}$$

für  $t \in [0, 1]$ . Das folgende Bild zeigt diese Kurve (rot), die Kontrollpunkte sowie ihre konvexe Hülle.<sup>7</sup>



<sup>7</sup>Die konvexe Hülle einer Menge von Punkten in  $\mathbb{R}^d$  ist die kleinste konvexe Menge, die alle diese Punkte enthält.

Das Beispiel illustriert folgende Sachverhalte, die auch allgemein gelten:

**Satz 16.5** (Eigenschaften von Bézier-Kurven). Es seien  $b_0, \dots, b_n \in \mathbb{R}^d$  gegebene Kontrollpunkte,  $n \in \mathbb{N}_0$ , und  $B$  die dadurch erzeugte Bézier-Kurve (16.7). Dann gilt:

- (i)  $B(0) = b_0$  und  $B(1) = b_n$ .
- (ii) Im Falle von  $n \geq 1$  gilt für die Ableitung  $B'(0) = n(b_1 - b_0)$  und  $B'(1) = n(b_n - b_{n-1})$ .
- (iii) Die Kurve  $B(t)$  verläuft im Intervall  $[0, 1]$  innerhalb der konvexen Hülle von  $b_0, \dots, b_n$ . (Man bezeichnet diese konvexe Hülle der Kontrollpunkte als **Kontrollpolygon**.)

**Quizfrage:** Wie sieht man die Aussagen (ii) und (iii) im Bild oben?

Für die Auswertung von Punkten  $B(t)$  entlang der Kurve (etwa zum Zeichnen) ist die Verwendung der Darstellung (16.7) mitsamt der Definition (16.1) für  $B_i^{(n)}$  numerisch ungünstig. Stattdessen verwendet man den **Algorithmus von de Casteljau**. Dieser nutzt die rekursive Darstellung (16.5) aus. Für  $n \geq 1$  gilt nämlich:

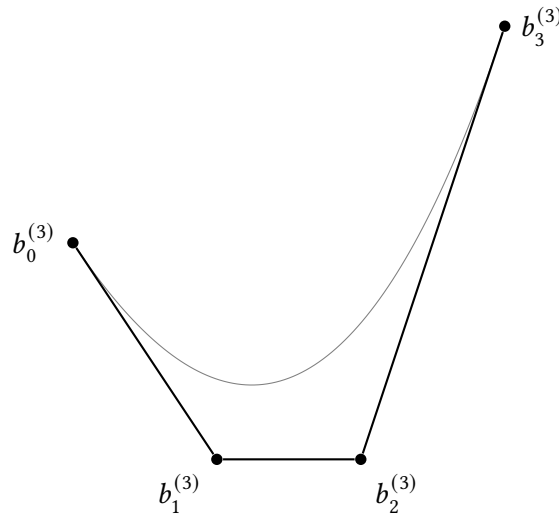
$$\begin{aligned}
 B(t) &= \sum_{i=0}^n B_i^{(n)}(t) b_i \\
 &= B_0^{(n)}(t) b_0 + \sum_{i=1}^{n-1} [B_i^{(n)}(t) b_i] + B_n^{(n)}(t) b_n \\
 &= (1-t) B_0^{(n-1)}(t) b_0 + \sum_{i=1}^{n-1} [t B_{i-1}^{(n-1)}(t) + (1-t) B_i^{(n-1)}(t)] b_i + t B_{n-1}^{(n-1)}(t) b_n \quad (\text{wegen (16.5)}) \\
 &= \sum_{i=0}^{n-1} B_i^{(n-1)}(t) \underbrace{[(1-t) b_i + t b_{i+1}]}_{\text{neue Kontrollpunkte durch Konvexkombination der vorherigen}} \quad (\text{durch Umsortieren}). \tag{16.8}
 \end{aligned}$$

Diese Darstellung sieht wieder genauso aus wie die, mit der wir begonnen haben, jedoch mit dem Grad  $n-1$  an Stelle von  $n$  und mit neuen Kontrollpunkten, die sich aus der Konvexkombination (anteiligen Mischung) von Paaren benachbarter Kontrollpunkte ergeben. Wir bezeichnen jetzt diese neuen Kontrollpunkte mit  $b_0^{(n-1)}, \dots, b_{n-1}^{(n-1)}$  und die ursprünglichen mit  $b_0^{(n)}, \dots, b_n^{(n)}$ . Durch wiederholte Anwendung dieses Arguments erhalten wir schließlich eine «Kurve» mit dem Bernstein-Polynom  $B_0^{(0)}$  (**Quizfrage:** Welche Funktion ist das?) und nur noch einem Kontrollpunkt  $b_0^{(0)}$ . Dieser ist der gesuchte Funktionswert  $B(t)$ .

Bevor wir den Algorithmus von de Casteljau angeben, wollen wir das Vorgehen noch grafisch anhand der Kurve aus **Beispiel 16.4** veranschaulichen.

**Beispiel 16.6** (Grafische Durchführung des Algorithmus von de Casteljau). Wir betrachten wieder

die durch die vier Kontrollpunkte  $b_0, \dots, b_3$  (die wir jetzt mit  $b_0^{(3)}, \dots, b_3^{(3)}$  bezeichnen) gegebene Bézier-Kurve  $B$  aus [Beispiel 16.4](#). Wir wollen diese an der Stelle  $t = 1/3$  auswerten.<sup>8</sup>



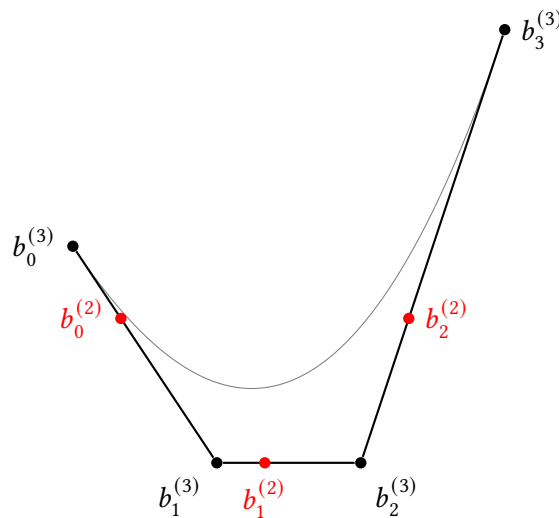
Nach (16.8) müssen zunächst die drei neuen Kontrollpunkte

$$b_0^{(2)} := \frac{2}{3}b_0^{(3)} + \frac{1}{3}b_1^{(3)}$$

$$b_1^{(2)} := \frac{2}{3}b_1^{(3)} + \frac{1}{3}b_2^{(3)}$$

$$b_2^{(2)} := \frac{2}{3}b_2^{(3)} + \frac{1}{3}b_3^{(3)}$$

konstruieren. Das heißt, dass wir die drei Strecken zwischen jeweils benachbarten Kontrollpunkten der Stufe 3 anteilig teilen und dort die drei neuen Kontrollpunkte der Stufe 2 einzeichnen:



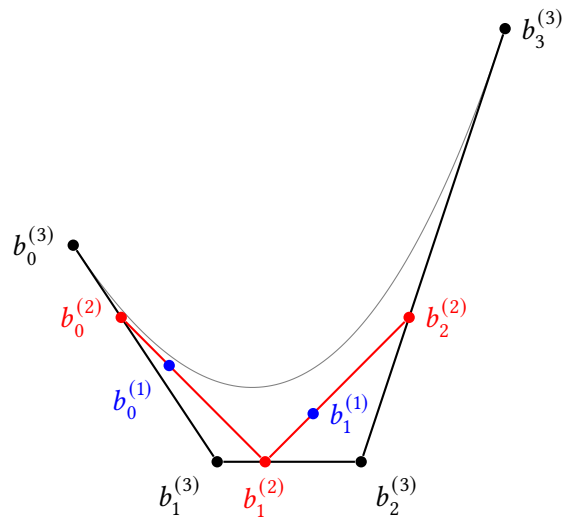
<sup>8</sup>Die Kurve ist zwar dünn in grau eingezeichnet, dies dient aber nur der Orientierung. Wir wollen ihre Funktionswerte, exemplarisch den für  $t = 1/3$ , ja gerade erst bestimmen.

Auf der nächste Stufe wiederholt sich das Ganze. Wir teilen also die Strecken zwischen den gerade hinzugefügten Kontrollpunkten der Stufe 2 wiederum anteilig und fügen die zwei neuen Kontrollpunkte

$$b_0^{(1)} := \frac{2}{3}b_0^{(2)} + \frac{1}{3}b_1^{(2)}$$

$$b_1^{(1)} := \frac{2}{3}b_1^{(2)} + \frac{1}{3}b_2^{(2)}$$

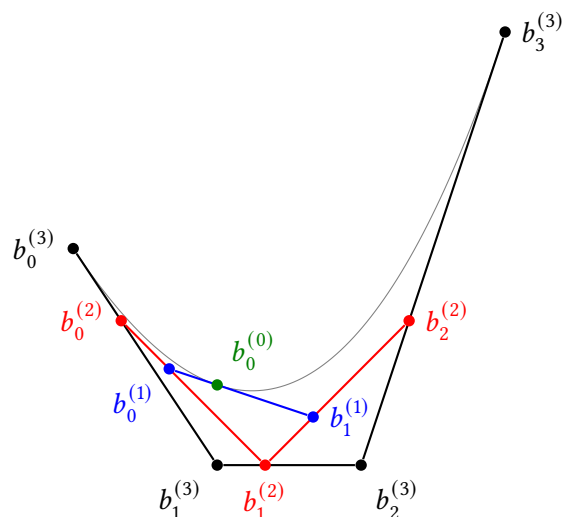
der Stufe 1 hinzu:



Schließlich wird auch die letzte Strecke wieder geteilt. Der neue Kontrollpunkt

$$b_0^{(0)} := \frac{2}{3}b_0^{(1)} + \frac{1}{3}b_1^{(1)}$$

der Stufe 0 ist der gesuchte Punkt  $B(t)$  auf der Kurve:





Eine ansprechende interaktive Version dieser grafischen Konstruktion finden Sie zum Beispiel unter <https://pomax.github.io/bezierinfo/#decasteljau>.

**Quizfrage:** Wie kann man mit Hilfe der grafischen Konstruktion begründen, dass die Kurve  $t \mapsto B(t)$  innerhalb der konvexen Hülle der Kontrollpunkte verläuft (Aussage (iii) aus Satz 16.5)?

Wir können das Verfahren nun formal angeben:

**Algorithmus 16.7** (Algorithmus von de Casteljau zur Auswertung einer Bézier-Kurve).

**Eingabe:** Kontrollpunkte  $b_0, \dots, b_n \in \mathbb{R}^d$  ( $n \in \mathbb{N}_0$ ) einer Bézier-Kurve  $B$

**Eingabe:** Auswertungsstelle  $t \in [0, 1]$

**Ausgabe:** Funktionswert  $y = B(t)$

```

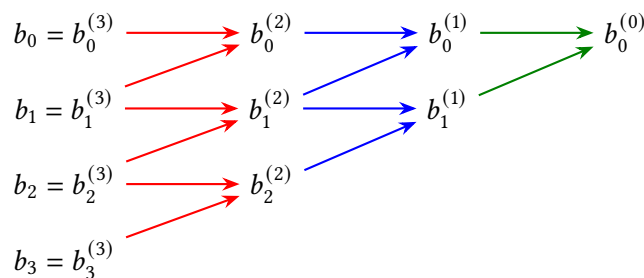
1: for  $i = 0, \dots, n$  do
2:   Setze  $b_i^{(n)} := b_i$ 
3: end for
4: for  $\ell = n - 1, \dots, 0$  do
5:   for  $i = 0, \dots, \ell$  do
6:     Setze  $b_i^{(\ell)} := (1 - t) b_i^{(\ell+1)} + t b_{i+1}^{(\ell+1)}$ 
7:   end for
8: end for
9: return  $y := b_0^{(0)}$ 

```

**Quizfrage:** Wieviele FLOPs benötigt ungefähr eine Auswertung  $B(t)$ ?

**Quizfrage:** Wie könnte man das Verfahren «vektorisieren», sodass die Schleifen über  $i$  jeweils wegfallen?

Die Struktur der Kontrollpunkte auf den verschiedenen Stufen und ihre Abhängigkeiten ähnelt dem Schema für die dividierten Differenzen, siehe nach Definition 15.8:



Unsere Variante von Algorithmus 16.7 berechnet die Kontrollpunkte Spalte für Spalte, wir könnten allerdings wie bei den dividierten Differenzen üblich auch auch Diagonale für Diagonale vorgehen.

**Quizfrage:** Wie kann das Hinzufügen einer Diagonale am unteren Ende grafisch gedeutet werden?

**Quizfrage:** Beim Vorgehen Diagonale für Diagonale erhalten wir als letzte Größe jeweils diejenige, die oben in der ersten Zeile steht, also  $b_0^{(3)}$ , dann  $b_0^{(2)}$ , anschließend  $b_0^{(1)}$  und schließlich  $b_0^{(0)}$ . Wie lassen sich diese Größen interpretieren?

Bézier-Kurven haben viele vorteilhafte Eigenschaften:

- (1) Auch für die Bestimmung der ersten und höherer Ableitungen lässt sich ein Verfahren ähnlich dem [Algorithmus 16.7](#) von de Casteljau angeben.
- (2) Bézier-Kurven sind affin invariant. Das heißt, dass das Bild einer Bézier-Kurve  $B_{n,GM}$  unter einer affinen Abbildung  $x \mapsto Ax + b$  wieder eine Bézier-Kurve ist, und zwar diejenige, die durch die affin transformierten Kontrollpunkte festgelegt wird.
- (3) Bestimmte Operationen lassen sich für Bézier-Kurven effizient ausführen, darunter die Teilung einer Kurve in zwei Kurven und die (scheinbare) Erhöhung des Grades<sup>9</sup> einer Kurve.
- (4) Bézier-Kurven lassen sich zu **Bézier-Flächen** verallgemeinern. In dieser Form haben sie große Bedeutung im *Computer Aided Design* (CAD).

## § 17 APPROXIMATION

Wie in der Einleitung zu [Kapitel 6](#) erwähnt, geht es bei der Approximation einer gegebenen Funktion  $f$  darum, eine andere Funktion  $g$  aus einer Menge von zur Verfügung stehenden Funktionen  $g$  auszuwählen, dass Approximationsfehler  $\|f - g\|$  in einer geeigneten Norm minimiert wird.

### § 17.1 TSCHEBYSCHOW-APPROXIMATION

Eine häufig verwendete Approximationsaufgabe ergibt sich aus dem Wunsch, den *maximalen punktweisen* Approximationsfehler in einem gegebenen Intervall  $[a, b]$  zu minimieren. Der Fehler wird also in der bereits bekannten Supremumsnorm gemessen:

$$\|f - p\|_{C([a,b])} := \sup_{x \in [a,b]} |f(x) - p(x)|.$$

Zur Approximation verwenden wir wieder Polynome mit gegebenem Höchstgrad  $n$ . Das führt auf die folgende Aufgabe: Zu einer gegebenen stetigen Funktion  $f: [a, b] \rightarrow \mathbb{R}$  finde ein Polynom  $p \in P_n$ , das die Aufgabe

$$\text{Minimiere } \|f - p\|_{C([a,b])}, \quad p \in P_n \tag{17.1}$$

löst. Diese Aufgabe heißt die **Tschebyschow-Approximationsaufgabe**<sup>10</sup>, häufig ungenau in der Form «**Tschebyscheff**» transkribiert (englisch: *Chebyshev approximation problem*).

<sup>9</sup>Dadurch erhält man mehr Kontrollpunkte, also auch mehr Möglichkeiten, die Kurve z. B. in einem Grafikprogramm zu verändern.

<sup>10</sup>Чебышёв

**Satz 17.1** (Existenz ~~und Eindeutigkeit~~ der Lösung der Tschebyschow-Approximationsaufgabe (17.1), vgl. Rannacher, 2017, Satz 2.8). Für jede stetige Funktion  $f: [a, b] \rightarrow \mathbb{R}$  und jeden Polynomgrad  $n \in \mathbb{N}_0$  besitzt die Tschebyschow-Approximationsaufgabe (17.1) eine Lösung.

*Beweis.* Es ist ausreichend, einen Minimierer innerhalb der Menge

$$\{p \in P_n \mid \|p\|_{C([a,b])} \leq 2 \|f\|_{C([a,b])}\} \quad (17.2)$$

zu suchen, denn jedes  $p \in P_n$  mit der Eigenschaft  $\|p\|_{C([a,b])} > 2 \|f\|_{C([a,b])}$  erfüllt

$$\begin{aligned} \|f - p\|_{C([a,b])} &\geq \|p\|_{C([a,b])} - \|f\|_{C([a,b])} \quad (\text{umgedrehte Dreiecksungleichung}) \\ &> \|f\|_{C([a,b])} \quad (\text{Voraussetzung}) \\ &= \|f - 0\|_{C([a,b])} \\ &\geq \inf_{p \in P_n} \|f - p\|_{C([a,b])} \end{aligned}$$

und ist damit sicher kein Minimierer. Die Menge (17.2) ist aber als beschränkte abgeschlossene Menge (Kugel) in einem endlich-dimensionalen Raum kompakt. Da die Zielfunktion in (17.1) stetig ist (**Quizfrage:** Warum ist das so?), existiert nach dem **Satz von Weierstraß** ein Minimierer, also eine Lösung von (17.1).  $\square$

Minimierer von (17.1) haben folgende interessante charakteristische Eigenschaft:

**Satz 17.2 (Alternantensatz)** (englisch: *equioscillation theorem*), vgl. Rannacher, 2017, Satz 2.19). Folgende Bedingungen sind für eine Funktion  $p \in P_n$  äquivalent:

(i)  $p$  ist eine Lösung von (17.1).

(ii) Es existieren  $m \geq n + 2$  Stellen  $a \leq \xi_1 < \xi_2 < \dots < \xi_m \leq b$  und  $s \in \{\pm 1\}$ , sodass für die Fehlerfunktion  $f - p$  gilt:

$$f(\xi_i) - p(\xi_i) = s(-1)^i \|f - p\|_{C([a,b])}, \quad i = 1, \dots, m. \quad (17.3)$$

Die Menge der Extremstellen  $\{\xi_1, \dots, \xi_m\}$  wird **Alternante** der Funktion  $f - p$  genannt.

Wir verzichten hier auf einen Beweis. **Quizfrage:** Was bedeutet die Eigenschaft (17.3)?

**Folgerung 17.3** (Eindeutigkeit der Lösung von (17.1), vgl. Rannacher, 2017, Korollar 2.2). Die Lösung der Tschebyschow-Approximationsaufgabe (17.1) ist eindeutig.

*Beweis.* Nehmen wir an,  $p_1$  und  $p_2$  seien beides Lösungen von (17.1). Wir kürzen die Fehlerfunktionen mit  $e_1 := f - p_1$  und  $e_2 := f - p_2$  ab. Wir können annehmen, dass  $e_1$  nicht die Nullfunktion ist, sonst wäre auch  $e_2$  die Nullfunktion (**Quizfrage:** Warum?) und  $p_1 = p_2 = f$  gezeigt. Es gilt also  $\|e_1\|_{C([a,b])} = \|e_2\|_{C([a,b])} > 0$ .

Es sei nun  $\lambda \in (0, 1)$  beliebig. Dann ist  $\|\lambda e_2\|_{C([a,b])} < \|e_1\|_{C([a,b])}$ . Daraus folgt, dass der Graph von  $\lambda e_2$  den Graphen von  $e_1$  mindestens  $n+1$  Mal schneidet. (**Quizfrage:** Klar? ) Die Funktion  $e_\lambda := e_1 - \lambda e_2$  hat also mindestens  $n+1$  verschiedene Nullstellen. Im Grenzübergang  $\lambda \nearrow 1$  fallen einige davon möglicherweise zusammen, es bleiben jedoch  $n+1$  Nullstellen (ihrer Vielfachheit entsprechend oft gezählt) von  $e_1 - e_2 = p_2 - p_1$ . Da jedoch  $p_2 - p_1 \in P_n$  ist, muss dies die Nullfunktion sein.  $\square$

Für die Lösung der Tschebyschow-Optimierungsaufgabe (17.1) gibt es Verfahren, die auf der Struktur der Lösung nach dem **Alternantensatz 17.2** aufbauen, darunter der **Algorithmus von Remez**. Eine Alternative ist die Formulierung als eine lineare Optimierungsaufgabe (mit eigentlich unendlich vielen Nebenbedingungen, die man dann in der Regel diskretisiert), siehe Vorlesung *Grundlagen der Optimierung*.

Ein Beispiel für die Tschebyschow-Approximation wird in **Abbildung 17.1** gezeigt.

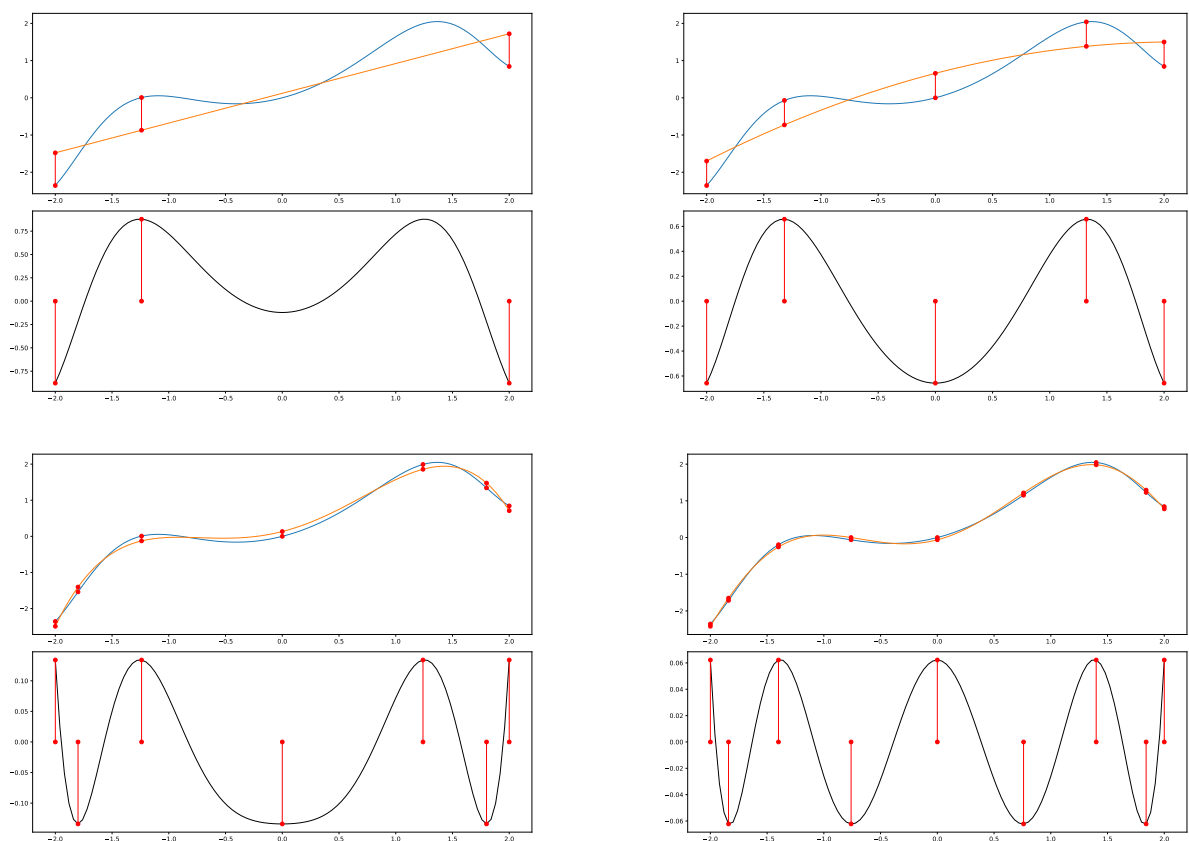


Abbildung 17.1: Illustration einer Tschebyschow-Approximation durch Polynome verschiedener Grade  $n \in \{1, 3, 5, 7\}$  auf  $[-2, 2]$ .

## § 17.2 ANWENDUNG DER TSCHEBYSCHOW-APPROXIMATION ZUR OPTIMIERUNG DER STÜTZSTELLEN

Wir kommen noch einmal zurück auf die Darstellung (15.25) des Approximationsfehlers aus Satz 15.11

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \ell(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \prod_{j=0}^n (x - x_j) \quad (15.25)$$

bei der Polynominterpolation einer Funktion  $f \in C^{n+1}([a, b])$ .

**Frage:** Wie sind die Stützstellen  $x_0, \dots, x_n \in [a, b]$  zu wählen, sodass das Knotenpolynom  $\ell(x) = \prod_{j=0}^n (x - x_j)$  minimale Ausschläge hat? Wir betrachten also die Optimierungsaufgabe

$$\text{Minimiere} \quad \sup_{x \in [a, b]} \left| \prod_{j=0}^n (x - x_j) \right|, \quad x_0, \dots, x_n \in [a, b]. \quad (17.4)$$

Die Funktion  $\prod_{j=0}^n (x - x_j)$  ist gleich  $x^{n+1} - p$  für ein Polynom  $p \in P_n$ , das von den Stützstellen abhängt. Daher können wir die Aufgabe (17.4) auch so lesen, dass wir die Tschebyschow-Approximation (17.1) der Funktion  $x \mapsto x^{n+1}$  auf  $[a, b]$  suchen, also

$$\text{Minimiere} \quad \sup_{x \in [a, b]} |x^{n+1} - p(x)|, \quad p \in P_n. \quad (17.5)$$

Die Lösung der Aufgabe (17.5) kann explizit angegeben werden. Es reicht dabei aus, das Intervall  $[a, b] = [-1, 1]$  zu betrachten. Für andere Intervalle kann das Ergebnis nachher übertragen werden. Die Darstellung der Lösung in Satz 17.5 verwendet die folgenden Funktionen:

**Definition 17.4** (Tschebyschow-Polynome). Es sei  $n \in \mathbb{N}_0$ . Das  $n$ -te **Tschebschow-Polynom** (erster Art) (englisch: **Chebyshev polynomial (of the first kind)**) ist definiert durch die Rekursion

$$T_0(x) := 1, \quad (17.6a)$$

$$T_1(x) := x, \quad (17.6b)$$

$$T_{n+1}(x) := 2x T_n(x) - T_{n-1}(x). \quad (17.6c)$$

Mit Hilfe dieser Definition können wir leicht zeigen, dass  $T_n$  ein Polynom vom Grad  $n$  ist, dessen führender Koeffizient im Falle von  $n \geq 1$  gleich  $2^{n-1}$  ist. Eine andere Darstellung ist

$$T_n(x) = \cos(n \arccos(x)) \quad \text{für } x \in [-1, 1]. \quad (17.7)$$

Die  $n$  Nullstellen von  $T_n$  sind

$$x_j = \cos\left(\frac{\pi}{2} \frac{2j+1}{n}\right), \quad j = 0, \dots, n-1 \quad (17.8)$$

und heißen **Tschebyschow-Knoten** (englisch: **Chebyshev nodes**).

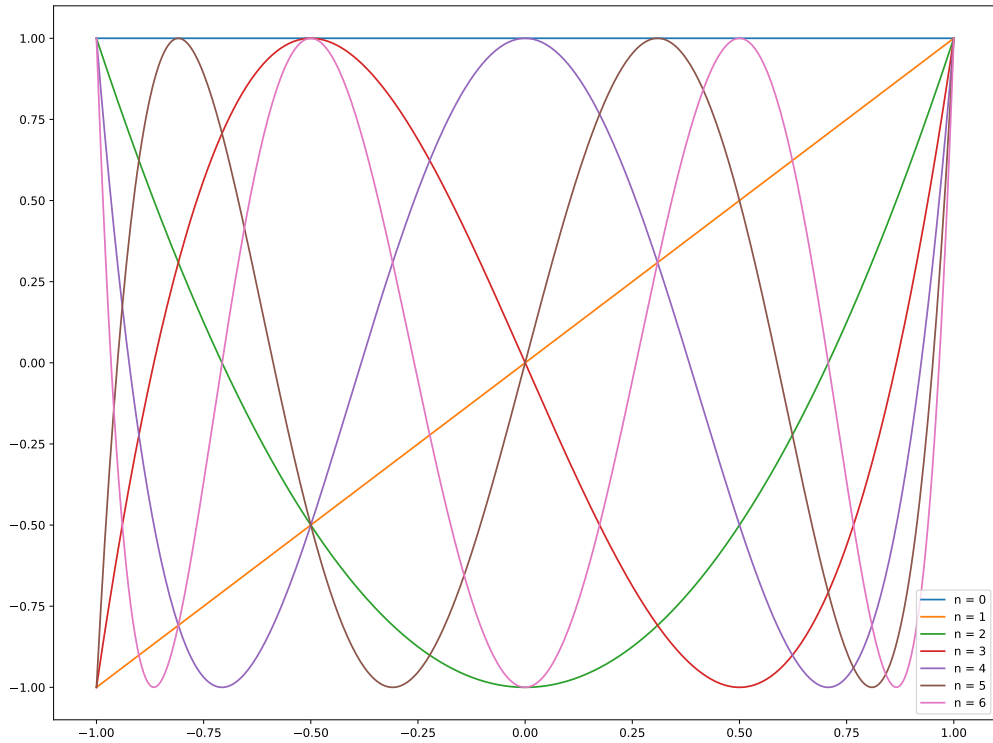


Abbildung 17.2: Illustration der Tschebyschow-Polynome verschiedener Grade  $n \in \{0, 1, \dots, 6\}$  auf  $[-1, 1]$ .

**Satz 17.5** (Lösung der Tschebyschow-Approximationsaufgabe (17.5) und optimale Stützstellen, vgl. Rannacher, 2017, Satz 2.20). Es sei  $n \in \mathbb{N}_0$ . Die eindeutige Lösung der Aufgabe (17.5) auf  $[a, b] = [-1, 1]$  ist gegeben durch

$$p(x) := x^{n+1} - \frac{1}{2^n} T_{n+1}(x). \quad (17.9)$$

Die Tschebyschow-Knoten (Nullstellen von  $T_{n+1}$ ) sind die optimalen Stützstellen von (17.4).

*Beweis.* Wir wissen, dass  $T_{n+1}$  den führenden Term  $2^n x^{n+1}$  besitzt, und kennen außerdem seine Nullstellen  $x_0, \dots, x_n$ , siehe (17.8). Daraus folgt, dass  $\frac{1}{2^n} T_{n+1}$  die Darstellung

$$\frac{1}{2^n} T_{n+1}(x) = \prod_{j=0}^n (x - x_j) = \ell(x)$$

besitzen muss. Aus der Darstellung (17.7) folgt, dass  $T_{n+1}$  im Intervall  $[-1, 1]$  genau  $(n+2)$  Mal seinen (positiven bzw. negativen) Extremwert annimmt, und zwar abwechselnd  $\pm 1$ . Dasselbe gilt für  $\frac{1}{2^n} T_{n+1}$  mit dem Extremwert  $\pm \frac{1}{2^n}$ . Mit anderen Worten: Diese  $(n+2)$  Extremstellen bilden eine Alternante der Funktion  $x \mapsto x^{n+1} - \frac{1}{2^n} T_{n+1}(x)$ , die ein Polynom in  $P_n$  ist. Der Alternantensatz 17.2 garantiert jetzt, dass  $x \mapsto x^{n+1} - \frac{1}{2^n} T_{n+1}(x)$  die Aufgabe (17.5) löst.

Es bleibt noch zu zeigen, dass die Nullstellen  $x_j$  aus (17.8) die optimale Stützstellen von (17.4) sind. Wir

haben

$$\begin{aligned} \sup_{x \in [-1,1]} \left| \prod_{j=0}^n (x - x_j) \right| &= \frac{1}{2^n} \sup_{x \in [-1,1]} |T_{n+1}(x)| \\ &= \sup_{x \in [-1,1]} \left| x^{n+1} - \left[ x^{n+1} - \frac{1}{2^n} T_{n+1}(x) \right] \right|, \quad \text{denn } x \mapsto x^{n+1} - \frac{1}{2^n} T_{n+1}(x) \text{ löst (17.5)} \\ &= \min_{p \in P_n} \sup_{x \in [-1,1]} |x^{n+1} - p(x)|. \end{aligned}$$

Wir wissen, dass  $x^{n+1} - p(x)$  irgendein Polynom in  $P_{n+1}$  mit führendem Koeffizienten 1 ist. Die nachstehend verwendeten Polynome  $\prod_{j=0}^n (x - y_j)$  mit Stützstellen  $y_j \in [-1, 1]$  sind ebenfalls Polynome in  $P_{n+1}$  mit führendem Koeffizienten 1 (aber möglicherweise nicht alle). Daher ist

$$\leq \min \left\{ \sup_{x \in [-1,1]} \left| \prod_{j=0}^n (x - y_j) \right| \mid y_j \in [-1, 1], j = 0, 1, \dots, n \right\}.$$

Das zeigt, dass  $\prod_{j=0}^n (x - x_j)$  tatsächlich optimale Stützstellen für (17.4) sind. □

**Bemerkung 17.6** (Bedeutung der Tschebyschow-Knoten). *Die Tschebyschow-Knoten sind diejenigen Stützstellen, die ohne Ansehen der zu interpolierenden Funktion  $f$  optimale (kleinstmögliche) Fehlerabschätzungen ermöglichen.*

Wir greifen nochmal die Runge-Funktion aus (15.13) auf und interpolieren sie durch Polynome wachsenden Grades, wobei wir als Stützstellen jetzt aber die jeweiligen Tschebyschow-Knoten verwenden. Dieses Mal können wir die gleichmäßige Konvergenz der Folge der Interpolationspolynome beobachten, siehe [Abbildung 17.3](#).

Ende der Woche 10

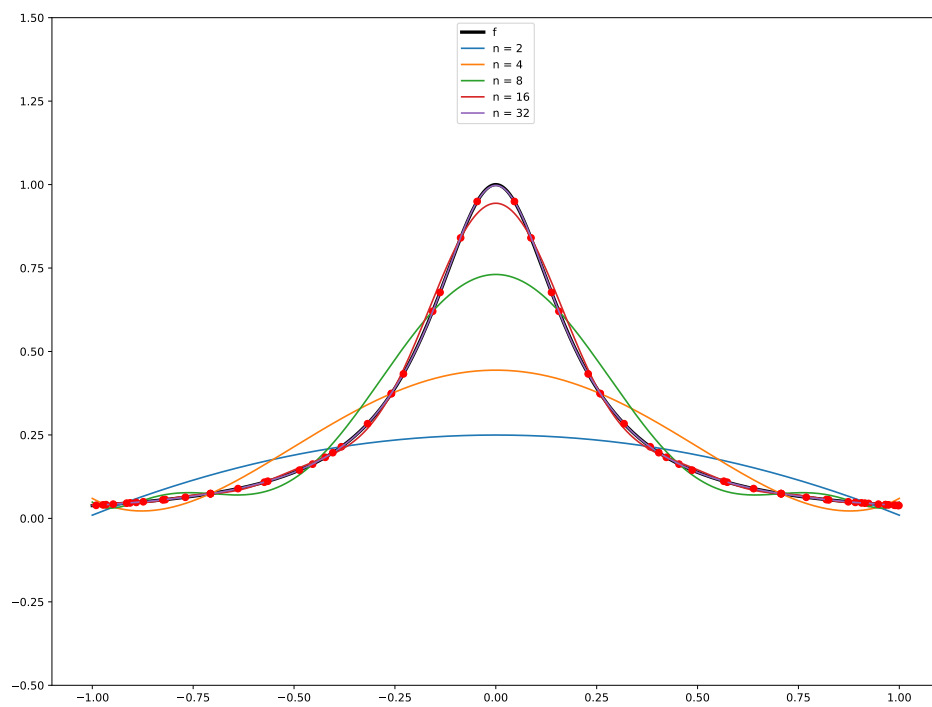


Abbildung 17.3: Interpolation der Runge-Funktion (schwarz) aus [Beispiel 15.13](#) durch Polynome mit Tschebyschow-Knoten auf  $[-1, 1]$  als Stützstellen.



## Kapitel 7 Numerische Integration

Die **numerische Integration** oder **numerische Quadratur** (englisch: *numerical integration, numerical quadrature*) umfasst numerische Verfahren zur approximativen Berechnung von Integralen. Wir betrachten hier nur Integrale  $\int_a^b f(x) dx$  über Funktionen  $f: [a, b] \rightarrow \mathbb{R}$  von einer Variablen. Genau wie in [Kapitel 6](#) (Interpolation und Approximation) wird es dabei wieder einen Approximationsfehler geben, den es abzuschätzen und zu kontrollieren gilt.

Die Notwendigkeit der numerischen Integration ergibt sich daraus, dass die Bestimmung einer Stammfunktion nur für «wenige» Funktionen möglich ist. (**Quizfrage:** Beispiel für eine Funktion, für die keine geschlossene Formel für eine Stammfunktion bekannt ist?) Außerdem könnte die zu integrierende Funktion nur an endlichen vielen Stellen bekannt sein (etwa aus Messungen). (**Quizfrage:** Können Sie ein praktisches Beispiel für diese Situation nennen?)

Alle Quadraturformeln werden die Darstellung

$$\int_a^b f(x) dx \approx \sum_{i=0}^n w_i f(x_i)$$

haben, wobei die Koeffizienten  $w_i$  als die **Gewichte** (englisch: *weights*) bezeichnet werden und die Zahlen  $x_i \in [a, b]$  als die **Stützstellen** (englisch: *integration points*). Eine Quadraturformel ist also eine gewichtete Summe (Linearkombination) von Funktionswerten.

Der **Genauigkeitsgrad** (englisch: *degree of exactness*) einer Quadraturformel ist definiert als die größte Zahl  $m \in \mathbb{N}_0$ , sodass alle Polynome vom Höchstgrad  $m$  exakt integriert werden. Eine Formel mit Genauigkeitsgrad 0 integriert also konstante Polynome exakt, aber nicht alle linearen Polynome. Eine Formel mit Genauigkeitsgrad 1 integriert alle konstanten und linearen Polynome exakt, aber nicht alle quadratischen Polynome usw.

### § 18 INTERPOLATORISCHE QUADRATURFORMELN

Es ist naheliegend, zur Konstruktion von Quadraturformeln folgende Strategie zu verwenden: Wir interpolieren die zu integrierende Funktion  $f$  durch ein Polynom  $p$  und integrieren dieses.<sup>1</sup> Wir verwenden zur Herleitung die Lagrangesche Darstellung des Interpolationspolynoms aus [Abschnitt 15.2](#)

---

<sup>1</sup>Dieselbe Idee (für die Differentiation an Stelle der Integration) lag schon den Finite-Differenzen-Formeln in [Abschnitt 15.6](#) zugrunde.

mit Stützstellen  $a \leq x_0 < \dots < x_n \leq b$ , vgl. [Gleichung \(15.9\)](#):

$$p(x) = \sum_{i=0}^n f(x_i) L_i^{(n)}(x). \quad (18.1)$$

Daraus ergibt sich die Quadraturformel

$$\sum_{i=0}^n w_i f(x_i) \quad \text{mit Gewichten } w_i := \int_a^b L_i^{(n)}(x) dx. \quad (18.2)$$

Die Gewichte hängen nur von den Stützstellen  $x_0, \dots, x_n$  ab. (**Quizfrage:** Warum verwenden wir nicht die Newtonsche Darstellung des Interpolationspolynoms aus [Abschnitt 15.3](#)?)

**Lemma 18.1** (Schranken für den Genauigkeitsgrad interpolatorischer Quadraturformeln). *Jede interpolatorische Quadraturformel mit  $n+1$  (paarweise verschiedenen) Stützstellen hat mindestens den Genauigkeitsgrad  $n$  und höchstens den Genauigkeitsgrad  $2n+1$ .*

*Beweis.* Wir wissen, dass das Interpolationspolynom (18.1) Polynome vom Höchstgrad  $\leq n$  exakt darstellt. Daraus folgt sofort, dass der Genauigkeitsgrad mindestens  $n$  ist. Um die obere Schranke zu zeigen, betrachten wir die Funktion

$$f(x) := \prod_{i=0}^n (x - x_i)^2,$$

ein Polynom vom Grad  $2(n+1)$ . Da  $f(x_i) = 0$  ist für alle  $i = 0, \dots, n$ , ergibt die Quadraturformel den Wert 0, während der exakte Wert des Integrals  $\int_a^b f(x) dx$  echt positiv ist.  $\square$

## § 18.1 ABGESCHLOSSENE NEWTON-COTES-FORMELN

Die Newton-Cotes-Formeln sind diejenigen interpolatorischen Quadraturformeln, die man erhält, wenn man die Stützstellen im Intervall  $[a, b]$  gleichverteilt. Die **(ab)geschlossenen Newton-Cotes-Formeln** (englisch: *closed Newton-Cotes formulas*) verwenden dabei Stützstellen inklusive der Randpunkte  $a$  und  $b$ , also

$$x_i := a + i h \quad \text{mit } h = \frac{b-a}{n}, \quad i = 0, \dots, n. \quad (18.3)$$

Dabei ist  $n \geq 1$ .

Die Berechnung der Gewichte für beliebige Intervalle  $[a, b]$  kann mit Hilfe der Koordinatentransformation  $x = a + t h$  für  $t \in [0, n]$  auf das Intervall  $[0, n]$  zurückgeführt werden. Es gilt

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j} = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{a + t h - (a + j h)}{a + i h - (a + j h)} = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - j}{i - j}$$

$n$	Gewichte	Name
1	$\frac{b-a}{2} [1 \quad 1]$	Trapezregel
2	$\frac{b-a}{6} [1 \quad 4 \quad 1]$	Simpson-Regel, Keplersche Fassregel
3	$\frac{b-a}{8} [1 \quad 3 \quad 3 \quad 1]$	3/8-Regel, «Pulcherrima»
4	$\frac{b-a}{90} [7 \quad 32 \quad 12 \quad 32 \quad 7]$	Milne-Regel, Boole-Regel
5	$\frac{b-a}{288} [19 \quad 75 \quad 50 \quad 50 \quad 75 \quad 19]$	6-Punkt-Regel
6	$\frac{b-a}{840} [41 \quad 216 \quad 27 \quad 272 \quad 27 \quad 216 \quad 41]$	Weddle-Regel

Tabelle 18.1: Abgeschlossene Newton-Cotes-Formeln auf  $[a, b]$ .

und daher

$$w_i = \int_a^b L_i^{(n)}(x) dx = h \int_0^n \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t-j}{i-j} dt, \quad i = 0, \dots, n. \quad (18.4)$$

**Quizfrage:** Wo kommt der Faktor  $h$  vor dem Integral her?

Die Gewichte der ersten abgeschlossenen Newton-Cotes-Formeln sind in [Tabelle 18.1](#) angegeben.

**Quizfrage:** Warum summieren sich die Gewichte eigentlich zu  $b - a$ ?

**Quizfrage:** Und warum sind die Gewichte alles *rationale* Vielfache von  $b - a$ ?

Die ersten beiden Quadraturformeln aus [Tabelle 18.1](#) lauten ausgeschrieben

$$\frac{b-a}{2} [f(a) + f(b)] \quad \text{Trapezregel,} \quad (18.5a)$$

$$\frac{b-a}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] \quad \text{Simpson-Regel.} \quad (18.5b)$$

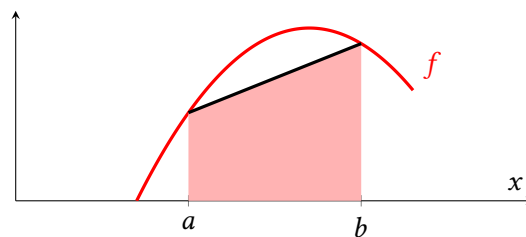


Abbildung 18.1: Veranschaulichung der Trapezregel (18.5a).

**Satz 18.2** (Genauigkeitsgrad der abgeschlossenen Newton-Cotes-Formeln). *Der Genauigkeitsgrad der abgeschlossenen Newton-Cotes-Formeln mit  $n + 1$  Stützstellen ist gleich  $n$  für ungerades  $n \in \mathbb{N}$  und gleich  $n + 1$  für gerades  $n \in \mathbb{N}$ .*

**Beachte:** Die abgeschlossenen Newton-Cotes-Formeln mit geradem  $n$  (also diejenigen mit einer ungeraden Anzahl von Stützstellen) sind also vorzuziehen.

*Beweis.* Wie oben in Lemma 18.1 bereits erwähnt, ist der Genauigkeitsgrad mindestens  $n$ . Es sei nun  $n \geq 2$  gerade. Wir betrachten nur den Fall  $[a, b] = [-1, 1]$ . Die Stützstellen sind dann symmetrisch um die mittlere Stützstelle bei 0 verteilt, und die Gewichte sind ebenfalls symmetrisch, also  $w_i = w_{n-i}$ . Wir betrachten den Integrationsfehler der Funktion  $x \mapsto x^{n+1}$ . (**Quizfrage:** Warum reicht es aus, Monome separat zu betrachten? Müssen wir nicht allgemeine Polynome vom Höchstgrad  $n+1$  ansehen?) Das exakte Integral  $\int_{-1}^1 x^{n+1} dx$  ist gleich 0. Der durch die Quadraturformel bestimmte Wert  $\sum_{i=0}^n w_i x_i^{n+1}$  ist ebenfalls gleich 0. (**Quizfrage:** Warum?)  $\square$

**Quizfrage:** Wie können wir uns davon überzeugen, dass das Resultat auch für allgemeine Intervalle  $[a, b]$  gilt und nicht nur für  $[-1, 1]$ ?

Ab  $n \geq 8$  sind einige der Gewichte in den abgeschlossenen Newton-Cotes-Formeln negativ. Das ist aus verschiedenen Gründen ungünstig:

- (1) Die Gefahr der Auslöschung erhöht sich durch negative Gewichte, da in der Praxis benachbarte Funktionswerte des Integranden  $f$  oft dasselbe Vorzeichen haben.
- (2) Für einen Integranden  $f \geq 0$  können sich für den Näherungswert des Integrals negative Werte ergeben. Das kann für bestimmte Anwendungen inakzeptabel sein, beispielsweise wenn  $f$  eine Stoffdichte repräsentiert und damit das Integral eine Masse. Äquivalent können wir das auch so ausdrücken, dass wir die Monotonie der exakten Integration

$$f \geq g \quad \Rightarrow \quad \int_a^b f(x) dx \geq \int_a^b g(x) dx$$

in der Quadraturformel verlieren, wenn negative Gewichte vorliegen.

Die abgeschlossenen Newton-Cotes-Formeln werden daher i. d. R. nur für  $1 \leq n \leq 7$  verwendet.

Wir geben nun noch eine Abschätzung des Quadraturfehlers, also des Approximationsfehlers bei der Quadratur, an. Die Abschätzung des Fehlers (vgl. (15.25), (15.26))

$$|f(x) - p(x)| \leq \sup_{x \in [a, b]} \frac{|f^{(n+1)}(x)|}{(n+1)!} \left| \prod_{i=0}^n (x - x_i) \right|$$

bei der Polynominterpolation kann durch Integration unmittelbar in eine Abschätzung des Quadraturfehlers überführt werden. Es gilt

$$\begin{aligned} & \left| \int_a^b f(x) dx - \sum_{i=0}^n w_i f(x_i) \right| \\ &= \left| \int_a^b f(x) - p(x) dx \right| \\ &\leq \int_a^b |f(x) - p(x)| dx \quad (\text{Dreiecksungleichung für Integrale}) \\ &\leq \sup_{x \in [a, b]} \frac{|f^{(n+1)}(x)|}{(n+1)!} \int_a^b \left| \prod_{i=0}^n (x - x_i) \right| dx \end{aligned}$$

Durch Substitution  $x = a + t h$  mit  $h = \frac{b-a}{n}$  wie in (18.4) folgt weiter:

$$\begin{aligned} &= \frac{\|f^{(n+1)}\|_{C([a,b])}}{(n+1)!} h \int_0^n \left| \prod_{i=0}^n h(t-i) \right| dt \\ &= \frac{\|f^{(n+1)}\|_{C([a,b])}}{(n+1)!} h^{n+2} \int_0^n \left| \prod_{i=0}^n (t-i) \right| dt \\ &=: C_n \|f^{(n+1)}\|_{C([a,b])} h^{n+2} \end{aligned}$$

mit

$$C_n := \frac{1}{(n+1)!} \int_0^n \left| \prod_{i=0}^n (t-i) \right| dt. \quad (18.6)$$

Wir fassen zusammen:

**Satz 18.3** (Fehlerabschätzung der abgeschlossenen Newton-Cotes-Formeln). Für  $f \in C^{n+1}([a, b])$  erfüllt die abgeschlossene Newton-Cotes-Formel mit  $n+1$  Stützstellen die Abschätzung

$$\left| \int_a^b f(x) dx - \sum_{i=0}^n w_i f(x_i) \right| \leq C_n \|f^{(n+1)}\|_{C([a,b])} h^{n+2} \quad (18.7)$$

für den Quadraturfehler, wobei  $C_n$  in (18.6) definiert und  $h = \frac{b-a}{n}$  ist.

Durch die Berechnung der Konstanten  $C_n$  können wir konkrete Abschätzungen angeben:

**Beispiel 18.4** (Fehlerabschätzung der abgeschlossenen Newton-Cotes-Formeln). Für die Trapezregel ( $n=1$ , Genauigkeitsgrad 1) gilt

$$C_1 = \frac{1}{2!} \int_0^1 |(t-0)(t-1)| dt = \frac{1}{12}$$

und damit

$$\left| \int_a^b f(x) dx - \frac{b-a}{2} [f(a) + f(b)] \right| \leq \frac{1}{12} \|f''\|_{C([a,b])} (b-a)^3.$$

Für die Simpson-Regel ( $n=2$ , Genauigkeitsgrad 3) gilt

$$C_2 = \frac{1}{3!} \int_0^2 |(t-0)(t-1)(t-2)| dt = \frac{1}{12}$$

und damit

$$\left| \int_a^b f(x) dx - \frac{b-a}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] \right| \leq \frac{1}{12} \|f^{(3)}\|_{C([a,b])} \frac{(b-a)^4}{2^4} = \frac{1}{192} \|f^{(3)}\|_{C([a,b])} (b-a)^4.$$

Diese Abschätzung ist aber leider nicht optimal. Mit anderen Techniken als im Beweis von Satz 18.3 kann man tatsächlich auch

$$\left| \int_a^b f(x) dx - \frac{b-a}{6} \left[ f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right] \right| \leq \frac{1}{2880} \|f^{(4)}\|_{C([a,b])} (b-a)^5$$

zeigen.

$n$	Gewichte	Name
0	$\frac{b-a}{1} [1]$	Rechteckregel, Mittelpunktsregel
1	$\frac{b-a}{2} [1 \quad 1]$	
2	$\frac{b-a}{3} [2 \quad -1 \quad 2]$	
3	$\frac{b-a}{24} [11 \quad 1 \quad 1 \quad 11]$	
4	$\frac{b-a}{20} [11 \quad -14 \quad 26 \quad -14 \quad 11]$	
5	$\frac{b-a}{1440} [611 \quad -453 \quad 562 \quad 562 \quad -453 \quad 611]$	
6	$\frac{b-a}{945} [460 \quad -954 \quad 2196 \quad -2459 \quad 2196 \quad -954 \quad 460]$	

Tabelle 18.2: Offene Newton-Cotes-Formeln auf  $[a, b]$ .

Die Fehlerterme bei abgeschlossenen Newton-Cotes-Formeln sind – entsprechende Glattheit des Integranden vorausgesetzt – von der Art  $(b-a)^p$ , wobei  $p = n+2$  für ungerades  $n \in \mathbb{N}$  und  $p = n+3$  für gerades  $n \in \mathbb{N}$  gilt, also  $p = \text{Genauigkeitsgrad} + 2$ .

## § 18.2 OFFENE NEWTON-COTES-FORMELN

Die **offenen Newton-Cotes-Formeln** (englisch: *open Newton-Cotes formulas*) verwenden ebenfalls gleichverteilte Stützstellen im Intervall  $[a, b]$ , lassen dabei aber die Randpunkte  $a$  und  $b$  aus. Es gilt also

$$x_i := a + (i+1)h \quad \text{mit } h = \frac{b-a}{n+2}, \quad i = 0, \dots, n \quad (18.8)$$

mit  $n \geq 0$ .

Die Berechnung der Gewichte für beliebige Intervalle  $[a, b]$  erfolgt analog zu den abgeschlossenen Formeln. Die Gewichte der ersten offenen Newton-Cotes-Formeln sind in [Tabelle 18.2](#) angegeben. Auch hier summieren sie sich wieder zu  $b-a$ .

Für den Genauigkeitsgrad der offenen Newton-Cotes-Formeln gilt dieselbe Aussage für wie die abgeschlossenen Formeln:

**Satz 18.5** (Genauigkeitsgrad der offenen Newton-Cotes-Formeln). *Der Genauigkeitsgrad der offenen Newton-Cotes-Formeln mit  $n+1$  Stützstellen ist gleich  $n$  für ungerades  $n \in \mathbb{N}$  und gleich  $n+1$  für gerades  $n \in \mathbb{N}$ .*

*Beweis.* Der Beweis kann ähnlich wie der Beweis für [Satz 18.2](#) geführt werden. □

Die offenen Formeln haben den Nachteil, dass bereits ab  $n=2$  negative Gewichte auftreten. Daher wird in der Praxis von den offenen Newton-Cotes-Formeln nur die Rechteckregel ( $n=0$ ) überhaupt

verwendet. (**Quizfrage:** Welchen Genauigkeitsgrad hat diese?) Die Regel für  $n = 1$  hat bei höherem Aufwand (2  $f$ -Auswertungen) denselben Genauigkeitsgrad wie die Rechteckregel.

**Beispiel 18.6** (Fehlerabschätzung der offenen Newton-Cotes-Formeln). Für die Rechteckregel ( $n = 0$ , Genauigkeitsgrad 1) kann man die Fehlerabschätzung

$$\left| \int_a^b f(x) dx - (b-a)f\left(\frac{a+b}{2}\right) \right| \leq \frac{1}{24} \|f''\|_{C([a,b])} (b-a)^3.$$

zeigen.

**Bemerkung 18.7** (Einschließung des Integralwertes). In bestimmten Fällen gelingt unter Verwendung zweier verschiedener Quadraturformeln eine Einschließung des echten Integralwertes. Beispielsweise gilt für konvexe Integranden, also solche mit  $f''(x) \geq 0$  für alle  $x \in [a, b]$ , die Abschätzung

$$(b-a)f\left(\frac{a+b}{2}\right) \leq \int_a^b f(x) dx \leq \frac{b-a}{2} [f(a) + f(b)]. \quad (18.9)$$

Die mit der Trapezregel und der Mittelpunktsregel erhaltenen Näherungswerte schließen also für konvexe Funktionen den echten Wert ein.

**Quizfrage:** Was gilt für konkave Funktionen, also solche mit  $f''(x) \leq 0$  auf  $[a, b]$ ?

## § 19 SUMMIERTE QUADRATURFORMELN

Wir haben bei den Newton-Cotes-Formeln in [Abschnitt 18](#) gesehen, dass die Erhöhung der Polynomordnung keine gute Strategie ist, um den Quadraturfehler zu senken. (**Quizfrage:** Warum nochmal?) Stattdessen kann man aber das Intervall  $[a, b]$  zunächst in  $N$  Teilintervalle  $I_k = [a_k, a_{k+1}]$ ,  $k = 0, \dots, N-1$ , mit den Punkten  $a_k = a + kH$  und der Länge  $H = \frac{b-a}{N}$  aufteilen

$$\int_a^b f(x) dx = \sum_{k=0}^{N-1} \int_{a_k}^{a_{k+1}} f(x) dx$$

und die Quadratur auf jedem Teilintervall dann z. B. mit einer abgeschlossenen Newton-Cotes-Formel einzeln durchführen. Die Gewichte an den Übergangsstellen zwischen zwei Intervallen addieren (verdoppeln) sich dabei effektiv, siehe [Beispiel 19.1](#).

Erfüllt die verwendete Quadraturformel auf dem Intervall  $[a_k, a_{k+1}]$  eine Abschätzung für den Fehler (vgl. (18.7) und [Beispiele 18.4](#) und [18.6](#))

$$\left| \int_{a_k}^{a_{k+1}} f(x) dx - \sum_{i=0}^n w_i^{(k)} f(x_i^{(k)}) \right| \leq C \|f^{(p+1)}\|_{C([a_k, a_{k+1}])} h^{p+2} \quad \text{mit } h = \frac{a_{k+1} - a_k}{n},$$

so kann man für den Fehler der summierten Quadraturformel zeigen:

$$\left| \int_a^b f(x) dx - \sum_{k=0}^{N-1} \sum_{i=0}^n w_i^{(k)} f(x_i^{(k)}) \right| \leq C(b-a) \|f^{(p+1)}\|_{C([a,b])} H^{p+1}, \quad (19.1)$$

vgl. [Rannacher, 2017](#), Gleichung (3.1.5).

**Beispiel 19.1** (Summierte Quadraturformeln). *Die summierten Varianten der ersten beiden abgeschlossenen Newton-Cotes-Formeln (Trapezregel und Simpson-Regel, vgl. (18.5)) und der ersten offenen Newton-Cotes-Formel (Rechteckregel) sind:*

$$\begin{aligned} & \sum_{k=0}^{N-1} \frac{a_{k+1} - a_k}{2} [f(a_k) + f(a_{k+1})] \\ &= \frac{b-a}{2N} \left[ f(a) + 2 \sum_{k=1}^{N-1} f(a_k) + f(b) \right] \end{aligned} \quad \text{summierte Trapezregel} \quad (19.2a)$$

mit der Fehlerschranke  $\frac{b-a}{12} \|f''\|_{C([a,b])} \left( \frac{b-a}{N} \right)^2$ ,

$$\begin{aligned} & \sum_{k=1}^n \frac{a_{k+1} - a_k}{6} \left[ f(a_k) + 4 f\left(\frac{a_k + a_{k+1}}{2}\right) + f(a_{k+1}) \right] \\ &= \frac{b-a}{6N} \left[ f(a) + 2 \sum_{k=1}^{N-1} f(a_k) + 4 \sum_{k=0}^{N-1} f\left(\frac{a_k + a_{k+1}}{2}\right) + f(b) \right] \end{aligned} \quad \text{summierte Simpson-Regel} \quad (19.2b)$$

mit der Fehlerschranke  $\frac{b-a}{2880} \|f^{(4)}\|_{C([a,b])} \left( \frac{b-a}{N} \right)^4$ ,

$$\begin{aligned} & \sum_{k=0}^{N-1} (a_{k+1} - a_k) f\left(\frac{a_k + a_{k+1}}{2}\right) \\ &= \frac{b-a}{N} \sum_{k=0}^{N-1} f\left(\frac{a_k + a_{k+1}}{2}\right) \end{aligned} \quad \text{summierte Rechteckregel} \quad (19.2c)$$

mit der Fehlerschranke  $\frac{b-a}{24} \|f''\|_{C([a,b])} \left( \frac{b-a}{N} \right)^2$ .



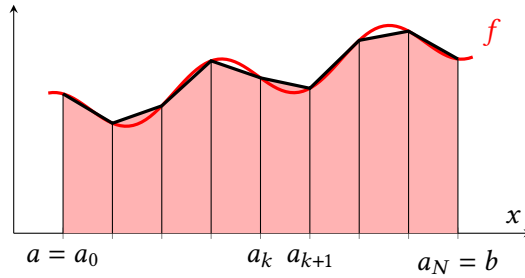


Abbildung 19.1: Veranschaulichung der summierten Trapezregel (19.2a).

## § 20 GAUSS-QUADRATURFORMELN

Die bisherigen Quadraturformeln basierten auf einer Gleichverteilung der Stützstellen. Wenn wir davon abweichen, können wir einen höheren Genauigkeitsgrad erzielen. Es wird sich herausstellen, dass wir für beliebiges  $n \in \mathbb{N}_0$  sogar die obere Schranke  $2n + 1$  aus Lemma 18.1 erreichen können. Es gibt genau eine Wahl der  $n + 1$  Stützstellen und Gewichte, mit der das realisiert wird. Das folgende Lemma charakterisiert diese Wahl mit Hilfe des zugehörigen Knotenpolynoms (15.10), also  $\ell(x) = \prod_{j=0}^n (x - x_j)$ .

**Lemma 20.1** (Bedingungen für das Erreichen des höchsten Genauigkeitsgrades). *Es sei  $n \in \mathbb{N}_0$ . Der Ansatz*

$$\sum_{i=0}^n w_i f(x_i) \quad \text{mit Gewichten } w_i := \int_a^b L_i^{(n)}(x) dx \quad (18.2)$$

*ist genau dann exakt für alle  $f \in P_{2n+1}$ , wenn für das Knotenpolynom  $\ell(x) = \prod_{j=0}^n (x - x_j) \in P_{n+1}$  gilt:*

$$\int_a^b z(x) \ell(x) dx = 0 \quad \text{für alle } z \in P_n. \quad (20.1)$$

Bevor wir dieses Lemma beweisen, weisen wir kurz darauf hin, dass der Ausdruck

$$\int_a^b p(x) \ell(x) dx$$

ein Innenprodukt<sup>2</sup>, genannt das  $L^2(a, b)$ -Innenprodukt (englisch:  $L^2(a, b)$ -inner product), auf dem Raum aller Polynome darstellt. Wir können die Bedingung (20.1) also so formulieren, dass das Knotenpolynom  $\ell$  im  $L^2(a, b)$ -Innenprodukt senkrecht auf allen Funktionen in  $P_n$  stehen muss – oder noch kürzer, dass  $\ell$  im  $L^2(a, b)$ -Innenprodukt senkrecht auf dem Raum  $P_n$  steht.

*Beweis.* Wir nehmen zunächst an, dass (20.1) gilt. Es sei  $f \in P_{2n+1}$  beliebig. Dann gibt es zwei Polynome  $z, p \in P_n$ , sodass

$$f(x) = z(x) \ell(x) + p(x)$$

<sup>2</sup>Genau wie das Euklidische Innenprodukt auf dem Vektorraum  $\mathbb{R}^n$  ist das  $L^2(a, b)$ -Innenprodukt symmetrisch, linear in beiden Funktionen, und es gilt  $\int_a^b p(x)^2 dx \geq 0$ , wobei der Wert 0 nur dann angenommen wird, wenn  $p$  das Nullpolynom ist.

gilt. Wir können  $z$  und  $p$  mit Hilfe einer Polynomdivision  $f/\ell$  mit Rest bestimmen. (**Quizfrage:** Klar?)  
Wir erhalten

$$\begin{aligned}
 \int_a^b f(x) \, dx &= \int_a^b z(x) \ell(x) + p(x) \, dx \\
 &= \int_a^b z(x) \ell(x) \, dx + \int_a^b p(x) \, dx \\
 &= \int_a^b p(x) \, dx \quad (\text{nach Voraussetzung (20.1)}) \\
 &= \sum_{i=0}^n p(x_i) \int_a^b L_i^{(n)}(x) \, dx \quad (\text{nach Darstellung von } p \in P_n \text{ mit Hilfe Lagrange-Basis}) \\
 &= \sum_{i=0}^n f(x_i) \int_a^b L_i^{(n)}(x) \, dx \quad (\text{wegen } f(x_i) = p(x_i)) \\
 &= \sum_{i=0}^n w_i f(x_i).
 \end{aligned}$$

Das heißt, dass alle Polynome  $f \in P_{2n+1}$  exakt integriert werden, was die eine Richtung des Beweises zeigt. (**Quizfrage:** Warum gilt eigentlich  $f(x_i) = p(x_i)$ ?)

Für die andere Richtung nehmen wir nun an, dass die Quadraturformel alle Polynome  $f \in P_{2n+1}$  exakt integriert. Es sei  $z \in P_n$  beliebig gewählt. Dann ist  $z\ell \in P_{2n+1}$  und wird exakt integriert. Das bedeutet:

$$\begin{aligned}
 \int_a^b z(x) \ell(x) \, dx &= \sum_{i=0}^n w_i z(x_i) \ell(x_i) \\
 &= \sum_{i=0}^n z(x_i) \ell(x_i) \int_a^b L_i^{(n)}(x) \, dx \\
 &= 0.
 \end{aligned}$$

Damit ist (20.1) gezeigt. (**Quizfrage:** Woraus folgt, dass sich im letzten Beweisschritt der Wert 0 ergibt?)

□

Das System (20.1) ist ein lineares Gleichungssystem für das Knotenpolynom, beschrieben durch die  $n+1$  unbekannten Stützstellen  $x_0, \dots, x_n$ . Es sieht zunächst so aus, als würde die Bedingung «für alle  $z \in P_n$ » eine unendliche Anzahl an Gleichungen erzeugen. Es ist jedoch ausreichend, diese Bedingung nur für  $n+1$  verschiedene Polynome  $z^{(0)}, \dots, z^{(n)}$  zu fordern, die eine Basis von  $P_n$  bilden. (**Quizfrage:** Klar?) Bei der Konstruktion dieser Basis ist folgende Überlegung hilfreich. Dabei ist jeweils  $n \in \mathbb{N}_0$ :

- $\ell$  ist ein Polynom in  $P_{n+1}$ , das im  $L^2(a, b)$ -Innenprodukt senkrecht auf allen Basiselementen  $z^{(0)}, \dots, z^{(n)}$  steht und das außerdem den führenden Koeffizienten 1 besitzt.

- Durch diese Bedingungen ist  $\ell \in P_{n+1}$  eindeutig festgelegt.
- Wenn wir die Basis  $z^{(0)}, \dots, z^{(n)}$  um  $\ell \in P_{n+1}$  als nächstes Basiselement  $z^{(n+1)}$  ergänzen, erhalten wir eine  $L^2(a, b)$ -orthogonale Folge von Polynomen, die bis zum  $n$ -ten Glied jeweils eine Basis von  $P_n$  bilden. Das gesuchte Knotenpolynom  $\ell \in P_{n+1}$  ist gerade das Element  $z^{(n+1)}$ .
- Die gesuchten Stützstellen  $x_0, \dots, x_n$  sind gerade die Nullstellen von  $z^{(n+1)}$ .

Zur Konstruktion der Basis  $z^{(n)}$  kann man z. B. das Gram-Schmidt-Verfahren auf die Monombasis  $\{1, x, x^2, \dots, x^n\}$  anwenden und anschließend das neue Folgenglied  $z^{(n+1)}$  so normieren, dass der führende Koeffizient gleich 1 wird.

Man kann zeigen, dass die Nullstellen von  $z^{(n+1)}$  alle verschieden sind, im Inneren des Intervalls  $[a, b]$  liegen und dass die Gewichte  $w_i = \int_a^b L_i^{(n)}(x) dx$  alle echt positiv sind. Die erzeugte Folge von Polynomen wachsenden Grades besteht gerade aus Vielfachen der **Legendre-Polynome** (englisch: *Legendre polynomials*). Man kann weiterhin zeigen (vgl. [Rannacher, 2017](#), Satz 3.5), dass die Werte der Quadraturformel für jeden stetigen Integranden gegen den wirklichen Integralwert konvergieren.

## § 21 QUADRATUR BASIEREND AUF RICHARDSON-EXTRAPOLATION

Wir betrachten abschließend eine weitere Idee, die Genauigkeit einer numerischen Approximation des Integrals  $\int_a^b f(x) dx$  zu erhöhen. Die Idee dazu ist allgemeiner Natur und wird auch bei anderen Aufgaben als der Quadratur eingesetzt. Nehmen wir an, ein numerisches Verfahren liefert uns für die eigentlich gesuchte Größe  $a^* \in \mathbb{R}$  einen Näherungswert  $a(h)$ , der einem Approximationsfehler unterliegt, wobei gilt:

$$a^* \in a(h) + C h^p + O(h^{p+1}). \quad (21.1)$$

Dabei ist  $h > 0$  ein Diskretisierungsparameter. Wenn wir dasselbe Verfahren nun zusätzlich auch beispielsweise für den Diskretisierungsparameter  $h/2$  durchführen, so gilt

$$a^* \in a(h/2) + C (h/2)^p + O(h^{p+1}) \quad (21.2)$$

und folglich

$$2^p a^* \in 2^p a(h/2) + C h^p + O(h^{p+1}). \quad (21.3)$$

Durch die Subtraktion von (21.1) und (21.3) können wir den führenden Fehlerterm  $C h^p$  eliminieren und erhalten die Aussage

$$a^* \in \frac{2^p a(h/2) - a(h)}{2^p - 1} + O(h^{p+1}). \quad (21.4)$$

Der Ausdruck

$$b(h) := \frac{2^p a(h/2) - a(h)}{2^p - 1} \quad (21.5)$$

definiert also ein neues Verfahren, dessen Fehlerterm um eine  $h$ -Potenz besser ist. Man bezeichnet dieses Vorgehen als **Richardson-Extrapolation** (englisch: *Richardson extrapolation*).

**Bemerkung 21.1.** Die gerade beschriebene Idee erhält man auch dadurch, dass man die Stützpaare  $(h^p, a(h))$  und  $((h/2)^p, a(h/2))$  durch ein lineares Interpolationspolynom verbindet und dieses an der Stelle 0 auswertet. Dieses Interpolationspolynom lautet

$$p(x) = a(h) + \frac{a(h) - a(h/2)}{h^p - (h/2)^p} (x - h^p).$$

Die Auswertung an der Stelle  $x = 0$  liefert also

$$\begin{aligned} p(0) &= a(h) + \frac{a(h) - a(h/2)}{h^p - (h/2)^p} (0 - h^p) \\ &= a(h) - h^p \frac{a(h) - a(h/2)}{h^p(1 - 2^{-p})} \\ &= a(h) - \frac{a(h) - a(h/2)}{1 - 2^{-p}} \\ &= a(h) - 2^p \frac{a(h) - a(h/2)}{2^p - 1} \\ &= \frac{2^p a(h/2) - a(h)}{2^p - 1}, \end{aligned}$$

was mit (21.5) übereinstimmt.

Die Richardson-Extrapolation lässt sich auch dann noch anwenden, wenn man neben den Diskretisierungsparametern  $h$  und  $h/2$  noch  $h/4$ ,  $h/8$  (die **Romberg-Folge**, englisch: *Romberg sequence*) usw. verwendet, um die Genauigkeit weiter zu erhöhen. Im Prinzip können wir auch beliebige andere Schrittweiten verwenden, die Romberg-Folge hat aber einen bestimmten Vorteil. (**Quizfrage:** Welcher Vorteil könnte das sein?)

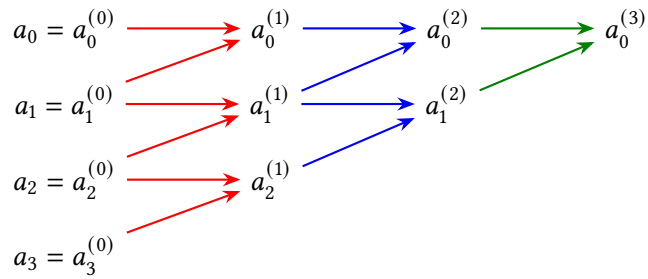
Nun noch zur Frage, wie man die Auswertung des Interpolationspolynoms an der Stelle 0 am besten macht. Wir könnten dazu das Interpolationspolynom z. B. in der Newtonschen Form aufstellen. Da wir aber nur an dessen Auswertung an der Stelle 0 interessiert ist, gibt es eine effizientere Form, das **Verfahren von Aitken-Neville** (englisch: *Aitken-Neville algorithm*).<sup>3</sup> Es bestimmt Größen analog zu den Newtonschen dividierten Differenzen, wobei jedoch an Stelle der Koeffizienten des Polynoms direkt der interessierende Funktionswert steht.

Gegeben die Fehlerordnung  $p \in \mathbb{N}$  und die Größen  $a_i^{(0)} := a(h_i)$ ,  $i = 0, \dots, n$ , die die Näherungswerte zu den Diskretisierungsparametern  $h_i$  sind, bestimmt das Verfahren sukzessive für  $\ell = 0, 1, \dots, n-1$  die Größen

$$a_i^{(\ell+1)} := \frac{h_{i+(\ell+1)}^p a_i^{(\ell)} - h_i^p a_{i+1}^{(\ell)}}{h_{i+(\ell+1)}^p - h_i^p} = a_i^{(\ell)} + \frac{a_i^{(\ell)} - a_{i+1}^{(\ell)}}{\left(\frac{h_{i+(\ell+1)}}{h_i}\right)^p - 1} \quad i = 0, \dots, n - (\ell + 1).$$

Die finale Größe  $a_0^{(n)}$  ist dann der gesuchte Funktionswert.

<sup>3</sup>Dieses ist dann empfehlenswert, wenn man von vornherein weiß, dass man nicht am Interpolationspolynom selbst, sondern nur an dessen Auswertung an einigen wenigen Stellen interessiert ist.



Um die Fehlerordnung eines mit Hilfe der Richardson-Extrapolation konstruierten Quadraturverfahrens analysieren zu können, benötigt man eine Darstellung des Fehlers in Potenzen von  $h$ . Eine solche liefert die **Euler-Maclaurinsche Summenformel**, die wir hier aber nicht besprechen.

Speziell als **Romberg-Quadratur** (englisch: *Romberg quadrature*) bezeichnet man das Verfahren bei Verwendung der Trapezregel und der Romberg-Folge  $h, \frac{h}{2}, \frac{h}{4}, \dots$ , an Diskretisierungsparametern.

Ende der Woche 11

# Kapitel 8 Iterative Lösungsverfahren für nichtlineare Gleichungssysteme

In diesem Kapitel beschäftigen wir uns mit Verfahren zur Lösung nichtlinearer Gleichungen. Diese können in der Form  $F(x) = 0$  geschrieben werden, wobei  $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$  eine Funktion ist. **Quizfrage:** Können alle nichtlinearen Gleichungen in dieser Form geschrieben werden?

## § 22 BISEKTIONSVERFAHREN

Zur Charakterisierung der Konvergenzgeschwindigkeit von Algorithmen führen wir folgende Begriffe ein:

**Definition 22.1** (Q-Konvergenzraten).

Es sei  $(x^{(k)}) \subseteq \mathbb{R}^n$  eine Folge und  $x^* \in \mathbb{R}^n$ .

(i) Die Folge  $(x^{(k)})$  konvergiert gegen  $x^*$  (mindestens) **Q-linear**, falls ein  $c \in (0, 1)$  existiert mit

$$\|x^{(k+1)} - x^*\|_2 \leq c \|x^{(k)} - x^*\|_2 \quad \text{für alle } k \in \mathbb{N} \text{ hinreichend groß.}$$

(ii) Die Folge  $(x^{(k)})$  konvergiert gegen  $x^*$  (mindestens) **Q-superlinear**, falls es eine Nullfolge  $\varepsilon^{(k)}$  gibt mit

$$\|x^{(k+1)} - x^*\|_2 \leq \varepsilon^{(k)} \|x^{(k)} - x^*\|_2 \quad \text{für alle } k \in \mathbb{N}.$$

(iii) Es gelte  $x^{(k)} \rightarrow x^*$ . Die Folge  $(x^{(k)})$  konvergiert gegen  $x^*$  (mindestens) **Q-quadratisch**, falls ein  $C > 0$  existiert mit

$$\|x^{(k+1)} - x^*\|_2 \leq C \|x^{(k)} - x^*\|_2^2 \quad \text{für alle } k \in \mathbb{N}.$$

Wir betrachten zunächst ein einfaches Verfahren im Eindimensionalen, also zur Lösung von  $f(x) = 0$  für eine stetige Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$ . Voraussetzung ist, dass zwei Punkte  $a^{(0)} < b^{(0)} \in \mathbb{R}$  bekannt sind, für die  $f(a^{(0)})$  und  $f(b^{(0)})$  verschiedenes Vorzeichen haben, also  $f(a^{(0)}) f(b^{(0)}) < 0$ . Dann liegt im Intervall  $[a^{(0)}, b^{(0)}]$  sicher eine Nullstelle von  $f$ . **Quizfrage:** Welchen mathematischen Satz für stetige Funktionen nutzt man hier aus?

Die Idee des **Bisektionsverfahrens** (englisch: *bisection method*) ist, den Funktionswert in der Intervallmitte  $(a^{(0)} + b^{(0)})/2$  auszuwerten und dann zu entscheiden, ob man mit dem linken oder dem

rechten Teilintervall als  $[a^{(1)}, b^{(1)}]$  weitermacht usw. Die Intervallmitte dient jeweils als aktuell beste Schätzung für die Nullstelle. Wir beenden das Verfahren, sobald wir höchstens noch  $\varepsilon$  von einer Nullstelle entfernt sind. Das führt zu folgendem Verfahren.

**Algorithmus 22.2** (Bisektionsverfahren für Funktionen  $f: \mathbb{R} \rightarrow \mathbb{R}$ ).

**Eingabe:** Startwerte  $a^{(0)} < b^{(0)}$ , sodass  $f(a^{(0)}) f(b^{(0)}) < 0$

**Eingabe:** Toleranz  $\varepsilon > 0$

**Ausgabe:** Näherung  $x$ , in deren  $\varepsilon$ -Umgebung sich eine Nullstelle von  $f$  befindet,

```

1: Setze  $k := 0$ 
2: repeat
3:   Setze  $x^{(k)} := (a^{(k)} + b^{(k)})/2$ 
4:   if  $f(a^{(k)}) f(x^{(k)}) < 0$  then                                // im linken Teilintervall liegt eine Nullstelle
5:     Setze  $a^{(k+1)} := a^{(k)}$ 
6:     Setze  $b^{(k+1)} := x^{(k)}$ 
7:   else                                                            // im rechten Teilintervall liegt eine Nullstelle
8:     Setze  $a^{(k+1)} := x^{(k)}$ 
9:     Setze  $b^{(k+1)} := b^{(k)}$ 
10:  end if
11:  Setze  $k := k + 1$ 
12: until  $b^{(k)} - a^{(k)} \leq \varepsilon$  or  $f(x^{(k)}) = 0$ 
13: return  $x^{(k)}$ 

```

Die Mengen  $[a^{(k)}, b^{(k)}]$  bilden eine Folge geschachtelter Intervalle, deren Länge sich von Iteration zu Iteration halbiert. Die Folge ihrer Mittelpunkte  $(x^{(k)})$  konvergiert daher gegen eine Nullstelle  $x^*$  von  $f$ , und es gilt

$$\|x^{(k)} - x^*\|_2 \leq \frac{b^{(k)} - a^{(k)}}{2} = \frac{b^{(0)} - a^{(0)}}{2^{k+1}} \rightarrow 0 \quad \text{für } k \rightarrow \infty. \quad (22.1)$$

Die Konvergenz ist Q-linear mit  $c = \frac{1}{2}$ , denn es gilt

$$\|x^{(k+1)} - x^*\|_2 \leq \frac{1}{2} \|x^{(k)} - x^*\|_2. \quad (22.2)$$

für alle  $k \in \mathbb{N}_0$ . (**Quizfrage:** Warum?)

## § 23 FIXPUNKTVERFAHREN

Es sei  $M \subset \mathbb{R}^n$  eine nichtleere, abgeschlossene Menge und  $T: M \rightarrow M$  eine **Kontraktion** (englisch: *contraction*), d. h. eine Lipschitz-stetige Abbildung

$$\|T(x) - T(y)\|_2 \leq L \|x - y\|_2 \quad \text{für alle } x, y \in M \quad (23.1)$$

mit Lipschitz-Konstante  $L < 1$ . Ein  $x \in M$  heißt **Fixpunkt** (englisch: *fixed point*) von  $T$ , wenn  $T(x) = x$  gilt. (**Quizfrage:** Was haben Fixpunkte mit Nullstellen zu tun?)

Unter einer **Fixpunkt-Iteration** versteht man die Iterationsvorschrift

$$x^{(k+1)} := T(x^{(k)}), \quad k = 0, 1, \dots \quad (23.2)$$

**Satz 23.1** (Banachscher Fixpunktsatz). *Unter den oben stehenden Voraussetzungen gilt:*

- (i) Die Funktion  $T$  besitzt in  $M$  genau einen Fixpunkt  $x^*$ .
- (ii) Für jeden beliebigen Startpunkt  $x^{(0)} \in M$  ist die durch (23.2) gegebene Folge wohldefiniert und konvergiert gegen  $x^*$ .
- (iii) Es gilt die **A-priori-Fehlerabschätzung**

$$\|x^{(k)} - x^*\|_2 \leq \frac{L^k}{1-L} \|x^{(1)} - x^{(0)}\|_2. \quad (23.3)$$

- (iv) Es gilt die **A-posteriori-Fehlerabschätzung**

$$\|x^{(k)} - x^*\|_2 \leq \frac{L}{1-L} \|x^{(k)} - x^{(k-1)}\|_2. \quad (23.4)$$

- (v) Es gilt die Abschätzung

$$\|x^{(k+1)} - x^*\|_2 \leq L \|x^{(k)} - x^*\|_2. \quad (23.5)$$

Man nennt (23.3) eine Abschätzung a-priori («im Vorhinein»), weil sie den Fehler  $\|x^{(k)} - x^*\|_2$  der  $k$ -ten Iterierten abschätzt, ohne dass man  $x^{(k)}$  dazu ausrechnen müsste. Im Gegensatz dazu ist (23.4) eine Abschätzung a-posteriori («im Nachhinein»), weil der Fehler  $\|x^{(k)} - x^*\|_2$  der  $k$ -ten Iterierten abgeschätzt wird, nachdem man  $x^{(k)}$  bereits bestimmt hat. Die Abschätzung (23.5) schließlich besagt, dass die Folge  $(x^{(k)})$  **Q-linear** in der  $\|\cdot\|_2$ -Norm gegen den Grenzwert  $x^*$  konvergiert.

*Beweis.* Zunächst zur Eindeutigkeit des Fixpunktes. Nehmen wir an,  $x^*$  und  $x^{**}$  seien zwei Fixpunkte von  $T$  in  $M$ . Dann gilt wegen (23.1)

$$\|T(x^*) - T(x^{**})\|_2 \leq L \|x^* - x^{**}\|_2 = L \|T(x^*) - T(x^{**})\|_2.$$

Die Annahme  $x^* \neq x^{**}$ , also  $T(x^*) \neq T(x^{**})$ , führt wegen  $L \in [0, 1)$  sofort zum Widerspruch.

Die Wohldefiniertheit der Folge  $(x^{(k)})$  folgt sofort daraus, dass  $T$  die Menge  $M$  wieder in  $M$  abbildet. Durch wiederholte Anwendung der Abschätzung (23.1) erhalten wir sofort

$$\|x^{(k+1)} - x^{(k)}\|_2 \leq L^k \|x^{(1)} - x^{(0)}\|_2. \quad (23.6)$$

Diese Abschätzung wird im Folgenden noch genutzt.



Um die Konvergenz der Folge zu sehen, zeigen wir ihre Cauchy-Eigenschaft. Es seien dazu  $k, m \in \mathbb{N}$ . Aus der Dreiecksungleichung folgt

$$\|x^{(k+m)} - x^{(k)}\|_2 \leq \sum_{\ell=0}^{m-1} \|x^{(k+\ell+1)} - x^{(k+\ell)}\|_2$$

Durch Anwendung der Abschätzung (23.6) erhalten wir weiter

$$\begin{aligned} \|x^{(k+m)} - x^{(k)}\|_2 &\leq \sum_{\ell=0}^{m-1} L^{k+\ell} \|x^{(1)} - x^{(0)}\|_2 = L^k \sum_{\ell=0}^{m-1} L^\ell \|x^{(1)} - x^{(0)}\|_2 \\ &\leq L^k \frac{1}{1-L} \|x^{(1)} - x^{(0)}\|_2 \quad (\text{geometrische Reihe}). \end{aligned} \quad (23.7)$$

Daraus folgt, dass wir für jedes  $\varepsilon > 0$  einen Index  $k(\varepsilon) \in \mathbb{N}$  finden, sodass  $\|x^{(k+m)} - x^{(k)}\|_2 \leq \varepsilon$  für alle  $k \geq k(\varepsilon)$  und alle  $m \in \mathbb{N}$  ist. Das bedeutet aber, dass  $(x^{(k)})$  eine Cauchy-Folge in der  $\|\cdot\|_2$ -Norm ist. Daraus folgt, dass die Folge in  $\mathbb{R}^n$  konvergiert, und weil alle Folgenglieder in der abgeschlossenen Menge  $M$  liegen, gehört auch der Grenzwert  $x^*$  zu  $M$ . Schließlich folgt aus der (Lipschitz-)stetigkeit von  $T$ :

$$x^* = \lim_{k \rightarrow \infty} x^{(k)} = \lim_{k \rightarrow \infty} T(x^{(k-1)}) = T(x^*),$$

d. h., der Grenzwert  $x^*$  ist der bereits als eindeutig erkannte Fixpunkt von  $T$  in  $M$ .

Bis hierher sind die Aussagen (i) und (ii) gezeigt. Aussage (iii), also die A-priori-Fehlerabschätzung (23.3), folgt sofort, wenn wir in (23.7)  $m \rightarrow \infty$  gehen lassen. (**Quizfrage:** Wie funktioniert dieses Argument genau?)

Ähnlich zu (23.7) können wir

$$\|x^{(k+m)} - x^{(k)}\|_2 \leq \sum_{\ell=0}^{m-1} L^{\ell+1} \|x^{(k)} - x^{(k-1)}\|_2 \leq \frac{1}{1-L} \|x^{(k)} - x^{(k-1)}\|_2$$

zeigen. Der Grenzübergang  $m \rightarrow \infty$  ergibt nun (23.4), also die Aussage (iv). Schließlich gilt auch

$$\|x^{(k+m)} - x^{(k)}\|_2 \leq L \|x^{(k+m-1)} - x^{(k-1)}\|_2,$$

und auch hier zeigt der Grenzübergang  $m \rightarrow \infty$  die gewünschte Abschätzung (23.5), also Aussage (v).  $\square$

**Quizfrage:** Ist die Wahl der  $\|\cdot\|_2$ -Norm im Banachschen Fixpunktsatz wesentlich, oder könnten wir auch eine andere Norm im  $\mathbb{R}^n$  verwenden?

## § 24 NEWTON-VERFAHREN

Wir untersuchen in diesem Abschnitt das Newton-Verfahren zur Lösung von  $F(x) = 0$ , wobei  $F$  als  $C^1$ -Funktion vorausgesetzt wird.

**Idee:** Es sei  $x^{(0)}$  die Schätzung einer Nullstelle von  $F$ . Wir legen im Punkt  $x^{(0)}$  die Tangente an die Funktion und bestimmen *deren* Nullstelle:

$$F(x^{(0)}) + F'(x^{(0)})(x - x^{(0)}) = 0 \quad \Leftrightarrow \quad x = x^{(0)} - F'(x^{(0)})^{-1}F(x^{(0)}).$$

Diese Nullstelle  $x$  dient als nächste Iterierte  $x^{(1)}$  usw.

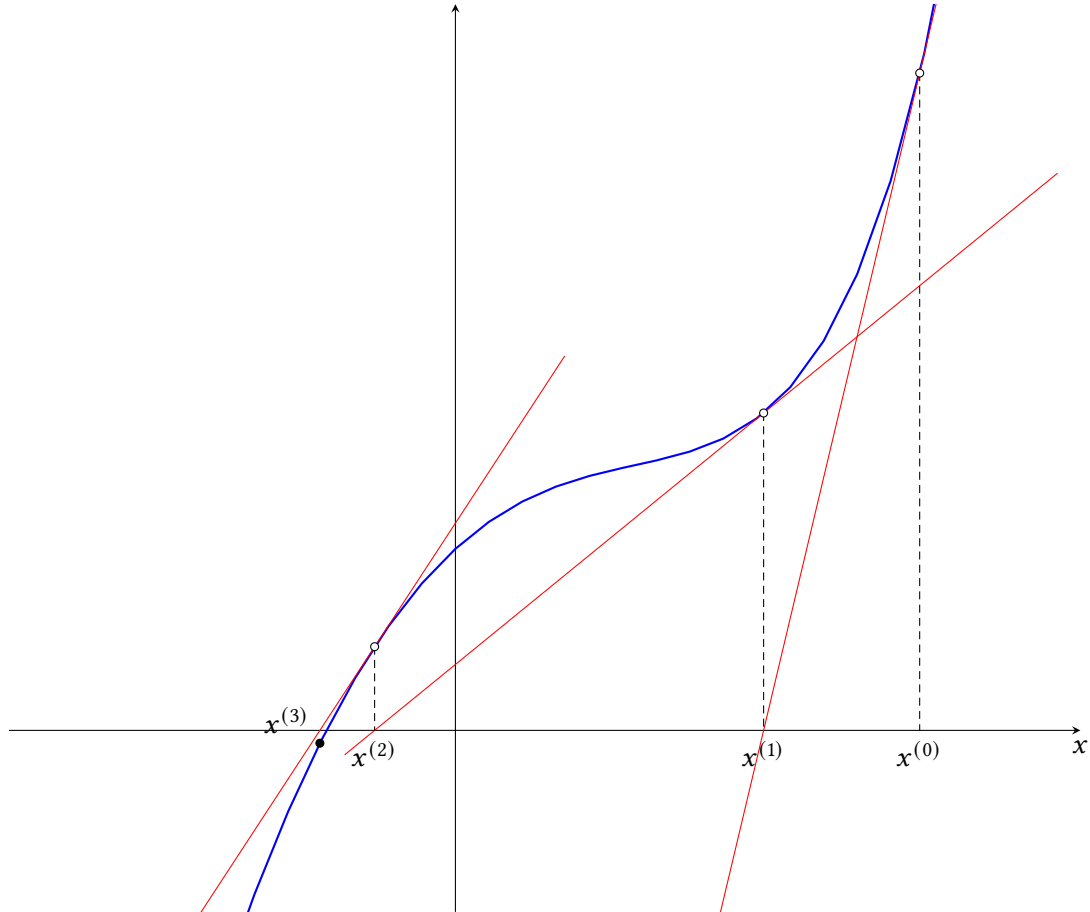


Abbildung 24.1: Illustration des Newton-Verfahrens zur Suche einer Nullstelle der Funktion  $F(x) = \exp(0.9x) - x^2$ .

Der Vektor  $F(x^{(k)})$  heißt dabei das **Residuum** zur Iterierten  $x^{(k)}$ , und  $F'(x^{(k)})$  ist die zugehörige Jacobimatrix:

$$F'(x) = \begin{pmatrix} \frac{\partial F_1(x)}{\partial x_1} & \dots & \frac{\partial F_1(x)}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial F_n(x)}{\partial x_1} & \dots & \frac{\partial F_n(x)}{\partial x_n} \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

**Algorithmus 24.1** (Lokales Newton-Verfahren).

**Eingabe:** Startschätzung  $x^{(0)} \in \mathbb{R}^n$

1: Setze  $k := 0$

```

2: while Abbruchkriterium nicht erfüllt do
3:   Löse das lineare Gleichungssystem  $F'(x^{(k)}) d^{(k)} = -F(x^{(k)})$  für die Newton-Richtung  $d^{(k)}$ 
4:   Setze  $x^{(k+1)} := x^{(k)} + d^{(k)}$ 
5:   Setze  $k := k + 1$ 
6: end while
    
```

Wir wollen die Konvergenzeigenschaften des Newton-Verfahrens (Algorithmus 24.1) untersuchen. Da wir in jeder Iteration in Zeile 3 ein lineares Gleichungssystem zu lösen haben, wird sich ein Konvergenzbeweis auch mit der Frage beschäftigen müssen, ob dieses überhaupt eindeutig lösbar ist. Daher folgen zunächst einige Resultate über die Invertierbarkeit von Matrizen.

## § 24.1 EINIGE HILFSRESULTATE

**Literatur:** Geiger, Kanzow, 1999, Kapitel 7, Lemma B.7 und B.8

Wir erinnern an die Operatornorm  $\|A\|_{2 \rightarrow 2}$  für Matrizen, siehe (2.3b).

**Lemma 24.2** (Banach-Lemma).

(i) Es sei  $M \in \mathbb{R}^{n \times n}$  mit  $\|M\|_{2 \rightarrow 2} < 1$ . Dann ist  $\text{Id} - M$  regulär (invertierbar), und es gilt

$$\|(\text{Id} - M)_{2 \rightarrow 2}^{-1}\| \leq \frac{1}{1 - \|M\|_{2 \rightarrow 2}}.$$

(ii) Es seien  $A, B \in \mathbb{R}^{n \times n}$  mit  $\|\text{Id} - BA\|_{2 \rightarrow 2} < 1$ . Dann sind  $A$  und  $B$  regulär, und es gilt

$$\|B^{-1}\|_{2 \rightarrow 2} \leq \frac{\|A\|_{2 \rightarrow 2}}{1 - \|\text{Id} - BA\|_{2 \rightarrow 2}} \quad \text{und} \quad \|A_{2 \rightarrow 2}^{-1}\| \leq \frac{\|B\|_{2 \rightarrow 2}}{1 - \|\text{Id} - BA\|_{2 \rightarrow 2}}.$$

Aussage (i) besagt, Matrizen «in der Nähe» der Einheitsmatrix invertierbar sind. Aussage (ii) besagt, dass wenn  $\text{Id} - BA$  klein ist, also  $B \approx A^{-1}$  gilt, notwendig  $A$  und  $B$  invertierbar sind.

**Beweis.** Aussage (i): Für  $x \in \mathbb{R}^n$  gilt

$$\|(\text{Id} - M)x\|_2 = \|x - Mx\|_2 \geq \|x\|_2 - \|Mx\|_2 \geq \underbrace{(1 - \|M\|_{2 \rightarrow 2})}_{>0} \|x\|_2.$$

Es folgt  $(\text{Id} - M)x \neq 0$  für  $x \neq 0$ , d. h.,  $\text{Id} - M$  ist injektiv und damit regulär.

Es sei nun  $y \in \mathbb{R}^n$  beliebig und  $x := (\text{Id} - M)^{-1}y$ . Für eine Abschätzung der Norm von  $(\text{Id} - M)^{-1}$  müssen wir  $\|x\|_2$  durch  $\|y\|_2$  abschätzen. Die Abschätzung oben zeigt

$$\begin{aligned} \|y\|_2 &\geq (1 - \|M\|_{2 \rightarrow 2}) \|x\|_2 \\ \Rightarrow \quad \|(\text{Id} - M)^{-1}\|_{2 \rightarrow 2} &= \max_{y \neq 0} \frac{\|(\text{Id} - M)^{-1}y\|_2}{\|y\|_2} \leq \frac{1}{1 - \|M\|_{2 \rightarrow 2}}. \end{aligned}$$

**Aussage (ii):** Es sei  $M = \text{Id} - BA$ , also  $\|M\|_{2 \rightarrow 2} < 1$ . Wegen **Aussage (i)** ist  $\text{Id} - M = \text{Id} - (\text{Id} - BA) = BA$  regulär, d. h.,  $A$  und  $B$  sind beide regulär. Weiter gilt

$$\begin{aligned} (\text{Id} - M)^{-1} &= (BA)^{-1} = A^{-1}B^{-1} \\ \Rightarrow B^{-1} &= A(\text{Id} - M)^{-1} \\ \Rightarrow \|B^{-1}\|_{2 \rightarrow 2} &\leq \|A\|_{2 \rightarrow 2} \|(\text{Id} - M)^{-1}\|_{2 \rightarrow 2} \stackrel{(i)}{\leq} \frac{\|A\|_{2 \rightarrow 2}}{1 - \|M\|_{2 \rightarrow 2}} = \frac{\|A\|_{2 \rightarrow 2}}{1 - \|\text{Id} - BA\|_{2 \rightarrow 2}}. \end{aligned}$$

Die andere Ungleichung folgt analog. □

Wir zeigen jetzt mit Hilfe des **Banach-Lemmas 24.2**, dass alle Jacobimatrizen «in der Nähe» einer invertierbaren Jacobimatrix ebenfalls invertierbar sind.

**Lemma 24.3.** *Es sei  $F$  eine  $C^1$ -Funktion,  $x^* \in \mathbb{R}^n$  und die Jacobimatrix  $F'(x^*)$  regulär. Dann existieren eine Umgebung  $B_\delta(x^*)$  und eine Konstante  $c > 0$ , sodass  $F'(x)$  für alle  $x \in B_\delta(x^*)$  regulär ist, und es gilt:*

$$\|F'(x)^{-1}\|_{2 \rightarrow 2} \leq c := 2 \|F'(x^*)^{-1}\|_{2 \rightarrow 2} \quad \text{für alle } x \in B_\delta(x^*).$$

*Beweis.* Da  $F'$  im Punkt  $x^*$  stetig ist, existiert ein  $\delta > 0$  mit

$$\|F'(x^*) - F'(x)\|_2 \leq \varepsilon = \frac{1}{2 \|F'(x^*)^{-1}\|_{2 \rightarrow 2}}$$

für alle  $x \in B_\delta(x^*)$ , also auch

$$\begin{aligned} \|\text{Id} - F'(x^*)^{-1} F'(x)\|_{2 \rightarrow 2} &= \|F'(x^*)^{-1} (F'(x^*) - F'(x))\|_{2 \rightarrow 2} \\ &\leq \|F'(x^*)^{-1}\|_{2 \rightarrow 2} \|F'(x^*) - F'(x)\|_{2 \rightarrow 2} \\ &\leq 1/2 < 1. \end{aligned}$$

Nach dem **Banach-Lemma 24.2**, **Aussage (ii)** [mit  $A = F'(x)$  und  $B = F'(x^*)^{-1}$ ] folgt, dass  $F'(x)$  für  $x \in B_\delta(x^*)$  regulär ist, und es gilt

$$\|F'(x)^{-1}\|_{2 \rightarrow 2} \leq \frac{\|F'(x^*)^{-1}\|_{2 \rightarrow 2}}{1 - \|\text{Id} - F'(x^*)^{-1} F'(x)\|_{2 \rightarrow 2}} \leq 2 \|F'(x^*)^{-1}\|_{2 \rightarrow 2} =: c.$$

□

Es folgt schließlich noch eine weitere Abschätzung.

**Lemma 24.4.** *Es sei  $F$  eine  $C^1$ -Funktion und  $x^* \in \mathbb{R}^n$ . Für alle  $\varepsilon > 0$  existiert  $\delta > 0$  mit*

$$\|F(x) - F(x^*) - F'(x^*)(x - x^*)\|_2 \leq \varepsilon \|x - x^*\|_2$$

*für alle  $\|x - x^*\|_2 \leq \delta$ .*

**Quizfrage:** Was würde die Aussage des Satzes bedeuten, wenn an Stelle von  $x$  der Punkt  $x^*$  stehen würde?

*Beweis.* Es sei  $\varepsilon > 0$  gegeben. Aus der Dreiecksungleichung ergibt sich

$$\begin{aligned} & \|F(x) - F(x^*) - F'(x)(x - x^*)\|_2 \\ & \leq \|F(x) - F(x^*) - F'(x^*)(x - x^*)\|_2 + \|F'(x^*) - F'(x)\|_{2 \rightarrow 2} \|x - x^*\|_2. \end{aligned}$$

Da  $F$  nach Voraussetzung in  $x^*$  diffbar ist, existiert  $\delta_1 > 0$  mit

$$\|F(x) - F(x^*) - F'(x^*)(x - x^*)\|_2 \leq \frac{\varepsilon}{2} \|x - x^*\|_2$$

für alle  $\|x - x^*\|_2 < \delta_1$ . Andererseits ist  $F'$  stetig in  $x^*$ , sodass  $\delta_2 > 0$  existiert mit

$$\|F'(x^*) - F'(x)\|_{2 \rightarrow 2} \leq \frac{\varepsilon}{2}$$

für alle  $\|x - x^*\|_2 < \delta_2$ . Mit  $\delta := \min\{\delta_1, \delta_2\}$  folgt die Behauptung.  $\square$

## § 24.2 DAS LOKALE NEWTON-VERFAHREN

Wir können nun einen lokalen Konvergenzsatz für [Algorithmus 24.1](#) (ohne Abbruchbedingung) beweisen. «Lokal» bedeutet dabei, dass es darauf ankommt, dass wir hinreichend nahe an einer Lösung starten.

**Satz 24.5** (Lokaler Konvergenzsatz für Newton-Verfahren).

Es sei  $F$  eine  $C^1$ -Funktion und  $x^* \in \mathbb{R}^n$  ein Punkt mit  $F(x^*) = 0$  und  $F'(x^*)$  regulär. Dann existiert eine Umgebung  $B_\delta(x^*)$  von  $x^*$ , sodass für jedes  $x^{(0)} \in B_\delta(x^*)$  gilt:

- (i) Das lokale Newton-Verfahren ist wohldefiniert und erzeugt eine Folge  $x^{(k)}$ , die gegen  $x^*$  konvergiert.
- (ii) Die Konvergenzrate ist  $Q$ -superlinear.
- (iii) Ist  $F'$  Lipschitz-stetig in  $B_\delta(x^*)$ , dann ist die Konvergenzrate sogar  $Q$ -quadratisch.

*Beweis.* **Aussage (i):** Nach [Lemma 24.3](#) existieren  $\delta_1 > 0$  und  $c > 0$ , sodass  $F'(x)$  für alle  $x \in B_{\delta_1}(x^*)$  regulär ist mit

$$\|F'(x)^{-1}\|_{2 \rightarrow 2} \leq c = 2 \|F'(x^*)^{-1}\|_{2 \rightarrow 2}. \quad (24.1)$$

Nach [Lemma 24.4](#) existiert zu  $\varepsilon = 1/(2c)$  ein  $\delta_2 > 0$  mit

$$\|F(x) - F(x^*) - F'(x)(x - x^*)\|_2 \leq \frac{1}{2c} \|x - x^*\|_2 \quad (24.2)$$

für alle  $x \in B_{\delta_2}(x^*)$ . Setze  $\delta := \min\{\delta_1, \delta_2\}$  und wähle  $x^{(0)} \in B_{\delta}(x^*)$ . Dann ist der Schritt  $x^{(1)} := x^{(0)} - F'(x^{(0)})^{-1}F(x^{(0)})$  wohldefiniert, und es gilt

$$\begin{aligned} \|x^{(1)} - x^*\|_2 &= \|x^{(0)} - x^* - F'(x^{(0)})^{-1}F(x^{(0)})\|_2 \\ &= \|F'(x^{(0)})^{-1} [F'(x^{(0)})(x^{(0)} - x^*) - F(x^{(0)}) + F(x^*)]\|_{2 \rightarrow 2} \\ &\leq \|F'(x^{(0)})^{-1}\|_{2 \rightarrow 2} \|F(x^{(0)}) - F(x^*) - F'(x^{(0)})(x^{(0)} - x^*)\|_2 \\ &\leq c \frac{1}{2c} \|x^{(0)} - x^*\|_2 = \frac{1}{2} \|x^{(0)} - x^*\|_2, \end{aligned}$$

also liegt auch  $x^{(1)}$  wieder in  $B_{\delta}(x^*)$ . Per Induktion ist  $x^{(k)}$  wohldefiniert, gehört zu  $B_{\delta}(x^*)$ , und  $x^{(k)} \rightarrow x^*$  zunächst Q-linear.

**Aussage (ii):** Wir stellen zunächst eine Gleichung für den Fehler auf:<sup>1</sup>

$$\begin{aligned} x^{(k+1)} - x^* &= x^{(k)} - x^* - F'(x^{(k)})^{-1}(F(x^{(k)}) - F(x^*)) \\ &= F'(x^{(k)})^{-1} [F'(x^{(k)})(x^{(k)} - x^*) - (F(x^{(k)}) - F(x^*))] \\ &= F'(x^{(k)})^{-1} \left[ F'(x^{(k)})(x^{(k)} - x^*) - \int_0^1 F'(x^{(k)} + t(x^* - x^{(k)}))(x^{(k)} - x^*) dt \right] \\ &= F'(x^{(k)})^{-1} \left[ \int_0^1 F'(x^{(k)}) - F'(x^{(k)} + t(x^* - x^{(k)})) dt \right] (x^{(k)} - x^*). \end{aligned}$$

Daraus erhalten wir folgende wichtige Abschätzung:

$$\begin{aligned} \|x^{(k+1)} - x^*\|_2 &\leq \|F'(x^{(k)})^{-1}\|_{2 \rightarrow 2} \int_0^1 \overbrace{\|F'(x^{(k)}) - F'(x^{(k)} + t(x^* - x^{(k)}))\|_{2 \rightarrow 2}}^{=: D^{(k)}(t)} dt \|x^{(k)} - x^*\|_2. \quad (24.3) \end{aligned}$$

Wegen  $x^{(k)} \rightarrow x^*$  gilt  $x^{(k)} + t(x^* - x^{(k)}) \rightarrow x^*$  gleichmäßig auf  $t \in [0, 1]$ . Außerdem ist  $F'$  stetig. Zu jedem  $\varepsilon > 0$  existiert also ein Index  $k_0 \in \mathbb{N}$  mit

$$\begin{aligned} \|D^{(k)}(t)\|_{2 \rightarrow 2} &\leq \varepsilon \quad \text{für alle } k \geq k_0 \text{ und alle } t \in [0, 1]. \\ \Rightarrow 0 &\leq \int_0^1 \|D^{(k)}(t)\|_{2 \rightarrow 2} dt \leq \varepsilon \quad \text{für alle } k \geq k_0. \end{aligned}$$

Das bedeutet aber:  $\int_0^1 \|D^{(k)}(t)\|_{2 \rightarrow 2} dt \rightarrow 0$ . Jetzt liefern (24.1) und (24.3):

$$\|x^{(k+1)} - x^*\|_2 \leq c \int_0^1 \|D^{(k)}(t)\|_{2 \rightarrow 2} dt \|x^{(k)} - x^*\|_2,$$

also die Q-superlineare Konvergenz.

**Aussage (iii):** Da  $x^{(k)}$  und  $x^{(k)} + t(x^* - x^{(k)})$  für alle  $t \in [0, 1]$  in  $B_{\delta}(x^*)$  liegen, können wir das Integral unter den stärkeren Voraussetzungen besser abschätzen:

$$\int_0^1 \|F'(x^{(k)}) - F'(x^{(k)} + t(x^* - x^{(k)}))\|_{2 \rightarrow 2} dt \leq \int_0^1 L t \|x^* - x^{(k)}\|_2 dt = \frac{L}{2} \|x^{(k)} - x^*\|_2.$$

Aus (24.3) erhalten wir nun:

$$\|x^{(k+1)} - x^*\|_2 \leq c \frac{L}{2} \|x^{(k)} - x^*\|_2^2.$$

□

<sup>1</sup>**Beachte:** Unter dem Integral stehen Matrizen.

**Bemerkung 24.6** (Zum lokalen Newton-Verfahren).

- (i) Das lokale Newton-Verfahren ([Algorithmus 24.1](#)) kann scheitern, denn  $F'(x^{(k)})$  muss nicht regulär sein, falls man außerhalb der (unbekannten) garantierten Konvergenzumgebung  $B_\delta(x^*)$  startet.
- (ii) Das sogenannte vereinfachte Newton-Verfahren, bei dem in [Zeile 3](#) von [Algorithmus 24.1](#) statt  $F'(x^{(k)})$  die feste (invertierbare) Matrix  $F'(x^{(0)})$  verwendet wird, konvergiert noch lokal  $Q$ -linear.

Ende der Woche 12

## Kapitel 9 Iterative Lösungsverfahren für lineare Gleichungssysteme

In diesem Kapitel geht es um **iterative Lösungsverfahren** für lineare Gleichungssysteme  $Ax = b$  mit  $A \in \mathbb{R}^{n \times n}$  und  $b \in \mathbb{R}^n$ . Diese kommen vor allem dann zum Einsatz, wenn die Bestimmung einer LR- oder Cholesky-Zerlegung von  $A$ , also ein direktes Lösungsverfahren (Kapitel 2) zu zeit- und/oder zu speicher-aufwändig ist. Das wiederum ist der Fall für hochdimensionale Aufgaben. Für diese ist die Matrix  $A$  oft dünn besetzt, die Faktoren der Zerlegung jedoch nicht. Man nennt dieses Phänomen des zusätzlichen Speicherbedarfs der Zerlegungsfaktoren **fill-in**. Zwar gibt es spezielle Pivotisierungstechniken, die den *fill-in* reduzieren, dennoch fährt man bei hochdimensionalen Aufgaben mit iterativen Lösungsverfahren oft besser.

Im Unterschied zu den direkten Lösungsverfahren aus Kapitel 4 erzeugen iterative Verfahren eine Folge  $(x^{(k)})$  von Näherungslösungen von  $Ax = b$ , die in einem vom konkreten Verfahren abhängigen Sinne zunehmend besser werden. Man kann also iterative Verfahren bereits vor der Konvergenz stoppen und erhält eine Näherungslösung, die besser als die Startschätzung  $x^{(0)}$  ist. Diese Tatsache lässt sich in vielen Algorithmen ausnutzen, in denen das Lösen linearer Gleichungssysteme als Teilschritt vorkommt. Beispielsweise gibt es inexacte Varianten des Newton-Verfahrens (Algorithmus 24.1), die es erlauben, das lineare Gleichungssystem  $F'(x^{(k)}) d^{(k)} = -F(x^{(k)})$  nur bis auf eine Toleranz zu lösen, die davon abhängt, welche Norm  $\|F(x^{(k)})\|$  die rechte Seite noch hat.

Ein weiterer Unterschied zwischen direkten und iterativen Lösungsverfahren ist, dass direkte Verfahren Zugriff auf die Einträge  $a_{ij}$  der Matrix benötigen, während iterative Verfahren oft mit Matrix-Vektor-Produkten  $x \mapsto Ax$  auskommen. Das bedeutet, dass wir, um das Gleichungssystem  $Ax = b$  zu lösen, bei Verwendung eines solchen iterativen Verfahrens die Matrix  $A$  nicht explizit benötigen. Man spricht deshalb auch von **matrix-freien Verfahren**.

Es stellt sich sicher die Frage, in welchen Situationen das relevant sein könnte. Um das zu beantworten, erinnern wir uns daran, dass Matrizen  $A$  und lineare Abbildungen  $x \mapsto L(x)$  in einer Eins-zu-Eins-Beziehung stehen.<sup>1</sup> Wir können aber durchaus die lineare Abbildung etwa in Form eines Computer-Codes vorliegen haben, ohne dass die zugehörige Matrix darin explizit auftaucht. (**Quizfrage:** Können Sie Beispiele dafür nennen?)

---

<sup>1</sup>Jede Matrix  $A$  repräsentiert eine lineare Abbildung  $x \mapsto Ax$ , und umgekehrt wird jede lineare Abbildung durch eine Matrix repräsentiert.



Eigenschaft	direkte Verfahren	iterative Verfahren
Zerlegung der Matrix wird berechnet	ja	nein
Zugriff auf Matrixeinträge benötigt	ja	oft nein
Matrix-Vektor-Produkte benötigt	nein	oft ja
kann zur inexakten Lösung verwendet werden	nein	ja

## § 25 STATIONÄRE VERFAHREN

Ein **stationäres Verfahren** (englisch: *stationary method*) entsteht dadurch, dass wir das lineare Gleichungssystem  $Ax = b$  als eine Fixpunktaufgabe umschreiben und dann eine Fixpunktiteration anwenden. Je nachdem, wie man die Fixpunktaufgabe formuliert, wird das daraus resultierende Verfahren unterschiedlich bezeichnet.

Das einfachste Beispiel ist die **Richardson-Iteration** (englisch: *Richardson iteration*). Sie basiert auf der Darstellung

$$x = x + b - Ax, \quad (25.1)$$

woraus sich die Iterationsvorschrift

$$\begin{aligned} x^{(k+1)} &:= x^{(k)} + b - Ax^{(k)} \\ &= b + (\text{Id} - A) x^{(k)} \end{aligned} \quad (25.2)$$

ergibt. Wir werden später die Konvergenz dieses und der noch folgenden Verfahren mit Hilfe des **Banachschen Fixpunktsatzes 23.1** untersuchen. Dabei spielt bekanntlich die Kontraktivität der Fixpunktfunktion  $T(x)$ , hier

$$T(x) := b + (\text{Id} - A) x,$$

die entscheidende Rolle. Charakteristisch für lineare Gleichungssysteme ist, dass die Fixpunktfunktion  $T$  affin-linear ist und daher bei der Frage der Kontraktivität der konstante, von der rechten Seite  $b$  abhängige Anteil unerheblich ist:

$$T(x) - T(y) = b + (\text{Id} - A) x - [b + (\text{Id} - A) y] = (\text{Id} - A) (x - y).$$

Die Kontraktivität wird also davon abhängen, dass eine gewisse Norm der Matrix  $N := \text{Id} - A$  nicht zu groß ist.

Verwenden wir an Stelle von (25.1) die Darstellung

$$Px = Px + b - Ax \quad (25.3)$$

mit einer regulären Matrix  $P$ , so ergibt sich daraus die Iterationsvorschrift

$$Px^{(k+1)} := Px^{(k)} + b - Ax^{(k)},$$

also

$$\begin{aligned} x^{(k+1)} &:= x^{(k)} + P^{-1}(b - Ax^{(k)}) \\ &= P^{-1}b + (\text{Id} - P^{-1}A) x^{(k)}. \end{aligned} \quad (25.4)$$

Dieses Verfahren bezeichnet man als **vorkonditionierte Richardson-Iteration** (englisch: *preconditioned Richardson iteration*) mit **Vorkonditionierer**  $P$ . Es kommt jetzt auf die Eigenschaften der Matrix  $N := \text{Id} - P^{-1}A$  an. **Quizfrage:** Können Sie das vorkonditionierte Richardson-Verfahren deuten als gewöhnliche Richardson-Iteration für ein modifiziertes Problem?

Eine andere Idee für die Herleitung stationärer Verfahren ist, die Matrix gemäß  $A = A_1 - A_2$  in zwei Summanden aufzuteilen, wobei  $A_1$  regulär sein soll. Die Darstellung

$$A_1 x = b + A_2 x \quad (25.5)$$

führt dann auf die Iterationsvorschrift

$$A_1 x^{(k+1)} = b + A_2 x^{(k)},$$

also

$$\begin{aligned} x^{(k+1)} &= A_1^{-1} b + A_1^{-1} A_2 x^{(k)} \\ &= A_1^{-1} b + A_1^{-1} (A_1 - A) x^{(k)} \\ &= A_1^{-1} b + (\text{Id} - A_1^{-1} A) x^{(k)} \\ &= x^{(k)} + A_1^{-1} (b - A x^{(k)}) \quad (\text{Update-Form}). \end{aligned} \quad (25.6)$$

Aufgrund der Aufteilung  $A = A_1 - A_2$  nennt man dies auch eine **Splitting-Iteration**. Der Vergleich mit (25.4) zeigt aber, dass Splitting-Iterationen und vorkonditionierte Richardson-Iterationen identisch sind, wobei  $A_1$  die Rolle des Vorkonditionierers  $P$  einnimmt. Lediglich die Idee der Herleitung ist verschieden.

**Bemerkung 25.1.** Jedes stationäre Verfahren zur Lösung von  $Ax = b$  kann in der Form einer vorkonditionierten Richardson-Iteration (25.4) bzw. äquivalent in der Form eines Splitting-Verfahrens (25.6) geschrieben werden. Ein solches Verfahren ist durch die Wahl von  $P = A_1$  eindeutig bestimmt.

Zwei bekannte Vertreter stationärer Iterationen sind das Jacobi-Verfahren und das Gauß-Seidel-Verfahren. Sie basieren auf der Zerlegung  $A = D + L + R$ , wobei  $D$  die Diagonale und  $L$  und  $R$  die strikte untere (linke) bzw. obere (rechte) Dreiecksmatrix von  $A$  bezeichnen:

$$D = \begin{bmatrix} a_{11} & 0 & \cdots & 0 \\ 0 & \ddots & \ddots & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & a_{nn} \end{bmatrix}, \quad L = \begin{bmatrix} 0 & \cdots & \cdots & 0 \\ a_{21} & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ a_{n1} & \cdots & a_{n,n-1} & 0 \end{bmatrix}, \quad R = \begin{bmatrix} 0 & a_{12} & \cdots & a_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & a_{n-1,n} \\ 0 & \cdots & \cdots & 0 \end{bmatrix}.$$

Beim **Jacobi-Verfahren** (auch **Gesamtschrittverfahren**, englisch: *Jacobi iteration*) wird  $A_1 = D$  gewählt und folglich  $A_2 = -(L + R)$ . Die Iteration (25.6) lautet daher

$$\begin{aligned} x^{(k+1)} &= D^{-1} b - D^{-1} (L + R) x^{(k)} \\ &= D^{-1} (b - (L + R) x^{(k)}). \end{aligned} \quad (25.7)$$

Für die Wohldefiniertheit muss natürlich vorausgesetzt werden, dass  $A_1 = D$  invertierbar ist, also dass auf der Diagonale von  $A$  keine Nullen stehen.

Eine Iteration des Jacobi-Verfahrens, wobei wir auf den Iterationszähler  $k$  verzichten und stattdessen die alte Iterierte  $x$  mit der neuen Iterierten überschreiben, kann wie folgt implementiert werden:

```

1: Setze  $y := x$ 
2: for  $i = 1, \dots, n$  do
3:   Setze  $x_i := \frac{1}{a_{ii}} \left( b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} y_j \right)$ 
4: end for

```

**Quizfrage:** Die obige Implementierung benötigt nun doch Zugriff auf die Elemente der Matrix  $A$ . Wie kann das Jacobi-Verfahren alternativ weitestgehend matrix-frei implementiert werden?

Beim **Gauß-Seidel-Verfahren** (auch **Einzelschrittverfahren**, englisch: *Gauss-Seidel iteration*) ist  $A_1 = D + L$  und folglich  $A_2 = -R$ . Die Iteration (25.6) lautet dann

$$\begin{aligned} x^{(k+1)} &= (D + L)^{-1} b - (D + L)^{-1} R x^{(k)} \\ &= (D + L)^{-1} (b - R x^{(k)}). \end{aligned} \quad (25.8)$$

**Quizfrage:** Was bedeutet nun die Voraussetzung, dass  $A_1 = D + L$  eine reguläre Matrix ist? **Quizfrage:** Wie können wir eine Iteration des Gauß-Seidel-Verfahrens implementieren?

**Quizfrage:** Wie sieht es mit der Parallelisierbarkeit des Jacobi- und des Gauß-Seidel-Verfahrens aus?

**Quizfrage:** Kann das Gauß-Seidel-Verfahren ebenfalls weitestgehend matrix-frei implementiert werden?

Das **symmetrische Gauß-Seidel-Verfahren** (englisch: *symmetric Gauss-Seidel iteration*) verwendet die Zerlegung  $A = A_1 - A_2$  mit  $A_1 = (D + L)D^{-1}(D + R)$  und  $A_2 = LD^{-1}R$ . (**Quizfrage:** Können Sie die Darstellung von  $A_2$  bestätigen?) Die Iteration (25.6) wird dann oft in zwei Teilschritten geschrieben:

$$\begin{aligned} x^{(k+1/2)} &= (D + L)^{-1} (b - R x^{(k)}), \\ x^{(k+1)} &= (D + R)^{-1} (b - L x^{(k+1/2)}). \end{aligned} \quad (25.9)$$

**Quizfrage:** Können Sie sich davon überzeugen, dass (25.9) in der Tat die Iteration (25.6) für das symmetrische Gauß-Seidel-Verfahren implementiert?

Schließlich erwähnen wir noch das **Überrelaxationsverfahren** (englisch: *successive overrelaxation*, kurz: SOR), das auf der Iterationsvorschrift

$$(\omega^{-1}D + L) x^{(k+1)} = b - (R + (1 - \omega^{-1})D) x^{(k)} \quad (25.10)$$

beruht, wobei oft  $\omega > 1$  ist. Es gilt also  $A_1 = \omega^{-1}D + L$  und  $A_2 = (\omega^{-1} - 1)D - R$ . Im Falle von  $0 < \omega < 1$  spricht man auch von **Unterrelaxation**. (**Quizfrage:** Was ist der Fall  $\omega = 1$ ?)

Vom Gauß-Seidel- und SOR-Verfahren gibt es auch Varianten, die rückwärts durch die Gleichungen gehen. Das entspricht der Vertauschung von  $L$  und  $R$ . Durch Kombination der Vorwärts- und Rückwärts-Varianten erhält man zusätzlich symmetrische Varianten. [Tabelle 25.1](#) fasst diese häufigen stationären Iterationsverfahren zusammen.

Name	Wahl von $P$ in (25.4)	Iterationsmatrix $N = \text{Id} - P^{-1}A$
Richardson	$\text{Id}$	$\text{Id} - A$
Jacobi	$D$	$-D^{-1}(L + R)$
Gauß-Seidel	$D + L$	$-(D + L)^{-1}R =: N_{\text{GSfwd}}$
Gauß-Seidel rückwärts	$D + R$	$-(D + R)^{-1}L =: N_{\text{GSbwd}}$
Gauß-Seidel symmetrisch	$(D + L)D^{-1}(D + R)$	$N_{\text{GSbwd}}N_{\text{GSfwd}}$
successive overrelaxation (SOR)	$\frac{1}{\omega}D + L$	$-(\frac{1}{\omega}D + L)^{-1}((1 - \frac{1}{\omega})D + R) =: N_{\text{SORfwd}}$
succ. overrelax. rückwärts	$\frac{1}{\omega}D + R$	$-(\frac{1}{\omega}D + R)^{-1}((1 - \frac{1}{\omega})D + L) =: N_{\text{SORbwd}}$
symm. succ. overrelax. (SSOR)	$\frac{\omega}{2-\omega}(\frac{1}{\omega}D + L)D^{-1}(\frac{1}{\omega}D + R)$	$N_{\text{SORbwd}}N_{\text{SORfwd}}$

Tabelle 25.1: Beispiele stationärer Iterationsverfahren.

Es bleibt noch die Frage, wann solche stationären Iterationsverfahren konvergieren. Wie oben bereits angedeutet, können wir das mit Hilfe des **Banachschen Fixpunktsatzes** 23.1 untersuchen. Die Fixpunkt-funktion des allgemeinen Splitting-Verfahrens (bzw. vorkonditionierten Richardson-Verfahrens) ist laut (25.4)

$$T(x) := P^{-1}b + (\text{Id} - P^{-1}A)x =: P^{-1}b + Nx. \quad (25.11)$$

Wir setzen im Fixpunktsatz  $M = \mathbb{R}^n$ , sodass  $T$  auf jeden Fall eine Selbstabbildung  $M \rightarrow M$  ist. (**Quizfrage:** Könnte nicht evtl. eine kleinere Menge  $M$  eine bessere Lipschitz-Konstante bringen?) Bei der Frage nach der Kontraktivität spielt der konstante Anteil  $P^{-1}b$  keine Rolle. Es kommt also nur darauf an, dass die Iterationsmatrix  $N := \text{Id} - P^{-1}A$  kleine Norm hat:

$$\|\text{Id} - P^{-1}A\|_{* \rightarrow *} < 1.$$

Dabei ist  $\|\cdot\|_{* \rightarrow *}$  die Operatornorm (Definition 2.2) zu einer beliebigen Vektornorm  $\|\cdot\|_*$  im  $\mathbb{R}^n$ .

Es stellt sich die Frage, wie wir die optimale Norm finden, um  $\|\cdot\|_{* \rightarrow *}$  und damit die Lipschitz-Konstante der Fixpunkt-funktion  $T$  möglichst klein zu machen. Tatsächlich ist es gar nicht erforderlich, die optimale Norm ausfindig zu machen, da wir die Kontraktivität auch mit Hilfe des Spektralradius charakterisieren können.

**Definition 25.2** (Spektralradius). Es sei  $N \in \mathbb{R}^{n \times n}$  eine beliebige Matrix. Der **Spektralradius** von  $N$  ist definiert als

$$\text{spr}(N) := \max\{|\lambda| \mid \lambda \in \mathbb{C} \text{ ist ein Eigenwert von } N\}, \quad (25.12)$$

also als der Betrag eines betragsgrößten Eigenwertes von  $N$ .

Die entscheidende Aussage zum Zusammenhang von Spektralradius und Operatornorm ist:

**Lemma 25.3** (Spektralradius).

(i) Es sei  $\|\cdot\|_*$  eine beliebige Vektornorm auf  $\mathbb{R}^n$  und  $\|\cdot\|_{* \rightarrow *}$  die dazugehörige Operatornorm. Dann gilt

$$\text{spr}(N) \leq \|N\|_{* \rightarrow *}. \quad (25.13)$$

(ii) Für jedes  $\varepsilon > 0$  gibt es eine Vektornorm  $\|\cdot\|_*$ , sodass für die dazugehörige Operatornorm gilt:

$$\|N\|_{* \rightarrow *} \leq \text{spr}(N) + \varepsilon. \quad (25.14)$$

*Beweis.* Es sei  $(\lambda, v)$  ein beliebiges Eigenpaar der Matrix  $N$ , also ist  $v \neq 0$ , und es gilt  $Nv = \lambda v$ . Es folgt

$$|\lambda| \|v\|_* = \|Nv\|_* \leq \|N\|_{* \rightarrow *} \|v\|_*.$$

Daraus folgt  $|\lambda| \leq \|N\|_{* \rightarrow *}$  für alle Eigenwerte  $\lambda$  von  $N$  und somit auch  $\text{spr}(N) \leq \|N\|_{* \rightarrow *}$ , also **Aussage (i)**. Den Beweis der **Aussage (ii)** findet man beispielsweise in [Rannacher, 2017](#), Hilfssatz 6.1.  $\square$

**Folgerung 25.4** (Spektralradius und Kontraktheit). *Die folgenden Aussagen sind äquivalent:*

- (i) Es gibt eine Vektornorm  $\|\cdot\|_*$  auf  $\mathbb{R}^n$ , sodass die Fixpunktfunktion (25.11) eine Kontraktion ist.
- (ii) Für den Spektralradius der Iterationsmatrix  $N$  gilt:  $\text{spr}(N) < 1$ .

*Beweis.* **Aussage (i):** Es sei  $\|\cdot\|_*$  eine Vektornorm, sodass die Fixpunktfunktion (25.11) eine Kontraktion ist. Das heißt, wir haben  $\|N\|_{* \rightarrow *} < 1$ . Aus (25.13) folgt dann  $\text{spr}(N) \leq \|N\|_{* \rightarrow *} < 1$ .

Nun sei umgekehrt  $\text{spr}(N) < 1$ . Wir wählen  $\varepsilon > 0$  so, dass auch  $\text{spr}(N) + \varepsilon < 1$  bleibt. Nach **Lemma 25.3 (ii)** gibt es dann eine Vektornorm  $\|\cdot\|_*$ , sodass gilt:  $\|N\|_{* \rightarrow *} \leq \text{spr}(N) + \varepsilon < 1$ . Damit gilt

$$\|T(x) - T(y)\|_* = \|N(x - y)\|_* \leq \|N\|_{* \rightarrow *} \|x - y\|_* < \|x - y\|_*,$$

d. h.,  $T$  ist eine Kontraktion bzgl. der Norm  $\|\cdot\|_*$ .  $\square$

**Bemerkung 25.5** (Stationäre Verfahren).

- (1) In jeder Iteration eines stationären Verfahrens ist ein lineares Gleichungssystem mit der verfahrensabhängigen Matrix  $P$  zu lösen.
- (2) Die Bezeichnung «Vorkonditionierung» und «vorkonditionierte Richardson-Iteration» ist etwas irreführend, weil die Wahl von  $P$  dazu genutzt wird, den Spektralradius  $\text{spr}(N) = \text{spr}(\text{Id} - P^{-1}A)$  zu reduzieren und nicht etwa die Konditionszahl.
- (3) Mit der Wahl von  $P$  verfolgt man zwei konkurrierende Ziele:
  - Der Spektralradius  $\text{spr}(\text{Id} - P^{-1}A)$  soll möglichst klein sein.
  - Gleichungssysteme mit  $P$  sollen möglichst einfach zu lösen sein.

Wir geben (ohne Beweise) noch einige Beispiele für Konvergenzresultate an:

- (1) Ist  $A \in \mathbb{R}^{n \times n}$  s. p. d., dann konvergiert das Gauß-Seidel-Verfahren.

(2) Ist  $A \in \mathbb{R}^{n \times n}$  eine **strikt diagonal-dominante Matrix** (vgl. [Abschnitt 10.1](#)), gilt also

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

so konvergieren das Jacobi- und das Gauß-Seidel-Verfahren.

(3) Für beliebiges  $A \in \mathbb{R}^{n \times n}$  mit regulärem Diagonalanteil  $D$  gilt  $\text{spr}(N_{\text{SORfwd}}) \geq |\omega - 1|$  für alle  $\omega \in \mathbb{R}$ . Daher ist nur der Parameterbereich  $\omega \in (0, 2)$  interessant. Falls  $A$  s. p. d. ist, dann konvergiert das SOR-Verfahren tatsächlich für alle  $\omega \in (0, 2)$ .

Um den Übergang zum folgenden [Abschnitt 26](#) zu motivieren, schreiben wir noch die ersten paar Iterationen der Richardson-Iteration explizit aus. Dabei verwenden wir neben den Iterierten  $x^{(k)}$  auch die Residuen (vgl. [\(11.2\)](#))

$$r^{(k)} := Ax^{(k)} - b.$$

Es gilt laut Iterationsvorschrift [\(25.2\)](#) des Richardson-Verfahrens

$$x^{(1)} = x^{(0)} + b - Ax^{(0)}$$

und daher

$$x^{(1)} - x^{(0)} = -r^{(0)}.$$

Das nächste Residuum ist also

$$\begin{aligned} r^{(1)} &= Ax^{(1)} - b \\ &= A(x^{(1)} - x^{(0)}) - b + Ax^{(0)} \\ &= -Ar^{(0)} + r^{(0)} \\ &= (\text{Id} - A) r^{(0)}. \end{aligned}$$

Für die folgende Iteration gilt laut Iterationsvorschrift [\(25.2\)](#)

$$\begin{aligned} x^{(2)} - x^{(0)} &= x^{(1)} - x^{(0)} + b - A(x^{(1)} - x^{(0)}) - Ax^{(0)} \\ &= -r^{(0)} - r^{(0)} + Ar^{(0)} \\ &= (-2 \text{Id} + A) r^{(0)} \end{aligned}$$

und für das Residuum

$$\begin{aligned} r^{(2)} &= A(x^{(2)} - x^{(0)}) - b + Ax^{(0)} \\ &= A(-2 \text{Id} + A) r^{(0)} + r^{(0)} \\ &= (\text{Id} - A)^2 r^{(0)}. \end{aligned}$$

Insgesamt können wir mit Hilfe eines Induktionsarguments folgende Relationen<sup>2</sup> zeigen:

$$r^{(k)} = p_k(A) r^{(0)}, \tag{25.15a}$$

$$x^{(k)} - x^{(0)} = q_{k-1}(A) r^{(0)}, \tag{25.15b}$$

<sup>2</sup> $p(A)$  ist dabei ein Polynom, in das die Matrix  $A$  eingesetzt wird. Ist  $p(x) = \sum_{j=0}^k \alpha_j x^j$ , so ist  $p(A) = \sum_{j=0}^k \alpha_j A^j$ , wobei  $A^0 = \text{Id}$  gilt,  $A^1 = A$ ,  $A^2 = AA$  usw.

wobei für die Polynome  $p_k, q_k \in P_k$  gilt:

$$p_k(x) = (1-x)^k \quad \text{für } k \geq 0, \quad (25.16a)$$

$$q_{k-1}(x) = \frac{p_k(x) - 1}{x} = \frac{(1-x)^k - 1}{x} \quad \text{für } k \geq 1. \quad (25.16b)$$

**Definition 25.6** (Krylov-Unterraum). Gegeben sei eine Matrix  $A \in \mathbb{R}^{n \times n}$  und ein Vektor  $r \in \mathbb{R}^n$ . Der Unterraum

$$K_k(A, r) := \text{span}\{r, Ar, A^2r, \dots, A^{k-1}r\} \quad (25.17)$$

heißt der  $k$ -te **Krylov-Unterraum** (englisch: **Krylov subspace**) zur Matrix  $A$  und zum Startvektor  $r$ .

Offenbar gilt für den Krylov-Unterraum

$$K_k(A, r) = \{p(A)r \mid p \in P_{k-1}\}.$$

Seine Elemente entstehen also durch Variation der Polynome von einem gewissen Höchstgrad. Wir können die Richardson-Iteration (25.2) demnach so verstehen, dass die Iterierten

$$r^{(k)} \in K_{k+1}(A, r^{(0)}), \quad (25.18a)$$

$$x^{(k)} \in x^{(0)} + K_k(A, r^{(0)}) \quad (25.18b)$$

liegen, wobei die jeweiligen Elemente des Krylov-Unterraumes durch die Polynome  $p_k$  und  $q_k$  festgelegt sind. Diese Polynome (25.16) stehen bereits fest, bevor die Iteration überhaupt begonnen wird, hängen also insbesondere nicht von der rechten Seite  $b$  und der Startschätzung  $x^{(0)}$  ab. Das begründet nochmals den Namen **stationäre Verfahren** für die in diesem Abschnitt betrachteten Verfahren.

Da sich das allgemeine stationäre Verfahren (25.4) auch als Richardson-Iteration für die modifizierte Aufgabe<sup>3</sup>  $P^{-1}Ax = P^{-1}b$  deuten lässt, können wir zeigen, dass an Stelle von (25.15)

$$r^{(k)} = p_k(P^{-1}A) P^{-1}r^{(0)}, \quad (25.19a)$$

$$x^{(k)} - x^{(0)} = q_{k-1}(P^{-1}A) P^{-1}r^{(0)} \quad (25.19b)$$

mit denselben Polynomen (25.16) gilt. Die relevanten Krylov-Unterräume haben dann die Gestalt

$$\begin{aligned} K_k(P^{-1}A, P^{-1}r) &= \text{span}\{P^{-1}r, P^{-1}AP^{-1}r, (P^{-1}A)^2P^{-1}r, \dots, (P^{-1}A)^{k-1}P^{-1}r\} \\ &= \{p(P^{-1}A) P^{-1}r \mid p \in P_{k-1}\}. \end{aligned}$$

## § 26 INSTATIONÄRE VERFAHREN

Die Überlegungen am Ende des vorhergehenden Abschnitts haben gezeigt, dass wir jedes stationäre Iterationsverfahren zur Lösung von  $Ax = b$  so verstehen können, dass die Iterierten und die Residuen in gewissen Krylov-Unterräumen liegen und dass die konkreten Elemente durch die Polynome  $p_k$  und  $q_k$

<sup>3</sup>Das Residuum ist noch immer als  $r = Ax - b$  definiert.

festgelegt werden, die von der rechten Seite  $b$  und der Startschätzung  $x^{(0)}$  unabhängig sind. **Instationäre Verfahren** hingegen bestimmen ihre Iterierten in Abhängigkeit des bisherigen Iterationsverlaufs. Die bekanntesten Vertreter instationärer Verfahren sind die **Krylov-Unterraumverfahren**, in denen dieselben Unterräume wie in (25.18) eine Rolle spielen. Verschiedene Verfahren unterscheiden sich darin, in welchen Räumen die Iterierten jeweils gewählt und wie diese festgelegt werden, wie also die Polynome in (25.15) zustande kommen. Auch für Krylov-Unterraumverfahren spielt wieder der Aspekt der «Vorkonditionierung» in der Praxis eine wichtige Rolle, auf den wir hier jedoch nicht eingehen werden.

Wir sprechen hier nur das bekannteste Krylov-Unterraumverfahren an, und zwar das **Verfahren der konjugierten Gradienten** (englisch: *conjugate gradient method*, kurz: CG-Verfahren). Dieses Verfahren ist gemacht für Matrizen  $A$ , die s. p. d. sind. Es kann dadurch definiert werden, dass wie in (25.18) die Iterierten

$$\begin{aligned} r^{(k)} &\in K_{k+1}(A, r^{(0)}), \\ x^{(k)} &\in x^{(0)} + K_k(A, r^{(0)}) \end{aligned}$$

gewählt werden und außerdem gefordert wird, dass  $x^{(k)}$  die durch  $A$  induzierte Norm (im Quadrat) des Fehlers  $\|x^{(k)} - x^*\|_A^2 = (x^{(k)} - x^*)^T A (x^{(k)} - x^*)$  in der zur Verfügung stehenden Menge minimiert. Diese Bedingung ist äquivalent dazu, dass das Residuum  $r^{(k)}$  senkrecht auf  $K_k(A, r^{(0)})$  steht. Das dabei (implizit) in jeder Iteration bestimmte Polynom besitzt folgende Minimierungseigenschaft:

$$\|x^{(k)} - x^*\|_A^2 = \min \{ \|p_k(A)(x^{(0)} - x^*)\|_A^2 \mid p_k \in P_k, p_k(0) = 1 \}.$$

Das CG-Verfahren kann wie in **Algorithmus 26.1** gezeigt implementiert werden.

**Algorithmus 26.1** (CG-Verfahren).

**Eingabe:** Startwert  $x^{(0)} \in \mathbb{R}^n$

**Eingabe:** rechte Seite  $b \in \mathbb{R}^n$

**Eingabe:** s. p. d. Matrix  $A$  (oder Matrix-Vektor-Produkte mit  $A$ )

**Ausgabe:** näherungsweise Lösung der Aufgabe  $Ax = b$

```

1: Setze  $k := 0$ 
2: Setze  $r^{(0)} := Ax^{(0)} - b$ 
3: Setze  $d^{(0)} := -r^{(0)}$ 
4: Setze  $\delta^{(0)} := -(r^{(0)})^T d^{(0)}$ 
5: while Abbruchkriterium nicht erfüllt do
6:   Setze  $q^{(k)} := Ad^{(k)}$ 
7:   Setze  $\alpha^{(k)} := \delta^{(k)} / ((d^{(k)})^T q^{(k)})$ 
8:   Setze  $x^{(k+1)} := x^{(k)} + \alpha^{(k)} d^{(k)}$ 
9:   Setze  $r^{(k+1)} := r^{(k)} + \alpha^{(k)} q^{(k)}$ 
10:  Setze  $d^{(k+1)} := -r^{(k+1)}$ 
11:  Setze  $\delta^{(k+1)} := -(r^{(k+1)})^T d^{(k+1)}$ 
12:  Setze  $\beta^{(k+1)} := \delta^{(k+1)} / \delta^{(k)}$ 
13:  Setze  $d^{(k+1)} := d^{(k+1)} + \beta^{(k+1)} d^{(k)}$ 
14:  Setze  $k := k + 1$ 
15: end while
16: return  $x^{(k)}$ 
```





# Index

- 1-Norm, 13
- 2-Norm, 13
- $C^k$ -Funktion, 18
- $L^2$ -Innenprodukt, 153
- $O$ , *siehe* Groß-O
- $o$ , *siehe* Klein-o
- $\infty$ -Norm, 13
- $k$ -mal stetig partiell differenzierbare Funktion, 18
- (Newtonsche) dividierten Differenzen, 122
- A-posteriori-Fehlerabschätzung der Fixpunktiteration, 160
- A-priori-Fehlerabschätzung der Fixpunktiteration, 160
- abgeschlossene Newton-Cotes-Formeln, 146
- Ableitung, 17
- absolute Konditionszahl, 27
- absolute partielle Konditionszahlen, 26
- Adjazenzmatrix, 8
- Algorithmus von de Casteljau, 134
- allgemeine linearen Gruppe, 38
- Alternante, 139
- Alternantensatz, 139
- Approximationsfehler, 123
- asymptotisch obere Schranke, 21
- asymptotisch vernachlässigbar, 22
- Ausgleichsproblem, 100
- Auslöschung, 33
- baryzentrische Darstellung, 120
- baryzentrische Gewichte, 120
- Basis einer Fließkommazahl, 47
- Berechnungsgraph, 59
- Bernstein-Polynom, 131
- binary64, *siehe* double precision
- Bisektionsverfahrens, 158
- Bézier-Kurve, 133
- Cauchy-Schwarz-Ungleichung, 13
- Cholesky-Zerlegung, 97
- Definitheit einer Norm, 12
- Designmatrix, 100
- diagonal-dominante Matrix, 96
- differentielle Sensitivitätsanalyse, 25
- differenzierbare Funktion, 17
- direkte Lösungsverfahren für lineare Gleichungssysteme, 72
- double precision, 49
- Dreiecksungleichung, 12
- dünne Singulärwertzerlegung, 41
- echte Addition, 33
- echte Subtraktion, 33
- Eigenpaar, 38
- Eigenvektor, 38
- Eigenwert, 38
- Eigenwertaufgabe, 10
- Einheitsmatrix, 14
- Einzelschrittverfahren, 171
- Elastizität, 29
- Euklidisches Innenprodukt, 13
- Euklidische Norm, 13
- Exponent, 47
- Extrapolation, 124
- Fehlerquadratsumme, 100
- Festkommazahlen, 47
- fill-in, 168
- Finite-Differenzen-Approximation, 129
- Fixpunkt, 159
- Fixpunkt-Iteration, 160
- Fließkommagitter, 47
- Fließkommazahl, 47
- Frobeniusmatrix, 74
- Ganzzahlen, 47
- Gauß-Seidel-Verfahren, 171
- Genauigkeitsgrad einer Quadraturformel, 145

- Gesamtschrittverfahren, 170  
Gewichte einer Quadraturformel, 145  
Givens-Rotation, 107  
Gram-Schmidt-Verfahren, 105  
Groß-O, 21, 22  
gut konditioniert im absoluten Sinne, 27  
gut konditioniert im relativen Sinne, 30
- Hermite-Interpolation, 129  
Hessematrix, 19  
Horner-Schema, 64  
Householder-Reflexion, 109  
Householder-Transformation, 109
- Instationäre Verfahren, 176  
Interpolation, 115  
Interpolationsaufgabe, 115  
Interpolationsbedingungen, 115  
iterative Lösungsverfahren für lineare Gleichungssysteme, 168
- Jacobi-Verfahren, 170  
Jacobimatrix, 17
- Kante eines Graphen, 8  
Klein-o, 22  
Kleinste-Quadrate-Aufgabe, 100  
Knoten eines Graphen, 8  
Knotenpolynom, 119  
Koeffizienten eines Polynoms, 115  
Konditionszahl einer Matrix, 35, 103  
konstante Interpolation, 115  
Kontraktion, 159  
Kontrollpolygon einer Bézier-Kurve, 134  
Kontrollpunkte einer Bézier-Kurve, 133  
konvexe Menge, 18  
Krylov-Unterraum, 175  
Krylov-Unterraumverfahren, 176  
kubische Interpolation, 115  
kubischer Spline, 128
- Lagrange-Interpolation, 129  
Lagrange-Polynom, 118  
Landau-Notation, 21  
Legendre-Polynome, 155  
lineare Interpolation, 115  
linke Dreiecksmatrix, *siehe* untere Dreiecksmatrix
- Linkssingulärvektoren, 39  
Lipschitz-Konstante, 20  
Lipschitz-stetige Funktion, 20  
LR-Zerlegung, 75  
LR-Zerlegung mit Pivotisierung, 79  
LR-Zerlegung mit Totalpivotisierung, 83  
LU-Zerlegung, *siehe* LR-Zerlegung
- Mantisse, 47  
Maschinengenauigkeit, 50  
Maschinenoperation, 52  
matrix-freie Verfahren, 168  
Metrik, 13  
Mittelwertsatz, 19  
Modellkalibrierung, 100  
Modellmatrix, 100  
modifiziertes Gram-Schmidt-Verfahren, 106  
Monome, 118
- Newton-Richtung, 163  
Newtonsche Basispolynom, 121  
Norm, 12  
Norm der gleichmäßigen Konvergenz, 124  
Normalengleichung(en), 101  
Normalgleichung(en), 101  
normalisierte Fließkommazahl, 47  
normierte QR-Zerlegung, 104  
Numerik, 7  
numerische Differentiation, 129  
numerische Integration, 145  
numerische Mathematik, 7  
numerische Quadratur, 145
- obere Dreiecksmatrix, 73  
offene  $\varepsilon$ -Kugel, 13  
offene  $\varepsilon$ -Umgebung, 13  
offene Newton-Cotes-Formeln, 150  
Operatornorm, 14  
orthogonale Gruppe, 38  
orthogonale Matrix, 38  
Orthonormalbasis, 38
- Parameteridentifikation, 100  
Parameterschätzung, 100  
Permutationsmatrix, 76  
Pivotelement, 77  
Polynome, 115  
positiv definite Matrix, 96

- positive Homogenität einer Norm, 12  
Potenzmethode, 10
- Q-lineare Konvergenz, 158, 160  
Q-quadratische Konvergenz, 158  
Q-superlineare Konvergenz, 158  
QR-Zerlegung, 104, 107  
quadratische Interpolation, 115
- rechte Dreiecksmatrix, *siehe* obere Dreiecksmatrix
- Rechtssingulärvektoren, 39  
reduzierte QR-Zerlegung, 104  
Regressionsanalyse, 100  
relative Konditionszahl, 30  
relative partielle Konditionszahlen, 29  
Residuenvektor, 100  
Residuum, 100, 162  
Restglied, 17  
Richardson-Extrapolation, 155  
Richardson-Iteration, 169  
Romberg-Folge, 156  
Romberg-Quadratur, 157  
Rundung, 50  
Rundung zur nächstgelegenen Fließkommazahl, 50
- Runge-Phänomen, 126  
Rückwärtsanalyse, 69  
rückwärtsstabiler Algorithmus, 69  
Rückwärtssubstitution, 73
- Satz von Weierstraß, 101, 139  
schlecht konditioniert im absoluten Sinne, 27  
schlecht konditioniert im relativen Sinne, 30  
Sensitivitätsanalyse, 24  
Signum, 16  
Singulärwerte, 39  
Singulärwertzerlegung, 39  
Spaltenpivotsuche, 78  
Spaltensummennorm, 15  
spd Matrix, 96  
Spektralnrm, 15  
Spektralradius, 172  
Spektralzerlegung einer reellen, symmetrischen Matrix, 38  
Spline, 128  
Splitting-Iteration, 170  
stabiler Algorithmus, 71
- stationäre Verfahren zur Lösung linearer Gleichungssysteme, 169  
strikt diagonal-dominante Matrix, 174  
Stützstellen, 115  
Stützstellen einer Quadraturformel, 145  
Stützwerte, 115  
Subadditivität einer Norm, 12  
Supremumsnorm, 124  
SVD, *siehe* Singulärwertzerlegung  
symmetrische Gauß-Seidel-Verfahren, 171
- total nichtnegative Matrix, 96  
totale Pivotsuche, 83  
Transpositionsmatrix, 76  
tridiagonale Matrix, 96  
Tschebschow-Polynom, 141  
Tschebyscheff-Approximation, *siehe* Tschebyschow-Approximation  
Tschebyschow-Approximationsaufgabe, 138  
Tschebyschow-Knoten, 141
- untere Dreiecksmatrix, 75  
Unterrelaxation, 171
- Vandermonde-Matrix, 116  
Vektor der Messwerte, 100  
Vektoriteration, 10  
verallgemeinertes Bernstein-Polynom, 131  
Verbindungsstrecke, 18  
Verfahren der konjugierten Gradienten, 176  
Verfahren von Aitken-Neville, 156  
Verfahrensfehler, 123  
Vertauschungsmatrix, 76  
volle QR-Zerlegung, 107  
Von-Mises-Iteration, 10  
Vorkonditionierer, 170  
vorkonditionierte Richardson-Iteration, 170  
Vorwärtsanalyse, 58  
Vorwärtselimination, 73  
vorwärtsstabiler Algorithmus, 66
- wissenschaftliches Rechnen, 7
- Zeilensummennorm, 15  
Zeilenvektoren, 13  
Zerlegung der Eins, 131  
zulässiger Bereich eines Fließkommagitters, 49  
Überrelaxationsverfahren, 171

# Literatur

- Bauer, F. L. (1974). «Computational graphs and rounding error». *SIAM Journal on Numerical Analysis* 11.1, S. 87–96. DOI: [10.1137/0711010](https://doi.org/10.1137/0711010).
- Beutelspacher, A. (2014). *Lineare Algebra. Eine Einführung in die Wissenschaft der Vektoren, Abbildungen und Matrizen*. 8. Aufl. Springer Fachmedien Wiesbaden. DOI: [10.1007/978-3-658-02413-0](https://doi.org/10.1007/978-3-658-02413-0).
- Bornemann, F. (2018). *Numerische lineare Algebra*. 2. Aufl. Springer Fachmedien Wiesbaden. DOI: [10.1007/978-3-658-24431-6](https://doi.org/10.1007/978-3-658-24431-6).
- Brin, S.; L. Page (1998). «The anatomy of a large-scale hypertextual Web search engine». *Computer Networks and ISDN Systems* 30.1-7, S. 107–117. DOI: [10.1016/s0169-7552\(98\)00110-x](https://doi.org/10.1016/s0169-7552(98)00110-x).
- Freund, R. W.; R. W. Hoppe (2007). *Stoer/Bulirsch: Numerische Mathematik 1*. 10. Aufl. Springer Berlin Heidelberg. DOI: [10.1007/978-3-540-45390-1](https://doi.org/10.1007/978-3-540-45390-1).
- Geiger, C.; C. Kanzow (1999). *Numerische Verfahren zur Lösung unrestringierter Optimierungsaufgaben*. New York: Springer. DOI: [10.1007/978-3-642-58582-1](https://doi.org/10.1007/978-3-642-58582-1).
- Golub, G.; C. van Loan (1983). *Matrix Computations*. Baltimore: Johns Hopkins University Press.
- Heuser, H. (2002). *Lehrbuch der Analysis. Teil 2*. 12. Aufl. Stuttgart: B.G.Teubner. DOI: [10.1007/978-3-322-96826-5](https://doi.org/10.1007/978-3-322-96826-5).
- Higham, N. J. (2002). *Accuracy and Stability of Numerical Algorithms*. Society for Industrial und Applied Mathematics. DOI: [10.1137/1.9780898718027](https://doi.org/10.1137/1.9780898718027).
- Higham, N. J.; M. R. Dennis; P. Glendinning; P. A. Martin; F. Santosa; J. Tanner, Hrsg. (2016). *The Princeton Companion to Applied Mathematics*. Princeton University Press. DOI: [10.1515/9781400874477](https://doi.org/10.1515/9781400874477).
- Kahan, W. (1966). «Numerical linear algebra». *Canadian Mathematical Bulletin* 9.05, S. 757–801. DOI: [10.4153/cmb-1966-083-2](https://doi.org/10.4153/cmb-1966-083-2).
- Rannacher, R. (2017). *Numerik 0: Einführung in die Numerische Mathematik*. Heidelberg University Publishing. DOI: [10.17885/HEIUP.206.281](https://doi.org/10.17885/HEIUP.206.281).
- Trefethen, L. N. (1992). «The definition of numerical analysis». *SIAM News*. URL: [https://people.maths.ox.ac.uk/trefethen/publication/PDF/1992\\_55.pdf](https://people.maths.ox.ac.uk/trefethen/publication/PDF/1992_55.pdf).
- Trefethen, L. N.; D. Bau (1997). *Numerical Linear Algebra*. SIAM. DOI: [10.1137/1.9780898719574](https://doi.org/10.1137/1.9780898719574).
- Tyrtshnikov, E. E. (1994). «How bad are Hankel matrices?» *Numerische Mathematik* 67.2, S. 261–269. DOI: [10.1007/s002110050027](https://doi.org/10.1007/s002110050027).