

# EINFÜHRUNG IN DIE NUMERISCHE MATHEMATIK

VORLESUNG VOM SS 2010

MARIO OHLBERGER

Institut für Numerische und Angewandte Mathematik  
Fachbereich Mathematik und Informatik  
Westfälische Wilhelms-Universität Münster

Dieses Skript beruht auf meinen Vorlesungen *Einführung in die Numerische Mathematik* und *Höhere Numerische Mathematik* vom Wintersemester 2007/2008 und Sommersemester 2008 an der Westfälische Wilhelms-Universität Münster.

Es besteht keine Garantie auf Richtigkeit und/oder Vollständigkeit des Manuskripts.

Mario Ohlberger

# Inhaltsverzeichnis

<b>0</b>	<b>Einleitung</b>	<b>1</b>
<b>1</b>	<b>Grundlagen</b>	<b>5</b>
1.1	Normierte Räume . . . . .	5
1.2	Operatoren . . . . .	7
1.3	Banachscher Fixpunktsatz . . . . .	10
1.4	Taylorreihe . . . . .	11
1.5	Approximationsfehler und Fehleranalyse . . . . .	13
<b>2</b>	<b>Lineare Gleichungssysteme</b>	<b>23</b>
2.1	Direkte Verfahren . . . . .	24
2.1.1	Gaußalgorithmus/LR-Zerlegung . . . . .	25
2.1.2	Gauß-Jordan Verfahren . . . . .	32
2.1.3	Cholesky Verfahren für SPD-Matrizen . . . . .	33
2.1.4	LR-Zerlegung für Tridiagonalmatrizen . . . . .	33
2.2	Überbestimmte Gleichungssysteme/Ausgleichsrechnung . . . . .	34
2.2.1	QR-Zerlegung nach Householder . . . . .	38
2.2.2	Singulärwertzerlegung einer Matrix . . . . .	41
2.2.3	Pseudoinverse einer Matrix . . . . .	43
2.3	Iterative Verfahren . . . . .	45
2.3.1	Gesamtschritt Verfahren (GSV)/ Jacobi Verfahren . . . . .	48
2.3.2	Einzelschritt Verfahren (ESV) / Gauß-Seidel-Verfahren . . . . .	52
2.4	Gradientenverfahren . . . . .	53
2.4.1	Eigentliches Gradientenverfahren . . . . .	55
2.4.2	<i>Conjugate Direction</i> Verfahren (CD) . . . . .	57
2.4.3	<i>Conjugate Gradient</i> Verfahren (CG) . . . . .	59
2.5	Zusammenfassung . . . . .	60
<b>3</b>	<b>Nichtlineare Gleichungen/ Nullstellensuche</b>	<b>63</b>
3.1	Verfahren in einer Raumdimension . . . . .	63
3.1.1	Intervallschachtelungsverfahren (ISV) . . . . .	63
3.1.2	Newton Verfahren . . . . .	64
3.1.3	Sekantenverfahren . . . . .	68
3.1.4	Zusammenfassung . . . . .	70
3.2	Konvergenzordnung von Iterationsverfahren . . . . .	71
3.2.1	Verfahren höher Ordnung ( $p = 3$ ) . . . . .	72
3.2.2	Newton-Verfahren für mehrfache Nullstellen . . . . .	73
3.3	Nichtlineare Gleichungssysteme . . . . .	75
3.3.1	Newton-Verfahren für nichtlineare Systeme . . . . .	75

<b>4</b>	<b>Eigenwertprobleme</b>	<b>77</b>
4.1	Grundbegriffe der linearen Algebra und theoretische Grundlagen . . . . .	77
4.2	Kondition des Eigenwertproblems . . . . .	80
4.3	Variationsprinzip für Eigenwerte hermitescher Matrizen . . . . .	83
4.4	Transformation auf Hessenberg-Form . . . . .	86
4.5	Eigenwertbestimmung für Hessenberg-Matrizen . . . . .	88
4.6	Vektoriteration für partielle Eigenwertprobleme . . . . .	91
4.7	Das QR-Verfahren . . . . .	93
<b>5</b>	<b>Approximation</b>	<b>95</b>
5.1	Allgemeine Approximation in normierten Räumen . . . . .	95
5.2	Der Satz von Weierstraß: Approximation durch Polynome . . . . .	100
5.3	Gleichmäßige Approximation / Tschebyschev Approximation . . . . .	104

# Abbildungsverzeichnis

1	Illustration der Vorgehensweise zur Lösung eines Anwendungsproblems . . . . .	1
2	Koordinatentransformation zur mathematischen Betrachtung des Wärmetransports in einem Draht. . . . .	2
1.1	Modellfehler . . . . .	14
1.2	Auswirkung des Datenfehlers. . . . .	15
1.3	Auswirkung des Datenfehlers. . . . .	15
2.1	Ausgleichsgerade . . . . .	37
2.2	Graph, Beispiel 2.33 . . . . .	50
3.1	Newton Verfahren, Beispiel 1 . . . . .	65
3.2	Newton Verfahren, Beispiel 2 . . . . .	66
3.3	Newton Verfahren, Beispiel 3 . . . . .	66
3.4	Sekantenverfahren, geometrische Interpretation . . . . .	69
4.1	Gerschgorinkreise: Beispiel, Radien nicht maßstabsgetreu!! . . . . .	83
5.1	Proximum: Beispiel . . . . .	95
5.2	Proximum: Beispiel 1) . . . . .	96
5.3	Proximum: Beispiel 3) . . . . .	96
5.4	Konvexe und streng konvexe Mengen. . . . .	97

# Kapitel 0

## Einleitung

Die Numerische Mathematik, oder auch Numerik genannt, beschäftigt sich mit der numerischen Lösung endlichdimensionaler Probleme, sowie mit der Approximation unendlichdimensionaler Probleme durch endlichdimensionale. Die Numerik ist somit eine mathematische Schlüsseldisziplin zur Behandlung von Anwendungsproblemen mit Hilfe des Computers.

Der Numerik geht stets die Modellierung voraus, deren Ziel es ist ein Anwendungsproblem in der mathematischen Sprache zu formulieren. Dieses Prozedere wird in der Abbildung 1 verdeutlicht.

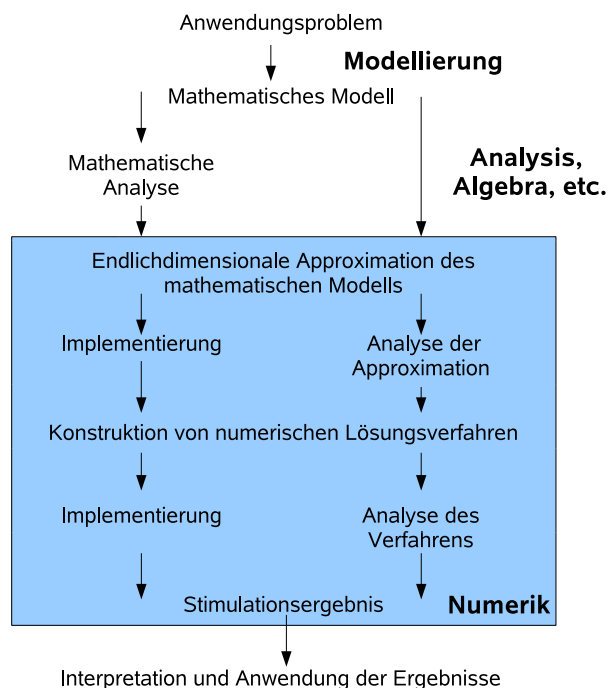


Abbildung 1: Illustration der Vorgehensweise zur Lösung eines Anwendungsproblems

Im folgenden wollen wir an einem einfachen Anwendungsbeispiel dieses Vorgehen skizzieren.

### Beispiel 0.1 (Berechnung des Wärmetransports in einem Draht)

#### Schritt 1: Modellierung

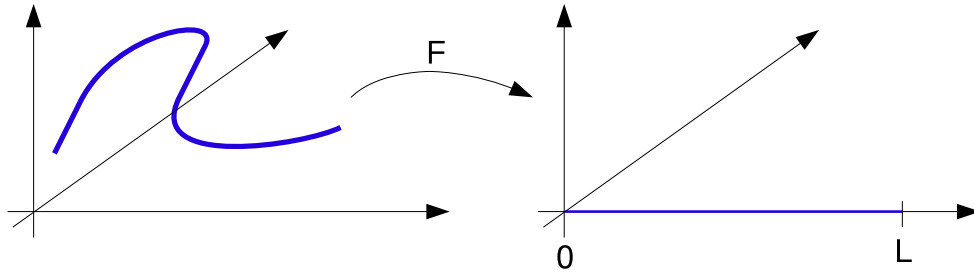


Abbildung 2: Koordinatentransformation zur mathematischen Betrachtung des Wärmetransports in einem Draht.

Wir betrachten den Wärmetransport in einem Draht. Nach einer Koordinatentransformation, wie sie in Abbildung 2 skizziert ist, können wir den Draht eindimensional durch ein Intervall  $I = [0, L]$  repräsentieren. Nach dem Fick'schen Gesetz gilt für die Wärmeleitung, dass der Wärmefluß  $q$  proportional zum Gradienten der Temperatur  $T$  ist, d.h.

$$q(x, t) = -\sigma \partial_x T(x, t), \forall x \in I, \forall t \in [0, T_{\max}].$$

Dabei bezeichnet  $\sigma > 0$  die Wärmeleitfähigkeit und ist eine Materialkonstante. Außerdem gilt für jedes Teilintervall  $[a, b] \in I$  und jeden Zeitabschnitt  $[t_1, t_2] \in [0, T_{\max}]$ :

$$\int_{[a,b]} T(x, t_2) dx - \int_{[a,b]} T(x, t_1) dx = \int_{[t_1, t_2]} q(a, t) dt - \int_{[t_1, t_2]} q(b, t) dt.$$

Sind die Temperatur und der Wärmefluß genügend oft differenzierbar, so folgt mit dem Hauptsatz der Differential- und Integralrechnung

$$\int_{[a,b]} \int_{[t_1, t_2]} \partial_t T(x, t) dx dt = - \int_{[a,b]} \int_{[t_1, t_2]} \partial_x q(x, t) dx dt.$$

Da dies für beliebige  $a, b, t_1, t_2$  gilt, folgt mit dem Hauptsatz der Variationsrechnung

$$\partial_t T(x, t) = -\partial_x q(x, t) = \sigma \partial_{xx} T(x, t), \quad \forall (x, t) \in I \times [0, T_{\max}].$$

Wir erhalten so als mathematisches Modell für die Wärmeleitung in einem Draht eine partielle Differentialgleichung, d.h. eine Gleichung, die die partiellen Ableitungen einer Funktion miteinander in Beziehung setzt. Eine mathematische Analyse zeigt, dass diese sogenannte Wärmeleitungsgleichung eine eindeutige Lösung  $T$  besitzt, falls man beispielsweise die Temperatur zum Zeitpunkt  $t = 0$  und an den Endpunkten  $x = 0, L$  vorgibt. Da die gesuchte Temperaturverteilung  $T$  eine differenzierbare Funktion in zwei Veränderlichen darstellt, handelt es sich hierbei um ein unendlichdimensionales Problem. Im nächsten Schritt wollen wir mit Hilfe von sogenannten Finite Differenzenverfahren eine endlichdimensionale Approximation angeben.

## Schritt 2: Endlichdimensionale Approximation

Zur Approximation der Wärmeleitungsgleichungen führen wir Zerlegungen  $T_h := \{x_i | x_i = hi, i = 0, \dots, N+1\}$  und  $J_k := \{t_n | t_n = kn, n = 0, \dots, M\}$  des Ortsintervalls  $[0, L]$  und des Zeitintervalls  $[0, T_{\max}]$  ein. Dabei ist  $h := L/(N+1)$  die Ortsschrittweite und  $k := T_{\max}/M$  die Zeitschrittweite der jeweiligen Zerlegung.

Die Idee der Finite Differenzen besteht darin, alle Ableitungen in der Wärmeleitungsgleichungen durch Differenzenquotienten zu ersetzen. Verwenden wir z.B.

$$\partial_t T(x, t_n) \approx (T(x, t_n) - T(x, t_{n-1}))/k$$

und

$$\partial_{xx}T(x_i, t) \approx (T(x_{i+1}, t) - 2T(x_i, t) + T(x_{i-1}, t))/h^2,$$

so erhalten wir für die Approximation  $T_i^n$  von  $T(x_i, t_n)$  die Gleichung

$$(T_i^n - T_i^{n-1})/k = \sigma(T_{i+1}^n - 2T_i^n + T_{i-1}^n)/h^2, \quad \forall i = 1, \dots, N, n = 1, \dots, M.$$

Dies ist eine lineare Gleichung für  $T_i^n$ , die jedoch mit den Linearen Gleichungen des selben Typs für die ebenfalls unbekannten Werte  $T_i^{n-1}$  und  $T_{i+1}^n, T_{i-1}^n$  gekoppelt ist. Für die Anfangs- und Randwerte verwenden wir die vorgegebenen Werte der Temperatur, d.h.

$$T_i^0 := T(x_i, 0), \quad T_0^n := T(0, t_n), \quad T_{N+1}^n := T(L, t_n).$$

Berechnet man sukzessive zunächst die Lösung zu den Zeitpunkten  $t_1, t_2, t_3, \dots$ , so erhält man für jeden dieser Zeitschritte  $t_n$  ein lineares Gleichungssystem mit  $N$  Gleichungen für die Unbekannten  $T_i^n$ ,  $i = 1, \dots, N$ . Wir haben das unendlichdimensionale Problem also durch das sukzessive Lösen von linearen Gleichungssystemen approximiert. Verfahren zur Approximation von Differentialgleichungen werden in Kapitel 6 vorgestellt und in der Vorlesung *Höhere Numerische Mathematik* im Wintersemester detailliert analysiert.

### Schritt 3: Numerische Lösung des endlichdimensionalen Problems

Im letzten Lösungsschritt geht es darum die linearen Gleichungssysteme numerisch zu lösen. Hierzu kann z.B. die Gaußelimination verwendet werden. Falls jedoch die Dimension  $N$  des Systems sehr groß wird, sind andere Verfahren besser geeignet. Im zweiten Kapitel der Vorlesung wenden wir uns daher der numerischen Lösung von linearen Gleichungssystemen zu.

Wäre in unserem Beispiel der Wärmeleitfähigkeitskoeffizient temperaturabhängig, d.h.  $\sigma = \hat{\sigma}(T(x, t))$ , mit einer nichtlinearen Funktion  $\hat{\sigma}$ , so wäre das resultierende Gleichungssystem nichtlinear. Solchen Problemen ist das Kapitel 3 gewidmet.

Wählt man zur Approximation der Differentialgleichungen andere Verfahren, wie z.B. Finite Elemente Verfahren, so ist auch die numerische Approximation von Funktionen durch Polynome und die numerische Berechnung von Integralen Bestandteil der Verfahren. Diese Themen werden in der Vorlesung *Höhere Numerische Mathematik* im Wintersemester detailliert analysiert.

Das Beispiel der Wärmeleitung in einem Draht verdeutlicht die Notwendigkeit von numerischen Methoden, wie z.B. die Lösung linearer, oder nichtlinearer Gleichungssysteme, oder die Approximation von Differentialgleichungen. Bevor wir uns diesen Methoden zuwenden, werden im nächsten Kapitel jedoch einige Grundlagen dargestellt. Hierbei handelt es sich zum einen um eine Zusammenstellung wichtiger Begriffe auf der Analysis und Linearen Algebra und zum anderen wird auf die numerische Approximation reeller Zahlen eingegangen und daraus resultierende Fehlerquellen diskutiert.





# Kapitel 1

## Grundlagen

Im folgenden Abschnitt werden wir Definitionen angeben, die wir in den weiteren Kapiteln benötigen.

### 1.1 Normierte Räume

Seien  $\mathbb{K} = \mathbb{R}$  oder  $\mathbb{K} = \mathbb{C}$  ein Körper und  $V$  ein Vektorraum über  $\mathbb{K}$ .

**Definition 1.1 (Norm)**

Eine Abbildung  $\|\cdot\| : V \longrightarrow \mathbb{R}$  heißt **Norm**, falls gilt

- (i)  $\|v\| > 0 \quad \forall v \in V \setminus \{0\}$ ,
- (ii)  $\|\lambda v\| = |\lambda| \|v\| \quad \forall \lambda \in \mathbb{K}, \forall v \in V$ ,
- (iii)  $\|v + w\| \leq \|v\| + \|w\| \quad \forall v, w \in V$ .

**Beispiel 1.2**

Sei  $V = \mathbb{R}^n, v = (v_1, \dots, v_n) \in \mathbb{R}^n$ . Dann ist

$$\|v\|_\infty := \max_{1 \leq i \leq n} |v_i|, \quad \|v\|_1 := \sum_{i=1}^n |v_i|,$$
$$\|v\|_2 = \left( \sum_{i=1}^n |v_i|^2 \right)^{\frac{1}{2}}, \quad \|v\|_p := \left( \sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}} \quad (1 \leq p < \infty).$$

**Beispiel 1.3**

Sei  $V = C^0(I), I = [a, b] \subset \mathbb{R}$ . Dann ist

$$\|v\|_\infty := \sup \{ |v(x)| \mid x \in I \},$$
$$\|v\|_p := \left( \int_a^b |v(x)|^p dx \right)^{\frac{1}{p}}.$$

**Definition 1.4 (Normierter Raum)**

Ein Vektorraum  $V$  zusammen mit einer Norm  $\|\cdot\|$ , geschrieben  $(V, \|\cdot\|)$ , heißt **normierter Raum**.

**Definition 1.5 (Banachraum)**

Eine Folge  $(u_n)_{n \in \mathbb{N}} \subset V$  konvergiert gegen  $u \in V$  :  $\Longleftrightarrow$   
 $\forall \varepsilon > 0 \exists N \forall n > N : \|u_n - u\| < \varepsilon.$

Eine Folge  $(u_n)_{n \in \mathbb{N}} \subset V$  heißt **Cauchy Folge** :  $\Longleftrightarrow$   
 $\forall \varepsilon > 0 \exists N \forall m, l > N : \|u_l - u_m\| < \varepsilon.$

Ein normierter Raum  $(V, \|\cdot\|)$  heißt **vollständig**, falls alle Cauchy-Folgen in  $V$  bzgl.  $\|\cdot\|$  in  $V$  konvergieren. Ein vollständiger normierte Raum heißt **Banachraum**.

**Beispiel 1.6**

$(\mathbb{R}^n, \|\cdot\|)$  ist ein Banachraum für alle  $\|\cdot\|$ ,

$(C^0(I), \|\cdot\|_\infty)$  ist ein Banachraum,  $(C^0(I), \|\cdot\|_p)$  ist dagegen nicht vollständig

**Satz 1.7**

Sei  $\dim V < \infty$ ,  $\|\cdot\|_a$  und  $\|\cdot\|_b$  zwei Normen. Dann existieren  $m, M \in \mathbb{R} : m \|v\|_a \leq \|v\|_b \leq M \|v\|_a \forall v \in V$ , d.h.  $\|\cdot\|_a$  und  $\|\cdot\|_b$  sind **äquivalente Normen**.

**Definition 1.8 (Skalarprodukt)**

Eine Abbildung  $\langle \cdot, \cdot \rangle : V \times V \longrightarrow \mathbb{C}$  heißt **Skalarprodukt**, falls gilt

- (i)  $\forall v \in V \setminus \{0\} : \langle v, v \rangle \geq 0$ ,
- (ii)  $\forall u, v \in V \langle u, v \rangle = \overline{\langle v, u \rangle}$ ,
- (iii)  $\forall u, v, w \in V \forall \alpha \in \mathbb{K} :$   
 $\langle \alpha u, v \rangle = \alpha \langle u, v \rangle$ ,  
 $\langle u + v, w \rangle = \langle u, w \rangle + \langle v, w \rangle.$

*Folgerung:*  $\langle u, \alpha v \rangle = \bar{\alpha} \langle u, v \rangle$ ,  $\langle u, v + w \rangle = \langle u, v \rangle + \langle u, w \rangle.$

**Satz 1.9 (Induzierte Norm)**

Sei  $\langle \cdot, \cdot \rangle$  ein Skalarprodukt, dann wird durch  $\|v\| := \sqrt{\langle v, v \rangle}$  eine Norm induziert.

**Definition 1.10 (Hilbertraum)**

Ein Vektorraum mit Skalarprodukt heißt **Prähilbertraum**, falls  $V$  mit der induzierten Norm **nicht** vollständig ist, sonst bezeichnet  $V$  einen **Hilbertraum**.

**Beispiel 1.11 (Cauchy-Schwarz-Ungleichung)**

$\forall u, v \in V : |\langle u, v \rangle| \leq \sqrt{\langle u, u \rangle \langle v, v \rangle}$ , Gleichheit  $\Longleftrightarrow u, v$  linear abhängig.

**Beispiel 1.12**

Sei  $V = \mathbb{R}^n$ ,  $\langle u, v \rangle := \sum_{i=1}^n u_i v_i$  ist ein Skalarprodukt und induziert die **euklidische Norm**

$$\|v\|_2 := \left( \sum_{i=1}^n |v_i|^2 \right)^{\frac{1}{2}}.$$

## 1.2 Operatoren

### Definition 1.13

$U, V$  normierte Vektorräume,  $D \subseteq U$ . Wir bezeichnen eine Abbildung  $T : D \rightarrow V$  als **Operator**. Dabei gilt:

- (i)  $T$  heißt **stetig** in  $u \in D$  :  $\Leftrightarrow$   
 $\forall \varepsilon > 0 \exists \delta > 0 \forall v \in D : \|u - v\|_U < \delta \implies \|T(u) - T(v)\|_V < \varepsilon.$
- (ii)  $T$  heißt **stetig** in  $D$  :  $\Leftrightarrow$   
 $T$  ist stetig für alle  $u \in D$ .
- (iii)  $T$  heißt **Lipschitz-stetig** :  $\Leftrightarrow$   
es existiert ein  $L > 0 \forall u, v \in D : \|T(u) - T(v)\|_V \leq L \|u - v\|_U.$

### Bemerkung 1.14

Es ist leicht zu sehen, dass aus (iii) (ii) folgt und aus (ii) folgt (i).

### Definition 1.15

$T$  heißt **linearer Operator** (oder einfach **linear**), falls  $\forall u, v \in V, \alpha \in \mathbb{K}$ :

- (i)  $T(u + v) = T(u) + T(v),$
- (ii)  $T(\alpha v) = \alpha T(v).$

### Bemerkung 1.16

Ist  $T$  linear, so schreibt man häufig  $Tu$  statt  $T(u)$ .

### Beispiel 1.17

$V = W = \mathbb{R}^n, A \in \mathbb{R}^{n \times n} : Tu = Au$  ist ein linearer Operator.

$V = C^0(I), W = \mathbb{R} : Tu := \int_a^b u(x) dx$  ist ein linearer Operator.

### Definition 1.18

$T : U \rightarrow V$  sei ein Operator.  $T$  heißt **beschränkt**, falls es ein  $C > 0$  gibt, so dass  
 $\forall u \in U : \|T(u)\|_V \leq C \|u\|_U.$

### Satz 1.19

Für einen linearen Operator  $T : U \longrightarrow V$  sind äquivalent:

- (i)  $T$  ist beschränkt,
- (ii)  $T$  ist Lipschitz-stetig,
- (iii)  $T$  ist stetig in 0.

*Beweis:* Siehe Übungsblatt 1

□

### Bemerkung 1.20

- (i)  $\dim U < \infty, \dim V < \infty$ , dann sind alle linearen Operatoren beschränkt und damit stetig.
- (ii) Auf unendlich-dimensionalen Vektorräumen existieren auch unbeschränkte lineare Operatoren.
- (iii) Die Aussage von Satz 1.19 (Seite 7) ist nur richtig für lineare Operatoren.  
Bsp:  $T : \mathbb{R} \longrightarrow \mathbb{R}, x \longmapsto x^2$  : es existiert keine Konstante  $C$  mit  $|x^2| < C|x| \ \forall x \in \mathbb{R}$ .

### Definition 1.21

Mit  $B(U, V)$  bezeichnen wir den Raum der beschränkten linearen Operatoren.  $B(U, V)$  ist ein Vektorraum. Durch

$$\|T\|_{U,V} := \sup_{\substack{u \in U \\ \|u\|_U = 1}} \|Tu\|_V$$

wird eine Norm auf  $B(U, V)$  definiert. Diese wird als die durch  $\|\cdot\|_U, \|\cdot\|_V$  induzierte **Operatornorm** bezeichnet.

### Folgerung 1.22

- (i)  $\|T\|_{U,V} = \sup_{\substack{u \in U \\ \|u\|_U = 1}} \|Tu\|_V$  (wegen Linearität von  $T$ ).
- (ii)  $\|Tu\|_V \leq \|T\|_{U,V} \|u\|_U$  und  $\|T\|_{U,V}$  ist die kleinste Konstante mit dieser Eigenschaft für alle  $u \in U$  (folgt aus der Definition).
- (iii)  $\|id\|_{U,U} = 1$ , dabei ist  $id \in B(U, U)$  mit  $id : u \longmapsto u$ .

### Beispiel 1.23

$U, V = \mathbb{R}^n$ , dann entspricht  $B(U, V)$  dem Raum der  $n \times n$  Matrizen. Daher wird die Operatornorm auch häufig Matrixnorm genannt.

### Satz 1.24

Die induzierte (Matrix-) Operatornorm ist submultiplikativ, d.h.  $\|A \circ B\| \leq \|A\| \cdot \|B\|$

*Beweis:* (gilt nur für die induzierte Matrixnorm)

$$\|(A \circ B)x\| = \|A(Bx)\| \stackrel{1.22ii}{\leq} \|A\| \|Bx\| \stackrel{1.22ii}{\leq} \|A\| \|B\| \|x\|.$$

$$\text{Sei } x \neq 0 \implies \frac{\|ABx\|}{\|x\|} \leq \|A\| \|B\| \quad \square$$

□

### Bemerkung 1.25

Die induzierten Operatornormen ergeben nicht alle Normen auf  $B(U, V)$ . Sei etwa  $A \in \mathbb{R}^{n \times n}$ ,  $A = (a_{ij})$ , dann wird durch  $\|A\| = \sup_{1 \leq i, j \leq n} |a_{ij}|$  eine Norm definiert, die nicht induziert ist.

### Beispiel 1.26

Die durch  $\|\cdot\|_1$  und  $\|\cdot\|_\infty$  induzierte Operatornormen werden in den Übungen behandelt.

Sei  $A : (\mathbb{R}^n, \|\cdot\|_2) \longrightarrow (\mathbb{R}^n, \|\cdot\|_2)$ , dann gilt  $\|A\|_{2,2} = \sqrt{\lambda_{\max}(A^*A)}$ , wobei  $\lambda_{\max}(B)$  für  $B \in \mathbb{R}^{n \times n}$  den betragsmäßig größten **Eigenwert (EW)** bezeichnet. Sei  $A = (a_{ij})$ , dann ist  $A^* = \overline{A}^\top$ . Diese Norm wird als **Spektralnorm** bezeichnet. Ist  $A \in \mathbb{R}^{n \times n}$ , dann ist  $\overline{A}^\top = A^\top$ .

*Beweis:* **Bemerkungen:**  $(A^*A)^* = A^*A \implies A^*A$  ist hermitesch  $\implies$  alle Eigenwerte sind reell.

Es gilt  $x^*(AA^*)x = (Ax)^*Ax = \langle Ax, Ax \rangle \geq 0$ . Also ist  $A^*A$  positiv definit und somit alle EW positiv.

Da  $A^*A$  hermitesch ist, existiert ein Matrix  $U \in \mathbb{C}^{n \times n}$  mit  $U^*U = id$  (d.h.  $U$  ist unitär) und

$$U^*(A^*A)U = \text{diag}(\lambda_1, \dots, \lambda_n) = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} =: D \quad (*)$$

Sei  $u_i$  die  $i$ -te Spalte von  $U$ , d.h.  $U = (u_1, \dots, u_n)$  und  $\|u_i\|_2 = 1 \forall i \in \{1, \dots, n\}$ , dann ist

$$A^*AU = UD,$$

da  $U^{-1} = U^*$  und somit  $A^*Au_i = \lambda_i u_i$ . Also sind die Vektoren  $u_i$  Eigenvektoren (EV) von  $A^*A$  zu den Eigenwerten  $\lambda_i$ . Aus (\*) folgt dann weiter  $u_i^* A^* A u_i = \lambda_i$ .

Sei  $x \in \mathbb{C}^n$  mit  $\|x\|_2 = 1$ . Setze  $y = U^*x$ , so folgt wegen  $x = Uy$ :

$$\begin{aligned} \|Ax\|_2^2 &= \langle Ax, Ax \rangle = \langle x, A^*Ax \rangle \stackrel{(*)}{=} \langle x, UDU^*x \rangle \\ &= \langle x, U Dy \rangle = \langle U^*x, Dy \rangle = \langle y, Dy \rangle \\ &= \sum_{i=1}^n \overline{y_i} \lambda_i y_i \leq \max_{1 \leq i \leq n} \lambda_i \sum_{i=1}^n y_i^2 \\ &= \lambda_{\max}(A^*A) \|y\|_2^2 = \lambda_{\max}(A^*A). \end{aligned}$$

da  $\|y\|_2 = \|U^*x\| = \|x\| = 1$  ( $U^*$  unitär).

Also gilt  $\|Ax\|_2 \leq \sqrt{\lambda_{\max}(A^*A)}$  für alle  $x$  mit  $\|x\|_2 = 1$  und somit folgt

$$\|A\|_{2,2} \leq \sqrt{\lambda_{\max}(A^*A)}.$$

Sei nun  $\lambda_i = \lambda_{\max}(A^*A)$  der größte Eigenwert und  $u_i$  der zugehörige Eigenvektor mit  $\|u_i\|_2 = 1$ . Da  $\|A\|_{2,2} = \sup_{\|x\|_2=1} \|Ax\|_2$  folgt  $\|A\|_{2,2} \geq \|Au_i\|_2$  und somit

$$\begin{aligned}\|A\|_{2,2}^2 &\geq \|Au_i\|_2^2 = \langle Au_i, Au_i \rangle = \langle u_i, A^* Au_i \rangle \\ &= \langle u_i, \lambda_i u_i \rangle = \lambda_i \langle u_i, u_i \rangle = \lambda_i \|u_i\|_2^2 = \lambda_i.\end{aligned}$$

Also folgt die Behauptung.  $\square$

### 1.3 Banachscher Fixpunktsatz

#### Definition 1.27 (Kontraktion)

Sei  $D \subset X$ ,  $X$  normierter Vektorraum,  $Y$  normierter Vektorraum. Dann heißt ein Operator  $T : D \rightarrow Y$  eine **Kontraktion**, falls  $T$  Lipschitz-stetig mit Lipschitz-Konstante  $0 < L < 1$  ist, d.h.  $\forall u, v \in D : \|T(u) - T(v)\|_Y \leq L \|u - v\|_X$ .

#### Definition 1.28 (Fixpunkt)

Sei  $T : D \rightarrow D$  ein Operator, der  $D$  auf sich selbst abbildet. Dann heißt  $\bar{u} \in D$  **Fixpunkt** von  $T$  in  $D$ , falls  $T(\bar{u}) = \bar{u}$ .

#### Satz 1.29 (Banachscher Fixpunktsatz)

Sei  $X$  ein Banachraum,  $D \subseteq X$  abgeschlossen,  $T : D \rightarrow D$  eine Kontraktion. Dann gilt:

- (i)  $T$  hat genau einen Fixpunkt  $\bar{u} \in D$ .
- (ii) Sei  $u_0 \in D$  beliebig und  $u_{k+1} := T(u_k)$ ,  $k = 0, 1, \dots \implies u_k \rightarrow \bar{u}$ .
- (iii)  $\|\bar{u} - u_k\| \leq L \|\bar{u} - u_{k-1}\|$  ( $k \geq 1$ ), d.h. der Fehler nimmt monoton ab.
- (iv)  $\|\bar{u} - u_k\| \leq \frac{L^k}{1-L} \|T(u_0) - u_0\|$  ( $k \geq 1$ ). (**a-priori Abschätzung**)
- (v)  $\|\bar{u} - u_k\| \leq \frac{L}{1-L} \|u_k - u_{k-1}\|$  ( $k \geq 1$ ). (**a-posteriori Abschätzung**)

*Beweis:* zu a): Wir zeigen zunächst, dass  $T$  höchstens einen Fixpunkt hat. Dazu nehmen wir zunächst an, dass  $T$  zwei Fixpunkte  $\bar{u}, \bar{v}$  hätte. Dann folgt aus  $T\bar{u} = \bar{u}$  und  $T\bar{v} = \bar{v}$ :

$$\|\bar{u} - \bar{v}\| = \|T\bar{u} - T\bar{v}\| \leq L \|\bar{u} - \bar{v}\|.$$

Da nach Voraussetzung  $L < 1$  ist, folgt  $\|\bar{u} - \bar{v}\| = 0$  und somit  $\bar{u} = \bar{v}$ .

Um die Existenz eines Fixpunktes zu beweisen, reicht es Teil b) zu zeigen.

zu b): Die Folge  $(u_k)_{k \in \mathbb{N}}$  sei für beliebiges  $u_0 \in D$  definiert durch die Fixpunktiteration

$$u_{k+1} = Tu_k.$$

Wir zeigen, dass  $(u_k)_{k \in \mathbb{N}}$  Cauchy-Folge ist. Die Konvergenz gegen ein  $\bar{u}$  folgt dann aus der Vollständigkeit von  $X$ . Mit vollständiger Induktion zeigen wir zunächst, dass für  $k \geq 1$  gilt

$$\|u_{k+1} - u_k\| \leq L^k \|T(u_0) - u_0\|. \quad (*)$$

Induktionsanfang ( $k = 1$ ):

Es ist  $\|u_1 - u_0\| = \|Tu_0 - u_0\| = L^0 \|Tu_0 - u_0\|$ .

Induktionsschritt ( $k \rightarrow k+1$ ):

$$\begin{aligned}
 \|u_{k+2} - u_{k+1}\| &= \|Tu_{k+1} - Tu_k\| \\
 &\leq L\|u_{k+1} - u_k\| \\
 \text{Ind. Vor. (*)} \quad &\leq LL^k\|T(u_0) - u_0\| \\
 &= L^{k+1}\|T(u_0) - u_0\|.
 \end{aligned}$$

Sei nun  $m < n$  so folgt hieraus

$$\begin{aligned}
 \|u_n - u_m\| &= \left\| \sum_{k=m}^{n-1} (u_{k+1} - u_k) \right\| \\
 &\leq \sum_{k=m}^{n-1} \|u_{k+1} - u_k\| \\
 &\leq \sum_{k=m}^{n-1} L^k \|T(u_0) - u_0\| \\
 &\leq L^m \|T(u_0) - u_0\| \sum_{k=0}^{n-1-m} L^k \\
 &\stackrel{\text{geom. Reihe}}{\leq} \frac{L^m}{1-L} \|T(u_0) - u_0\|.
 \end{aligned}$$

Da  $L < 1$  ist, konvergiert die rechte Seite gegen Null für  $m \rightarrow \infty$  und wir haben somit gezeigt, dass  $(u_k)_{k \in \mathbb{N}}$  Cauchy-Folge ist. Dies beweist b) und somit auch a).

Die Teile c)-d) werden in den Übungen behandelt. □

### Bemerkung 1.30

Satz 1.29 (iii) (Seite 10) gibt eine a-priori Schranke, die man nutzen kann, um einen Index  $k_0$  zu bestimmen mit  $\|\bar{u} - u_{k_0}\| \leq TOL$  für eine gegebene Toleranz  $TOL > 0$ :

Sei  $TOL$  gegeben, O.B.d.A  $TOL < 1$

$$\begin{aligned}
 \|u_{k_0} - \bar{u}\| &\leq \frac{L^{k_0}}{1-L} \|T(u_0) - u_0\| \leq TOL \\
 \iff L^{k_0} &\leq (1-L) \frac{TOL}{\|T(u_0) - u_0\|} \\
 \iff k_0 \log L &\leq \log(1-L) + \log TOL - \log(\|T(u_0) - u_0\|) \\
 \iff k_0 &\geq \frac{\log(1-L) + \log TOL - \log(\|T(u_0) - u_0\|)}{\log L},
 \end{aligned}$$

da  $0 < L < 1$  und somit  $\log L < 0$ . Meistens ist dies eine Überschätzung des Aufwands.

Satz 1.29 (iv) (Seite 10) kann als Abbruchkriterium während der Iteration benutzt werden, d.h. man bricht ab, falls  $\frac{L}{1-L} \|u_k - u_{k-1}\| < TOL$  ist.

## 1.4 Taylorreihe



**Definition 1.31**

Sei  $C^0(I), I = (a, b)$  der Raum der **stetigen Funktionen auf  $I$** . Mit  $C^m(I) := \{f : I \rightarrow \mathbb{R} \mid f, f', f'', \dots, f^{(m)} \text{ ex. und sind stetig}\}$  bezeichnen wir den Raum der  **$m$ -mal stetig differenzierbaren Funktionen**.

Kurzschreibweise:  $C^m(a, b)$  statt  $C^m((a, b))$ . Mit der Definition  $C^\infty(I) := \bigcap_{m \in \mathbb{N}} C^m(I)$  folgt dann

$$C^\infty(I) \subset \dots \subset C^m(I) \subset \dots \subset C^0(I).$$

**Satz 1.32 (Taylorreihe mit Lagrange Restglied)**

Seien  $f \in C^{m+1}(a, b)$  und  $x_0 \in (a, b)$  fest. Dann existiert für jedes  $x \in (a, b)$  ein  $\xi$  zwischen  $x_0$  und  $x$  mit

$$f(x) = \sum_{k=0}^m \frac{1}{k!} f^{(k)}(x_0)(x - x_0)^k + R_m(x),$$

mit  $R_m(x) := \frac{1}{(m+1)!} f^{(m+1)}(\xi)(x - x_0)^{m+1}$ .

**Satz 1.33 (Taylorreihe mit Integralrestterm)**

Seien  $f \in C^{m+1}(a, b)$ ,  $x_0 \in (a, b)$  fest. Dann gilt für jedes  $x \in (a, b)$

$$f(x) = \sum_{k=0}^m \frac{1}{k!} f^{(k)}(x_0)(x - x_0)^k + R_m(x),$$

mit  $R_m(x) := \frac{1}{m!} \int_{x_0}^x f^{(m+1)}(t)(x - t)^m dt$ .

*Beweis:* Für beide Sätze siehe Analysis I. □

**Folgerung 1.34 (Häufig verwendete Form)**

Seien  $f \in C^{m+1}(x_0 - h_0, x_0 + h_0)$ ,  $x_0 \in \mathbb{R}, h_0 > 0$ . Sei  $|h| \leq h_0$ , dann existiert eine Abbildung  $\omega_m : (-h_0, h_0) \rightarrow \mathbb{R}$  mit  $\lim_{h \rightarrow 0} \omega_m(h) = 0$ , so dass gilt

$$f(x_0 + h) = f(x_0) + \sum_{k=1}^m \frac{f^{(k)}(x_0)}{k!} h^k + \omega_m(h) h^m.$$

*Beweis:* Wende Satz 1.32 (Seite 12) an mit  $x = x_0 + h$ , d.h. es existiert ein  $\xi$  mit  $|\xi| < |h|$  und

$$\begin{aligned} f(x_0 + h) - \sum_{k=0}^{m-1} \frac{f^{(k)}(x_0)}{k!} h^k &= \frac{f^{(m)}(x_0 + \xi)}{m!} h^m \\ &= \frac{f^{(m)}(x_0)}{m!} h^m - \frac{f^{(m)}(x_0)}{m!} h^m + \frac{f^{(m)}(x_0 + \xi)}{m!} h^m \\ &= \frac{f^{(m)}(x_0)}{m!} h^m + \omega_m(h) h^m \end{aligned}$$

mit  $\omega_m(h) = \frac{f^{(m)}(x_0+\xi) - f^{(m)}(x_0)}{m!}$ .

Da  $|\xi| < |h|$  und  $f^{(m)}$  stetig, folgt  $\lim_{h \rightarrow 0} \frac{f^{(m)}(x_0+\xi) - f^{(m)}(x_0)}{m!} = 0$ .

□

### Definition 1.35

Die Funktion  $f \in C^1(x_0 - h_0, x_0 + h_0)$  ist in **erster Näherung** gleich  $f(x_0) + f'(x_0)h$  in einer Umgebung um  $x_0$ , d.h. es existiert ein  $\bar{\omega} : (-h_0, h_0) \rightarrow \mathbb{R}$  mit  $\frac{|\bar{\omega}(h)|}{|h|} \rightarrow 0$  und  $f(x_0 + h) = f(x_0) + f'(x_0)h + \bar{\omega}(h)$ .

**Notation:**  $f(x_0 + h) \stackrel{\bullet}{=} f(x_0) + f'(x_0)h$ .

### Definition 1.36 (Landau Symbole)

Seien  $g, h : \mathbb{R} \rightarrow \mathbb{R}$ . Dann schreiben wir:

(i)  $g(t) = O(h(t))$  für  $t \rightarrow 0 \iff$  es eine Konstante  $C > 0$  und ein  $\delta > 0$  gibt, so dass  $|g(t)| \leq C |h(t)| \quad \forall |t| < \delta$ .

(ii)  $g(t) = o(h(t))$  für  $t \rightarrow 0 \iff$  es ein  $\delta > 0$  und ein  $c : (0, \delta) \rightarrow \mathbb{R}$  gibt, so dass  $|g(t)| \leq c(|t|) |h(t)| \quad \forall |t| < \delta$  und  $c(t) \rightarrow 0$  für  $t \rightarrow 0$ .

**Beispiel:**  $f \in C^1(\mathbb{R})$ , dann ist  $f(x) - (f(x_0) + f'(x_0)(x - x_0)) = o(|x - x_0|)$  wegen Folgerung 1.34 (Seite 12) mit  $h = x - x_0$  und  $m = 1$ .

Ist  $f \in C^2(\mathbb{R})$ , dann ist  $f(x) - (f(x_0) + f'(x_0)(x - x_0)) = O(|x - x_0|^2)$  wegen Satz 1.32 (Seite 12), da  $f''$  beschränkt in einer Umgebung von  $x_0$ , d.h.  $|f''(\xi)| < C$ .

## 1.5 Approximationsfehler und Fehleranalyse

**Problem:** Ein Stahlseil der Länge  $L = 1$  sei an seinen Endpunkten so befestigt, dass es (fast) straff gespannt erscheint. Nun soll die Auslenkung des Seils berechnet werden, wenn sich in der Mitte des Seils ein Seiltänzer befindet.

1. **Modellfehler:** Wir gehen davon aus, dass sich das Seil als Graph einer Funktion  $y : (0, 1) \rightarrow \mathbb{R}$  beschreiben lässt, welche die sogenannte **potentielle Gesamtenergie**:

$$E(y) = \frac{c}{2} \int_0^1 \frac{y'(t)^2}{\sqrt{1 + y'(t)^2}} dt - \int_0^1 f(t)y(t) dt$$

minimiert.

Dabei ist  $c$  eine Materialkonstante und  $f$  die Belastungsdichte.

2. Zur Vereinfachung (Abb 1.1) nehmen wir an, dass  $|y'(t)| \ll 1$ . Dann können wir das Funktional  $E$  vereinfachen zu:

$$\bar{E}(y) = \frac{c}{2} \int_0^1 y'(t)^2 dt - \int_0^1 f(t)y(t) dt.$$

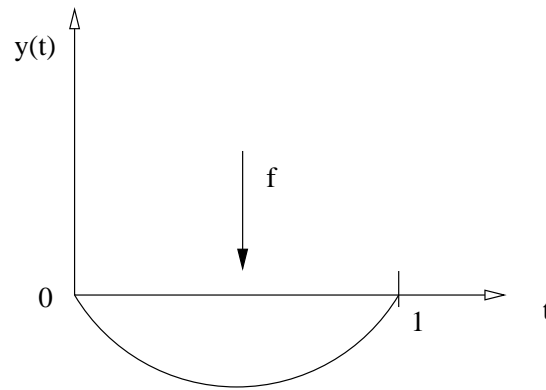


Abbildung 1.1: Modellfehler

Dabei sind eine Reihe von Effekten vernachlässigt worden. Dies führt zu Modellfehlern, die jedoch in dieser Vorlesung nicht weiter betrachtet werden. Wir nehmen an, dass die Minimierung von  $\bar{E}$  das zu lössende Problem sei: Als notwendige und hinreichende Bedingung für die Minimierung von  $\bar{E}$  erhält man durch Variation

$$\left( \frac{d}{d\alpha} \bar{E}(y + \alpha\varphi) \Big|_{\alpha=0} = 0 \quad \forall \text{ "zul" assigge " } \varphi \right)$$

die Differentialgleichung

$$-cy''(t) = f(t), \quad \forall t \in (0, 1)$$

mit Randwerten  $y(0) = y(1) = 0$ .

3. **Datenfehler:**  $c$  ist eine Materialkonstante, die vom Material des Seils abhängt (aber auch von Temperatur und Luftfeuchtigkeit). Der Wert für  $c$  kann nur durch Experimente bestimmt werden, und das ist zwangsläufig fehlerbehaftet. Daher muss sichergestellt werden, dass sowohl  $y$  als auch das numerische Verfahren nicht sensitiv vom konkreten Wert für  $c$  abhängen.

Beispiel: Betrachten wir die Differentialgleichung  $u'(t) = (c - u(t))^2$ ,  $u(0) = 1$   $c > 0$ , so folgt durch Substitution

$$v(t) = \frac{1}{c - u(t)} \implies v'(t) = \frac{u'(t)}{(c - u(t))^2} = 1 \implies u(t) = \frac{1 + tc(c - 1)}{1 + t(c - 1)}.$$

Studieren wir das **Verhalten von  $u$  in Abhängigkeit von  $c$** , so sehen wir:

$$c = 1 : u'(t) = 0 \implies u \equiv 1.$$

$$c > 1 : u' > 0, \text{ d.h. } u \text{ monoton wachsend und } \lim_{t \rightarrow \infty} u(t) = c.$$

$$c < 1 : u' > 0, \lim_{t \rightarrow t_0} u(t) = \infty \text{ für } t_0 = \frac{1}{1-c} > 0.$$

Durch Messfehler oder auch Approximationsfehler kann leicht  $c > 1$  oder  $c < 1$  eintreten, und man erhält qualitativ unterschiedliche Ergebnisse.

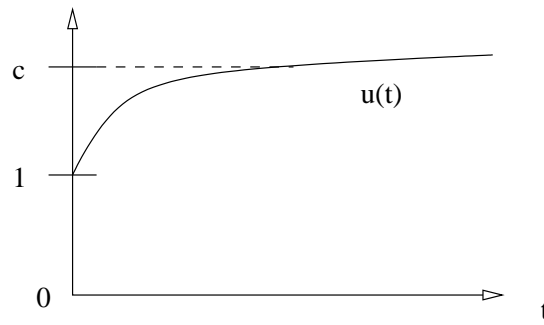


Abbildung 1.2: Auswirkung des Datenfehlers.

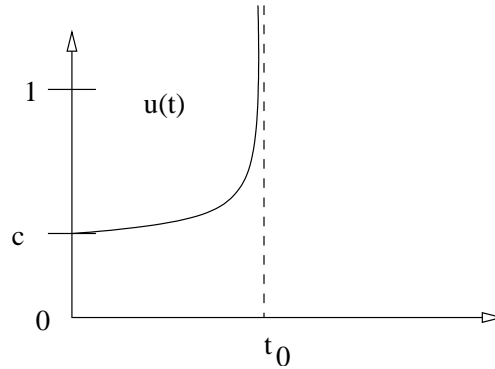


Abbildung 1.3: Auswirkung des Datenfehlers.

4. **Diskretisierungsfehler:** Zurück zu unserem Seiltänzerproblem. In der Numerik müssen wir Ableitungen durch etwas Berechenbares ersetzen, auch können wir  $y(t)$  nicht für alle  $t$  bestimmen. Sei  $N \in \mathbb{N}, x_i := ih, i = 0, \dots, N+1$  mit  $h = \frac{1}{N+1}$ , so approximieren wir auch hier

$$y''(x_i) \approx \frac{1}{h^2} (y(x_{i+1}) - 2y(x_i) + y(x_{i-1})).$$

Setze:  $f_i \equiv f(x_i)$  und sei  $y_i \approx y(x_i)$ , dann ist eine Finite-Differenzen Approximation von  $-cy''(t) = f(t), t \in (0, 1), y(0) = y(1) = 0$  gegeben durch

$$-\frac{c}{h^2} (y_{i+1} - 2y_i + y_{i-1}) = f_i, \quad i = 1, \dots, N, \quad y_0 = y_{N+1} = 0.$$

Die Differenz  $|y_i - y(x_i)|$  ist der Diskretisierungsfehler im Punkt  $x_i$ . Es muss untersucht werden, wie sich dieser Fehler verhält wenn die Gitterweite  $h$  gegen 0 geht.

5. **Lösungsfehler/Abbruchfehler:**

Setze  $A = \begin{pmatrix} 2 & -1 & & 0 \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & -1 \\ 0 & & -1 & 2 \end{pmatrix} \in \mathbb{R}^{N \times N}$  und  $F = (f_i)_{i=1}^N \in \mathbb{R}^N$ , so ergibt sich aus der

Finite Differenzen Diskretisierung das diskrete Problem: Finde  $y_h \in \mathbb{R}^N$  mit

$$\frac{c}{h^2} A y_h = F.$$

Zur Lösung verwenden wir die Identität

$$Dy_h = Dy_h - Ay_h + \frac{h^2}{c}F, \text{ mit } D = \text{diag}(2) = \begin{pmatrix} 2 & & 0 \\ & \ddots & \\ 0 & & 2 \end{pmatrix}.$$

Sei  $y_h^0$  ein beliebiger Startwert (z.B:  $y_h^0 = 0$ ). Zur Berechnung von  $y_h$  betrachten wir folgende Iterationsvorschrift:

$$Dy_h^{n+1} := Dy_h^n - Ay_h^n + \frac{h^2}{c}F$$

bzw.

$$y_h^{n+1} = y_h^n - D^{-1} \left( Ay_h^n + \frac{h^2}{c}F \right).$$

Es muss gezeigt werden, dass  $y_h^n \rightarrow y_h$  für  $n \rightarrow \infty$ .

In der Praxis können wir nur bis zu einem endlichen Wert  $n_0 \in \mathbb{N}$  rechnen. Das heißt die Lösung eines Problems wird  $y_h^{n_0}$  sein. Der Abbruchfehler in der Norm  $\|\cdot\|$  ist  $\|y_h^{n_0} - y_h\|$ .

6. **Rundungsfehler:** Auf einem Rechner kann nur eine endliche Teilmenge von  $\mathbb{R}$  bearbeitet werden. Daher wird nicht  $y_h^{n_0}$  berechnet, sondern die Approximation in dieser endlichen Menge.

### Definition 1.37 (Gleitkommazahl)

Eine Gleitkommazahl zur Basis  $b \in \mathbb{N}$  ist eine Zahl  $a \in \mathbb{R}$  der Form

$$a = \pm [m_1 b^{-1} + \dots + m_r b^{-r}] b^{\pm [e_{s-1} b^{s-1} + \dots + e_0 b^0]}. \quad (*)$$

Man schreibt  $\pm a = 0, m_1 \dots m_r b^{\pm E}$  mit  $E = [e_{s-1} b^{s-1} + \dots + e_0 b^0]$  und  $m_i \in \{0, \dots, b-1\}$ ,  $E \in \mathbb{N}$ ,  $r, s \in \mathbb{N}$  abhängig von der Rechnerarchitektur.

### Bemerkung:

1. Diese Darstellung ermöglicht die gleichzeitige Speicherung sehr unterschiedlich großer Zahlen, wie etwa die Lichtgeschwindigkeit  $c \approx 0.29998 \cdot 10^9 \frac{m}{s}$  oder Elektronenruhemasse  $m_0 \approx 0.911 \cdot 10^{-30}$  kg.
2. Als Normierung nimmt man für  $a \neq 0$  an, dass  $m_1 \neq 0$  ist.
3. Für Computer ist  $b = 2$  üblich, für Menschen  $b = 10$ .

### Definition 1.38 (Maschinenzahlen)

Zu geg.  $(b, r, s)$  sei  $A = A(b, r, s)$  die Menge der  $a \in \mathbb{R}$  mit einer Darstellung  $(*)$ .

$A(b, r, s)$  ist endlich mit größtem und kleinstem positiven Element  $a_{\max} = (1 - b^{-r}) \cdot b^{b^s-1}$ ,  $a_{\min} = b^{-b^s}$ .

Zur Speicherung einer Zahl  $a \in D = [-a_{\max}, -a_{\min}] \cup [a_{\min}, a_{\max}]$  wird eine Rundungsfunktion  $rd: D \rightarrow A$  mit  $rd(a) = \min_{\bar{a} \in A} |\bar{a} - a|$  definiert.

$rd(a)$  wird gespeichert als:<sup>1</sup>

$$\boxed{\pm \mid m_1 \mid \cdots \mid m_r \mid \pm \mid e_0 \mid \cdots \mid e_{s-1} \mid} \quad rd(a) = 0, \underbrace{m_1, \dots, m_r}_{\text{Mantisse } M}, b^{\pm E} \text{ mit } E \text{ als Exponent.}$$

Die heutigen PC benutzen 52 Bits für die Mantisse und 11 Bits für den Exponent; die  $\pm$  werden mit 1 (negativ) und 0 (positiv) dargestellt.

Für  $a \in (-a_{\min}, a_{\min})$  wird in der Regel  $rd(a) = 0$  gesetzt (“underflow”).

Für  $|a| > a_{\max}$  wird von “overflow” geredet. Viele Compiler setzen  $a = NaN$  (not a number) und die Rechnung muss abgebrochen werden.

### Satz 1.39 (Rundungsfehler)

Der absolute Rundungsfehler, der durch Rundung verursacht wird, kann abgeschätzt werden durch

$$|a - rd(a)| \leq \frac{1}{2} b^{-r} \cdot b^E,$$

wobei  $E$  der Exponent von  $a$  ist (in der  $(*)$  Darstellung). Für den relativen Rundungsfehler gilt für  $a \neq 0$

$$\frac{|rd(a) - a|}{|a|} \leq \frac{1}{2} b^{-r+1}.$$

Die Zahl  $eps := \frac{1}{2} b^{-r+1}$  heißt Maschinengenauigkeit.

*Beweis:*  $rd(a)$  weicht maximal eine halbe Einheit in der letzten Mantissenstelle von  $a$  ab. Also  $|a - rd(a)| \leq \frac{1}{2} b^{-r} b^E$ .

Aufgrund der Normalisierung  $m_1 \neq 0$  folgt  $|a| \geq b^{-1} b^E$  und weiter

$$\frac{|rd(a) - a|}{|a|} \leq \frac{\frac{1}{2} b^{-r} b^E}{b^{-1} b^E} = \frac{1}{2} b^{-r+1}. \quad \square$$

Setzt man  $\varepsilon := \frac{rd(a)-a}{a}$ , so folgt  $|\varepsilon| \leq eps$  und  $rd(a) = \varepsilon a + a = a(1 + \varepsilon)$ .  $\square$

### Definition 1.40 (Maschinenoperation)

Die Grundoperation  $\star \in \{+, -, \times, /\}$  wird ersetzt durch  $\otimes$ . In der Regel gilt:

$$a \otimes b = rd(a \star b) = (a \star b)(1 + \varepsilon)$$

mit  $|\varepsilon| \leq eps$ .

**Bemerkung:** Die Verknüpfungen  $\otimes$  erfüllen **nicht** das Assoziativ- bzw. Distributivgesetz.

### Beispiel 1.41

Berechne das Integral  $I_k := \int_0^1 \frac{x^k}{x+5} dx$ .

---

<sup>1</sup>Jedes Kästchen entspricht einem Bit

(A) Es gilt

$$I_0 = \ln(6) - \ln(5)$$

und

$$I_k + 5I_{k-1} = \frac{1}{k} \quad (k \geq 1), \text{ da}$$

$$\int_0^1 \frac{x^k}{x+5} + 5 \frac{x^{k-1}-1}{x+5} = \int_0^1 x^{k-1} dx = \frac{1}{k}.$$

Bei einer Berechnung mit nur 3 Dezimalstellen ( $r = 3, b = 10$ ) ergibt sich:

$$\begin{aligned} \bar{I}_0 &= 0.182 \cdot 10^0 \\ \bar{I}_1 &= 0.900 \cdot 10^{-1} \\ \bar{I}_2 &= 0.500 \cdot 10^{-1} \\ \bar{I}_3 &= 0.833 \cdot 10^{-1} \\ \bar{I}_4 &= -0.166 \cdot 10^0 \end{aligned}$$

Dabei bezeichnet  $\bar{I}_k$  den berechneten Wert unter Berücksichtigung der Rundungsfehler. Die Berechnung ist fehlerhaft. Offensichtlich sind die  $I_k$  monoton fallend, da  $I_k \searrow 0$  ( $k \rightarrow \infty$ ), aber es gibt widersprüchliche Ergebnisse (siehe  $I_3$ ). Auf einem Standard PC ergab:  $\bar{I}_{21} = -0.158 \cdot 10^{-1}$  und  $\bar{I}_{39} = 8.960 \cdot 10^{10}$ .

Dies ist ein Beispiel für **Fehlerfortpflanzung**, da der Fehler in  $I_{k-1}$  mit 5 multipliziert wird, um  $I_k$  zu berechnen.

(B) Berechnet man die Werte  $I_k$  exakt, so ergibt sich bei einer Rundung auf drei Dezimalstellen  $I_9 = I_{10}$  und eine Rückwärtsiteration  $I_{k-1} = \frac{1}{5} \left( \frac{1}{k} - I_k \right)$  ergibt:

$$\begin{aligned} \bar{I}_4 &= 0.343 \cdot 10^{-1} \\ \bar{I}_3 &= 0.431 \cdot 10^{-1} \\ \bar{I}_2 &= 0.500 \cdot 10^{-1} \\ \bar{I}_1 &= 0.884 \cdot 10^{-1} \\ \bar{I}_0 &= 0.182 \cdot 10^0 \end{aligned}$$

Hier tritt **Fehlerdämpfung** auf.

### Beispiel 1.42

Zu lösen ist das LGS

$$\begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0.86419999 \\ 0.14400001 \end{pmatrix} =: b.$$

Die exakte Lösung ist  $\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 0.9911 \\ -0.4870 \end{pmatrix}$ .

Durch Messfehler oder auch Rundung erhalten wir eine rechte Seite

$$\bar{b} = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix}$$

Dann ist die Lösung  $\begin{pmatrix} \tilde{x}_1 \\ \tilde{x}_2 \end{pmatrix} = \begin{pmatrix} 2 \\ -2 \end{pmatrix}$ , d.h. wir erhalten ca. 100% Abweichung.

Dies bedeutet, dass kleine Änderungen der Eingabedaten zu großen Änderungen der Lösungen führen können. In diesem Kontext führen wir den Begriff der Kondition eines Problems ein.

### Definition

Eine numerische Aufgabe (z.B. effizientes Lösen eines LGS oder Integrals) heißt **gut konditioniert**, falls kleine Änderungen der Eingabedaten zu kleinen Änderungen der Lösung führen; sonst heißt das Problem **schlecht konditioniert**.

Präzisieren wir: Was ist eine numerische Aufgabe? Was heißt klein?

Die Matrix

$$A := \begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix}$$

sollte schlecht konditioniert sein.

Im folgenden 2 Ansätze:

1. Für einfache Probleme.
2. Für etwas komplexere Probleme.

### Definition 1.43

Sei  $f : U \rightarrow \mathbb{R}^n$  mit  $U \subset \mathbb{R}^m$  und sei  $x_0 \in U$  vorgegeben. Dann versteht man unter der Aufgabe  $(f, x_0)$  die effektive Berechnung von  $f$  an der Stelle  $x_0$ . Dabei sind  $x_0$  die Eingabedaten.

Beispiel:  $Ax = b$ ,  $(f, b)$  mit  $f(b) = A^{-1}b$

### Satz 1.44

Sei  $x_0 = (x_1, \dots, x_m)$  und  $x_0 + \Delta x \in U$  eine Störung der Eingabedaten mit  $\|\Delta x\| \ll 1$ . Falls  $f : U \rightarrow \mathbb{R}$  (d.h.  $n = 1$ ) einmal stetig differenzierbar, so ist der Ergebnisfehler  $\Delta f(x_0) = f(x_0) - f(x_0 + \Delta x)$  in erster Näherung gleich

$$\sum_{j=1}^m \frac{\partial f}{\partial x_j}(x_0) \Delta x_j = \nabla f(x_0) \Delta x.$$

Für den relativen Fehler gilt in erster Näherung

$$\frac{\Delta f(x_0)}{f(x_0)} \doteq \sum_{j=1}^m \left( \frac{\partial f}{\partial x_j}(x_0) \frac{x_j}{f(x_0)} \right) \frac{\Delta x_j}{x_j}.$$

### Definition 1.45 (Konditionszahlen I)

Wir nennen den Faktor  $k_j := \frac{\partial f}{\partial x_j}(x_0) \frac{x_j}{f(x_0)}$  (relative) **Konditionszahl**.



*Beweis:* (Beweis von Satz 1.44)

Wie in der Folgerung 1.34 (Seite 12) kann man hier den Satz von Taylor anwenden:

$$f(x_0 + \Delta x) = f(x_0) + \nabla f(x_0) \Delta x + \bar{\omega}(\|\Delta x\|)$$

mit  $\bar{\omega}(\|\Delta x\|) = o(\|\Delta x\|) \implies$  Behauptung für den absoluten Fehler.

**Bemerkung:**  $k_j$  beschreibt, wie der relative Fehler in den Eingabedaten  $x_j$  verstärkt bzw. abgeschwächt wird.  $\square$

#### Definition 1.46

Wir nennen das Problem  $(f, x_0)$  **gut konditioniert**, falls alle  $k_j$  ( $j = 1, \dots, m$ ) klein sind, sonst **schlecht konditioniert**.

#### Beispiel 1.47 (Arithmetische Operationen)

$$(i) \ f(x_1, x_2) = x_1 x_2, \ k_1 = \frac{\partial f}{\partial x_1}(x_1, x_2) \frac{x_1}{f(x_1, x_2)} = \frac{x_2 x_1}{x_1 x_2} = 1$$

Analog für  $k_2$  ergibt sich ebenfalls 1  $\implies$  Multiplikation ist gut konditioniert.

(ii) Division ist gut konditioniert.

(iii) Addition  $f(x_1, x_2) = x_1 + x_2$ :

$$k_j = 1 \frac{x_j}{x_1 + x_2} = \frac{x_j}{x_1 + x_2}.$$

$k_j$  wird beliebig groß, wenn  $x_1 x_2 < 0$  und  $x_1$  und  $x_2$  betragsmäßig gleich groß sind. Das heißt, in diesem Fall ist die Addition schlecht konditioniert, ansonsten ist sie gut konditioniert.

(iv) Subtraktion ist schlecht konditioniert, falls  $x_1 x_2 > 0$  und  $x_1$  und  $x_2$  betragsmäßig gleich groß sind.

**Beispiel:**  $(n = 3)x = 0.9995 \quad y = 0.9984 \quad rd(x) = 0.1 \cdot 10^1 \quad rd(y) = 0.998 \cdot 10^0$  Dann gilt für  $\circledast = -$

$$x \circledast y = rd(1 - 0.998) - rd(0.2 \cdot 10^{-2}) = 0.2 \cdot 10^{-2}$$

Der absolute Fehler beträgt  $x \circledast y - (x - y) = 0.0001$

Der relative Fehler beträgt  $\frac{x \circledast y - (x - y)}{(x - y)} = 0.82$

Das Problem wird als Auslöschung bezeichnet.

Bei komplexeren Problemen (etwa  $n > 1$ ) betrachten wir einen anderen Ansatz:

**Definition 1.48**

Das Problem  $(f, x_0)$  ist **wohlgestellt** in

$$B_\delta(x_0) := \{x \in U \mid \|x - x_0\| < \delta\}$$

falls es eine Konstante  $L_{abs} \geq 0$  gibt, mit

$$\|f(x) - f(x_0)\| \leq L_{abs} \|x - x_0\| \quad (*)$$

für alle  $x \in B_\delta(x_0)$ . Gibt es keine solche Konstante, so heißt das Problem **schlecht gestellt**.

Sei im folgenden  $L_{abs}(\delta)$  die kleinste Zahl mit der Eigenschaft (\*).

Analog sei  $L_{rel}(\delta)$  die kleinste Zahl mit

$$\frac{\|f(x) - f(x_0)\|}{\|f(x_0)\|} \leq L_{rel}(\delta) \frac{\|x - x_0\|}{\|x_0\|}.$$

**Definition 1.49 (Konditionszahlen II)**

Wir definieren  $K_{abs} := \lim_{\delta \searrow 0} L_{abs}(\delta)$  die **absolute Konditionszahl** und  $K_{rel} := \lim_{\delta \searrow 0} L_{rel}(\delta)$  die **relative Konditionszahl**.

**Bemerkung:** Falls  $f$  differenzierbar, so gilt

$$K_{rel} = \|f'(x_0)\| \frac{\|x_0\|}{\|f(x_0)\|}.$$

**Beachte:**  $f'(x_0)$  ist eine Matrix und  $\|f'(x_0)\|$  eine Matrixnorm.  $K_{rel}$  hängt von der Wahl der Normen ab.

**Beispiel 1.50 (Konditionierung eines LGS)**

Zu lösen ist  $Ax = b$ , d.h.  $f(b) = A^{-1}b$  und  $f'(x) = A^{-1}$ .

Damit folgt  $K_{abs} = \|A^{-1}\|$ ; und hieraus mit  $Ax = b$  und der Submultiplikativität der zugeordneten Norm:

$$K_{rel} = \|A^{-1}\| \frac{\|b\|}{\|A^{-1}b\|} = \|A^{-1}\| \frac{\|Ax\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|A\| \cdot \|x\|}{\|x\|} = \|A^{-1}\| \cdot \|A\|.$$

Wir definieren entsprechend die Kondition der Matrix  $A$  durch

$$\text{cond}(A) := \|A^{-1}\| \cdot \|A\|.$$

Beachte, dass ein  $x \in \mathbb{R}^m$  existiert mit  $\|Ax\| = \|A\| \|x\|$ , d.h.  $\text{cond}(A)$  ist eine gute Abschätzung für die Konditionierung vom Problem  $(f, b)$

Mit  $A$  wie in Beispiel 1.42 gilt:  $\text{cond}(A) = \|A^{-1}\| \|A\| \approx 10^9$ . Das Problem ist also schlecht konditioniert.



## Kapitel 2

# Lineare Gleichungssysteme

Wir werden in diesem Kapitel Probleme der Form

$$Ax = b$$

betrachten, wobei  $A \in \mathbb{R}^{n \times n}$  und  $x, b \in \mathbb{R}^n$ . Es gibt im Wesentlichen 2 Klassen von Verfahren

1. Direkte Verfahren
2. Iterative Verfahren

Aus der Schule (und den Lineare Algebra-Vorlesungen) ist uns ein direkte Verfahren bekannt, das Gaußsche Eliminationsverfahren. Für kleine Gleichungssysteme eignet sich dieses Verfahren, jedoch kann das Verfahren für  $n \gg 1000$  sehr ineffizient werden, da das Verfahren einen Rechenaufwand der Ordnung  $n^3$  hat. Aus diesem Grund werden wir andere Verfahren kennenlernen, mit denen man schneller ans Ziel kommen kann.

Wir werden Probleme folgender Art behandeln:

- (A) Geg:  $A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n$   
Ges:  $x \in \mathbb{R}^n$  mit  $Ax = b$  (falls eine Lösung existiert).
- (B) Geg:  $A \in \mathbb{R}^{n \times n}, b_1, \dots, b_l \in \mathbb{R}^n$   
Ges:  $x_i \in \mathbb{R}^n$  mit  $Ax_i = b_i$  ( $i = 1, \dots, l$ ) (falls Lösungen existieren).
- (C) Geg:  $A \in \mathbb{R}^{n \times n}$   
Ges:  $A^{-1}$  (falls die Inverse existiert).

Es sind äquivalent

- (i)  $\exists! x \in \mathbb{R}^n : Ax = b$ .
- (ii)  $Ax = 0 \iff x = 0$ .
- (iii)  $\det(A) \neq 0$ .
- (iv) 0 ist kein Eigenwert von  $A$ .
- (v)  $A$  ist regulär, d.h.  $\exists B \in \mathbb{R}^{n \times n}$  mit  $AB = BA = E_n$ . Dabei ist  $B = A^{-1}$  und  $x = A^{-1}b$  ist die eindeutige Lösung von  $Ax = b$ .

Alle diese Probleme sind äquivalent, aber es existieren Verfahren, die besonders geeignet für eines dieser Probleme sind.

## Verfahren

1. Direkte Verfahren liefern die exakte Lösung  $x$  nach endlich vielen Schritten (bis auf Rundungsfehler). Beispiele dafür sind der *Gaußalgorithmus* mit Aufwand  $O(n^3)$  und die *Cramersche Regel* mit Aufwand  $O(n!)$ . Der minimale theoretische Aufwand liegt bei  $O(n^2)$ , aber es existiert kein direktes Verfahren mit dieser Komplexität.

Der Vorteil direkter Verfahren ist, dass  $A^{-1}$  in der Regel mitbestimmt wird und somit der Aufwand für (A), (B) und (C) ungefähr gleich groß ist.

Ein Nachteil ist, dass während der laufenden Berechnung keine Näherung vorliegt, d.h. das Resultat steht erst nach Abarbeitung des Algorithmus, also erst nach  $n$  Schritten, fest. Je nach Anwendung sind diese Verfahren viel zu aufwändig und deshalb besonders ungeeignet für Problem (A), wenn  $n$  sehr groß ist.

2. Iterative Verfahren liefern nach endlich vielen Schritten eine beliebig genaue Approximation der Lösung (bis auf Rundungsfehler).

Der Vorteil liegt darin, dass man in der Lage ist, die Lösung so genau zu bestimmen, wie es nötig ist. Häufig hat man bereits eine brauchbare Lösung nach  $k \ll n$  Schritten

### Satz 2.1 (Störungssatz für lineare Gleichungssysteme)

Sei  $A \in \mathbb{R}^{n \times n}$  regulär und  $\|\cdot\|$  die induzierte Matrixnorm. Sei  $\Delta A \in \mathbb{R}^{n \times n}$  gegeben mit  $\|\Delta A\| < \frac{1}{\|A^{-1}\|}$  und sei  $b \in \mathbb{R}^n$  und  $\Delta b \in \mathbb{R}^n$ . Dann ist  $A + \Delta A$  regulär und es gilt

$$\frac{\|x - \bar{x}\|}{\|x\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\Delta A\|}{\|A\|}} \left( \frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right).$$

Dabei ist  $Ax = b$  und  $\bar{x}$  die Lösung des von  $(A + \Delta A)\bar{x} = b + \Delta b$ .

**Bemerkung:**  $\text{cond}(A) := \|A\| \|A^{-1}\|$  ist der entscheidende Verstärkungsfaktor für den relativen Fehler.

## 2.1 Direkte Verfahren

Idee: Hat  $A$  eine einfache Gestalt, so lässt sich  $x$  leicht bestimmen.

### Beispiel 2.2 (Dreiecksmatrizen)

Sei  $A \in \mathbb{R}^{n \times n}$  eine obere Dreiecksmatrix ( $\Delta$ -Matrix), d.h.  $a_{ij} = 0$  für  $i > j$ , oder

$$A = \begin{pmatrix} * & \cdots & * \\ & \ddots & \vdots \\ 0 & & * \end{pmatrix}.$$

Dann gilt  $\det(A) = \prod_{i=1}^n a_{ii}$ , d.h.  $A$  ist regulär  $\iff a_{ii} \neq 0 \forall i \in \{1, \dots, n\}$ . Ist  $A$  regulär, so ist  $Ax = b$  lösbar. Aus

$$b_i = \sum_{u=1}^n a_{iu}x_u = \sum_{u=i}^n a_{iu}x_u$$

erhalten wir den Algorithmus:

$$\begin{aligned} i = n : & \quad x_n = \frac{b_n}{a_{nn}}, \\ i < n : & \quad x_i = \frac{1}{a_{ii}} \left( b_i - \sum_{u=i+1}^n a_{iu}x_u \right). \end{aligned}$$

Frage: Kann eine beliebige reguläre Matrix  $A$  so umgeformt werden, dass sie obere  $\Delta$ -Gestalt hat? D.h. gesucht ist  $\tilde{A} \in \mathbb{R}^{n \times n}$  mit oberer  $\Delta$ -Gestalt,  $\tilde{b} \in \mathbb{R}^n$ , so dass  $\tilde{A}x = \tilde{b}$  dieselbe Lösung hat wie  $Ax = b$ .

Eine Lösung dieses Problems liefert der Gaußalgorithmus:

### 2.1.1 Gaußalgorithmus/LR-Zerlegung

Der Algorithmus startet mit

$$(A, b) = (A^{(0)}, b^{(0)}) = \left( \begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ a_{21} & a_{22} & \cdots & a_{2n} & b_2 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} & b_n \end{array} \right)$$

und führt durch sukzessive Manipulation auf  $(A^{(p)}, b^{(p)})$ ,  $p = 1, \dots, n-1$ , so dass aus  $A^{(p-1)}x = b^{(p-1)}$  folgt  $A^{(p)}x = b^{(p)}$ . Um  $(A^{(1)}, b^{(1)})$  zu berechnen, wird zur  $i$ -ten Zeile für  $i = 2, \dots, n$  das  $a_{i1}^{(0)}/a_{11}^{(0)}$ -fache der ersten Zeile hinzuaddiert. Wir erhalten somit

$$\begin{aligned} (A^{(1)}, b^{(1)}) &= \left( \begin{array}{cccc|c} a_{11} & a_{12} & \cdots & a_{1n} & b_1 \\ 0 & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} & b_2^{(1)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & a_{n2}^{(1)} & \cdots & a_{nn}^{(1)} & b_n^{(1)} \end{array} \right) \\ &\quad \downarrow \\ (A^{(p-1)}, b^{(p-1)}) &= \left( \begin{array}{cccc|cccc|c} a_{11} & \cdots & \cdots & \cdots & \cdots & \cdots & a_{1n} & b_1 \\ 0 & a_{22}^{(1)} & \cdots & \cdots & \cdots & \cdots & a_{2n}^{(1)} & b_2^{(1)} \\ \vdots & 0 & \ddots & \cdots & \cdots & \cdots & a_{in}^{(i-1)} & b_i^{(i-1)} \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & 0 & a_{pp}^{(p-1)} & \cdots & a_{pn}^{(p-1)} & b_p^{(p-1)} \\ \vdots & \vdots & \vdots & \vdots & a_{(p+1)(p+1)}^{(p-1)} & \cdots & a_{(p+1)n}^{(p-1)} & b_{(p+1)}^{(p-1)} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 0 & a_{n(p+1)}^{(p-1)} & \cdots & a_{nn}^{(p-1)} & b_n^{(p-1)} \end{array} \right) \end{aligned}$$

Wir erhalten schließlich  $(A^{(n-1)}, b^{(n-1)})$ , wobei  $A^{(n-1)}$  eine obere  $\Delta$ -Matrix ist und  $Ax = b \iff A^{(n-1)}x = b^{(n-1)}$ .

Unsere Rechnung setzt voraus, dass stets gilt  $a_{pp}^{(p-1)} \neq 0$  gilt, ansonsten müssen zuerst Zeilen vertauscht werden. Weder das Vertauschen von Zeilen noch der Eliminationsschritt verändern die Lösung. Kann in einem Schritt  $a_{pp}^{(p-1)} \neq 0$  nicht erreicht werden, nachdem man sämtliche Zeilen vertauscht hat, so bedeutet dies, dass  $A$  singulär ist. Das Zeilenvertauschen wird als **Pivotisierung** bezeichnet. Häufig wird die Zeile ausgesucht mit

$$|a_{kp}^{(p-1)}| = \max_{p \leq i \leq n} |a_{ip}^{(p-1)}|$$

und wird als **Teilpivotisierung** oder **Spaltenpivotisierung** bezeichnet.

### Beispiel 2.3

Sei  $A = \begin{pmatrix} \varepsilon & 1 \\ 1 & 1 \end{pmatrix}$  und  $b = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \implies x = \begin{pmatrix} \frac{1}{1-\varepsilon} \\ \frac{1-2\varepsilon}{1-\varepsilon} \end{pmatrix} \approx \begin{pmatrix} 1 \\ 1 \end{pmatrix}$  für  $\varepsilon \ll 1$ .

Ohne Pivotisierung folgt:  $(A^{(1)}, b^{(1)}) = \left( \begin{array}{cc|c} \varepsilon & 1 & 1 \\ 0 & 1 - \frac{1}{\varepsilon} & 2 - \frac{1}{\varepsilon} \end{array} \right)$

Es folgt  $x_2 = \frac{2-\varepsilon^{-1}}{1-\varepsilon^{-1}} \approx 1$  und  $x_1 = (1 - x_2)\varepsilon^{-1} \approx 0$ , da auf einem Computer  $rd(2 - \varepsilon^{-1}) = -\varepsilon^{-1}$ ,  $rd(1 - \varepsilon^{-1}) = -\varepsilon^{-1}$  berechnet werden.

Mit Pivotisierung folgt hingegen nach Zeilentausch:

$$\left( \begin{array}{cc|c} 1 & 1 & 2 \\ \varepsilon & 1 & 1 \end{array} \right)$$

und schließlich nach der Elimination

$$\left( \begin{array}{cc|c} 1 & 1 & 2 \\ 0 & 1 - \varepsilon & 1 - 2\varepsilon \end{array} \right)$$

Hier folgt also  $x_2 = \frac{1-2\varepsilon}{1-\varepsilon} \approx 1$  und  $x_1 = 2 - x_2 \approx 1$ , da auf einem Computer  $rd(1 - 2\varepsilon) = 1$  für sehr kleines  $\varepsilon$  gilt.

Das äquivalente Problem

$$\begin{pmatrix} 1 & \varepsilon^{-1} \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \varepsilon^{-1} \\ 2 \end{pmatrix}$$

kann durch Spaltenpivotisierung nicht gelöst werden. Hier muss **total pivoting** benutzt werden, d.h. sind die Matrixeinträge sehr unterschiedlich groß, so müssen auch die Spalten vertauscht werden. Das Vertauschen der Spalten ist jedoch umständlich und wird selten angewandt. Man vertauscht die Spalten nur dann, wenn man keine andere Wahl hat.

Wir fassen das Gaußverfahren, wie folgt zusammen.

**Algorithmus 2.4 (Gaußverfahren)**

Setze  $q_i = i$  ( $i = 1, \dots, n$ ).

Für  $p = 1, \dots, n-1$ :

Wähle  $j \in \{p, \dots, n\}$  mit  $|a_{q_j p}| = \max_{k=p, \dots, n} |a_{q_k p}|$ .

Vertausche die Zeilen  $q_j \longleftrightarrow q_p$  [Spaltenpivotisierung]

Für  $k = p+1, \dots, n$ :

Falls  $a_{q_p p} = 0 \implies$  Abbruch.

Setze  $l = \frac{a_{q_k p}}{a_{q_p p}}$  [Multiplikationsfaktor]

Setze  $a_{q_k p} = l$  [Speichere  $l$  statt  $a_{q_k p} = 0$ ]

Für  $j = p+1, \dots, n$ :

Setze  $a_{q_k j} = a_{q_k j} - l \cdot a_{q_p j}$  [Matrix  $A^{(p)}$ ]

Setze  $b_{q_k} = b_{q_k} - l \cdot b_{q_p}$  [Vektor  $b^{(p)}$ ]

Die Lösung von  $Ax = b$  wird anschließend durch Rückwärtseinsetzen wie folgt gelöst:

Setze  $x_n = b_{q_n} / a_{q_n n}$ .

Für  $k = n-1, \dots, 1$ :

$$x_k = \left( b_{q_k} - \sum_{i=k+1}^n a_{q_k i} x_i \right) / a_{q_k k}.$$

**Bemerkungen:**

- (i) Der Aufwand des Algorithmus liegt bei  $\frac{1}{3}n^3 + O(n^2)$ .
- (ii) Anstelle der entstehenden Nullen wird der Multiplikationsfaktor  $l$  gespeichert.
- (iii) Die Matrix  $A$  und der Vektor  $b$  werden überschrieben. Es ist deshalb ratsam, eine Kopie der Vektoren zu machen.
- (iv) Die Pivotisierung wird als Vektor gespeichert, und die Zeilenvertauschung im Speicher wird nicht durchgeführt.

**Formaler Zugang:****1) Uminterpretation der Pivotisierung:**

Im  $i$ -ten Schritt des Gaußalgorithmus werden Zeilen vertauscht  $i \longleftrightarrow k$ ,  $k > i$ . Zur Umformulierung



betrachten wir folgende Matrix.

$$P_{ik} := \begin{pmatrix} 1 & & & & & & \\ & \ddots & & & & & \\ & & 1 & & & & \\ & & & 0 & \cdots & 1 & \leftarrow i \\ & & & \vdots & \ddots & \vdots & \\ & & & 1 & \cdots & 0 & \leftarrow k \\ & & & & & & 1 \\ & & & & & & \ddots \\ & & & & & & & 1 \end{pmatrix}$$

### Lemma 2.5

Für die Matrizen  $P_{ik}$  gilt:

- (i)  $B = P_{ik}A$  entspricht der Matrix nach der Vertauschung der  $i$ -ten und  $k$ -ten Zeile.
- (ii)  $B = AP_{ik}$  entspricht der Matrix nach der Vertauschung der  $i$ -ten und  $k$ -ten Spalte.
- (iii)  $P_{ik}^2 = E_n$ , d.h.  $P_{ik}^{-1} = P_{ik}$ .

*Beweis:* Durch nachrechnen. □

### Definition 2.6

Eine Matrix  $P \in \mathbb{R}^{n \times n}$  heißt **Permutationsmatrix**, falls  $P$  durch Zeilenvertauschungen aus der Einheitsmatrix  $E_n$  entsteht.

**Bemerkung:**  $P_{ik}$  ist eine Permutationsmatrix. Ist  $P = P_{l_m, k_m} \cdots P_{l_1, k_1}$ , eine Permutationsmatrix, so gilt  $P^{-1} = P_{l_1, k_1} \cdots P_{l_m, k_m}$ .

## 2) Berechnung im Gaußverfahren:

Wir betrachten die Matrix:

$$L_i := \left( \begin{array}{ccc|ccc} 1 & & & 0 & & 0 \\ & \ddots & & \vdots & & \\ & & 1 & & & \\ & & & l_{i+1,i} & \ddots & \\ & & & \vdots & & \\ 0 & & & l_{n,i} & 0 & 1 \end{array} \right)$$

mit  $l_{ji} := \frac{a_{ji}}{a_{ii}}$ ,  $j = i+1, \dots, n$ .

### Lemma 2.7

- (i) Sei  $B = L_i A = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$ , wobei  $b_i \in \mathbb{R}^n$  Zeilenvektoren sind. Dann gilt
 
$$b_j = a_j \quad (j = 1, \dots, i) \quad \text{und} \quad b_j = a_j + l_{ji} a_i \quad (j = i+1, \dots, n).$$

(ii)  $L_i^{-1} = 2E_n - L_i$ , also

$$L_i^{-1} := \left( \begin{array}{ccc|ccc} 1 & & & 0 & & 0 \\ & \ddots & & \vdots & & \\ & & 1 & & & \\ & & -l_{i+1,i} & \ddots & & \\ & & \vdots & & \ddots & \\ 0 & & -l_{n,i} & 0 & & 1 \end{array} \right)$$

d.h. für  $B = L_i^{-1}A = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}$  gilt

$$b_j = a_j, \quad (j = 1, \dots, i) \quad \text{und} \quad b_j = a_j - l_{ji}a_i, \quad (j = i+1, \dots, n).$$

*Beweis:* Durch nachrechnen. □

### Folgerung 2.8

Die Transformation von  $A$  auf obere  $\triangle$ -Gestalt kann geschrieben werden als

$$R = L_{n-1}^{-1}P_{n-1} \cdots L_1^{-1}P_1A.$$

Dabei ist  $P_i = P_{ij}$  für ein  $j \geq i$  und  $R$  eine obere  $\triangle$ -Matrix ist.

### Satz 2.9 (LR-Zerlegung)

Sei  $A \in \mathbb{R}^{n \times n}$  eine reguläre Matrix. Dann gilt:

- Es existiert eine Permutationsmatrix  $P$ , eine untere  $\triangle$ -Matrix  $L$  mit Diagonalelementen 1 und eine obere  $\triangle$ -Matrix  $R$  mit

$$PA = LR.$$

- Es gilt: Ist  $A = LR = MS$ , wobei  $L, M$  untere  $\triangle$ -Matrizen mit Diagonalelementen 1 sind, und  $R, S$  obere  $\triangle$ -Matrizen sind, so folgt  $L = M$ ,  $R = S$ .

**Bemerkung:** Ist  $PA = LR$  gegeben, so kann man  $Ax = b$  lösen, indem man in zwei Schritten durch Vorwärts-, bzw. Rückwärtseinsetzen löst:

(a) Löse  $Lz = Pb$ ,

(b) löse  $Rx = z$ .

Dies gilt, da

$$\begin{aligned} Ax = b &\iff PAx = Pb, \\ &\iff LRx = Pb, \\ &\iff LRx = Lz, \\ &\iff Rx = z. \end{aligned}$$

I.A. wird  $L$  in den frei werdenden Stellen von  $A$  gespeichert (die 1 auf der Diagonalen muss man nicht speichern). Also

$$A^{(n-1)} := \begin{pmatrix} & \tilde{R} \\ \tilde{L} & \end{pmatrix}.$$

Dabei ist  $\tilde{R}$  die obere Dreiecksmatrix von  $R$  und  $\tilde{L}$  die untere Dreiecksmatrix von  $L$  ohne die Diagonale.  $L$  und  $R$  benötigen also zusammen  $n^2$  Speicherstellen.

*Beweis:* (von 2.9)

Nach dem Gaußalgorithmus gilt:  $R = L_{n-1}^{-1}P_{n-1}\dots L_1^{-1}P_1A$ , wobei  $R$  obere  $\Delta$ -Matrix ist  $\implies P_1L_1\dots P_{n-1}L_{n-1}R = A$ .

Definiere Permutationsmatrix  $P := P_{n-1}\dots P_1$  und  $L := PP_1L_1\dots P_{n-1}L_{n-1}$

$$\implies P^{-1}LR = A \implies LR = PA$$

Noch zu zeigen:  $L$  ist untere  $\Delta$ -Matrix mit Diagonalelementen 1.

Es ist:

$$\begin{aligned} L &= PP_1L_1\dots P_{n-1}L_{n-1} = P_{n-1}\dots \underbrace{P_1P_1}_{=E_n}L_1P_2\dots P_{n-1}L_{n-1} \\ &= P_{n-1}\dots P_2L_1P_2\dots P_{n-1}L_{n-1}. \end{aligned}$$

Wir setzen:

$$\tilde{L}_1 := P_2L_1P_2 \text{ und für } p = 2, \dots, n-1: \tilde{L}_p := P_{p+1}\tilde{L}_{p-1}L_pP_{p+1}.$$

Dann hat  $\tilde{L}_p$  die Gestalt:

$$\tilde{L}_p = \begin{pmatrix} 1 & & & & & & 0 \\ * & \ddots & & & & & \\ \vdots & \ddots & \ddots & & & & \\ \vdots & & * & \ddots & & & \\ \vdots & & \vdots & 0 & \ddots & & \\ \vdots & & \vdots & \vdots & \ddots & \ddots & \\ * & \dots & * & 0 & \dots & 0 & 1 \end{pmatrix}$$

Außerdem gilt, dass  $P_q\tilde{L}_pP_q$  (mit  $q > p$ ) ebenfalls diese Gestalt, da  $P_q = P_{qk}$ ,  $q < k$ .

Also folgt

$$\begin{aligned} L &= P_{n-1}\dots P_3P_2L_1P_2L_2P_3\dots P_{n-1}L_{n-1} \\ &= P_{n-1}\dots P_3\tilde{L}_1L_2P_3\dots P_{n-1}L_{n-1} \\ &= P_{n-1}\dots P_4\tilde{L}_2L_3P_4\dots P_{n-1}L_{n-1} \quad . \\ &\vdots \\ &= \tilde{L}_{n-1} \end{aligned}$$

Somit ist  $L$  untere  $\Delta$ -Matrix.

Zur Eindeutigkeit:

Es gilt:  $L^{-1}$  hat ebenfalls untere  $\Delta$ -Gestalt mit Diagonalelementen 1 (betrachte  $L^{-1}L = E_n$ ). Analog folgt, dass  $S^{-1}$  eine obere  $\Delta$ -Matrix ist.

Also ist  $L^{-1}M$  untere  $\Delta$ -Matrix mit Diagonalelementen 1 und  $RS^{-1}$  hat obere  $\Delta$ -Gestalt. Aus  $LR = MS$  folgt  $RS^{-1} = L^{-1}M = E_n$  und hieraus mit der Eindeutigkeit der Inversen:  $R = S \wedge L = M$ .

□

## Weitere Anwendungen der LR-Zerlegung

- (a) Determinantenberechnung einer Matrix  $A$ .

Hat  $R$  obere/untere  $\triangle$ -Gestalt, so gilt

$$\det(R) = \prod_{i=1}^n r_{ii}.$$

Aus der  $LR$  Zerlegung folgt

$$R = L_{n-1}^{-1} P_{n-1} \dots L_1^{-1} P_1 A$$

und somit

$$\det R = \det(L_{n-1}^{-1}) \det(P_{n-1}) \dots \det(L_1^{-1}) \det(P_1) \det(A).$$

Weiter gilt:

$$\det(L_i^{-1}) = 1 \quad \text{und} \quad \det(P_i) = \det(P_{ik}) = \begin{cases} 1 & : i = k \\ -1 & : i \neq k \end{cases}.$$

Also folgt:

$$\begin{aligned} \det(A) &= \begin{cases} \det(R) & : \text{gerade Anzahl von Zeilenvertauschungen} \\ -\det(R) & : \text{ungerade Anzahl von Zeilenvertauschungen} \end{cases} \\ &= \begin{cases} \prod_{i=1}^n r_{ii} & : \text{gerade Anzahl von Zeilenvertauschungen} \\ -\prod_{i=1}^n r_{ii} & : \text{ungerade Anzahl von Zeilenvertauschungen} \end{cases} \end{aligned}$$

- (b) Bestimmung von  $\text{Rang}(A) = \#$  der linear unabhängigen Zeilenvektoren bei einer nicht unbedingt quadratischen Matrix.

Ist im  $p$ . Schritt  $a_{pp}^{(p)} = 0$ , so müssen Zeilen und eventuell auch Spalten vertauscht werden, was den Rang der Matrix nicht verändert. Ist dies nicht möglich, so hat  $A^{(p)}$  die Gestalt

$$A^{(p)} = \left( \begin{array}{c|c} * & * \\ \hline 0 & 0 \end{array} \right)$$

Dabei sind die ersten  $p$  Zeilenvektoren linear unabhängig, aber alle weiteren Zeilenvektoren sind linear abhängig. Es folgt  $\text{Rang}(A) = \text{Rang}(A^{(p)}) = p$ .

**Achtung:** Aufgrund von Rundungsfehlern kann dieses Verfahren numerisch zu falschen Ergebnissen führen:

(c) Berechnung der Umkehrmatrix  $A^{-1}$  einer Matrix  $A$ .

1. Ansatz: Sei  $e_i$  der  $i$ . Einheitsvektor. Löse  $Ax^{(i)} = e_i$  für  $i = 1, \dots, n \implies A^{-1} = (x_1, \dots, x_n)$  mittels LR-Zerlegung mit  $Lz^{(i)} = Pe_i$  und  $Rx^{(i)} = z^{(i)}$

2. Berechnung durch simultane Elimination

$$\begin{array}{c|cc} & 1 & 0 \\ \hline A & & \\ & \ddots & \\ & 0 & 1 \end{array} \longrightarrow \begin{array}{ccc|ccc} \text{Vorwärtselimination} & & & & & \\ \hline r_{11} & \cdots & * & * & & 0 \\ & \ddots & \vdots & \vdots & \ddots & \\ 0 & & r_{nn} & * & \cdots & * \end{array} \longrightarrow$$

$$\begin{array}{cc|c} \text{Rückwärtselimination} & & \\ \hline r_{11} & 0 & \\ & \ddots & * \\ 0 & r_{nn} & \end{array} \longrightarrow \begin{array}{cc|c} & & \\ \hline 1 & 0 & \\ & \ddots & \\ 0 & 1 & \end{array} A^{-1}$$

### 2.1.2 Gauß-Jordan Verfahren

Diese Methode beruht darauf, durch Matrixumformungen von  $Ax = b$  zu  $Bb = x$  mit  $B = A^{-1}$  überzugehen. Die Idee des Verfahren ist folgende: Ist  $a_{pq} \neq 0$ , so kann die  $p$ -te Gleichung nach  $x_q$  aufgelöst werden:

$$x_q = -\frac{a_{p1}}{a_{pq}}x_1 - \dots - \frac{a_{pq-1}}{a_{pq}}x_{q-1} + \frac{1}{a_{pq}}b_q - \frac{a_{pq+1}}{a_{pq}}x_{q+1} - \dots - \frac{a_{pn}}{a_{pq}}x_n.$$

Durch Einsetzung von  $x_q$  in die anderen Gleichungen ( $j \neq p$ )

$$\sum_{k=1}^{q-1} \left[ a_{jk} - \frac{a_{jq}a_{pk}}{a_{pq}} \right] x_k + \frac{a_{jq}}{a_{pq}}b_q + \sum_{k=q+1}^n \left[ a_{jk} - \frac{a_{jq}a_{jk}}{a_{pq}} \right] x_k = b_j.$$

Man erhält also eine Matrix  $\tilde{A}$  mit

$$\tilde{A} \begin{pmatrix} x_1 \\ \vdots \\ b_q \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ x_q \\ \vdots \\ b_n \end{pmatrix}.$$

Kann dieser Schritt z.B. mit  $p = q$   $n$ -mal durchgeführt werden, so ergibt sich

$$\tilde{A} \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \implies \tilde{A} = A^{-1}.$$

Dies entspricht einem Algorithmus ohne Pivotisierung, d.h.  $a_{ii} \neq 0$ , Varianten mit Pivotisierung sind ebenfalls möglich (siehe z.B. Stoer, Numerische Mathematik I, Springer, 1989, Abschnitt 4.2).

### 2.1.3 Cholesky Verfahren für SPD-Matrizen

Sei  $A \in \mathbb{R}^{n \times n}$  eine symmetrische positive definite Matrix. Dann existiert ein unterer  $\triangle$ -Matrix

$Q \in \mathbb{R}^{n \times n}$ , so dass gilt

$$A = QQ^T.$$

Dabei ist  $Q = LD^{1/2}$  mit  $D^{1/2} = \text{diag}(\sqrt{r_{11}}, \dots, \sqrt{r_{nn}})$ .

Die Existenz einer solchen Zerlegung sieht man wie folgt ein:

$$A = LR \implies A = LD\tilde{R} = LDL^T = LD^{1/2}(LD^{1/2})^T = QQ^T,$$

wobei  $\tilde{R} = L^T$  aus der Symmetrie von  $A$  folgt. Die positive Definitheit der Matrix  $A$  ist notwendig, damit  $D^{1/2}$  wohldefiniert ist.

Unter Ausnutzung der Symmetrie erhält man folgenden Algorithmus, der die untere Dreiecksmatrix von  $Q$  anstelle der unteren Dreiecksmatrix von  $A$  abspeichert und  $D^{-1/2}$  in einem Vektor  $d$ :

#### Algorithmus 2.10 (Cholesky Verfahren)

```

Für  $i = 1, \dots, n$ :
  Für  $j = i, \dots, n$ :
    Setze  $u := a_{ij}$ 
    Für  $k = i - 1, \dots, 1$ :
      Setze  $u := u - a_{jk}a_{ik}$ 
    Falls  $i = j$ ,
      Setze  $d_i := 1/\sqrt{u}$  (Abbruch, falls  $u \leq 0$ )
    Sonst
      Setze  $a_{ji} := d_i u$ 

```

Dieser Algorithmus hat aufgrund der Symmetrie den halben Aufwand im Vergleich zum Gauß-Algorithmus.

### 2.1.4 LR-Zerlegung für Tridiagonalmatrizen

Sei  $A$  eine Tridiagonalmatrix

$$A = \begin{pmatrix} \alpha_1 & \gamma_1 & & 0 \\ \beta_1 & \ddots & \ddots & \\ & \ddots & \ddots & \gamma_{n-1} \\ 0 & & \beta_{n-1} & \alpha_n \end{pmatrix}$$

$A$  kann mittels LR-Zerlegung zerlegt werden. Diese Zerlegung kann explizit in Abhängigkeit von  $\alpha, \beta$  und  $\gamma$  hingeschrieben werden (siehe Übungsaufgabe).

## 2.2 Überbestimmte Gleichungssysteme/Ausgleichsrechnung

**Problem:** Gegeben sind  $m$  Messdaten (zum Beispiel Zeit und Konzentration)  $(x_1, y_1), \dots, (x_m, y_m)$  und Funktionen  $u_1, \dots, u_n$ , ( $n, m \in \mathbb{N}$ ,  $n \leq m$ ).

**Gesucht:** Linearkombination  $u(x) = \sum_{i=1}^n c_i u_i(x)$ , welche die mittlere Abweichung minimiert, also:

$$\Delta_2 := \left( \sum_{j=1}^m (u(x_j) - y_j)^2 \right)^{\frac{1}{2}} = \inf_{c_1, \dots, c_n \in \mathbb{R}} \left( \sum_{j=1}^m \left( \sum_{i=1}^n (c_i u_i(x_j)) - y_j \right)^2 \right)^{\frac{1}{2}}$$

Dieses Problem wird als das **Gaußsche Ausgleichsproblem** oder als die Methode der kleinsten Quadrate (least squares) bezeichnet.

**Bemerkung:** Das **Tschebyscheffsche Ausgleichsproblem**, bei dem bezüglich der Maximumnorm minimiert wird, d.h.

$$\Delta_\infty := \inf_{c_1, \dots, c_n \in \mathbb{R}} \max_{i=1, \dots, m} |c_i u_i(x_j) - y_j|,$$

ist deutlich schwieriger.

Sei:

$$\begin{aligned} c &= (c_1, \dots, c_n)^\top \in \mathbb{R}^n \text{ (der gesuchte Lösungsvektor),} \\ x &= (x_1, \dots, x_m)^\top \in \mathbb{R}^m, \quad y = (y_1, \dots, y_m)^\top \in \mathbb{R}^m, \\ A &= (a_{ij}) \in \mathbb{R}^{m \times n} \text{ mit } a_{ij} = u_i(x_j). \end{aligned}$$

Dann ist das Ausgleichsproblem äquivalent zur Minimierung des Funktionals

$$F(c) := \|Ac - y\|_2. \quad (AGP)$$

**Bemerkung:** Sind  $m = n$ ,  $u_1, \dots, u_n$  linear unabhängig und  $x_1, \dots, x_m$  paarweise verschieden, so ist  $A$  regulär und  $c = A^{-1}y$  ist das gesuchte Minimum. Im Allgemeinen ist jedoch  $n \ll m$ , so dass  $\text{Rang}(A) \leq n$  folgt. In solchen Fällen erwarten wir, dass (AGP) einen Minimierer hat, jedoch  $Ac = y$  entweder keine oder sehr viele Lösungen hat.

### Satz 2.11 (Normalengleichung)

Sei  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$  ( $n \leq m$ ) gegeben. Dann existiert mindestens eine Lösung  $\bar{x} \in \mathbb{R}^n$  des Ausgleichsproblems ( $Ax = b$ ) mit kleinstem Fehlerquadrat, d.h.  $\bar{x}$  minimiert  $F(x) = \|Ax - b\|_2$ . Dies ist äquivalent dazu, dass  $\bar{x}$  die **Normalengleichung**

$$A^\top Ax = A^\top b$$

löst. ist  $\text{Rang}(A) = n$  (d.h. maximal), so ist die Lösung eindeutig bestimmt, andernfalls ist jede weitere Lösung von der Form

$$x = \bar{x} + y$$

mit  $y \in \text{Kern}(A)$ . In diesem Fall wird meistens die Lösung  $x_A$  mit minimaler 2-Norm gesucht, d.h.

$$\|x_A\| = \inf \left\{ \|x\|_2 \mid x \text{ Lösung des Ausgleichsproblems} \right\}.$$

Diese Lösung ist eindeutig.

**Lemma 2.12**

Sei  $A \in \mathbb{R}^{m \times n}$ , dann gelten für  $A^\top A \in \mathbb{R}^{n \times n}$  folgende Eigenschaften:

- (i)  $A^\top A$  ist symmetrisch.
- (ii)  $A^\top A$  ist positiv semidefinit. Falls  $\text{Rang}(A) = n$  ist, so ist  $A^\top A$  positiv definit.
- (iii)  $\text{Kern}(A^\top A) = \text{Kern}(A)$ .
- (iv)  $\mathbb{R}^m = \text{Bld}(A) \oplus \text{Kern}(A^\top)$ .
- (v)  $A^\top A$  und  $AA^\top$  haben dieselben (positiven, reellen) Eigenwerte und es gilt  $\dim \text{Kern}(A^\top A - \lambda I) = \dim \text{Kern}(AA^\top - \lambda I)$  für alle Eigenwerte  $\lambda > 0$ .
- (vi)  $r = \text{Rang}(A) = \text{Rang}(A^\top A) = \text{Rang}(AA^\top) = \text{Rang}(A^\top)$   
 $= \left| \left\{ \lambda > 0 \mid \lambda \text{ EW von } A^\top A \right\} \right|.$

*Beweis:* (Lemma 2.12(iv))

Es gilt  $\mathbb{R}^m = \text{Bld}(A) \oplus \text{Bld}(A)^\perp$ , d.h. es ist zu zeigen:  $\text{Bld}(A)^\perp = \text{Kern}(A^\top)$

$$\begin{aligned} \text{Sei } y \in \text{Bld}(A)^\perp, \text{ d.h. } \forall z \in \text{Bld}(A) : \langle y, z \rangle &= 0 \\ \iff \forall x \in \mathbb{R}^n : \langle y, Ax \rangle &= 0 \iff \forall x \in \mathbb{R}^n : \langle A^\top y, x \rangle = 0 \\ \iff A^\top y &= 0 \iff y \in \text{Kern}(A^\top). \end{aligned}$$

□

*Beweis:* (Satz 2.11)

Wir zeigen zunächst die Äquivalenz des Minimierungsproblems mit dem Lösen der Normalengleichung. Sei  $\bar{x}$  Lösung von  $A^\top A \bar{x} = A^\top b$ . Dann folgt

$$\begin{aligned} \|b - Ax\|_2^2 &= \|b - A\bar{x} + A(\bar{x} - x)\|_2^2 \\ &= \langle b - A\bar{x} + A(\bar{x} - x), b - A\bar{x} + A(\bar{x} - x) \rangle \\ &= \langle b - A\bar{x}, b - A\bar{x} \rangle + 2 \langle b - A\bar{x}, A(\bar{x} - x) \rangle + \langle A(\bar{x} - x), A(\bar{x} - x) \rangle \\ &= \|b - A\bar{x}\|_2^2 + \|A(\bar{x} - x)\|_2^2 + 2 \langle A^\top (b - A\bar{x}), \bar{x} - x \rangle \\ &\geq \|b - A\bar{x}\|_2^2 \end{aligned}$$

Also gilt für alle  $x \in \mathbb{R}^n$  :  $\|b - A\bar{x}\|_2 \leq \|b - Ax\|_2$ .

Sei nun umgekehrt  $\bar{x}$  eine Lösung des Minimierungsproblems, so folgt

$$\begin{aligned} 0 &= \frac{\partial}{\partial x_i} (F(x)) \Big|_{x=\bar{x}} = \frac{\partial}{\partial x_i} \left( \sum_{j=1}^m \left( \sum_{k=1}^n a_{jk} x_k - b_j \right)^2 \right) \Big|_{x=\bar{x}} \\ &= \sum_{j=1}^n a_{ji} 2 \left( \sum_{k=1}^n a_{jk} \bar{x}_k - b_j \right) = 2 \left( \sum_{j=1}^m a_{ji} \sum_{k=1}^n a_{jk} \bar{x}_k - \underbrace{\sum_{j=1}^m a_{ji} b_j}_{a_{ij}^\top} \right) \\ &= 2(A^\top A \bar{x} - A^\top b)_i. \end{aligned}$$

Also folgt  $A^\top A \bar{x} = A^\top b$  und somit ist  $\bar{x}$  Lösung der Normalengleichung.



Insbesondere kann also die Existenz einer Lösung des AGPs durch das Lösen der Normalgleichung gezeigt werden.

Zur Lösung der Normalgleichung: Es ist  $b \in \mathbb{R}^m = \text{Bld}(A) \oplus \text{Kern}(A^\top)$ . Also können wir  $b = s + r$  zerlegen mit  $s \in \text{Bld}(A)$ ,  $r \in \text{Kern}(A^\top)$ . Zu  $s \in \text{Bld}(A)$  existiert ein  $\bar{x} \in \mathbb{R}^n$  mit  $A\bar{x} = s$ . Es folgt:

$$A^\top A\bar{x} = A^\top s + A^\top r = A^\top(r + s) = A^\top b,$$

d.h.  $\bar{x}$  ist Lösung der Normalgleichung.

Ist  $\text{Rang}(A) = n$ , so folgt, dass  $A^\top A$  positiv definit (Lemma 2.12(ii)) und somit auch regulär ist. Insbesondere folgt daraus, dass  $\bar{x} = (A^\top A)^{-1} A^\top b$  die eindeutige Lösung der Normalgleichung ist.

Ist  $\text{Rang}(A) < n$  und seien  $x_1, x_2$  Lösungen der Normalgleichung. Dann gilt  $b = Ax_i + (b - Ax_i) \in \text{Bld}(A) \oplus \text{Kern}(A^\top)$ .

Da die Zerlegung  $\mathbb{R}^m = \text{Bld}(A) \oplus \text{Kern}(A^\top)$  eindeutig ist, folgt  $Ax_i = s = A\bar{x}$ , also  $A(x_i - \bar{x}) = 0$ , d.h.  $x_i - \bar{x} \in \text{Kern}(A)$ .

Wir definieren die Lösungsmenge  $K$  durch

$$K := \left\{ x \in \mathbb{R}^n \mid x \text{ Lösung des AGPs und } \|x\|_2 \leq \|\bar{x}\|_2 \right\}.$$

Dann ist  $K$  kompakt und da die Norm  $\|\cdot\|_2$  stetig ist, nimmt sie ihr Minimum auf  $K$  an. Sei also  $x_A$  die Lösung des AGPs mit

$$\|x_A\|_2 = \inf \left\{ \|x\|_2 \mid x \in K \right\} =: \rho.$$

Sind nun  $x_1, x_2 \in K$  Lösungen mit  $\|x_1\|_2 = \|x_2\|_2 = \rho$ , so folgt  $\frac{x_1+x_2}{2} \in K$  und daher

$$\rho \leq \left\| \frac{x_1 + x_2}{2} \right\|_2 \leq \frac{1}{2} \|x_1\|_2 + \frac{1}{2} \|x_2\|_2 = \rho.$$

Wir erhalten somit  $\left\| \frac{x_1+x_2}{2} \right\|_2 = \rho$  und es folgt

$$\begin{aligned} \rho^2 &= \left\| \frac{x_1+x_2}{2} \right\|_2^2 = \frac{1}{4} \langle x_1 + x_2, x_1 + x_2 \rangle \\ &= \frac{1}{4} \left( \|x_1\|_2^2 + 2 \langle x_1, x_2 \rangle + \|x_2\|_2^2 \right) \\ &= \frac{1}{4} (\rho^2 + 2 \langle x_1, x_2 \rangle + \rho^2) = \frac{1}{2} \rho^2 + \frac{1}{2} \langle x_1, x_2 \rangle, \\ \implies &\langle x_1, x_2 \rangle = \rho^2, \\ \implies &\|x_1 - x_2\|_2^2 = \|x_1\|_2^2 - 2 \langle x_1, x_2 \rangle + \|x_2\|_2^2 = 2\rho^2 - 2\rho^2 = 0, \\ \implies &x_1 = x_2. \end{aligned}$$

□

### Beispiel 2.13 (Ausgleichsgerade)

**Gegeben:** Messdaten:

$x_i$	-2	-1	0	1	2
$y_i$	1/2	1/2	2	7/2	7/2

**Gesucht:** Ausgleichsgerade (linear fit<sup>1</sup>)  $u(x) = bx + a$  mit

$$\left( \sum_{j=1}^5 (bx_j + a - y_j)^2 \right)^{\frac{1}{2}} = \min_{(\hat{a}, \hat{b}) \in \mathbb{R}^2} \left( \sum_{j=1}^5 (\hat{b}x_j + \hat{a} - y_j)^2 \right)^{\frac{1}{2}}.$$

<sup>1</sup>englischer Ausdruck der Ausgleichsgerade

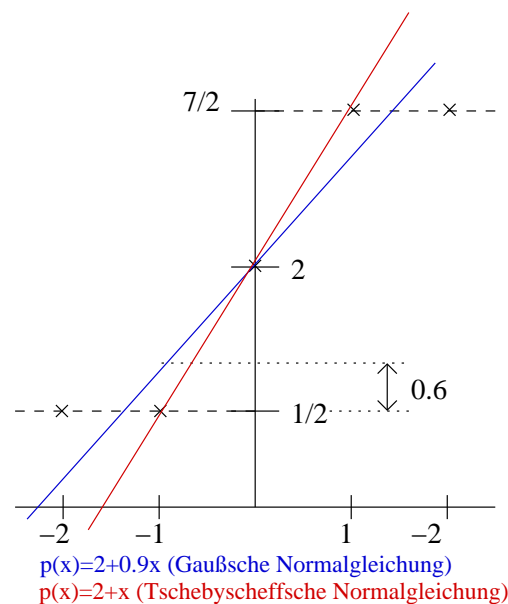


Abbildung 2.1: Ausgleichsgerade

Nach Satz 2.12 ist dann  $(a, b)^\top$  die Lösung der Normalengleichung mit

$$A = \begin{pmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix}; \quad c = \begin{pmatrix} 1/2 \\ 1/2 \\ 2 \\ 7/2 \\ 7/2 \end{pmatrix}$$

und folglich  $A^\top A = \begin{pmatrix} 5 & 0 \\ 0 & 10 \end{pmatrix}$  und  $A^\top c = \begin{pmatrix} 10 \\ 9 \end{pmatrix}$

Da  $\text{Rang}(A) = 2$  ist, ist die Normalengleichung eindeutig lösbar. Aus  $A^\top A \begin{pmatrix} a \\ b \end{pmatrix} = A^\top c$  folgt  $\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 2 \\ 0.9 \end{pmatrix}$ , d.h.  $u(x) = 2 + 0.9x$  ist die Ausgleichsgerade.

Berechnet man die Abweichung, so folgt  $\Delta_2 = \sqrt{0.9} < 1$ ,  $\Delta_\infty = 0.6$ .

Die Lösung des Tschebyscheffs-Problems ist gegeben durch  $u(x) = 2 + x$ . Hier erhält man  $\Delta_2 = 1$ ;  $\Delta_\infty = \frac{1}{2}$ .

**Bemerkung:** Physikalisch könnte z.B. ein nichtlinearer Zusammenhang der Form  $u(x) = \frac{a}{1+bx}$  sinnvoller sein.

Mithilfe der Transformation  $\hat{u}(x) = \frac{1}{u(x)} = \frac{1}{a} + \frac{b}{a}x = \hat{a} + \hat{b}x$  lassen sich jedoch solche Probleme oft auf die Berechnung der linearen Ausgleichsgeraden zurückführen.

**Bemerkung:** Das Ausgleichsproblem kann durch Lösen der Normalengleichung (etwa LR-Zelegung) behandelt werden. Allerdings ist dies numerisch nicht unbedingt der beste Zugang, da  $\text{cond}(A^\top A)$  sehr viel größer als  $\text{cond}(A)$  ist. Beispiel: gilt  $\text{Rang}(A) = n$ ,  $A \in \mathbb{R}^{n \times n}$ , dann ist  $\text{cond}(A^\top A) \sim \text{cond}(A)^2$ .

### 2.2.1 QR-Zerlegung nach Householder

**Idee:** Sei  $\text{Rang}(A) = n$ ,  $A \in \mathbb{R}^{m \times n}$ . Anstelle einer LR-Zerlegung sei

$$A = QR$$

mit einer oberen  $\Delta$ -Matrix  $R \in \mathbb{R}^{m \times n}$  und einer orthogonalen Matrix  $Q \in \mathbb{R}^{m \times m}$  gegeben, d.h.  $Q^{-1} = Q^\top$ . Dann folgt

$$A = QR = Q \begin{pmatrix} * & & * \\ & \ddots & \\ 0 & & * \\ \hline & & 0 \end{pmatrix} = Q \begin{pmatrix} \tilde{R} \\ 0 \end{pmatrix},$$

wobei  $\tilde{R}$  eine reguläre obere  $\Delta$ -Matrix ist. Durch Einsetzen erhalten wir

$$A^\top A = (QR)^\top QR = R^\top Q^\top QR = R^\top R,$$

$$A^\top b = R^\top Q^\top b,$$

d.h. es gilt  $A^\top A \bar{x} = A^\top b \iff R^\top R \bar{x} = R^\top Q^\top b$ .

Definiere  $\tilde{c}$  durch  $c := Q^\top b = \begin{pmatrix} c_1 \\ \vdots \\ c_n \\ c_{n+1} \\ \vdots \\ c_m \end{pmatrix}$ ;  $\tilde{c} = \begin{pmatrix} c_1 \\ \vdots \\ c_n \end{pmatrix} \in \mathbb{R}^n$ ,

so folgt aus  $R^\top = \left( \underbrace{\tilde{R}^\top}_n \mid \underbrace{0}_{m-n} \right)$ :  $R^\top c = \begin{pmatrix} \tilde{R}^\top \tilde{c} \\ 0 \end{pmatrix}$ .

$\tilde{R}^\top$  ist regulär. Sei also  $\bar{x} \in \mathbb{R}^n$  die Lösung von  $\tilde{R} \bar{x} = \tilde{c}$ , so ist  $\bar{x}$  leicht zu berechnen, da  $\tilde{R}$  obere  $\Delta$ -Matrix ist.

Da  $R \bar{x} = \begin{pmatrix} \tilde{c} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$  folgt  $R^\top R \bar{x} = R^\top c = R^\top Q^\top b$  und somit ist  $\bar{x}$  die Lösung der Normalengleichung.

Konzentrieren wir uns also auf die Berechnung einer  $QR$  Zerlegung.

### QR-Zerlegung nach Householder

Sei  $A \in \mathbb{R}^{m \times n}$  ( $m \geq n$ ) mit  $\text{Rang}(A) = n$ .

**Ziel:** Finde obere  $\Delta$ -Matrix  $R \in \mathbb{R}^{m \times n}$  und eine orthogonale Matrix  $Q \in \mathbb{R}^{m \times m}$  mit  $A = QR$ .

**Definition 2.14 (Dyadisches Produkt)**

Seien  $u, v \in \mathbb{R}^m$  (Spaltenvektoren). Dann heißt die Matrix

$$A = uv^\top = \begin{pmatrix} u_1 \\ \vdots \\ u_m \end{pmatrix} (v_1, \dots, v_m)$$

das **dyadische Produkt** von  $u, v$ .

Es gilt  $A \in \mathbb{R}^{m \times m}$  und  $a_{ij} = u_i v_j$  ( $1 \leq i, j \leq m$ ).

**Beachte:**  $\langle u, v \rangle = u^\top v = (u_1, \dots, u_m) \begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix} \in \mathbb{R}.$

**Folgerung 2.15**

Seien  $u, v \in \mathbb{R}^m$ ,  $A = uv^\top$  und  $w \in \mathbb{R}^m$ . Dann gilt

- (i)  $Aw = \langle v, w \rangle u$ ,
- (ii)  $A^2 = \langle u, v \rangle A$ .

*Beweis:*

$$(i) \quad (Aw)_i = \sum_{k=1}^m u_i v_k w_k = \langle v, w \rangle u_i.$$

$$(ii) \quad (A^2)_{ij} = (AA)_{ij} = \sum_{k=1}^m u_i v_k u_k v_j = \left( \sum_{k=1}^m v_k u_k \right) u_i v_j = \langle v, u \rangle (A)_{ij}.$$

□

**Definition 2.16 (Householder Matrix)**

Sei  $v \in \mathbb{R}^m$ ,  $v \neq 0$ . Die Matrix  $H(v) = \mathbb{I} - 2 \frac{vv^\top}{\|v\|_2^2}$  heißt **Householder Matrix**.

Wir setzen  $H(0) = \mathbb{I}$ .

**Folgerung 2.17**

Sei  $v \in \mathbb{R}^m$ , dann gilt:

- (i)  $H(v)$  ist symmetrisch.
- (ii)  $H(v)$  ist orthogonal, d.h.  $H(v) = H(v)^\top = H(v)^{-1}$ .

*Beweis:*

$$(i) \quad \text{Es ist } H(v)_{ij} = \delta_{ij} - 2 \frac{v_i v_j}{\|v\|_2^2} = \delta_{ji} - 2 \frac{v_j v_i}{\|v\|_2^2} = H(v)_{ji}^\top. \text{ Dabei bezeichnet } \delta_{ij} \text{ das Kronecker Symbol.}$$

(ii) Wegen (i) bleibt zu zeigen, dass  $H(v)^2 = \mathbb{I}$  gilt. Es ist

$$\begin{aligned} H(v)^2 &= \left( \mathbb{I} - 2 \frac{vv^\top}{\|v\|_2^2} \right) \left( \mathbb{I} - 2 \frac{vv^\top}{\|v\|_2^2} \right) = \mathbb{I} - 4 \frac{vv^\top}{\|v\|_2^2} + 4 \frac{(vv^\top)^2}{\|v\|_2^4} \\ &= \mathbb{I} - \frac{4}{\|v\|_2^2} \left( vv^\top - \frac{(vv^\top)^2}{\|v\|_2^2} \right) \\ &= \mathbb{I} - \frac{4}{\|v\|_2^2} \left( vv^\top - \frac{\langle v, v \rangle vv^\top}{\|v\|_2^2} \right) \quad (\text{wegen 2.15(ii)}) \\ &= \mathbb{I}. \end{aligned}$$

□

### Satz 2.18

Sei  $a \in \mathbb{R}^m$  und  $u := a \pm \|a\|_2 e_k \in \mathbb{R}^m$  ( $1 \leq k \leq m$ ). Dann gilt

$$H(u)a = \mp \|a\|_2 e_k.$$

*Beweis:* Im Fall  $u = 0$  gilt aufgrund der Definition von  $u$   $a = \mp \|a\|_2 e_k$ . Also folgt  $H(u)a = \mathbb{I}a = \mp \|a\|_2 e_k$ .

Sei also  $u \neq 0$ , etwa  $u = a - \|a\|_2 e_k$ . Dann folgt

$$H(u) = \mathbb{I} - \frac{uu^\top}{h} \quad \text{mit} \quad h = \frac{1}{2} \|u\|_2^2.$$

Und weiter

$$\begin{aligned} h &= \frac{1}{2} \langle a - \|a\|_2 e_k, a - \|a\|_2 e_k \rangle = \frac{1}{2} \left( \|a\|_2^2 - 2 \|a\|_2 \langle a, e_k \rangle + \|a\|_2^2 \right) \\ &= \|a\|_2^2 - \|a\|_2 \langle a, e_k \rangle, \end{aligned}$$

$$\begin{aligned} H(u)a &= \left( \mathbb{I} - \frac{1}{h} uu^\top \right) a = a - \frac{1}{h} (uu^\top) a \\ &= a - \frac{1}{h} \langle u, a \rangle u \quad (\text{Folgerung 2.15(i)}) \\ &= a - \frac{1}{h} \langle a - \|a\|_2 e_k, a \rangle u \\ &= a - \frac{1}{h} \left( \|a\|_2^2 - \|a\|_2 \langle e_k, a \rangle \right) u \\ &= a - u = \|a\|_2 e_k \quad (\text{Definition von } u). \end{aligned}$$

□

### Verfahren: QR Zerlegung

Sei  $A \in \mathbb{R}^{m \times n}$  mit  $A = (a_1, \dots, a_n) = (a_1^{(0)}, \dots, a_n^{(0)})$  gegeben.

#### Schritt 1:

Setze  $u^{(0)} = (a_1^{(0)}) - \left\| a_1^{(0)} \right\|_2 e_1 \in \mathbb{R}^m$  und  $Q_1 = H(u^{(0)})$ , dann folgt

$$Q_1 A = R^{(1)} = \begin{pmatrix} * & \cdots & * \\ 0 & & \\ \vdots & A^{(1)} & \\ 0 & & \end{pmatrix} \quad \text{mit } A^{(1)} \in \mathbb{R}^{m-1 \times n-1} = (a_2^{(1)}, \dots, a_n^{(1)}).$$

#### Schritt 2:

Setze  $u^{(1)} = a_2^{(1)} - \|a_2^{(1)}\| e_1 \in \mathbb{R}^{m-1}$  und  $Q_2 := \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & & & \\ \vdots & & H(u^{(1)}) & \\ 0 & & & \end{pmatrix}$ .

Dann ist  $Q_2 \in \mathbb{R}^{m \times m}$ ,  $H(u^{(1)}) \in \mathbb{R}^{m-1 \times m-1}$  und es folgt

$$Q_2 Q_1 A = Q_2 R^{(1)} = \begin{pmatrix} * & \cdots & \cdots & 0 \\ 0 & * & \cdots & * \\ \vdots & 0 & & \\ \vdots & \vdots & A^{(2)} & \\ 0 & 0 & & \end{pmatrix}.$$

Iterativ erhalten wir so nach  $n$  Schritten:

$$Q_n \cdots Q_1 A = R^{(n)} = \begin{pmatrix} * & \cdots & * \\ & \ddots & \\ 0 & & * \\ \hline & & 0 \end{pmatrix} = R$$

und  $Q := Q_1 \cdots Q_n$  ist orthogonal, da alle  $Q_i$  orthogonal sind. Ausserdem gilt  $A = QR$ , da  $Q_i^2 = \mathbb{I}$  und somit  $QQ_n \cdots Q_1 = \mathbb{I}$ .

Wir fassen die bisherigen Ergebnisse zusammen:

$$\begin{aligned} \text{Ausgleichsproblem} &\iff \text{Normalengleichung} \\ &\iff \tilde{R}\bar{x} = \tilde{c} \text{ mit } A = QR = \begin{pmatrix} \tilde{R} \\ 0 \end{pmatrix}, R \\ &\quad \text{und } R \text{ ist reguläre obere } \triangle\text{-Matrix,} \\ &\quad \text{falls } \text{Rang}(A) = n \text{ gilt.} \end{aligned}$$

Betrachten wir nun die Kondition der Matrix  $R$ . Falls  $A \in \mathbb{R}^{n \times n}$  und  $\text{Rang}(A) = n$  ist, gilt

$$\begin{aligned} \text{cond}_2(A) &= \|A\|_2 \|A^{-1}\|_2 \\ &= \|QR\|_2 \|R^{-1}Q^{-1}\|_2 \\ &= \|QR\|_2 \|R^{-1}Q^\top\|_2 \\ &= \|R\|_2 \cdot \|R^{-1}\|_2 \end{aligned}$$

Also folgt  $\implies \text{cond}_2(A) = \text{cond}_2(R)$ .

### 2.2.2 Singulärwertzerlegung einer Matrix

Die  $QR$ -Zerlegung liefert eine Möglichkeit, (AGP) numerisch zu lösen, falls  $\text{Rang}(A) = n$  ist. Für Probleme mit  $\text{Rang}(A) \leq n$  betrachten wir nun die Singulärwertzerlegung einer Matrix.

**Satz 2.19 (Singulärwertzerlegung)**

Sei  $A \in \mathbb{R}^{m \times n}$ ,  $\text{Rang}(A) = r$ ,  $p = \min\{m, n\}$ . Dann existieren orthogonale Matrizen  $U = (u_1, \dots, u_m) \in \mathbb{R}^{m \times m}$  und  $V = (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$  mit  $U^\top AV = \Sigma \in \mathbb{R}^{m \times n}$  mit

$$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_p),$$

wobei  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_p = 0$ .

D. h.  $\Sigma$  hat die Form

$$\Sigma = \left( \begin{array}{ccc|c} \sigma_1 & & & 0 \\ & \ddots & & \\ & & \sigma_r & 0 \\ \hline 0 & & & 0 \end{array} \right),$$

mit  $\text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$ .

Die Werte  $\sigma_1, \dots, \sigma_r$  heißen **singuläre Werte** von  $A$ . Sie entsprechen gerade den Wurzeln aus den Eigenwerten von  $A^\top A$  bzw.  $AA^\top$  (Bem: nach 2.12(v) haben  $A^\top A$  und  $AA^\top$  dieselben positiven und reellen Eigenwerte).

*Beweis:* Eindeutigkeit: Seien  $U^\top AV = \Sigma$  und  $U, V$  orthogonal, dann gelten  $Av_i = \sigma_i u_i$ , da  $AV = U\Sigma$  und  $A^\top u_i = \sigma_i v_i$ , da  $A^\top U = V\Sigma$ . Daraus ergibt sich

$$A^\top Av_i = \sigma_i A^\top u_i = \sigma_i^2 v_i \implies \sigma_i^2 \text{ ist Eigenwert von } A^\top A.$$

Analog folgt

$$AA^\top u_i = \sigma_i^2 u_i \implies \sigma_i^2 \text{ ist Eigenwert von } AA^\top.$$

Da nach 2.12(v) die Eigenwerte von  $A^\top A$  und  $AA^\top$  übereinstimmen, folgt die Eindeutigkeit.

Existenz: Sei  $\sigma_1 := \|A\|_2 = \max_{\|x\|_2=1} \|Ax\|_2$  (Bemerkung: Da  $\|A\|_2 = \sqrt{\lambda_{\max}(A^\top A)}$ , ist  $\sigma_1$  ein guter Kandidat).

Dann existieren  $x_1 \in \mathbb{R}^n, y_1 \in \mathbb{R}^m$  mit  $\|x_1\|_2 = 1, \|y_1\|_2 = 1$  und  $Ax_1 = \sigma_1 y_1$ . Sei  $V = (x_1, \dots, x_n) \in \mathbb{R}^{n \times n}$  eine orthonormale Basis (ONB) des  $\mathbb{R}^n$  und  $U_1 = (y_1, \dots, y_m) \in \mathbb{R}^{m \times m}$  eine ONB des  $\mathbb{R}^m$ .

Dann folgt:

$$A_1 := U_1^\top AV_1 = \left( \begin{array}{c|c} \sigma_1 & w^\top \\ \hline 0 & B \end{array} \right), \quad B \in \mathbb{R}^{(m-1) \times (n-1)}, \quad w \in \mathbb{R}^{n-1}.$$

Da  $U_1, V_1$  orthogonal, gilt:

$$\|A_1\|_2 = \|A\|_2 = \sigma_1.$$

Desweiteren gilt

$$A_1 \begin{pmatrix} \sigma_1 \\ w \end{pmatrix} = \begin{pmatrix} \sigma_1^2 + w^\top w \\ Bw \end{pmatrix} = \begin{pmatrix} \sigma_1^2 + \|w\|_2^2 \\ Bw \end{pmatrix}$$

und daher

$$\begin{aligned} \sigma_1^2 = \|A_1\|_2^2 &= \left( \max_x \frac{\|A_1 x\|_2}{\|x\|_2} \right)^2 \geq \frac{1}{\|(\sigma_1, w)^\top\|_2^2} \left\| A_1 \begin{pmatrix} \sigma_1 \\ w \end{pmatrix} \right\|_2^2 \\ &= \frac{1}{(\sigma_1^2 + \|w\|_2^2)} \left( (\sigma_1^2 + \|w\|_2^2)^2 + \underbrace{\|Bw\|_2^2}_{\geq 0} \right) \\ &\geq \frac{1}{\sigma_1^2 + \|w\|_2^2} (\sigma_1^2 + \|w\|_2^2)^2 = \sigma_1^2 + \|w\|_2^2 \end{aligned}$$

Insgesamt folgt also

$$\sigma_1^2 \geq \sigma_1^2 + \|w\|_2^2 \implies \|w\|_2^2 = 0 \implies w = 0$$

und wir erhalten

$$U_1^\top AV_1 = \left( \begin{array}{c|c} \sigma & 0 \\ \hline 0 & B \end{array} \right).$$

Die Aussage des Satzes folgt nun durch Induktion. □

### Lösung des Ausgleichsproblems mit Singulärwertzerlegung:

**Gesucht:**  $x \in \mathbb{R}^n$  mit  $\|Ax - b\|_2 = \inf_{z \in \mathbb{R}^n} \|Az - b\|_2$ ,  $A \in \mathbb{R}^{m \times n}$ ,  $b \in \mathbb{R}^m$ ,  $m \geq n \geq r = \text{Rang}(A)$ .

Sei  $U^\top AV = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ , so folgt

$$\begin{aligned} \|Ax - b\|_2^2 &= \langle Ax - b, Ax - b \rangle \stackrel{\text{da } U^\top = \text{orth.}}{=} \langle U^\top(Ax - b), U^\top(Ax - b) \rangle \\ &\stackrel{VV^\top = \mathbb{I}}{=} \langle U^\top AV(V^\top x) - U^\top b, U^\top AV(V^\top x) - U^\top b \rangle \\ &= \langle \Sigma V^\top x - U^\top b, \Sigma V^\top x - U^\top b \rangle \\ &= \|\Sigma V^\top x - U^\top b\|_2^2 \\ &= \sum_{i=1}^r (\sigma_i(V^\top x)_i - u_i^\top b)^2 + \sum_{i=r+1}^m (u_i^\top b)^2 \\ &\geq \sum_{i=r+1}^m (u_i^\top b)^2 \end{aligned}$$

#### Folgerung 2.20

$x \in \mathbb{R}^n$  ist genau dann Lösung des Ausgleichsproblems, wenn

$$V^\top x = \left( \frac{u_1^\top b}{\sigma_1}, \dots, \frac{u_r^\top b}{\sigma_r}, \alpha_{r+1}, \dots, \alpha_n \right), \text{ mit } \alpha_i \in \mathbb{R} \text{ beliebig.}$$

Ist  $x$  Lösung des AGPs, so ist  $\|x\|_2^2 \stackrel{V^\top = \text{orth.}}{=} \|V^\top x\|_2^2 = \sum_{i=1}^r \left( \frac{u_i^\top b}{\sigma_i} \right)^2 + \sum_{i=r+1}^n \alpha_i^2$ .

Also ist  $\|x\|_2$  minimal, g.d.w.  $\alpha_{r+1} = \dots = \alpha_n = 0$  ist. D.h. die eindeutige Lösung des AGPs mit minimaler 2-Norm ist gegeben durch

$$x = V \left( \frac{u_1^\top b}{\sigma_1}, \dots, \frac{u_r^\top b}{\sigma_r}, 0, \dots, 0 \right)^\top = \sum_{i=1}^r \frac{u_i^\top b}{\sigma_i} v_i.$$

### 2.2.3 Pseudoinverse einer Matrix

**Ziel:** Verallgemeinerung der Inversen einer regulären Matrix auf beliebige Matrizen  $A \in \mathbb{R}^{m \times n}$  mit

Hilfe der Singulärwertzerlegung.



**Definition 2.21 (Pseudoinverse)**

Zu  $A \in \mathbb{R}^{m \times n}$  mit  $\text{Rang}(A) = r$  sei

$$\Sigma := U^\top A V = \text{diag}(\sigma_1, \dots, \sigma_{\min\{m,n\}})$$

eine Singulärwertzerlegung von  $A$ . Wir definieren die  $(n \times m)$ -Matrix  $\Sigma^+$  durch

$$\Sigma^+ := \left( \begin{array}{ccc|c} 1/\sigma_1 & & & 0 \\ & \ddots & & \\ & & 1/\sigma_r & 0 \\ \hline & & 0 & 0 \end{array} \right),$$

dann heißt  $A^+ := V \Sigma^+ U^\top \in \mathbb{R}^{n \times m}$  **Pseudoinverse** oder **Penrose Inverse** von  $A$ .

**Bemerkung:** Die singulären Werte einer Matrix sind eindeutig bestimmt, die orthogonalen Matrizen  $U$  und  $V$  jedoch nicht. Die Eindeutigkeit der Pseudoinversen muß also noch gezeigt werden.

**Satz 2.22**

Sei  $A \in \mathbb{R}^{m \times n}$  mit  $\text{Rang}(A) = r$  gegeben. Dann gilt

(i) Ist  $A^+ \in \mathbb{R}^{n \times m}$  eine Pseudoinverse zu  $A$ , so ist

$$AA^+ = (AA^+)^\top, \quad A^+A = (A^+A)^\top, \quad AA^+A = A, \quad A^+AA^+ = A^+.$$

(ii) Durch die “Penrose Bedingung”

$$AB = (AB)^\top, \quad BA = (BA)^\top, \quad ABA = A, \quad BAB = B \quad (*)$$

ist eine Matrix  $B \in \mathbb{R}^{n \times m}$  eindeutig bestimmt.

Insbesondere ist die Pseudoinverse wohldefiniert.

(iii) Sind  $m \geq n$ ,  $b \in \mathbb{R}^m$  und  $x_A$  die eindeutige Lösung des Ausgleichsproblems, so ist  $x_A = A^+b$ .

(iv) Ist  $m = n$  und  $\text{Rang}(A) = n$ , so ist  $A^+ = A^{-1}$ .

(v) Ist  $m \geq n$  und  $\text{Rang}(A) = n$ , so ist  $A^+ = (A^\top A)^{-1} A^\top$ .

(vi) Sind  $\sigma_1 \geq \dots \geq \sigma_r$  die singulären Werte von  $A$ , so ist  $\|A\|_2 = \sigma_1$  und  $\|A^+\|_2 = \frac{1}{\sigma_r}$ .

*Beweis:* (siehe Übungsaufgabe)

**Bemerkung:** Für die Pseudoinverse gilt auch  $(A^+)^+ = A$ , sowie  $(A^\top)^+ = (A^+)^\top$ , jedoch gilt im Allgemeinen nicht  $(AB)^+ = B^+A^+$ .

**Definition 2.23 (Kondition einer singulären Matrix)**

Für eine beliebige Matrix  $A \in \mathbb{R}^{m \times n}$  definieren wir die Kondition von  $A$  bezüglich der Spektralnorm durch

$$\text{cond}_2(A) = \|A\|_2 \|A^+\|_2 = \frac{\sigma_1}{\sigma_r}.$$

## 2.3 Iterative Verfahren

Wir wollen uns nun iterativen Verfahren zur Lösung linearer Gleichungssysteme zuwenden.

**Idee:** Formuliere  $Ax = b$  äquivalent als Fixpunktgleichung, so dass die Kontraktionsbedingung (vgl. Satz 1.29 auf Seite 10) erfüllt ist.

**Ansatz:** Wir zerlegen  $A = M - N$  mit einer regulären Matrix  $M$  (i.A. ist  $M$  vorgegeben und  $N := M - A$ ). Wir erhalten:

$$\begin{aligned} Ax = b &\iff (M - N)x = b \iff Mx - Nx = b \iff Mx = Nx + b \\ &\iff x = M^{-1}Nx + M^{-1}b. \end{aligned}$$

Wir definieren  $T := M^{-1}N = \mathbb{I} - M^{-1}A$ ,  $c := M^{-1}b$  sowie eine Abbildung  $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$  durch  $F(x) = Tx + c$ .

Dann gilt:  $x$  ist die Lösung von  $Ax = b$ , g.d.w.  $x$  ein Fixpunkt von  $F$  ist, d.h. wenn gilt  $x = F(x) = Tx + c$ .

**Iterationsverfahren:** Sei ein Startwert  $x^0 \in \mathbb{R}^n$  gegeben, dann definieren wir die Iteration für  $k \in \mathbb{N}$  durch

$$x^{k+1} = F(x^k),$$

d.h.  $x^{k+1}$  wird berechnet durch folgende Schritte:

- 1)  $My^k = r^k$ , mit dem Residuum  $r^k := b - Ax^k$ ,
- 2)  $x^{k+1} = x^k + y^k$ .

**Bemerkung:** Aus den Lösungsschritten erkennt man, dass solche Iterationsverfahren nur dann von Interesse sind, wenn die Defektgleichung  $My^k = r^k$  leicht zu lösen ist. Dies ist z.B. der Fall, falls  $M$  eine Diagonalmatrix oder eine obere  $\Delta$ -Matrix ist.

### Satz 2.24

Die Folge  $(x^k)_{k \in \mathbb{N}}$  sei durch eine Fixpunktiteration zu  $x = F(x) = Tx + c$  gegeben. Sei  $\|\cdot\|$  eine Norm auf dem  $\mathbb{R}^n$ , so dass für die induzierte Matrixnorm gilt

$$q := \|T\| < 1.$$

Dann konvergiert  $x^k$  gegen ein  $x$  mit  $Ax = b$  und es gelten die Fehlerabschätzungen

$$\|x - x^k\| \leq \frac{q^k}{1 - q} \|x^1 - x^0\|,$$

bzw.

$$\|x - x^k\| \leq \frac{q}{1 - q} \|x^k - x^{k-1}\|.$$

*Beweis:* Der Beweis folgt mit dem Banachschen Fixpunktsatz 1.29 da gilt

$$\begin{aligned} \|F(y_1) - F(y_2)\| &= \|Ty_1 + c - (Ty_2 + c)\| \\ &= \|T(y_1 - y_2)\| \leq \|T\| \|y_1 - y_2\| = q \|y_1 - y_2\|. \end{aligned}$$

Aus  $q < 1$  folgt also, dass  $F$  eine Kontraktion ist.

**Problem:** Die Kontraktionsbedingung ist abhängig von der Norm, die Konvergenz nicht, da alle Normen auf  $\mathbb{R}^n$  äquivalent sind.

Gesucht: Notwendige und hinreichende Bedingung für die Konvergenz.

**Definition 2.25 (Spektralradius)**

Für eine Matrix  $B \in \mathbb{R}^{n \times n}$  definieren wir den **Spektralradius**  $\rho(B)$  durch

$$\rho(B) := \max \{ |\lambda| \mid \lambda \in \mathbb{C} \text{ ist Eigenwert von } B \}.$$

**Bemerkung:** Eine Matrix  $B$  hat in  $\mathbb{C}$  die Eigenwerte  $\lambda_1, \dots, \lambda_n$  (falls Vielfachheit zugelassen wird). Es existiert eine reguläre Matrix  $U \in \mathbb{C}^{n \times n}$  mit

$$U^{-1}BU = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix},$$

also eine Ähnlichkeitstransformation (Jordansche Normalform).

**Lemma 2.26**

Für  $B \in \mathbb{R}^{n \times n}$  gilt

- (i)  $\rho(B) \leq \|B\|$  für jede induzierte Matrixnorm.
- (ii)  $\forall \varepsilon > 0$  gibt es eine induzierte Norm  $\|\cdot\|$  auf  $\mathbb{C}^{n \times n}$  mit  $\|B\| \leq \rho(B) + \varepsilon$ .

*Beweis:*

- (i) Seien  $\lambda$  ein Eigenwert von  $B$ ,  $u \in \mathbb{C}^n \setminus \{0\}$  der zugehörige Eigenvektor und  $\|\cdot\|$  eine Norm auf  $\mathbb{C}^n$ . Dann folgt

$$\begin{aligned} Bu = \lambda u &\implies |\lambda| \|u\| = \|\lambda u\| = \|Bu\| \leq \|B\| \|u\| \\ &\implies |\lambda| \leq \|B\| \text{ für alle EWe } \implies \rho(B) \leq \|B\|. \end{aligned}$$

- (ii) Sei  $U \in \mathbb{C}^{n \times n}$  regulär mit  $U^{-1}BU = \begin{pmatrix} \lambda_1 & & r_{ij} \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$ .

Für  $\delta > 0$  sei  $D_\delta = \text{diag}(\delta^0, \dots, \delta^{n-1})$ . Für  $G \in \mathbb{C}^{n \times n}$  gilt dann  $(D_\delta^{-1}GD_\delta) = g_{ij}\delta^{j-i}$ . Also folgt

$$(UD_\delta)^{-1}B(UD_\delta) = D_\delta^{-1}(U^{-1}BU)D_\delta = \begin{pmatrix} \lambda_1 & \delta r_{12} & \delta^2 r_{13} & \cdots & \delta^{n-1} r_{1n} \\ 0 & \lambda_2 & \delta r_{23} & \cdots & \delta^{n-2} r_{2n} \\ \vdots & & \ddots & \ddots & \vdots \\ \vdots & & & \ddots & \delta r_{n-1,n} \\ 0 & \cdots & \cdots & \cdots & \lambda_n \end{pmatrix}.$$

Sei  $\varepsilon > 0$  vorgegeben. Dann existiert ein  $\delta > 0$  mit  $\sum_{k=i+1}^n \delta^{k-i} |r_{ik}| \leq \varepsilon$  für alle  $i \in \{1, \dots, n-1\}$ .

Setze  $\|x\|_\delta = \|(UD_\delta)^{-1} x\|_\infty$ . Diese ist eine Norm auf  $\mathbb{C}^n$ , da  $(UD_\delta)^{-1}$  regulär. Es gilt:

$$\|B\|_\delta = \sup_{x \in \mathbb{C}^n \setminus \{0\}} \frac{\|Bx\|_\delta}{\|x\|_\delta} = \sup_{x \in \mathbb{C}^n \setminus \{0\}} \frac{\|(UD_\delta)^{-1} Bx\|_\infty}{\|(UD_\delta)^{-1} x\|_\infty}.$$

Mit  $w := (UD_\delta)^{-1} x$  bzw.  $x = (UD_\delta) w$  gilt:

Wenn  $x$  über  $\mathbb{C}^n$  läuft, dann läuft auch  $w$  über den ganzen  $\mathbb{C}^n$ , da  $(UD_\delta)^{-1}$  regulär, d.h.  $\sup_{x \neq 0} \dots = \sup_{w \neq 0} \dots$ . Also folgt

$$\begin{aligned} \|B\|_\delta &= \sup_{w \neq 0} \frac{\|(UD_\delta)^{-1} B(UD_\delta)w\|_\infty}{\|w\|_\infty} \\ &\leq \left\| (UD_\delta)^{-1} B (UD_\delta) \right\|_\infty \\ &= \max_i \left\{ |\lambda_i| + \underbrace{\sum_{k=i+1}^n \delta^{k-1} \cdot |r_{ik}|}_{\leq \varepsilon} \right\} \quad (\text{Zeilensummennorm}) \\ &= \max_i |\lambda_i| + \varepsilon = \rho(B) + \varepsilon \end{aligned}$$

Aus Lemma 2.26 folgt folgender Satz über die Konvergenz geometrischer Matrixfolgen. □

### Satz 2.27

Sei  $T \in \mathbb{C}^{n \times n}$  eine reguläre Matrix. Dann sind äquivalent

- (i)  $T^\nu \rightarrow 0$  für  $\nu \rightarrow \infty$ .
- (ii) Für  $u \in \mathbb{C}^n$  gilt:  $T^\nu u \rightarrow 0$  für  $\nu \rightarrow \infty$ .
- (iii)  $\rho(T) < 1$ .
- (iv) Es existiert eine Norm auf  $\mathbb{C}^n$ , so dass für die induzierte Matrixnorm  $\|T\| < 1$  gilt.

*Beweis:*

(i)  $\implies$  (ii): Es gilt  $\|T^\nu u\| \leq \|T^\nu\| \|u\| \rightarrow 0$ , da  $T^\nu \rightarrow 0$ . Also folgt  $T^\nu u \rightarrow 0$  (für  $\nu \rightarrow \infty$ ).

(ii)  $\implies$  (iii): Annahme:  $\rho(T) \geq 1$ , d.h. es existiert ein EW  $\lambda \in \mathbb{C}$  mit  $|\lambda| \geq 1$ . Sei  $u \in \mathbb{C}^n \setminus \{0\}$  der zugehörige EV, dann gilt  $T^\nu u = T^{\nu-1}(Tu) = T^{\nu-1}(\lambda u) = \lambda T^{\nu-1}u = \dots = \lambda^\nu u$ . Weiter folgt  $\|T^\nu u\| = \|\lambda^\nu u\| = |\lambda|^\nu \|u\| \not\rightarrow 0$ , da  $|\lambda| \geq 1$ . Dies ist ein Widerspruch zu (ii).

(iii)  $\implies$  (iv): Lemma 2.26.

(iv)  $\implies$  (i): Es gilt  $\|T^\nu\| \stackrel{\text{Submultipl.}}{\leq} \|T\|^\nu \rightarrow 0$ , da  $\|T\| < 1$ . □

### Folgerung 2.28

Das Iterationsverfahren konvergiert genau dann wenn  $\rho(T) < 1$ .

**Bemerkung:** (Ohne Beweis)

Aus den bisher gezeigten Sätzen folgt:

Um eine Dezimalstelle im Fehler zu gewinnen, müssen  $K \sim -\frac{\ln 10}{\ln(\rho(T))}$  Schritte durchgeführt werden.

Das heißt für  $\rho(T) \sim 1$  ist  $-\ln(\rho(T)) \sim 0$  und  $K$  sehr groß. Somit muß das Ziel bei der Wahl der regulären Matrix  $M$  sein:

(a)  $\rho(T) = \rho(\mathbb{I} - M^{-1}A)$  möglichst klein.

(b) Das Gleichungssystem  $My^k = r^k$  muss leicht zu lösen sein.

Dies sind widersprüchliche Forderungen: Optimal für Bedingung (a) wäre  $M = A \implies \rho(T) = 0$ , aber dann ist die Bedingung (b) nicht erfüllt.

Auf der anderen Seite ist (b) erfüllt, falls  $M$  eine Diagonalmatrix ist. Dies führt uns zu folgendem Verfahren.

### 2.3.1 Gesamtschritt Verfahren (GSV)/ Jacobi Verfahren

Sei  $A$  regulär und  $a_{ii} \neq 0$  für alle  $i = 1, \dots, n$ . Setze

$$M := D = \text{diag}(a_{11}, \dots, a_{nn}) \implies T = \mathbb{I} - D^{-1}A.$$

Dies führt zu folgender Iterationsvorschrift:

Sei ein Startvektor  $x^0 \in \mathbb{R}^n$  gegeben.

Iteration:

$$\begin{aligned} x^{k+1} &= (\mathbb{I} - D^{-1}A)x^k + D^{-1}b \\ &= D^{-1}(Dx^k - Ax^k + b), \\ \implies x_i^{k+1} &= \frac{1}{a_{ii}} \left( b_i - \sum_{l \neq i} a_{il}x_l^k \right), \quad i = 1, \dots, n. \end{aligned}$$

**Satz 2.29 (Hinreichende Bedingung für die Konvergenz des Jacobi-Verfahrens)**

Falls entweder

$$(a) \max_i \left( \sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} \right) < 1 \quad (\text{starkes Zeilensummenkriterium})$$

oder

$$(b) \max_i \left( \sum_{k \neq i} \frac{|a_{ki}|}{|a_{ii}|} \right) < 1 \quad (\text{starkes Spaltensummenkriterium})$$

gilt, so konvergiert das Jacobi-Verfahren.

Falls (a) oder (b) erfüllt ist, heißt die Matrix  $A$  **stark diagonal dominant**, da in diesem Fall die Beträge der Diagonaleinträge größer sind als die Summe der Zeilen bzw. Spalten der Matrix.

*Beweis:* Gelte (a), so folgt

$$\begin{aligned}
\rho(T) \leq \|T\|_\infty &= \|\mathbb{I} - D^{-1}A\|_\infty \\
&= \max_i \left( \sum_{k=1}^n \left| \delta_{ik} - \frac{|a_{ik}|}{|a_{ii}|} \right| \right) \\
&= \max_i \left( \sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} \right) < 1 \text{ (wegen (a))}.
\end{aligned}$$

Gelte (b), so folgt

$$\rho(T) \leq \|T\|_1 = \max_i \left( \sum_{k \neq i} \frac{|a_{ki}|}{|a_{ii}|} \right) < 1 \text{ (wegen (b))}. \quad \square$$

Die starke Diagonaldominanz ist nur eine hinreichende Bedingung. Betrachten wir folgendes Beispiel.

### Beispiel 2.30

Bei der Diskretisierung  $-\partial_{xx}u = f$  in §1.5 musste ein LGS  $Ax = b$  gelöst werden mit

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

Für  $A$  gilt  $\sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} = 1$  für  $i = 2, \dots, n-1$  und für  $i = 1, n$  gilt  $\sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} < 1$ . Für dieses Beispiel ist also die starke Diagonaldominanz nicht erfüllt. Wir wollen nun eine schwächere Bedingung herleiten, die auch Matrizen eines solchen Typs mit beinhaltet.

### Definition 2.31 (Zerlegbare Matrizen)

Eine Matrix  $A = (a_{ik})$  heißt **zerlegbar**, falls es eine Zerlegung von  $N := \{1, \dots, n\}$  in zwei Teilmengen  $N_1, N_2 \subset N$  gibt mit  $a_{ik} = 0 \ \forall (i, k) \in N_1 \times N_2$ .

(Zerlegung heißt:  $N_1 \neq \emptyset$ ,  $N_2 \neq \emptyset$ ,  $N = N_1 \cup N_2$ ,  $N_1 \cap N_2 = \emptyset$ ).

### Lemma 2.32

Für  $A \in \mathbb{R}^{n \times n}$  sind äquivalent

(i)  $A$  ist zerlegbar.

(ii) Der zugehörige gerichtete Graph

$$G(A) := (\text{Knoten } P_1, \dots, P_n, \text{ gerichtete Kanten } \overline{P_j P_k} \iff a_{jk} \neq 0)$$

ist nicht zusammenhängend, d.h. es existieren Knoten  $P_j$  und  $P_k$ , so dass kein Pfad zwischen ihnen existiert. Es gibt also keine Folge  $l_0, \dots, l_L \in \{1, \dots, N\}$  mit  $l_0 = j$ ,  $l_L = k$  und  $a_{l_i l_{i+1}} \neq 0$  für alle  $i = 0, \dots, L$ .

(iii) Es existiert eine Permutationsmatrix  $P \in \mathbb{R}^{n \times n}$  mit

$$PAP^\top = \begin{pmatrix} A_{11} & 0 \\ A_{21} & A_{22} \end{pmatrix}$$

mit  $A_{11} \in \mathbb{R}^{p \times p}$ ,  $A_{22} \in \mathbb{R}^{q \times q}$ ,  $A_{21} \in \mathbb{R}^{q \times p}$  und  $p + q = n$ .

*Beweis:* Wir zeigen  $(i) \iff (ii)$ . Die Äquivalenz mit (iii) wird in den Übungen behandelt.

$(i) \implies (ii)$ : Sei  $A$  zerlegbar, d.h. es existieren  $N_1 \neq \emptyset$ ,  $N_2 \neq \emptyset$ ,  $N = N_1 \cup N_2$ ,  $N_1 \cap N_2 = \emptyset$  und es gilt  $a_{ik} = 0 \ \forall (i, k) \in N_1 \times N_2$ . Sei  $(i, k) \in N_1 \times N_2$ . Annahme:  $G(A)$  ist zusammenhängend.

Dann existiert eine Folge  $l_0, \dots, l_L \in \{1, \dots, N\}$  mit  $l_0 = i$  und  $l_L = k$  und  $a_{l_i l_{i+1}} \neq 0$  für alle  $i = 0, \dots, L$ . Nach Voraussetzung ist  $l_0 = i \in N_1$ . Da  $a_{l_0 l_1} \neq 0$ , ist  $l_1 \notin N_2 \implies l_1 \in N_1$  und induktiv folgt somit  $l_L = k \in N_1$ . Dies ist ein Widerspruch zur Annahme  $k \in N_2$ .

$(ii) \implies (i)$ : Sei nun  $G(A)$  nicht zusammenhängend, existiere etwa kein Pfad zwischen  $P_j$  und  $P_k$ . Setze  $N_1 := \{l \mid \text{es existiert ein Pfad zwischen } P_j \text{ und } P_l\} \cup \{j\}$ ,  $N_2 := N \setminus N_1$ . Dann gilt  $k \in N_2$ . Also folgt  $N_1 \neq \emptyset$ ,  $N_2 \neq \emptyset$ ,  $N = N_1 \cup N_2$ ,  $N_1 \cap N_2 = \emptyset$ . Sei  $(i, l) \in N_1 \times N_2$ . Wäre  $a_{il} \neq 0$ , so würde ein Pfad von  $P_j$  zu  $P_l$  existieren und somit wäre  $l \in N_1$ . Dies ist ein Widerspruch und folglich gilt  $a_{il} = 0$  für alle  $(i, l) \in N_1 \times N_2$ .

□

### Beispiel 2.33

Wir betrachten die Matrix

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

Dann ist der Graph  $G(A)$  wie in Abb. 2.2 gegeben. Es gelten:

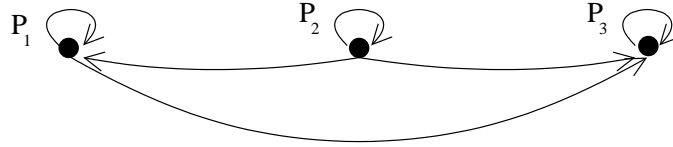


Abbildung 2.2: Graph, Beispiel 2.33

(i)  $N_1 = \{1, 3\}$ ,  $N_2 = \{2\}$  ist eine geeignete Zerlegung.

(ii)  $G(A)$  nicht zusammenhängend, da es keine Verbindung zwischen  $P_1$  und  $P_2$  existiert. (Siehe Abbildung 2.2).

(iii)

$$P = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}; \quad PAP^\top = \left( \begin{array}{c|cc} 2 & 0 & 0 \\ \hline 1 & 2 & 2 \\ 2 & 0 & 2 \end{array} \right).$$

### Beispiel 2.34

Sei  $A$  eine Tridiagonalmatrix ohne Nullen auf der Diagonalen und den Nebendiagonalen, d.h.

$$A = \begin{pmatrix} * & * & & 0 \\ * & \ddots & \ddots & \\ & \ddots & \ddots & * \\ 0 & & * & * \end{pmatrix}$$

Dann ist  $A$  unzerlegbar, da der Graph zusammenhängend ist. (Man kommt von jedem inneren Knoten zu den Nachbarn rechts und links. Siehe Übungsaufgabe)

**Satz 2.35 (Schwachtes Zeilensummenkriterium)**

Sei  $A \in \mathbb{R}^{n \times n}$  unzerlegbar und erfülle das schwache Zeilensummenkriterium, d.h.

$$\max_i \left( \sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} \right) \leq 1$$

und es existiert ein  $r \in \{1, \dots, n\}$  mit

$$\sum_{k \neq r} \frac{|a_{rk}|}{|a_{rr}|} < 1.$$

Dann kann das Jacobi-Verfahren angewendet werden und es konvergiert für alle Startvektoren  $x^0 \in \mathbb{R}^n$ .

*Beweis:* Wir zeigen zunächst, dass gilt  $|a_{ii}| > 0$ . Dazu nehmen wir an, dass gilt  $\sum_{k \neq i} |a_{ik}| \leq |a_{ii}|$ . Da  $A$  unzerlegbar ist, folgt  $\sum_{k \neq i} |a_{ik}| > 0$ , und somit  $|a_{ii}| > 0$ . Insbesondere kann also das Jacobi-Verfahren angewendet werden.

Analog zum Beweis von Satz 2.29 können wir zeigen:  $\rho(T) \leq 1$  mit  $T = M^{-1}N$ . Wir müssen also noch zeigen, dass  $\rho(T) \neq 1$  ist.

Annahme: Es existiert ein Eigenwert  $\lambda \in \mathbb{C}$  mit  $|\lambda| = 1$ . Sei  $v \in \mathbb{C}^n$  der zugehörige Eigenvektor mit  $\|v\|_\infty = 1$ , d.h.  $\exists s \in \{1, \dots, n\}$  mit  $|v_s| = 1$ . Aus  $Tv = \lambda v$  folgt für  $i = 1, \dots, n$ :

$$\lambda v_i = \sum_{k=1}^n t_{ik} v_k = \sum_{k=1}^n \left( \delta_{ik} - \frac{a_{ik}}{a_{ii}} \right) v_k = \sum_{k \neq i} \frac{a_{ik}}{a_{ii}} v_k.$$

Also folgt

$$|v_i| = |\lambda| |v_i| \leq \frac{1}{|a_{ii}|} \sum_{k \neq i} |a_{ik}| |v_k|. \quad (*)$$

Da  $G(A)$  zusammenhängend ist, existiert ein Pfad zwischen  $P_s$  und  $P_r$ , d.h.  $l_0 = s, \dots, l_L = r$  und  $a_{l_i, l_{i+1}} \neq 0$  ( $i = 0, \dots, L-1$ ). Mit  $(*)$  folgt also:

$$\begin{aligned} |v_r| &\leq \frac{1}{|a_{rr}|} \sum_{k \neq r} |a_{rk}| |v_k| \leq \frac{\|v\|_\infty}{|a_{rr}|} \sum_{k \neq r} |a_{rk}| < \|v\|_\infty, \\ |v_{l_{L-1}}| &\stackrel{*}{\leq} \frac{1}{|a_{l_{L-1}, l_{L-1}}|} \left( \sum_{k \neq l_{L-1}; k \neq l_L} |a_{l_{L-1}, k}| |v_k| + |a_{l_{L-1}, l_L}| |v_{l_L}| \right) \\ &< \frac{\|v\|_\infty}{|a_{l_{L-1}, l_{L-1}}|} \left( \sum_{k \neq l_{L-1}} |a_{l_{L-1}, k}| \right) \leq \|v\|_\infty, \\ &\vdots \\ \|v\|_\infty = |v_s| &= |v_{l_0}| < \|v\|_\infty. \end{aligned}$$

Dies ist ein Widerspruch und somit muss  $\rho(T) < 1$  sein. Mit Folgerung 2.28 folgt dann die Behauptung.



Als nächstes Iterationsverfahren wählen wir  $M$  als die untere Dreiecksmatrix von  $A$ . Wir erhalten  $\square$  das Einzelschritt Verfahren.

### 2.3.2 Einzelschritt Verfahren (ESV) / Gauß-Seidel-Verfahren

Sei  $A$  regulär mit  $a_{ii} \neq 0$ . Wir zerlegen  $A$  additiv in  $A = L + D + R$  und setzen

$$M = L + D = \begin{pmatrix} a_{11} & & 0 \\ \vdots & \ddots & \\ a_{1n} & \cdots & a_{nn} \end{pmatrix}.$$

Da  $a_{ii} \neq 0$  ist, ist  $M$  regulär und  $N = M - A = -R$ .

Dies führt zu folgender Iterationsvorschrift:

Sei ein Startvektor  $x^0 \in \mathbb{R}^n$  gegeben.

Iteration (ESV):

$$\begin{aligned} x^{k+1} &= (L + D)^{-1} (b - Rx^{(k)}) \\ \Rightarrow x_i^{k+1} &= \frac{1}{a_{ii}} \left( b_i - \sum_{l=1}^{i-1} a_{il} x_l^{(k+1)} - \sum_{l=i+1}^n a_{il} x_l^{(k)} \right), \quad i = 1, \dots, n. \end{aligned}$$

Vergleiche mit (GSV):

$$x_i^{k+1} = \frac{1}{a_{ii}} \left( b_i - \sum_{l=1}^{i-1} a_{il} x_l^{(k)} - \sum_{l=i+1}^n a_{il} x_l^{(k)} \right), \quad i = 1, \dots, n.$$

#### Satz 2.36

Die Matrix  $A$  erfülle das starke Zeilensummenkriterium. Dann konvergiert das Einzelschrittverfahren.

*Beweis:* Setze  $T := M^{-1}N$  und  $q = \max_i \left( \sum_{k \neq i} \frac{|a_{ik}|}{|a_{ii}|} \right) < 1$ .

Wir zeigen:  $\|T\|_\infty \leq q < 1$  und folglich ist  $T$  Kontraktion.

Es ist  $\|T\|_\infty = \sup_{\|x\|_\infty=1} \|Tx\|_\infty$ . Sei also  $x \in \mathbb{R}^n$  mit  $\|x\|_\infty = 1$ . Setze  $y := Tx$ . Zu zeigen  $\|y\|_\infty \leq q$ , d.h.  $|y_i| \leq q$  ( $1 \leq i \leq n$ ).

Induktion:

$$\text{I.A. } y_1 = -\frac{1}{a_{11}} \left( \sum_{k=2}^n a_{1k} x_k \right)$$

$$\Rightarrow |y_1| \leq \frac{1}{|a_{11}|} \sum_{k=2}^n |a_{1k}| \underbrace{|x_k|}_{\leq 1} \leq \frac{1}{|a_{11}|} \sum_{k \neq 1} |a_{1k}| \leq q.$$

$$\text{I.S. } y_i = -\frac{1}{a_{ii}} \left( \sum_{k=i+1}^n a_{ik} x_k - \sum_{k=1}^{i-1} a_{ik} y_k \right)$$

$$\Rightarrow |y_i| \leq \frac{1}{|a_{ii}|} \left( \sum_{k=i+1}^n |a_{ik}| \underbrace{|x_k|}_{\leq 1} + \sum_{k=1}^{i-1} |a_{ik}| \underbrace{|y_k|}_{\leq q} \right) \leq q.$$

Wie für das Gesamtschrittverfahren erhalten wir auch hier folgenden Konvergenzsatz.  $\square$

**Satz 2.37 (Ohne Beweis)**

Sei  $A \in \mathbb{R}^{n \times n}$  unzerlegbar und erfülle das schwache Zeilensummenkriterium, dann konvergiert das Gauß-Seidel-Verfahren.

Darüberhinaus können wir für das Gauß-Seidel-Verfahren auch folgenden Satz zeigen, der für das Jacobi-Verfahren nicht gilt.

**Satz 2.38**

$A$  sei symmetrisch und positiv definit, dann konvergiert das Gauß-Seidel-Verfahren.

*Beweis:* Sei  $A = L + D + R$ . Da  $A$  symmetrisch ist gilt:  $R = L^\top$  bzw.  $R^\top = L$ ,  $M = L + D$ ,  $N = -R$ ,  $T = M^{-1}N = -(L + D)^{-1}R$ .

Sei  $\lambda \neq 0$  ein Eigenwert von  $T$  in  $\mathbb{C}$ . Sei  $x \in \mathbb{C}^n$  der zugehörige Eigenvektor mit  $\|x\|_2 = 1$ . Dann gilt:  $-(L + D)^{-1}Rx = \lambda x$ , bzw.  $-Rx = \lambda Dx + \lambda Lx = \lambda Dx + \lambda R^\top x$ . Es folgt

$$D = A - L - R = A - R^\top - R \implies -Rx = \lambda(A - R^\top - R)x + \lambda R^\top x = \lambda Ax - \lambda Rx.$$

Setze  $\alpha := \langle Ax, x \rangle > 0$ , da  $A$  positiv definit,  $\sigma := \langle Rx, x \rangle = \sigma_1 + i\sigma_2 \in \mathbb{C}$ ,  $\delta := \langle Dx, x \rangle$ . Dann folgt

$$-\sigma = \lambda\alpha - \lambda\sigma = \lambda(\alpha - \sigma)$$

Hieraus folgt  $\alpha - \sigma \neq 0$ , da sonst  $-\sigma = 0$  und somit  $\alpha = \sigma = 0$  wäre. Dies ist ein Widerspruch zur positiven Definitheit. Also folgt

$\lambda = -\frac{\sigma}{\alpha - \sigma} \implies |\lambda|^2 = \frac{\sigma\bar{\sigma}}{(\alpha - \sigma)(\alpha - \bar{\sigma})} = \frac{\sigma_1^2 + \sigma_2^2}{(\alpha - \sigma_1)^2 + \sigma_2^2}$ , da  $\alpha \in \mathbb{R}$  und  $\alpha - \sigma = \alpha - \sigma_1 - i\sigma_2$ . Weiter haben wir

$$\begin{aligned} \alpha = \langle Ax, x \rangle &= \langle R^\top x, x \rangle + \langle Dx, x \rangle + \langle Rx, x \rangle \\ &= \delta + 2\sigma_1 \implies \delta = \alpha - 2\sigma_1, \end{aligned}$$

$$\begin{aligned} (\alpha - \sigma_1)^2 &= (\delta + \sigma_1)^2 = \delta^2 + 2\delta\sigma_1 + \sigma_1^2 = \delta(\alpha - 2\sigma_1) + 2\delta\sigma_1 + \sigma_1^2 \\ &= \delta\alpha + \sigma_1^2 \geq \mu\delta + \sigma_1^2, \end{aligned}$$

wobei  $\mu > 0$  der kleinste Eigenwert von  $A$  ist.

$$\delta = \langle Dx, x \rangle = \sum_{i=1}^n a_{ii}x_i\bar{x}_i \geq \min_i a_{ii} \|x\|_2^2 = \min_i a_{ii} =: \underline{\delta}$$

$$\implies |\lambda|^2 = \frac{\sigma_1^2 + \sigma_2^2}{(\alpha - \sigma_1)^2 + \sigma_2^2} \leq \frac{\sigma_1^2 + \sigma_2^2}{\mu\underline{\delta} + \sigma_1^2 + \sigma_2^2} < 1.$$

da  $\mu\underline{\delta} > 0$ . Also folgt insgesamt  $|\lambda| < 1$  und somit  $\rho(T) < 1$ .

□

## 2.4 Gradientenverfahren

**Generalvereinbarung 2.39**

In diesem Abschnitt gelte stets

- 1)  $A \in \mathbb{R}^{m \times m}$  mit  $A = A^T$ ,
- 2)  $A$  sei positiv definit und
- 3)  $b \in \mathbb{R}^m$ .

**Ziel:** Löse das lineare Gleichungssystem  $Ax = b$ .

Dazu wollen wir das Problem zunächst in ein Minimierungsproblem überführen.

**Definition 2.40 (Minimierungsaufgabe)**

Sei  $F(x) := \langle A^{-1}(Ax - b), Ax - b \rangle$ ,  $x \in \mathbb{R}^m$ .  
 $x$  heißt Lösung der Minimierungsaufgabe (M), g.d.w.

$$(M) \quad F(x) = \min_{y \in \mathbb{R}^m} F(y).$$

**Lemma 2.41**

$A, b$  seien mit den Eigenschaften der Generalvereinbarung 2.39 gegeben. Dann sind äquivalent:

- 1)  $x \in \mathbb{R}^m$  löst  $Ax = b$ ,
- 2)  $x \in \mathbb{R}^m$  löst (M).

*Beweis:*  $A$  positiv definit  $\implies A^{-1}$  positiv definit. Also gilt  $F(y) \geq 0 \quad \forall y \in \mathbb{R}^m$ .

1)  $\Rightarrow$  2): Gilt  $Ax = b$ , so folgt  $F(x) = \langle A^{-1}0, 0 \rangle = 0$ .  
 Folglich löst  $x$  die Minimierungsaufgabe (M).

2)  $\Rightarrow$  1):  $x$  löst (M)  $\implies \nabla F(x) = 0 \iff 2(Ax - b) = 0 \iff Ax - b = 0 \implies$  1).

□

**Idee der Gradientenverfahren 2.42**

Aus Lemma 2.41 folgt:  $Ax = b \iff F(x) = 0$ .

**Idee:** Sei  $(z_n)_{n \in \mathbb{N}}$  eine Folge im  $\mathbb{R}^m$  definiert durch

$$z_{n+1} := z_n + \alpha_n t_n \quad n = 1, 2, \dots$$

mit Koeffizienten  $\alpha_n \in \mathbb{R}$  und Richtungsvektoren  $t_n \in \mathbb{R}^m \setminus \{0\}$ .  
 Wähle  $\alpha_n, t_n$  so, dass  $F(z_n) \longrightarrow 0 \quad (n \rightarrow \infty)$ .

**Ansatz für die Wahl von  $\alpha_n$ :**

$$\alpha_n := \beta_n \frac{\langle t_n, r_n \rangle}{\langle At_n, t_n \rangle}$$

mit  $\beta_n \in \mathbb{R}$  und  $r_n := b - Az_n$  der "Residuenvektor".

Dann gilt:

$$\begin{aligned}
 (*) \quad \left| \begin{aligned}
 F(z_n) - F(z_{n+1}) &= \langle A^{-1}r_n, r_n \rangle - \langle A^{-1}(r_n - A\alpha_n t_n), r_n - A\alpha_n t_n \rangle \\
 &= 2\alpha_n \langle t_n, r_n \rangle - \alpha_n^2 \langle At_n, t_n \rangle \\
 &= 2\beta_n \frac{\langle t_n, r_n \rangle^2}{\langle At_n, t_n \rangle} - \beta_n^2 \frac{\langle t_n, r_n \rangle^2}{\langle At_n, t_n \rangle} \\
 &= (2 - \beta_n) \underbrace{\beta_n \frac{\langle t_n, r_n \rangle^2}{\langle At_n, t_n \rangle}}_{\geq 0} \stackrel{0 \leq \beta_n \leq 2}{\geq} 0.
 \end{aligned} \right.
 \end{aligned}$$

$\Rightarrow$  Für  $0 \leq \beta_n \leq 2$  ist  $F(z_n) \leq F(z_{n+1})$ .

$\Rightarrow (F(z_n))_{n \in \mathbb{N}}$  ist monoton fallend.

Also folgt:  $\lim_{n \rightarrow \infty} F(z_n) = \lambda \geq 0$  existiert.

$\Rightarrow (F(z_n) - F(z_{n-1}))_{n \in \mathbb{N}}$  ist Nullfolge.

Man nennt  $\beta_n$  den Relaxationsparameter der Gradientenverfahren.

Die Gleichung (\*) zeigt:

Für  $\beta_n = 1$  wird  $F(z_n) - F(z_{n+1})$  maximal und  $F$  nimmt auf der Geraden  $z_n + \alpha_n t_n$  ein Minimum an.

#### Definition 2.43 (Allgemeines Gradientenverfahren)

Seien  $(\beta_n)_{n \in \mathbb{N}}$  und  $(t_n)_{n \in \mathbb{N}}$  gegeben mit  $\beta_n \in [0, 2]$  und  $t_n \in \mathbb{R}^m$ . Dann heißt die Folge  $(z_n)_{n \in \mathbb{N}}$ ,  $z_n \in \mathbb{R}^m$  Lösung des Gradientenverfahren mit Startwert  $z_1 \in \mathbb{R}^m$  wenn gilt:

$$r_1 = b - Az_1$$

und für  $n = 1, 2, \dots$  gilt

$$\alpha_n = \beta_n \frac{\langle t_n, r_n \rangle}{\langle At_n, t_n \rangle},$$

$$z_{n+1} = z_n + \alpha_n t_n,$$

$$r_{n+1} = b - Az_{n+1} = r_n - \alpha_n At_n.$$

#### Definition 2.44 (Konvergenz)

Ein Gradientenverfahren heißt konvergent, falls gilt

$$r_n \longrightarrow 0 \quad (n \rightarrow \infty)$$

Dies ist äquivalent zu

$$\begin{aligned}
 & z_n \longrightarrow A^{-1}b \quad (n \rightarrow \infty) \\
 \iff & z_n \longrightarrow x \quad (n \rightarrow \infty) \text{ mit } Ax = b.
 \end{aligned}$$

### 2.4.1 Eigentliches Gradientenverfahren

**Definition 2.45 (Eigentliches Gradientenverfahren)**

Das Gradientenverfahren 2.43 mit  $t_n = r_n$  und  $\beta_n = 1$  heißt eigentliches Gradientenverfahren. Die Richtungsvektoren werden in Richtung des Gradienten von  $F$  gewählt:

$$r_n = b - Az_n = -\frac{1}{2}\nabla F(z_n).$$

**Satz 2.46 (Konvergenz)**

Das eigentliche Gradientenverfahren 2.45 ist konvergent.

*Beweis:* Es gilt  $\forall x \in \mathbb{R}^m : \langle Ax, x \rangle \stackrel{\text{Scharzsche Ungleichung}}{\leq} \left( \sum_{i,j=1}^m a_{ij}^2 \right)^{\frac{1}{2}} \langle x, x \rangle$ .

Setze  $k := \left( \sum_{i,j=1}^m a_{ij}^2 \right)^{\frac{1}{2}}$ . Dann gilt  $\forall x \in \mathbb{R}^m \setminus \{0\}$ :

$$\frac{\langle Ax, x \rangle}{\langle x, x \rangle} \leq k.$$

Mit (\*) aus 2.42 folgt:

$$F(z_n) - F(z_{n+1}) = \frac{\langle r_n, r_n \rangle^2}{\langle Ar_n, r_n \rangle} \geq \frac{1}{k} \langle r_n, r_n \rangle.$$

Aus 2.42 wissen wir, dass  $F(z_n) - F(z_{n+1}) \rightarrow 0$  für  $(n \rightarrow \infty)$ . Also folgt für  $k \neq 0$ :  $r_n \rightarrow 0$  ( $n \rightarrow \infty$ ). □

**Bemerkung 2.47**

Je zwei aufeinanderfolgende Residuenvektoren des eigentlichen Gradientenverfahrens stehen senkrecht aufeinander:

$$\begin{aligned} \langle r_n, r_{n+1} \rangle &= \langle r_n, r_n - \alpha_n A t_n \rangle \\ &= \left\langle r_n, r_n - \frac{\langle r_n, r_n \rangle}{\langle Ar_n, r_n \rangle} Ar_n \right\rangle \\ &= \langle r_n, r_n \rangle - \langle r_n, r_n \rangle = 0. \end{aligned}$$

**Definition 2.48 (Gradientenverfahren bezüglich der kanonischen ON-Basis des  $\mathbb{R}^m$ )**

Wähle  $\beta_n = 1$  und  $t_{i+jm} = e_i$  für  $i = 1, \dots, m$ ;  $j = 0, 1, 2, \dots$ , wobei  $e_i \in \mathbb{R}^m$  der  $i$ -te Einheitsvektor ist. Dann folgt mit  $n = i + jm$ :

$$\begin{aligned} \alpha_n = \frac{\langle e_i, r_n \rangle}{\langle Ae_i, e_i \rangle} &= \frac{1}{a_{ii}} \langle e_i, b - Az_n \rangle \\ &= \frac{1}{a_{ii}} \left( b_i - \sum_{l=1}^m a_{il} z_{n,l} \right) \end{aligned}$$

$$\implies z_{n+1} = z_n + \frac{1}{a_{ii}} \left( b_i - \sum_{l=1}^m a_{il} z_{n,l} \right) e_i.$$

D.h.  $z_{n+1}$  und  $z_n$  unterscheiden sich nur in der  $i$ -ten Komponente.

Das ist das Einschrittverfahren (siehe Numerik I).

**Satz 2.49 (Konvergenz)**

Das Gradientenverfahren 2.48 ist konvergent.

*Beweis:* Setze  $\gamma := \max_{i=1,\dots,m} a_{ii}$ . Dann ist mit (\*) aus 2.4

$$F(z_n) - F(z_{n+1}) \stackrel{n=i+jm}{=} \frac{\langle e_i, r_n \rangle^2}{\langle Ae_i, e_i \rangle} \geq \frac{1}{\gamma} (r_{n,i})^2 \geq 0.$$

Da nach 2.4  $F(z_n) - F(z_{n+1}) \rightarrow 0$ , folgt  $r_{n,i} \rightarrow 0 \forall i = 1, \dots, m$  und somit  $r_n \rightarrow 0$  ( $n \rightarrow \infty$ ).  $\square$

**2.4.2 Conjugate Direction Verfahren (CD)****Definition 2.50 (A-orthogonal)**

Ein System von  $k$  Vektoren  $q_1, \dots, q_k \in \mathbb{R}^m$  ( $k \leq m$ ) heißt A-orthogonal oder A-konjugiert, wenn gilt:

$$\begin{aligned} \langle Aq_i, q_j \rangle &= 0 \quad i \neq j; \quad i, j = 1, \dots, k, \\ \langle Aq_i, q_i \rangle &\neq 0 \quad \forall i = 1, \dots, k. \end{aligned}$$

Für  $k = m$  bilden A-orthogonale Systeme eine Basis des  $\mathbb{R}^m$ . Wir setzen  $\langle x, y \rangle_A = \langle Ax, y \rangle$ ,  $\|x\|_A = \sqrt{\langle x, x \rangle_A} = \sqrt{\langle Ax, x \rangle}$ .

**Definition 2.51 (cd-Verfahren) "conjugate direction"**

Sei  $\{q_i\}_{i=1}^m$  A-orthogonal. Wähle im allgemeinen Gradientenverfahren 2.43:

$$t_n = q_n \text{ und } \beta_n = 1 \text{ für } n = 1, \dots, m.$$

**Satz 2.52 (Konvergenz der cd-Methode)**

Das cd-Verfahren ist für beliebige Startvektoren  $z_1 \in \mathbb{R}^m$  ein endliches Verfahren. Es gilt:

$$r_{m+1} = 0$$

und somit

$$z_{m+1} = A^{-1}b = x.$$

*Beweis:* Laut Gradientenverfahren folgt induktiv für  $1 \leq n \leq m+1$ :

$$r_n = r_{n-1} - \alpha_{n-1} A t_{n-1} = r_{n-2} - \alpha_{n-2} A t_{n-2} - \alpha_{n-1} A t_{n-1} = \dots = r_{n-l} - \sum_{i=n-l}^{n-1} \alpha_i A t_i.$$

Für  $l = n - j$ ,  $1 \leq j \leq n$  folgt

$$r_n = r_j - \sum_{i=j}^{n-1} \alpha_i A q_i.$$

$$\begin{aligned} \implies \langle q_j, r_n \rangle &\stackrel{\text{A-orthogonal}}{=} \langle q_j, r_j \rangle - \alpha_j \langle Aq_j, q_j \rangle \\ &\stackrel{\text{Def. von } a_j}{=} \langle q_j, r_j \rangle - \langle q_j, r_j \rangle = 0. \end{aligned}$$

$$\implies \langle q_j, r_{m+1} \rangle = 0 \quad \forall j = 1, \dots, m.$$

$$\implies r_{m+1} = 0, \text{ da } \{q_i\}_{i=1}^m \text{ Basis des } \mathbb{R}^m.$$

□

### 2.4.3 Conjugate Gradient Verfahren (CG)

**Definition 2.53 (cg-Verfahren) "conjugate gradient"**

Wähle  $\beta_n = 1 \quad \forall n = 1, \dots, m+1$ ,  $t_1 = r_1$  und

$$t_n = r_n + \gamma_{n-1} t_{n-1}$$

mit  $\gamma_{n-1} := -\frac{\langle Ar_n, t_{n-1} \rangle}{\langle At_{n-1}, t_{n-1} \rangle}$  für  $n = 2, 3, \dots$ . Idee: Das  $A$ -orthogonale System wird mit Hilfe des Schmitd'schen Orthogonalisierungsverfahrens bzgl.  $\langle \cdot, \cdot \rangle_A$  schrittweise aufgebaut.

**Satz 2.54 (Konvergenz des cg-Verfahrens)**

Sei  $z_1 \in \mathbb{R}^m$  und  $t_1 = r_1 = b - Az_1$ . Dann lautet das cg-Verfahren für  $n = 1, \dots, l$  mit  $l \leq m$ :

$$\left\{ \begin{array}{l} a_n = \frac{\langle t_n, r_n \rangle}{\langle At_n, t_n \rangle} \\ z_{n+1} = z_n + a_n t_n \\ r_{n+1} = b - Az_{n+1} \\ \gamma_n = -\frac{\langle Ar_{n+1}, t_n \rangle}{\langle At_n, t_n \rangle} \\ t_{n+1} = r_{n+1} + \gamma_n t_n \end{array} \right.$$

Es erfüllt die Gleichungen

a)  $\langle At_i, t_j \rangle = 0 \quad 1 \leq j \leq i-1$

b)  $\langle r_i, r_j \rangle = 0 \quad 1 \leq j \leq i-1$

c)  $\langle t_i, r_j \rangle = \langle r_j, r_j \rangle \quad 1 \leq j \leq i$

und  $l \leq m$  ist so gewählt, dass gilt  $r_{l+1} = 0$ .

*Beweis:* Folgt aus linearer Algebra und Satz 2.52. □

**Bemerkung 2.55**

Für Gleichungssysteme resultierend aus Diskretisierungsverfahren, wie z.B. der Finit Elemente Methode, ist  $m$  so groß, dass man in der Praxis weniger als  $m$  Schritte iterieren wird. Dass dies Sinn macht zeigt die folgende Fehlerabschätzung.

**Satz 2.56**

Sei  $\kappa(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$  die Kondition von  $A \in \mathbb{R}^{n \times n}$ . Dann gilt für das cg-Verfahren für  $Ax = b$ :

$$\|z_n - x\|_A \leq 2 \left( \frac{\sqrt{\kappa(A)} - 1}{\sqrt{\kappa(A)} + 1} \right)^n \|z_1 - x\|_A$$

wobei  $\|y\|_A := \sqrt{\langle Ay, y \rangle}$  definiert ist.  
(ohne Beweis)

**Bemerkung 2.57**



- 1) Vorkonditionierung: Anstelle von  $Ax = b$  löst man  $C^{-1}A(C^{-1})^T y = C^{-1}b$ , wobei  $y = C^T x$  ist und  $C$  auch folgende Anforderungen erfüllt
- 1)  $C \in \mathbb{R}^{n \times n}$  regulär.
  - 2) Gleichungssysteme mit  $C$  sollen einfach zu lösen sein.
  - 3)  $\kappa(C^{-1}A(C^{-1})^T)$  sollt möglich nahe an 1 liegen.
- 2) Ist  $A$  nicht symmetrisch und positiv definit, so kann man z.B.  $A^T Ax = A^T b$  anstelle von  $Ax = b$  lösen, denn  $A^T A$  ist positiv definit und symmetrisch.  
Dieser Ansatz führt auf das bicg-Verfahren.
- 3) Für eine gegebene Toleranz  $TOL$ , kann

$$\langle r_n, r_n \rangle \leq TOL$$

als Abbruchbedingung verwendet werden.

## 2.5 Zusammenfassung

Wir haben in diesem Kapitel verschiedene Methoden zur Lösung linearer Gleichungssysteme der Form

$$Ax = b$$

kennengelernt und näher betrachtet,

- (a) für  $A \in \mathbb{R}^{n \times n}$  regulär
- (b) für  $A \in \mathbb{R}^{m \times n}$  mit  $m \geq n$  (d.h. überstimmt)

### Verfahren:

- (a) Direkte Verfahren oder Direktlöser (LR-, Cholesky, QR-Zerlegung)
- (b) Iterative Verfahren oder Iterativlöser (Jacobi-, Gauß-Seidel-Verfahren)

### Vor-/Nachteile

- **Direkte Verfahren:** Die Lösung wird bis auf Rundungsfehler exakt berechnet. Sie sind für große Gleichungssysteme sehr langsam (die Anzahl der arithmetischen Operationen liegen in  $O(n^3)$ ); es tritt ein *fill-in* Problem auf: In Anwendungen ist  $A$  häufig dünn besetzt, d.h. pro Zeile sind nur  $k$  viele Einträge ungleich Null, wobei  $k$  unabhängig von  $n$ , und diese Einträge sind unstrukturiert verteilt. Eine typische Zeile sieht dann so aus:

$$\begin{pmatrix} * & & & & 0 \\ 0 & \ddots & & & \\ \hline * & 0 \cdots 0 & * & * 0 \cdots 0 & * \\ \hline & & & \ddots & \\ & & & & * \end{pmatrix}$$

Bei Zerlegungsverfahren werden Nulleinträge u.a. durch Einträge ungleich Null ersetzt, der Speicheraufwand zum Speichern von  $A$  liegt in der Regel bei  $O(NK) = O(N)$ ; nach der Zerlegung wächst der Speicheraufwand bis auf  $O(N^2)$ .

- **Iterative Verfahren:** Die Lösung wird nur näherungsweise bestimmt und die Geschwindigkeit der Konvergenz hängt von  $\rho(T)$  (Spektralradius) ab.

Allerdings ist der zusätzliche Speicheraufwand sehr klein (das Gauß-Seidel-Verfahren hat gar kein zusätzlicher Speicheraufwand). Iterative Verfahren benötigen ein gutes Abbruchkriterium, z.B.  $\|Ax^{(k)} - b\| < TOL$ .

### Beschleunigung/Stabilisierung

Das Hauptproblem: Einfluss durch Rundungsfehler und ihre Reduzierung.

**Beispiel:** Die Pivotisierung bei der LR-Zerlegung. Die Kondition des Problems ist proportional zur Kondition von  $A$ :  $cond(A) = \|A\| \|A^{-1}\|$ . Durch **Vorkonditionierung** kann versucht werden, die Kondition des Problems zu verkleinern.

**Beispiel:** Wähle  $C_1, C_2$  reguläre Matrizen. Dann gilt:

$$Ax = b \iff \underbrace{C_1 A C_2}_{:=\tilde{A}} \underbrace{C_2^{-1} x}_{:=\tilde{x}} = \underbrace{C_1 b}_{:=\tilde{b}}$$

mit  $\tilde{A} := C_1 A C_2$ ,  $\tilde{x} := C_2^{-1} x$ ,  $\tilde{b} := C_1 b$  und  $C_1, C_2$  so gewählt, dass  $cond(\tilde{A}) < cond(A)$  gilt.

### Beschleunigung durch Relaxation:

Ein Iterationsverfahren hat die Gestalt

$$r^k = b - Ax^k, \quad My^k = r^k, \quad x^{k+1} = x^k + y^k$$

Statt  $y^k$  zu korrigieren, wählt man einen Relaxationsparameter  $w > 0$  und setzt

$$x^{k+1} = x^k + wy^k.$$

Dies führt auf das **SOR-Verfahren**<sup>2</sup>.

---

<sup>2</sup>Succesive Over Relaxation.



## Kapitel 3

# Nichtlineare Gleichungen/ Nullstellensuche

In diesem Kapitel wollen wir uns mit der Berechnung von Nullstellen einer gegebenen Funktion

$$f : \mathbb{R}^n \longrightarrow \mathbb{R}^n$$

beschäftigen. Da man die Lösung linearer und nichtlinearer Gleichungen auch als Nullstellensuche einer geeigneten Funktion  $f$  auffassen kann, ist diese Fragestellung eine direkte Verallgemeinerung der Lösung linearer Gleichungsprobleme, die in Kapitel 2 behandelt wurde. Bei der Bestimmung von Nullstellen unterscheiden wir zwei Problemstellungen.

**Problem A:** Gesucht ist ein  $x^* \in \mathbb{R}^n$  mit  $f(x^*) = 0$ .

**Problem B:** Gesucht sind alle (größtes, kleinstes)  $x^* \in \mathbb{R}^n$  mit  $f(x^*) = 0$ .

**Beispiel:**

$f(x) = Ax - b$ ,  $A \in \mathbb{R}^{m \times n}$  (siehe Kap. 2).

$f(x) = ax^2 + bx + c$ . Hier sind alle Nullstellen explizit berechenbar.

$f(x) = \cos(x)$ . Gesucht  $x^* \in [1, 2] \implies x^* = \frac{\pi}{2}$ .

Wir beschäftigen uns im wesentlichen mit Verfahren zur Lösung vom Problem **A** für  $n = 1$  und

$$f : [a, b] \longrightarrow \mathbb{R}$$

glatt.

Wir nehmen an, dass es ein  $x^* \in [a, b]$  gibt mit  $f(x^*) = 0$ , etwa  $f \in C^0(a, b)$  und  $f(a)f(b) < 0$ . Alle Verfahren konstruieren eine Folge  $(x^{(k)})_{k \in \mathbb{N}}$  mit  $x^{(k)} \longrightarrow x^*$ ,  $f(x^*) = 0$ . Direkte Verfahren existieren nur in Spezialfällen.

### 3.1 Verfahren in einer Raumdimension

#### 3.1.1 Intervallschachtelungsverfahren (ISV)

Voraussetzungen: Sei  $f \in C^0(a, b)$  mit  $a < b$  und  $f(a)f(b) < 0$ .

Verfahren: Setze  $a_0 := a$ ,  $b_0 := b$  und  $x^{(0)} := \frac{1}{2}(a_0 + b_0)$ .

Für  $n = 0, \dots, N$ :

Falls  $f(a_n)f(x^{(n)}) = 0$ , dann Abbruch.

Falls  $f(a_n)f(x^{(n)}) < 0 \implies a_{n+1} := a_n; b_{n+1} := x^{(n)}$

Falls  $f(a_n)f(x^{(n)}) > 0 \implies a_{n+1} := x^{(n)}; b_{n+1} := b_n$

Setze  $x^{(n+1)} := \frac{1}{2}(a_{n+1} + b_{n+1})$ .

### Satz 3.1 (Konvergenz von ISV)

Seien  $(a_n)_{n \in \mathbb{N}}, (b_n)_{n \in \mathbb{N}}, (x^{(n)})_{n \in \mathbb{N}}$  durch das Intervallschachtelungsverfahren definiert. Dann gilt

$$\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n = \lim_{n \rightarrow \infty} x^{(n)} = x^* \text{ und } f(x^*) = 0.$$

Es gilt die Fehlerabschätzung:

$$|x^{(n)} - x^*| \leq 2^{-(n+1)} |b - a|.$$

*Beweis:* Es gilt  $(a_n)_{n \in \mathbb{N}}$  monoton wachsend und  $(b_n)_{n \in \mathbb{N}}$  monoton fallend und  $0 < b_n - a_n = 2^{-n}(b - a)$ ,  $a_n < b, a < b_n$ . Daraus folgt  $(a_n)_{n \in \mathbb{N}}, (b_n)_{n \in \mathbb{N}}$  konvergieren und  $\lim_{n \rightarrow \infty} a_n = \lim_{n \rightarrow \infty} b_n =: x^*$ . Da  $f \in C^0$  folgt weiter  $f(x^*) = \lim_{n \rightarrow \infty} f(a_n) = \lim_{n \rightarrow \infty} f(b_n)$ .

Nach Konstruktion gilt stets  $f(a_n)f(b_n) < 0$ , also folgt

$$f(x^*)^2 = \lim_{n \rightarrow \infty} (f(a_n)f(b_n)) \leq 0 \implies f(x^*) = 0.$$

Da  $x^* \in (x^{(n)}, b_n)$  oder  $x^* \in (a_n, x^{(n)})$  folgt auch

$$|x^{(n)} - x^*| \leq \min \left\{ |x^{(n)} - b_n|, |x^{(n)} - a_n| \right\} = \frac{1}{2} |b_n - a_n| \leq 2^{-n-1}(b - a).$$

**Bemerkung:** Das Verfahren ist sehr robust aber auch sehr langsam: Man benötigt etwa 3 Schritte, um eine Dezimalzahlstelle Genauigkeit zu gewinnen.

Aufwand: Eine Auswertung von  $f$  pro Schritt.

Das Verfahren wird in der Regel eingesetzt, um grob ein Intervall  $[a, b]$  zu bestimmen, in dem eine Nullstelle liegt. Zur genaueren Berechnung der Nullstellen werden dann effizientere Verfahren eingesetzt.

### 3.1.2 Newton Verfahren

**Idee:** Sei  $x^{(k)}$  eine gegebene Approximation von  $x^*$ , d.h.  $h := x^* - x^{(k)} \neq 0$  ist klein. Dann gilt mit der Taylorentwicklung:

$$0 = f(x^*) = f(x^{(k)} + h) = f(x^{(k)}) + f'(x^{(k)})h + \mathcal{O}(h^2).$$

Unter Vernachlässigung der Terme höherer Ordnung folgt

$$h = -\frac{f(x^{(k)})}{f'(x^{(k)})}$$

Daher sollte  $x^{(k+1)} = x^{(k)} + h = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$  eine bessere Approximierung von  $x^*$  sein als  $x^{(k)}$ .

Verfahren: Sie Startwert  $x^{(0)}$  gegeben. Setze iterativ:

$$x^{(k+1)} := x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}.$$

**Aufwand:** Eine Auswertung von  $f$  und  $f'$  pro Schritt, d.h. es muß gelten  $f \in C^1$  und  $f'$  muss bekannt sein.

### Geometrische Interpretation

Sei  $l(x)$  die Linearisierung von  $f$  an der Stelle  $x^{(k)}$ , d.h.  $l(x^{(k)}) = f(x^{(k)})$ ,  $l'(x^{(k)}) = f'(x^{(k)})$  und  $l(x) = ax + b$ . Dann folgt aus der Taylorentwicklung  $l(x) = f(x^{(k)}) + f'(x^{(k)})(x - x^{(k)})$ .

Statt Nullstellen von  $f$  zu suchen, definieren wir  $x^{(n+1)}$  Nullstelle von  $l$

$$0 = f(x^{(k)}) + f'(x^{(k)})(x^{(k+1)} - x^{(k)}).$$

Dies ist gerade das Newton Verfahren.

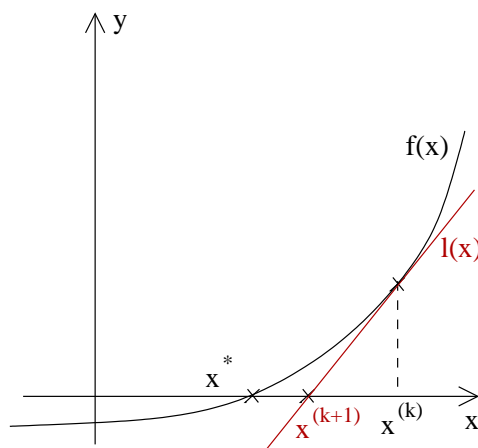


Abbildung 3.1: Newton Verfahren, Beispiel 1

Abbildung 3.1 veranschaulicht das Newton Verfahren. Mit dieser Anschauung sehen wir direkt ein, dass das Newton Verfahren nicht für alle Startwerte  $x^{(0)}$  konvergieren wird. Abbildung 3.2 zeigt eine solche Situation. Hier gilt  $|x^{(k)}| \rightarrow \infty$ . In Abbildung 3.3 gibt es zwei Nullstellen, aber es ist nicht klar, welche der beiden Nullstellen gefunden wird.

### Satz 3.2 (Konvergenz des Newton-Verfahrens)

Sei  $f \in C^2(a, b)$  und es existiere ein  $x^* \in (a, b)$  mit  $f(x^*) = 0$ . Sei  $m := \min_{a \leq x \leq b} |f'(x)| > 0$  und  $M := \max_{a \leq x \leq b} |f''(x)|$ . Sei  $\rho > 0$  so gewählt, dass  $B_\rho(x^*) := \{x \mid |x - x^*| < \rho\} \subset [a, b]$  und  $q := \frac{M}{2m}\rho < 1$ . Dann konvergiert das Newton-Verfahren für jeden Startwert  $x^{(0)} \in B_\rho(x^*)$ . Es gilt die a-priori Fehlerschranke

$$(a) \quad |x^{(k)} - x^*| \leq \frac{M}{2m} |x^{(k-1)} - x^*| \leq \frac{2m}{M} q^{2^k}$$

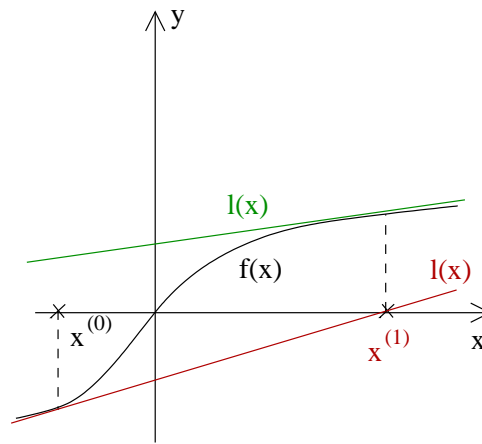


Abbildung 3.2: Newton Verfahren, Beispiel 2

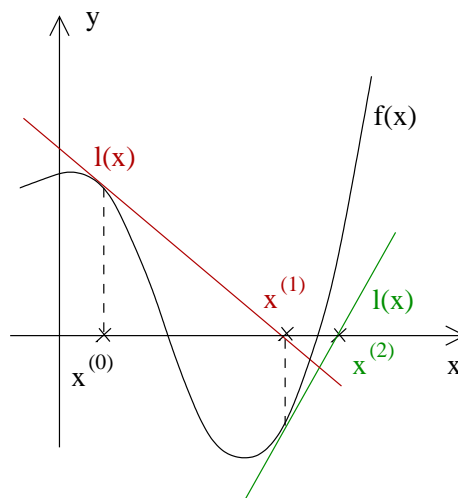


Abbildung 3.3: Newton Verfahren, Beispiel 3

und die a-posteriori Fehlerschranke

$$(b) \quad |x^{(k)} - x^*| \leq \frac{1}{m} |f(x^{(k)})| \leq \frac{M}{2m} |x^{(k)} - x^{(k-1)}|^2$$

**Bemerkung:** Aus dem Mittelwertsatz folgt  $\left| \frac{f(x) - f(y)}{x - y} \right| = |f'(\xi)| \geq m \quad \forall x, y \in B_\rho(x^*); x \neq y \implies |x - y| \leq \frac{1}{m} |f(x) - f(y)|$ . Folglich ist  $x^*$  die einzige Nullstelle in  $B_\rho(x^*)$  und  $x^*$  ist **einfache** Nullstelle, d.h.  $f(x^*) = 0$  und  $f'(x^*) \neq 0$ .

*Beweis:* Nach Taylorentwicklung (Satz 1.33) gilt

$$(1) \quad f(y) = f(x) + f'(x)(y - x) + R(y, x) \text{ mit}$$

$$R(y, x) = \int_x^y f''(\xi)(y - \xi) d\xi = (y - x)^2 \int_0^1 f''(x + s(y - x))(1 - s) ds.$$

Also folgt für alle  $x, y \in B_\rho(x^*)$

$$(2) \quad |R(y, x)| \leq |y - x|^2 M \int_0^1 (1 - s) ds = \frac{M}{2} |y - x|^2$$

Für  $x \in B_\rho(x^*)$  setze  $\Phi(x) := x - \frac{f(x)}{f'(x)}$ . Dann folgt

$$\begin{aligned} |\Phi(x) - x^*| &= \left| \left( x - x^* \right) - \frac{f(x)}{f'(x)} \right| = \left| -\frac{1}{f'(x)} [f(x) + (x^* - x)f'(x)] \right| \\ &\stackrel{(1)}{=} \left| \frac{1}{f'(x)} R(x^*, x) \right| = \frac{1}{|f'(x)|} |R(x^*, x)| \stackrel{(2)}{\leq} \frac{1}{m} |x - x^*|^2 \frac{M}{2} \end{aligned}$$

Also folgt für  $x \in B_\rho(x^*)$

$$(3) \quad \begin{aligned} |\Phi(x) - x^*| &\leq \frac{M}{2m} |x - x^*|^2, \\ &\leq \frac{M}{2m} \rho^2 =: q\rho < \rho, \text{ da } q = \frac{M}{2m}\rho < 1. \end{aligned}$$

Insbesondere folgt  $x^{(k)} \in B_\rho(x^*)$ , falls  $x^{(0)} \in B_\rho(x^*)$ .

Sei  $\rho^{(k)} := \frac{M}{2m} |x^{(k)} - x^*|$ , so gilt wegen (3)

$$\begin{aligned} \rho^{(k)} &= \frac{M}{2m} |\Phi(x^{(k-1)}) - x^*| \leq \frac{M}{2m} \left( \frac{M}{2m} |x^{(k-1)} - x^*|^2 \right) \\ &= (\rho^{(k-1)})^2 \leq \dots \leq (\rho^{(0)})^{2^k} \\ \implies |x^{(k)} - x^*| &\leq \frac{2m}{M} \rho^{(k)} \leq \frac{2m}{M} (\rho^{(0)})^{2^k} = \frac{2m}{M} \left( \frac{M}{2m} |x^{(0)} - x^*| \right)^{2^k} \\ &\leq \frac{2m}{M} \underbrace{\left( \frac{M}{2m} \rho \right)^{2^k}}_{=q} = \frac{2m}{M} q^{2^k}. \end{aligned}$$

Dies ist die a-priori Abschätzung. Da  $q < 1$  ist, folgt  $q^{(2^k)} \rightarrow 0 \implies x^{(k)} \rightarrow x^*$ .

Für die a-posteriori Abschätzung benutzen wir nochmal (1) mit  $y = x^{(k)}$  und  $x = x^{(k-1)}$ . Es folgt

$$\begin{aligned} f(x^{(k)}) &= f(x^{(k-1)}) + (x^{(k)} - x^{(k-1)}) f'(x^{(k-1)}) + R(x^{(k)}, x^{(k-1)}) \\ &= R(x^{(k)}, x^{(k-1)}), \text{ da } x^{(k)} = x^{(k-1)} - \frac{f(x^{(k-1)})}{f'(x^{(k-1)})} \end{aligned}$$

und somit

$$\begin{aligned} |x^{(k)} - x^*| &\stackrel{MWS}{\leq} \frac{1}{m} |f(x^{(k)}) - f(x^*)| = \frac{1}{m} |f(x^{(k)})|, \\ &\stackrel{(2)}{=} \frac{1}{m} R(x^{(k)}, x^{(k-1)}) \leq \frac{M}{2m} |x^{(k)} - x^{(k-1)}|^2. \end{aligned}$$

□

**Bemerkung:** Falls  $x^{(0)} \in B_\rho(x^*)$ , so konvergiert das Newton-Verfahren sehr schnell. Sei zum Beispiel  $q = \frac{1}{2}$ , so gilt nach 10 Iterationen  $|x^{(10)} - x^*| \leq \frac{2m}{M} q^{1024} \sim \frac{2m}{M} 10^{-303}$ .

Vergleichen wir dies mit dem Intervallschachtelungsverfahren, so folgt mit dem selben Startintervall:

$$\begin{aligned} |b - a| = \rho &= q \frac{2m}{M} = \frac{m}{M}, \text{ falls } q = \frac{1}{2} \text{ gilt. Nach 10 Schritten gilt also} \\ |x^{(10)} - x^*| &= 2^{-11} |b - a| = 2^{-11} \frac{m}{M} \sim 10^{-4} \frac{2m}{M}. \end{aligned}$$



**Folgerung 3.3**

Für  $f \in C^2(\mathbb{R})$  existiert für jede einfache Nullstelle  $x^*$  eine Umgebung  $U$  um den Wert  $x^*$ , so dass das Newton-Verfahren für alle  $x^{(0)} \in U$  konvergiert und  $|x^{(k)} - x^*| \leq q^{2^k}$  für  $q < 1$ .

*Beweis:* Übungsaufgabe. □

**Offene Fragen:**

- (a) Wie kann  $\rho$  effektiv bestimmt werden?
- (b) Kann die Berechnung von  $f'$  umgangen werden?
- (c) Was passiert, falls  $f'(x^*) = 0$  ist, d.h falls  $x^*$  mehrfache Nullstelle ist?
- (d) Gibt es noch schnellere Verfahren?

**Kombination von Newton Verfahren und Intervallschachtelung**

**Idee:** ISV einsetzen, um ausreichend nahe an eine Nullstelle  $x^*$  zu kommen, so dass das Newton Verfahren schnell konvergiert.

**Verfahren:** Seien  $a < b$  gegeben mit  $f(a)f(b) < 0$ .

Setze  $x := \frac{1}{2}(a + b)$ ,  $\tilde{a} := a$ ,  $\tilde{b} := b$ ,  $f_0 := f(x)$ ,  $f_a := f(a)$ .

Solange  $|f_0| > TOL$ :

$$\left[ \begin{array}{l} \text{Falls } f_a f_0 < 0, \text{ dann } \tilde{b} := x, \text{ sonst } (\tilde{a} := x; f_a := f_0) \\ x := x - \frac{f_0}{f'(x)} \\ f_1 := f(x) \\ \text{Falls } (|f_1| > |f_0| \text{ oder } x \notin (a, b)), \text{ dann} \\ \quad \left[ a := \tilde{a}, b := \tilde{b}, x := \frac{1}{2}(a + b), f_1 := f(x) \right] \\ f_0 := f_1 \end{array} \right.$$

**3.1.3 Sekantenverfahren**

Ein Nachteil des Newton-Verfahrens ist die Auswertung von  $f'(x^{(k)})$ .

**Idee:** Ersetze die Ableitung durch einen Differenzenquotient

$$f'(x^{(k)}) \sim \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}.$$

**Sekantenverfahren:**

$$x^{(k+1)} = x^{(k)} - f(x^{(k)}) \frac{x^{(k)} - x^{(k-1)}}{f(x^{(k)}) - f(x^{(k-1)})}.$$

**Geometrische Interpretation:**

Die Sekante an  $f$  durch die Punkte  $x^{(k)}, x^{(k-1)}$  ist gegeben durch

$$\frac{y - f(x^{(k)})}{x - x^{(k)}} = \frac{f(x^{(k-1)}) - f(x^{(k)})}{x^{(k-1)} - x^{(k)}},$$

wobei  $y = y(x)$  die Gerade durch die Punkte  $(x^{(k-1)}, f(x^{(k-1)}))$ ,  $(x^{(k)}, f(x^{(k)}))$  ist, d.h.

$$y(x) = \frac{f(x^{(k-1)}) - f(x^{(k)})}{x^{(k-1)} - x^{(k)}} (x - x^{(k)}) + f(x^{(k)}).$$

$x^{(k+1)}$  wird also als Nullstelle der Sekante wählen.

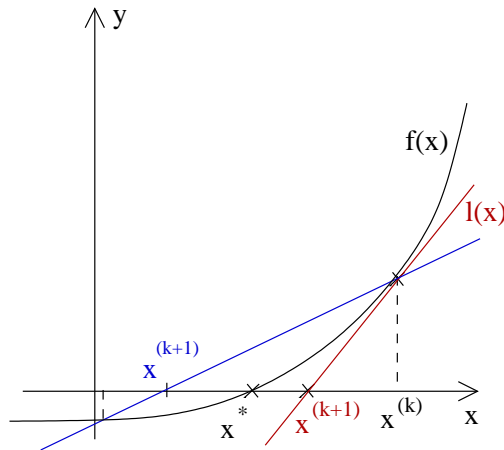


Abbildung 3.4: Sekantenverfahren, geometrische Interpretation

**Bemerkung:** Das Verfahren erfordert 2 Startwerte  $x^{(0)}$ ,  $x^{(1)}$  (siehe Abb. 3.4). Das Verfahren erfordert eine Auswertung von  $f$  pro Iterationsschritt (speichern von  $f(x^{(k-1)})$ ).

**Satz 3.4 (Konvergenz des Sekantenverfahrens)**

Sei  $f \in C^2(a, b)$  mit  $x^* \in [a, b]$ ,  $f(x^*) = 0$  und  $m := \min_{a \leq x \leq b} |f'(x)| > 0$ ,  $M := \max_{a \leq x \leq b} |f''(x)|$ . Sei

$q := \frac{M}{2m} \rho < 1$  für  $\rho > 0$ . Seien  $x^{(0)}$ ,  $x^{(1)} \in B_\rho(x^*)$ ,  $x^{(0)} \neq x^{(1)}$

und  $(x^{(k)})_{k \in \mathbb{N}}$  die Folge definiert durch das Sekantenverfahren.

Dann gilt  $x^{(k)} \in B_\rho(x^*)$  und  $x^{(k)} \rightarrow x^*$  für  $k \rightarrow \infty$ . Es gelten folgende Fehlerabschätzungen.

(a) A-priori Fehlerschranke:

$$|x^{(k)} - x^*| \leq \frac{2m}{M} q^{\gamma_k},$$

wobei  $(\gamma_k)_{k \in \mathbb{N}}$  die Folge der Fibonacci-Zahlen ist, d.h.  $\gamma_0 = \gamma_1 = 1$ ;  $\gamma_{k+1} = \gamma_k + \gamma_{k-1}$ .

(b) A-posteriori Fehlerschranke:

$$|x^{(k)} - x^*| \leq \frac{1}{m} |f(x^{(k)})| \leq \frac{M}{2m} |x^{(k)} - x^{(k-1)}| \cdot |x^{(k)} - x^{(k-2)}|.$$

*Beweis:* Übungsaufgabe. □

**Folgerung 3.5 (Konvergenz des Sekantenverfahrens)**

Für  $f \in C^2(\mathbb{R})$  existiert eine Umgebung  $U$  um jede einfache Nullstelle, so dass  $x^{(k)} \rightarrow x^*$  für  $x^{(0)}$ ,  $x^{(1)} \in U$  und es gilt  $|x^{(k)} - x^*| \leq \frac{2m}{M} \tilde{q}^{\alpha^k}$  mit  $\tilde{q} < 1$  und  $\alpha = \frac{1}{2}(1 + \sqrt{5}) \approx 1.618$ .

*Beweis:* Zunächst zeigt man analog zu Satz 3.3, dass es ein  $\rho > 0$  gibt, so dass die Voraussetzungen von Satz 3.4 auf  $U = B_\rho(x^*)$  erfüllt sind. Mit der A-priori-Abschätzung von Satz 3.4 ist noch zu zeigen, dass gilt  $q^{\gamma_k} \leq \tilde{q}^{\alpha^k}$  mit  $\tilde{q} < 1$  geeignet gewählt.

**Ansatz:**  $\gamma_k = \lambda^k$  mit  $\lambda > 0$ . Aus  $\gamma_k - \gamma_{k-1} - \gamma_{k-2} = 0$  folgt

$$\lambda^{k-2}(\lambda^2 - \lambda - 1) = 0 \iff \lambda^2 - \lambda - 1 = 0 \implies \lambda_{1/2} = \frac{1}{2}(1 \pm \sqrt{5}).$$

D.h. die allgemeine Lösung der Differenzengleichung  $\gamma_k - \gamma_{k-1} - \gamma_{k-2} = 0$  ist

$$\gamma_k = c_1 \lambda_1^k + c_2 \lambda_2^k$$

für  $c_1, c_2 \in \mathbb{R}$ .

Mit  $\gamma_0 = \gamma_1 = 1$  folgt  $c_1 = \frac{\lambda_1}{\sqrt{5}}$ ,  $c_2 = -\frac{\lambda_2}{\sqrt{5}}$ . Also  $\gamma_k = \frac{1}{\sqrt{5}}(\lambda_1^{k+1} - \lambda_2^{k+1})$ . Aus  $|\lambda_2| < |\lambda_1|$  folgt  $\gamma_k \geq \frac{\lambda_1}{2\sqrt{5}} \lambda_1^k$  und somit

$$\begin{aligned} q^{\gamma_k} &\leq q^{\frac{\lambda_1}{2\sqrt{5}}(\lambda_1^k)}, \text{ da } q < 1 \\ &= \tilde{q}^{(\alpha^k)} \end{aligned}$$

mit  $\tilde{q} := q^{\frac{\lambda_1}{2\sqrt{5}}} < 1$  und  $\alpha = \lambda_1$ . □

### 3.1.4 Zusammenfassung

Für die betrachteten Verfahren gilt:

**ISV:**  $|x^{(k)} - x^*| \leq \frac{(b-a)}{2} \left(\frac{1}{2}\right)^k$ , eine Auswertung von  $f$  pro Schritt.

**Newton:**  $|x^{(k)} - x^*| \leq \frac{2m}{M} q^{2^k}$ , je eine Auswertung von  $f$  und  $f'$  pro Schritt.

**Sekanten:**  $|x^{(k)} - x^*| \leq \frac{2m}{M} q^{1.618^k}$ , eine Auswertung von  $f$  pro Schritt.

**Annahme:** Die Auswertungen von  $f$  und  $f'$  seien gleich aufwändig. Dann hat das Newton-Verfahren pro Schritt den doppelten Aufwand.

#### Vergleich der Verfahren:

Definiere  $z^{(k)} := x^{2^k}$  beim ISV und Sekanten-Verfahren. Dann ist auch hier der Aufwand in einem Schritt durch zwei Auswertung von  $f$  gegeben. Es gilt:

**ISV:**  $|z^{(k)} - x^*| \leq \frac{(b-a)}{2} \left(\frac{1}{2}\right)^{2^k} \leq \frac{(b-a)}{2} \left(\frac{1}{4}\right)^k$ .

**Sekanten:**  $|z^{(k)} - x^*| \leq \frac{2m}{M} \tilde{q}^{(\lambda_1^{2^k})} = \frac{2m}{m} \tilde{q}^{(2.618^k)}$ .

Bei gleichen Aufwand konvergiert das Sekanten-Verfahren also schneller als das Newton- oder IS-Verfahren.

#### Problem beim Sekanten-Verfahren

Die Fehleranalyse wurde ohne Berücksichtigung von Rundungsfehlern gemacht. Seien  $x^{(k)}$ ,  $x^{(k-1)}$  sehr nah an  $x^*$ , dann liegen auch  $f(x^{(k)})$ ,  $f(x^{(k-1)})$  nahe zusammen und haben gleiches Vorzeichen. Da die Differenz in diesem Fall schlecht konditioniert ist, kann es hier zu Auslöschungen kommen.

## 3.2 Konvergenzordnung von Iterationsverfahren

Wir betrachten allgemeine Iterationsverfahren auf einem Banach-Raum  $X$  der Form

$$x^{(k+1)} = \Phi(x^{(k)}, \dots, x^{(k-j)}), \quad k \geq j \geq 0 \quad (1)$$

mit Startwerten  $x^{(0)}, \dots, x^{(j)}$  und Iterationsfunktion  $\Phi : X^j \rightarrow X$ .

### Definition 3.6 Lokale Konvergenz

Das Iterationsverfahren (1) **konvergiert lokal** gegen ein  $x^* \in X$ , falls es eine Umgebung  $U$  von  $x^*$  gibt, so dass für alle  $x^{(0)}, \dots, x^{(j)} \in U$  die Folge  $(x^{(k)})_{k \in \mathbb{N}}$  gegen  $x^*$  konvergiert.

### Definition 3.7 (Konvergenzordnung)

Sei  $X$  ein Banach-Raum und  $(x^{(k)})_{k \in \mathbb{N}}$  Folge in  $X$ , die gegen ein  $x^* \in X$  konvergiert. Die Folge hat mindestens die **Konvergenzordnung**  $p \geq 1$  falls

$$\limsup_{k \rightarrow \infty} \frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|^p} = c$$

mit  $c < 1$  für  $p = 1$  und  $c < \infty$  für  $p > 1$ .

Die **Ordnung ist genau  $p$** , falls  $c \neq 0$ .

Der Fall  $p = 1$  heißt **lineare Konvergenz** und falls für ein  $p \geq 1$  gilt

$$\limsup_{k \rightarrow \infty} \frac{\|x^{(k+1)} - x^*\|}{\|x^{(k)} - x^*\|^p} = 0,$$

so spricht man von **superlinearer Konvergenz**.

**Beispiel:** Das Newton-Verfahren konvergiert quadratisch ( $p = 2$ ), da  $|x^{(k+1)} - x^*| \leq c \cdot |x^{(k)} - x^*|^2$ . Das ISV Verfahren konvergiert linear, da  $|x^{(k+1)} - x^*| \leq \frac{1}{2} |x^{(k)} - x^*|$ .

### Satz 3.8 (Konvergenzordnung von Iterationsverfahren)

Sei  $I = [a, b]$ ,  $\Phi : I \rightarrow I$ ,  $\Phi \in C^p(I)$ . Sei  $x^{(0)} \in I$ ,  $x^{(k+1)} = \Phi(x^{(k)})$ ,  $k \geq 0$  und es gelte  $x^{(k)} \rightarrow x^*$  mit  $x^* = \Phi(x^*)$ . Dann gilt:

(i)  $p = 1$  :  $x^{(k)} \rightarrow x^*$  mind. linear  $\iff |\Phi'(x^*)| < 1$ .

(ii)  $p \geq 2$  :  $x^{(k)} \rightarrow x^*$  mind. mit der Ordnung  $p$ , g.d.w.

$$\Phi^{(\nu)}(x^*) = 0, \quad (1 \leq \nu \leq p-1).$$

*Beweis:* (i) Es existiert eine Zwischenstelle  $\xi_k$  zwischen  $x^k$  und  $x^*$ , so dass gilt

$$\lim_{k \rightarrow \infty} \frac{\|\Phi(x^{(k)}) - x^*\|}{\|x^{(k)} - x^*\|} = \lim_{k \rightarrow \infty} |\Phi'(\xi_k)| = |\Phi'(x^*)|.$$

Also folgt die Behauptung.

(ii) “ $\implies$ ”: Annahme:  $\exists j < p$  mit  $\Phi^{(j)}(x^*) \neq 0$  ( $j$  minimal gewählt). Dann folgt mit Taylorentwicklung:

$$\begin{aligned} x^{(k+1)} - x^* &= \Phi(x^{(k)}) - \Phi(x^*) \\ &= \sum_{\nu=1}^{p-1} \frac{\Phi^{(\nu)}(x^*)}{\nu!} (x^{(k)} - x^*)^\nu + \Phi^{(p)}(\xi_k) \frac{(x^{(k)} - x^*)^p}{p!} \end{aligned}$$

mit  $\xi_k$  zwischen  $x^{(k)}$  und  $x^*$ . Mit der Annahme folgt

$$x^{(k+1)} - x^* = \frac{\Phi^{(j)}(x^*)}{j!} (x^{(k)} - x^*)^j \cdot \underbrace{\left[ 1 + (x^{(k)} - x^*) \sum_{\nu=j+1}^p a_\nu (x^{(k)} - x^*)^{\nu-j-1} \right]}_{:=b}$$

mit  $a_\nu := \frac{j! \Phi^{(\nu)}(x^*)}{\nu! \Phi^{(j)}(x^*)}$ . Für große  $k$  gilt  $|b| \geq \frac{1}{2}$ , da  $x^{(k)} \rightarrow x^*$  und  $a_\nu$  unabhängig von  $k$  beschränkt. Wir erhalten also

$$|x^{(k+1)} - x^*| \geq \frac{1}{2} |x^{(k)} - x^*|^j \frac{|\Phi^{(j)}(x^*)|}{j!}.$$

Da  $x^{(k)} \rightarrow x^*$  mit Ordnung  $p$ , gilt

$$\infty > c \geq \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} \geq \frac{1}{2j!} |\Phi^{(j)}(x^*)| \cdot |x^{(k)} - x^*|^{j-p} \rightarrow \infty.$$

Dies ist ein Widerspruch, da nach Annahme  $j - p < 0$ .

“ $\Leftarrow$ ”: Es gilt

$$\begin{aligned} x^{(k+1)} - x^* &= \sum_{\nu=1}^{p-1} \frac{\Phi^{(\nu)}(x^*)}{\nu!} (x^{(k)} - x^*)^\nu + \frac{\Phi^{(p)}(\xi_k)}{p!} (x^{(k)} - x^*)^p \\ &= \frac{\Phi^{(p)}(\xi_k)}{p!} (x^{(k)} - x^*)^p, \text{ da } \Phi^{(\nu)}(x^*) = 0 \text{ } (1 \leq \nu \leq p) \\ \implies \frac{|x^{(k+1)} - x^*|}{|x^{(k)} - x^*|^p} &= \frac{1}{p!} \Phi^{(p)}(\xi_k) \xrightarrow{k \rightarrow \infty} \frac{1}{p!} \Phi^{(p)}(x^*), \end{aligned}$$

da  $\xi_k$  zwischen  $x^{(k)}$  und  $x^*$  und  $x^{(k)} \xrightarrow{k \rightarrow \infty} x^*$ . □

### 3.2.1 Verfahren höher Ordnung ( $p = 3$ )

**1. Idee:** Konstruiere Verfahren 3. Ordnung aus einer Linearkombination von 2 Verfahren zweiter Ordnung. Seien  $\Phi_0$ ,  $\Phi_1$  Iterationsverfahren, die quadratisch konvergieren. Betrachte

$$\Phi_s(x) = (1 - s)\Phi_0(x) + s\Phi_1(x).$$

Dann gilt  $\Phi'_s(x^*) = (1 - s)\Phi'_0(x^*) + s\Phi'_1(x^*) = 0$ , da nach Satz 3.7  $\Phi'_0(x^*) = \Phi'_1(x^*) = 0$ .

Bestimmt man  $s$  so, dass  $\Phi''_s(x^*) = 0$  gilt, so konvergiert nach Satz 3.7  $\Phi_s$  mindestens mit Ordnung 3. (Beispiel: siehe Übungsblatt).

**2. Idee: (Verbessertes Newton-Verfahren)**

Ansatz:  $\Phi(x) = x - g(x)f(x) - h(x)f(x)^2$  mit  $g(x) \neq 0$ ,  $h(x) \neq 0$  in einer Umgebung von  $x^*$ . Dann gilt  $\Phi(x^*) = x^* \iff f(x^*) = 0$ .

Idee: Bestimme  $g$ ,  $h$ , so dass  $\Phi'(x^*) = \Phi''(x^*) = 0$  für eine einfache Nullstelle  $x^*$ .

Sei  $x^*$  einfache Nullstelle, so gilt

$$\Phi'(x) = 1 - f'(x)g(x) - f(x)g'(x) - 2f(x)f'(x)h(x) - h'(x)f^2(x).$$

$$\text{Also folgt } \Phi'(x^*) = 0 \stackrel{f(x^*)=0}{\iff} 1 - f'(x^*)g(x^*) = 0 \implies g(x) = \frac{1}{f'(x)}.$$

Analog folgt für  $\Phi''$  mit dieser Wahl von  $g$ :

$$\begin{aligned} \Phi''(x) &= -f'(x)(g'(x) + h(x)f'(x) + 2h(x)f'(x)) - f(x)(\dots) - \underbrace{(g(x)f'(x))'}_{\substack{=1 \\ =0}} \\ &= \frac{f''(x)}{f'(x)} - 2h(x)f'(x) - f(x)(\dots). \end{aligned}$$

Aus

$$0 = \Phi''(x^*) = \frac{f''(x^*)}{f'(x^*)} - 2h(x^*)f'(x^*)^2$$

folgt also für die Wahl von  $h$ :

$$h(x) = \frac{1}{2f'(x)^2} \frac{f''(x)}{f'(x)} = \frac{f''(x)}{2f'(x)^3}.$$

### Folgerung 3.9 (Verbessertes Newton-Verfahren)

Das Verfahren, definiert durch

$$\Phi(x) = x - \frac{f(x)}{f'(x)} - \frac{1}{2} \frac{f(x)^2 f''(x)}{f'(x)^3}$$

konvergiert in der Umgebung einer einfachen Nullstelle  $x^*$  von  $f$  mit mindestens dritter Ordnung.

### 3.2.2 Newton-Verfahren für mehrfache Nullstellen

#### Definition 3.10 (Ordnung und Vielfachheit einer Nullstelle)

Sei  $f \in C^n(I)$  ( $n \geq 1$ ).  $f$  hat eine Nullstelle  $x^*$  der Ordnung  $(n-1)$  bzw. der Vielfachheit  $n$  gdw.  $f^{(\nu)}(x^*) = 0$ ,  $0 \leq \nu \leq n-1$  und  $f^{(n)}(x^*) \neq 0$ .

#### Satz 3.11

Sei  $f \in C^{n+1}(a,b)$ ,  $n \geq 1$  und  $x^* \in (a,b)$  eine  $n$ -fache Nullstelle von  $f$  mit  $f^{(k)}(x) \neq 0$  ( $x \neq x^*$ )  $0 \leq k \leq n-1$  und  $f^{(n)}(x) \neq 0$  für  $x \in (a,b)$ . Dann existiert eine Umgebung von  $x^*$ , so dass das Newton-Verfahren mindestens linear konvergiert.

$$\text{Beweis: Setze } \Phi(x) = \begin{cases} x - \frac{f(x)}{f'(x)} & : x \neq x^* \\ x^* & : x = x^* \end{cases}.$$

Zu zeigen:  $\Phi \in C^1(a, b)$  und  $|\Phi(x)| < 1$ . Zusammen mit Satz 3.8 folgt dann die Behauptung des Satzes. Es gilt

$$\lim_{\substack{x \rightarrow x^* \\ x \neq x^*}} \Phi(x) = x^* - \lim_{x \rightarrow x^*} \frac{f(x)}{f'(x)} \stackrel{\text{L'Hôpital}}{=} x^* - \lim_{x \rightarrow x^*} \frac{f^{(n-1)}(x)}{f^{(n)}(x)} = x^* + \frac{0}{f^{(n)}(x^*)} = x^*.$$

Weiter folgt mit L'Hopital und  $\Phi'(x) = 1 - \frac{f'(x)^2 - f(x)f''(x)}{f'(x)^2} = \frac{f(x)f''(x)}{f'(x)^2}$

$$\lim_{x \rightarrow x^*} \Phi'(x) = 1 - \frac{1}{n} \implies |\Phi'(x^*)| < 1.$$

**Idee:** Modifiziere das Verfahren je nach Vielfachheit: □

Wähle  $\Phi_\alpha(x) = x - \alpha \frac{f(x)}{f'(x)}$ ,  $\alpha \in \mathbb{R}$  fest

$$\implies \Phi'_\alpha(x^*) = 1 - \frac{\alpha}{n} = 0 \iff \alpha = n$$

### Folgerung 3.12

Sei  $f \in C^{n+1}(a, b)$ ,  $n \geq 1$  und  $x^*$  eine  $n$ -fache Nullstelle, dann konvergiert das Verfahren

$$x^{(k+1)} = x^{(k)} - n \frac{f(x^{(k)})}{f'(x^{(k)})}$$

quadratisch gegen  $x^*$ .

### 3.3 Nichtlineare Gleichungssysteme

**Problem:**  $D \subset \mathbb{R}^n$  abgeschlossen und konvex ( $n \geq 2$ ),  
 $f : D \rightarrow \mathbb{R}^n$ , d.h.  $f(x) = (f_1(x), \dots, f_n(x))^T$ ,  $x \in D$ .

**Gesucht:**  $x^* \in D$  mit  $f(x^*) = 0$ , d.h.  $f_i(x^*) = 0$  ( $i = 1, \dots, n$ ).

**Annahme:** Es existiere ein  $x^* \in D$  mit  $f(x^*) = 0$ .

#### 3.3.1 Newton-Verfahren für nichtlineare Systeme

**Idee:** Iteriere analog zum skalaren Fall mit

$$x^{(k)} = x^{(k)} - Df(x^{(k)})^{-1} f(x^{(k)}).$$

Dabei ist  $Df(x)$  die Jacobi-Matrix von  $f$  und  $Df(x)^{-1}$  die Inverse.

#### Algorithmus: (Newton-Verfahren für Systeme)

Sei  $x^{(0)} \in \mathbb{R}^n$  gegeben. Für  $k \geq 0$  iteriere

1) Löse Defektgleichung:  $Df(x^{(k)}) y^{(k)} = -f(x^{(k)})$

2) Setze neu Iterierte:  $x^{(k+1)} = x^{(k)} + y^{(k)}$

**Problem:** In jedem Schritt muss ein  $n \times n$  LGS gelöst werden. Häufig wird  $Df$  auch für einige Schritte festgehalten. Mit der LR-Zerlegung können diese Schritte dann sehr effizient gelöst werden.

#### Satz 3.13

Sei  $f_i \in C^2(\mathbb{R})$  und  $Df(x^*)$  regulär (d.h.  $x^*$  einfache Nullstelle).

Dann existiert ein  $\rho > 0$ , so dass für alle  $x^{(0)} \in B_\rho(x^*) := \{x \mid \|x - x^*\| < \rho\}$  gilt:  $x^{(k)} \in B_\rho(x^*)$

und  $x^{(k)} \xrightarrow{k \rightarrow \infty} x^*$  mindestens quadratisch.

*Beweis:* Analog zum Satz 3.2

□





# Kapitel 4

## Eigenwertprobleme

**Definition 4.1 (Eigenwertproblem)**

Sei  $A \in \mathbb{K}^{n \times n}$  mit  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$ . Dann heißt  $\lambda \in \mathbb{K}$  Eigenwert von  $A$ , wenn ein Eigenvektor  $x \in \mathbb{K}^n$  mit  $x \neq 0$  existiert mit der Eigenschaft  $Ax = \lambda x$ .

Vollständiges Eigenwertproblem: Finde alle EW (und EV) von  $A$

Partielles Eigenwertproblem: Finde einzelne EW (und EV) von  $A$  (z.B. den größten und kleinsten EW)

### 4.1 Grundbegriffe der linearen Algebra und theoretische Grundlagen

**Notation und Grundlagen 4.2**

- 1) Für  $x, y \in \mathbb{K}^n$  bezeichne  $\langle x, y \rangle$  und  $\|x\|$  das euklidische Skalarprodukt und die euklidische Norm.
- 2)  $\|A\|$  bezeichnet die Spektralnorm von  $A$ . Bemerke:  $\rho(A) \leq \|A\|_*$  für alle Normen  $\|\cdot\|_*$ .
- 3) Für  $\mathbb{K} = \mathbb{C}$  ist  $A^H := \bar{A}^T$ .  $A$  heißt hermitesch, falls  $A^H = A$ .
- 4) Die EWe von  $A$  sind die Nullstellen des charakteristischen Polynoms

$$\rho_A(\lambda) := \det(A - \lambda I).$$

- 5) Ist ein EW  $\lambda$  bekannt, so findet man alle EVen zu  $\lambda$  durch Lösung des singulären homogenen Gleichungssystems

$$(A - \lambda I)x = 0$$

Umgekehrt bestimmt ein EV  $x \neq 0$  den zugehörigen EW  $\lambda$  durch den Rayleigh Quotienten

$$R(x) = \frac{\langle Ax, x \rangle}{\|x\|^2},$$

d.h.  $\lambda = R(x)$ .

- 6)  $\sigma(A) := \{\lambda \mid \lambda \text{ EW von } A\}$  heißt Spektrum von  $A$ .

7) Das charakteristische Polynom besitzt die Darstellung

$$\rho_A(x) = \prod_{i=1}^s (x - \lambda_i)^{\sigma_i},$$

wobei  $\lambda_i$  die paarweise verschiedenen Eigenwerte von  $A$  sind. Es gilt  $\sum_{i=1}^s \sigma_i = n$  und  $\sigma_i$  heißt algebraische Vielfachheit von  $\lambda_i$ .

Die Eigenvektoren zu  $\lambda_i$  (Vereinigt mit dem Nullvektor) bilden den sogenannten Eigenraum  $E_i := \text{Kern}(A - \lambda_i I)$ . Ist  $\rho_i := \dim(E_i)$ , so heißt  $\rho_i$  geometrische Vielfachheit von  $\lambda_i$ .

### Ähnlichkeitstransformationen 4.3

Ist  $T \in \mathbb{K}^{n \times n}$  regulär, so heite  $B := T^{-1}AT$  Ähnlichkeitstransformation und  $B$  ähnlich zu  $A$ .

1) Ähnliche Matrizen besitzen dieselben EWe  $\lambda$ , denn mit  $y := T^{-1}x$  folgt:

$$T^{-1}ATy = T^{-1}Ax = \lambda T^{-1}x = \lambda y.$$

D.h.  $\lambda$  ist EW zu  $B$  mit EV  $y = T^{-1}x$ .

2) Ähnliche Matrizen besitzen dasselbe char. Polynom, denn

$$\begin{aligned} \det(T^{-1}AT - \lambda I) &= \det(T^{-1}(A - \lambda I)T) \\ &= \det(T^{-1}) \det(A - \lambda I) \det(T) \\ &= \det(A - \lambda I). \end{aligned}$$

### Idee einer numerischen Methode 4.4

Wende eine Folge von Ähnlichkeitstransformationen an, um  $A$  in einfachen Gestalt zu transformieren, d.h.

$$\begin{aligned} A^{(0)} &:= A, \\ A^{(i)} &:= T_i^{-1}A^{(i-1)}T_i, \quad i = 1, 2, 3, \dots \end{aligned}$$

und geeignete Matrizen  $T_i$ .

### Satz 4.5 (Satz von Schur)

Zu jeder Matrix  $A \in \mathbb{C}^{n \times n}$  existiert eine unitäre Matrix  $U \in \mathbb{C}^{n \times n}$  (d.h.  $U^H U = I$ ) mit

$$U^H A U = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}.$$

*Beweis:* (Induktion über  $n$ )

Ind. Anf:  $n = 1$  ✓

Ind. Vor: Für  $A \in \mathbb{C}^{(n-1) \times (n-1)}$  gilt die Behauptung.

Ind. Beh: Sei  $A \in \mathbb{C}^{n \times n}$ . Zeige:  $\exists U$  unitär mit  $U^H A U = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$ .

Sei  $\lambda \in \mathbb{C}$  EW von  $A$  und  $z \in \mathbb{C}^n \setminus \{0\}$  ein EV mit  $\|z\| = 1$ .

Wir können  $z$  zu einer Orthonormalbasis auf  $\mathbb{C}^n$  ergänzen, d.h. wir können eine unitäre Matrix  $\hat{V}$  finden, so dass

$$V = (z, \hat{V}) \text{ mit } V^H V = I.$$

Sei  $B = V^H A V$ . Wegen  $V B e_1 = A V e_1 = A z = \lambda z = \lambda V e_1$

$\implies B e_1 = \lambda e_1$ , d.h. die erste Spalte von  $B$  ist ein  $\lambda$ -faches des ersten Einheitsvektors.

$$\text{Also } B = V^H A V = \left( \begin{array}{c|c} \lambda & b \\ \hline 0 & C \end{array} \right).$$

Nach Ind. Vor. ex. eine unitäre Matrix  $W$  mit  $W^H C W = T$  und  $T$  hat obere Dreiecksgestalt.

$$\implies A = V \left( \begin{array}{c|c} \lambda & b \\ \hline 0 & W^H C W \end{array} \right) V^H = \underbrace{V \left( \begin{array}{c|c} 1 & 0 \\ \hline 0 & W \end{array} \right)}_{=:U} \left( \begin{array}{c|c} \lambda & b \\ \hline 0 & T \end{array} \right) \underbrace{\left( \begin{array}{c|c} 1 & 0 \\ \hline 0 & W^H \end{array} \right) V^H}_{:=U^H}.$$

Daraus folgt die Behauptung. □

#### Folgerung 4.6 (Schur für hermitesche Matrix)

Sei  $A \in \mathbb{C}^{n \times n}$  hermitesch. Dann ex. eine unitäre Matrix  $U = (u_1, \dots, u_n)$  mit

$$U^H A U = \text{diag}(\lambda_1, \dots, \lambda_n).$$

$\lambda_1, \dots, \lambda_n$  sind EW von  $A$ , die reell sind und  $u_i$  die EVen zu  $\lambda_i$ .

Insbesondere ist  $A$  eine Matrix mit  $n$  linear unabhängigen zu einander orthogonalen Eigenvektoren.

*Beweis:*  $A$  hermitesch  $\implies A^H = A$  und somit

$$(U^H A U)^H = U^H A (U^H)^H = U^H A U.$$

$\implies U^H A U$  ist selbst wieder hermitesch.

$\implies$  Beh. mit 4.5. □

## 4.2 Kondition des Eigenwertproblems

### Satz 4.7 (Gerschgorinscher Kreissatz)

Sei  $A \in \mathbb{C}^{n \times n}$ . Für  $i = 1, \dots, n$  definiere die sogenannten Gerschgorin Kreise

$$G_i := \{z \in \mathbb{C} \mid |z - a_{ii}| \leq r_i\}, \quad r_i := \sum_{j=1, j \neq i}^n |a_{ij}|.$$

Dann gilt:

- (i) Ist  $\lambda$  EW von  $A$ , so ist  $\lambda \in \bigcup_{i=1}^n G_i$ , d.h.  $\sigma(A) \subset \bigcup_{i=1}^n G_i$ .
- (ii) Hat die Vereinigung  $\hat{G}$  von  $m$  Kreisen  $G_i$  einen leeren Schnitt mit den restlichen  $n-m$  Kreisen, so enthält  $\hat{G}$  genau  $m$  EWe von  $A$  (gezählt mit ihren algebraischen Vielfachheiten).

*Beweis:*

zu (i): Sei  $\lambda$  EW von  $A$  mit EV  $x \neq 0$ .

Aus  $Ax = \lambda x$  folgt:

$$(\lambda - a_{ii})x_i = \sum_{j=1, j \neq i}^n a_{ij}x_j \quad \forall i = 1, \dots, n.$$

Für  $i \in \{1, \dots, n\}$  mit  $|x_i| = \max_{i=1, \dots, n} |x_j|$  folgt:

$$|\lambda - a_{ii}| = \left| \sum_{j=1, j \neq i}^n \frac{a_{ij}x_j}{x_i} \right| \leq \sum_{j=1, j \neq i}^n |a_{ij}|$$

$$\implies \lambda \in G_i \subset \bigcup_{j=1}^n G_j.$$

zu (ii): Wir setzen  $D = (a_{ii}\delta_{ij})_{i,j=1, \dots, m}$  und betrachten

$$B(t) := D + t(A - D), \quad 0 \leq t \leq 1$$

mit Gerschgorin Kreisen

$$G_i(t) := \left\{ z \in \mathbb{C} \mid |z - a_{ii}| \leq t \cdot \sum_{j=1, j \neq i}^n |a_{ij}| \right\} \quad i = 1, \dots, n.$$

Es ist  $B(0) = D$  und  $B(1) = A$ . Die EW von  $B(t)$  sind die Nullstellen von  $\rho_{B(t)}$  und hängen daher stetig von  $t$  ab. Wende (i) auf  $B(t)$  an und lasse  $t$  von 0 nach 1 laufen.

Dabei wird der Radius der Kreise bei festem Mittelpunkt immer größer.

Die Anzahl der EWe in einem Kreis  $G_i(t)$  kann sich erst dann ändern, wenn dieser einen anderen Kreis trifft. Daraus folgt die Beh.

□

### Folgerung 4.8

Da  $A$  und  $A^H$  dieselben Eigenwerte besitzen, gilt der Satz 4.7 auch mit

$$G'_i := \{z \in \mathbb{C} \mid |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ji}|\}$$

$$\implies \sigma(A) \subset \left( \bigcup_{i=1}^n G_i \right) \cap \left( \bigcup_{i=1}^n G'_i \right).$$

#### Lemma 4.9

Seien  $A, B \in \mathbb{C}^{n \times n}$ . Ist  $\lambda$  ein Eigenwert von  $A$ , aber kein Eigenwert von  $B$ , so gilt

$$1 \leq \|(\lambda I - B)^{-1}(A - B)\|.$$

*Beweis:* Sei  $Ax = \lambda x$ ,  $x \neq 0 \implies \lambda x - Bx = (A - B)x$   
 $\implies x = (\lambda I - B)^{-1}(A - B)x.$

Also ist 1 ein EW von  $(\lambda I - B)^{-1}(A - B)$ , d.h.

$$1 \leq \rho((\lambda I - B)^{-1}(A - B)) \leq \|(\lambda I - B)^{-1}(A - B)\|.$$

□

#### Folgerung 4.10

Lemma 4.9 gilt insbesondere für  $B = D = (a_{ij}\delta_{ij})_{i,j}$ . Da die Spektralnorm kleiner als z.B. die Zeilensummennorm ist, folgt dann

$$1 \leq \max_{i=1, \dots, n} \sum_{k \neq j} \left| \frac{a_{jk}}{\lambda - a_{ii}} \right|$$

$$\implies \exists j \text{ mit } |\lambda - a_{jj}| = \sum_{k \neq j} |a_{jk}| \implies \lambda \in G_j \subset \bigcup_{i=1}^n G_i$$

#### Satz 4.11 (Kondition des Eigenwertproblems)

Sei  $B \in \mathbb{C}^{n \times N}$  diagonalisierbar, d.h. es ex eine reguläre Matrix  $P \in \mathbb{C}^{n \times n}$  mit

$$P^{-1}BP = D = \text{diag}(\lambda_1(B), \dots, \lambda_n(B)).$$

Dann gibt es zu jedem EW  $\lambda_j(B)$  einen EW  $\lambda_j(A)$  von  $A \in \mathbb{C}^{n \times n}$  mit

$$|\lambda_j(A) - \lambda_j(B)| \leq \kappa(P) \|A - B\|.$$

Dabei ist  $\kappa(P) = \frac{\lambda_{\max}(P)}{\lambda_{\min}(P)}$  die Kondition von  $P$ .

*Beweis:*

$$\begin{aligned} \|(\lambda I - B)^{-1}\| &= \|P(\lambda I - D)^{-1}P^{-1}\| \leq \|(\lambda I - D)^{-1}\| \kappa(P) \\ &\leq \max_j \frac{1}{|\lambda - \lambda_j(B)|} \kappa(P) \\ &= \frac{1}{\min_j |\lambda - \lambda_j(B)|} \kappa(P). \end{aligned}$$

Aus Lemma 4.7 folgt:

$$1 \leq \|(\lambda I - B)^{-1}\| \|A - B\|.$$

Also folgt:

$$1 \leq \frac{1}{\min_j |\lambda - \lambda_j(B)|} \kappa(P) \|A - B\|$$

$$\implies \min_j |\lambda - \lambda_j(B)| \leq \kappa(P) \|A - B\|$$

$$\implies \exists \lambda = \lambda_j(A) \text{ mit } |\lambda_j(A) - \lambda_j(B)| \leq \kappa(P) \|A - B\|.$$

□

Für unitäre Matrizen  $U$  gilt  $\kappa(U) = 1$ . Also ist nach Folgerung 4.6 das Eigenwertproblem für hermitesche Matrizen gut konitioniert.

#### Beispiel 4.12 (Anwendung der Gerschgorin-Kreise)

$$A = \begin{pmatrix} 0,9 & 0 & 0 \\ 0 & 0,4 & 0 \\ 0 & 0 & 0,2 \end{pmatrix} + 10^{-5} \begin{pmatrix} 0,1 & 0,4 & -0,2 \\ -0,1 & 0,5 & 0,1 \\ 0,2 & 0,1 & 0,3 \end{pmatrix}$$

Nach Satz 4.11 erwartet man, dass die EW einer Matrix  $A$  mit kleinen Außendiagonalelementen ungefähr mit dem Diagonalelementen übereinstimmen.

Die 3 Gerschgorinkreise sind disjunkt, daher besitzt  $A$  die EW  $\lambda_1, \lambda_2, \lambda_3$  mit

$$\begin{aligned} |\lambda_1 - (0,9 + 0,1 \cdot 10^{-5})| &\leq 0,6 \cdot 10^{-5}, \\ |\lambda_2 - (0,4 + 0,5 \cdot 10^{-5})| &\leq 0,2 \cdot 10^{-5}, \\ |\lambda_3 - (0,2 + 0,3 \cdot 10^{-5})| &\leq 0,3 \cdot 10^{-5}. \end{aligned}$$

Diese Abschätzungen können noch wesentlich verbessert werden.

Sei  $P := \text{diag}(10^5, 1, 1)$ . Dann ist

$$P^{-1}AP = \begin{pmatrix} 0,9 & 0 & 0 \\ 0 & 0,4 & 0 \\ 0 & 0 & 0,2 \end{pmatrix} + \begin{pmatrix} 0,1 \cdot 10^{-5} & 0,1 \cdot 10^{-10} & -0,2 \cdot 10^{-10} \\ -0,1 & 0,5 \cdot 10^{-5} & 0,1 \cdot 10^{-5} \\ 0,2 & 0,1 \cdot 10^{-5} & 0,3 \cdot 10^{-5} \end{pmatrix}.$$

Der erste Gerschgorin Kreis ist noch disjunkt zu den beiden anderen, die nicht mehr disjunkt sind.

Also Folgt für  $\lambda_1$ :

$$|\lambda_1 - (0,9 + 0,1 \cdot 10^{-5})| \leq 0,6 \cdot 10^{-10}.$$

Entsprechend kann die Abschätzung für  $\lambda_2, \lambda_3$  verbessert werden (siehe Abbildung 4.1).

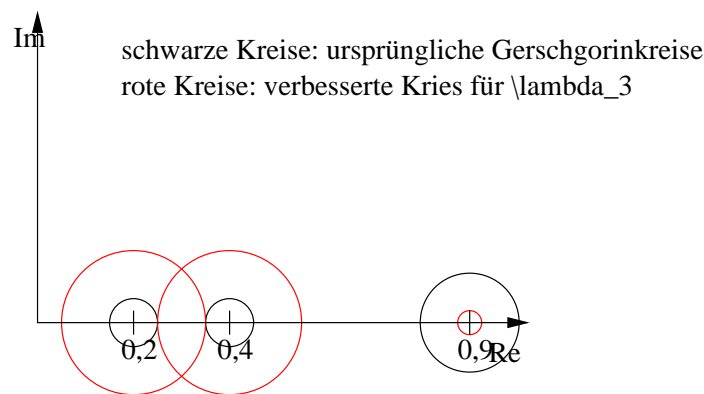


Abbildung 4.1: Gerschgorinkreise: Beispiel, Radien nicht maßstabsgetreu!!

### 4.3 Variationsprinzip für Eigenwerte hermitescher Matrizen

#### Satz 4.13 (Rayleighsches Maximumsprinzip)

Sei  $A \in \mathbb{K}^{n \times n}$  hermitesch. Die EW von  $A$  seien  $\lambda_1 \geq \dots \geq \lambda_n$ . Sei  $U = (u_1, \dots, u_n)$  unitäre Matrix mit  $U^H A U = \text{diag}(\lambda_1, \dots, \lambda_n) =: \Lambda$ .

Für  $j = 1, \dots, n$  definiere den  $(n+1-j)$ -dimensionalen Teilraum

$$M_j := \{x \in \mathbb{K}^n \mid \langle u_i, x \rangle = 0 \forall i = 1, \dots, j-1\} = (\text{span}(u_1, \dots, u_{j-1}))^\perp.$$

Dann gilt:

$$\lambda_j = \max_{x \in M_j \setminus \{0\}} R(x)$$

mit dem Rayleigh-Quotienten  $R(x) := \frac{\langle Ax, x \rangle}{\langle x, x \rangle}$ .

*Beweis:* Sei  $j \in \{1, \dots, n\}$  und  $x \in M_j \setminus \{0\}$  beliebig.

setze  $y := U^H x$ . Dann folgt:

$$y_i = \langle u_i, x \rangle = 0, \quad i = 1, \dots, j-1.$$

Wegen  $\lambda_i \leq \lambda_j$  für  $i = j, \dots, n$  gilt:

$$R(x) = \frac{\langle Ax, x \rangle}{\langle x, x \rangle} = \frac{\langle U \Lambda U^H x, x \rangle}{\langle x, x \rangle} = \frac{\langle \Lambda U^H x, U^H x \rangle}{\langle x, x \rangle} = \frac{\langle \Lambda y, y \rangle}{\langle y, y \rangle} = \frac{\sum_{i=1}^n \lambda_i |y_i|^2}{\sum_{i=1}^n |y_i|^2} \leq \lambda_j$$

und somit

$$\sup_{x \in M_j \setminus \{0\}} R(x) \leq \lambda_j.$$

Andererseits ist  $u_j \in M_j \setminus \{0\}$  und  $R(u_j) = \lambda_j$ .

$\Rightarrow \max_{x \in M_j \setminus \{0\}} R(x) = \lambda_j$ .

□

#### Bemerkung 4.14

Definiert man  $f: \mathbb{R} \rightarrow \mathbb{R}$  bei festem  $x \in \mathbb{K}^n \setminus \{0\}$  durch

$$f(\lambda) = \frac{1}{2} \|Ax - \lambda x\|^2 = \frac{1}{2} \lambda^2 \|x\|^2 - \lambda \langle Ax, x \rangle + \frac{1}{2} \|Ax\|^2$$



so nimmt  $f$  in  $\lambda = R(x)$  sein Minimum an.

Ist daher  $x$  näherungsweise ein EV von  $A$ , so ist  $R(x)$  eine gute Näherung des zugehörigen Eigenwertes.

**Satz 4.15 (Courantsches Minimum-Maximum Prinzip)**

Sei  $A \in \mathbb{K}^{n \times n}$  hermitesch mit EWen  $\lambda_1 \geq \dots \geq \lambda_n$ . Für  $j = 1, \dots, n$  definiere

$$\mathcal{N}_j := \{N_j \subset \mathbb{K}^n \mid N_j \text{ ist linearer Teilraum der Dimension } n+1-j\}$$

Dann gilt:

$$\lambda_j = \min_{N_j \in \mathcal{N}_j} \max_{x \in N_j \setminus \{0\}} R(x), \quad j = 1, \dots, n.$$

*Beweis:* Sei  $U = (u_1, \dots, u_n)$  die unitäre Matrix von EVen zu den EWen  $\lambda_1, \dots, \lambda_n$  von  $A$ . Sei  $j \in \{1, \dots, n\}$ . Definiere  $L_j = \text{span}(u_1, \dots, u_j)$  und wähle  $N_j \in \mathcal{N}_j$  beliebig.

Wegen

$$\begin{aligned} \dim(L_j \cap N_j) &= \dim(L_j) + \dim(N_j) - \dim(L_j \cup N_j) \\ &= n+1 - \underbrace{\dim(L_j \cup N_j)}_{\leq n} \geq 1 \end{aligned}$$

existiert ein  $x \in L_j \cap N_j$  mit  $x \neq 0$ .

Da  $x \in L_j$ , folgt:  $x = \sum_{i=1}^j \alpha_i u_i$ .

$$\implies R(x) = \frac{\sum_{i=1}^j \lambda_i |\alpha_i|^2}{\sum_{i=1}^j |\alpha_i|^2} \geq \lambda_j, \text{ da } \lambda_i \geq \lambda_j \text{ für } i = 1, \dots, j$$

$$\implies \min_{N_j \in \mathcal{N}_j} \max_{x \in N_j \setminus \{0\}} R(x) \geq \lambda_j$$

Wählt man andererseits  $N_j = M_j$ , so gilt nach Satz 4.15

$$\max_{x \in M_j \setminus \{0\}} R(x) = \lambda_j.$$

□

**Folgerung 4.16**

Seien  $A, B \in \mathbb{K}^{n \times n}$  hermitesch und gelte  $\lambda_1(A) \geq \dots \geq \lambda_n(A)$  sowie  $\lambda_1(B) \geq \dots \geq \lambda_n(B)$ .

Dann gilt die Abschätzung:

$$|\lambda_j(A) - \lambda_j(B)| \leq \|A - B\|, \quad j = 1, \dots, n.$$

*Beweis:*  $A, B$  hermitesch  $\implies E = A - B$  hermitesch. Sei  $x \in \mathbb{K}^n \setminus \{0\}$ .

Dann gilt:

$$R_E(x) = \frac{\langle (A - B)x, x \rangle}{\langle x, x \rangle} \leq \|A - B\|.$$

Somit folgt

$$(*) \quad R_A(x) \leq R_B(x) + \|A - B\|.$$

Sei  $j \in \{1, \dots, n\}$  und  $N_j \in \mathcal{N}_j$  (vgl. 4.15), so folgt aus (\*):

$$\min_{N_j \in \mathcal{N}_j} \max_{x \in N_j \setminus \{0\}} R_A(x) \leq \min_{N_j \in \mathcal{N}_j} \max_{x \in N_j \setminus \{0\}} R_B(x) + \|A - B\|.$$

Mit dem Courantschen Min-Max Prinzip folgt  $\lambda_j(A) \leq \lambda_j(B) + \|A - B\|$ .

Vertauscht man die Rollen von  $A$  und  $B$ , so folgt auch

$$\begin{aligned} \lambda_j(B) &\leq \lambda_j(A) + \|A - B\| \\ \implies |\lambda_j(B) - \lambda_j(A)| &\leq \|A - B\|. \end{aligned}$$

□

## 4.4 Transformation auf Hessenberg-Form

### Definition 4.17 (Hessenberg-Matrizen)

Eine Matrix  $A = (a_{ij})_{i,j} \in \mathbb{R}^{n \times n}$  heißt eine (obere) Hessenberg-Matrix, wenn  $a_{ij} = 0$  für  $1 \leq j \leq i - 2$ , d.h.

$$A = \begin{pmatrix} * & \dots & \dots & * \\ * & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ 0 & & * & * \end{pmatrix}$$

Eine Hessenberg-Matrix heißt unreduziert, falls sämtliche Subdiagonalelemente  $a_{i+1,i}$ ,  $i = 1, \dots, n - 1$  ungleich 0 sind.

Eine symmetrische Hessenberg-Matrix ist somit eine symmetrische Tridiagonalmatrix.

### Erinnerung an die Numerik I 4.18 (Householder-Matrix)

Ist  $x \in \mathbb{R}^n \setminus \{0\}$  und definiert man  $u := x + \text{sign}(x_1) \|x\| e_1 \in \mathbb{R}^n$ , so ist durch  $P = I - \frac{2uu^T}{\langle u, u \rangle} = I - \beta uu^T$  mit  $\beta := \frac{2}{\langle u, u \rangle} = \frac{1}{\|x\|(\|x\| + |x_1|)}$  eine Householder Matrix definiert mit  $Px = -\text{sgn}(x_1) \|x\| e_1$ .

### Satz 4.19 (Householder-Transformationen auf Hessenberg-Form)

Zu einer Matrix  $A \in \mathbb{R}^{n \times n}$  existieren  $n - 2$  Householder-Matrizen  $P_1, \dots, P_{n-2}$ , so dass

$$P_{n-2} \dots P_1 A P_1 \dots P_{n-2}$$

eine zu  $A$  orthogonal-ähnliche Hessenberg-Matrix ist.

*Beweis:* (Induktion über  $k = 1, \dots, n - 2$ )

Angenommen es seien schon  $k - 1$  Householder-Matrizen  $P_1, \dots, P_{k-1}$  bestimmt, so dass

$$P_{k-1} \dots P_1 A P_1 \dots P_{k-1} = \left( \begin{array}{c|c} H_k & B_k \\ \hline 0 & a_k \mid C_k \end{array} \right) \left\{ \begin{array}{l} k \\ n - k \end{array} \right\},$$

wobei  $H \in \mathbb{R}^{k \times k}$  eine Hessenberg-Matrix,  $B_k \in \mathbb{R}^{k \times k}$ ,  $C \in \mathbb{R}^{(n-k) \times (n-k)}$  und  $a_k \in \mathbb{R}^{n-k}$  ist.

Für  $P_k$  machen wir folgenden Ansatz:

$$P_k = \text{diag}(I_k, \tilde{P}_k) = \left( \begin{array}{c|c} I_k & 0 \\ \hline 0 & \tilde{P}_k \end{array} \right)$$

mit der Identität  $I_k \in \mathbb{R}^{k \times k}$  und einer  $(n - k) \times (n - k)$  Householder-Matrix  $\tilde{P}_k$ . Dann folgt:

$$\begin{aligned} P_k P_{k-1} \dots P_1 A P_1 \dots P_{k-1} P_k &= \left( \begin{array}{c|c} I_k & 0 \\ \hline 0 & \tilde{P}_k \end{array} \right) \left( \begin{array}{c|c} H_k & B_k \\ \hline 0 & a_k \mid C_k \end{array} \right) \left( \begin{array}{c|c} I_k & 0 \\ \hline 0 & \tilde{P}_k \end{array} \right) \\ &= \left( \begin{array}{c|c} H_k & B_k \tilde{P}_k \\ \hline 0 & \tilde{P}_k a_k \mid \tilde{P}_k C_k \tilde{P}_k \end{array} \right) = \left( \begin{array}{c|c} H_{k+1} & B_{k+1} \\ \hline 0 & a_{k+1} \mid C_{k+1} \end{array} \right) \left\{ \begin{array}{l} k \\ n - k \end{array} \right\} \end{aligned}$$

mit einer Hessenberg-Matrix  $H_{k+1} \in \mathbb{R}^{(k+1) \times (k+1)}$ , wenn die Householder-Matrix  $\tilde{P}_k$  so gewählt wird, dass  $\tilde{P}_k a_k$  ein Vielfaches des ersten Einheitsvektors in  $\mathbb{R}^{n-k}$  ist (siehe 4.18).  $\square$

**Algorithmus 4.20 (Householder-Transformation auf Hessenberg-Form)**

**Input:**  $A \in \mathbb{R}^{n \times n}$

**Für**  $k = 1, \dots, n-2$

**Falls**  $a_k = (a_{k+1,k}, \dots, a_{n,k})^T \neq 0$

**dann** Berechne Householder Matrix  $\tilde{P}_k$  durch

$$u_k := (a_{k+1,k} + \text{sign}(a_{k+1,k} \|a_k\|), a_{k+2,k}, \dots, a_{n,k})^T$$

$$\beta_k := \|a_k\|^{-1} (\|a_k\| + |a_{k+1,k}|)^{-1}$$

$$\tilde{P}_k := I_{n-k} - \beta_k u_k (u_k)^T \text{ und berechne } A := P_k A P_k$$

**sonst** setze  $P_k := I$

**Output:** Die Ausgangsmatrix  $A$  wird in  $n-2$  Schritten mit orthogonal-ähnlichen Transformationen in eine Hessenberg-Matrix  $P^T A P$ ,  $P = P_1 \cdot \dots \cdot P_{n-2}$  überführt.

**Bemerkung:**

Da orthogonale Ähnlichkeitstransformationen die Symmetrie erhalten, wird eine symmetrische Matrix  $A$  durch  $n-2$  Schritte auf eine Tridiagonalmatrix transformiert.

## 4.5 Eigenwertbestimmung für Hessenberg-Matrizen

### Grundlegende Idee 4.21

Bestimme das charakteristische Polynom  $\rho_A(\lambda)$  einer Hessenberg-Matrix  $A$ , sowie die Ableitung  $\rho'_A(\lambda)$ , so dass man die Eigenwerte durch Nullstellensuchen, etwa mit dem Newtonverfahren, berechnen kann.

### Satz 4.22 (Berechnung von $\rho_A(\lambda), \rho'_A(\lambda)$ für Tridiagonalmatrizen)

Sei  $T \in \mathbb{R}^{n \times n}$  eine symmetrische Tridiagonalmatrix, d.h.

$$T = \begin{pmatrix} b_1 & c_1 & & & 0 \\ c_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & c_{n-1} \\ 0 & & & c_{n-1} & b_n \end{pmatrix}$$

und  $\rho_A(\lambda) = \det(T - \lambda I)$ .

Für  $k = n, \dots, 0$  seien  $f_k(\lambda), g_k(\lambda)$  definiert durch:

$$\begin{aligned} f_n(\lambda) &= 1, & g_n(\lambda) &= 0, \\ f_{n-1}(\lambda) &= b_1 - \lambda, & g_{n-1}(\lambda) &= -1, \\ f_{n-i-1}(\lambda) &= (b_{i+1} - \lambda)f_{n-i}(\lambda) - |c_i|^2 f_{n-i+1}(\lambda), \\ g_{n-i-1}(\lambda) &= -f_{n-i}(\lambda) + (b_{i+1} - \lambda)g_{n-i}(\lambda) - |c_i|^2 g_{n-i+1}(\lambda). \end{aligned}$$

Dann gilt:

$$\begin{aligned} f_0(\lambda) &= \rho_A(\lambda), \\ g_0(\lambda) &= \rho'_A(\lambda). \end{aligned}$$

*Beweis:* Wir zeigen durch vollst. Induktion über  $i$ , dass  $f_{n-i}(\lambda)$  die Determinante des  $i$ -ten Hauptminors von  $T - \lambda I$  ist:

Für  $i = 0, 1$  ist die Aussage richtig.

Sei die Aussage richtig für alle  $j$  mit  $1 \leq j \leq i, 2 \leq i$ .

Dann folgt mit Determinanten-Entwicklungssatz

$$\begin{aligned} \det \begin{pmatrix} b_1 - \lambda & c_1 & & & 0 \\ c_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & c_{i-1} \\ \hline & & & c_{i-1} & b_i - \lambda \\ \hline 0 & & & c_i & b_{i+1} - \lambda \end{pmatrix} &= \\ &= (b_{i+1} - \lambda)f_{n-i}(\lambda) - c_i c_i f_{n-(i-1)}(\lambda) \\ &= f_{n-i-1}(\lambda). \end{aligned}$$

Da  $f'_{n-i}(\lambda) = g_{n-i}(\lambda)$  für  $i = 0, \dots, n$ , folgt die Behauptung. □

### Bemerkung 4.23

Das Verfahren lässt sich einfach erweitern auf allgemeine Tridiagonalmatrizen  $T = \text{tridiag}(a, b, c)$ , falls im Algorithmus  $|c_i|^2$  durch  $a_i c_i$  ersetzt wird.

**Beispiel 4.24**

Sei  $T$  gegeben durch

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

Dann ist

$$f_3(\lambda) = 1, f_2(\lambda) = 1 - \lambda, f_1(\lambda) = \lambda^2 - 4\lambda - 1, f_0(\lambda) = -\lambda^3 + 9\lambda^2 - 3\lambda - 21.$$

Aus  $f_0(\lambda) = 0$  ergeben sich die Eigenwerte.

$$\text{Weiter ist } g_3(\lambda) = 0, g_2(\lambda) = -1, g_1(\lambda) = 2\lambda - 4, g_0(\lambda) = -3\lambda^2 + 18\lambda - 3.$$

**Bemerkung 4.25**

Gilt  $c_k \neq 0$  für  $1 \leq k \leq n-1$ , so haben die Polynome  $f_{n-i}$   $i$  reelle einfache Nullstellen  $\lambda_j^{(i)}, j = 1, \dots, i$  und die Nullstellen von  $f_{n-i}$  trennen die Nullstellen von  $f_{n-i-1}$ . Daher bilden die Polynome  $(f_n, \dots, f_0)$  eine Sturmsche Kette.

Dadurch gilt für ein beliebiges Intervall  $[a, b]$  mit  $f_0(a)f_0(b) \neq 0$ :

Ist  $n_a$  die Anzahl von Vorzeichenwechsel der Sequenz

$$(f_n(a), \dots, f_0(a))$$

und  $n_b$  die Anzahl von Vorzeichenwechsel der Sequenz

$$(f_n(b), \dots, f_0(b)),$$

so besitzt  $f_0$  auf  $[a, b]$  genau  $n_a - n_b$  Nullstellen. Mit Intervallhalbierung kann man dann alle Nullstellen in  $[a, b]$  finden.

**Motivaton 4.26 (Bestimmung von  $\rho_A(\lambda)$  für unreduzierte Hessenberg-Matrizen)**

Sei  $H = (h_{ij})_{i,j} \in \mathbb{K}^{n \times n}$  unreduzierte Hessenberg-Matrix, d.h.  $h_{i+1,i} \neq 0 \forall i = 1, \dots, n-1$ .

Vorschlag von Hyman (1957)

Betrachte das lineare Gleichungssystem

$$(H - \lambda I)x = -c e_1, c \in \mathbb{K}.$$

Setzte man  $x_n = 1$ , so kann man durch Rückwärtseinsetzen nacheinander  $x_{n-1}, x_{n-2}, \dots, x_1$  berechnen und schließlich  $c$  bestimmen.

Andererseits kann  $x_n$  mit der Cramerschen Regel berechnet werden durch:

$$1 = x_n = \frac{(-1)^n c h_{2,1} h_{3,2} \dots h_{n,n-1}}{\det(H - \lambda I)}.$$

Also folgt

$$\rho_H(\lambda) = \det(H - \lambda I) = (-1)^n c h_{2,1} h_{3,2} \dots h_{n,n-1}.$$

Mit diesem Vorgehen erhalten wir folgenden Algorithmus.

**Satz 4.27 (Verfahren nach Hyman)**

Sei  $H \in \mathbb{K}^{n \times n}$  unreduzierte Hessenberg-Matrix. Für  $H$  mit charakteristischem Polynom

$$\rho_H(\lambda) = (-1)^n h_{2,1} h_{3,2} \dots h_{n,n-1} \varphi(\lambda)$$

liefert der Algorithmus

Input:  $\lambda \in \mathbb{C}$

$$h_{1,0} = 1, x_n = 1, y_n = 0$$

$$x_{n-i} = \frac{1}{h_{n-i+1,n-i}} \left( \lambda x_{n-i+1} - \sum_{j=n-i+1}^n h_{n-i+1,j} x_j \right)$$

$$y_{n-i} = \frac{1}{h_{n-i+1,n-i}} \left( x_{n-i+1} + \lambda y_{n-i+1} - \sum_{j=n-i+1}^n h_{n-i+1,j} y_j \right)$$

für  $1 \leq i \leq n$  das Ergebnis

$$\begin{aligned} x_0 &= \varphi(\lambda) \quad (= c), \\ y_0 &= \varphi'(\lambda). \end{aligned}$$

Ist  $\lambda$  ein Eigenwert von  $H$ , so ist  $x = (x_1, \dots, x_n)^T$  ein zugehöriger Eigenvektor.

*Beweis:* Folgt aus 4.26 und  $y_{n-i} = x'_{n-i}$ . □

## 4.6 Vektoriteration für partielle Eigenwertprobleme

**Definition 4.28 (Vektoriteration nach von Mises)**

Sei  $A \in \mathbb{C}^{n \times n}$  eine diagonalisierbare Matrix mit dominanten Eigenwert  $\lambda_1$ , d.h.

$$(D - EW) \quad |\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|.$$

Dann erhält man eine Folge von Approximationen  $\lambda_k, k = 1, 2, \dots$  von  $\lambda_1$  durch folgenden Algorithmus:

Input:  $z^0 \in \mathbb{C}^n$  mit  $\|z^0\| = 1$  und  $l \in \{1, \dots, n\}$

Für  $k = 1, 2, \dots$  berechne

$$\begin{aligned} \tilde{z}^k &= Az^{k-1} \\ z^k &= \frac{1}{\|\tilde{z}^k\|} \tilde{z}^k \\ \lambda^k &= \frac{(Az^k)_l}{z_l^k} \end{aligned}$$

**Satz 4.29 (Konvergenz der Vektoriteration nach von Mises)**

Sei  $A \in \mathbb{C}^{n \times n}$  diagonalisierbar mit einer Basis  $\{u_1, \dots, u_n\}$  von normierten Eigenvektoren zu den Eigenwerten  $\lambda_1, \dots, \lambda_n$ . Sei  $\lambda_1$  ein dominanter Eigenwert von  $A$ , d.h. es gelte  $(D - EW)$ . Der Startwert  $z^0 \in \mathbb{C}$  habe eine nichttriviale Komponente in Richtung  $u_1$ , d.h.

$$z^0 = \sum_{i=1}^n \alpha_i u_i \quad \text{mit } \alpha_1 \neq 0.$$

Dann gilt für  $z^k, \lambda^k$  aus der Vektoriteration nach von Mises (Def. 4.28), falls  $u_{1,l} \neq 0$  ist:

- 1)  $\|z^k - \sigma_k u_1\| = \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \quad (k \rightarrow \infty)$  mit  $\sigma_k := \frac{\lambda_1^k \alpha_1}{|\lambda_1^k \alpha_1|}$ , also  $|\sigma_k| = 1$ .
- 2)  $\lambda^k - \lambda_1 = \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \quad (k \rightarrow \infty)$ .

*Beweis:* Für die Iterierten  $z^k$  zeigt man durch vollständige Induktion, dass

$$z^k = \frac{A^k z^0}{\|A^k z^0\|}.$$

Mit  $z^0 = \sum_{i=1}^n \alpha_i u_i$  folgt weiter

$$A^k z^0 = \sum_{i=1}^n \alpha_i \lambda_i^k u_i = \lambda_1^k \alpha_1 \left( u_1 + \sum_{i=2}^n \frac{\alpha_i}{\alpha_1} \left( \frac{\lambda_i}{\lambda_1} \right)^k u_i \right).$$

Wegen  $(D - EW)$  folgt dann

$$A^k z^0 = \lambda_1^k \alpha_1 \left( u_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \right) \quad (k \rightarrow \infty).$$



und hieraus

$$z^k = \frac{\lambda_1^k \alpha_1 \left( u_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \right)}{|\lambda_1^k \alpha_1| \left\| u_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \right\|} = \sigma_k u_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right).$$

Also ist 1) gezeigt. Weiter gilt

$$\begin{aligned} \lambda^k &= \frac{(Az^k)_l}{z_l^k} \\ &= \frac{(A^{k+1}z^0)_l \|A^k z^0\|}{\|A^k z^0\| (A^k z^0)_l} \\ &= \frac{\lambda_1^{k+1} \left( \alpha_1 u_{1,l} + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1}\right)^{k+1} u_{i,l} \right)}{\lambda_1^k \left( \alpha_1 u_{1,l} + \sum_{i=2}^n \alpha_i \left(\frac{\lambda_i}{\lambda_1}\right)^k u_{i,l} \right)} \\ &\stackrel{u_{1,l} \neq 0}{=} \lambda_1 + \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right) \quad (k \rightarrow \infty). \end{aligned}$$

□

### Bemerkung 4.30

1) Die Konvergenz der Vektoriteration nach von Mises ist also umso besser, je weiter der dominante Eigenwert  $\lambda_1$  von den anderen Eigenwerten entfernt ist.

2) Variante für  $A \in \mathbb{C}^{n \times n}$  hermitesch:

Ist  $A$  hermitesch, so erhält man eine bessere Näherung von  $\lambda_1$ , wenn man zur Berechnung von  $\lambda^k$  den Rayleigh-Quotienten verwendet, d.h.

$$\lambda^k := R_A(z^k).$$

In diesem Fall gilt:

$$\lambda^k - \lambda_1 = \mathcal{O}\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right) \quad (k \rightarrow \infty).$$

### Definition 4.31 (Inverse Iteration nach Wielandt)

Sei  $A \in \mathbb{C}^{n \times n}$  eine diagonalisierbare Matrix mit

$$|\lambda_1| \geq |\lambda_2| \geq \dots > |\lambda_n|,$$

so erhält man eine Näherung des Kehrwertes  $\frac{1}{\lambda_n}$  des kleinsten Eigenwertes von  $A$ , indem man in der Vektoriteration nach von Mises  $A$  durch  $A^{-1}$  ersetzt.

### Bemerkung 4.32 (Inverse Iteration mit Diagonal-Shift)

Ist  $A \in \mathbb{C}^{n \times n}$  diagonalisierbar und  $\mu \neq \lambda_i$  für alle  $i = 1, \dots, n$  eine gute Näherung eines Eigenwertes  $\lambda_j$  mit

$$|\lambda_j - \mu| \ll |\lambda_i - \mu| \quad \forall i \neq j,$$

so kann die Näherung  $\mu$  durch inverse Iteration für die Matrix  $B := A - \mu I$  verbessert werden (Diagonal-Shift).

## 4.7 Das QR-Verfahren

**Ziel:** Verwende Ähnlichkeitstransformationen, um  $A$  sukzessive auf obere Dreiecksform zu transformieren:

$$\begin{aligned} A^0 &:= A, \\ A^i &:= T_i^{-1} A^{i-1} T_i, \quad i = 1, 2, \dots \end{aligned}$$

Beim QR-Verfahren wird  $T_i$  unitär gewählt!

**Definition 4.33 (QR-Verfahren)**

Sei  $A \in \mathbb{C}^{n \times n}$ . Dann ist das QR-Verfahren definiert durch:

$$\begin{aligned} A^0 &:= A, \\ A^i &:= Q^i R^i, & (QR\text{-Zerlegung von } A^i) \\ A^{i+1} &:= R^i Q^i. \end{aligned}$$

Hierbei ist  $Q^i$  unitär und  $R^i$  obere Dreiecksmatrix (siehe Numerik I, Kapitel 2.2). Wegen

$$A^{i+1} = R^i Q^i = (Q^i)^H A^i Q^i$$

sind alle Iterierten  $A^i$  ähnlich zu  $A$ .

**Satz 4.34 (Konvergenz des QR-Verfahrens)**

Die Eigenwerte von  $A \in \mathbb{C}^{n \times n}$  seien betragsmäßig getrennt, d.h.

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|.$$

Dann gilt für die Diagonalelemente  $a_{jj}^i$  der Matrizen  $A^i$  des QR-Verfahrens 4.33:

$$\{\lim_{i \rightarrow \infty} a_{jj}^i \mid j = 1, \dots, n\} = \{\lambda_1, \dots, \lambda_n\}.$$

Weiter gilt

$$\lim_{i \rightarrow \infty} a_{jk}^i = 0 \text{ für } j > k.$$

*Beweis:* Siehe z.B. Stoer, Bulirsch [? ].

□

**Bemerkung 4.35**

Da die Berechnung der QR-Zerlegung für allgemeine Matrizen  $A$  sehr aufwendig ist, bringt man in den Anwendungen  $A$  zunächst durch Housholder-Transformation auf Hessenberg-Form. Die Berechnung der QR-Zerlegung einer Hessenberg-Matrix kann dann mit Hilfe der sogenannten Givens-Rotation durchgeführt werden.

**Algorithmus 4.36 (QR-Verfahren für Hessenberg-Matrizen)**

Sei  $A \in \mathbb{R}^{n \times n}$  eine Hessenberg-Matrix. Der folgende Algorithmus bestimmt in Schritt 1 (implizit) die QR-Zerlegung  $A = QR$  und überschreibt in Schritt 2  $A$  mit  $A = RQ$ .

1. Für  $k = 1, \dots, n-1$ :

Bestimme  $c_k = \cos \Phi_k$  und  $s_k = \sin \Phi_k$  mit

$$\begin{pmatrix} c_k & -s_k \\ s_k & c_k \end{pmatrix} \begin{pmatrix} a_{k,k} \\ a_{k+1,k} \end{pmatrix} = \begin{pmatrix} * \\ 0 \end{pmatrix}.$$

Für  $j = k, \dots, n$  setze

$$\begin{pmatrix} a_{k,j} \\ a_{k+1,j} \end{pmatrix} := \begin{pmatrix} c_k & -s_k \\ s_k & c_k \end{pmatrix} \begin{pmatrix} a_{k,j} \\ a_{k+1,j} \end{pmatrix}.$$

2. Für  $k = 1, \dots, n-1$ :

Für  $j = k, \dots, n$  setze

$$(a_{j,k}, a_{j,k+1}) = (a_{j,k}, a_{j,k+1}) \begin{pmatrix} c_k & s_k \\ -s_k & c_k \end{pmatrix}.$$

**Bemerkung 4.37 (LR-Verfahren)**

Ersetzt man in Definition 4.33 die QR-Zerlegung durch die LR-Zerlegung, so erhält man das LR-Verfahren. Unter geeigneten Voraussetzungen gilt

$$\lim_{i \rightarrow \infty} A^i = \lim_{i \rightarrow \infty} R^i = \begin{pmatrix} \lambda_1 & & * \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix}$$

und

$$\lim_{i \rightarrow \infty} L^i = I.$$

Nachteil:

- 1) Eventuell ist Pivotisierung notwendig und gilt nur  $P^i A^i = L^i R^i$  mit einer Permutationsmatrix  $P^i$ , so ist die Konvergenz nicht gesichert.
- 2)  $L^i$  ist nicht unitär und das Verfahren konvergiert schlechter als das QR-Verfahren.

Vorteil:

Für Hessenberg-Matrizen  $A \in \mathbb{R}^{n \times n}$  ist die Berechnung der LR-Zerlegung mit dem Gauß-Algorithmus etwa doppelt so schnell wie die Berechnung der QR-Zerlegung.

# Kapitel 5

## Approximation

### 5.1 Allgemeine Approximation in normierten Räumen

**Definition 5.1 (Beste Approximation/Proximum)**

Sei  $(X, \|\cdot\|)$  ein normierter Raum und  $T \subset X$  eine beliebige Teilmenge. Zu einem  $v \in X$  definiere

$$J_v(u) := \|v - u\|.$$

Dann heißt ein  $u \in T$  beste Approximation oder Proximum von  $v$  in  $T$ , g.d.w.

$$J_v(u) = \inf_{w \in T} J_v(w).$$

Die Zahl  $E_v(T) := \inf_{w \in T} J_v(w)$  heißt Minimalabstand von  $v \in X$  zur Teilmenge  $T$ .

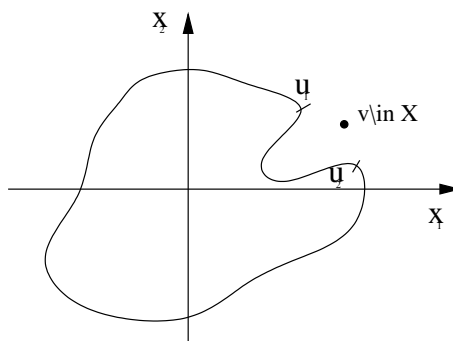


Abbildung 5.1: Proximum: Beispiel

**Beispiel 5.2**

- 1) Sei  $X = \mathbb{R}^2$  und  $\|\cdot\| = \|\cdot\|_2$  die euklidische Norm. Sei  $T := \{x \in \mathbb{R}^2 \mid \|x\| \leq 1\}$ . Dann existiert zu jedem  $v \in X$  eine beste Approximation  $u \in T$ :  
 $\implies$  Hier existiert zu jedem  $v \in X$  genau ein Proximum.

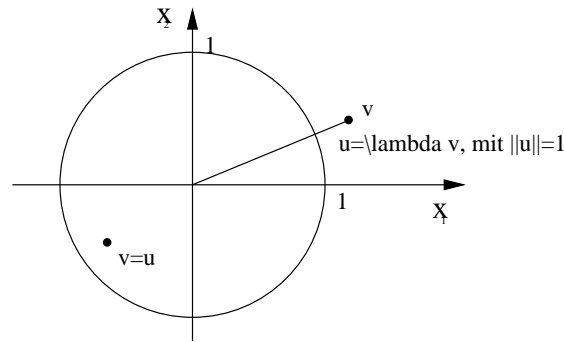


Abbildung 5.2: Proximum: Beispiel 1)

2) Sei  $X = \mathbb{R}^2$ ,  $\|\cdot\| = \|\cdot\|_2$  und  $T := \{x \in \mathbb{R}^2 \mid \|x\| < 1\}$ .

Ist  $v \in X \setminus T$ , so existiert keine beste Approximation  $u \in T$  von  $V$ , denn  $E_v(T) = \|v\| - 1$  und für alle  $w \in T$  gilt  $\|w - v\| > \|v\| - 1$

$\implies$  Hier existiert für  $x \in X \setminus T$  kein Proximum.

3) Sei  $X = \mathbb{R}^2$ ,  $\|\cdot\| = \|\cdot\|_\infty$ , d.h.  $\|x\|_\infty = \max\{|x_1|, |x_2|\}$  und  $T := \{(x_1, 0) \mid x_1 \in \mathbb{R}\}$ . Sei  $v = (0, 1) \in X$ . Dann gilt  $u \in [-1, 1] \times \{0\} \implies u$  ist Proximum, da  $J_v(u) = 1 = E_v(T)$

$\implies$  Hier existieren unendlich viele beste Approximationen.

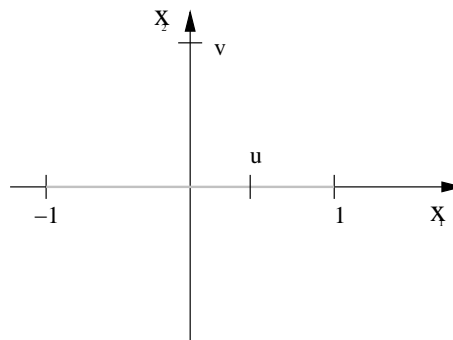


Abbildung 5.3: Proximum: Beispiel 3)

### Satz 5.3 (Existenz eines Proximums)

Sei  $T \subset X$  eine kompakte Teilmenge. Dann existiert zu jedem  $v \in X$  ein Proximum  $u \in T$ .

*Beweis:* Sei  $(u_n)_{n \in \mathbb{N}}$  eine Minimalfolge in  $T$  für  $v \in X$ , d.h.  $\lim_{n \rightarrow \infty} J_v(u_n) = E_v(T)$ . Da  $T$  kompakt ist, enthält  $(u_n)_{n \in \mathbb{N}}$  eine Teilfolge, die in  $T$  konvergiert, d.h.  $\exists u \in T$  mit

$$\lim_{j \rightarrow \infty} u_{n_j} = u.$$

Zu zeigen:  $u$  ist Proximum, d.h.  $J_v(u) = E_v(T)$ .

Es gilt:

$$\begin{aligned} \|v - u\| &\leq \|v - u_{n_j}\| + \|u_{n_j} - u\| \\ &\stackrel{\downarrow (j \rightarrow \infty)}{\leq} E_v(T) \stackrel{\downarrow (j \rightarrow \infty)}{\leq} 0 \end{aligned}$$

$$\implies \|v - u\| \leq E_v(T).$$

Da  $E_v(T) = \inf_{w \in T} J_v(w) = \inf_{w \in T} \|v - w\|$ , folgt:

$$J_v(u) = \|v - u\| = E_v(T).$$

□

**Definition 5.4 (Konvexe Teilmengen)**

$T \subset X$  heißt *konvex*, g.d.w.

$$K_{u_1, u_2} := \{\lambda u_1 + (1 - \lambda)u_2 \mid \lambda \in (0, 1)\} \subset T$$

für alle  $u_1, u_2 \in T$ .

$T \subset X$  heißt *streng konvex*, g.d.w.

$$K_{u_1, u_2} \subset \overset{\circ}{T} \quad \forall u_1, u_2 \in T$$

mit  $u_1 \neq u_2$  (vgl. Abbildung 5.4).

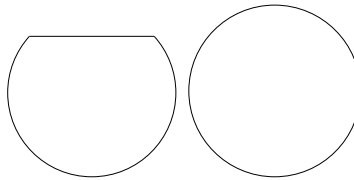


Abbildung 5.4: Konvexe und streng konvexe Mengen.

**Satz 5.5 (Eindeutigkeit von Proxima)**

Sei  $T \subset X$  kompakt und streng konvexe Teilmenge des normierten Raumes  $X$ . Dann gibt es zu jedem  $v \in X$  genau ein Proximum  $u \in T$ .

*Beweis:* Die Existenz ist klar nach Satz 5.3.

Eindeutigkeit: Seien  $u_1, u_2$  mit  $u_1 \neq u_2$  Proxima von  $v \in X$  in  $T$ . Dann gilt:

$$\left\| \frac{1}{2}(u_1 + u_2) - v \right\| \leq \frac{1}{2} \underbrace{\|u_1 - v\|}_{=E_v(T)} + \frac{1}{2} \underbrace{\|u_2 - v\|}_{=E_v(T)} = E_v(T).$$

$$T \text{ konvex} \implies \left\| \frac{1}{2}(u_1 + u_2) - v \right\| = E_v(T).$$

Da  $T$  streng konvex ist, folgt  $\frac{1}{2}(u_1 + u_2) \in \overset{\circ}{T}$ . Also existiert ein  $\tilde{\lambda} \in (0, 1)$ , so dass

$$\tilde{u} = \frac{1}{2}(u_1 + u_2) + \tilde{\lambda}(v - \frac{1}{2}(u_1 + u_2)) \in T$$

ist. Dann gilt:

$$\begin{aligned} \|\tilde{u} - v\| &= \left\| \frac{1}{2}(1 - \tilde{\lambda})(u_1 + u_2) - (1 - \tilde{\lambda})v \right\| \\ &= (1 - \tilde{\lambda}) \left\| \frac{1}{2}(u_1 + u_2) - v \right\| \\ &= (1 - \tilde{\lambda})E_v(T) \\ &< E_v(T). \end{aligned}$$

Dies ist ein Widerspruch zur Definition von  $E_v(T) \implies u_1 = u_2$ . □

Für Anwendungen ist vor allem der Fall wichtig, dass  $T$  ein endlich dimensionaler Teilraum von  $X$  ist.

**Satz 5.6 (Fundamentalsatz der Approximationstheorie in normierten Räumen)**

Sei  $T \subset X$  ein endlichdimensionaler linearer Teilraum des normierten Vektorraums  $X$ . Dann existiert zu jedem  $v \in X$  ein Proximum  $u \in T$ .

*Beweis:* Sei  $(u_n)_{n \in \mathbb{N}}$  Minimalfolge für  $v \in X$ .

Wir zeigen zunächst:  $(u_n)_{n \in \mathbb{N}}$  ist beschränkt.

Es gilt:

$$E_v(T) \leq \|v - u_n\| \leq E_v(T) + 1 \quad \forall n \geq N.$$

Also ist  $\|u_n\| \leq \|v - u_n\| + \|v\| \leq E_v(T) + 1 + \|v\| =: K_1 \quad \forall n \geq N$ .

Sei nun  $K_2 \geq \|u_n\|$  für  $n < N$  und setze  $K = \max\{K_1, K_2\}$ , so folgt  $\|u_n\| \leq K \quad \forall n \in \mathbb{N}$ .

Da  $T$  endlich dimensionaler linearer Teilraum ist, existiert also eine Teilfolge  $(u_{n_j})_{j \in \mathbb{N}}$ , die gegen ein  $u \in T$  konvergiert.

Analog zum Beweis von Satz 5.3 zeigt man, dass  $u$  ein Proximum ist. □

**Definition 5.7 (Streng normierter Raum)**

Sei  $(X, \|\cdot\|)$  ein normierter Raum.  $\|\cdot\|$  heißt strenge Norm und  $X$  streng normiert, g.d.w.

$$(\|f + g\| = \|f\| + \|g\| \text{ für } f, g \in X \text{ mit } f, g \neq 0) \implies (\exists \lambda \in \mathbb{C} \text{ mit } g = \lambda f).$$

**Satz 5.8 (Eindeutigkeit in streng normierten Räumen)**

Ist  $(X, \|\cdot\|)$  ein streng normierter Raum und  $T \subset X$  ein endlichdimensionaler linearer Teilraum, so existiert zu jedem  $v \in X$  genau ein Proximum  $u \in T$ .

*Beweis:* Ist  $v \in T$ , so ist  $u = v$  das eindeutige Proximum.

Sei also  $v \in X \setminus T$ . Sind  $u_1, u_2$  verschiedene Proxima zu  $v$ , so gilt wie in Beweis von Satz 5.5:

$$E_v(T) \leq \left\| v - \frac{1}{2}(u_1 + u_2) \right\| \leq \frac{1}{2} \|v - u_1\| + \frac{1}{2} \|v - u_2\| = E_v(T)$$

also  $\|(v - u_1) + (v - u_2)\| = \|v - u_1\| + \|v - u_2\|$

Da  $\|\cdot\|$  strenge Norm ist,  $\exists \lambda \in \mathbb{C}$  mit  $v - u_1 = \lambda(v - u_2)$

$\implies (1 - \lambda)v = u_1 - \lambda u_2$ .

Da  $v \notin T$ , folgt  $\lambda = 1$  und somit  $0 = u_1 - u_2$ . Dies ist ein Widerspruch zur Annahme. □



## 5.2 Der Satz von Weierstraß: Approximation durch Polynome

### Motivation 5.9

Bei der Approximation durch Polynome wurde in der Analysis immer von der Regularität der Funktion  $f$  gebrauch gemacht. So z.B. bei der Approximation von  $f$  durch die Taylorreihe.

In diesem Kapitel stellen wir uns die Frage, ob eine beliebig gute Polynomapproximation auch für Funktionn  $f \in C([a, b])$  möglich ist.

### Satz 5.10 (Approximationssatz von Weierstraß)

Sei  $X = C([a, b])$  und  $\|\cdot\| = \|\cdot\|_\infty$  für  $|a|, |b| < \infty$ .

Dann gibt es zu jedem  $f \in X$  und  $\varepsilon > 0$  ein  $n \in \mathbb{N}$  und ein  $p \in \mathbb{P}_n$ , so dass  $\|f - p\|_\infty < \varepsilon$  ist.

*Beweis:* (Konstruktiv!)

Ohne Einschränkung sei  $[a, b] = [0, 1]$  (sonst Transformation).

Wir zeigen, dass die Folge der Bernstein-Polynome

$$(B_n f)(x) := \sum_{i=0}^n f\left(\frac{i}{n}\right) \binom{n}{i} x^i (1-x)^{n-i}$$

für  $n \rightarrow \infty$  auf  $[0, 1]$  gleichmäßig gegen  $f$  konvergieren.

Definiere  $q_{ni} := \binom{n}{i} x^i (1-x)^{n-i}$ , dann ist

$$1 = (x + (1-x))^n = \sum_{i=0}^n \binom{n}{i} x^i (1-x)^{n-i} = \sum_{i=0}^n q_{ni}.$$

Also folgt:

$$\begin{aligned} f(x) - (B_n f)(x) &= \sum_{i=0}^n \left(f(x) - f\left(\frac{i}{n}\right)\right) q_{ni}(x) \\ \implies |f(x) - (B_n f)(x)| &\leq \sum_{i=0}^n \left|f(x) - f\left(\frac{i}{n}\right)\right| q_{ni}(x) \quad \forall x \in [0, 1]. \end{aligned}$$

Da  $f$  gleichmäßig stetig ist, existiert zu jedem  $\varepsilon > 0$  ein  $\delta$ , so dass  $\forall x, i$

$$\left|x - \frac{i}{n}\right| < \delta \implies \left|f(x) - f\left(\frac{i}{n}\right)\right| < \frac{\varepsilon}{2}.$$

Sei  $x \in [0, 1]$ . Setze

$$N_{<} := \{i \in \{0, \dots, n\} \mid \left|x - \frac{i}{n}\right| < \delta\},$$

$$N_{\geq} := \{i \in \{0, \dots, n\} \mid \left|x - \frac{i}{n}\right| \geq \delta\}.$$

Wir erhalten:

$$\sum_{i \in N_{<}} \left|f(x) - f\left(\frac{i}{n}\right)\right| q_{ni}(x) \leq \frac{\varepsilon}{2} \sum_{i \in N_{<}} q_{ni}(x) \leq \frac{\varepsilon}{2}.$$

und

$$\begin{aligned}
 \sum_{i \in N_{\geq}} \left| f(x) - f\left(\frac{i}{n}\right) \right| q_{ni}(x) &\leq \sum_{i \in N_{\geq}} \left| f(x) - f\left(\frac{i}{n}\right) \right| q_{ni}(x) \frac{(x - \frac{i}{n})^2}{\delta^2} \\
 &\leq \frac{2 \|f\|_{\infty}}{\delta^2} \sum_{i \in N_{\geq}} q_{ni}(x) (x - \frac{i}{n})^2 \\
 &\leq \frac{2 \|f\|_{\infty}}{\delta^2} \sum_{i=0}^n q_{ni}(x) (x^2 - 2x \frac{i}{n} + (\frac{i}{n})^2).
 \end{aligned}$$

Es ist

$$1) \sum_{i=0}^n q_{ni}(x) x^2 = x^2,$$

$$2) \sum_{i=0}^n q_{ni}(x) 2x \frac{i}{n} = 2x \cdot x \underbrace{\sum_{i=1}^n \binom{n-1}{i-1} x^{i-1} (1-x)^{(n-1)-(i-1)}}_{=1} = 2x^2,$$

3)

$$\begin{aligned}
 \sum_{i=0}^n q_{ni}(x) \left(\frac{i}{n}\right)^2 &= \frac{x}{n} \sum_{i=1}^n (i-1) \binom{n-1}{i-1} x^{i-1} (1-x)^{(n-1)-(i-1)} + \frac{x}{n} \\
 &= \frac{x^2}{n} (n-1) \underbrace{\sum_{i=2}^n \binom{n-2}{i-2} x^{i-2} (1-x)^{(n-2)-(i-2)}}_{=1} + \frac{x}{n} \\
 &= x^2 \left(1 - \frac{1}{n}\right) + \frac{x}{n} = x^2 + \frac{x}{n} (1-x).
 \end{aligned}$$

Mit 1), 2) und 3) folgt:

$$\begin{aligned}
 \sum_{i \in N_{\geq}} \left| f(x) - f\left(\frac{i}{n}\right) \right| q_{ni}(x) &\leq \frac{2 \|f\|_{\infty}}{\delta^2} \underbrace{(x^2 - 2x + x^2)}_{=0} + \underbrace{\frac{x}{n} (1-x)}_{\leq \frac{1}{4n}} \\
 &\leq \frac{2 \|f\|_{\infty}}{\delta^2} \frac{1}{4n} \\
 &< \frac{\varepsilon}{2} \text{ f\"ur } n > \frac{M}{\delta^2 \varepsilon}.
 \end{aligned}$$

Insgesamt folgt also

$$|f(x) - (B_n f)(x)| \leq \frac{\varepsilon}{2} + \frac{\varepsilon}{2} = \varepsilon \quad \forall x \in [0, 1].$$

□

### Bemerkung 5.11

Satz 5.10 zeigt, dass sich  $f \in C([a, b])$  in folgender Weise entwickeln lässt:

$$f(x) = (B_1 f)(x) + [(B_2 f)(x) - (B_1 f)(x)] + \dots + [(B_n f)(x) - (B_{n-1} f)(x)] + \dots$$

Die Reihe konvergiert gleichmäßig, lässt sich aber im allgemeinen nicht zu einer Potenzreihe umordnen!

**Satz 5.12 (Fehlerabschätzung für die Approximation mit Bernsteinpolynomen)**

Seien die Voraussetzungen von Satz 5.10 erfüllt und sei

$$w_f(\delta) := \sup_{|x'-x''| \leq \delta, x', x'' \in [a, b]} |f(x') - f(x'')|$$

das Stetigkeitsmodul von  $f$  bezgl.  $\delta$ . Dann gilt die Abschätzung:

$$|f(x) - (B_n f)(x)| \leq \frac{5}{4} w_f\left(\frac{1}{\sqrt{n}}\right).$$

*Beweis:* Sei  $\lambda = \lambda(x', x'', \delta) := \left\lceil \frac{|x' - x''|}{\delta} \right\rceil$  das größte Ganze.

Mit der Definition von  $w_f(\delta)$  folgt:

$$\delta_1 \leq \delta_2 \implies w_f(\delta_1) \leq w_f(\delta_2).$$

Damit folgt:

$$|f(x') - f(x'')| \leq w_f(|x' - x''|) \leq w_f((\lambda + 1)\delta).$$

Aus  $w_f(\mu\delta) \leq \mu w_f(\delta)$  für  $\mu \in \mathbb{N}$  folgt:

$$|f(x') - f(x'')| \leq (\lambda + 1)w_f(\delta).$$

Setze  $N_{\geq} := \{i \in \{0, \dots, n\} \mid \lambda(x, \frac{i}{n}, \delta) \geq 1\}$  und  $N_{<}$  entsprechend. Jetzt gehen wir analog zum Beweis von Satz 5.10 vor:

$$\begin{aligned} |f(x) - (B_n f)(x)| &\leq \sum_{i=0}^n \left| f(x) - f\left(\frac{i}{n}\right) \right| q_{ni} \\ &\leq w_f(\delta) \sum_{i=0}^n \left(1 + \lambda\left(x, \frac{i}{n}, \delta\right)\right) q_{ni}(x). \end{aligned}$$

Da  $\lambda(x, \frac{i}{n}, \delta) = 0$  für alle  $i \in N_{<}$  gilt, folgt weiter:

$$\begin{aligned} |f(x) - (B_n f)(x)| &\leq w_f(\delta) \left(1 + \sum_{i \in N_{\geq}} \lambda\left(x, \frac{i}{n}, \delta\right) q_{ni}(x)\right) \\ &\leq w_f(\delta) \left(1 + \frac{1}{\delta} \sum_{i \in N_{\geq}} \left|x - \frac{i}{n}\right| q_{ni}(x)\right) \\ &\stackrel{\substack{|x - \frac{i}{n}|}{\delta} \geq 1 \quad \forall i \in N_{\geq}}}{\leq} w_f(\delta) \left(1 + \frac{1}{\delta^2} \sum_{i \in N_{\geq}} \left(x - \frac{i}{n}\right)^2 q_{ni}(x)\right) \end{aligned}$$

Analog zum  
Beweis von 5.10

$$\leq w_f(\delta) \left(1 + \frac{1}{4n\delta^2}\right).$$

Wählen wir  $\delta = \frac{1}{\sqrt{n}}$ , so folgt die Abschätzung des Satzes. □

**Bemerkung 5.13**

- 1) Abhängig vom Stetigkeitsmodul kann die Schranke in Satz 5.12 beliebig langsam konvergieren. Bei höheren Anforderungen an die Stetigkeit von  $f$  kann andererseits eine schnellere Konvergenz erwartet werden.
- 2) In der Praxis ist die Approximation durch Bernstein-Polynome nicht von Bedeutung. Im nächsten Abschnitt werden wir wirkungsvollere Verfahren kennen lernen!

### 5.3 Gleichmäßige Approximation / Tschebyschev Approximation

#### Motivation 5.14

In §5.2 haben wir gesehen, dass sich jede stetige Funktion  $f \in C^0([a, b])$  durch Polynome  $p_n \in \mathbb{P}_n$  approximieren lassen. Die Frage, welches Polynom  $p \in \mathbb{P}_n$  das Proximum zu  $f$  in  $T = \mathbb{P}_n$  ist, wurde jedoch nicht beantwortet. Ist  $\|\cdot\| = \|\cdot\|_\infty$ , so beantworten wir diese Frage in diesem Abschnitt.  $\|\cdot\|_\infty$  wird auch Tschebyschev-Norm genannt.

#### Definition 5.15 (Best Approximating Polynomial (BAP))

Sei  $X = C^0(I)$ ,  $I \subset \mathbb{R}$  beschränktes Intervall ausgestattet mit der  $\infty$ -Norm  $\|\cdot\| = \|\cdot\|_\infty$ .

Dann heißt  $p_n \in \mathbb{P}_n$  Best Approximating Polynomial (BAP) vom Grad  $\leq n$  von  $f \in C^0(I)$ , g.d.w.  $p_n$  Proximum von  $f$  in  $\mathbb{P}_n$  ist.

#### Bemerkung 5.16

- 1) Der Fundamentalsatz der Approximationstheorie 5.6 liefert die Existenz eines (BAP), da  $\mathbb{P}_n$  ein endlichdimensionaler linearer Teilraum von  $C^0(I)$  ist.
- 2) Eindeutigkeit des (BAP) können wir zunächst nicht erwarten, da  $\|\cdot\|_\infty$  keine strenge Norm ist.

#### Satz 5.17 (Charakterisierung von BAP)

Sei  $f \in C^0(I)$ ,  $I = [a, b]$  beschränkt und  $p \in \mathbb{P}_n$ . Es gebe  $n + 2$  Punkte  $a \leq x_0 < x_1 < \dots < x_n < x_{n+1} \leq b$ , so dass

- 1)  $|f(x_i) - p(x_i)| = J_f(p) = \|f - p\|_\infty$  für  $i = 0, \dots, n + 1$ ,
- 2)  $f(x_{i+1}) - p(x_{i+1}) = -(f(x_i) - p(x_i))$  für  $i = 0, \dots, n$ .

Dann ist  $p$  BAP vom Grad  $\leq n$  an  $f$ .

*Beweis:* Sei  $p^* \in \mathbb{P}_n$  und  $M := \{x \in I \mid |f(x) - p^*(x)| = J_f(p^*)\}$ .

Ist  $p^*$  kein BAP, so existiert ein  $\bar{p} \in \mathbb{P}_n$ ,  $\bar{p} \neq 0$ , so dass  $p^* + \bar{p}$  BAP ist (nach Satz 5.6). Dann gilt:

$$|(f(x) - p^*(x)) - \bar{p}(x)| = |f(x) - (p^*(x) + \bar{p}(x))| < |f(x) - p^*(x)| \quad \forall x \in M$$

$$\implies (f(x) - p^*(x))\bar{p}(x) > 0 \quad \forall x \in M \quad (*)$$

Sei nun  $p \in \mathbb{P}_n$  und es gelten die Voraussetzungen 1) und 2).

Dann kann es kein  $\bar{p} \neq 0$ ,  $\bar{p} \in \mathbb{P}_n$  geben, so dass (\*) erfüllt ist. Denn dazu müsste  $\bar{p}$  in  $[a, b]$  mindestens  $(n + 1)$ -mal das Vorzeichen wechseln (wegen 2)), also mindestens  $n + 1$  Nullstellen besitzen. Nach dem Fundamentalsatz der Algebra ist das nicht möglich.

Da es zu  $p$  kein  $\bar{p}$  gibt, so dass (\*) gilt, muß  $p$  bereits BAP sein.

□

#### Satz 5.18 (Eindeutigkeit eines Proximums in $\mathbb{P}_n$ )

Sei  $U := \mathbb{P}_n$  Unterraum von  $C^0([a, b])$ . Dann ist das Proximum  $p \in U$  an ein  $f \in C^0([a, b])$  eindeutig bestimmt.

*Beweis:* Seien  $p_1$  und  $p_2$  Proxima aus  $U$  an  $f \in C^0([a, b])$ .

Dann ist auch  $\frac{1}{2}(p_1 + p_2)$  Proximum (Satz 5.5) und nach dem Charakterisierungssatz 5.17 existiert eine sogenannte Alternante der Länge  $n + 2$ , d.h.

$$f(x_i) - \frac{1}{2}(p_1(x_i) + p_2(x_i)) = \sigma(-1)^i E_f(U).$$

Also ist  $\frac{1}{2}(f(x_i) - p_1(x_i)) + \frac{1}{2}(f(x_i) - p_2(x_i)) = \sigma(-1)^i E_f(U)$ .

Da  $p_1, p_2$  Proxima sind, gilt  $|f(x_i) - p_k(x_i)| \leq E_f(U)$   $k = 1, 2$  und  $i = 0, \dots, n + 1$ . Daraus folgt

$$f(x_i) - p_1(x_i) = f(x_i) - p_2(x_i) \quad \forall i = 0, \dots, n + 1$$

$$\implies p_1(x_i) = p_2(x_i) \quad \forall i = 0, \dots, n + 1.$$

Da  $U$  Polynome vom maximalen Grad  $n$  sind, folgt hieraus  $p_1 - p_2 \equiv 0$ , da  $p_1 - p_2$  mindestens  $n + 2$  Nullstellen hat.

□

Der Charakterisierungssatz 5.17 bildet die Grundlage für ein Verfahren zur Approximation des BAP zu einem  $f \in C^0([a, b])$ .

### (Austauschalgorithmus von Remez) 5.19

Sei  $f \in C^0([a, b])$ . Dann erzeugt der folgende Algorithmus von Remez eine Folge von Polynomen  $p^{(k)} \in \mathbb{P}_n$  mit

$$\lim_{k \rightarrow \infty} \|p^{(k)} - p_n\|_{\infty} = 0,$$

wobei  $p_n$  das eindeutig bestimmte BAP von  $f$  in  $\mathbb{P}_n$  ist.

Als Abbruchkriterium kann die Bedingung

$$\Delta - \delta \leq \varepsilon$$

gewählt werden, wobei  $\delta, \Delta$  bezüglich  $p^{(k)}$  wie folgt definiert sind

$$\delta := \min_{i=0, \dots, n+1} |f(x_i) - p(x_i)|,$$

$$\Delta := \max_{i=0, \dots, n+1} |f(x_i) - p(x_i)|.$$

(Tatsächlich kann gezeigt werden, dass gilt  $\delta \leq E_f(U) \leq \Delta$ .)

#### Remez-Algorithmus:

Input: Startzerlegung  $I^{(0)}$  von  $[a, b]$ , d.h.

$$I^{(0)} := \{x_0, \dots, x_{n+1}\}$$

mit  $a \leq x_0 < x_1 < \dots < x_{n+1} \leq b$ .

Iteration für  $k = 0, 1, 2, \dots$

1. Schritt: Zu  $I^{(k)}$  bestimme Polynom  $p^{(k)} \in \mathbb{P}_n$ , so dass  $I^{(k)}$  Alternante für  $f - p^{(k)}$  ist und für alle  $x_i \in I^{(k)}$  gilt:

$$\left| f(x_i) - p^{(k)}(x_i) \right| = \left| \xi^{(k)} \right|$$

mit einer Konstante  $\xi^{(k)} \in \mathbb{R}$ .

Setze man  $p^{(k)}(x) := \sum_{j=0}^n \alpha_j^{(k)} x^j$ , so führen diese Forderungen auf das Gleichungssystem:

$$(-1)^i \xi^{(k)} + \sum_{j=0}^n \alpha_j^{(k)} x_i^j = f(x_i), \quad i = 0, \dots, n+1$$

für die Unbekannten  $\xi^{(k)}, \alpha_0^{(k)}, \dots, \alpha_n^{(k)}$ .

2. Schritt: Setze  $r^{(k)}(x) := f(x) - p^{(k)}(x)$  und bestimme  $y \in [a, b]$  mit der Eigenschaft:

$$r^{(k)}(y) = \max_{x \in [a, b]} r^{(k)}(x).$$

Ersetze einen Punkt  $x_i \in I^{(k)}$  durch  $y$  und erhalte  $I^{(k+1)}$  gemäß folgender Austauschvorschrift:

1. Fall:  $x_j < y < x_{j+1}$  für eine  $j \in \{0, \dots, n\}$ .  
Falls  $\text{sign}(r^{(k)}(x_j)) = \text{sign}(r^{(k)}(y))$ , ersetze  $x_j$  durch  $y$ , sonst ersetze  $x_{j+1}$  durch  $y$ .
2. Fall:  $y < x_0$ .  
Falls  $\text{sign}(r^{(k)}(x_0)) = \text{sign}(r^{(k)}(y))$ , ersetze  $x_0$  durch  $y$ , sonst ersetze  $x_{n+1}$  durch  $y$ .
3. Fall:  $y > x_{n+1}$ .  
Falls  $\text{sign}(r^{(k)}(x_{n+1})) = \text{sign}(r^{(k)}(y))$ , ersetze  $x_{n+1}$  durch  $y$ , sonst ersetze  $x_0$  durch  $y$ .

Bemerkung: Gilt  $y = x_i$  für ein  $x_i \in I^{(k)}$ , so ist  $p^{(k)}$  bereits BAP. (Charakterisierungssatz 5.17)

Beweis: Skizze: Man zeigt:

- 1) Das Gleichungssystem in Schritt 1 ist stets eindeutig lösbar.
- 2) Es gilt  $\xi^{(k+1)} > \xi^{(k)} \quad \forall k = 0, 1, 2, \dots$
- 3)  $\lim_{m \rightarrow \infty} I^{(k_m)} = \{\bar{x}_0, \dots, \bar{x}_{n+1}\}$  für eine Teilfolge.
- 4)  $\lim_{m \rightarrow \infty} p^{(k_m)} = p$  für diese Teilfolge.
- 5)  $p$  ist BAP an  $f$ .

□

Im folgenden wollen wir der Frage nachgehen, wie für  $f(x) = x^n$  in  $[-1, 1]$  das BAP in  $\mathbb{P}_{n-1}$  aussieht. Dies führt uns zu den Tschebyschev Polynomen 1. Art.

**Definition 5.20 (Tschebyschev Polynome 1. Art)***Durch die Rekursion*

$$T_0(x) = 1, T_1(x) = x$$

$$T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x), n \geq 1$$

*werden die sogenannten Tschebyschev Polynome 1. Art definiert.***Satz 5.21**

Die Tschebyschev-Polynome 1. Art haben die Darstellung  $T_n(x) = \cos(n \arccos(x))$ ,  $x \in [-1, 1]$ ,  $n \in \mathbb{N}$ . Es gelten die Eigenschaften

- 1)  $|T_n(x)| < 1 \quad \forall x \in [-1, 1]$ ,
- 2)  $T_n(x)$  hat in  $[-1, 1]$  die Extrempunkte

$$x_i^{(n)} = \cos\left(\frac{i\pi}{n}\right), T_n(x_i^{(n)}) = (-1)^i \quad i = 0, \dots, n,$$

- 3)  $T_n$  hat  $n$  einfache Nullstellen in  $[-1, 1]$

$$\tilde{x}_i^{(n)} = \cos\left(\frac{2i-1}{2n}\pi\right) \quad i = 1, \dots, n,$$

- 4) Zwischen je zwei Nullstellen von  $T_{n+1}$  liegt eine Nullstelle von  $T_n$ .

*Beweis: (ohne Beweis)*

□

**Satz 5.22**

Sei  $f(x) = x^n \in C^0([-1, 1])$ . Dann ist

$$p_{n-1} := x^n - \frac{1}{2^{n-1}} T_n(x)$$

BAP zu  $f$  in  $\mathbb{P}_{n-1}$ .

*Beweis: 1.) Zeige  $p_{n-1} \in \mathbb{P}_{n-1}$ .*

Unter Verwendung der Rekursionsformel zeigt man induktiv, dass  $T_n(x)$  die Form

$$T_n(x) = 2^{n-1} x^n + \sum_{i=1}^{n-1} \alpha_i x^i$$

hat. Dann folgt:

$$p_{n-1} = x^n - \frac{1}{2^{n-1}} T_n(x) = -\frac{1}{2^{n-1}} \sum_{i=1}^{n-1} \alpha_i x^i \implies p_{n-1} \in \mathbb{P}_{n-1}.$$



2.) Zeige  $p_{n-1}$  ist BAP zu  $f$ .

Es gilt  $x^n - p_{n-1}(x) = \frac{1}{2^{n-1}}T_n(x)$ . Nach Satz 5.21 (2) ist  $\{x_i^{(n)} := \cos(\frac{i\pi}{n}) \mid i = 0, \dots, n\}$  eine Alternante der Länge  $(n-1) + 2$  für  $p_{n-1}(x)$  und es gilt nach 5.21 1), 2)

$$\left| \frac{1}{2^{n-1}}T_n(x_n^{(i)}) \right| = \frac{1}{2^{n-1}} = \max_{x \in [-1,1]} \left| \frac{1}{2^{n-1}}T_n(x) \right|.$$

Nach dem Charakterisierungssatz 5.17 muß  $p_{n-1}$  BAP zu  $f$  sein. □

### Satz 5.23

Die Tschebyschev-Polynome bilden bezüglich der Gewichtsfunktion  $w(x) := \frac{1}{\sqrt{1-x^2}}$  ein Orthogonalsystem. Es gilt:

$$\int_{-1}^1 T_n(x)T_m(x)w(x)dx = \begin{cases} \pi & \text{für } n = m = 0 \\ \frac{\pi}{2} & n = m \neq 0 \\ 0 & n \neq m \end{cases}.$$

(ohne Beweis)

### Definition 5.24 (Tschebyschev-Entwicklung)

Sei  $f$  stetig in  $[-1, 1]$ . Dann heißen

$$a_k(f) := \frac{2}{\pi} \int_{-1}^1 f(x)T_k(x) \frac{1}{\sqrt{1-x^2}} dx, \quad k = 0, 1, 2, \dots$$

Tschebyschev-Koeffizienten von  $f$  und die formal gebildete Reihe

$$S_f^n(x) = \frac{a_0(f)}{2} + \sum_{k=1}^n a_k(f) \cdot T_k(x)$$

heißt Tschebyschev-Entwicklung von  $f$ .

### Entwicklungssatz 5.25

Sei  $f \in C^2([-1, 1])$ . Dann konvergiert die Tschebyscheventwicklung für  $x \in [-1, 1]$  gleichmäßig gegen  $f$ , d.h.

$$S_f^\infty := \lim_{n \rightarrow \infty} S_f^n = f.$$

Für die Tschebyschev-Koeffizienten gilt die Abschätzung:

$$|a_k(f)| \leq \frac{c}{k^2}, \quad k = 1, 2, 3, \dots$$

*Beweis:* (ohne Beweis) □

### Folgerung 5.26

Die Koeffizientenabschätzung in Satz 5.25 zeigt, dass für  $f \in C^2([-1, 1])$  eine gute  $N$  durch  $S_f^n$  gegeben ist.

Diese Möglichkeit der Approximation bietet sich dann an, wenn die Koeffizienten  $a_k(f)$  einfach zu berechnen sind.

# Index

- LR*-Verfahren, 92
- QR*-Verfahren, 91
- Ähnlichkeitstransformation, 76
- Gaußalgorithmus**, 23
- Pseudoinverse**, 41
- Singulärwertzerlegung**, 39
  
- A-konjugiert, 56
- A-orthogonal, 56
- Abschneidefehler, 119
- Alternante, 103
- Anfangswertproblem, 115
  - höherer Ordnung -, 124
  - Stetigkeitssatz für das -, 117
- Approximation
  - Finite-Differenzen, 14
- Approximationssatz von Weierstraß, 98
- Arithmetische Operationen, 19
- Ausgleichsproblem, 41
- Ausgleichsrechnung, 31
- Austauschalgorithmus von Remez, 105
- AWP, 115
  
- Banachraum, 6, 10
- Banachscher Fixpunktsatz, 10
  - a-posteriori Abschätzung, 10
  - a-priori Abschätzung, 10
  - Fixpunkt, 10
  - Kontraktion, 10
  - TOL, 10
  - Toleranz, 10
- BAP, 102
- Bernstein-Polynome, 98
- Besselsche Ungleichung, 112
- best approximatin polynomial, 102
- Beste Approximation, 93
  
- Cauchy-Schwarz-Ungleichung, 6
- cd-Verfahren, 56
- cg-Verfahren, 57
- charakteristisches Polynom, 75
- Cholesky Verfahren, 31
- Courantsches Minimum-Maximum Prinzip, 82
  
- Cramersche Regel, 22
- Crank-Nicholson Verfahren, 123
  
- Diagonalmatrix, 77
- Diskretisierungsfehler, 119
  - konvergent, 119
- dyadisches Produkt, 36
  
- Eigenraum, 76
- Eigentliches Gradientenverfahren, 53
- Eigenvektor, 9, 75
- Eigenwert, 9, 21, 75
- Eigenwertproblem, 75
- einfache Nullstelle, 64
- Einschrittverfahren, 119
- Einzelschritt Verfahren, 50
- erster Näherung, 12
- Explizites Euler Verfahren, 120, 121
- Explizites Verfahren, 119
  - Abschneidefehler, 119
  - Verfahren, 119
- Explizites Verfahren
  - Euler -, 120
  
- Fehleranalyse, 12
  - Abbruchfehler, 14
  - Approximationsfehler, 12, 13
  - Datenfehler, 13
  - Modellfehler, 12
  - potentielle Energie, 12
  - Rundungsfehler, 15
- Fehlerdämpfung, 17
- Fehlerfortpflanzung, 17
- Feinheit, 119
- Fibonacci-Zahlen, 67
- Finite-Differenzen, 14
- Fixpunkt, 10
- Folge, 6
  - Cauchy Folge, 6
  - Konvergenz, 6
  
- Gauß-Jordan Verfahren, 30
- Gaußalgorithmus, 22

- Pivotisierung, 24
- Spaltenpivotisierung, 24
- Teilpivotisierung, 24
- total pivoting, 24
- Gaußsches Ausgleichsproblem, 32
- Gaußverfahren, 25
- Gerschgorin Kreise, 76
- Gerschgorinscher Kreissatz, 76
- Gesamtschritt Verfahren, 46
- gestörtes AWP, 117
- gewöhnlicher Differentialgleichungen, 115
- Gitter, 119
- Gleitkommazahl, 15
  - eps, 16
  - Exponent, 16
  - Mantisse, 16
  - Maschinenoperation, 16
  - overflow, 16
  - Rundungsfehler, 16
  - underflow, 16
- Globaler Fehler, 119
- Gradientenverfahren, 51
- Grounwalls Lemma, 118
  - diskrete Version, 118
- Haarscher Raum, 103
- hermitesch, 9, 75
- Hessenberg-Matrizen, 84
- Heun-Verfahren, 122
- Hilbertraum, 6
- Householder Matrix, 37
- Householder-Matrix, 84
- Implizites Euler Verfahren, 121, 122
- Intervallschachtelung, 61
- Inverse Iteration mit Diagonal-Shift, 90
- Inverse Iteration nach Wielandt, 90
- Jaccobi-Matrix, 73
- Jacobi Verfahren, 46
  - Diagonaldominanz, 46
  - starkes Spaltensummenkriterium, 46
  - starkes Zeilensummenkriterium, 46
- Jacobi-Verfahren, 49
- Kondition, 57
- Konditioniert, 18
  - gut konditioniert, 18, 19
  - schlecht konditioniert, 18, 19
- Konditionszahlen, 18, 20
  - absolute Konditionszahl, 20
  - relative Konditionszahl, 18, 20
- Kontraktion, 10
- Konvegenzordnung
  - lineare Konvergenz, 69
  - super lineare Konvergenz, 69
- Konvergenz
  - lokale Konvergenz, 69
- Konvergenzordnung, 69, 119
- Konvexe Teilmengen, 95
- Kronecker Symbol, 37
- Landau Symbole, 12
  - $O(n)$ , 12
  - $o(n)$ , 12
- least squares, 32
- Legendre Polynome, 113
- linear abhängig, 6
- Lineare Gleichungssysteme, 21
- lineare Konvergenz, 69
- linearer Operator, 7
- Lipschitz-stetig, 7
- LR-Zerlegung, 27, 28
- Maschinenoperation, 16
- Maschinenzahlen, 15
- Matrix
  - Householder Matrix, 37
  - obere Dreiecksmatrix, 22
  - orthogonal, 35
  - regulär, 21, 23
  - singulär, 24
  - unitär, 9
  - zerlegbar, 47
- Matrixnorm, 8
- Mehrschrittverfahren, 119, 122
- Minimierungsaufgabe, 52
- Mittelpunktregel, 123
- mittlere Abweichung, 32
- Mittwersatz, 64
- Newton Verfahren, 62
  - für  $n \geq 2$ , 73
  - für mehrfache Nullstellen, 71
- Newton-Verfahren für mehrfache Nullstellen, 71
- nicht zusammenhängend, 47
- Nichtlineare Gleichungen, 61
- Nichtlineare Gleichungssysteme, 73
- Norm, 5
  - äquivalente Normen, 6
  - euklidische Norm, 7
  - induzierte Norm, 6

- Matrixnorm, 8
- Operatornorm, 8
- Spektralnorm, 9
- Normalengleichung, 32, 112
- normierter Raum
  - Hilbertraum, 6
  - Prähilbertraum, 6
- Nullstelle
  - Vielfachheit, 71
- Nullstellensuche, 61
- Operator, 7
  - beschränkt, 7
  - linear, 7
  - Lipschitz-stetig, 7
  - Matrixnorm, 8
  - Operatornorm, 8
  - Raum der beschränkten linearen, 8
  - Stetigkeit, 7
- Operatornorm, 8
- Ordnung der Nullstelle, 71
- Orthonormalsysteme, 112
- Peano
  - Peanoscher Satz, 117
- Penrose Inverse, 41
- Permutationsmatrix, 26
- Picard-Lindelöf, 117
  - lokale Version, 116
- Pivotisierung, 24
- positiv definit, 9
- Prä-Hilbertraum, 111
- Prähilbertraum, 6
- Proximum, 93
- QR-Zerlegung, 35
  - QR-Zerlegung nach Householder, 36
- Raum, 5
  - normierter Raum, 5
- Rayleigh Quotienten, 75
- Rayleighsches Maximumsprinzip, 82
- Relaxation, 59
- Relaxationsparameter, 53
- Residuenvektor, 52
- Rundungsfehler
  - absoluter Rundungsfehler, 16
  - relativer Rundungsfehler, 16
- Runge-Kutta-Verfahren, 122
- Satz von Bauer-Fike, 79
- Satz von Schur, 80
- schlecht gestellt, 20
- Schrittweitenvektor, 119
- schwache Zeilensummenkriterium, 50
- Schwaches Zeilensummenkriterium, 48
- Sekantenverfahren, 66
- Shooting-Verfahren, 125
- singuläre Werte, 39
- Skalarprodukt, 6
- SOR-Verfahren, 59
- Spaltenpivotisierung, 24
- Spektralnorm, 75
- Spektralradius, 44
- Störungssatz, 22
- starkes Spaltensummenkriterium, 46
- starkes Zeilensummenkriterium, 46
- Streng normierter Raum, 96
- submultiplikativ, 9
- superlineare Konvergenz, 69
- Taylorreihe, 11
  - Raum der stetigen diferenzierbaren Funktionen, 11
- Taylorreihe mit Integralrestterm, 11
- Taylorreihe mit Lagrange Restglied, 11
- Teilpivotisierung, 24
- Tridiagonalmatrix, 31, 48
- Tschebyscheffsches Ausgleichsproblem, 32
- Tschebyschev Polynome 1. Art, 108
- Tschebyschev Systeme, 103
- Tschebyschev-Entwicklung, 109
- Tschebyschev-Koeffizienten, 109
- Tschebyschev-Norm, 102
- Vektoriteration nach von Mises, 89
- Verfahren höher Ordnung, 70
- Verfahren in einer Raumdimension, 61
- Verfahren nach Hyman, 88
- Vollständiges ONS, 112
- Vollständigkeitsrelation, 113
- Vorkonditionierung, 59
- wohlgestellt, 20
- zerlegbar, 47
- zulässig, 117