

Skript zur Vorlesung

Numerisches Praktikum

Wintersemester 2005/2006

Universität Konstanz
Dr. Eberhard Luik

private Mitschrift

Stand: 17. Februar 2006
www.meidert.net/uni

Achtung:

Dies ist kein offizielles Skript, sondern nur eine private Mitschrift. Ich kann daher keine Gewähr für die Richtigkeit und Vollständigkeit übernehmen. Vor allem können die Nummerierungen zum Teil von den in den Vorlesungen verwendeten abweichen. Falls jemand einen Fehler entdeckt, so möge er/sie mir bitte eine eMail schicken - vielen Dank!

Frieder Meidert (uni@meidert.net)

Inhaltsverzeichnis

1	Numerisches Rechnen	1
a	Zahlen und ihre Darstellung	1
b	Operationen mit Gleitpunktzahlen	2
c	Algorithmen	4
d	Das Hornerschema	5
2	Interpolation	6
a	Problemstellung	6
b	Existenz und Eindeutigkeit	6
c	Der Neville-Algorithmus	7
d	Der Interpolationsfehler	8
e	Tschebyscheff-Polynome	9
f	Konvergenzfragen	11
g	Weitere Interpolationsarten	12
3	Nullstellenbestimmung	14
a	Bisektionsverfahren	14
b	Newton-Verfahren	14
c	Konvergenz des Newton-Verfahrens	15
d	Das Sekanten-Verfahren	19
e	Anwendung aus dem Bereich der Finanzmathematik	20
4	Lineare Gleichungssysteme	22
a	Einleitung	22
b	Gestaffelte Systeme	22
c	Die LR-Zerlegung	23
d	Die Cholesky-Zerlegung	24
e	Die QR-Zerlegung von A	26
f	Die Gauß-Elimination	28
g	Matrizentechnische Deutung der Gauß-Elimination	30

5	Das lineare Ausgleichsproblem	32
	a Problemstellung	32
	b Die Normalgleichungen	33
6	Lineare Optimierung	35
	a Problemstellung	35
	b Grafische Lösung	36
	c Normalformen	36
	d Charakterisierung des zulässigen Bereichs	37
	e Ecken des zulässigen Bereichs	38
	f Basislösungen	40
	g Lösung der linearen Optimierungsaufgabe	41
	h Der Simplex-Algorithmus	41
	i Zweiphasenmethode	43
7	Numerische Lösung von Anfangswertaufgaben	44
	a Einleitung und Beispiele	44
	b Einteilung der Näherungsverfahren	45
	c Einschnitt-Verfahren	46
	d Konsistenz und Konvergenz	48
	e Schrittweitensteuerung	51
	f Stabilität	52
8	Numerische Integration	54
	a Einleitung	54
	b Interpolatorische Quadraturformeln	55
	c Zusammengesetzte Quadraturformeln	56
	d Das Prinzip der Extrapolation	58
	e Das Romberg-Verfahren	58
	f Gauß-Quadraturformeln	60
9	Eigenwertaufgaben	63
	a Eigenwerte und Eigenvektoren	63

b	Beispiel aus der Physik: lineare Kette	63
c	Lokalisierung der Eigenwerte	65
d	Transformation auf obere Hessenberg-Gestalt	66
e	Das QR-Verfahren	67
f	Das Verfahren von Hyman	68
g	Eigenwerte und Eigenvektoren in Matlab	70
h	Vektornormen und Matrixnormen	71
i	Störempfindlichkeit bei Eigenwertaufgaben	74
10	Iterative Verfahren	77
a	Iteration	77
b	Indirekte Verfahren für lineare Gleichungssysteme	78
c	Nichtlineare Gleichungssysteme, mehrdimensionales Newton- verfahren	82
11	Minimierung	85
a	Einleitung	85
b	Nichlinearer Ausgleich	85
c	Höhenlinien	86
d	Abstiegsverfahren	86
e	Differentialgleichungsmethode	88
12	Stabilität und Störungsfragen	89
a	Einleitung	89
b	Eine allgemeine Fehlerdarstellung	90
c	Stabilität und Störuneempfindlichkeit einer Nullstelle	91
d	Störempfindlichkeit linearer Gleichungssysteme	92

1. Numerisches Rechnen

a. Zahlen und ihre Darstellung

Sei $x \in \mathbb{R} \setminus \{0\}$, $g \in \mathbb{N}$, $g \geq 2$.

$$x = \sigma \cdot g^N \cdot \sum_{v=1}^{\infty} x_v \cdot g^{-v} \quad (1)$$

mit $\sigma \in \{1, -1\}$ (Vorzeichen), $N \in \mathbb{Z}$, $x_v \in \{0, 1, \dots, g-1\}$

Zusatzforderung:

- $x_1 \neq 0$
- zu jedem $v \in \mathbb{N}$ gibt es ein $n > v$ mit $x_n \neq g-1$.

Damit wird die Darstellung 1 eindeutig (normalisierte Darstellung).

Schreibweise $\sigma \cdot 0, x_1, x_2, x_3 \dots g^N$ (Gleitkommazahl)

BEISPIELE:

$g = 10$: Dezimalsystem: Ziffern $0, \dots, 9$

$g = 2$: Dualsystem: Ziffern $0, 1$

$g = 16$: Hexadezimalsystem: Ziffern $0, 1, \dots, 9, A, B, C, D, E, F$

$$\text{z.B. } x = 11,1875 \begin{cases} \text{dezimal} & 0,111875 \cdot 10^2 \\ \text{dual} & 0,10110011 \cdot 2^4 \\ \text{hexadezimal} & 0,B3 \cdot 16^1 = 16 \cdot \left(\frac{11}{16} + \frac{3}{256}\right) \end{cases}$$

$$11,1875 = 1 \cdot 8 + 0 \cdot 4 + 1 \cdot 2 + 1 \cdot 1 + 0 \cdot \frac{1}{2} + 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{8} + 1 \cdot \frac{1}{16} = 16 \cdot \left(1 \cdot \frac{1}{2} + 0 \cdot \frac{1}{4} \dots\right)$$

Computer: $g = 2$

Darstellbar sind dann die 0 und Zahlen der Form $\sigma \cdot 2^N \cdot \sum_{v=1}^t x_v 2^{-v}$ mit $\sigma \in \{1, -1\}$, $N^- \leq N^+$, $x_{nu} \in \{0, 1\}$, $x_1 = 1$ und einem festen t .

z.B. 64 Bit: 1 Bit: Vorzeichen, 11 Bit Hochzahl, 52 Bits: Mantisse (x_1, \dots, x_{52}) (IEEE-Standard)

DEFINITION:

Die Menge M der Zahlen der Form (2) heißt **Maschinenzahlen** (endliche Menge!)

Beispiel: nicht von der Form (2): 0,1

Ist x eine Maschinenzahl $\neq 0$, so gilt:

$$g^{N^- - 1} \leq |x| < g^{N^+}$$

$y \in \mathbb{R}$ mit $|y| \geq g^{N^+}$: Exponentenüberlauf

$y \in \mathbb{R}$ mit $|y| < g^{N^- - 1}$: wird (in der Regel) durch Null ersetzt.

Beispiel: $\sum_{n=1}^{\infty} \frac{1}{n}$ konvergiert im PC!

b. Operationen mit Gleitpunktzahlen

Sei $x \in \mathbb{R}$ mit $|y| < g^{N^+}$, $x = \sigma \cdot g^N \cdot \sum_{v=1}^{\infty} \cdot g^{-v}$

$$rd(x) := \begin{cases} \sigma \cdot g^N \cdot \sum_{v=1}^t x_v g^{-v} & , \text{ falls } x_{t+1} < \frac{g}{2} \text{ abrunden} \\ \sigma \cdot g^N \cdot \left(\sum_{v=1}^t x_v g^{-v} + g^{-t} \right) & , \text{ falls } x_{t+1} \geq \frac{g}{2} \text{ aufrunden} \end{cases}$$

Sei \bar{x} ein Näherungswert.

DEFINITION:

- (a) $x - \bar{x}$ heißt **absoluter Fehler**,
- (b) $\frac{x - \bar{x}}{x}$ heißt **relativer Fehler** (für $x \neq 0$).

SATZ:

- (a) $|rd(x) - x| \leq 0,5g^{N-t}$,
- (b) $\left| \frac{rd(x) - x}{x} \right| \leq 0,5g^{-t+1} (2^{-t})$ im Zweiersystem

DEFINITION:

Die Zahl $\tau := 0,5 \cdot g^{-t+1}$ heißt **relative Rechengenauigkeit**.

$$1 + \frac{1}{2}$$

$$1 + \frac{1}{2^2}$$

⋮

$$1 + \frac{1}{2^{52}} > 1$$

$$1 + \frac{1}{2^{-53}} = 1 \quad (2^{-53} \text{ kann noch exakt dargestellt werden})$$

⇒ $2^{-52} \cong$ relative Rechengenauigkeit

$$g^{N^-} < 2^{-53}$$

Setze $\varepsilon := \frac{rd(x)-x}{x} \Rightarrow rd(x) = x(1 + \varepsilon)$ mit $|\varepsilon| \leq \tau$.

VERKNÜPFUNG VON GLEITKOMMAZAHLEN:

$$\diamond \in \{+, \cdot, -, /\}$$

$x, y \in M$: dann ist im Allgemeinen $x \diamond y \notin M$

⇒ Runden: $rd(x \diamond y) = (x \diamond y)(1 + \varepsilon_1)$ mit $|\varepsilon| \leq \tau$.

BEISPIEL: :

$$g := 19, t := 3$$

$$x = 0,123 \cdot 10^6$$

(Maschinenzahl)

$$y = 0,456 \cdot 10^2$$

$$x + y = 0,1230456 \cdot 10^6 \notin M$$

$$rd(x + y) = 0,123 \cdot 10^6 = x$$

$x, y \in \mathbb{R}$:

$$rd(x) = x(1 + \varepsilon_1) \text{ mit } |\varepsilon_1| \leq \tau$$

$$rd(y) = y(1 + \varepsilon_2) \text{ mit } |\varepsilon_2| \leq \tau$$

$$rd(x) + rd(y) = x(1 + \varepsilon_1) + y(1 + \varepsilon_2) = x + \varepsilon_1x + y + \varepsilon_2y$$

$$rd(x) + rd(y) - (x + y) = \varepsilon_1x + \varepsilon_2y$$

(absoluter Fehler)

$$\frac{rd(x)+rd(y)-(x+y)}{(x+y)} = \frac{\varepsilon_1x}{x+y} + \frac{\varepsilon_2y}{x+y}$$

(relativer Fehler)

Fall $\text{sign}(x) = \text{sign}(y)$:

$$\left| \frac{rd(x)+rd(y)-(x+y)}{(x+y)} \right| \leq \left| \frac{\varepsilon_1x}{x+y} \right| + \left| \frac{\varepsilon_2y}{x+y} \right| = |\varepsilon_1| \underbrace{\left| \frac{x}{x+y} \right|}_{<1} + |\varepsilon_2| \underbrace{\left| \frac{y}{x+y} \right|}_{<1} \leq |\varepsilon_1| + |\varepsilon_2| \leq 2\tau, \text{ also unpro-}$$

blematisch!

Fall $x \approx -y \Rightarrow$ großer relativer Fehler.

ACHTUNG: Problematisch ist die Subtraktion zweier ungefähr gleicher Zahlen. (siehe Aufgabe I, 1a)

$$rd(rd(x) \diamond rd(y)) = x \diamond y + F.$$

$$rd(rd(x) + rd(y)) = (x(1 + \varepsilon_1) + y(1 + \varepsilon_2))(1 + \varepsilon_3) = \dots \leq x + y + (|x| + |y|) \cdot 2\tau(|x| + |y|) \cdot \tau^2 = x + y + F \text{ mit } |F| \leq (|x| + |y|)(\tau^2 + 2\tau).$$

ACHTUNG: Im Computer kann man nicht beliebig klammern:

$$(x + y) + z \neq x + (y + z) \text{ (wegen der Rundung)}$$

c. Algorithmen

Algorithmus: Beschreibung einer Methode zur Lösung einer gegebenen Aufgabenstellung

- Bestimmtheit
- Endlichkeit (nach endlich vielen Schritten am Ziel)
- Allgemeingültigkeit

Auswahlkriterien unter Algorithmen: Rechenzeit, Stabilität.

DEFINITION:

Sei A ein Algorithmus. Die Abbildung $T_A : \mathbb{N} \rightarrow \mathbb{N}$, welche jeder Anzahl von Eingabedaten die Anzahl der Rechenoperationen zuordnet, heißt **Komplexität von A** .

BEISPIEL:

$$1) T_A(n) = \frac{1}{2}n(n + 3) = \frac{1}{2}n^2 + \frac{3}{2}n = \sigma(n^2) \text{ (für } n \rightarrow \infty)$$

$$2) p(t) = a_0 + a_1t + \dots + a_nt^n, \xi \in \mathbb{R}, \text{ Berechne } p(\xi).$$

A_1 : Eingabe $n, (a_0, \dots, a_n), \xi$

$p := a_0$.

Für $i := 1, \dots, n$:

$b = a_i$

Für $j = 1, \dots, i$:

$b = b \cdot \xi$

{Ende der j-Schleife}

$p = p + b$

{Ende der i-Schleife}

Ausgabe p .

n Additionen, $(1 + 2 + \dots + n) = \frac{n(n+1)}{2}$ Multiplikationen

$$\Rightarrow T_{A_1}(n) = \frac{1}{2}n^2 + \frac{3}{2}n = \sigma(n^2).$$

d. Das Hornerschema

$p(t) = a_0^{(0)} + a_1^{(0)}t + \dots + a_n^{(0)}t^n; \xi \in \mathbb{R};$ gesucht: $p(\xi)$

Setze $a_n^{(1)} = a_n^{(0)}$

A_2 :

Für $i = 1, \dots, n$:

$$a_{n-1}^{(1)} := a_{n-i}^{(0)} + a_{n+1-i}^{(1)} \cdot \xi.$$

Behauptung: $a_0^{(1)} = p(\xi)$ (vgl. Übungsaufgaben)

n Additionen, n Multiplikationen

$$T_{A_2}(n) = 2n = \sigma(n).$$

2. Interpolation

a. Problemstellung

Gegeben seien $(t_i, s_i) \in \mathbb{R}^2$, $i = 0, \dots, m$, t_0, \dots, t_m paarweise verschieden.

Gesucht ist ein $p(t) := \sum_{j=0}^m a_j t^j$ mit $p(t_i) = s_i$ ($i = 0, \dots, m$) (1)

Interpolationsaufgabe nach Lagrange:

Die (t_i, s_i) heißen **Stützpaare**.

Anwendung:

$f \in C([a, b]; \mathbb{R})$, $a \leq t_0 < t_1 < \dots < t_m \leq b$, $s_i := f(t_i)$

Gesucht: Interpolationspolynom.

Fall $m = 1$: klar (Gerade)

Fall $m = 2$: (Parabel)

b. Existenz und Eindeutigkeit

SATZ:

Die Interpolationsaufgabe nach Lagrange besitzt genau eine Lösung.

BEWEIS:

(i) Eindeutigkeit:

Seien p und q Lösungen der Interpolationsaufgabe nach Lagrange.

Bezeichnung: Π_m sei der Raum der Polynome vom Grad $\leq m$. D.h. $p, q \in \Pi_m$, und erfüllen (1).

Setze $r(t) := p(t) - q(t)$. Dann gilt $r \in \Pi_m$. Es ist $r(t_i) = p(t_i) - q(t_i) = 0$ für $i = 0, \dots, m \Rightarrow r$ hat mindestens $m + 1$ paarweise verschiedene Nullstellen $\Rightarrow r \equiv 0 \Rightarrow p = q$.

(ii) Existenz:

Wir setzen $l_i(t) := \prod_{j=0, j \neq i}^m \left(\frac{t-t_j}{t_i-t_j} \right)$ für $i = 1, \dots, m$, wohldefiniert, da Nenner stets $\neq 0$.

Es gilt $l_i \in \Pi_m$ für $i = 1, \dots, m$

$$l_i(t_k) = \begin{cases} 0 & , i \neq k \\ 1 & , i = k \end{cases}$$

Die l_i heißen **Lagrange-Grundpolynome**.

$$\boxed{p(t) := \sum_{i=0}^m s_i l_i(t)} ; p \in \Pi_m \quad (2)$$

Es gilt $p(t_k) = \sum_{i=0}^m s_i l_i(t_k) = s_k l_k(t_k) = s_k$ für $k = 0, \dots, m$. □

(2) heißt **Interpolationsformel nach Lagrange**.

BEMERKUNG:

Das Interpolationspolynom ist abhängig von der Stützstellenwahl.

Gegeben ein ξ , $p(\xi) = ?$.

c. Der Neville-Algorithmus

Seien $(t_i, s_i) \in \mathbb{R}^2$ mit t_i paarweise verschieden für $i = 0, \dots, m$. $p(t)$ sei das Interpolationspolynom nach Lagrange.

Gesucht $p(\xi)$ für ein gegebenes $\xi \in \mathbb{R}$.

$p_{m-1}(t)$ sei das Interpolationspolynom zu den Stützpaaren $(t_0, s_0), \dots, (t_{m-1}, s_{m-1})$

$q_{m-1}(t)$ sei das Interpolationspolynom zu den Stützpaaren $(t_1, s_1), \dots, (t_m, s_m)$

LEMMA:

Es gilt:

$$p(t) = \frac{(t-t_m)p_{m-1}(t) - (t-t_0)q_{m-1}(t)}{t_0-t_m} =: r(t)$$

BEWEIS:

$r \in \Pi_m$ klar.

$$r(t_i) = \frac{(t_i-t_m)p_{m-1}(t_i) - (t_i-t_0)q_{m-1}(t_i)}{t_0-t_m} = \frac{(t_i-t_m)s_i - (t_i-t_0)s_i}{t_0-t_m} = s_i \text{ für } 1 \leq i \leq m-1,$$

$$r(t_0) = \frac{(t_0-t_m)s_0 - 0}{t_0-t_m} = s_0,$$

$$r(t_m) = \frac{-(t_m-t_0)s_m}{t_0-t_m} = s_m.$$

$$\Rightarrow r(t) = p(t)$$

□

BEZEICHNUNG:

$p_{i_0, i_1, \dots, i_k} \in \Pi_k$ sei das Interpolationspolynom zu den Stützpaaren $(t_{i_0}, s_{i_0}), (t_{i_1}, s_{i_1}), \dots, (t_{i_k}, s_{i_k})$.

d. Der Interpolationsfehler

$f : [a, b] \rightarrow \mathbb{R}; a \leq t_0 < t_1 < \dots < t_m \leq b$ ($s_i = f(t_i)$)
 $p(t)$ sei das Interpolationspolynom zu f .

SATZ:

Sei $f \in C^{m+1}[a, b]$, sei $z \in [a, b]; \alpha := \min\{z, t_0\}, \beta := \max\{z, t_m\}$.
 Es gibt ein $\xi \in (a, b)$ mit

$$f(z) - p(z) = \frac{f^{(m+1)}(\xi)}{(m+1)!} \underbrace{\prod_{i=0}^m (z - t_i)}_{=: w(z)}$$

BEWEIS:

$$\text{Setze } \Phi(t) := f(t) - p(t) - \frac{w(t)}{w(z)}(f(z) - p(z))$$

$$\Phi^{(m+1)}(t) = f^{(m+1)}(t) - 0 - \frac{(m+1)!}{w(z)}(f(z) - p(z))$$

Für $z \in \{t_0, \dots, t_m\}$ klar.Sei deshalb $z \notin \{t_0, \dots, t_m\}$

$$\Phi(t) := f(t) - p(t) - \frac{w(t)}{w(z)}(f(z) - p(z))$$

$$\Phi^{(m+1)}(t) = f^{(m+1)}(t) - \frac{(m+1)!}{w(z)}(f(z) - p(z))$$

Sei ξ mit $\Phi^{(m+1)}(\xi) = 0$.

$$0 = f^{(m+1)}(\xi) - \frac{(m+1)!}{w(z)}(f(z) - p(z)) \Rightarrow \text{Behauptung.}$$

$$\Phi(t_k) = \underbrace{f(t_k)}_{=s_i} - \underbrace{p(t_k)}_{=s_i} - \overbrace{\frac{w(t_k)}{w(z)}(f(z) - p(z))}^{=0} = 0 \text{ für } k = 0, \dots, m.$$

$$\Phi(z) = 0$$

Φ besitzt in $[\alpha, \beta]$ mindestens $m + 2$ Nullstellen.

Φ' besitzt in (α, β) mindestens $m + 1$ Nullstellen.

\vdots

$\Phi^{(m+1)}$ besitzt in (α, β) mindestens eine Nullstelle. □

Fehlerschranken:

zu $g \in C[a, b]$ sei $\|g\| := \max\{|g(t)| : t \in [a, b]\}$.

$$|f(z) - p(z)| \leq \frac{\|f^{(m+1)}\|}{(m+1)!} |w(z)|$$

$$\|f - p\| \leq \frac{\|f^{(m+1)}\|}{(m+1)!} \|w\|$$

$\|w\|$ ist abhängig von der Stützstellenwahl t_0, \dots, t_m

→ günstige Stützstellenwahl??

(d.h. so, dass $\|w\|$ minimal wird)

$$w(t) = t^{m+1} + a_m t^m + \dots + a_0 \quad Q_n := \{p \in \Pi_n \mid p(t) = t^n + a_{n-1} + \dots + a_0\}.$$

e. Tschebyscheff-Polynome

Zu jedem $n \in \mathbb{N}$ definiert man

$$T_n : [-1, 1] \rightarrow \mathbb{R}, T_n(t) := \cos(n \cdot \arccos(t))$$

LEMMA:

Für $n \geq 1$ gilt:

$$T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t)$$

$$T_{n+1}(t) = \cos((n+1) \cdot \arccos(t)) =$$

$$= \cos(n \cdot \arccos(t)) \cdot \cos(\arccos(t)) - \sin(n \cdot \arccos(t)) \cdot \sin(\arccos(t))$$

$$T_{n-1}(t) = \cos((n-1) \cdot \arccos(t)) =$$

$$= \cos(n \cdot \arccos(t)) \cdot \cos(\arccos(t)) + \sin(n \cdot \arccos(t)) \cdot \sin(\arccos(t))$$

$$T_{n+1} + T_{n-1} = 2 \cdot \cos(n \cdot \arccos(t)) \cdot t = 2tT_n(t) \quad \square$$

SATZ:

Damit sind wegen $T_0 = 1$ und $T_1 = t$ alle $T_n, n \in \mathbb{N}$ Polynome vom Grad n .

$$T_n(t) = 2^{n-1}t^n + a_{n-1}t^{n-1} + \dots + a_0.$$

BEWEIS:

Induktion mit vorangegangenem Lemma.

BEZEICHNUNG:

Die T_n heißen **Tschebyscheff-Polynome 1. Art**.

$$T_0(t) = 1, T_1(t) = t, T_2(t) = 2t^2 - 1, T_3(t) = 4t^3 - 3t.$$

EIGENSCHAFTEN:

- (1) Symmetrie: $T_n(-t) = (-1)^n T_n(t)$,
- (2) $T_n(1) = 1, T_n(-1) = (-1)^n, |T_n(t)| \leq 1 \forall t \in [-1, 1] \Rightarrow \|T_n\|_{[-1,1]} = 1$,
- (3) Nullstellen: T_n hat n einfache Nullstellen in $[-1, 1]$:
 $t_i = \cos\left(\frac{2i+1}{2n} \cdot \pi\right)$ für $i = 0, \dots, n-1$,
- (4) Extremstellen: T_n hat in $[-1, 1]$ genau $(n+1)$ Extremstellen:
 $z_i := \cos\left(\frac{i}{n}\pi\right)$ für $i = 0, 1, \dots, n$
 $T_n(z_i) = (-1)^i$ für $i = 0, \dots, n$,
- (5) Minimalitätseigenschaft: sei $\widehat{T}_n := \frac{T_n}{2^{n-1}} \Rightarrow \widehat{T}_n(t) = t^n + c_{n-1}t^{n-1} + \dots$, d.h. $\widehat{T}_n \in \mathbb{Q}_n$

SATZ:

Unter allen Polynomen aus \mathbb{Q}_n hat \widehat{T}_n die minimale Norm.

BEWEISIDEE:

Sei $q \in \mathbb{Q}_n$ mit $\|q\| < \|\widehat{T}_n\|$

$$r(t) := \widehat{T}_n(t) - q(t) \Rightarrow r \in \Pi_{n-1}$$

$$r(z_0) = \widehat{T}_n(z_0) - q(z_0) > 0$$

$$r(z_1) = \widehat{T}_n(z_1) - q(z_1) < 0$$

$$r(z_2) = \widehat{T}_n(z_2) - q(z_2) > 0$$

⋮

$\Rightarrow r$ hat Grad $\leq n-1$ mit mindestens n Nullstellen, da $r \neq 0$ (da $\|q\| < \|\widehat{T}_n\|$) ist dies ein Widerspruch.

Wähle bei der Interpolation die Nullstellen t_0, \dots, t_m von $T_{n+1}(t)$.

Dann gilt (bezüglich $[-1, 1]$):

$$\|f - p\|_{[-1,1]} \leq \frac{\|f^{(m+1)}\|_{(-1,1)}}{(m+1)!} \cdot \frac{1}{2^m} \cdot c.$$

f. Konvergenzfragen

Sei $f \in C([a, b])$; zu $m \in \mathbb{N}$ seien $a \leq t_{m_0} < t_{m_1} < \dots < t_{m_m} \leq b$.

Weiter sei $p_m(t)$ das Interpolationspolynom zu f bezüglich t_{m_0}, \dots, t_{m_m} .

FRAGE:

- (1) Gilt gleichmäßige Konvergenz? Also $\|p_m - f\| \rightarrow 0$ für $m \rightarrow \infty$?
- (2) Gilt punktweise Konvergenz? Also $p_m(x) \rightarrow f(x)$ für alle $x \in [a, b]$?

Wir ordnen die Stützstellen in einer Matrix an:

$$S = \begin{pmatrix} t_{00} & & & & \\ & t_{11} & & & \\ & \vdots & & \ddots & \\ & t_{m0} & \dots & \dots & t_{mm} \\ & \vdots & & & \end{pmatrix} \text{- Stützstellenmatrix}$$

SATZ: (Faber)

Zu jeder Stützstellenmatrix S gibt es ein $f \in C[a, b]$, so dass die Folge $(p_m)_{m \in \mathbb{N}}$ der Interpolationspolynome nicht gleichmäßig konvergiert.

SATZ: (Marcinkiewicz)

Zu jedem $f \in C[a, b]$ gibt es eine Stützstellenmatrix S , so dass die Folge $(p_m)_{m \in \mathbb{N}}$ der Interpolationspolynome gleichmäßig gegen f konvergiert.

Aber: S ist im Allgemeinen nicht bekannt!

BEISPIEL:

$f(x) = \exp(x)$ im Intervall $[a, b]$.

Sei $p_m(x)$ das Interpolationspolynom zu f bezüglich $a \leq t_{m_0} < \dots < t_{m_m} \leq b$.

BEHAUPTUNG: $\|p_m - f\| \rightarrow 0$ für $m \rightarrow \infty$.

BEWEIS:

Für den Interpolationsfehler gilt $\|p_m - f\|_{[a,b]} \leq \frac{\|f_{[a,b]}^{(m+1)}\|}{(m+1)!} \|w\|_{[a,b]}$.

$$|w(t)| := \prod_{i=0}^m |(t - t_i)| \leq (b - a)^{m+1} \Rightarrow \|w\|_{[a,b]} \leq (b - a)^{m+1}$$

$$\|f^{(m+1)}\|_{[a,b]} = \|\exp\|_{[a,b]} = e^b.$$

Somit $\|p_m - f\|_{[a,b]} \leq \frac{e^b (b-a)^{m+1}}{(m+1)!} \rightarrow 0$ für $m \rightarrow \infty$ □

g. Weitere Interpolationsarten

HERMITE INTERPOLATION:

Gegeben seien $(t_i, s_i^{(0)}, s_i^{(1)}) \in \mathbb{R}^3$, t_i paarweise verschieden.

Gesucht ist ein $p \in \Pi_{2m+1}$ mit $(i = 0, \dots, m)$

- $p(t_i) = s_i^{(0)}$ für $i = 0, \dots, m$;
- $p'(t_i) = s_i^{(1)}$ für $i = 0, \dots, m$.

SATZ:

Es gibt genau ein solches Polynom.

BEWEIS:

Eindeutigkeit: ähnlich wie bei Lagrange.

Existenz:

$$p(t) = \sum_{k=0}^m s_k^{(0)} l_k^2(t) [1 - 2l_k'(t_k)(t - t_k)] + \sum_{k=0}^m s_k^{(1)} l_k^2(t)(t - t_k)$$

Verallgemeinerung:

auch höhere Ableitungen \rightarrow verallgemeinertes Hermite-Interpolationspolynom.

RATIONALE INTERPOLATION:

$$R_{m,n}(t) := \frac{a_0 + a_1 t + \dots + a_m t^m}{b_0 + b_1 t + \dots + b_n t^n}$$

Gegeben seien $m + n + 1$ Paare $(t_i, s_i) \in \mathbb{R}^2$, t_i paarweise verschieden.

Gesucht ist eine rationale Funktion $R_{m,n}$ mit

$$R_{m,n}(t_i) = s_i \text{ für } i = 0, \dots, m + n$$

ACHTUNG: Existenz nicht immer gesichert!

TRIGONOMETRISCHE INTERPOLATION:

Eine Funktion der Form

$$s_m(t) := a_0 + \sum_{k=1}^m a_k \cos(k \cdot t) + \sum_{k=1}^m b_k \sin(k \cdot t)$$

heißt trigonometrisches Polynom vom Grad m

Gegeben seien $2m + 1$ Paare $(t_i, s_i) \in \mathbb{R}^2$ mit $-\pi \leq t_0 < \dots < t_m < \pi$.
Gesucht ist ein $S_m(t)$ mit

$$S_m(t_i) = s_i \text{ für } i = 0, \dots, 2m$$

Es gibt genau ein solche S_m .

SPLINE-INTERPOLATION:

Sei $a_0 = t_0 < t_1 \dots < t_m = b$

DEFINITION:

Unter einer (kubischen) **Splinefunktion** versteht man eine Funktion $h : [a, b] \rightarrow \mathbb{R}$ mit folgenden Eigenschaften:

- 1) $h \in C^2[a, b]$;
- 2) auf jedem Teil-Intervall $[t_i, t_{i+1}]$ gilt $g|_{[t_i, t_{i+1}]} \in \Pi_3$.

Sei nun $(t_i, s_i) \in \mathbb{R}$ für $i = 0, \dots, m$

Gesucht ist eine kubische Splinefunktion mit $h(t_i) = s_i$ für $i = 0, \dots, m$.

Es gibt eine Lösung, aber keine eindeutige Lösung.

Zusatzforderung:

$$h^{(2)}(t_0) = h^{(2)}(t_m) = 0.$$

3. Nullstellenbestimmung

a. Bisektionsverfahren

Sei $f \in C[a, b]$. Gesucht: $\xi \in [a, b] = I_0$ mit $f(\xi) = 0$.

Es gelte $f(a) \cdot f(b) < 0$.

Zwischenwertsatz \Rightarrow es gibt ein $\xi \in (a, b)$ mit $f(\xi) = 0$.

Bilde $s = \frac{a+b}{2}$.

Dann gilt genau einer der folgenden 3 Fälle:

$$(1) f(s) = 0 \Rightarrow s = \xi;$$

$$(2) f(a) \cdot f(s) < 0 \Rightarrow \xi \in (a, s) =: I_1;$$

$$(3) f(b) \cdot f(s) < 0 \Rightarrow \xi \in (s, b) =: I_1.$$

Fortgesetzte Intervall-Halbierung: I_n

Sei $u_n \in I_n$: dann gilt $|u_n - \xi| \leq \frac{b-a}{2^n} \rightarrow 0$ für $n \rightarrow \infty$.

\Rightarrow Eindeutigkeit (**natürlicher kubischer Spline**)

b. Newton-Verfahren

HERLEITUNG:

$f \in C^2[a, b]$, sei $\xi \in [a, b]$ mit $f(\xi) = 0$ und $f'(t) \neq 0$ in $[a, b]$.

Sei $x_0 \in [a, b]$ ein Näherungswert an ξ .

Ziel: Berechnung eines besseren Näherungswertes x_1 .

Taylor:

$$0 = f(\xi) = f(x_0) + f'(x_0)(\xi - x_0) + \underbrace{\frac{f^{(2)}(\eta)}{2!}(\xi - x_0)^2}_{\text{vernachlässigbar}}$$

$$\Rightarrow 0 \approx f(x_0) + f'(x_0)(\xi - x_0)$$

$$\Rightarrow 0 = f(x_0) + f'(x_0)(x_1 - x_0)$$

$$\Rightarrow x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$$

Wiederhole diesen Vorgang:

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)} \text{ für } i = 0, 1, \dots$$

Diese Rechenvorschrift heißt **Newton-Verfahren**, x_0 heißt **Startwert**.

FRAGEN:

- durchführbar? (also $x_i \in [a, b]$ für alle i ?)
- Konvergenz gegen ξ ?
- Konvergenzgeschwindigkeit?
- Wahl des Startwertes?

Iterationsverfahren:

$\Phi : [a, b] \rightarrow [a, b]$ stetig.

Zu einem Startwert $x_0 \in [a, b]$ berechnet man

$$x_{i+1} = \Phi(x_i) \quad (i = 0, 1, \dots)$$

(Φ heißt dann **Iterationsfunktion**)

Ein $\xi \in [a, b]$ heißt Fixpunkt von Φ , falls $\Phi(\xi) = \xi$ gilt.

FRAGE:

$\lim_{i \rightarrow \infty} x_i = \xi$?

→ Banachscher Fixpunktsatz.

Newton-Verfahren: $\Phi(t) = t - \frac{f(t)}{f'(t)}$

c. Konvergenz des Newton-Verfahrens

DEFINITION:

- (1) Konvergiert die Folge $(x_i)_{i \in \mathbb{N}}$ für jeden Startwert $x_0 \in [a, b]$, so heißt das Iterationsverfahren **global konvergent**.

- (2) Konvergiert die Folge $(x_i)_{i \in \mathbb{N}}$ nur für alle Startwerte $x_0 \in U$ (U eine Umgebung um ξ), so heißt das Verfahren **lokal konvergent**.

DEFINITION:

Sei U eine Umgebung um ξ . Gilt für alle Startwerte $x_0 \in U$ für die Iterationsfolge eine Ungleichung der folgenden Form:

$$|x_{i+1} - \xi| \leq C|x_i - \xi|^p$$

($p \geq 1$, C eine Konstante - falls $p = 1$, so gelte $C < 1$ - , unabhängig von x_0)
so nennt man das Verfahren ein Verfahren der Konvergenzordnung p .

$p = 1$: lineare Konvergenz

$p = 2$: quadratische Konvergenz

SATZ:

Sei $f \in C^2[a, b]$; $\xi \in [a, b]$ mit $f(\xi) = 0$ und $f'(t) \neq 0$ für alle $t \in [a, b]$.

Das Newton-Verfahren liefert für jeden Startwert $x_0 \in [a, b]$ die Ungleichung

$$|x_1 - \xi| \leq C|x_0 - \xi|^2$$

mit einer von x_0 unabhängigen Konstanten $c > 0$.

BEWEIS:

$$0 = f(\xi) = f(x_0) + f'(x_0)(\xi - x_0) + \frac{f^{(2)}(\eta)}{2!}(\xi - x_0)^2$$

$$\Rightarrow \xi - \underbrace{\left(x_0 + \frac{f(x_0)}{f'(x_0)}\right)}_{=x_1} = -\frac{f^{(2)}(\eta)}{2!f'(x_0)}(\xi - x_0)^2.$$

$$\Rightarrow |\xi - x_1| = \frac{|f^{(2)}(\eta)|}{2|f'(x_0)|}|\xi - x_0|^2$$

$$C_1 := \max\{|f^{(2)}(t)| : t \in [a, b]\} < \infty$$

$$C_2 := \min\{|f'(t)| : t \in [a, b]\} > 0$$

$$|\xi - x_1| \leq \underbrace{\frac{C_1}{2C_2}}_{=:C} |\xi - x_0|^2 = C|\xi - x_0|^2 \quad \square$$

SATZ:

Sei $f \in C^2[a, b]$; $\xi \in (a, b)$ mit $f(\xi) = 0$ und $f'(\xi) \neq 0$. Dann gilt:

- (1) Das Newton-Verfahren besitzt Konvergenzordnung 2
- (2) Es gibt eine Umgebung U von ξ so, dass für jeden Startwert $x_0 \in U$ das Newton-Verfahren gegen ξ konvergiert.

BEWEIS:

Sei C wie oben, sei $\varepsilon := \min\{\frac{1}{2C}, |\xi - a|, |\xi - b|\} > 0$

$U := (\xi - \varepsilon, \xi + \varepsilon)$

Sei $x_0 \in U \Rightarrow |x_1 - \xi| \leq C \underbrace{|x_0 - \xi|^2}_{< \varepsilon} < C \cdot \varepsilon^2 \leq C \cdot \frac{1}{2C} \cdot \varepsilon = \frac{\varepsilon}{2}$

$\Rightarrow x_1 \in U \Rightarrow x_i \in U$ für $i = 1, 2, \dots$

\Rightarrow Konvergenzordnung 2.

Es gilt

$$|x_{i+1} - \xi| \leq C|x_i - \xi|^2 < \left(\frac{1}{2}\right)^{i+1} \cdot \varepsilon \rightarrow 0 \text{ für } i \rightarrow \infty.$$

Sei $f : [a, b] \rightarrow \mathbb{R}$ mit

(V1) es gibt ein $\xi \in [a, b]$ mit $f(\xi) = 0$,

(V2) $f \in C^1[a, b]$ und $f'(t) \neq 0$ für alle $t \in [a, b]$,

(V3) $f'(s)(t - s) + f(s) \leq f(t)$ für alle $s, t \in [a, b]$,

(V4) Für jeden Startwert $x_0 \in [a, b]$ gilt $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$.

LEMMA:

Sei $f \in C^2[a, b]$ mit $f^{(2)}(t) \geq 0$ für alle $t \in [a, b]$. Dann folgt daraus (V3).

BEWEIS:

$$f(t) = f(s) + f'(s)(t - s) + \underbrace{\frac{f^{(2)}(\eta)}{2} (t - s)^2}_{\geq 0} > f(s) + f'(s)(t - s) \quad \square$$

SATZ: (monotone Konvergenz)

Seien (V1) - (V4) erfüllt und $x_0 \in [a, b]$ ein beliebiger Startwert.

Dann gilt:

- (i) $\xi \leq x_{i+1} \leq x_i$ für $i = 1, 2, 3, \dots$ und $f'(t) > 0, t \in [a, b]$,
- (ii) $x_i \leq x_{i+1} \leq \xi$ für $i = 1, 2, 3, \dots$ und $f'(t) < 0, t \in [a, b]$,
- (iii) $\lim_{i \rightarrow \infty} x_i = \xi$.

BEWEIS:

Es gibt genau eine Nullstelle.

- (i) Sei $f'(t) > 0$ in $[a, b]$.

$$x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} \Leftrightarrow 0 = f'(x_0)(x_1 - x_0) + f(x_0) \leq f(x_1).$$

Falls $x_1 = \xi \Rightarrow x_i = \xi$ für $i = 1, 2, \dots$

Sei $x_1 \neq \xi$. Nachdem Mittelwertsatz folgt:

$$\frac{\underbrace{f(x_1) - f(\xi)}_{>0}}{x_1 - \xi} \stackrel{=0}{=} \underbrace{f'(\eta)}_{>0} \Rightarrow x_1 - \xi > 0 \Rightarrow x_1 > \xi.$$

Induktionsanfang: ✓

Es gelte $\xi \leq x_n \leq x_{n-1} \leq \dots \leq x_1$ mit $f(x_n) \geq 0$

$$0 = f'(x_n)(x_{n+1} - x_n) + f(x_n) \leq f(x_{n+1}) \quad (\text{nach (V3)})$$

\Rightarrow mit dem Mittelwertsatz $\frac{f(x_{n+1}) + f(\xi)}{x_{n+1} - \xi} = f''(\eta) \Rightarrow$ (wie oben) $\xi \leq x_{n+1}$.

$$x_{n+1} = x_n - \underbrace{\frac{f(x_n)}{f'(x_n)}}_{>0} \leq x_n.$$

Induktionsschritt: ✓

- (ii) Die Folge $(x_i)_{i \in \mathbb{N}}$ ist monoton fallend und beschränkt und damit konvergent.

$$\lim_{i \rightarrow \infty} x_i =: \gamma \quad \Phi(t) = t - \frac{f(t)}{f'(t)} \text{ stetig.}$$

$$\Phi(\gamma) = \Phi(\lim_{i \rightarrow \infty} x_i) = \lim_{i \rightarrow \infty} \Phi(x_i) = \lim_{i \rightarrow \infty} x_{i+1} = \gamma.$$

Also ist γ eine Nullstelle von $f \Rightarrow \gamma = \xi$ (da f streng monoton wachsend).

ANWENDUNG:

Berechnung von \sqrt{a} für $a > 1$

$$\begin{aligned} f(t) &= t^2 - a \\ f'(t) &= 2t \\ f''(t) &= 2 \end{aligned}$$

$$\begin{aligned} 1 &< \sqrt{a} < a \\ f'(t) &> 0 \text{ im Intervall } [1, a] \\ f''(t) &> 0 \text{ im Intervall } [1, a] \end{aligned}$$

$$\begin{aligned} x_1 &= x_0 - \frac{f(x_0)}{f'(x_0)} = x_0 - \frac{x_0^2 - a}{2x_0} = \frac{x_0}{2} + \frac{a}{2x_0} \\ &\leq \frac{a}{2} + \frac{a}{2} = a. \\ \text{bzw. } &\geq \frac{1}{2} + \frac{1}{2} = 1. \end{aligned}$$

$\Rightarrow x_1 \in [1, a]$ für jeden Startwert $x_0 \in [1, a]$.

\Rightarrow (V1) - (V4) sind erfüllt. Für jedes $x_0 \in [1, a]$ konvergiert also das Newton-Verfahren gegen \sqrt{a} .

d. Das Sekanten-Verfahren

IDEE:

Ersetze im Newton-Verfahren die Tangente durch eine Sekante:

Seien x_0, x_1 zwei Startwerte:

$$\begin{aligned} f'(x_i) &\approx \frac{f(x_i) - f(x_{i-1})}{x_i - x_{i-1}} \\ x_{i+1} &= x_i - \frac{f(x_i)(x_i - x_{i-1})}{f(x_i) - f(x_{i-1})} \cdot \frac{f(x_i)x_{i-1} - f(x_{i-1})x_i}{f(x_i) - f(x_{i-1})}. \end{aligned}$$

BEMERKUNG:

Vorteil:

es muss nur $f(x)$ ausgewertet werden und nicht $f'(x) \Rightarrow$ schneller

Nachteil:

Konvergenzordnung: $p = \frac{1+\sqrt{5}}{2} \approx 1,6 \dots$

REGULA FALSI VERFAHREN:

Seien $x_0, y_0 \in [a, b]$ zwei Startwerte mit $f(x_0) \cdot f(y_0) < 0$.

$$s := \frac{f(x_0)y_0 - f(y_0)x_0}{f(x_0) - f(y_0)} \quad (\text{Sekantenschritt})$$

$$1. \quad f(s) = 0 \quad \checkmark$$

$$2. \quad f(x_0) \cdot f(s) < 0 \Rightarrow x_1 = x_0, y_1 = s \quad (\text{Bisektionsschritt})$$

$$3. \quad f(y_0) \cdot f(s) < 0 \Rightarrow x_1 = s, y_1 = y_0$$

e. Anwendung aus dem Bereich der Finanzmathematik

EFFEKTIVZINS VON SPARVERTRÄGEN (PRÄMIENBEGÜNSTIGTE SPARVERTÄGE):

Ein Sparer zahlt zu Beginn eines Jahres einen festen Betrag R ein.

Verzinsung: $p\%$. Dauer n Jahre.

Nach n Jahren bekommt der Sparer das Kapital ausbezahlt und eine Prämie von 14% auf die Sparbeträge.

FRAGE:

Effektivzins: welcher Zinssatz wäre anzusetzen, damit der gleiche Auszahlungsbetrag zustande kommt, wenn keine Prämie bezahlt wird?

Hier hat man

für die Spareinlage zu Beginn des ersten Jahres:	$R(1 + \frac{p}{100})^n = Rs^n$
für die Spareinlage zu Beginn des zweiten Jahres:	$R(1 + \frac{p}{100})^{n-1} = Rs^{n-1}$
⋮	⋮
für die Spareinlage zu Beginn des n -ten Jahres:	Rs

zusammen

$$R(s^n + s^{n-1} + \dots + s) + 0,14 \cdot R$$

Keine Prämie, dafür aber $q\%$ Zins. ($r = 1 + \frac{q}{100}$)

$$K_2 = R(r^n + \dots + r)$$

Wie ist r zu wählen, damit $K_1 = K_2$ gilt?

O.B.d.A. $R = 1$.

$$r^n + r^{n-1} + \dots + r = \underbrace{s^n + s^{n-1} + \dots + s + 0,14 \cdot n}_{=:c>0}$$

$$f(t) = t^n + \dots + t - c = 0$$

(positive Nullstellen für $t \in (0, \infty)$)

$$f'(t) = nt^{n-1} + \dots + 1 > 0 \quad \text{(V2) } \checkmark$$

$$f''(t) = n(n-1)t^{n-2} + \dots + 2 > 0 \quad \text{(V3) } \checkmark$$

$$f(0) = -c < 0$$

$$f(c) = c^n + \dots + c - c = c^n + \dots + c^2 > 0$$

\Rightarrow (nach dem Zwischenwertsatz) gibt es eine Nullstelle r in $(0, c)$ (V1) ✓

$$\begin{aligned} f'(x_0)(r - x_0) + f(x_0) &\leq f(r) = 0 \\ \Rightarrow r &\leq x_0 - \frac{f(x_0)}{f'(x_0)} = x_1 \Rightarrow x_1 \in (0, \infty) \end{aligned} \quad \text{(V4) ✓}$$

Also ist der Satz über die monotone Konvergenz beim Newton-Verfahren erfüllt.
Für jeden Startwert $x_0 \in (0, \infty)$ konvergiert das Newton-Verfahren gegen r .

Wahl des Startwertes: s .

4. Lineare Gleichungssysteme

a. Einleitung

$$Ax = b: \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix}.$$

Wir setzen eine eindeutige Lösung voraus, d.h. $A \neq 0$.

- **direkte Verfahren:** liefern nach endlich vielen Schritten die exakte Lösung.
Aber: Rundungsfehler!
ca. n^3 Operationen.
- **indirekte Verfahren:** iterative Verfahren: Formelfehler und Rundungsfehler.
ca. $k \cdot n^2$ Operationen \Rightarrow insbesondere für Gleichungssysteme mit großes n interessant.

b. Gestaffelte Systeme

(1) A obere Dreiecksmatrix:

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ & \ddots & \vdots \\ 0 & & a_{nn} \end{pmatrix} \text{ mit } a_{ii} \neq 0 \text{ für } i = 1, \dots, n.$$

$$x_n = \frac{b_n}{a_{nn}},$$

$$x_{n-1} = \frac{b_{n-1} - a_{n-1,n}x_n}{a_{n-1,n-1}},$$

\vdots

$$x_k = \frac{b_k - a_{k,k+1}x_{k+1} - \dots - a_{k,n}x_n}{a_{kk}} \text{ für } k = n, \dots, 1.$$

„Rückwärtsauflösen“

(2) A eine untere Dreiecksmatrix:

$$\begin{pmatrix} a_{11} \\ \vdots & \ddots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \text{ mit } a_{ii} \neq 0 \text{ für } i = 1, \dots, n.$$

$$x_1 = \frac{b_1}{a_{11}},$$

$$\vdots$$

$$x_k = \frac{b_k - a_{k1}x_1 - \dots - a_{k,k-1}x_{k-1}}{a_{kk}} \text{ für } k = 1, \dots, n.$$

„Vorwärtsauflösen“

c. Die LR-Zerlegung

Hier versucht man eine Zerlegung der Form $A = L \cdot R$ mit

$$L = \begin{pmatrix} 1 & & 0 \\ \vdots & \ddots & \\ l_{n1} & \dots & 1 \end{pmatrix}, R = \begin{pmatrix} r_{11} & \dots & r_{1n} \\ & \ddots & \vdots \\ 0 & & r_{nn} \end{pmatrix}.$$

$$Ax = b \Leftrightarrow L \cdot \underbrace{(Rx)}_y = b.$$

1. Löse das lineare Gleichungssystem $Ly = b$ durch Vorwärtsauflösen.
2. Löse das lineare Gleichungssystem $Rx = y$ durch Rückwärtsauflösen.

SATZ: Ist A regulär, so kann man durch vorherige Zeilenvertauschungen erreichen, dass die LR-Zerlegung existiert, also:

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

mit einer Permutationsmatrix P .

$$Ax = b \Leftrightarrow PAx = Pb$$

BERECHNUNG DER LR-ZERLEGUNG (VERFAHREN VON CROUT)

$$\underbrace{\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}}_A = \underbrace{\begin{pmatrix} 1 & & \\ \vdots & \ddots & \\ l_{n1} & \dots & 1 \end{pmatrix}}_L \cdot \underbrace{\begin{pmatrix} r_{11} & \dots & r_{1n} \\ & \ddots & \vdots \\ 0 & & r_{nn} \end{pmatrix}}_R$$

\Rightarrow

1. $a_{1k} = r_{1k}$ für $k = 1, \dots, n$.

$$2. a_{j1} = l_{j1} \cdot r_{11} \Rightarrow l_{j1} = \frac{a_{j1}}{r_{11}} \text{ für } j = 2, \dots, n.$$

3. usw.

ALLGEMEIN:

$$r_{ik} = a_{jk} - l_{j1}r_{1k} - \dots - l_{j,j-1}r_{j-1,k} \text{ für } k = j, \dots, n$$

$$l_{ij} = \frac{a_{ij} - l_{i1}r_{1j} - \dots - l_{i,j-1}r_{j-1,j}}{r_{jj}} \text{ für } i = j + 1, \dots, n$$

RECHENAUFWAND:

- für LR-Zerlegung benötigt man $\approx \frac{n^3}{3}$ Punktoperationen,
- ein gestaffeltes Gleichungssystem benötigt man $\approx \frac{n^2}{2}$ Punktoperationen.

ANMERKUNG:

Man vermeidet die Berechnung von A^{-1} , denn man benötigt n^3 Punktoperationen und das Ergebnis ist instabil (d.h. starke Verfälschung durch z.B. Rundungsfehler).

d. Die Cholesky-Zerlegung

$A \in \mathbb{R}^{n \times n}$ heißt positiv definit, falls gilt

- (1) A ist symmetrisch
- (2) $x^T A x > 0$ für alle $x \in \mathbb{R}^n \setminus \{0\}$.

SATZ:

Sei A positiv definit. Dann existiert eine Zerlegung der Form

$$A = L \cdot L^T$$

mit einer eindeutig bestimmten unteren Dreiecksmatrix $L = \begin{pmatrix} l_{11} & & 0 \\ \vdots & \ddots & \\ l_{n1} & \dots & l_{nn} \end{pmatrix}$ und

$l_{ii} > 0$ für $i = 1, \dots, n$.

Diese Zerlegung heißt **Cholesky-Zerlegung**.

BEWEIS:

vollständige Induktion nach n :

$$n = 1: A = (a_{11}) \text{ positiv definit} \Rightarrow a_{11} > 0 \Rightarrow L = (\sqrt{a_{11}}); LL^T = A \quad \checkmark$$

Behauptung sei richtig für $A_{n-1} \in \mathbb{R}^{(n-1) \times (n-1)}$

$$A_{n-1} = L_{n-1} \cdot L_{n-1}^T$$

Sei nun $A \in \mathbb{R}^{n \times n}$ positiv definit.

$$A = \left(\begin{array}{c|c} A_{n-1} & b \\ \hline b^T & a_{nn} \end{array} \right) \Rightarrow A_{n-1} \text{ ist positiv definit und } A_{n-1} = L_{n-1} \cdot L_{n-1}^T.$$

ANSATZ:

$$L = \left(\begin{array}{c|c} L_{n-1} & 0 \\ \hline c^T & \alpha \end{array} \right)$$

$$A = L \cdot L^T$$

$$\begin{aligned} \Rightarrow L \cdot L^T &= \left(\begin{array}{c|c} L_{n-1} & 0 \\ \hline c^T & \alpha \end{array} \right) \cdot \left(\begin{array}{c|c} L_{n-1}^T & c \\ \hline 0 & \alpha \end{array} \right) \\ &= \left(\begin{array}{c|c} L_{n-1} \cdot L_{n-1}^T & L_{n-1} \cdot c \\ \hline c^T \cdot L_{n-1}^T & c^T \cdot c\alpha^2 \end{array} \right) = \left(\begin{array}{c|c} A_{n-1} & b \\ \hline b^T & a_{nn} \end{array} \right) \end{aligned}$$

$$\Rightarrow L_{n-1} \cdot c = b. \text{ Da } L_{n-1} \text{ vollen Rang hat, ist sie invertierbar und es folgt } L_{n-1}^{-1} b.$$

$$c^T c + \alpha^2 = a_{nn} \Rightarrow \alpha^2 = a_{nn} - c^T c$$

$$0 < \det A = \det L \cdot \det L^T = (\det L)^2 = \underbrace{(\det L_{n-1})^2}_{>0} \cdot \alpha^2 \Rightarrow \alpha^2 > 0.$$

$$\Rightarrow \alpha = \sqrt{a_{nn} - c^T c}.$$

BERECHNUNG DER CHOLESKY-ZERLEGUNG:

Spaltenweises Vorgehen liefert für $j = 1, \dots, n$:

$$l_{jj} = \sqrt{a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2}$$

$$l_{ij} = \frac{1}{l_{jj}} \left(a_{ij} - \sum_{k=1}^{j-1} l_{ik} l_{jk} \right) \text{ für } i = j+1, \dots, n.$$

RECHENAUFWAND:

- für Cholesky-Zerlegung benötigt man $\approx \frac{n^3}{6}$ Punktoperationen.

e. Die QR-Zerlegung von A

Hier zerlegt man A in der Form $A = Q \cdot R$ mit einer orthogonalen Matrix Q und einer oberen Dreiecksmatrix R .

$$Ax = b \Leftrightarrow QRx = b \Rightarrow Rx = Q^T b \text{ (} \rightarrow \text{ Rückwärtsauflösung)}$$

SATZ:

Die QR-Zerlegung existiert für jede Matrix $A \in \mathbb{R}^{n \times n}$.

HOUSEHOLDER-MATRIZEN:

$h^T = (0, \dots, 0, h_k, \dots, h_n) \in \mathbb{R}^n$ mit $\|h^T\|_2 = 1$

$$h \cdot h^T = \begin{pmatrix} 0 & \dots & 0 \\ \vdots & h_k^2 & \dots & h_k h_n \\ & \vdots & \ddots & \vdots \\ 0 & h_n h_k & \dots & h_n^2 \end{pmatrix}$$

DEFINITION:

Eine Matrix der Form $H = I - 2h \cdot h^T$ heißt **Householder-Matrix**.

EIGENSCHAFTEN:

(a) $H^T = H$;

(b) H ist eine orthogonale Matrix:

$$H^T H = (I - 2hh^T)^T (I - 2hh^T) = (I - 2hh^T)(I - 2hh^T) = I - 4hh^T + 4h \underbrace{h^T h}_{=1} h^T = I.$$

ANMERKUNG:

Seien Q_1, Q_2 orthogonale Matrizen, dann folgt $Q_1 Q_2$ ist ebenfalls eine orthogonale Matrix.

KONSTRUKTION DER QR-ZERLEGUNG:

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nm} \end{pmatrix} =: a_1.$$

1. Schritt:

Es gibt eine Householder-Matrix H_1 mit $H_1 a_1 = \sigma e_1$

$$e_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^n$$

$$\sigma^2 = \sigma e_1^T \cdot \sigma e_1 = (H_1 a_1)^T H_1 a_1 = a_1^T \underbrace{H^T H}_{=I} a_1 = a_1^T a_1 = \|a_1\|_2^2 \Rightarrow \sigma = \pm \|a_1\|_2$$

$$\sigma e_1 = H_1 a_1 = (I - 2h_1 h_1^T) a_1 = a_1 - 2h_1 \underbrace{h_1^T a_1}_{\in \mathbb{R}} = a_1 - 2(h_1^T a_1) h_1$$

$$2(h_1^T a_1) h_1 = a_1 - \sigma e_1.$$

$$\text{Wegen } \|h_1\|_2 = 1 \text{ folgt } h_1 = \frac{a_1 - \sigma e_1}{\|a_1 - \sigma e_1\|_2}.$$

Aus Stabilitätsgründen setzt man $\sigma = -\text{sign}(a_{11}) \|a_1\|_2$.

$$\Rightarrow h_1 = \frac{a_1 + \text{sign}(a_{11}) \|a_1\|_2 e_1}{\|a_1 + \text{sign}(a_{11}) \|a_1\|_2 e_1\|}.$$

$$\text{Nun gilt mit } h_1: H_1 \cdot A = \begin{pmatrix} \sigma & \dots \\ 0 & \tilde{A}_1 \\ \vdots & \\ 0 & \end{pmatrix}$$

2. Schritt:

Es gibt eine Householder-Matrix \tilde{H}_2 ($(n-1) \times (n-1)$ -Matrix) mit $\tilde{H}_2 a_2 = \gamma \tilde{e}_1$

$$\tilde{e}_1 = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{n-1}.$$

Setze $H_2 := \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & \tilde{H}_2 \end{array} \right)$; H_2 ist eine $n \times n$ -Householder-Matrix; $h_2 = (0, x, \dots, x)$.

Nun gilt $H_2 \cdot H_1 \cdot A = \left(\begin{array}{c|c} \sigma & \\ \hline \gamma & \\ \hline & \tilde{A}_2 \end{array} \right)$.

⋮

$n - 1$. Schritt:

$$\Rightarrow \underbrace{H_{n-1}H_{n-2} \dots H_2H_1}_{=:Q^T} \cdot A = \begin{pmatrix} x & \dots & x \\ & \ddots & \vdots \\ 0 & & x \end{pmatrix}$$

$$\Rightarrow Q = H_1^T H_2^T \dots H_{n-1}^T = H_1 H_2 \dots H_{n-1}.$$

RECHENAUFWAND:

- für QR-Zerlegung benötigt man $\approx \frac{2n^3}{3}$ Punktoperationen, dafür ist der Algorithmus stabil.

Die QR-Zerlegung bzw. die LR-Zerlegung ist besonders günstig, wenn mehrere Gleichungssysteme mit der selben Koeffizientenmatrix zu lösen sind.

f. Die Gauß-Elimination

$Ax = b$.

$$(A, b) = \begin{pmatrix} a_{11} & \dots & a_{1n} & b_1 \\ \vdots & \ddots & \vdots & \vdots \\ a_{n1} & \dots & a_{nn} & b_n \end{pmatrix} \rightarrow \text{elementare Zeilenumformungen} \rightarrow \begin{pmatrix} \gamma_{11} & \dots & \gamma_{1n} & c_1 \\ & \ddots & \vdots & \vdots \\ 0 & & \gamma_{nn} & c_n \end{pmatrix}$$

$$R := \begin{pmatrix} \gamma_{11} & \dots & \gamma_{1n} \\ & \ddots & \vdots \\ 0 & & \gamma_{nn} \end{pmatrix}$$

$Ax = b \Leftrightarrow Rx = c \rightarrow$ Rückwärtsauflösen.

Elementare Zeilenumformungen:

- (1) Vertauschen von Zeilen;
- (2) Addition der Vielfachen einer Zeile zu einer anderen Zeile.

Es sei $a_{11} \neq 0$. Mit den Operationen

$$l_{i1} := \frac{a_{i1}}{a_{11}} \text{ für } i = 2, \dots, n$$

$$a_{ik}^{(1)} := a_{ik} - l_{i1}a_{1k} \text{ für } i, k = 2, \dots, n$$

$$b_i^{(1)} = b_i - l_{i1}b_1$$

$$\Rightarrow \begin{pmatrix} a_{11} & \dots & \dots & a_{1n} & b_1 \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} & b_2^{(1)} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & a_{n2}^{(1)} & \dots & a_{nn}^{(1)} & b_n^{(1)} \end{pmatrix}$$

Ein Eliminationsschritt:

$$\Rightarrow x_1 = \frac{1}{a_{11}} \left(b_1 - \sum_{k=2}^n a_{1k}x_k \right)$$

Nach $n - 1$ solcher Schritte erhalte:

$$\begin{pmatrix} a_{11} & \dots & \dots & a_{1n} & b_1 \\ 0 & a_{22}^{(1)} & \dots & a_{2n}^{(1)} & b_2^{(1)} \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & \dots & 0 & a_{nn}^{(n-1)} & b_n^{(n-1)} \end{pmatrix}$$

DEFINITION:

Die Einträge $a_{11}, a_{22}^{(1)}, \dots, a_{nn}^{(n-1)}$ heißen **Pivot-Elemente**.

SATZ:

Ist A regulär, so kann durch Zeilentausch erreicht werden, dass alle Pivot-Elemente von Null verschieden sind (vgl. Übungsaufgabe).

RECHENAUFWAND:

- für die Gauß-Elimination benötigt man $\approx \frac{n^3}{3}$ Punktoperationen.

PIVOT-STRATEGIEN

- BETRAGSMAXIMALE SPALTENPIVOTWAHL:

betragsgrößtes Element: $\max_{i \geq k} |a_{ik}^{(k-1)}| =: |a_{pk}^{(k-1)}|$.

Vertausche vor dem k -ten Eliminationsschritt die k -te und p -te Zeile.

g. Matrizentechnische Deutung der Gauß-Elimination

$$G_1 = \begin{pmatrix} 1 & & & 0 \\ -l_{21} & 1 & & \\ \vdots & & \ddots & \\ -l_{n1} & 0 & & 1 \end{pmatrix}.$$

Der 1. Eliminationsschritt bedeutet dann $G_1 A x = G_1 b$.

$$G_2 = \begin{pmatrix} 1 & & & & \\ & 1 & & & \\ & -l_{32} & 1 & & \\ & \vdots & & \ddots & \\ & -l_{n2} & & & 1 \end{pmatrix}.$$

Der 2. Eliminationsschritt bedeutet dann $G_2 G_1 A x = G_2 G_1 b$.

⋮

Nach dem $n - 1$ -sten Eliminationsschritt:

$$G_{n-1} \cdot \dots \cdot G_1 \cdot A = R.$$

GAUSS-JORDAN-ELIMINATION

$$A = \begin{pmatrix} a_{11} & \dots & a_{1k} & \dots & a_{1m} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{i1} & \dots & a_{ik} & \dots & a_{im} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nk} & \dots & a_{nm} \end{pmatrix} \rightarrow JA = \begin{pmatrix} a_{11}^{(1)} & \dots & 0 & \dots & a_{1m}^{(1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{i1}^{(1)} & \dots & 1 & \dots & a_{im}^{(1)} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ a_{n1}^{(1)} & \dots & 0 & \dots & a_{nm}^{(1)} \end{pmatrix}$$

$$J = \begin{pmatrix} 1 & & & -l_{1i} & & \\ & \ddots & & \vdots & & \\ & & & l_{ii} & & \\ & & & \vdots & \ddots & \\ & & & -l_{ni} & & 1 \end{pmatrix} \text{ mit } l_{ii} = \frac{1}{a_{ik}}, l_{ij} = \frac{a_{jk}}{a_{ik}}, i \neq j.$$

Sei $A \in \mathbb{R}^{n \times n}$ invertierbar. Berechne A^{-1} :

$$\begin{aligned} & \left(\begin{array}{ccc|cc} a_{11} & \dots & a_{1n} & 1 & 0 \\ \vdots & \ddots & \vdots & & \ddots \\ a_{n1} & \dots & a_{nn} & 0 & 1 \end{array} \right) \rightarrow \left(\begin{array}{cccc|ccc} 1 & x & \dots & x & x & 0 & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & & \ddots \\ 0 & x & \dots & x & x & 0 & 1 \end{array} \right) \\ & \rightarrow \dots \rightarrow \left(\begin{array}{cc|ccc} 1 & 0 & x & \dots & 0 \\ & \ddots & \vdots & \ddots & \vdots \\ 0 & 1 & x & \dots & x \end{array} \right) \end{aligned}$$

5. Das lineare Ausgleichsproblem

a. Problemstellung

GEGEBEN: Messergebnisse (z.B. $(t_i, s_i), i = 1, \dots, m$);
Theoriefunktion (z.B. $G(t, \underbrace{\alpha_1, \dots, \alpha_n}_{=\alpha}) = G(t, \alpha)$); $m > n$.

AUFGABE: Bestimme die Parameter $\alpha_1, \dots, \alpha_n$ anhand der Messergebnisse.

$G(t_i, \alpha_1, \dots, \alpha_n) = s_i$ für $i = 1, \dots, m$ ist optimal, aber im Allgemeinen nicht möglich, da $m > n$.

aber: $G(t_i, \alpha_1, \dots, \alpha_n) \approx s_i$ für $i = 1, \dots, m$ kann „möglichst gut“ erreicht werden. Bestimme also $\alpha_1, \dots, \alpha_n$ so, dass dies erreicht wird,

d.h. $\underbrace{\sum_{i=1}^n (G(t_i, \alpha_1, \dots, \alpha_n) - s_i)^2}_{\text{Fehlerquadratsumme}}$ minimieren.

DEFINITION:

Dies ist ein allgemeines Ausgleichsproblem.

BEISPIEL:

$G(t, a, b, c) = \frac{a}{1 + \exp(b - ct)}$ (logistische Kurve, nicht linear)

DEFINITION:

Ein Ausgleichsproblem heißt **linear**, falls

$$G(t, \alpha_1, \dots, \alpha_n) = \sum_{i=1}^n \alpha_i g_i(t)$$

mit gewissen Grundfunktionen $g_i(t)$.

BEISPIEL:

$g_i(t) = t^{i-1}$. (Polynomausgleich)
 $\sum_{i=1}^m \left(\sum_{j=1}^n \alpha_j g_j(t_i) - s_i \right)^2$.

Setze $a_{ij} := g_j(t_i)$ für $j = 1, \dots, n; i = 1, \dots, m$.

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix}.$$

$$\text{Setze } b = \begin{pmatrix} s_1 \\ \dots \\ s_n \end{pmatrix}.$$

$$\text{Dann gilt } \sum_{i=1}^m \left(\sum_{j=1}^n \alpha_j g_j(t_i) - s_i \right)^2 = \|A\alpha - b\|_2^2.$$

Beim Polynomausgleich gilt:

$$A = \begin{pmatrix} 1 & t_1 & \dots & t_1^{m-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & t_m & \dots & t_m^{m-1} \end{pmatrix}.$$

b. Die Normalgleichungen

$A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ gegeben. Gesucht ist ein $x \in \mathbb{R}^n$ mit $\|Ax - b\|_2^2 \rightarrow \min$.

$F: \mathbb{R}^n \rightarrow \mathbb{R}^+$, $F(x) = \|Ax - b\|_2^2 = (Ax - b)^T \cdot (Ax - b) = x^T A^T A x - 2(A^T b)^T x + b^T b$.
 $\Rightarrow F(x)$ ist differenzierbar.

$$\nabla F(x) = 2A^T A x - 2A^T b$$

Notwendige Bedingung für ein lokales Minimum: $\nabla F = 0$.

$$\Rightarrow A^T A x = A^T b$$

(Normalgleichungen)

SATZ:

Jede Lösung der Normalgleichungen ist ein globales Minimum von F .

BEWEIS:

Sei \bar{x} eine Lösung der Normalgleichungen: $A^T A \bar{x} = A^T b$.

Dann gilt für jedes $x \in \mathbb{R}^n$:

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

$$= F(\bar{x}) + 2 \underbrace{\langle A^T A \bar{x} - A^T b, x - \bar{x} \rangle}_{=0} + \underbrace{\|A(x - \bar{x})\|_2^2}_{\geq 0} = \quad (\text{da } \langle Ax, y \rangle = \langle x, A^T y \rangle)$$

$$\geq F(\bar{x}).$$

$\Rightarrow \bar{x}$ ist ein globales Minimum.

□

KOROLLAR:

A hat den Rang n . Dann besitzt F genau ein globales Minimum.

BEWEIS:

$$x^T A^T A x = (Ax)^T (Ax) = \|Ax\|_2^2 \geq 0.$$

$$\|Ax\|_2^2 = 0 \Leftrightarrow Ax = 0 \Rightarrow x = 0.$$

$\Rightarrow A^T A$ positiv definit.

□

BEMERKUNG:

Man kann die Cholesky-Zerlegung verwenden.

6. Lineare Optimierung

a. Problemstellung

Maximum bzw. Minimum einer Funktion unter einer Nebenbedingung. Maximum bzw. Minimum einer linearen Funktion unter einer linearen Nebenbedingung → lineare Optimierung.

BEISPIEL:

Ein Landwirt besitzt 80 ha Ackerfläche und baut Gerste und Weizen an. Der Arbeitseinsatz pro ha Gerste beträgt 6 Stunden, pro ha Weizen 10 Stunden. Der Landwirt hat maximal 560 Stunden Zeit. Die Kosten für pro ha Gerste betragen 30 €, pro ha Weizen 100 €; es stehen maximal 4700 € zur Verfügung. Der Gewinn beträgt pro ha Gerste 40 € und pro ha Weizen 50 €.

FRAGE: wie ist die Ackerfläche aufzuteilen, um einen möglichst hohen Gewinn zu erzielen?

Sei x_1 die Fläche für Gerste und x_2 die Fläche für Weizen.

$$\begin{aligned} \Rightarrow x_1 + x_2 &\leq 80 && \text{(Anbaufläche)} \\ 6x_1 + 10x_2 &\leq 560 && \text{(Arbeit)} \\ 30x_1 + 100x_2 &\leq 4700 && \text{(Kapital)} \\ x_1 \geq 0, x_2 &\geq 0. \end{aligned}$$

Gewinn: $g(x_1, x_2) = 40x_1 + 50x_2 \rightarrow \max$.

MATRIXSCHREIBWEISE:

$$P = \begin{pmatrix} 1 & 1 \\ 6 & 10 \\ 30 & 100 \end{pmatrix}, x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, b = \begin{pmatrix} 80 \\ 560 \\ 4700 \end{pmatrix}, q = \begin{pmatrix} 40 \\ 50 \end{pmatrix}.$$

$\Rightarrow g(x) := q^T x \rightarrow \max$ unter den Nebenbedingungen:

$$Px \leq b, x \geq 0. \quad x \leq y: \Leftrightarrow x_i \leq y_i, i = 1, \dots, n.$$

ALLGEMEINER FALL:

$$q \in \mathbb{R}^n, b \in \mathbb{R}^m, P \in \mathbb{R}^{m \times n}, x \in \mathbb{R}^n.$$

$$g(x) = q^T x \rightarrow \max \text{ unter } Px \leq b, x \geq 0. \quad \text{(Normalform 1)}$$

Diese Aufgabenstellung heißt ein **lineares Optimierungsproblem (lineares Programm)**. $g(x) = q^T x$ heißt **Zielfunktion**.

Die Menge $K = \{x \in \mathbb{R}^n : Px \leq b \text{ und } x \geq 0\}$ heißt **zulässiger Bereich**.

b. Grafische Lösung

siehe Mathe-Heft Klasse 9 ;-).

⇒ das Maximum liegt in einer Ecke!

c. Normalformen

NORMALFORM 1: $g(x) = q^T x \rightarrow \max$ unter $Px \leq b, x \geq 0$. (siehe oben)

NORMALFORM 2: $\max g(x) = q^T x$ unter $Px + s = b, x \geq 0, s \geq 0$.

$$s = \begin{pmatrix} s_1 \\ \vdots \\ s_m \end{pmatrix}$$

s heißt **Schlupfvariable**.

NORMALFORM 3: $A \in \mathbb{R}^{m \times (m+n)}$ mit $rk(A) = m; b \in \mathbb{R}^m, c \in \mathbb{R}^{m+n}$
 $\max g(z) = c^T z$ unter $Az = b, z \geq 0$.
 $A = (P, I); c = (q, 0 \dots 0)^T, z = (x_1, \dots, x_n, s_1, \dots, s_m)^T$.

ANMERKUNG:

- (1) $\min g(x)$ ist äquivalent zu $\max -g(x)$;
- (2) Eine Nebenbedingung der Form

$$a_{k1}x_1 + \dots + a_{kn}x_n \geq b_k$$

wird zu

$$-a_{k1}x_1 - \dots - a_{kn}x_n \leq -b_k;$$

- (3) unterliegt eine Komponente x_k keiner Vorzeichen-Beschränkung, dann $x_k = x_{k1} - x_{k2}$ mit $x_{k1} \geq 0$ und $x_{k2} \geq 0$.

Alle linearen Optimierungsprobleme lassen sich also auf die oben genannten Normalformen zurückführen.

d. Charakterisierung des zulässigen Bereichs

$$K = \{x \in \mathbb{R}^n : Px \leq b \text{ und } x \geq 0\} \quad (\text{Normalform 1})$$

$$K = \{x \in \mathbb{R}^{n+m} : Az = b \text{ und } z \geq 0\} \quad (\text{Normalform 2})$$

KONVEXE MENGEN:

DEFINITION:

Eine Menge $M \subset \mathbb{R}^k$ heißt **konvex**, falls für alle $x_1, x_2 \in M$ und alle $0 \leq \lambda \leq 1$ der Punkt $x_\lambda := \lambda x_1 + (1 - \lambda)x_2 \in M$ gilt.

EIGENSCHAFTEN:

(1) Der Durchschnitt beliebig vieler konvexer Mengen ist konvex.

(2) Seien $x_1, \dots, x_r \in M \Rightarrow \sum_{i=1}^r \mu_i x_i \in M$ für alle $\mu_i \geq 0$ mit $\sum_{i=1}^r \mu_i = 1$. (konvexe Linearkombination)

(3) Seien $y_1, \dots, y_r \in \mathbb{R}^k$. Dann ist die Menge $K = K(y_1, \dots, y_r) := \{x \in \mathbb{R}^k : x = \sum_{i=1}^r \mu_i y_i, \mu_i \geq 0 \text{ mit } \sum_{i=1}^r \mu_i = 1\}$ konvex. (**konvexe Hülle** der Punkte y_1, \dots, y_r)

DEFINITION:

Sei M eine konvexe Menge. Ein $x \in M$ heißt **Ecke** (oder **Extrempunkt**), wenn $x = \mu x_1 + (1 - \mu)x_2$ mit $0 < \mu < 1$, $x_1, x_2 \in M$, folgt, dass $x = x_1 = x_2$.

HYPEREBENEN UND HALBRÄUME

DEFINITION:

Sei $a^T = (\alpha_1, \dots, \alpha_k) \in \mathbb{R}^k \setminus \{0\}$ und $\gamma \in \mathbb{R}$. Die Menge $H := \{x \in \mathbb{R}^k : a^T x = \gamma\}$ heißt **Hyperebene**.

Eine Hyperebene zerteilt den \mathbb{R}^k in zwei Halbräume

$$H^+ := \{x \in \mathbb{R}^k : a^T x \geq \gamma\}$$

$$H^- := \{x \in \mathbb{R}^k : a^T x \leq \gamma\}.$$

LEMMA:

H, H^-, H^+ sind konvex und abgeschlossen.

ZULÄSSIGER BEREICH K :

Jede Nebenbedingung liefert eine Hyperebene bzw. einen Halbraum.
 K ist endlicher Durchschnitt von Hyperebenen bzw. Halbräumen. \Rightarrow

SATZ:

K ist konvex und abgeschlossen.

ANMERKUNG:

Sei $K \neq \emptyset$ und K beschränkt, dann gibt es eine optimale Lösung.

e. Ecken des zulässigen Bereichs

(Normalform 3: $\max g(z) = c^T z$ unter $Az = b, z \geq 0$)

$$A = (a^{(1)}, a^{(2)}, \dots, a^{(m+n)}).$$

Sei $x \in K$, dann definiere $I(x) := \{i \in \{1, \dots, m+n\} : x_i > 0\}$.

SATZ: (Charakterisierungssatz)

Äquivalent sind (für $x \in K$)

- (1) x ist eine Ecke von K ;
- (2) Die Spalten $a^{(i)}$ mit $i \in I(x)$ sind linear unabhängig.

BEWEIS:

(i) \Rightarrow (ii):

sei x eine Ecke. Ohne Einschränkung der Allgemeinheit sei $I(x) = \{1, \dots, r\}$.

$$b = Ax = \sum_{i=1}^{m+n} x_i a^{(i)} = \sum_{i=1}^r x_i a^{(i)}.$$

Annahme: $a^{(1)}, \dots, a^{(r)}$ sind linear abhängig: dann gibt es ein $(\alpha_1, \dots, \alpha_r) \neq (0, \dots, 0)$

$$\text{mit } \sum_{i=1}^r \alpha_i a^{(i)} = 0. \quad (*)$$

Wegen $x_i > 0$ für $i = 1, \dots, r$ gibt es ein $\varepsilon > 0$ mit $x_j \pm \varepsilon \alpha_j \geq 0$ für $j = 1, \dots, r$.

$$\bar{\alpha} = (\alpha_1, \dots, \alpha_r, 0, \dots, 0)$$

$$y_+ = (x_1 + \varepsilon \alpha_1, \dots, x_r + \varepsilon \alpha_r, 0, \dots, 0) \geq 0$$

$$y_- = (x_1 - \varepsilon\alpha_1, \dots, x_r - \varepsilon\alpha_r, 0, \dots, 0) \geq 0$$

Es gilt $y_+ \neq y_-$; $x = \frac{1}{2}y_+ + (1 - \frac{1}{2})y_-$.

$$Ay_{\pm} = Ax \pm \varepsilon A\bar{a} = Ax + \varepsilon \underbrace{\sum_{i=1}^r \alpha_i a^{(i)}}_{=0(\text{nach } (*))} = b$$

$\Rightarrow y_+, y_- \in K \Rightarrow x$ ist keine Ecke.

(ii) \Rightarrow (i):

zu $x \in K$ seien $a^{(1)}, \dots, a^{(r)}$ linear unabhängig.

Sei $x = \lambda y + (1 - \lambda)z$, $0 < \lambda < 1$, $y, z \in K$.

$\Rightarrow I(x) = I(y) \cup I(z)$.

$$Ay = Az = b \Rightarrow 0 = Ay - Az = A(y - z) = \sum_{j=1}^{m+n} (y_j - z_j)a^{(j)} = \sum_{j=1}^r (y_j - z_j)a^{(j)}$$

$\Rightarrow (y_j - z_j) = 0$ für alle $j \in \{1, \dots, r\}$ (da $a^{(1)}, \dots, a^{(r)}$ linear unabhängig)

$\Rightarrow y = z \Rightarrow x$ ist eine Ecke.

□

SATZ: (Existenzsatz)

Der zulässige Bereich $K \subset \mathbb{R}^{m+n}$, $K \neq \emptyset$, hat Ecken.

BEWEIS:

$$I := \{|I(z)| : z \in K\} \subset \{0, \dots, n + m\}$$

$\gamma = \min I$. Sei $x \in K$ mit $|I(x)| = \gamma$.

$\gamma = 0$: $\Rightarrow x = (0, \dots, 0)$ ist eine Ecke.

$\gamma > 0$: Ohne Einschränkung sei $I(x) = \{1, \dots, \gamma\}$.

Betrachte die Spalten $a^{(1)}, \dots, a^{(\gamma)}$. Annahme: diese Spalten sind linear abhängig.

$$\Rightarrow \sum_{j=1}^{\gamma} \alpha_j a^{(j)} = 0 \text{ mit } (\alpha_1, \dots, \alpha_{\gamma}) \neq (0, \dots, 0).$$

$$\lambda : \min \left\{ \frac{x_j}{|\alpha_j|} : \alpha_j \neq 0, 1 \leq j \leq \gamma \right\} = \frac{x_k}{|\alpha_k|} \text{ für ein } k. \quad (\text{O.E. } |\alpha_k| = \alpha_k)$$

$$\text{Setze } \bar{x} = (x_1 - \lambda\alpha_1, x_2 - \lambda\alpha_2, \dots, x_{\gamma} - \lambda\alpha_{\gamma}, 0, \dots, 0) \geq 0.$$

$$A\bar{x} = Ax - \lambda \underbrace{\sum_{j=1}^{\gamma} \alpha_j a^{(j)}}_{=0} = Ax = b.$$

$\Rightarrow \bar{x} \in K$ und $\bar{x}_k = 0$.

$\Rightarrow |I(\bar{x})| < \gamma \rightarrow$ Widerspruch.

Also $a^{(1)}, \dots, a^{(\gamma)}$ sind linear unabhängig. $\Rightarrow x$ ist eine Ecke.

□

SATZ: (Darstellungssatz)

Sei $K \neq \emptyset$ und beschränkt.

Zu jedem $x \in K$ gibt es Ecken $z^{(1)}, \dots, z^{(l)}$ mit

$$x = \sum_{j=1}^l \lambda_j z^{(j)}; \lambda_j \geq 0, \sum_{j=1}^l \lambda_j = 1.$$

f. Basislösungen

DEFINITION:

Sei $A \in \mathbb{R}^{m \times (n+m)}$ mit $\text{rk}(A) = m$ und $B(a^{(i_1)}, \dots, a^{(i_m)})$ eine Teil-Matrix mit Rang $B = m$.

Ein $x \in \mathbb{R}^{m+n}, x \geq 0$ heißt **Basispunkt** zu B , falls gilt:

$$x_j = 0 \text{ für } j \notin \{i_1, \dots, i_m\}$$

$$\sum_{j=1}^m x_{i_j} a^{(i_j)} = b$$

Die **Komponenten** x_{i_1}, \dots, x_{i_m} heißen **Basisvariable**.

Ein $x \in \mathbb{R}^{m+n}$ heißt **Basispunkt** (bzw. **Basislösung**), wenn es eine Teil-Matrix B von A gibt, so dass x Basispunkt zu B ist.

SATZ: (Äquivalenzsatz)

Sei $\text{rk}(A) = m$ und $x \in \mathbb{R}^{m+n}$. Dann sind äquivalent:

- (1) x ist eine Ecke von K ;
- (2) x ist eine Basislösung.

BEWEIS:

(1) \Rightarrow (2):

Sei x eine Ecke mit $I(x) = \{i_1, \dots, i_p\}$.

Nach dem Charakterisierungssatz sind $a^{(i_1)}, \dots, a^{(i_p)}$ linear unabhängig.

($p \leq m$).

$$B = (a^{(i_1)}, \dots, a^{(i_p)}, b^{(i_{p+1})}, \dots, b^{(i_m)}).$$

(ergänze $a^{(i_1)}, \dots, a^{(i_p)}$ linear unabhängig zu $a^{(i_1)}, \dots, a^{(i_p)}, b^{(i_{p+1})}, \dots, b^{(i_m)}$)

(2) \Rightarrow (1):

Sei x ein Basispunkt zu $B = (a^{(i_1)}, \dots, a^{(i_m)}) \Rightarrow a^{(i_1)}, \dots, a^{(i_m)}$ sind linear unabhängig.
 $I(x) \subset \{i_1, \dots, i_m\} \Rightarrow a^{(j)}, j \in I(x)$ sind linear unabhängig $\Rightarrow x$ ist eine Ecke. \square

KOROLLAR:

Es gibt maximal $\binom{n+m}{m}$ Ecken.

g. Lösung der linearen Optimierungsaufgabe

SATZ: (Hauptsatz)

Sei $K \neq \emptyset$ und beschränkt. Dann nimmt die Zielfunktion ihr Maximum in einer Ecke an.

BEWEIS:

K ist kompakt; die stetige Funktion $g(z)$ nimmt also ihr Maximum an in einem $\bar{z} \in K$.

$\Rightarrow \bar{z} = \sum_{j=1}^l \lambda_j x^{(j)}$, wobei $x^{(j)}$ Ecken sind.

$$c^T x^{(k)} := \max\{c^T x^{(j)} : j = 1, \dots, l\} \leq \max\{c^T z : z \in K\} = c^T \bar{z} = \sum_{j=1}^l \lambda_j \underbrace{c^T x^{(j)}}_{\leq c^T x^{(k)}} \leq c^T x^{(k)}$$

$$\Rightarrow c^T x^{(k)} = c^T \bar{z}. \quad \square$$

Ziel: Systematisches Berechnen von Ecken. \Rightarrow **Simplexverfahren.**

h. Der Simplex-Algorithmus

$g(x) := q^T x + \alpha^{(0)} \rightarrow \max$ unter den Nebenbedingungen

$$Px + s = b$$

$$x \geq 0$$

$$s \geq 0$$

(Normalform 2)

Zusatzforderung: $b \geq 0$

$$\left(\begin{array}{c|c|c} P & I_m & b \\ \hline q^T & 0^T & \alpha^{(0)} \end{array} \right) = \left(\begin{array}{ccc|cc|c} p_{11} & \dots & p_{1n} & 1 & 0 & b_1 \\ \vdots & \ddots & \vdots & & \ddots & \vdots \\ p_{m1} & \dots & p_{mn} & 0 & 1 & b_m \\ \hline q_1 & \dots & q_n & 0 & \dots & 0 & \alpha^{(0)} \end{array} \right) = H^{(0)}$$

$$\text{Eine Basislösung ist } z^{(0)} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ b_1 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} x_1 \\ \vdots \\ x_n \\ s_1 \\ \vdots \\ s_m \end{pmatrix}$$

Suche nun nach einer neuen Basislösung (Ecke) mit höherem Wert der Zielfunktion.

Gauß-Jordan-Schritt auf $H^{(0)}$ liefert eine neue Basislösung aus $H^{(1)}$. Gauß-Jordan-Schritt muss so erfolgen, dass $b^{(1)} \geq 0$ (sonst ist die abgelesene Lösung nicht im zulässigen Bereich).

Wahl des Pivots: $p_{rs} > 0; q_s > 0$.

$$b_i^{(1)} = b_i^{(0)} - \frac{p_{is}}{p_{rs}} b_r \geq 0$$

$$\Rightarrow b_i^{(0)} p_{rs} \geq p_{is} b_r$$

$$\frac{b_i^{(0)}}{p_{is}} \geq \frac{b_r}{p_{rs}}$$

$$H^{(i)} = \left(\begin{array}{cccc|c} a_{11}^{(i)} & \dots & a_{1n}^{(i)} & \dots & a_{1,n+m}^{(i)} & b_1^{(i)} \\ \vdots & & \ddots & & \vdots & \vdots \\ a_{m1}^{(i)} & \dots & a_{mn}^{(i)} & \dots & a_{m+n,mn}^{(i)} & b_m^{(i)} \\ \hline c_1^{(i)} & \dots & a_{mn}^{(i)} & \dots & c_{m+n}^{(i)} & -\alpha^{(i)} \end{array} \right) \text{ mit Basislösung } z^{(i)}.$$

SATZ: (OPTIMALITÄTSPRÜFUNG)

Gilt $c_j^{(i)} \leq 0$ für $j = 1, \dots, m+n$, so gilt $z^{(i)}$ eine optimale Lösung.

SATZ: (NICHTLÖSBARKEIT)

Gibt es ein $c_k^{(i)} < 0$ mit $a_{jk}^{(i)} \leq 0$ für $j = 1, \dots, m$, dann gibt es keine Lösung (also kein Maximum).

SATZ:

Gibt es ein $c_k^{(i)} > 0$ und ein $a_{rk}^{(i)} > 0$, so findet man eine Basislösung $z^{(i+1)}$ mit

$$g(z^{(i+1)}) \geq g(z^{(i)}).$$

i. Zweiphasenmethode

Falls $b \geq 0$ nicht erfüllt ist, dann ist $z^{(0)} = \begin{pmatrix} 0 \\ b \end{pmatrix}$ keine zulässige Lösung!

⇒ Anlaufrechnung, um eine zulässige Basislösung zu bekommen.

PHASE 1:

Bestimme eine zulässige Basislösung.

OBdA sei $b_i < 0$ für $i = 1, \dots, p$ und $b_i \geq 0$ für $i = p + 1, \dots, m$.

Führe zusätzliche Variablen ξ_1, \dots, ξ_p ein.

$$-\sum_{k=1}^n p_{ik}x_k - s_i + \xi_i = -b_i \text{ für } i = 1, \dots, p$$

$$x \geq 0, s \geq 0, \xi \geq 0.$$

$$\left(\begin{array}{ccc|cccc|ccc} -p_{11} & \dots & -p_{1n} & -1 & & & 0 & 1 & 0 & -b_1 \\ & & & & \ddots & & & & \ddots & \vdots \\ & & & & & -1 & & 0 & 1 & -b_p \\ & & & & & & 1 & 0 & \dots & b_{p+1} \\ & & & & & & & \vdots & \vdots & \vdots \\ p_{m1} & \dots & p_{mn} & 0 & & & 1 & 0 & \dots & b_m \\ \hline q_1 & \dots & q_m & 0 & & & 0 & 0 & \dots & -\alpha \end{array} \right)$$

Basislösung ist $z^{(0)} = (x_1, \dots, x_m, s_1, \dots, s_p, s_{p+1}, \dots, s_m, \xi_1, \dots, \xi_p) = (0, \dots, 0, 0, \dots, 0, b_{p+1}, \dots, b_m, -b_1,$

Zielfunktion:

$$h(\xi) := -\sum_{i=1}^p \xi_i \rightarrow \max.$$

(x, s, ξ) ist eine optimale Lösung des Hilfsproblems $\Leftrightarrow \xi = (0, \dots, 0)$

(x, s, ξ) ist eine optimale Lösung des Hilfsproblems $\Rightarrow (x, s)$ ist eine zulässige Lösung des Ausgangsproblems.

PHASE 2:

Führe mit (x, s) das Simplex-Verfahren durch.

7. Numerische Lösung von Anfangswertaufgaben

a. Einleitung und Beispiele

$$\dot{x} = F(t, x); x(t_0) = y^{(0)}$$

$$\left[\begin{array}{l} \dot{x}_1 = f_1(t, x_1, \dots, x_n); x_1(t_0) = y_1^{(0)} \\ \vdots \\ \dot{x}_n = f_n(t, x_1, \dots, x_n); x_n(t_0) = y_n^{(0)} \end{array} \right]$$

Wir setzen die eindeutige Lösbarkeit voraus.

BEISPIELE:

(1) Biologie: Wachstum

Wachstum ist proportional zum Umfang der Population.

$x(t)$: Umfang der Population zum Zeitpunkt t .

$$\dot{x}(t) = \lambda x(t) \text{ für } \lambda > 0.$$

$$\dot{x} = F(x).$$

Wachstum wird durch die Umwelt begrenzt.

$$\dot{x} = Rx \left(1 - \frac{x}{K}\right), R > 0, K > 0$$

(Verhulstgleichung)

(2) Physik / Chemie: Radioaktiver Zerfall

$$\dot{x}(t) = k \cdot x(t) \text{ für } k > 0.$$

(3) Physik:

(a) Fallschirmspringer

$v(t)$: Fallgeschwindigkeit zum Zeitpunkt t

t_1 : Reißleine wird gezogen; t_2 : Fallschirm vollständig offen.

$$\dot{v}(t) = g - \frac{k(t) \cdot v(t)^2}{m};$$

$$k(t) = \begin{cases} k_1(t) & , \text{ für } 0 \leq t \leq t_1 \\ k_1(t) + \frac{(t-t_1)(k_2-k_1)}{t_2-t_1} & , \text{ für } t_1 \leq t < t_2 \\ k_2(t) & , \text{ für } t_2 \leq t \end{cases} \quad k(1) < k_2$$

$$\dot{x} = F(t, x).$$

(b) Pendel (ohne Reibung)

$$\ddot{x} = -\lambda \sin(x)$$

(4) Räuber-Beute-Modell:

 $y(t)$: Räuber $z(t)$: Beute

$$\dot{y} = \alpha \cdot z \cdot y - \beta \cdot y$$

$$\dot{z} = -\gamma \cdot z \cdot y + \delta \cdot z$$

(Differentialgleichungssystem)

Es reicht, Algorithmen für autonome Anfangswertaufgaben 1. Ordnung zu entwickeln:

$$\dot{x} = F(x), x(t_0) = y^{(0)}$$

denn

(1) nicht autonome Differentialgleichungen \Rightarrow autonome Differentialgleichungen:

$$z := \begin{pmatrix} t \\ x \end{pmatrix} \Rightarrow \dot{z} = \begin{pmatrix} \dot{t} \\ \dot{x} \end{pmatrix} = \begin{pmatrix} 1 \\ \dot{x} \end{pmatrix} = \begin{pmatrix} 1 \\ F(z) \end{pmatrix} = G(z) \text{ (autonom)}$$

(2) Differentialgleichungen höherer Ordnung \Rightarrow Differentialgleichung 1. Ordnung:

$$\begin{aligned} \ddot{x} = F(x) &\rightarrow z_1 := x, z_2 := \dot{x} \\ \Rightarrow \dot{z}_1 = \dot{x} = z_2, \dot{z}_2 = \ddot{x} = F(x) = F(z_1). \end{aligned}$$

$$\Rightarrow \dot{z} = G(z) = \begin{pmatrix} z_2 \\ F(x) \end{pmatrix}.$$

b. Einteilung der Näherungsverfahren

$$\dot{x} = F(x), x(t_0) = y^{(0)}$$

Es sei $x(t)$ die exakte Lösung.

KONTINUIERLICHE VERFAHREN liefern eine (mindestens stetige) Funktion.

- Picard-Iteration:

$$\varphi_0(t) = y^{(0)}$$

$$\varphi_{n+1}(t) = y^{(0)} + \int_{t_0}^t F(\varphi_n(s)) ds.$$

DISKRETE VERFAHREN liefern in Gitterpunkten $t_0 < t_1 < t_2 < \dots < t_k < \dots$ Näherungswerte $y^{(k)} \approx x(t_k)$.

- (1) EINSCHRITT-VERFAHREN: $y^{(k+1)}$ wird aus $y^{(k)}$ berechnet.
- (2) MEHRSCHRITT-VERFAHREN: $y^{(k+1)}$ wird berechnet aus $y^{(k)}, y^{(k-1)}, \dots, y^{(k-l+1)}$ (l -Schritt-Verfahren).
Man benötigt hier l Startwerte.

c. Einschritt-Verfahren

Gegeben sei $h > 0$ (Schrittweite)

$$t_k := t_0 + kh.$$

Ein Einschritt-Verfahren basiert auf einer Verfahrensfunktion $V(h, x) \in C(\mathbb{R}^{n+1}; \mathbb{R})$.

$y^{(0)}$ erhält man aus der Anfangsbedingung.

$$y^{(k+1)} = y^{(k)} + h \cdot V(h, y^{(k)}).$$

DAS EULER-CAUCHY-VERFAHREN

Hier wählt man als Verfahrensfunktion:

$$V(h, x) = F(x)$$

$$\rightarrow \boxed{y^{(k+1)} = y^{(k)} + hF(y^{(k)})}$$

$$\dot{x} = F(t, x), \dot{z} = G(z).$$

Ohne Rückführung auf autonome Differentialgleichung:

$$y^{(k+1)} = y^{(k)} + hF(t_k, y^{(k)}).$$

Matlab:

```
function z = ecv(funk, y, h)
```

```
z = zeros(size(y))
```

$K = \text{feval}(\text{funkt}, y)$

(liefert funkt an der Stelle y)

$z = y + h \cdot K$

EXPLIZITE RUNGE-KUTTA-VERFAHREN

werden dargestellt durch eine VERFAHRENSMATRIX:

$$\begin{pmatrix} \beta_{21} \\ \beta_{31} & \beta_{32} \\ \vdots & & \ddots \\ \beta_{s1} & & & \beta_{s,s-1} \\ \gamma_1 & \dots & \dots & \gamma_{s-1} & \gamma_s \end{pmatrix} \text{ mit } \beta_{ij} \in \mathbb{R}, \gamma_i \in \mathbb{R}.$$

$$K^{(1)}(h, x) = F(x)$$

$$K^{(2)}(h, x) = F(x + h \cdot \beta_{21} K^{(1)}(h, x))$$

\vdots

$$K^{(j)} = F\left(x + h \cdot \sum_{i=1}^{j-1} \beta_{ij} K^{(i)}(h, x)\right) \text{ für } j = 1, \dots, s$$

$$V(h, x) = \sum_{j=1}^s \gamma_j K^{(j)}(h, x)$$

$$y^{(k+1)} = y^{(k)} + hV(h, y^{(k)})$$

explizites Runge-Kutta-Verfahren der Stufe s .

BEISPIELE:

(1) ($\bar{1}$): Euler-Cauchy-Verfahren ($s = 1$).

(2) $s = 2$:

• $\begin{pmatrix} \frac{1}{2} & \\ 0 & 1 \end{pmatrix}$: Halbschritt-Verfahren.

• $\begin{pmatrix} 1 & \\ \frac{1}{2} & \frac{1}{2} \end{pmatrix}$: zweistufiges Verfahren von Heun.

$$K^{(1)}(h, y^{(k)}) = F(y^{(k)})$$

$$K^{(2)}(h, y^{(k)}) = F\left(y^{(k)} + hK^{(1)}(h, y^{(k)})\right)$$

$$\Rightarrow y^{(k+1)} = y^{(k)} + \frac{h}{2} \left(K^{(1)}(h, y^{(k)}) + K^{(2)}(h, y^{(k)}) \right)$$

(3) $s = 3$

• $\begin{pmatrix} \frac{1}{2} & & \\ -1 & 2 & \\ \frac{1}{4} & 0 & \frac{3}{4} \end{pmatrix}$: Verfahren von Kutta.

(4) $s = 4$:

$$\bullet \begin{pmatrix} \frac{1}{2} & & & \\ 0 & \frac{1}{2} & & \\ 0 & 0 & 1 & \\ \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{pmatrix} : \text{klassisches Runge-Kutta-Verfahren.}$$

BEMERKUNG:

Bei allen guten Verfahren gilt $\sum_{i=1}^s \gamma_j = 1$.

IMPLIZITE RUNGE-KUTTA-VERFAHREN

werden dargestellt durch eine VERFAHRENSMATRIX:

$$\begin{pmatrix} \beta_{11} & \dots & \beta_{1s} \\ \vdots & \ddots & \vdots \\ \beta_{s1} & \dots & \beta_{ss} \\ \gamma_1 & \dots & \gamma_s \end{pmatrix} \text{ mit } \beta_{ij} \in \mathbb{R}, \gamma_i \in \mathbb{R}, \beta_{ij} \neq 0 \text{ f\u00fcr mindestens ein } j \geq i.$$

$$K^{(j)}(h, x) = F(x + h \cdot \sum_{i=1}^s \beta_{ji} K^{(i)}(h, x)) \text{ f\u00fcr } j = 1, \dots, s.$$

BEISPIEL:

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix} (s = 1)$$

$$K^{(1)} = F(\underbrace{y^{(k)} + hK^{(1)}}_{y^{(k+1)}})$$

$$y^{(k+1)} = y^{(k)} + h \cdot K^{(1)} = y^{(k)} + hF(y^{(k+1)})$$

implizites Euler-Verfahren (R\u00fcckw\u00e4rts-Euler)

d. Konsistenz und Konvergenz

DISKRETISIERUNGSFEHLER:

$$\dot{x} = F(x), x(t_0) = y^{(0)}.$$

$$\Delta(h, t_0, y^{(0)}) := \begin{cases} \frac{x(t_0+h) - x(t_0)}{h} & , \text{ falls } h \neq 0; \\ F(y^{(0)}) & , \text{ falls } h = 0. \end{cases}$$

Δ hei\u00dft wahre **Zuwachsfunktion**.

Einschritt-Verfahren: $y^{(1)} = y^{(0)} + hV(h, y^{(0)})$.

DEFINITION:

Die Funktion $\tau(h, t_0, y^{(0)}) := \Delta(h, t_0, y^{(0)}) - V(h, y^{(0)})$ heißt **lokaler Diskretisierungsfehler**.

$$x(t_1) - y^{(1)} = x(t_0 + h) - y^{(1)} = x(t_0 + h) - y^{(0)} - hV(h, y^{(0)}) = h \left(\frac{x(t_0+h) - x(t_0)}{h} - V(h, y^{(0)}) \right) = h \cdot \tau(h, t_0, y^{(0)}).$$

$x(t_1) - y^{(1)}$: Schrittfehler = Schrittweite · Diskretisierungsfehler.

BEZEICHNUNG:

B_N sei die Menge der Funktionen $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$ mit

- F ist stetig und beschränkt;
- es existieren alle partiellen Ableitungen der Ordnung N von F ; diese sind stetig und beschränkt.

DEFINITION:

Das von der Verfahrensfunktion $V(h, x)$ definierte Verfahren heißt **konsistent**, falls gilt

$$\lim_{h \rightarrow 0} \tau(h, t_0, y^{(0)}) = 0$$

für jeden Anfangswert (t_0, y_0) und jedes $F \in B_1$.

Das von $V(h, x)$ erzeugte Verfahren hat die **Konsistenzordnung** p , falls gilt

$$\tau(h, t_0, y^{(0)}) = O(h^p) \text{ für } h \rightarrow 0$$

für jeden Anfangswert $(t_0, y^{(0)})$ und jedes $F \in B_p$

Aus Konsistenz folgt also $V(0, y^{(0)}) = \lim_{h \rightarrow 0} V(h, y^{(0)}) = F(y^{(0)}) = V(0, y^{(0)})$

RUNGE-KUTTA-VERFAHREN:

$$V(h, x) = \sum_{i=1}^s \gamma_i K^{(i)}(h, x)$$

$$\text{Konsistenz: } F(x) = V(0, x) = \sum_{i=1}^s \gamma_i K^{(i)}(0, x) = \underbrace{\left(\sum_{i=1}^s \gamma_i \right)}_{\text{muss 1 sein}} F(x)$$

Fazit:

$$\text{Runge-Kutta-Verfahren konsistent} \Rightarrow \sum_{i=1}^s \gamma_i = 1.$$

BEISPIEL: Konsistenzordnung des Euler-Cauchy-Verfahrens (im skalaren Fall):
 $\dot{x} = f(x), x(t_0) = y^{(0)}$ (Taylorentwicklung)

$$\begin{aligned} x(t_0 + h) &= x(t_0) + \dot{x}(t_0) \cdot h + \ddot{x}(\xi) \frac{h^2}{2} \\ \dot{x}(t_0) &= f(x(t_0)) = f(y^{(0)}) \\ \ddot{x}(\xi) &= f'(x(\xi)) \cdot \dot{x}(\xi) = f'(x(\xi)) \cdot f(x(\xi)) \end{aligned}$$

$$\Rightarrow \Delta(h, t_0, y^{(0)}) = \frac{x(t_0+h) - x(t_0)}{h} = \dot{x}(t_0) + \ddot{x}(\xi) \cdot \frac{h}{2} = f(y^{(0)}) + \frac{h}{2} f'(x(\xi)) \cdot f(x(\xi))$$

$$\begin{aligned} V(h, y^{(0)}) &= f(y^{(0)}) \\ \Rightarrow \tau(h, t_0, y^{(0)}) &= \frac{h}{2} \cdot f'(x(\xi)) \cdot f(x(\xi)) = O(h) \end{aligned} \quad (\text{wegen } f \in B_1)$$

\Rightarrow Konsistenzordnung 1.

KONVERGENZ:

$$\dot{x} = F(x), x(t_0) = y^{(0)}.$$

Sei $x(t)$ die exakte Lösung. Wir berechnen mit einem Einschrittverfahren Näherungen in Gitterpunkten $t_i := t_0 + i \cdot h$ ($h > 0$ Schrittweite), $y_h^{(i)}$

Intervall: $[t_0, b]$.

Sei $R = R(h)$ der Index, welcher $t_{R-1} < b \leq t_R$ erfüllt.

$$e_h^{(i)} := x(t_i) - y_h^{(i)} \text{ für } i = 0, 1, \dots, R.$$

$$E(h) := \max\{\|e_h^{(i)}\| : i = 0, \dots, R\} \text{ (maximaler Fehler in } [t_0, t_R])$$

DEFINITION:

(1) Das Einschrittverfahren heißt **konvergent**, falls gilt

$$\lim_{h \rightarrow 0} E(h) = 0$$

für jedes $F \in B_1$, jeden Anfangswert $(t_0, y^{(0)})$ und jedes Intervall $[t_0, b]$.

(2) Das Einschrittverfahren hat die Konvergenzordnung p ($p \in \mathbb{N}_1$), falls

$$E(h) = O(h^p)$$

für jedes $F \in B_p$, jeden Anfangswert $(t_0, y^{(0)})$ und jedes Intervall $[t_0, b]$.

ANMERKUNG:

Ein konsistentes explizites Runge-Kutta-Verfahren ist auch konvergent.

Formelfehler: $E(h)$: h klein $\Rightarrow E(h)$ klein.

Rundungsfehler: $R(h)$: h klein $\Rightarrow R(h)$ groß.

Optimales h ??

e. Schrittweitensteuerung

$$\underbrace{t_k, y^k}_{\text{sei akzeptiert}} \rightarrow t_{k+1}, y^{(k+1)}$$

sei akzeptiert

Formelfehler = Fortpflanzungsfehler (Fehler, der schon in $t_k, y^{(k)}$ enthalten ist) + Schrittfehler.

IDEE: Kontrolliere den Schrittfehler.

Betrachte die Anfangswertaufgaben $\dot{x} = F(x)$, $x(t_k) = y^{(k)}$.

Exakte Lösung sei $u(t)$. $y_h^{(k+1)} = y^{(k)} + hV(h, y^{(k)})$.

Das Einschrittverfahren habe die Konsistenzordnung p .

$$u(t_h + h) - y_h^{(k+1)} = d \cdot h^{p+1} + O(h^{p+2}).$$

Wenn h klein ist, dann kann $O(h^{p+2})$ vernachlässigt werden. Man nennt deshalb $d \cdot h^{p+1}$ das **Fehlerhauptglied**.

Berechne zusätzlich zwei Schritte mit Schrittweite $\frac{h}{2}$:

$$z^{(k+1)} = y^{(k)} + \frac{h}{2} V\left(\frac{h}{2}, y^{(k)}\right)$$

$$z_{\frac{h}{2}}^{(k+1)} = z^{(k+1)} + \frac{h}{2} V\left(\frac{h}{2}, z^{(k+1)}\right)$$

SATZ:

Das Einschrittverfahren besitze die Konsistenzordnung p und $f \in B_{p+1}$. Dann gilt

$$u(t_k + h) - y_h^{(k+1)} = \frac{y_{\frac{h}{2}}^{(k+1)} - y_h^{(k+1)}}{1 - 2^{-p}} + O(h^{p+2}).$$

KOROLLAR:

Setzt man

$$y^{(k+1)} := y_h^{(k+1)} + \frac{y_{\frac{h}{2}}^{(k+1)} - y_h^{(k+1)}}{1 - 2^{-p}} = \frac{2^p y_{\frac{h}{2}}^{(k+1)} - y_h^{(k+1)}}{2^{p-1}},$$

so gilt: $u(t_k + h) - y^{(k+1)} = O(h^{p+2})$ (Konsistenzordnung wird um 1 größer).

Praktische Durchführung der Schrittweitensteuerung am Beispiel des klassischen Runge-Kutta-Verfahrens (also Konsistenzordnung 4) für nichtautonome Anfangswertaufgaben (skalar)

$$\text{sei gegeben } \dot{x} = f(t, x), x(t_0) = y^{(0)}$$

Sei $h > 0$ eine gegebene Schrittweite, $0 < \varepsilon_1 < \varepsilon_2$ (Genauigkeitsschranken). Sei $y^{(k)}$ eine akzeptierte Näherung zum Zeitpunkt t_k .

(1) Berechne mit dem aktuellen h die Näherungen $y_h^{(k+1)}$ und $y_{\frac{h}{2}}^{(k+1)}$.

(2) Berechne die Fehlerschätzung $E := \frac{16}{15} \left\| y_{\frac{h}{2}}^{(k+1)} - y_h^{(k+1)} \right\|$

(3) Gilt $E > \varepsilon_2$, so setze $h_{\text{neu}} := \frac{h}{2}$ und gehe zu (1).

(4) Ist $E \leq \varepsilon_2$, so setze gemäß dem obigen Korollar

$$y^{(k+1)} = \frac{16 y_{\frac{h}{2}}^{(k+1)} - y_h^{(k+1)}}{15} \text{ und } t_{k+1} := t_k + h.$$

$y^{(k+1)}$ wird zum Zeitpunkt t_k akzeptiert.

(5) Gilt $E < \varepsilon_1$, so setze $h = 2h$.

f. Stabilität

Test-Differentialgleichung: $\dot{x} = \lambda x, x(0) = 0, \lambda < 0$.

Die exakte Lösung: $x(t) = a \cdot \exp(\lambda t); \lim_{t \rightarrow \infty} x(t) = 0$.

Löse diese Anfangswertaufgabe mit einem Einschrittverfahren: Gitter $t_k = t_0 + k \cdot h$, Näherungswert $y^{(k)}$.

FORDERUNG: $\lim_{k \rightarrow \infty} y^{(k)}$ muss 0 sein (sonst ist das Verfahren schlecht oder die Schrittweite zu groß)

Z.B. Euler-Cauchy-Verfahren:

$$y^{(k)} = y^{(k-1)} + h \cdot \lambda y^{(k-1)} = (1 + \lambda h)y^{(k-1)} = (1 + \lambda h)^2 y^{(k-2)} = \dots = (1 + \lambda h)^k y^{(0)}.$$

Forderung: $(1 + \lambda h)^k y^{(0)} \rightarrow 0$ für $k \rightarrow \infty$.

$\Rightarrow |1 + \lambda h| < 1$ oder $\lambda h \in (-2, 0)$. \rightarrow Einschränkung an die Schrittweite.

Z.B. $\lambda = -10 \Rightarrow h < 0,2$.

Man setzt $W(z) := 1 + z$.

$$\Rightarrow y^{(k)} = W(\lambda h)^k \cdot a.$$

$W(z)$ heißt **Stabilitätsfunktion** des Euler-Cauchy-Verfahrens.

Allgemein wird jedem Runge-Kutta-Verfahren $\beta = \begin{pmatrix} \beta_{11} & \dots & \beta_{1s} \\ \vdots & \ddots & \vdots \\ \beta_{s1} & \dots & \beta_{ss} \\ \gamma_1 & \dots & \gamma_s \end{pmatrix}$ eine Stabili-

tätsfunktion $W_\beta(z)$ zugeordnet.

Dies ist die Funktion, welche $y^{(k)} = W_\beta(\lambda h)^k \cdot y^{(0)}$ (also $y^{(k)} = W_\beta(\lambda h) \cdot y^{(k-1)}$) erfüllt.

FORDERUNG:

$$\lim_{k \rightarrow \infty} y^{(k)} = 0 \Rightarrow |W_\beta(\lambda h)| < 1.$$

DEFINITION:

Die Menge $S_\beta := \{z \in \mathbb{R} : |W_\beta(z)| < 1\}$ heißt **Stabilitätsbereich** von β .

NACHTEIL der expliziten Runge-Kutta-Verfahren: der Stabilitätsbereich ist relativ klein, d.h. es gibt eine Einschränkung an die Schrittweite h .

IMPLIZITES EULER-VERFAHREN:

$$y^{(k)} = y^{(k-1)} + hf(y^{(k)}) = y^{(k-1)} + h\lambda y^{(k)}.$$

$$\Rightarrow y^{(k)} = \frac{1}{1-h\lambda} y^{(k-1)}.$$

$$\Rightarrow \text{Stabilitätsfunktion } W_\beta(z) = \frac{1}{1-z}.$$

$$|W_\beta(z)| < 1 \text{ f+r alle } z \in (-\infty, 0).$$

Also gibt es keine Einschränkung an die Schrittweite h .

8. Numerische Integration

a. Einleitung

Sei $f \in C[a, b]$. Dann existiert $\int_a^b f(x) dx$.

Ist $F(x)$ eine Stammfunktion zu $f(x)$, so gilt $\int_a^b f(x) dx = F(b) - F(a)$.

Z.B. günstig, falls f ein Polynom ist.

PROBLEMATISCH:

- Stammfunktion ist nicht explizit bekannt (z.B. $f(x) = e^{-x^2}$),
- Stammfunktion wird kompliziert (z.B. bei rationalen Funktionen),
- f ist nicht explizit bekannt, sondern nur $f(x_i)$ für $i = 0, \dots, m$ (z.B. bei Messergebnissen).

Deshalb braucht man Näherungsmethoden:

$$\int_a^b f(x) dx \approx \sum_{i=0}^m w_i f(x_i) \quad (\text{Quadratursumme})$$

Die x_i heißen **Stützstellen**, die w_i heißen **Gewichte**.

FORMELFEHLER (bzw. Quadraturfehler):

$$R[f] := \int_a^b f(x) dx - \sum_{i=0}^m w_i f(x_i)$$

$R : C[a, b] \rightarrow \mathbb{R}$ ist ein lineares Funktional (also $R[\alpha f + \beta g] = \alpha R[f] + \beta R[g]$).

Eine Quadratursumme (Quadraturformel) heißt **exakt** für ein $f \in C[a, b]$, falls $R[f] = 0$ gilt. Eine Quadraturformel hat den **Exaktheitsgrad** $n \in \mathbb{N}$, falls gilt

$$\begin{aligned} R[p] &= 0 \text{ für alle } p \in \Pi_n \text{ und} \\ R[q] &\neq 0 \text{ für ein } q \in \Pi_{n+1} \end{aligned}$$

Äquivalent dazu:

$$R[1] = R[x] = \dots = R[x^n] = 0 \text{ und } R[x^{n+1}] \neq 0$$

GÜNSTIGE WAHL EINER QUADRATURSUMME:

- *interpolatorische Quadraturformel*: Integrand f wird durch ein Polynom p approximiert (z.B. Interpolation):

$$\int_a^b f(x) dx \approx \int_a^b p(x) dx$$

- *zusammengesetzte Quadraturformel*: das Intervall $[a, b]$ wird in Teilintervalle zerlegt, auf jedem Teilintervall verwendet man eine eigene Quadratursumme
- *Romberg-Verfahren*: durch geeignete Kombination von zwei Quadratursummen erhält man eine neue (vielleicht bessere) Quadraturformel.

AFFINE TRANSFORMATIONEN:

Wie erhält man aus einer Formel für ein Intervall $[a, b]$ eine Formel für ein Intervall $[c, d]$?

$$\int_a^b f(x) dx \approx \sum_{i=0}^m w_i f(x_i).$$

Sei nun $g \in C[c, d]$; gesucht ist $\int_c^d g(s) ds$.

$$\varphi : [a, b] \rightarrow [c, d], \varphi(t) = \alpha t + \beta \text{ mit } \varphi(a) = c \text{ und } \varphi(b) = d. \Rightarrow \varphi(t) = \frac{d-c}{b-a}t + \frac{bc-ad}{b-a}.$$

$$f := g \circ \varphi \in C[a, b] \Rightarrow \int_c^d g(s) ds = \int_a^b g(\varphi(t)) \cdot \varphi'(t) dt = \alpha \int_a^b g(\alpha t + \beta) dt$$

$$\approx \alpha \cdot \sum_{i=0}^m w_i g(\alpha x_i + \beta) = \sum_{i=0}^m \bar{w}_i g(y_i) \text{ mit } \bar{w}_i := \alpha \cdot w_i.$$

b. Interpolatorische Quadraturformeln

Sei $a \leq t_0 < t_1 < \dots < t_m \leq b$ und $f \in C[a, b]$. Bilde das Lagrange-Interpolationspolynom $p(t)$ zu diesen Stützstellen $p(t) = \sum_{i=0}^m l_i f(t_i)$.

$$\int_a^b f(t) dt \approx \int_a^b p(t) dt = \int_a^b \left(\sum_{i=0}^m l_i(t) f(t_i) \right) dt = \sum_{i=0}^m \left[f(t_i) \cdot \int_a^b l_i(t) dt \right] = \sum_{i=0}^m w_i f(t_i) \text{ mit } w_i = \int_a^b l_i(t) dt.$$

SATZ:

Eine interpolatorische Quadraturformel (mit $m + 1$ Stützstellen) hat mindestens den Exaktheitsgrad m .

LEMMA:

Für eine interpolatorische Quadraturformel gilt:

$$\sum_{i=0}^m w_i = b - a.$$

BEISPIELE:

- *Mittelpunktsregel* ($m = 0, t_0 = \frac{a+b}{2}$):

$$\int_a^b f(t) dt = (b - a) \cdot f\left(\frac{a+b}{2}\right) + R[f] \text{ mit } R[f] = \frac{(b-a)^3}{24} f^{(2)}(\xi) \text{ mit einem } \xi \in (a, b).$$
 Man kann zeigen, dass der Exaktheitsgrad $M = 1$ ist.
- *Trapezregel* ($m = 1, t_0 = a, t_1 = b$):

$$\int_a^b f(t) dt = \frac{b-a}{2} (f(a) + f(b)) + R[f] \text{ mit } R[f] = -\frac{(b-a)^3}{12} f^{(2)}(\eta).$$
 Wieder lässt sich zeigen, dass der Exaktheitsgrad $M = 1$ ist.
- *Keplersche Fassregel* ($m = 2, t_0 = a, t_1 = \frac{a+b}{2}, t_2 = b$):

$$\int_a^b f(t) dt = \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) + R[f].$$
- *Newton-Cotes-Formel* (m beliebig, $t_i = a + i \cdot \frac{b-a}{m}$ (äquidistant))

c. Zusammengesetzte Quadraturformeln

Zerlege das Intervall $[a, b]$ in N Teilintervalle. $a = x_0 < x_1 < \dots < x_N = b$.

$$\int_a^b f(x) dx = \sum_{i=1}^N \int_{x_{i-1}}^{x_i} f(x) dx.$$

Verwende auf jedem Teilintervall eine einfache Quadraturformel (z.B. Trapezregel, Keplersche Fassregel, ...).

BEISPIEL (zusammengesetzte Trapezregel):

$$\int_{x_{i-1}}^{x_i} f(x) dx = \frac{x_i - x_{i-1}}{2} (f(x_{i-1}) + f(x_i)) + R_i[f].$$

Fall: alle Teilintervalle gleich: $x_i = a + i \cdot \frac{b-a}{N}$.

$$\Rightarrow \frac{x_i - x_{i-1}}{2} (f(x_{i-1}) + f(x_i)) + R_i[f] = \frac{b-a}{2N} (f(x_{i-1}) + f(x_i)) + R_i[f] \text{ mit } R_i[f] = -\frac{(b-a)^3}{12N^3} f^{(2)}(\xi_i).$$

$$\text{Also: } \int_a^b f(x) dx = \frac{b-a}{2N} \sum_{i=1}^N (f(x_{i-1}) + f(x_i)) + R[f] = \frac{b-a}{N} \left(\frac{1}{2} f(a) + \frac{1}{2} f(b) + \sum_{i=1}^{N-1} f(x_i) \right) + R[f].$$

Bei der zusammengesetzten Trapezregel gilt, falls $f \in C^2[a, b]$: $R[f] = -\frac{(b-a)^3}{12N^2} f^{(2)}(\xi)$ mit $\xi \in (a, b)$

FEHLERSCHRANKE:

$$|R[f]| \leq \frac{(b-a)^3}{12N^2} \|f^{(2)}\|_{[a,b]} \rightarrow 0 \text{ f\u00fcr } N \rightarrow \infty.$$

$$R[f] = \sum_{i=1}^N R_i[f] = \sum_{i=1}^N -\frac{(b-a)^3}{12N^3} f^{(2)}(\xi_i) = -\frac{(b-a)^3}{12N^3} \cdot \sum_{i=1}^N f^{(2)}(\xi_i) = \frac{1}{N} \sum_{i=1}^N f^{(2)}(\xi_i) = f^{(2)}(\eta).$$

FEHLERABSCH\u00c4TZUNG (zusammengesetzte Mittelpunktsregel):

$$R[f] = \frac{(b-a)^3}{24N^2} f^{(2)}(\xi) \text{ mit } \xi \in (a, b), \text{ falls } f \in C^2[a, b].$$

ANWENDUNG:

Sei $f \in C^2[a, b]$ und $\varepsilon > 0$. Bestimme einen N\u00e4herungswert J f\u00fcr $I = \int_a^b f(x) dx$ mit

$$|J - I| < \varepsilon.$$

W\u00e4hle dazu N so, dass $\frac{(b-a)^3}{12N^2} \|f^{(2)}\| < \varepsilon$ gilt. Berechne dann diesen N\u00e4herungswert mit Hilfe der zusammengesetzten Trapezregel mit N Teilintervallen.

ZUSAMMENGESetzte KEPLERSche FASSREGEL (SIMPSONREGEL):

$$\text{Teilintervalle: } R_i[f] = -\frac{(b-a)^5}{2880N^4} f^{(4)}(\xi_i), \text{ falls } f \in C^4[a, b].$$

Gesamt:

$$R[f] = -\frac{(b-a)^5}{2880N^4} f^{(4)}(\xi).$$

d. Das Prinzip der Extrapolation

Sei $\varphi : [-a, a] \rightarrow \mathbb{R}$ mit $\varphi(h) = \varphi(-h)$ für alle $h \in [-a, a]$.

Seien $0 < h_m < \dots < h_1 < h_0$ gegeben und seien $\varphi(h_0), \dots, \varphi(h_m)$ bekannt.

Gesucht: $\varphi(0)$

IDEE 1: bestimme das Interpolationspolynom \bar{p} zu den Stützpaaren $(h_i, \varphi(h_i))$.

$\rightarrow \bar{p}(0) \approx \varphi(0)$.

IDEE 2: bestimme das Interpolationspolynom \bar{q} zu den Stützpaaren $(\pm h_i, \varphi(h_i))$.

$\rightarrow \bar{q}(0) \approx \varphi(0)$.

IDEE 3: Bilde das Interpolationspolynom zu den Stützpaaren $(h_i^2, \varphi(h_i))$ für $i = 0, \dots, m$.

$\Rightarrow p(0) \approx \varphi(0)$

(beste Idee)

Durchführung mit Neville-Schema.

Sei nun $\varphi(h) = \underbrace{\tau_0}_{=\varphi(0)} + \tau_1 h^2 + \tau_2 h^4 + \tau_3 h^6 + \dots$

Sei $h_0 > h_1, h_1 = \alpha \cdot h_0$.

$p_0 = \tau_0 + \tau_1 h_0^2 + \tau_2 h_0^4 + \dots$

$p_1 = \tau_0 + \tau_1 h_1^2 + \tau_2 h_1^4 + \dots$

$$p_{0,1} = \frac{(0-h_1^2)p_0 - (0-h_0^2)p_1}{h_0^2 - h_1^2} = \frac{h_0^2(\tau_0 + \tau_1 h_1^2 + \tau_2 h_1^4 + \dots) - h_1^2(\tau_0 + \tau_1 h_0^2 + \tau_2 h_0^4 + \dots)}{h_0^2 - h_1^2} = \tau_0 - \tau_2 h_0^2 h_1^2 - \dots = \tau_0 + \bar{\tau}_2 h_0^4 + \bar{\tau}_3 h_0^6 + \dots \text{ mit } \bar{\tau}_2 = -\tau_2 \alpha^2.$$

FAZIT:

$$p_0 - \tau_0 = \tau_1 h_0^2 + \tau_2 h_0^4 + \dots$$

$$p_{01} - \tau_0 = \bar{\tau}_2 h_0^4 + \bar{\tau}_3 h_0^6 + \dots$$

$$p_{012} - \tau_0 = \bar{\tau}_3 h_0^6 + \dots$$

e. Das Romberg-Verfahren

$T(h) = \frac{h}{2} \cdot (f(a) + f(b)) + 2 \sum_{i=1}^{N-1} f(x_i)$ für $h = \frac{b-a}{N}$ (zusammengesetzte Trapezregel mit N Teilintervallen)

BEACHTEN: $T(h)$ ist nur für gewisse h definiert.

SATZ (SUMMENFORMEL VON EULER-MACLAURIN):

$$T(h) = \underbrace{\int_a^b f(t) dt}_{\tau_0} + \tau_1 h^2 + \tau_2 h^4 + \tau_3 h^6 + \dots$$

NEVILLE-SCHEMA:

Romberg-Tableau				
h_0^2	$p_0 = T(h_0) = T_0^0$	T_1^0	\dots	T_i^0
h_1^2	$p_0 = T_0^1$	T_1^1	\dots	
h_2^2	$p_0 = T_0^2$	T_2^1		
\vdots	\vdots	\vdots		
h_i^2	$p_0 = T_0^i$			

T_0^0 : Trapezregel

T_0^1 : Trapezregel mit zwei Teilintervallen

T_0^2 : Trapezregel mit 4 Teilintervallen

T_0^i : Trapezregel mit 2^i Teilintervallen

$$T_j^k = \frac{4^j \cdot T_{j-1}^{k+1} - T_{j-1}^k}{4^j - 1}$$

$$T_j^k = p_{k \dots k+j}(0); h_{k+j}^2 = \frac{h_k^2}{4^j}$$

PRAKTISCHE DURCHFÜHRUNG DES ROMBERG-VERFAHRENS:

$M(h_i)$: zusammengesetzte Mittelpunktsregel mit 2^i Teilintervallen.

LEMMA:

$$T(h_{i+1}) = \frac{1}{2}(T(h_i) + M(h_i)).$$

NUMERISCHES VERFAHREN:

Romberg-Tableau wird zeilenweise aufgebaut:

- maximale Zeilenzahl N vorgeben, so dass Algorithmus auch abbricht, falls das Romberg-Verfahren nicht konvergiert.

- T_0^0 wird berechnet ($T_0^0 = T(h_0)$)
- i -te Zeile: berechne $M(h_{i-1})$
 $T_0^i = T(h_i) = \frac{1}{2}(T(h_{i-1}) + M(h_{i-1})) = \frac{1}{2}(T_0^{i-1} + M(h_{i-1}))$
- für die restlichen Einträge in der i -ten Zeile berechne T_j^k für $j+k=i$: mit der Rekursionsformel.
- Abbruchbedingung: sei $\varepsilon > 0$ als Genauigkeitsschranke gegeben.
 $|T_{i-1}^1 - T_{i-1}^0| < \varepsilon \Rightarrow$ akzeptiere T_i^0
 $|T_{i-1}^1 - T_{i-1}^0| \geq \varepsilon \Rightarrow i = i + 1$, gehe zum dritten Punkt, falls $i \leq N$.

f. Gauß-Quadraturformeln

Maximaler Exaktheitsgrad bei $m + 1$ Stützstellen $t_0 < t_1 < \dots < t_m$?

Bekannt: interpolatorische Quadraturformel \Rightarrow Exaktheitsgrad $\geq m$.

LEMMA:

Eine Quadraturformel mit $m + 1$ Stützstellen hat höchstens den Exaktheitsgrad $2m + 1$.

BEWEIS:

$$q(t) := \prod_{i=0}^m (t - t_i)^2; q \in \Pi_{2m+2}$$

$$\Rightarrow q(t) \geq 0 \text{ und } \int_a^b q(t) dt < 0.$$

$$\sum_{j=0}^m w_j \underbrace{q(t_j)}_{=0} = 0$$

$$\Rightarrow R[q] \neq 0.$$

□

LEGENDRE-POLYNOME

$C[a, b]$; $f, g \in C[a, b]$:

$\langle f, g \rangle := \int_{-1}^1 f(t)g(t) dt$ definiert ein Skalarprodukt auf $C[-1, 1]$.

Zwei Polynome p, q heißen orthogonal, falls $\langle p, q \rangle = 0$.

SATZ:

Es gibt eine Folge $(P_n)_{n \in \mathbb{N}}$ orthogonaler Polynome mit positivem Hauptkoeffizienten und $\deg P_n = n$ (d.h. $\langle P_i, P_k \rangle = 0$ für $i \neq k$).

(ohne Beweis - Beweis z.B. mit Basis von $C[a, b]$ und Gram-Schmidt)

ANMERKUNG:

$\langle r, P_k \rangle = 0$ für alle $r \in \Pi_{j-1}$

Diese Polynome P_n heißen Legendre-Polynome.

Für die Legendre-Polynome gilt

$$P_{n+1}(t) = \frac{2n+1}{n+1}t \cdot P_n(t) - \frac{n}{n+1}P_{n-1}(t)$$

und $P_0(t) = 1; P_1(t) = t$.

(Bemerkung: mit dem Skalarprodukt $\int_{-1}^1 \frac{1}{\sqrt{1-t^2}} f(t)g(t) dt$ erhält man die Tschebyscheff-Polynome)

SATZ:

P_n hat in $(-1, 1)$ n einfache Nullstellen.

BEWEIS:

Annahme: P_n hat nur $l < n$ Nullstellen $-1 < t_1 < \dots < t_l < 1$ ungerader Ordnung.

$$q(t) = \prod_{j=1}^l (t - t_j) \Rightarrow \deg q < n.$$

$$P_n(t) \cdot q(t) \geq 0 \Rightarrow \int_{-1}^1 P_n(t)q(t) dt > 0$$

$$\text{Andererseits gilt } \langle q, P_n \rangle = 0, \text{ da } \deg q < \deg P_n. \langle q, P_n \rangle = 0 = \int_{-1}^1 P_n(t)q(t).$$

\Rightarrow Widerspruch; es gibt also mindestens n Nullstellen ungerader Ordnung, also n einfach Nullstellen in $(-1, 1)$.

□

Verwende nun die Nullstellen $-1 < t_0 < \dots < t_m < 1$ von P_{m+1} und bilde damit die interpolatorische Quadraturformel.

DEFINITION:

Diese Formeln heißen **Gauß-Quadraturformeln**.

SATZ:

Die Gauß-Quadraturformel mit $m + 1$ Stützstellen hat den Exaktheitsgrad $M = 2m + 1$.

BEWEIS:

Sei $q \in \Pi_{2m+1}$.

$\Rightarrow q = r \cdot P_{m+1} + s$ mit $r, s \in \Pi_m$

$\Rightarrow R[q] = R[r \cdot P_{m+1}] + R[s]$ mit $R[s] = 0$, da interpolatorische Quadraturformel und $s \in \Pi_m$.

$$R[r \cdot P_{m+1}] = \underbrace{\int_{-1}^1 r(t)P_{m+1}(t) dt}_{=(r, P_{m+1})=0, \text{ da } r \in \Pi_m} - \underbrace{\sum_{i=0}^m w_i r(t_i) \cdot \underbrace{P_{m+1}(t_i)}_{=0}}_{=0} = 0.$$

□

LEMMA:

Bei Gauß-Quadraturformeln sind alle Gewichte positiv (vgl. Übungsaufgabe auf Blatt 11).

BEISPIEL:

$m = 2$: (vgl. Übungsaufgabe 19 (a)).

9. Eigenwertaufgaben

a. Eigenwerte und Eigenvektoren

Sei $A \in \mathbb{R}^{n \times n}$. Ein $\lambda \in \mathbb{C}$ heißt Eigenwert von A , falls ein $x \in \mathbb{C}^n \setminus \{0\}$ gibt mit $Ax = \lambda x$. x heißt dann ein Eigenvektor.

Eigenwerte sind genau die Nullstellen vom charakteristischen Polynom $p(\lambda) := \det(\lambda I - A)$. In der Numerik vermeidet man allerdings diesen Weg (weil instabil!, vgl. Abschnitt 12).

Sei Q eine invertierbare Matrix. Dann sind A und $B := QAQ^{-1}$ ähnlich und haben damit dieselben Eigenwerte. (Eigenvektoren: x ist Eigenvektor zu $A \Rightarrow y := Qx$ ist Eigenvektor zu B)

$$B = \begin{pmatrix} b_{11} & \dots & * \\ & \ddots & \vdots \\ 0 & & b_{nn} \end{pmatrix} \Rightarrow b_{11}, \dots, b_{nn} \text{ sind die Eigenwerte.}$$

In der Numerik verwendet man Ähnlichkeitstransformationen, insbesondere orthogonale Ähnlichkeitstransformationen (da dann keine Probleme mit den Inversen Matrizen auftreten). Sei Q eine orthogonale Matrix, dann gilt $Q^{-1} = Q^T$, also $B = QAQ^T$.

b. Beispiel aus der Physik: lineare Kette

N Massen M , die durch $N + 1$ Federn linear verbunden sind; die beiden äußeren Federn werden fest eingespannt.



c_i = Federkonstante der i -ten Feder

l_0 = Länge der Federn (alle gleich lang)

L = Endpunkt-Anfangspunkt

$x_i(t)$: Ort der i -ten Masse zum Zeitpunkt t

$\dot{x}_i(t)$: Geschwindigkeit der i -ten Masse zum Zeitpunkt t

$\ddot{x}_i(t)$: Beschleunigung der i -ten Masse zum Zeitpunkt t

$F_i(t)$: Federkraft der i -ten Feder zum Zeitpunkt t

$K_i(t)$: Kraft, die auf die i -te Kugel wirkt zum Zeitpunkt t

Bewegungsgleichung (für die i -te Kugel): $M \cdot \ddot{x}_i(t) = K_i(t)$

$$K_i(t) = F_i(t) - F_{i-1}(t)$$

$$\text{Hook'sches Gesetz: } F_i(t) = c_i \cdot (x_{i+1}(t) - x_i(t) - l_0)$$

$$\begin{aligned} \Rightarrow \ddot{x}_i(t) &= \frac{1}{M} K_i(t) = \frac{1}{M} [c_i (x_{i+1}(t) - x_i(t) - l_0) - c_{i-1} (x_i(t) - x_{i-1}(t) - l_0)] \\ &= \frac{1}{M} [c_{i-1} x_{i-1}(t) - (c_i + c_{i-1}) x_i(t) + c_i x_{i+1}(t) + (c_{i-1} - c_i) l_0] \text{ für } i = 1, \dots, N \end{aligned}$$

$$x(t) = \begin{pmatrix} x_0(t) \\ \vdots \\ x_N(t) \end{pmatrix}; A = \frac{1}{M} \begin{pmatrix} c_0 + c_1 & -c_1 & & 0 \\ -c_1 & c_1 + c_2 & \ddots & \\ & \ddots & \ddots & -c_{N-1} \\ 0 & & -c_{N-1} & c_{N-1} + c_N \end{pmatrix};$$

$$b = \frac{1}{M} \begin{pmatrix} (c_0 - c_1) l_0 \\ \vdots \\ (c_{N-2} - c_{N-1}) l_0 \\ (c_{N-1} - c_N) l_0 + c_N L \end{pmatrix} \Rightarrow \ddot{x}(t) = b - Ax(t)$$

TRANSFORMATION AUF EIN HOMOGENES SYSTEM:

RUHELAGE: $0 = b - Ax$: sei \bar{x} die Lösung von $Ax = b$.

$$\Rightarrow z(t) = \bar{x} \text{ (spezielle Lösung)}$$

Sei $x(t)$ eine Lösung von $\ddot{x}(t) = b - Ax(t)$

$$\text{Setze } y(t) := x(t) - \bar{x}$$

$$\Rightarrow \ddot{y}(t) = \ddot{x}(t) = b - Ax(t) = A\bar{x} - Ax(t) = -A(x(t) - \bar{x}) = -Ay(t).$$

LÖSUNGEN DES HOMOGENEN SYSTEMS:

Falls alle $c_i > 0$, so ist A positiv definit (vgl. Übungsaufgabe 9c); das ist hier der Fall, da c_i Federkonstanten sind.

$\Rightarrow A$ hat lauter positive Eigenwerte $0 < \lambda_1 \leq \dots \leq \lambda_n$.

Seien $v^{(1)}, \dots, v^{(n)}$ zugehörige (linear unabhängige) Eigenvektoren.

BEHAUPTUNG:

Alle Lösungen von $\dot{y} = -Ay$ sind von der Form

$$y(t) = \sum_{j=1}^N (\alpha_j \cos(\sqrt{\lambda_j} \cdot t) \cdot v^{(j)} + \beta_j \sin(\sqrt{\lambda_j} \cdot t) \cdot v^{(j)}).$$

Beweis siehe Übungsblatt.

Die Lösungen des inhomogenen Systems sind dann

$$x(t) = \sum_{j=1}^N (\alpha_j \cos(\sqrt{\lambda_j} \cdot t) \cdot v^{(j)} + \beta_j \sin(\sqrt{\lambda_j} \cdot t) \cdot v^{(j)}) + \bar{x} \text{ mit } \alpha_j, \beta_j \in \mathbb{R}.$$

ANFANGSBEDINGUNGEN:

$x_i(0)$ für $i = 1, \dots, N$ und

$\dot{x}_i(0)$ für $i = 1, \dots, N$ legen α_j und β_j eindeutig fest.

NUMERISCHE PROBLEME:

- lineares Gleichungssystem mit A positiv definit (z.B. Cholesky-Zerlegung);
- Berechnung der Eigenwerte und Eigenvektoren von A ;
- zwei lineare Gleichungssysteme zur Berechnung der α_j bzw. der β_j .

c. Lokalisierung der Eigenwerte

$$\text{Sei } A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} \in \mathbb{R}^{n \times n}$$

$$K_i := \{z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{k=1, k \neq i}^n |a_{ik}|\} \text{ (Kreis in der komplexen Ebene um } a_{ii} \text{ mit Radius } \sum_{k=1, k \neq i}^n |a_{ik}|)$$

DEFINITION:

Diese K_i heißen Gerschgorin-Kreise.

SATZ: (von Gerschgorin)

$\bigcup_{j=1}^{K_j}$ enthält alle Eigenwerte von A .

(Beweis vgl. Übungsaufgabe)

BEISPIEL:

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

d. Transformation auf obere Hessenberg-Gestalt

DEFINITION:

Eine $(n \times n)$ -Matrix der Gestalt
$$\begin{pmatrix} * & \dots & \dots & * \\ * & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ 0 & & * & * \end{pmatrix}$$
 heißt obere **Hessenberg-Matrix**.

SATZ:

Sei $A \in \mathbb{R}^{n \times n}$. Dann gibt es eine orthogonale Matrix Q , so dass $B = QAQ^T$ eine obere Hessenberg-Matrix ist.

Q ist Produkt von $n - 2$ Householder-Matrizen (siehe QR-Zerlegung, 4.e.)

BEWEISIDEE:

$A = \left(\begin{array}{c|c} a_{11} & b \\ \hline c & A \end{array} \right)$; es gibt eine $(n-1) \times (n-1)$ Householder-Matrix \tilde{H}_1 mit $\tilde{H}_1 c = be_1$
($e_1 \in \mathbb{R}^{n-1}$).

Setze $H_1 := \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & \tilde{H}_1 \end{array} \right)$; diese ist ebenfalls eine Householder-Matrix.

Dann gilt $H_1 A = \left(\begin{array}{c|c} a_{11} & b \\ \hline \sigma & \bar{A}_1 \\ \hline 0 & \end{array} \right)$.

$H_1 A H_1^T = H_1 A H_1 = \left(\begin{array}{c|c} a_{11} & b \\ \hline \sigma & \bar{A}_1 \\ \hline 0 & \end{array} \right) \cdot \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & \tilde{H}_1 \end{array} \right) = \left(\begin{array}{c|ccc} a_{11} & b & & \\ \hline \sigma & x & \dots & x \\ & \vdots & & \vdots \\ 0 & x & \dots & x \end{array} \right)$.

KOROLLAR:

Sei $A \in \mathbb{R}^{n \times n}$ symmetrisch.

Dann gibt es eine orthogonale Matrix Q , so dass $QAQ^T = B$ mit B eine symmetri-

sche Tridiagonalmatrix, also $B = \begin{pmatrix} * & * & & 0 \\ * & \ddots & \ddots & \\ & \ddots & \ddots & * \\ 0 & & * & * \end{pmatrix}$

BEWEISIDEE:

$$H_1 A H_1 = \left(\begin{array}{c|c} a_{11} & c^T \\ \hline \sigma & \tilde{H}_1 A_1 \\ \hline 0 & \end{array} \right) \left(\begin{array}{c|c} 1 & 0 \\ \hline 0 & \tilde{H}_1 \\ \hline \end{array} \right) = \left(\begin{array}{c|c} a_{11} & \sigma & 0 \\ \hline \sigma & & \\ \hline 0 & & \bar{A}_1 \\ \hline \end{array} \right), \text{ wegen } c^T \tilde{H}_1 = c^T \tilde{H}_1^T =$$

$$(\tilde{H}_1 c)^T = (\sigma e_1)^T = \begin{pmatrix} \sigma & 0 & \dots & 0 \end{pmatrix}.$$

e. Das QR-Verfahren

Sei $A \in \mathbb{R}^{n \times n}$.QR-Zerlegung: $A = QR$.Verwende diese orthogonale Matrix Q zur Ähnlichkeitstransformation: $Q^T A Q = Q^T Q R Q = R Q$, d.h. RQ ist ähnlich zu QR . RQ heißt die **QR-Transformierte** von A .

SATZ:

Sei B eine obere Hessenberg-Matrix mit $B = QR$.Dann ist RQ ebenfalls eine obere Hessenberg-Matrix.

BEWEIS:

vgl. Übungsaufgabe 23.

KOROLLAR:

Sei B eine symmetrische Tridiagonalmatrix. Dann ist die QR-Transformierte ebenfalls eine symmetrische Tridiagonalmatrix.

BEWEIS:

$$B = QR, C = Q^T B Q, C^T = (Q^T B Q)^T = Q B^T Q^T = Q^T B Q = C.$$

 $\Rightarrow C$ ist symmetrisch und C ist eine obere Hessenberg-Matrix (Satz). $\Rightarrow C$ ist eine symmetrische Tridiagonalmatrix.

QR-ALGORITHMUS

 $A \in \mathbb{R}^{n \times n}$; $A^{(1)} := A$. Für $k = 1, 2, 3, \dots$:

$$A^{(k)} = Q^{(k)} \cdot R^{(k)}$$

$$A^{(k+1)} := R^{(k)} Q^{(k)}$$

(QR-Zerlegung)
(QR-Transformierte)

IN DER PRAXIS:

1. Stufe Transformation auf Hessenberg-Gestalt B ;
2. Stufe QR-Algorithmus mit $A^{(1)} = B$ (spart Rechenarbeit).

KONVERGENZ:

1. Alle Eigenwerte betragsmäßig verschieden:

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$$

$$A^{(k)} = \begin{pmatrix} a_{11}^{(k)} & * & \dots & * \\ a_{21}^{(k)} & \ddots & & \vdots \\ & \ddots & \ddots & * \\ 0 & & a_{n,n-1}^{(k)} & a_{nm}^{(k)} \end{pmatrix}.$$

Konvergenz falls $a_{21}^{(k)}, \dots, a_{n,n-1}^{(k)} \rightarrow 0$.

Es gilt $a_{j+1,j}^{(k)} \leq c \left| \frac{\lambda_{j+1}}{\lambda_j} \right|^k$ für $j = 1, \dots, n-1$.

Wegen $|\lambda_{j+1}| < |\lambda_j|$ gilt $\left| \frac{\lambda_{j+1}}{\lambda_j} \right|^k \rightarrow 0$ für $k \rightarrow \infty$.
 $\Rightarrow a_{j+1,j}^{(k)} \rightarrow 0$ für $k \rightarrow \infty$ und $j = 1, \dots, n-1$.

Aber Konvergenz kann langsam sein!

2. \rightarrow Möglichkeiten zur Konvergenzbeschleunigung!
3. Reduktion.
4. Komplexe Eigenwerte.

f. Das Verfahren von Hyman

Sei $B = \begin{pmatrix} b_{11} & \dots & \dots & b_{1n} \\ b_{21} & \ddots & & \vdots \\ & \ddots & \ddots & \vdots \\ 0 & & b_{n,n-1} & b_{nm} \end{pmatrix}$ eine obere Hessenberg-Matrix.

Sei $p(\lambda) := \det(B - \lambda I)$ das charakteristische Polynom.

ZIEL: Sei $\mu \in \mathbb{R}$, berechne $p(\mu)$ (ohne $p(\lambda)$ explizit auszurechnen).

Sei zusätzlich $b_{21} \neq 0, \dots, b_{n,n-1} \neq 0$.

ANMERKUNG: Eine Matrix mit diesen Eigenschaften heißt unzerlegbar.

Sei $\mu \in \mathbb{R}$. Betrachte das lineare Gleichungssystem

$$(B - \mu I)x = \alpha \cdot e_1 \quad (\alpha \text{ unbekannt})$$

$$\begin{aligned} (b_{11} - \mu)x_1 + b_{12}x_2 + \dots + b_{1n}x_n &= \alpha \\ b_{21}x_1 + (b_{22} - \mu)x_2 + \dots + b_{2n}x_n &= 0 \\ \vdots \quad \ddots \quad \vdots & \\ b_{n-1,n-2}x_{n-2} + (b_{n-1,n-1} - \mu)x_{n-1} + b_{n-1,n}x_n &= 0 \\ b_{n,n-1}x_{n-1} + (b_{nn} - \mu)x_n &= 0 \end{aligned}$$

ANMERKUNG:

Sei x ein Eigenvektor zur oberen Hessenberg-Matrix B . Dann gilt für $\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} : x_n \neq 0$.

(denn aus $(B - \lambda I)x = 0$ folgt dann $x_n = x_{n-1} = \dots = x_1 = 0$)

Setze also $x_n = 1$.

$$\Rightarrow x_{n-1} = \frac{-(b_{nn} - \mu)}{b_{n,n-1}} = x_{n-1}(\mu).$$

$$\Rightarrow x_{n-2} = \dots = x_{n-2}(\mu).$$

\vdots

$$\Rightarrow x_1 = \dots = x_1(\mu)$$

$$\Rightarrow \alpha = \dots = \alpha(\mu)$$

LEMMA:

$$\text{Es gilt } \alpha(\mu) = \frac{(-1)^{n+1}}{b_{21}b_{32} \dots b_{n,n-1}} p(\mu).$$

BEWEIS:

Cramersche Regel für $(B - \mu I)x = \alpha e_1$ ergibt:

$$x_n = \frac{\det C}{\det(B - \mu I)} \text{ mit } C = \begin{pmatrix} b_{11} - \mu & \dots & b_{1,n-1} & \alpha \\ b_{21} & \ddots & & 0 \\ & \ddots & \ddots & \vdots \\ 0 & & b_{n,n-1} & 0 \end{pmatrix}.$$

$$\det C = (-1)^{n+1} \alpha \cdot b_{21} b_{32} \cdot \dots \cdot b_{n,n-1}$$

(Entwicklungssatz Laplace)

$$\Rightarrow 1 = x_n = \frac{\det C}{p(\mu)} = \frac{(-1)^{n+1} \cdot \alpha \cdot b_{21} \cdot \dots \cdot b_{n,n-1}}{p(\mu)}$$

$$\Rightarrow \alpha(\mu) = \frac{(-1)^{n+1}}{b_{21} \cdot \dots \cdot b_{n,n-1}} \cdot p(\mu).$$

□

Nun Nullstellenverfahren, z.B. Bisektionsverfahren.

$\alpha(\mu)$ ist differenzierbar (nach μ). BERECHNUNG VON α' :

$$(B - \mu I)x = \alpha \cdot e_1$$

(α unbekannt)

$$(b_{11} - \mu)x_1(\mu) + b_{12}x_2(\mu) + \dots + b_{1n}x_n(\mu) = \alpha(\mu)$$

$$b_{21}x_1(\mu) + (b_{22} - \mu)x_2(\mu) + \dots + b_{2n}x_n(\mu) = 0$$

$$\vdots \quad \ddots \quad \vdots$$

$$b_{n-1,n-2}x_{n-2}(\mu) + (b_{n-1,n-1} - \mu)x_{n-1}(\mu) + b_{n-1,n}x_n(\mu) = 0$$

$$b_{n,n-1}x_{n-1}(\mu) + (b_{nn} - \mu)x_n(\mu) = 0$$

\Rightarrow

$$(b_{11} - \mu)x_1'(\mu) - x_1(\mu) + b_{12}x_2'(\mu) + \dots + b_{1n}x_n'(\mu) = \alpha'(\mu)$$

$$b_{21}x_1'(\mu) + (b_{22} - \mu)x_2'(\mu) - x_2(\mu) + \dots + b_{2n}x_n'(\mu) = 0$$

$$\vdots \quad \ddots \quad \vdots$$

$$b_{n,n-1}x_{n-1}'(\mu) - \underbrace{x_n(\mu)}_{=1} = 0$$

Unbekannt: $x_i'(\mu)$ und $\alpha'(\mu)$. Durch Auflösen erhält man diese.

\Rightarrow Newton-Verfahren ist auch anwendbar.

g. Eigenwerte und Eigenvektoren in Matlab

$d = \text{eig}(A)$: d ein Vektor, der die Eigenwerte enthält;

$[V,D] = \text{eig}(A)$: D eine Diagonalmatrix, die die Eigenwerte enthält, die Spalten von V enthalten die Eigenvektoren;

$[P,H] = \text{hess}(A)$: Hessenberg-Gestalt;

$[Q,R] = \text{qr}(A)$: QR-Zerlegung.

h. Vektornormen und Matrixnormen

VEKTORNORMEN AUF \mathbb{K}^n :

$\|\cdot\| : \mathbb{K}^n \rightarrow \mathbb{R}$ heißt Vektornorm auf \mathbb{K}^n , falls gilt

- (1) $\|x\| > 0$ für alle $x \in \mathbb{K}^n \setminus \{0\}$;
- (2) $\|\lambda x\| = |\lambda| \|x\|$ für alle $x \in \mathbb{K}^n$ und alle $\lambda \in \mathbb{K}$;
- (3) $\|x + y\| \leq \|x\| + \|y\|$ für alle $x, y \in \mathbb{K}^n$.

BEISPIELE:

- $\|x\|_\infty := \max\{|x_i| \mid i = 1, \dots, n\}$;
- $\|x\|_2 := \sqrt{\sum_{i=1}^n |x_i|^2}$;
- $\|x\|_1 := \sum_{i=1}^n |x_i|$.

MATRIXNORMEN:

$N(\cdot) : \mathbb{K}^{n \times n} \rightarrow \mathbb{R}$ heißt eine Matrixnorm, falls gilt:

- (1) $N(A) > 0$ für alle $x \in \mathbb{K}^{n \times n} \setminus \{0\}$;
- (2) $N(\lambda x) = |\lambda| N(A)$ für alle $x \in \mathbb{K}^{n \times n}$ und alle $\lambda \in \mathbb{K}$;
- (3) $N(A + B) \leq N(A) + N(B)$ für alle $A, B \in \mathbb{K}^{n \times n}$.

Eine Matrix heißt **submultiplikativ**, falls gilt:

- (1) $N(A \cdot B) \leq N(A) \cdot N(B)$ für alle $A, B \in \mathbb{K}^{n \times n}$.

Die Matrixnorm $N(\cdot)$ und die Vektornorm $\|\cdot\|$ heißen **verträglich**, falls gilt $\|Ax\| \leq N(A) \cdot \|x\|$ für alle $A \in \mathbb{K}^{n \times n}$ und alle $x \in \mathbb{K}^n$.

BEISPIELE:

$$A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}$$

- $N_E(A) := \sqrt{\sum_{i=1}^n \sum_{j=1}^n |a_{ij}|^2}$: **Euklidische Norm** oder **Schur-Norm**;
- $N_Z(A) := \max\{\sum_{k=1}^n |a_{ik}| : i = 1, \dots, n\}$: **Zeilensummennorm**;
- $N_S(A) := \max\{\sum_{i=1}^n |a_{ik}| : k = 1, \dots, n\}$: **Spaltensummennorm**.

Verträglich sind: N_E mit $\|\cdot\|_2$, N_Z mit $\|\cdot\|_\infty$, N_S mit $\|\cdot\|_1$.

ANMERKUNG:

Zu jeder Vektornorm $\|\cdot\|$ wird durch

$$\text{lub}(A) := \max\{\|Ax\| : x \in \mathbb{K}^n \text{ mit } \|x\| = 1\}$$

eine *submultiplikative, verträgliche* Matrixnorm definiert.

ANMERKUNG:

- (1) *lub* bedeutet **least upper bound**.
- (2) $\text{lub}(A)$ ist die kleinste Konstante c , welche $\|Ax\| \leq c \cdot \|x\|$ erfüllt.
- (3) $\text{lub}(A)$ heißt auch die von $\|\cdot\|$ induzierte Matrixnorm.

BEISPIELE:

- $\|x\|_\infty \Rightarrow \text{lub}_\infty(A) = N_Z(A)$;
- $\|x\|_1 \Rightarrow \text{lub}_1(A) = N_S(A)$;
- $\|x\|_2 \Rightarrow \text{lub}_2(A) = \sqrt{\rho(A \cdot A^T)} \neq N_E(A)$.

DEFINITION:

Sei $A \in \mathbb{K}^{n \times n}$. Dann heißt $\sigma(A) := \{\lambda \in \mathbb{C} : \lambda \text{ ist ein Eigenwert von } A\}$ das **Spektrum** von A .

$\rho(A) := \max\{|\lambda| : \lambda \in \sigma(A)\}$ heißt **Spektralradius** von A .

LEMMA:

Es gilt $\text{lub}(A) \geq \rho(A)$.

BEWEIS:

Sei $\lambda \in \sigma(A)$ und v ein zugehöriger Eigenvektor mit $\|v\| = 1$. Dann gilt $\|Av\| = \|\lambda v\| = |\lambda| \|v\| = |\lambda| \leq \rho(A)$ (weil $\rho = \max\{|\lambda|\}$).

 $\Rightarrow \text{lub}(A) \geq \rho(A)$.

□

LEMMA:

Sei $A \in \mathbb{R}^{n \times n}$ symmetrisch. Dann gilt $\text{lub}_2(A) = \rho(A)$.

BEWEIS:

Es gibt eine orthogonale Matrix Q mit $Q^T A Q = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} =: D$

Sei $x \in \mathbb{R}^n$ und $\|x\|_2 = 1 = \|x\|_2^2$. Setze $y = Qx$ ($\Rightarrow x = Qy$).

$\|y\|_2^2 = \|Q^T y\|_2^2 = (x^T Q) Q^T x = x^T x = \|x\|_2^2$.

$$\|Ax\|_2^2 = \|AQy\|_2^2 = (y^T Q^T A^T) A Q y = y^T Q^T A^T Q Q^T A Q y = y^T D^T D y = \|Dy\|_2^2 = \sum_{i=1}^n (\lambda_i y_i)^2 = \sum_{i=1}^n \lambda_i^2 y_i^2 \leq \rho(A) \cdot \sum_{i=1}^n y_i^2 = \rho(A)^2 \cdot \|y\|_2^2 = \rho(A)^2.$$

$\Rightarrow \|Ax\|_2 \leq \rho(A) \Rightarrow \text{lub}_2(A) \leq \rho(A)$.

Mit vorangegangenem Lemma gilt also $\text{lub}_2(A) = \rho(A)$.

□

KOROLLAR:

$$A = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix} \Rightarrow \text{lub}_2(A) = \max\{|a_{ii}| : i = 1, \dots, n\}.$$

Sei A invertierbar, so heißt die Größe

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

die **Konditionszahl** einer Matrix (es gilt immer $\text{cond}(A) \geq 1$).

i. Störempfindlichkeit bei Eigenwertaufgaben

LEMMA:

Seien $A, B \in \mathbb{R}^{n \times n}$. Sei $\lambda \in \mathbb{C}$ ein Eigenwert von A , aber *kein* Eigenwert von B . Dann gilt:

$$1 \leq \text{lub}((\lambda I - B)^{-1}(A - B)) \leq \text{lub}((\lambda I - B)^{-1})\text{lub}(A - B)$$

BEWEIS:

Sei v ein Eigenvektor zu λ mit $\|v\| = 1$.

$$(A - B)v = (\lambda I - B)v \Rightarrow (\lambda I - B)^{-1}(A - B)v = v.$$

$$\|(\lambda I - B)^{-1}(A - B)v\| = \|v\| = 1$$

$$\Rightarrow \text{lub}((\lambda I - B)^{-1}(A - B)) \geq 1.$$

□

SATZ:

Sei $B \in \mathbb{R}^{n \times n}$ diagonalisierbar mit Eigenwerten $\lambda_1(B), \dots, \lambda_n(B)$ (d.h. $B = PDP^{-1}$ mit $D = \text{diag}(\lambda_1(B), \dots, \lambda_n(B))$).

Sei $A \in \mathbb{R}^{n \times n}$ eine beliebige Matrix, so gibt es zu jedem Eigenwert λ von A einen Eigenwert $\lambda_j(B)$

$$|\lambda - \lambda_j(B)| \leq \text{cond}_2(P) \cdot \text{lub}_2(A - B)$$

BEWEIS:

Sei λ kein Eigenwert von B . Dann gilt $\text{lub}_2((\lambda I - B)^{-1}) = \text{lub}_2((\lambda PP^{-1} - PDP^{-1})^{-1}) = \text{lub}_2((P(\lambda I - D)P^{-1})^{-1}) = \text{lub}_2(P(\lambda I - D)^{-1}P^{-1}) \leq \text{lub}_2(P) \cdot \text{lub}_2((\lambda I - D)^{-1}) \cdot \text{lub}_2(P^{-1}) = \text{cond}_2(P) \cdot \text{lub}_2((\lambda I - D)^{-1}) = \text{cond}_2(P) \cdot \max \left\{ \left| \frac{1}{\lambda - \lambda_i(B)} \right| : i = 1, \dots, n \right\} = \text{cond}_2(P) \cdot \frac{1}{|\lambda - \lambda_j(B)|}$
(für ein bestimmtes $j \in \{1, \dots, n\}$).

Nach vorangegangenem Lemma gilt $1 \leq \text{cond}_2(P) \frac{1}{|\lambda - \lambda_j(B)|} \text{lub}_2(A - B)$.

□

ALSO:

$\text{lub}_2(P)$ klein \Rightarrow kleine Störung in B bedeutet, dass die Eigenwerte nahe beieinander liegen.

$\text{lub}_2(P)$ groß \Rightarrow kleine Störung in B kann total verschiedene Eigenwerte zur Folge haben.

LEMMA:

Ist Q eine orthogonale Matrix, so gilt

- (1) $\text{lub}_2(Q) = 1$;
- (2) $\text{cond}_2(Q) = 1$.

BEWEIS:

$$\text{lub}_2(Q) = \max\{\|Qx\|_2 : x \in \mathbb{K}^n, \|x\|_2 = 1\}.$$

Sei $\|x\|_2 = 1$.

$$\Rightarrow \|Qx\|_2^2 = (Qx)^T(Qx) = x^T Q^T Q x = x^T x = \|x\|_2^2 = 1 \Rightarrow \|Qx\|_2 = 1.$$

$$\Rightarrow \text{lub}_2(Q) = 1.$$

$$\text{cond}_2(Q) = \text{lub}_2(Q) \cdot \text{lub}_2(Q^{-1}) = \text{lub}_2(Q) \cdot \text{lub}_2(Q^T).$$

$$Q^T \text{ ist ebenfalls orthogonal} \Rightarrow \text{lub}_2(Q^T) = 1.$$

$$\Rightarrow \text{cond}_2(Q) = 1.$$

□

KOROLLAR:

Ist $B \in \mathbb{R}^{n \times n}$ symmetrisch und $A \in \mathbb{R}^{n \times n}$ beliebig, so gibt es zu jedem Eigenwert λ von A einen Eigenwert $\lambda_j(B)$ von B mit $|\lambda - \lambda_j(B)| \leq \text{lub}_2(A - B)$.

BEWEIS:

Es gibt eine orthogonale Ähnlichkeitstransformation.

Rest folgt mit obigem Lemma aus obigem Satz.

□

Bei symmetrischen Matrizen ist die Eigenwert-Aufgabe also störungsempfindlich.

STÖREMPFINDLICHKEIT BEI ÄHNLICHKEITSTRANSFORMATIONEN

$B = P^{-1}AP$ und Störung $A + \Delta A$.

Ähnlichkeitstransformation:

$$P^{-1}(A + \Delta A)P = \underbrace{P^{-1}AP}_{=B} + \underbrace{P^{-1}\Delta AP}_{=\Delta B}$$

$$\Rightarrow \text{lub}_2(\Delta B) = \text{lub}_2(P^{-1}\Delta AP) \leq \text{lub}_2(P^{-1})\text{lub}_2(P)\text{lub}_2(\Delta A) = \text{cond}_2(P) = \text{lub}_2(\Delta A)$$

Beim QR-Verfahren orthogonale Ähnlichkeitstransformation (Q ist orthogonal):
 $\Rightarrow \text{lub}_2(\Delta B) \leq \text{lub}_2(\Delta A)$

FAZIT: QR-Verfahren ist störungsempfindlich.

10. Iterative Verfahren

a. Iteration

Sei $D \subset \mathbb{R}^n$ und $\Phi : D \rightarrow D$. Ein $\xi \in D$ heißt Fixpunkt von Φ , falls $\Phi(\xi) = \xi$ gilt.

Fixpunkte werden iterativ berechnet. Startwert $x^{(0)} \in D$ berechne dann $x^{(i+1)} = \Phi(x^{(i)})$.

FRAGE: $\lim_{i \rightarrow \infty} x^{(i)} = \xi$??

Φ stetig: wenn Konvergenz $\Rightarrow \lim_{i \rightarrow \infty} x^{(i)}$ ist Fixpunkt.

NEWTON-VERFAHREN: $x^{(i+1)} = x^{(i)} - \frac{f(x^{(i)})}{f'(x^{(i)})}$; $\Phi(x) = x - \frac{f(x)}{f'(x)}$.

DEFINITION:

Sei (X, d) ein metrischer Raum. Eine Abbildung $\Phi : X \rightarrow X$ heißt kontrahierend, falls es ein $0 \leq L < 1$ gibt mit

$$d(\Phi(s), \Phi(t)) \leq L \cdot d(s, t) \text{ für alle } s, t \in X.$$

KONTRAKTIONSSATZ (BANACHSCHER FIXPUNKTSATZ):

Sei (X, d) vollständig und $\Phi : X \rightarrow X$ kontrahierend. Dann gilt

- (1) Es gibt genau einen Fixpunkt $\xi \in X$.
- (2) Für jeden beliebigen Startwert $x^{(0)} \in X$ konvergiert die Folge $x^{(i+1)} = \Phi(x^{(i)})$ gegen ξ .

NEWTON-VERFAHREN: $X := [a, b]$, $d(s, t) = |s - t|$.

ANMERKUNG:

Sei $\Phi : [a, b] \rightarrow [a, b]$ stetig differenzierbar.

$L := \|\Phi'\|_{[a,b]} = \max\{|\Phi'(s)| : s \in [a, b]\}$.

$$\Rightarrow |\Phi(s) - \Phi(t)| \leq L \cdot |s - t|.$$

(Denn aus Mittelwertsatz: $\frac{|\Phi(s) - \Phi(t)|}{|s - t|} = |\Phi'(\xi)| \leq \|\Phi'\|_{[a,b]}$)

b. Indirekte Verfahren für lineare Gleichungssysteme

$Ax = b \rightarrow$ FIXPUNKTFORMEN: $x = \Phi(x) = Tx + g$.

LEMMA:

Sei $A \in \mathbb{R}^{n \times n}$ regulär. Dann kann man durch Zeilenumformungen erreichen, dass alle Diagonalelemente von Null verschieden sind.

BEWEIS:

Induktion nach n :

- Sei A regulär. $n = 1$: $A = (a_{11})$ mit $a_{11} \neq 0 \Rightarrow$ Aussage klar.

- $n - 1 \Rightarrow n$: $A = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix}$

$$\det A = \sum_{i=1}^n (-1)^{i+1} a_{i1} \cdot \det A_{i1} \quad (\text{Entwicklung nach Laplace, 1. Spalte})$$

Es gibt mindestens ein $k \in \{1, \dots, n\}$ mit $a_{k1} \cdot \det A_{k1} \neq 0 \Rightarrow a_{k1} \neq 0$ und $\det A_{k1} \neq 0$.

Vertausche erste und k -te Zeile $\Rightarrow \tilde{A} = \left(\begin{array}{c|ccc} a_{k1} & * & \dots & * \\ * & & & \\ \vdots & & & \\ * & & & \end{array} \right)$

\tilde{A}_{k1} entsteht aus A_{k1} durch Zeilenvertauschung. $\Rightarrow \det \tilde{A}_{k1} \neq 0$.

\Rightarrow Behauptung

□

Sei nun $A \in \mathbb{R}^{n \times n}$, A regulär mit $a_{ii} \neq 0$ für $i = 1, \dots, n$.

$$\begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix} + \begin{pmatrix} 0 & \dots & 0 \\ & \ddots & \vdots \\ a_{n1} & & 0 \end{pmatrix} + \begin{pmatrix} 0 & & a_{1n} \\ \vdots & \ddots & \\ 0 & \dots & 0 \end{pmatrix}$$

$A = L + D + U$

Daraus ergeben sich zwei Verfahren:

$$1. Ax = b \Rightarrow (L + D + U)x = b \Rightarrow Dx + (L + U)x = b \Leftrightarrow Dx = -(L + U)x + b \Leftrightarrow x = -D^{-1}(L + U)x + D^{-1}b.$$

$$\text{Also } x = Tx + g \text{ mit } T = -D^{-1}(L + U) \text{ und } g = D^{-1}b.$$

$$\text{Nun: } x^{(0)} \in \mathbb{R}^n: x^{(i+1)} = Tx^{(i)} + g.$$

Dieses Verfahren heißt **Gesamtschrittverfahren** oder **Gauß-Jacobi-Verfahren**.

$$2. Ax = b \Leftrightarrow (L + U + D)x = b \Leftrightarrow (L + D)x + Ux = b \Leftrightarrow (L + D)x = -Ux + b \Rightarrow x = -(L + D)^{-1}Ux + (L + D)^{-1}b.$$

$$\text{Also } x = Tx + g \text{ mit } T = -(L + D)^{-1}U \text{ und } g = (L + D)^{-1}b.$$

$$\text{Nun: } x^{(0)} \in \mathbb{R}^n: x^{(i+1)} = Tx^{(i)} + g.$$

Dieses Verfahren heißt **Einzelschrittverfahren** oder **Gauß-Seidel-Verfahren**.

Praktische Durchführung des Einzelschrittverfahrens:

$(L + D)x^{(i+1)} = -Ux^{(i)} + b$. Löse dieses lineare Gleichungssystem (durch Vorwärtsauflösen).

Gesamtschrittverfahren:

$$\begin{pmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \rightarrow \begin{pmatrix} x_1^{(i+1)} \\ \vdots \\ x_n^{(i+1)} \end{pmatrix}$$

Einzelschrittverfahren:

$$\begin{pmatrix} x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \rightarrow \begin{pmatrix} x_1^{(i+1)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \rightarrow \begin{pmatrix} x_1^{(i+1)} \\ x_2^{(i+1)} \\ x_3^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \dots$$

ZUSAMMENFASSUNG:

$$Ax = b \Leftrightarrow x = \Phi(x) = Tx + g$$

$$A = L + D + U$$

$$T_G = -D^{-1}(L + U)$$

$$T_E = -(L + D)^{-1}U$$

Gesamtschritt

Einzelschritt

AUFWAND:

pro Iterationsschritt: n^2 Punktoperationen $x^{(i+1)} = Tx^{(i)} + g$
 gesamt: $k \cdot n^2$ (k Schritte)

direkte Verfahren: $O(n^3)$ (indirekte Verfahren also nur bei großen n sinnvoll)

KONVERGENZ:

$\Phi(x) = Tx + g$ - wann ist hier Φ kontrahierend?

$$\|\Phi(x) - \Phi(y)\| \leq L\|x - y\| \text{ mit } L < 1.$$

$$\|Tx + g - (Ty + g)\| = \|T(x - y)\| \leq N(T) \cdot \|x - y\| \leftarrow N(T) < 1.$$

($\|\cdot\|$ und $N(\cdot)$ seien verträglich).

Gilt $N(T) < 1$, so erhält man Konvergenz der Iteration $x^{(i+1)} = Tx^{(i)} + g$ für jeden Startwert $x^{(0)} \in \mathbb{K}^n$.

GENAU-DANN-BEDINGUNG:

Sei $\rho(T)$ der Spektralradius.

SATZ:

$x^{(i+1)} = Tx^{(i)} + g$ ist konvergent (gegen den Fixpunkt) für jeden Startwert $\Leftrightarrow \rho(T) < 1$.

BEWEIS:

„ \Rightarrow “ Sei \bar{x} ein Fixpunkt und $x^{(0)} \in \mathbb{K}^n$ ein Startwert.

$$x^{(i)} - \bar{x} = Tx^{(i-1)} - T\bar{x} = T(x^{(i-1)} - \bar{x}) = T^2(x^{(i-2)} - \bar{x}) = \dots = T^i(x^{(0)} - \bar{x}).$$

\Rightarrow Sei λ ein Eigenwert von T und y ein zugehöriger Eigenvektor. Wähle $x^{(0)} := y + \bar{x}$

$$\rightarrow x^{(i)} = Tx^{(i-1)} + g$$

Nach Voraussetzung gilt $\lim_{i \rightarrow \infty} (x^{(i)} - \bar{x}) = 0$.

$$\Rightarrow \lim_{i \rightarrow \infty} (T^i(x^{(0)} - \bar{x})) = \lim_{i \rightarrow \infty} (T^i(x^{(0)} - \bar{x})) = \lim_{i \rightarrow \infty} (T^i y) = \lim_{i \rightarrow \infty} (\lambda^i y) = 0$$

$$\Rightarrow |\lambda| < 1 \Rightarrow \rho(T) < 1.$$

„ \Leftarrow “ Sei $\rho(T) < 1$. Dann gibt es eine Vektornorm $\|\cdot\|$, so dass für die zugehörige lub-Norm gilt $\rho(T) \leq \text{lub}(T) < 1$. \Rightarrow Behauptung.

□

HINREICHENDE BEDINGUNG:

1. für Gesamtschrittverfahren: $T = T_G = -D^{-1}(L + U)$.

$$\begin{pmatrix} \frac{1}{a_{11}} & & \\ & \cdots & \\ & & \frac{1}{a_{nn}} \end{pmatrix} \begin{pmatrix} 0 & & a_{1n} \\ & \ddots & \\ a_{n1} & & 0 \end{pmatrix}$$

Zeilensummennorm $N_Z(T_G) < 1$

$$N_Z(T_G) = \max\left\{\frac{1}{|a_{ii}|} \cdot \sum_{k=1, k \neq i}^n |a_{ik}| : i = 1, \dots, n\right\} < 1.$$

$$\sum_{k=1, k \neq i}^n |a_{ik}| < |a_{ii}| \text{ für jedes } i = 1, \dots, n \quad (\text{Zeilensummenkriterium})$$

Entsprechend liefert die Spaltensummennorm $N_S(T_G)$ die Bedingung

$$\sum_{k=1, k \neq j} |a_{kj}| \leq |a_{jj}| \quad (\text{Spaltensummenkriterium})$$

2. für Einzelschrittverfahren:

SATZ:

Zeilensummen- und Spaltensummenkriterium sind auch hinreichend für die Konvergenz des Einzelschritt-Verfahrens.

BEWEIS:

für Zeilensummenkriterium:

Zeige $N_Z(T_E) \leq N_Z(T_G) (< 1 \text{ nach Voraussetzung})$

$$N_Z(\cdot) = \text{lub}_{\infty}(\cdot)$$

Sei $y \in \mathbb{K}^n$ und $z := T_E y$

Wir zeigen durch vollständige Induktion $|z_k| \leq \frac{1}{|a_{kk}|} \sum_{j=1, j \neq k}^n |a_{kj}| \|y\|_{\infty} \leq N_Z(T_G) \|y\|_{\infty}$.

$$(L + D)z = -Uy \quad (\text{Vorwärtsaufösen})$$

$$z_k = -\frac{1}{a_{kk}} \cdot \left(\sum_{j=1}^{k-1} a_{kj} z_j + \sum_{j=k+1}^n a_{kj} y_j \right)$$

$$\text{Induktionsanfang: } |z_1| < \frac{1}{|a_{11}|} \cdot \sum_{j=2}^n |a_{1j}| \underbrace{\|y\|_{\infty}}_{|y_j|} \leq \left(\frac{1}{|a_{11}|} \cdot \sum_{j=2}^n |a_{1j}| \right) \|y\|_{\infty} \leq N_Z(T_G) \|y\|_{\infty}$$

$$\begin{aligned}
\text{Induktionsschluss: } |z_k| &\leq \frac{1}{|a_{kk}|} \left(\sum_{j=1}^{k-1} |a_{kj}| |z_j| + \sum_{j=k+1}^n |a_{kj}| |y_j| \right) \\
&\leq \frac{1}{|a_{kk}|} \left(\left(\sum_{j=1}^{k-1} |a_{kj}| \right) \cdot \underbrace{N_Z(T_G)}_{<1, n.V.} \|y\|_\infty + \left(\sum_{j=k+1}^n |a_{kj}| \right) \|y\|_\infty \right) \\
&\leq \frac{1}{|a_{kk}|} \left(\sum_{j=1, j \neq k}^n |a_{kj}| \right) \cdot \|y\|_\infty \leq N_Z(T_G) \|y\|_\infty \\
&\Rightarrow \|T_E y\|_\infty = \|z\|_\infty \leq N_Z(T_G) \|y\|_\infty \\
&\Rightarrow \text{lub}_\infty(T_E) = N_Z(T_E) \leq N_Z(T_G).
\end{aligned}$$

Analog für Spaltensummenkriterium.

□

KONVERGENZGESCHWINDIGKEIT

hängt ab vom Spektralradius.

Methoden zur Konvergenzbeschleunigung

- Relaxationsverfahren „Verschieben der Eigenwerte“- hängt davon ab, wie die Eigenwerte verteilt sind

c. Nichtlineare Gleichungssysteme, mehrdimensionales Newtonverfahren

Sei $U \subset \mathbb{R}^n$ und $F : U \rightarrow \mathbb{R}^n$ hinreichend glatt.

Gesucht ist ein $x \in U$ mit $F(x) = 0$, d.h.

$$f_1(x_1, \dots, x_n) = 0$$

⋮

$$f_n(x_1, \dots, x_n)$$

$$\text{Sei } DF(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1}(x) & \dots & \frac{\partial f_1}{\partial x_n}(x) \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1}(x) & \dots & \frac{\partial f_n}{\partial x_n}(x) \end{pmatrix} \text{ die Funktionalmatrix.}$$

Taylor-Entwicklung:

$$F(x) = F(y) + DF(y) \cdot (x - y) + \text{Rest}$$

$$\Rightarrow 0 = F(y) + DF(y)(x - y) \Rightarrow x = y - DF(y)^{-1} \cdot F(y), \text{ falls } \det(DF(y)) \neq 0.$$

Wähle nun ein $x^{(0)} \in U$ und bilde

$$x^{(i)} = x^{(i-1)} - [DF(x^{(i-1)})]^{-1} \cdot F(x^{(i-1)})$$

(mehrdimensionales Newton-Verfahren).

PRAKTISCHE DURCHFÜHRUNG:

$$DF(x^{(i-1)}) \cdot (x^{(i)} - x^{(i-1)}) = -F(x^{(i-1)}) \quad (\text{gesucht: } x^{(i)}; \text{ sei } v := x^{(i)} - x^{(i-1)})$$

Löse also das lineare Gleichungssystem $DF(x^{(i-1)}) \cdot v = F(x^{(i-1)})$ und setze dann $x^{(i)} = v + x^{(i-1)}$

ABBRUCHBEDINGUNGEN:

$$\|F(x^{(i)})\|_{\infty} < \varepsilon \text{ und } \|x^{(i)} - x^{(i-1)}\|_{\infty} < \varepsilon.$$

MATLAB:

$Ax = b$ wird mit dem Befehl „ $x = b \setminus A;$ “ gelöst.

In jedem Schritt muss $A_k = DF(y^{(k)})$ berechnet werden; das sind n^2 Funktionsauswertungen. Deshalb gibt es einfachere Formen:

VEREINFACHTE FORM

Wähle für A_k eine feste Matrix, z.B. $A_k := A_0 = DF(y^{(0)})$ mit $y^{(0)}$ ist Startwert.

Löse dann $A_0 \cdot v = F(y^{(k)})$ z.B. mit LR-Zerlegung (die Zerlegung muss nur einmal gemacht werden).

Diese Methode heißt **Parallel-Verfahren**.

VERANSCHAULICHUNG IM FALL $n = 1$:

$$y^{(k+1)} = y^{(k)} - \frac{f(y^{(k)})}{f'(y^{(0)})}$$

KONVERGENZ DES NEWTON-VERFAHRENS:

SATZ: (lokale Konvergenz)

Sei $U \subset \mathbb{R}^n$ offen und $F : U \rightarrow \mathbb{R}^n$ zweimal stetig differenzierbar. Sei $\bar{u} \in U$ mit $F(\bar{u}) = 0$ und $DF(\bar{u})$ invertierbar.

Dann gilt:

- (1) Es gibt eine Kugel $K_\rho(\bar{u}) := \{x \in U : \|x - \bar{u}\|_2 \leq \rho\}$ mit $F(x) \neq 0$ für alle $x \in K_\rho(\bar{u}) \setminus \{\bar{u}\}$ (d.h. \bar{u} ist eine isolierte Nullstelle).
- (2) Für jeden Startwert $y^{(0)} \in K_\rho(\bar{u})$ konvergiert das Newton-Verfahren gegen \bar{u} .
- (3) es gilt quadratische Konvergenz.

In der Praxis nicht nützlich, weil die Nullstelle und ρ nicht bekannt sind.

11. Minimierung

a. Einleitung

Gesucht ist ein (globales oder lokales) Minimum von

$$h : \mathbb{R}^n \supset D \rightarrow \mathbb{R}.$$

Gesucht ist also ein $\bar{z} \in \mathbb{R}^n$ mit

- $h(\bar{z}) \leq h(x)$ für alle $x \in D$ (globales Minimum)
- $h(\bar{z}) \leq h(x)$ für alle $x \in U, U \subset D$ (lokales Minimum)

Wir beschränken uns auf die Suche der lokalen Minima.

b. Nichlinearer Ausgleich

$$h(z) = h(z_1, \dots, z_n) = \sum_{j=1}^m (s_j - G(t_j, z_1, \dots, z_n))^2 \quad (\text{vergleiche Abschnitt 5.a.})$$

(t_j, s_j) für $j = 1, \dots, m$, waren die Messwerte.

BEISPIELE:

- (1) Physik: Messungen (t_i, s_i) mit $i = 0, \dots, 20$.

$$G(t, \alpha_1, \alpha_2, \alpha_3, \alpha_4) = \alpha_1 \cdot \sin(\alpha_2 t + \alpha_3) + \alpha_4.$$

- (2) Biologie: Enzymaktionen $X \xrightarrow{\text{Enzym}} P$ (X : Substrat, P : Produkt)
Reaktionsgeschwindigkeit v hängt von der Konzentration des Substrats x ab.

$$v(x) = \frac{\mu \cdot x}{K+x}$$

Messungen (x_i, v_i) mit $i = 1, \dots, m$.

$$G(x_i, \mu, K) := \frac{\mu \cdot x_i}{K+x_i}.$$

Sei $h(x) = h(x_1, \dots, x_n)$ hinreichend glatt. NOTWENDIGE BEDINGUNG FÜR EIN LOKALES MINIMUM \bar{z} :

$$\left. \begin{array}{l} \frac{\partial}{\partial x_1} h(\bar{z}_1, \dots, \bar{z}_n) = 0 \\ \vdots \\ \frac{\partial}{\partial x_n} h(\bar{z}_1, \dots, \bar{z}_n) = 0 \end{array} \right\} \Rightarrow \text{grad } h(\bar{z}) = 0.$$

HINREICHENDE BEDINGUNG:
Hess $h(\bar{z})$ positiv definit.

c. Höhenlinien

$n = 2; h(z_1, z_2); h : \mathbb{R}^2 \rightarrow \mathbb{R}.$

Zu $c \in \mathbb{R}$ heißt $H_c := \{(z_1, z_2) \in \mathbb{R}^2 : h(z_1, z_2) = c\}$ **Höhenlinie** zum Niveau c .

ZEICHNEN VON HÖHENLINIEN durch Lösen von Anfangswertaufgaben:

$$\dot{z}_1 = h_{z_2}(z_1, z_2), z_1(0) = \alpha$$

$$\dot{z}_2 = -h_{z_1}(z_1, z_2), z_2(0) = \beta$$

Für die Lösung $(z_1(t), z_2(t))$ dieser Anfangswertaufgabe gilt

$$(z_1(t), z_2(t)) \in H_c \text{ mit } c = h(\alpha, \beta)$$

(vgl. Übungsaufgabe 27).

d. Abstiegsverfahren

$h : D \subset \mathbb{R}^n \rightarrow \mathbb{R}$: gesucht ist ein lokales Minimum.

h hinreichend glatt, differenzierbar, ...

Notwendig: $\text{grad}(h(\bar{z})) = 0$.

Sei $z \in D$ gegeben. Gibt es ein $z^* \in D$ mit $h(z^*) < h(z)$? Gibt es einen Vektor $p \in \mathbb{R}^n$ mit $h(z + tp) < h(z)$ für alle $t \in (0, t_0]$?

Dann heißt p eine **Abstiegsrichtung**.

$$z + tp = \begin{pmatrix} z_1 + tp_1 \\ \vdots \\ z_n + tp_n \end{pmatrix}.$$

Sei $I \supset (0, t_0]$. Setze $\varphi : I \rightarrow \mathbb{R} : \varphi(t) := h(z + tp)$ (bei festen z, p).

$$\psi : I \rightarrow \mathbb{R}^n: \psi(t) = \begin{pmatrix} z_1 + tp_1 \\ \vdots \\ z_n + tp_n \end{pmatrix} \Rightarrow \psi'(t) = \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix}.$$

$h : \mathbb{R}^n \rightarrow \mathbb{R}$.

$V(t) = h(\psi(t)) \Rightarrow \varphi$ ist differenzierbar (als Verknüpfung von Funktionen) und es gilt $\varphi'(t) = \text{grad } h(\psi(t)) \cdot \psi'(t)$

$$= \begin{pmatrix} h_{z_1}(\psi(t)) & \dots & h_{z_n}(\psi(t)) \end{pmatrix} \cdot \begin{pmatrix} p_1 \\ \vdots \\ p_n \end{pmatrix} = \sum_{i=1}^n h_{z_i}(\psi(t)) \cdot p_i = \sum_{i=1}^n h_{z_i}(z + tp) \cdot p_i.$$

$\varphi'(0) = \sum_{i=1}^n h_{z_i}(z) \cdot p_i$. $\varphi'(0) < 0 \Rightarrow p$ ist eine Abstiegsrichtung.

Gilt $\text{grad } h(z) \neq 0$, so ist $p := -\text{grad } h(z)$ eine Abstiegsrichtung, denn dann gilt

$$\varphi'(0) = -\sum_{i=1}^n h_{z_i}(z)^2 < 0.$$

BASISALGORITHMUS:

Gesucht ist ein $\bar{z} \in D \subset \mathbb{R}^n$ mit $\text{grad } h(\bar{z}) = 0$.

- 1) Sei $z^{(0)} \in D$ gegeben; $k = 0$
- 2) Berechne $\text{grad } h(z^{(k)})$; falls $\text{grad } h(z^{(k)}) = 0$, beende diese Schleife. Ansonsten wähle eine Abstiegsrichtung $p^{(k)}$ und eine Schrittweite $t_k > 0$ mit $h(z^{(k)} + t_k \cdot p^{(k)}) < h(z^{(k)})$.
- 3) Setze dann $z^{(k+1)} := z^{(k)} + t_k p^{(k)}$; $k = k + 1$ und gehe zu 2).

Ein Abstiegsverfahren besteht aus einer Richtungs- und aus einer Schrittweitenstrategie.

RICHTUNGSSTRATEGIE:

z.B. $p^{(k)} = -\text{grad } h(z^{(k)})$

(negativer Gradient)

→ **Gradientenverfahren**

SCHRITTWEITENSTRATEGIE:

Betrachte $\varphi(t) := h(z^{(k)} + tp^{(k)})$ ($\varphi : (0, t_0) \rightarrow \mathbb{R}, t_0$).

t_k als erstes Minimum von $\varphi(t)$.

→ **exakte Schrittweite.**

KONVERGENZ (FÜR DAS GRADIENTENVERFAHREN):

SATZ:

Sei $h \in \mathcal{C}^2(D; \mathbb{R})$, D offen und $z^{(0)} \in D$ der Startwert.

Die Niveaumenge $L_0 := \{x \in D : h(x) \leq h(z^{(0)})\}$ sei kompakt und besitze genau ein $\bar{z} \in L_0$ mit $\text{grad } h(\bar{z}) = 0$.

Dann konvergiert das Gradientenverfahren (mit Startwert $z^{(0)}$) gegen \bar{z} .

(vgl. J. Werner - Numerische Mathematik, Band 2).

e. Differentialgleichungsmethode

Idee: Nullstellen des Gradienten als stationäre Punkte einer Differentialgleichung.

$\dot{x} = F(x)$: sei \bar{y} eine Nullstelle von F , dann ist $x(t) = \bar{y}$ eine Lösung der Differentialgleichung (\rightarrow stationäre Punkte).

$$\begin{aligned} \dot{z} = -\text{grad } h(z)^T, \quad x(0) = z^{(0)}, \text{ also} \quad & \begin{array}{l} \dot{z}_1 = -h_{z_1}(z) \\ \vdots \\ \dot{z}_n = -h_{z_n}(z) \end{array} \end{aligned} \quad (*)$$

Sei $\Psi(t)$ die Lösung dieser Anfangswertaufgabe.

$\varphi(t) = h(\Psi(t))$ ist differenzierbar und es gilt

$$\varphi'(t) = \text{grad}(h(\Psi(t))) \cdot \varphi'(t) = -\text{grad } h(\Psi(t)) \cdot \text{grad } h(\Psi(t)) = \sum_{i=1}^n h_{z_i}(\Psi(t))^2 \leq 0.$$

$\Rightarrow \varphi(t)$ ist monoton fallend.

NUMERISCHES VERFAHREN:

- (1) Wähle Startwert $z^{(0)} \in D$; feste Schrittweite σ ; $k := 1$.
- (2) Berechne mit einem Näherungsverfahren (für Differentialgleichungen; am Besten implizite Verfahren) Näherungswerte der Anfangswertaufgabe (*) zum Zeitpunkt $k \cdot \sigma$.

$$z^{(k)} \approx \Psi(k \cdot \sigma)$$

- (3) Gilt $\|\text{grad } h(z^{(k)})\| < \varepsilon$ (für $\varepsilon > 0$ gegeben), so beende das Verfahren und akzeptiere $z^{(k)}$ als Näherungswert für \bar{z} .
Ansonsten setze $k = k + 1$ und gehe zu (2).

12. Stabilität und Störungsfragen

a. Einleitung

Im ganzen Kapitel bedeutet „klein“ nicht klein, sondern „betragsmäßig klein“.

STABILITÄT

Sei $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$. Gesucht ist ein \bar{x} mit $F(\bar{x}) = 0$. Bekannt ist ein y mit $F(y) \approx 0$.

FRAGE: Ist y ein guter Näherungswert an \bar{x} ?

BEZEICHNUNG: Die Größe $F(y)$ heißt **Defekt**.

Gute Stabilität: kleiner Defekt \Rightarrow gute Näherung.

Schlechte Stabilität: kleiner Defekt \Rightarrow gute Näherung.

STÖREMPFINDLICHKEIT

F ist nicht exakt bekannt, sondern nur eine Approximation \bar{F} . Gesucht ist ein \bar{x} mit $F(\bar{x}) = 0$. Bekannt ist ein y mit $\bar{F}(y) = 0$.

FRAGE: Folgt aus $F - \bar{F}$ klein auch $\bar{x} - y$ klein.

BEZEICHNUNG: **Störungsempfindlichkeit** bedeutet, dass aus einer kleiner Störung folgt, dass y ein guter Näherungswert an \bar{x} ist.

Die Störung wird durch einen p -dimensionalen Parameter-Vektor modelliert:

$$G : \mathbb{R}^{n+p} \rightarrow \mathbb{R}^n, (x, \mu) \mapsto G(x, \mu)$$

z.B. $p(x) = F(x)$ ein Polynom vom Grad n :

$$G(x, \mu) = G(x, \mu_0, \dots, \mu_n) = p(x) + \sum_{i=0}^n \mu_i x^i.$$

Das ungestörte Problem ist ein Ausgangsproblem für ein bestimmtes $\bar{\mu}$. Dann gilt $G(x, \bar{\mu}) = F(x)$.

AUSGANGSPROBLEM: $F(\bar{x}) = 0$ (hier $\bar{\mu} = 0$).

Es gibt ein $\bar{\mu}$ mit $G(\bar{x}, \bar{\mu}) = F(\bar{x}) = 0$

Berechnet ist $G(y, \eta) = 0$; Störung $\bar{\mu} - \eta$.

Störungsempfindlich: $\bar{\mu} - \eta$ klein $\Rightarrow y - \bar{x}$ klein.

b. Eine allgemeine Fehlerdarstellung

Gegeben sei eine stetig differenzierbare Funktion $G : \mathbb{R}^{n+p} \rightarrow \mathbb{R}^n, (u, v) \mapsto G(u, v)$.

$$G(u, v) = \begin{pmatrix} g_1(u_1, \dots, u_n, v_1, \dots, v_n) \\ \vdots \\ g_n(u_1, \dots, u_n, v_1, \dots, v_n) \end{pmatrix}$$

Für jede Komponente g_j gilt nach dem Mittelwertsatz:

$$\begin{aligned} & g_j(x, \mu) - g_j(y, \eta) \\ &= D_u g_j(y + t_j(x - y), \eta + t_j(\mu - \eta))(x - y) + D_v g_j(y + t_j(x - y), \eta + t_j(\mu - \eta))(\mu - \eta) \\ &= \sum_{i=1}^n \frac{\partial}{\partial u_i} g_j(y + t_j(x - y), \eta + t_j(\mu - \eta))(x_i - y_i) + \sum_{i=1}^p \frac{\partial}{\partial v_i} g_j(y + t_j(x - y), \eta + t_j(\mu - \eta))(\mu_i - \eta_i) \end{aligned}$$

mit $t_j \in (0, 1)$, D_u, D_v die partiellen Ableitungen nach u bzw. v .

Kompakt:

$$G(x, \mu) - G(y, \eta) = D_u G \cdot (x - y) + D_v G \cdot (\mu - \eta).$$

$D_u G$ hat die Komponente ($n \times n$ -Matrix)

$$\frac{\partial}{\partial u_i} g_j(\dots) \text{ für } i = 1, \dots, n; j = 1, \dots, n.$$

$D_v G$ hat die Komponente ($n \times p$ -Matrix)

$$\frac{\partial}{\partial v_j} g_i(\dots) \text{ für } i = 1, \dots, n; j = 1, \dots, p.$$

$p = 0$:

$$G(x) - G(y) = DG \cdot (x - y).$$

STABILITÄTSPROBLEM:

$F(\bar{x}) = 0$, Näherungswert y .

$$F(\bar{x}) - F(y) = DF \cdot (\bar{x} - y).$$

Problem: Zwischenstellen in DF sind nicht bekannt.

Deshalb wird $DF \approx DF(\bar{x})$ gesetzt.

$$\Rightarrow \underbrace{F(\bar{x}) - F(y)}_{=0} \approx DF(\bar{x}) \cdot (\bar{x} - y)$$

Sei nun $DF(\bar{x})$ invertierbar. Dann folgt $\bar{x} - y \approx -[DF(\bar{x})]^{-1} \cdot F(y)$.

FAZIT: Norm($[DF(\bar{x})]^{-1}$) klein \Rightarrow gute Stabilität.

STÖRUNGSPROBLEM:

$$0 = G(\bar{x}, \bar{\mu}) - G(y, \eta) \approx D_u G(\bar{x}, \bar{\mu}) \cdot (\bar{x} - y) + D_v G(\bar{x}, \bar{\mu}) \cdot (\bar{\mu} - \eta).$$

Sei $D_u G(\bar{x}, \bar{\mu})$ invertierbar.

$$(\bar{x} - y) \approx -[D_u G(\bar{x}, \bar{\mu})]^{-1} D_v G(\bar{x}, \bar{\mu}) \cdot (\bar{\mu} - \eta).$$

FAZIT: $\text{Norm}([D_u G(\bar{x}, \bar{\mu})]^{-1} D_v G(\bar{x}, \bar{\mu}))$ klein \Rightarrow störungsunempfindlich.

c. Stabilität und Störuneempfindlichkeit einer Nullstelle

Sei $f : \mathbb{R} \rightarrow \mathbb{R}$ stetig und differenzierbar, sei \bar{x} eine Nullstelle mit $f'(\bar{x}) \neq 0$:

$$y - \bar{x} \approx \frac{1}{f'(\bar{x})} f(y).$$

$$\text{Z.B. } |f'(\bar{x})| \approx 10^k, |f(y)| = 10^d.$$

$$\Rightarrow |y - \bar{x}| \approx 10^{d-k}.$$

BEACHTEN:

Das Stabilitätsverhalten in den einzelnen Nullstellen kann sehr unterschiedlich sein.

STÖRUNGSPROBLEM

$$p(\bar{x}) = 0$$

Berechnet: $f(\bar{x}, \bar{\mu}) = 0, f(y, \eta) = 0.$

Annahme: $f(x, \mu) = p(x) + \mu q(x)$ mit $p, q \in C^1[a, b].$

$\bar{\mu} = 0 \Rightarrow$ ungestörtes Problem.

$$f_x(x, \mu) = p'(x) + \mu q'(x) \Rightarrow f_x(\bar{x}, \bar{\mu}) = f_x(\bar{x}, 0) = p'(\bar{x})$$

$$f_\mu(x, \mu) = q(x) \Rightarrow f_\mu(\bar{x}, 0) = q(\bar{x}).$$

$$\Rightarrow \bar{x} - y \approx \frac{q(\bar{x})}{p'(\bar{x})} \cdot \eta$$

DEFINITION:

Die GröÙte $\left| \frac{q(\bar{x})}{p'(\bar{x})} \right|$ heißt **Konditionszahl** der Nullstelle $\bar{x}.$

FAZIT:

- kleine Konditionszahl \Rightarrow störungsempfindlich;
- große Konditionszahl \Rightarrow störungsempfindlich.

BEISPIEL: (vgl. Buch von Schwarz)

$$p(x) = \prod_{k=1}^{12} (x - k) = x^{12} \pm \dots - 6926634x^7 \pm \dots (= p_{12}x^{12} + \dots + p_1x + p_0)$$

Störung nur im j -ten Koeffizienten: $\mu \cdot q(x) = \mu \cdot x^j$.

Konditionszahlen:

$$\frac{q(k)}{p'(k)} = \frac{|p_j| \cdot k^j}{(k-1)!(12-k)!} = c_{k,j}$$

(z.B. Konditionszahlen bei $j = 7$ relativ groß, bei $j = 1$ relativ klein)

Störung bei $j = 7$: $p_7^* := p_7 + \eta p_7 = 6926634,001 \Rightarrow \eta \approx 1,44 \cdot 10^{-10}$.

$k = 9$: $9 - y_9 \approx c_{9,7} \cdot \eta \Rightarrow y_9 \approx 8,980 \dots$

d. Störempfindlichkeit linearer Gleichungssysteme

$Ax = b, A \in \mathbb{R}^{n \times n}, \det A \neq 0$

exakte Lösung: $\bar{x}: A\bar{x} = b$.

STÖRUNG IN DER RECHTEN SEITE:

$$F(x, \mu) = Ax - b - \mu, \text{ mit } \mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} \in \mathbb{R}^n.$$

$$D_x F(\bar{x}, 0) = A, D_\mu F(x, \mu) = -I \Rightarrow D_\mu F(\bar{x}, 0) = -I.$$

$$\Rightarrow \bar{x} - y \approx -[D_x F(\bar{x}, 0)]^{-1} D_\mu F(\bar{x}, 0) \cdot \eta = -A^{-1} \eta.$$

Hier gilt sogar: $\bar{x} - y = -A^{-1} \eta$, da $D_x F$ und $D_\mu F$ konstant sind (und dadurch die

Zwischenwerte nicht abweichen); vgl. 12.b, Störungsproblem.

$$\Rightarrow \|\bar{x} - y\| \leq \text{lub}(A^{-1})\eta \quad (\text{absoluter Fehler})$$

$$A\bar{x} = b \Rightarrow \|b\| \leq \text{lub}(A)\|\bar{x}\| \Rightarrow \frac{1}{\|\bar{x}\|} \leq \text{lub}(A) \cdot \frac{1}{\|b\|}$$

$$\Rightarrow \frac{\|\bar{x} - y\|}{\|\bar{x}\|} \leq \underbrace{\text{lub}(A^{-1}) \cdot \text{lub}(A)}_{\text{cond}(A)} \cdot \frac{\|\eta\|}{\|b\|} \quad (\text{relativer Fehler})$$