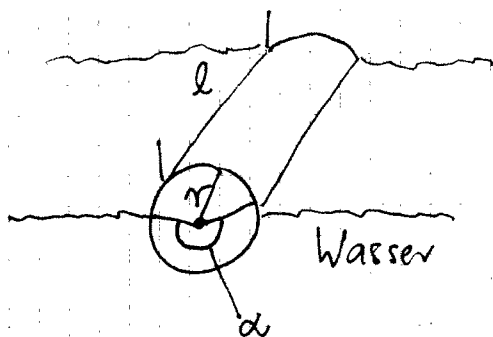


I Gleichungen

Ziel: Lösen von Glch.

Bsp 1: Berechne Temperatur T aus der Energie E in
Circuit Breaker

Bsp 2: Schwimmender Balken



Länge : l
 Radius : r
 spez. Gewicht Holz : ρ_1
 spez. " Wasser : ρ_0
 Winkel : α
 Gewicht : G

▷ Archimedes :

$$G = \pi r^2 l \rho_1 = \left[\frac{r^2}{2} \alpha + r^2 \sin\left(\pi - \frac{\alpha}{2}\right) \cos\left(\pi - \frac{\alpha}{2}\right) \right] l \rho_0$$

$$| \cdot \frac{2}{r^2 \rho_0}$$

$$\Rightarrow \alpha - \sin(\alpha) = 2\pi s \quad \text{mit } s = \frac{\rho_1}{\rho_0} \in (0, 1)$$

▷ Fixpunktform

$$\alpha = g(\alpha)$$

 g ist bekannt; hier:

$$\alpha = \sin(\alpha) + 2\pi s$$

$$g(\alpha) = \sin(\alpha) + 2\pi s$$

1. Iteration1.1 Fixpunkte

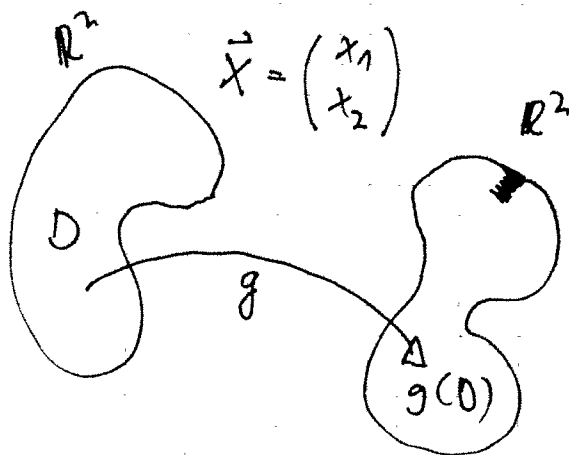
Def. bereich

Wertebereich

Gegeben : Fkt. $g : x \in D \mapsto g(x) \in g(D)$

1 dim: $D = [a, b]$

2 dim:



▷ Def: $\xi \in D$ heisst Fixpkt.

Falls: $\xi = g(\xi)$

▷ Fixpktiteration:

Wähle $x_0 \in D$

$$x_{k+1} = g(x_k) \quad \text{für } k = 0, 1, 2, \dots$$

Man benötigt: $g(D) \subset D$

Klar: Falls die Folge $\{x_k\}$ konvergiert gilt:

$$\xi = \lim_{k \rightarrow \infty} x_k \quad \text{ist}$$

Fixpkt von g .

▷ zum Bsp. 2: Tannenholz: $g = \frac{s_1}{s_0} = \frac{5}{6}$

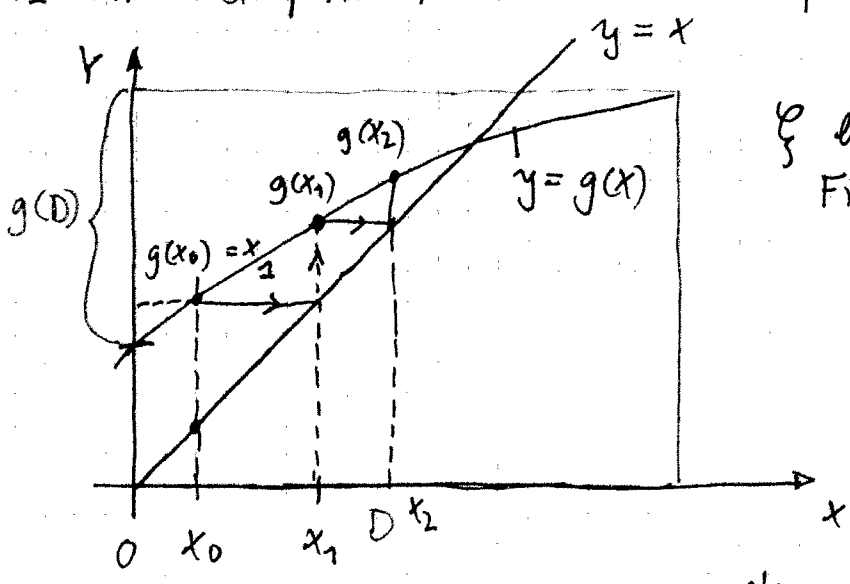
$$g(\alpha) = 2\pi \cdot \frac{5}{6} + \sin(\alpha)$$

Startwert: z. Bsp. $\alpha = 0$

$\Rightarrow \alpha_1 = 5.236 \dots$
 $\alpha_2 = 4,370 \dots$
 \vdots
 $\alpha_{23} = 4,314216594$
 $\alpha_{24} = 4,314216594$
 $\alpha_{24} = \text{Fixpunkt}$
 $\alpha \cong 247^\circ$

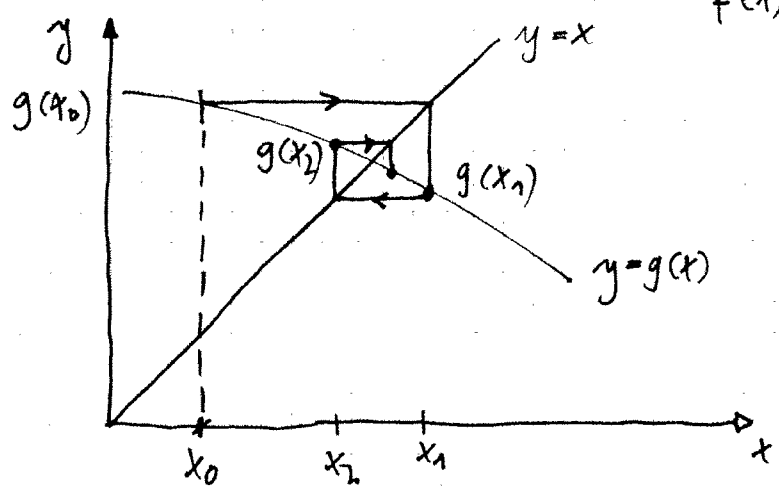
▷ 1dim: Graphisch

$f'(x) > 0$



ξ heisst anziehender Fixpkt.

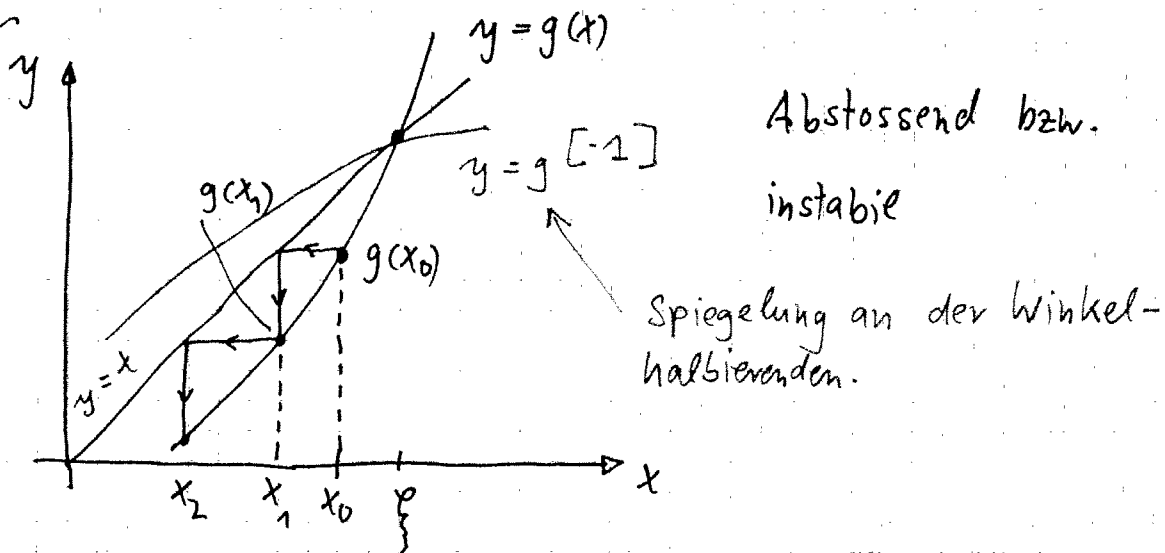
$f'(x) < 0$



ξ ist auch anziehend

$\xi = \lim_{k \rightarrow \infty} x_k$

Aber:



Fixpunkt: Die Iteration $x_{k+2} = g(x_k)$ konvergiert nicht gegen ξ .

Hingegen konvergiert $x_{k+1} = g^{-1}(x_k)$ für $k=0,1,\dots$

$g(x_{k+1}) = x_k$ (Rückwärtsiteration)

1.2. Fehlertheorie

Absoluten Fehler von x_k bezg. ξ

$$e_k := x_k - \xi$$

wobei $x_{k+1} = g(x_k)$

Idee: e_{k+1} durch e_k ausdrücken!

Taylorentwicklung im Punkt ξ

$$\begin{aligned} \Rightarrow e_{k+1} &= x_{k+1} - \xi = g(x_k) - \xi \\ &= g(\xi) + g'(\xi)e_k + \frac{1}{2}g''(\xi)e_k^2 - g(\xi) \end{aligned}$$

$$\Rightarrow e_{k+1} = g'(\xi) e_k + \sigma(e_k^2)$$

▷ Asymptotisches Fehlergesetz

Ann: $x_k \rightarrow \xi$, $\{x_k\}$ konvergiert

$\Rightarrow \xi$ ist attraktiver Fixpunkt:

$$e_k \rightarrow 0 \text{ mit } e_{k+1} = g'(\xi) e_k + \sigma(e_k^2)$$

▷ Def: Lineare Konvergenzgeschwindigkeit

Sei $\{x_k\}$ konvergent mit $\lim_{k \rightarrow \infty} x_k = \xi$,

$$e_k = x_k - \xi$$

$\{x_k\}$ heißt linear konvergent falls

$$\left| \frac{e_{k+1}}{e_k} \right| \leq K < 1 \text{ für } k > N.$$

$$\lim_{k \rightarrow \infty} \left| \frac{e_{k+1}}{e_k} \right| = q < 1 \text{ mit } q := \text{Konvergenz-faktor}$$

▷ Zusammenfassung:

ξ Fixpunkt von g , d.h. $\xi = g(\xi)$, x_0 Startwert
genügend nahe bei ξ . Fixpunktiteration

$$x_{k+1} = g(x_k)$$

▷ 2 Hauptfälle:

(i) $|g'(\xi)| < 1 \Leftrightarrow \xi$ attraktiv

(ii) $|g'(\xi)| > 1 \Leftrightarrow \xi$ abstoßend

▷ Bemerkung: Sei ξ abstoßend, d.h. $|g'(\xi)| > 1$

dann gilt:

$$|g^{[-1]}'| = \left| \frac{1}{g'(\xi)} \right| < 1$$

ξ ist anziehend von Funktion $g^{[-1]}$

1.3. Banachscher Fixpunktsatz

▷ Satz: Annahme:

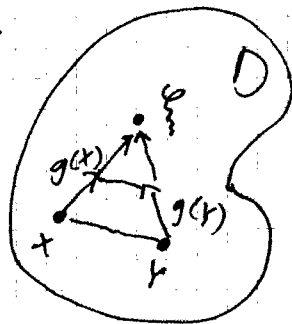
- (1) Sei $g(D) \subset D$ (Selbstabbildung)
- (2) Sei D abgeschlossen (Ränder sind noch in D enthalten)

g ist stetig in D

- (3) g ist kontrahierend, d.h. es gibt eine Zahl L mit $0 < L < 1$, mit $x, y \in D$

$$|g(y) - g(x)| \leq L|x - y| \quad (*)$$

Bsp: 2 dim:



Behauptung:

g hat genau einen Fixpunkt ξ . Die Fixpunktiteration konvergiert nach ξ für jeden beliebigen Startwert $x_0 \in D$

Es gilt:

$$|x_2 - \xi| \leq \frac{L^2}{1-L} |x_0 - x_1|$$

▷ Bemerkung 1:

(*) Existenz Lipschitzbedingung, $L < 1$ garantiert Kontraktion, L Existenz Lipschitzkonstante.

$$L = \max_{\substack{x, y \in D \\ x \neq y}} \left| \frac{g(y) - g(x)}{y - x} \right|$$

Maximale Sekantensteigung.

▷ Bemerkung 2:

Konvergenzgeschwindigkeit nahe bei ξ ist $L \cong g'(\xi)$

1.4. Newtonverfahren

Zweite Standardform einer Gleichung

$$f(x) = 0$$

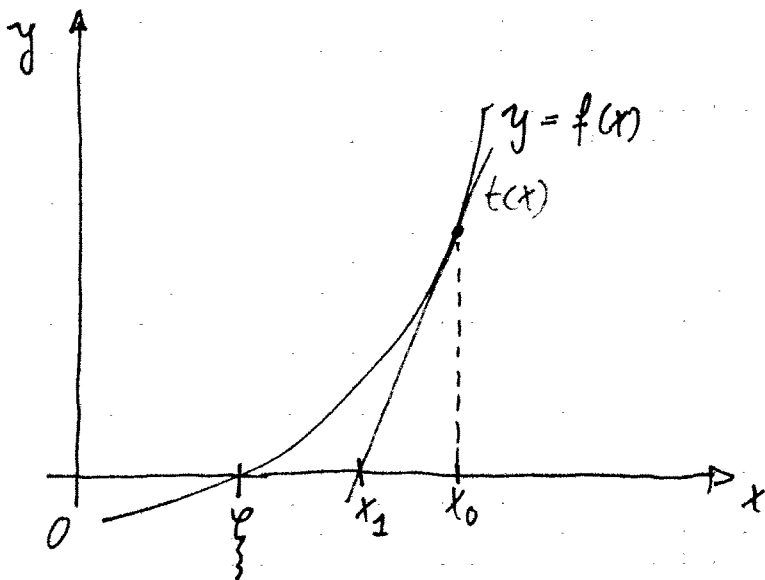
Gegeben: f ist differenzierbar

$$f: x \mapsto f(x)$$

Gesucht: Nullstelle ξ ; d. h. $f(\xi) = 0$

Nichtlineare Gleichung.

Idee: ersetze f durch lineare Funktion



$$t(x) = f(x_0) + f'(x_0)(x - x_0)$$

x_1 sodass $t(x_1) = 0$ linearisieren !!!

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

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

Voraussetzung: $f'(x_0) \neq 0$

25.02.08

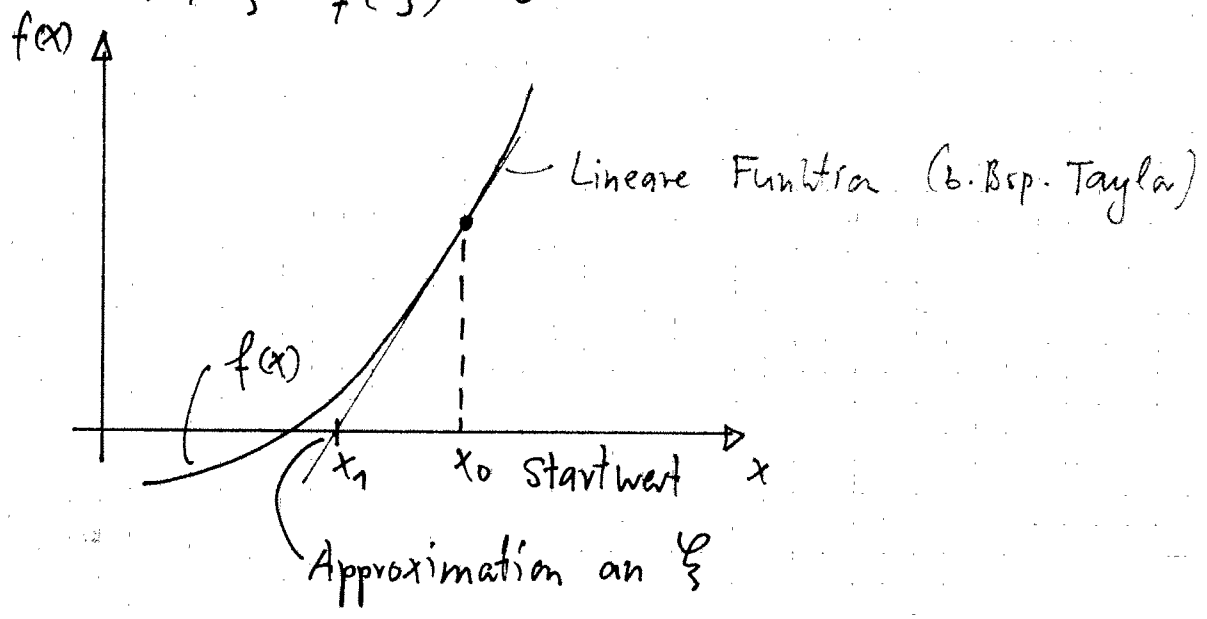
Newtonverfahren

1) $x \stackrel{?}{=} \cos(x)$; $f(x) = x - \cos(x)$

Nst von $f(x)$

2) $x^2 = 2$; $x = \sqrt{2}$; $f(x) = x^2 - 2$

Nst ξ $f(\xi) = 0$



Newton - Algorithmus

$$x_{h+1} = x_h - \frac{f(x_h)}{f'(x_h)} \quad ; \quad x_0 \text{ ist Schätzer}$$

$$h = 0, 1, 2, \dots$$

Bsp: $f(x) = x^2 - 2$; $f'(x) = 2x$

$$x_{h+1} = x_h - \frac{x_h^2 - 2}{2x_h} = \frac{1}{2}x_h + \frac{1}{x_h} \quad \forall x_h \neq 0$$

$$x_0 = 2$$

$$x_1 = \frac{3}{2} = 1,5$$

$$x_4 = \underline{1,414213562374}$$

$$x_2 = \underline{1,4166666\dots}$$

$$x_3 = \underline{1,41421568\dots}$$

Wie verringert sich der Fehler allgemein? (Konvergenztheorie)

$$x_k = \xi + e_k \quad ; \quad e_k : \text{Fehler}$$

$$x_{k+1} = \xi + e_{k+1}$$

$$x_{k+1} \stackrel{\uparrow}{=} \left(\xi + e_k \right) - \frac{f(\xi + e_k)}{f'(\xi + e_k)}$$

Verfahren

$$\stackrel{\uparrow}{=} \xi + \frac{f''(\xi)}{2f'(\xi)} e_k^2 + O(e_k^3) = \xi + e_{k+1}$$

Taylor Entw.

$$\Rightarrow e_{k+1} = \frac{f''(\xi)}{2f'(\xi)} e_k^2 \quad \left. \vphantom{\frac{f''(\xi)}{2f'(\xi)}} \right\} e_k = 10^{-3} \quad , \quad e_{k+1} \approx 10^{-6}$$

▷ Definition:

Falls für den Fehler $e_k \xrightarrow{k \rightarrow \infty} 0$ eines Verfahrens gilt

$$e_{k+1} \leq C e_k^p$$

mit $C = \text{const}$ (unabhängig von k) und $p \geq 1$,

so heißt das Verfahren konvergent mit Ordnung p .

Bsp: Newton : $p = 2$

allg. Fixpunktiteration : $p = 1$ ($C < 1$: Kontraktionsbedingung)

Abbruchkriterium:

Eigentlich : $|x_{n+1} - \xi| < \delta$; δ : Toleranz

Wir kennen aber ξ nicht!

Ersatz: $|x_{n+1} - x_n| < \delta$ (absolute Fehler)

$$\frac{|x_{n+1} - x_n|}{|x_n|} < \delta \quad (\text{relative Fehler})$$

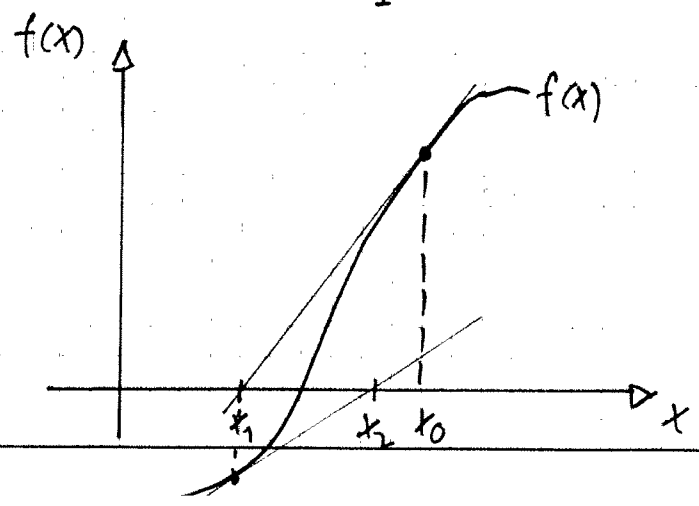
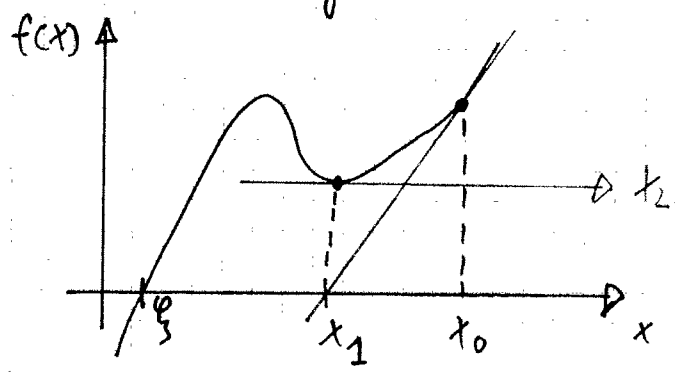
Kombination:

$$|x_{n+1} - x_n| < \delta_{rel} |x_{n+1}| + \delta_{abs}$$

Zusätzlich : Iteration $< 10 \dots 15$

Konvergenzprobleme:

Newton konvergiert, falls der Schätzer x_0 "nahe genug" bei der Lösung ist. In der Praxis ist das heikel.



(schlechte Zeichnung siehe Skript)

zyklisch!

Oft ist $f'(x)$ mühsam zu berechnen:

Strategien:

1) Wir approximieren $f'(x) \approx \frac{f(x+h) - f(x)}{h}$ mit h klein.

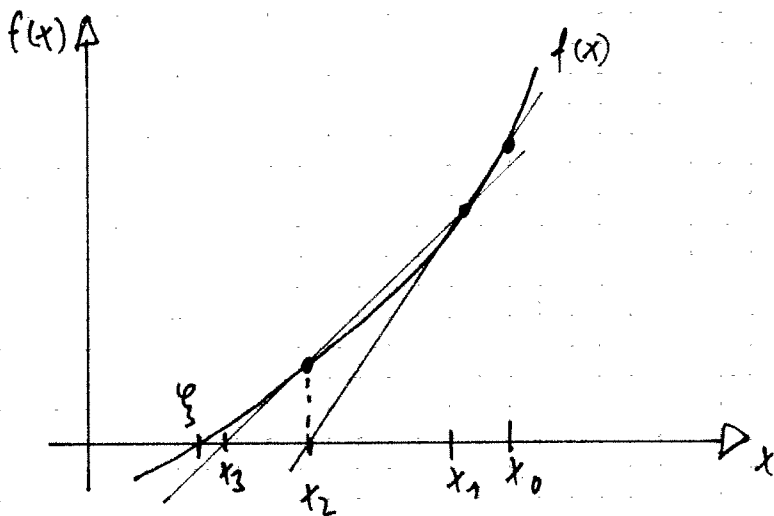
2) Oder $f'(x) \approx f'(x_0)$ "Quasi-Newton"

für alle h die gleiche Ableitung. $p=1$

3) $f'(x_n) = \frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$ "Sekantenverfahren"

Trick: Wir benutzen die Auswertung $f(x_{n-1})$ aus einem früheren Schritt.

Dafür brauchen wir zwei Startwerte. $x_0, x_{-1} = x_0 - h$



Es lässt sich zeigen $p = \frac{\sqrt{5}+1}{2} \approx 1,61 \dots$
Goldene Schnitt!

Also: $p_{\text{quasi}} < p_{\text{sek.}} < p_{\text{Newton}}$

Zusätzlich ist Sekanten Verf. effizienter als Newton Verf.

Effizienz = $\frac{\text{Fehlerreduktion}}{\text{Fehlerauswertung}}$
pro Schritt

hier: $E = p \frac{1}{m}$

m : Funktionsauswertung!

Newton $m=2$ $p=2$ $E=\sqrt{2}$

Sekante $m=1$ $p=1,6$ $E=1,6$

Newton in N Dimensionen

Analog: Nullstelle von $\vec{f}: \mathbb{R}^N \rightarrow \mathbb{R}^N$;

$$\begin{cases} \underline{x} \mapsto \underline{f}(\underline{x}) \\ \underline{\xi} = (\xi_1, \xi_2, \dots, \xi_N) \end{cases}$$

d.h. $\underline{f}(\underline{\xi}) = 0$ also

$$\begin{aligned} f_1(\underline{\xi}) &= 0 \\ f_2(\underline{\xi}) &= 0 \\ &\vdots \\ f_N(\underline{\xi}) &= 0 \end{aligned}$$

Bsp: $x^2 - y^2 = 0$

$$xy + 1 = 0$$

$$f_1(x, y) = x^2 - y^2$$

$$f_2(x, y) = xy + 1$$

$\underline{\xi} \equiv (x, y)$: Lösung: $\underline{\xi} = (1, -1)$ oder $(-1, 1)$

Approximation durch lineare Fkt. (Taylor) \leadsto Newton

$$f_i(\underline{x}_0 + \Delta \underline{x}) \stackrel{\text{Taylor}}{\approx} f_i(\underline{x}_0) + \sum_{j=1}^N \frac{\partial f_i}{\partial x_j}(\underline{x}_0) \Delta x_j + o(\|\Delta \underline{x}\|^2)$$

$$= \underline{f} + \underline{J} \Delta \underline{x} = 0$$

$i = 1, \dots, N$

\underline{x}_0 Schätzwert

$\stackrel{!}{=} 0$ (schwierig)
 $\Delta x = ?$

$\stackrel{!}{=} 0$ (leicht)
 $\Delta x = ?$

Jacobi-Matrix: $\underline{J} = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots \\ \frac{\partial f_2}{\partial x_1} & \dots & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$

Df

Newton : $\left\{ \begin{array}{l} \text{d) Löse lineares Gleich. system} \\ \text{für } k=0,1,2,\dots \\ \text{ii } J(\underline{x}_k) \Delta \underline{x} = -f(\underline{x}_k) \text{ nach } \Delta \underline{x} \\ \text{iii } \underline{x}_{k+1} = \underline{x}_k + \Delta \underline{x} \end{array} \right.$

\underline{x}_0 Schätzer

Schwierig : J ist teuer zu rechnen
typischer Weise mit Differenzenquotient

Gleich. system

Merke: ein nichtlineares $\checkmark \rightarrow$ wird zu einer Iteration von lin. Gleich. systemen

Lineare Gleichungssysteme (LGS)

Die Lsg. von LGS ist Grundbaustein von sehr vielen Berechnungen.

Abstrakt : $Ax = b$

$A \in \mathbb{R}^{N \times N}$; $x, b \in \mathbb{R}^N$

Erinnere : A regulär ($\det \neq 0$)

\Leftrightarrow genau eine Lösung \underline{x}

A singular ($\det = 0$)

\Leftrightarrow entweder keine oder unendlich viele Lsg.

Klassischer Algorithmus: Gauß - Elimination

numerische Varianten: LR - Zerlegung

Wir zerlegen ("faktorisieren") $A = L \cdot R$, $L, R \in \mathbb{R}^{N \times N}$

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

Linke untere Dreiecksmatrix

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

Rechte obere Dreiecksmatrix

(Englisch: LU - Decompositia)

Wozu? $Ax = b$

$$\iff \underbrace{(LR)}_Y x = b$$

Wir lösen $LY = b$ (leicht!)

$$\begin{pmatrix} 1 & & & 0 \\ & * & & \\ & & \ddots & \\ & & & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \begin{pmatrix} b \end{pmatrix} \text{ "Vorwärts" einsetzen.}$$

$Rx = Y$ (auch leicht)

$$\begin{pmatrix} * & & & \\ & * & & \\ & & \ddots & \\ 0 & & & * \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_N \end{pmatrix} \text{ "Rückwärts" einsetzen}$$

Das ist übrigens genau was Gauss macht!

Die Zerlegung ist nützlich falls mehrere rechte Seiten gelöst werden muss.

$$Ax = b_i \text{ mit } i = 1, 2, \dots, J$$

$$A = \begin{array}{c|c} \alpha & q^T \\ \hline p & \tilde{A} \end{array} \stackrel{!}{=} \begin{array}{c|c} 1 & -o- \\ \hline l & \tilde{L} \end{array} \stackrel{R}{=} \begin{array}{c|c} s & r^T \\ \hline 0 & \tilde{R} \\ 1 & \end{array} =$$

$p, q \in \mathbb{R}^N$

Wir iterieren über die Zeilen (Gauß)

$$= \begin{array}{c|c} s & r^T \\ \hline s \cdot l & l r^T + \tilde{L} \tilde{R} \end{array}$$

$$\underbrace{l, r \in \mathbb{R}^N}_{?} \quad \underbrace{s \in \mathbb{R}}_{?}$$

Es folgt: $s := \alpha$ $r := q$
 $l := \frac{1}{s} p$

Matlab!

Problem $\alpha \neq 0$ notwendig

Ansonsten Zeilen vertauschen!

$$\rightsquigarrow PA = L \cdot R$$

↑
Vertauschungsmatrix

$$PA = P \begin{pmatrix} a_{1T} \\ a_{2T} \\ \vdots \\ a_{nT} \end{pmatrix} = \begin{pmatrix} a_{i_1}^T \\ a_{i_2}^T \\ \vdots \end{pmatrix}$$

2) P ist orthogonal

$$P^{-1} = P^T$$

$$I = PP^{-1} = PP^T$$

▷ Satz: Für jede Matrix $A \in \mathbb{R}^{n \times n}$ existieren Permutationsmatrizen und Faktoren L und R mit

$$PA = LR$$

MATLAB: $[L, R, P] = \text{lu}(A)$ liefert die Matrizen L, P, und R

Bemerkungen:

i) A singular, Rang $r < n$

Es gilt: $PA = LR$

mit $\begin{pmatrix} R \\ 0 \ 0 \end{pmatrix} \Bigg\}^r$

ii) Determinante

= Volumen des Parallelepipedes das von den Spaltenvektoren aufgespannt wird.

$$\det(A) = \det(PLR) = \underbrace{\det(P)}_{(-1)^v} \underbrace{\det(L)}_1 \det(R)$$

$$= (-1)^v \cdot \prod_{i=1}^n r_{ii} \quad ; \quad v \equiv \text{Anzahl der Vertausch. } P$$

iii) Vorteile der LR Zerlegung:

— Speicherung von L, R und P an Ort und Stelle in A
 zusammen mit Vektor der Länge n für P

— Zerlegungsaufwand:

$$\mu_{LR} = \frac{n^3}{3} + O(n^2)$$

Multiplik./Divis.

$$\mu_{\text{solve}} = n^2$$

$$Ax = b$$

$$x = A^{-1}b \quad \text{Aufwand ebenfalls } n^2$$

Mehrere rechte Seiten, mit derselben Matrix:

$$X = [x_1, x_2, \dots, x_L] ; B = [b_1, b_2, \dots, b_L]$$

(iv) Ist A dünnbesetzte Matrix, so ist A^{-1} vollbesetzt
 (sparse)

sparse: Anzahl von 0 verschiedenen Elementen
 ist $k \cdot n$

Bei dünnbesed. Matrizen sind i.a. die Faktoren

L, R auch dünnbesetzt.

(E-Technik:

f kompakte Träger, Fouriertransf. geht von $-\infty, \infty$)

2.2. Kondition, Norm

$$A\underline{x} = \underline{b} \quad ; \quad \underline{b}, \underline{x} \in \mathbb{R}^n, \quad A \in \mathbb{R}^n$$

A kann fast singulär sein.

Im Allgemeinen ist die Lösung \tilde{x} dann ungenau.

$$\text{Residuum: } \underline{r} := A\tilde{x} - \underline{b}$$

\tilde{x} kann ungenau sein auch wenn \underline{r} klein ist.

$$\text{Bsp: } n=2 : A = \begin{pmatrix} 137 & 100 \\ 100 & 73 \end{pmatrix}$$

$$A^{-1} = \begin{pmatrix} 73 & -100 \\ -100 & 137 \end{pmatrix}$$

$$\text{Mit: } \underline{b} = \begin{pmatrix} 4,30 \\ 3,10 \end{pmatrix} \Rightarrow \underline{x} = \begin{pmatrix} 3,9 \\ -5,3 \end{pmatrix}$$

Kleine Störung bei \underline{b} :

$$\tilde{\underline{b}} = \begin{pmatrix} 4,31 \\ 3,10 \end{pmatrix} = \underline{b} + \Delta, \quad \Delta = \begin{pmatrix} 0,01 \\ 0,00 \end{pmatrix} \sim 2/100$$

$$\text{Lösung für } \tilde{\underline{b}} : \tilde{x} = \begin{pmatrix} 4,63 \\ -6,20 \end{pmatrix} \sim \text{Änderung } 20\%$$

Jemand behauptet \tilde{x} löse $A\underline{x} = \underline{b}$.

$$\text{Kontrolle: Residuum } \underline{r} = A\tilde{x} - \underline{b} = \begin{pmatrix} 0,01 \\ 0 \end{pmatrix} \text{ scheint ok.}$$

Trotzdem ist \bar{x} um 20% falsch.

Zum Verständnis brauchen wir die Vektornorm, Matrixnorm und Kondition.

Vektornormen:

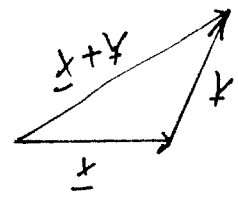
Sei $x \in \mathbb{R}^{n \times n}$

▷ Def: Die Zahl $\|x\| \in \mathbb{R}$ heißt eine Norm falls gilt:

- i) $\|x\| \geq 0 \quad \forall x \in \mathbb{R}^n, " = 0 " \iff x = \underline{0}$
- ii) $\alpha \in \mathbb{R} : \cancel{\| \alpha x \|}$

$\| \alpha x \| = |\alpha| \| x \| \quad \forall x \in \mathbb{R}^n$

iii) $\| x + y \| \leq \| x \| + \| y \|$



▷ Bem: $\| \cdot \|$ ist Abstandsbegriff zwischen Pkten x und y

Bsp: Höldernorm, p-Norm

$\underline{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$

▷ Def: $\| \underline{x} \|_p := \left(\sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}} ; p \geq 1$

3 Spezialfälle:

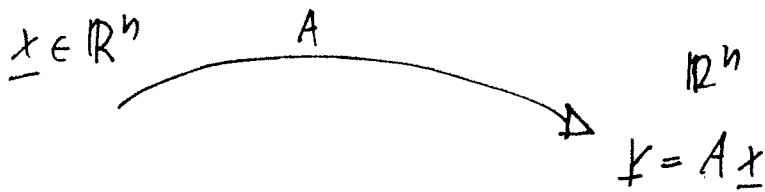
$p=2 : \| \underline{x} \|_2 = \sqrt{\left(\sum_{i=1}^n x_i^2 \right)}$ "Euklidische Norm"

$p=1 : \| \underline{x} \|_1 = \sum_{i=1}^n |x_i|$ "Summennorm"

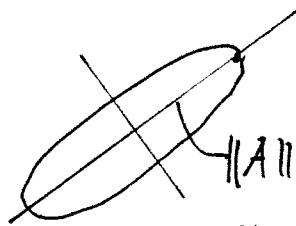
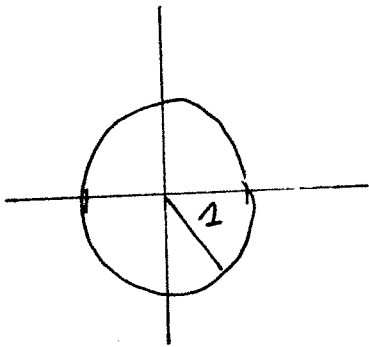
$p=\infty : \| \underline{x} \|_\infty = \max_{k=1, \dots, n} |x_k|$ "Maximumnorm"

Die einer Vektornorm zugehörige Matrixnorm

▷ Def: $\|A\| := \max_{\underline{x} \neq 0} \frac{\|A\underline{x}\|}{\|\underline{x}\|} = \max_{\|\underline{x}\|=1} \|A\underline{x}\|$



$p=2$



die grösste Streckung

1) Hauptsatz: $\|A\underline{x}\| \leq \|A\| \|\underline{x}\|$ (1)

▷ Satz: $A = (a_{ij})$

$$\|A\|_1 = \max_j \sum_{i=1}^n |a_{ij}|$$

maximale
Kolumnen
Spalten Betrags-Summe

$$\|A\|_\infty = \max_i \sum_{j=1}^n |a_{ij}|$$

maximale
Zeilen — " —

$p=2$

$$\|A\|_2 = |\lambda_{\max}|, \text{ falls } A = A^T$$

allgemein:

$$\|A\|_2 = \sigma_{\max} \text{ (grösster Singulärwert)}$$

Kondition einer Matrix

▷ Einfluss Datenfehler!

Sei $A \in \mathbb{R}^{n \times n}$, regulär, $b \neq 0$

LGS:

$$\left. \begin{array}{l} \text{ungestört: } \underline{Ax} = \underline{b} \quad (2) \\ \text{gestört: } A(\underline{x} + \Delta \underline{x}) = \underline{b} + \Delta \underline{b} \end{array} \right\} -$$

$$\Rightarrow A \Delta \underline{x} = \Delta \underline{b}$$

$$\Delta \underline{x} = A^{-1} \Delta \underline{b} \quad (3) \quad \text{absolute Fehler}$$

Wende (1) auf (2) u. (3) an:

$$\|\underline{b}\| \leq \|A\| \|\underline{x}\|$$

$$\|\Delta \underline{x}\| \leq \|A^{-1}\| \|\Delta \underline{b}\|$$

$$\|\underline{b}\| \cdot \|\Delta \underline{x}\| \leq \|A\| \|A^{-1}\| \|\underline{x}\| \cdot \|\Delta \underline{b}\|$$

$$\underbrace{\frac{\|\Delta \underline{x}\|}{\|\underline{x}\|}}_{\text{rel. Fehler}} \leq \|A\| \|A^{-1}\| \underbrace{\frac{\|\Delta \underline{b}\|}{\|\underline{b}\|}}_{\substack{\text{rel. Fehl.} \\ \text{rechte Seite}}}$$

▷ Def: $\text{cond}(A) = \|A\| \|A^{-1}\|$
Kondition von A bezüglich $\|\cdot\|$.

▷ Satz:

$$\frac{\|\Delta \underline{x}\|}{\|\underline{x}\|} \leq \underbrace{\text{cond}(A)}_{\substack{\text{maximale Verstärkungsfaktor} \\ \text{mit dem} \\ \text{der relative Fehler von } \underline{b} \\ \text{sich auf} \\ \text{den rel. Fehler der Lösung überträgt}}} \frac{\|\Delta \underline{b}\|}{\|\underline{b}\|} \quad (4)$$

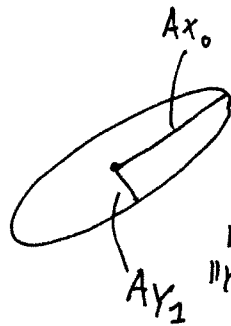
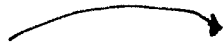
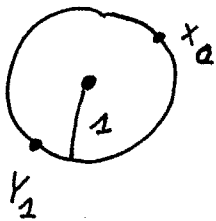
maximale Verstärkungsfaktor mit dem der relative Fehler von \underline{b} sich auf den rel. Fehler der Lösung überträgt

▷ Bem: i)

$\text{cond}(A)$ hängt von Norm ab

$$\text{cond}(A) = \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|y\|=1} \|Ay\|}$$

$$\underline{x} \mapsto A\underline{x}$$



$$\min_{\|y\|=1} \|Ay\| = \frac{1}{\|A^{-1}\|}$$

z. Bsp: $A = A^T$

$$\|\cdot\|_2 \Rightarrow \text{cond}_2(A) = \frac{\max |\lambda|}{\min |\lambda|}$$

$\lambda :=$ Eigenwerte

ii) MATLAB: $\text{cond}(A, p)$

Es wird die exakte $\text{cond}(A)$ berechnet.

Fehleranalyse

Rundungsfehler bei LR-Zerlegung (Gauss) wirken sich aus wie eine Störung der rechten Seite b

exakt: $A\underline{x} = \underline{b} \in \mathbb{R}^b$; $A \in \mathbb{R}^{n \times n}$

$\tilde{\underline{x}} = \underline{x} + \Delta\underline{x}$ numerisches Resultat

Residuum: $\underline{r} = A\tilde{\underline{x}} - \underline{b} \cong \Delta\underline{b}$

relatives Residuum: $\frac{\|\underline{r}\|}{\|\underline{b}\|} = \frac{\|\Delta\underline{b}\|}{\|\underline{b}\|}$

eps := Maschinengenauigkeit

eps := min(ε) mit 1 + ε > 1 (auf Rechner)

i. a. erreicht ein guter Algorithmus unabhängig von n und von cond(A)

(5) $\frac{\|\underline{r}\|}{\|\underline{b}\|} \leq \rho \cdot \text{eps}$; $\rho \approx 1$

(4), (5) \Rightarrow $\frac{\|\Delta\underline{x}\|}{\|\underline{x}\|} \leq \text{cond}(A) \cdot \rho \cdot \text{eps}$

▷ Bem:

iii) $\text{cond}(A) \geq 1$

$J = A^{-1}A$

$1 = \|J\| = \|A^{-1}\| \cdot \|A\| \leq \|A^{-1}\| \cdot \|A\|$

(iv) Fehler bei Multiplikation "exakte" Matrix mal

"gestörter" Vektor:

Geg: T reg. $T \in \mathbb{R}^{n \times n}$; $\underline{v} \in \mathbb{R}^n$

Ges: $\underline{w} = T\underline{v}$; $\underline{w} \in \mathbb{R}^n$

Ann: $\tilde{\underline{v}} = \underline{v} + \Delta\underline{v}$

Wie groß ist der Fehler von $\tilde{\underline{w}} = T\tilde{\underline{v}}$

Absolute Fehler: $\Delta\underline{w} = \tilde{\underline{w}} - \underline{w} = T\Delta\underline{v}$

$$\frac{\|\Delta\underline{w}\|}{\|\underline{w}\|} \leq \text{cond}(T) \frac{\|\Delta\underline{v}\|}{\|\underline{v}\|} \quad \text{"Vorkonditionierer" (Aufgabe)}$$

Wenn immer man mit Matrizen multipliziert, sollten Matrizen mit kleiner Konditionszahl gewählt werden.

2.3. QR-Zerlegung

Ausgleichsrechnung nach kleinsten Quadraten.

Vorbereitung: orthogonale Matrizen

$Q \in \mathbb{R}^{m \times m}$ heißt orthogonal $\Leftrightarrow Q^{-1} = Q^T$

Eigenschaften:

Normtreue in $\|\cdot\|_2$: $\|Q\underline{x}\|_2 = \|\underline{x}\|_2$; $\forall \underline{x}$

Bsp: Drehung, Spiegelungen

$\det(Q) = \pm 1$

- Euklidische Norm ist natürliche Norm

$\|Q\|_2 = 1, \text{cond}(Q) = 1$

Satz: QR-Zerlegung

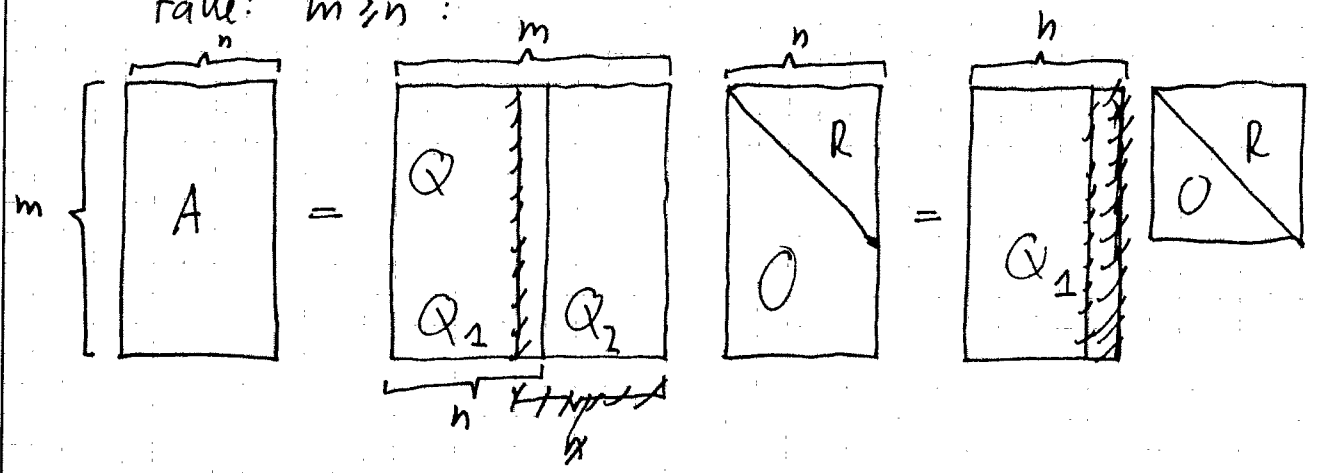
$A \in \mathbb{R}^{m \times n}$, beliebige Rechtecksmatrix

Dann existiert eine orthogonale Matrix

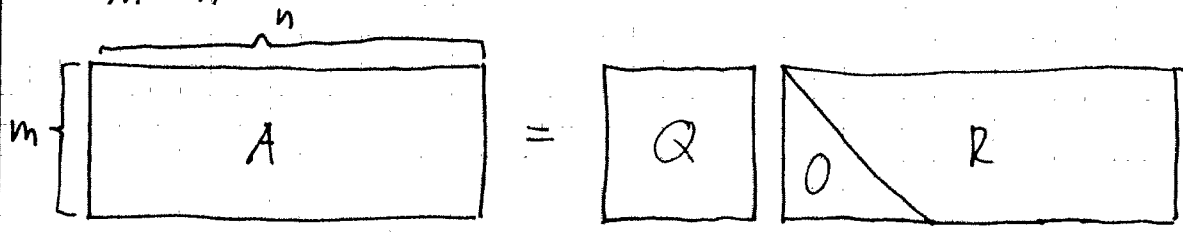
$Q \in \mathbb{R}^{m \times m}$ mit:

$A = QR$

Fälle: $m \geq n$:



$m < n$:



Beweis: Konstruktiv, Rechenverfahren

Idee: $R = Q^T A$

Wirkung der Givensrotation auf A

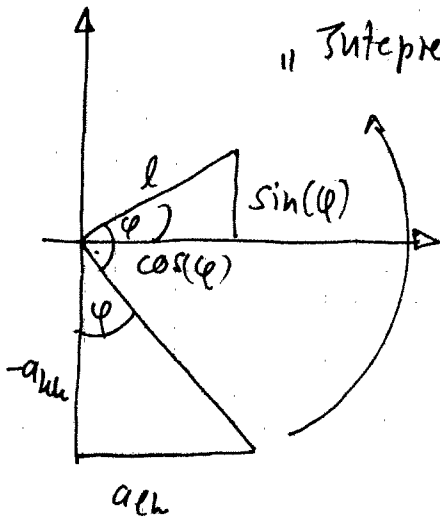
$$A_1 = Q_{ke}^T(\varphi) A \quad \left(\begin{array}{c} \underline{a}_1^T \\ \underline{a}_2^T \\ \vdots \\ \underline{a}_m^T \end{array} \right); A \in \mathbb{R}^{m \times n}$$

$$A_1 = \begin{pmatrix} \underline{a}_1^T & & & \\ \underline{a}_2^T & & & \\ \underline{a}_k^T \cos(\varphi) + \underline{a}_e^T \sin(\varphi) & & & \\ \underline{a}_{k+1}^T & & & \\ -\underline{a}_k^T \sin(\varphi) + \underline{a}_e^T \cos(\varphi) & & & \\ \vdots & & & \\ \underline{a}_m^T & & & \end{pmatrix} \begin{array}{l} \text{--- Zeile } k \\ \\ \text{--- Zeile } e \end{array}$$

$$(A_1)_{ek} = -a_{ek} \sin(\varphi) + a_{kk} \cos(\varphi) \stackrel{!}{=} 0$$

$$\varphi = \text{angle}(a_{kk} + i a_{ek}) \quad \text{"Matlab"}$$

"Interpretation"



QR:

Kolonnenweise alle Elemente $\neq 0$ des unteren Dreiecks von A "zu Null" machen.

$$R = Q_{m-1, n-1}^T \cdots Q_{23}^T Q_{1n}^T \cdots \\ \cdots Q_{13}^T Q_{12}^T(\varphi_{12}) A$$

▷ Anwendung: Least squares approximation

Bsp: Ansatz: linear

$$K(t) = \frac{a_0}{2} + \sum_{k=1}^N a_k \cos\left(\frac{2\pi}{p} kt\right) + b_k \sin\left(\frac{2\pi}{p} kt\right)$$

Finde $a_0, a_1, \dots, a_N, b_1, \dots, b_N$ so dass es passt

Typisches Problem:

Lineares Modell für eine Messreihe

Geg: m Zeitpunkte $t_k, k = 1, 2, \dots, m$ exakt

m Messwerte $f_k, k = 1, 2, \dots, m$ fehler behaftet

Betrachte "lineares Modell"

Ansatz:
$$\phi(t) := \sum_{l=1}^n \phi_l(t) c_l$$

* Basisfkt $l = 1, \dots, n$

das die Messreihe approximiert

$$\phi(t_k) \approx f_k$$

$K(t) = \phi(t)$, Bsp.

Basisfkt: $n = 2N + 1$

$$\phi_l = \begin{pmatrix} \dots \\ \dots \\ \dots \end{pmatrix}$$

Ges: $c = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}$, sodass "Abstand" klein

Lösung: Betrachte Residuenvektor

$$\underline{r} = \begin{pmatrix} r_1 \\ \vdots \\ r_m \end{pmatrix} \in \mathbb{R}^m; \quad r_k = \phi(t_k) - f_k$$

Modellmatrix:

$$A = \left(\phi_l(t_k) \right)_{\substack{l=1, \dots, m \\ k=1, \dots, n}}$$

$$\underline{r} = A\underline{c} - \underline{f} \quad ; \quad \underline{f} = \begin{pmatrix} f_1 \\ \vdots \\ f_m \end{pmatrix}$$

$$\underline{c} \in \mathbb{R}^n; \quad \underline{f} \in \mathbb{R}^m$$

Somit suche $\underline{c} \in \mathbb{R}^n$ mit $\|\underline{r}\| \stackrel{!}{=} \text{minimal}$ in geeigneter Norm. \otimes

17.03.2008

▷ Methode der kleinsten Quadrate heisst $\|\cdot\| := \|\cdot\|_2 =$
 $= \sum_{k=1}^m r_k^2$; Gauss., least squares

\otimes Lösen mit der QR-Faktorisierung

$$\begin{aligned} \|\underline{r}\|_2^2 &= \|A\underline{c} - \underline{f}\|_2^2 \\ &= \|QR\underline{c} - \underline{f}\|_2^2 \\ &= \underbrace{\|R\underline{c} - Q^T\underline{f}\|_2^2}_{\text{Residuum } \tilde{\underline{r}}} \stackrel{!}{=} \min \end{aligned}$$

$$\tilde{\underline{r}} := \begin{matrix} \begin{array}{|c|} \hline \begin{array}{c} \text{R} \\ \hline 0 \end{array} \\ \hline 0 \end{array} & \cdot & \begin{array}{|c|} \hline \underline{c} \\ \hline \end{array} & - & \begin{array}{|c|} \hline \underline{f} \\ \hline \end{array} & = & \begin{array}{|c|} \hline \tilde{\underline{r}}_{\text{oben}} \\ \hline \tilde{\underline{r}}_{\text{unten}} = (Q^T\underline{f})_{\text{unten}} \\ \hline \end{array} \end{matrix}$$

$Q^{-1} = Q^T$ $Q^T \underline{f}$

$$\text{Aus } \|\tilde{r}\|_2^2 = \underbrace{\|\tilde{r}_{\text{oben}}\|_2^2}_{=0} + \underbrace{\|\tilde{r}_{\text{unten}}\|_2^2}_{(Q^T f)_{\text{unten}} \text{ ist fest}} \stackrel{!}{=} \min$$

$$\tilde{r}_{\text{oben}} \stackrel{!}{=} 0$$

$$R_{\underline{c}} - (Q^T f)_{\text{oben}} = 0$$

$$R_{\underline{c}} = (Q^T f)_{\text{oben}} \quad \text{"Lösen durch Rückwärtseinsetzen"}$$

MATLAB:

$$A \in \mathbb{R}^{m \times n} ; \underline{f} \in \mathbb{R}^m$$

$$m \geq n$$

Befehl: $\underline{c} = A \backslash \underline{f}$ löst $\|A_{\underline{c}} - \underline{f}\|_2^2 \stackrel{!}{=} \min$ mit QR-Zerlegung.

2.4. Eigenwertprobleme

Rückblick:

$$A \in \mathbb{R}^{n \times n} ; \text{Finde } \underline{x} \neq 0 \text{ mit } A\underline{x} = \lambda \underline{x} ; \lambda \in \mathbb{C}$$

"Eigenwertproblem"

\underline{x} heißt Eigenvektor (EV)

λ -"- Eigenwert (EW)

▷ Charakteristisches Polynom: $(A - \lambda I) \underline{x} = 0 \quad (1)$

$$\Rightarrow \det(A - \lambda I) = 0$$

$P(\lambda) := \text{char. Polynom}$

$$P(\lambda) = \lambda^n - \lambda^{n-1} \text{tr}(A) + \dots +$$

$$= \sum_{k=1}^n a_{kk}$$

Allgemein $\lambda_k \in \mathbb{C}$.

Lin. Alg. \Rightarrow bestimme $p(\lambda)$ — MATLAB: $c = \text{poly}(A)$
Nullstellen bestimmen — $\lambda = \text{roots}(c)$

unbrauchbar !!

Grund: λ_k sind schlecht bestimmt durch Koeff. eines Polynoms.

Bsp: $A = \begin{pmatrix} 1 & \epsilon \\ \epsilon & 1 \end{pmatrix}$; $p(\lambda) = \lambda^2 - 2\lambda + 1 - \epsilon^2$

Sei MATLAB $\epsilon = 5 \cdot 10^{-9}$

$$\tilde{p}(\lambda) = \lambda^2 - 2\lambda + 1$$

↑
mitt Rechner

$$p(\lambda) = 0 \Rightarrow \lambda_{1/2} = 1 \pm \epsilon$$
 "Beide sind auf Rechner darstellbar"

$$\tilde{p}(\lambda) = 0 \Rightarrow \lambda_{1/2} = 1$$

"Nie Nullstellen von Polynomen ausrechnen, immer in Matrizen umwandeln und λ_k bestimmen."

▷ EW können zusammenfallen:

mehrfache EW : Anzahl des Auftretens nennt man die algebraische Vielfachheit.

▷ Diagonalform einer Matrix:

Lemma: Zu verschiedenen EW gehören linear unabhängige EV.

Ann: A hat n verschiedene EW.

$$\text{EW: } \lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n$$

$$\text{EV: } \underline{x}_1, \underline{x}_2, \underline{x}_3, \dots, \underline{x}_n \Rightarrow \text{linear unabhängig}$$

$$A \underline{x}_k = \lambda_k \underline{x}_k \quad ; \quad k = 1, \dots, n$$

$$A \underbrace{\begin{bmatrix} | & & & | \\ \vdots & & & \vdots \\ \underline{x}_1 & | & \underline{x}_2 & | \dots & | \underline{x}_n \\ \vdots & & & \vdots \\ | & & & | \end{bmatrix}}_T = \underbrace{\begin{bmatrix} | & & & | \\ \underline{x}_1 & | & \underline{x}_2 & | \dots & | \underline{x}_n \\ | & & & | \end{bmatrix}}_D \begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \dots & \\ 0 & & & \lambda_n \end{pmatrix}$$

$$AT = T D$$

Satz: $A \in \mathbb{R}^{n \times n}$; n linear unabhängige EV \underline{x}_k

$$AT = TD \Rightarrow A = TDT^{-1}$$

$$D = T^{-1}AT$$

} Ähnlichkeitstransformationen

MATLAB: $[T, D] = \text{eig}(A)$

▷ Eigenraum: (ER)

$\underline{x}_1, \underline{x}_2$ zwei Eigenvektoren zum selben EW λ :

$$\begin{array}{l|l} A\underline{x}_1 = \lambda \underline{x}_1 & \cdot \alpha_1 \\ A\underline{x}_2 = \lambda \underline{x}_2 & \cdot \alpha_2 \end{array}$$

$$A(\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2) = \lambda(\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2)$$

⇒ $\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2$ ist auch EV

∀ α_1, α_2

Fälle: $\underline{x}_1, \underline{x}_2$ linear unabhängig

$\{\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2\}$ ~~ist~~ bilden 2-dim. Eigenraum

$\underline{x}_1, \underline{x}_2$ linear abhängig

$$\underline{x}_2 = \gamma \underline{x}_1$$

$$\Rightarrow \{\alpha_1 \underline{x}_1 + \alpha_2 \underline{x}_2\} = \{c_1 \underline{x}_1\}$$

1.-dim. Raum

Zusammenfassung:

▷ EV sind höchstens auf freien Faktor bestimmt.
z.Bsp. $\| \cdot \| = 1$.

▷ alle EV zum Eigenwert λ bilden einen linearen Raum: ER zu λ

▷ Ein einfacher EW ($\text{alg.-V} = 1$) hat immer eine 1-Dim ER.

Def: Die dim. des ER zu einem mehrfachen EW heißt die geom. Vielf. des EW (= Anzahl d. linear unabh. EV)

Satz: geometrische Vielfachheit \leq algebraische Vielfachheit

Bsp: $A = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$; $\lambda_1, \lambda_2 = 1$ alg. $V = 2$
geom. $V = 2$

$A = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$; $\lambda_1 = \lambda_2 = 1$ alg. $V = 2$

es existiert nur 1 EV; $\underline{x}_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \rightarrow$ geom $V = 1$

Satz: Matrix ist diagonalisierbar, wenn für jeden EW gilt geom $V =$ algebr. V .

Folgerung: $A = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$ ist nicht diagonalisierbar.

MATLAB: $[T, D] = \text{eig}(A)$

Falls A nicht diagonalisierbar:

- ein k -facher EW λ wird in k einfache EW aufgespalten.

$$|\tilde{\lambda}_\ell - \lambda| \leq \text{const}(\text{eps})^{\frac{1}{k}}$$

\Rightarrow Damit wird A "diagonalisierbar"

Die zugehörigen EV; 3fach



▷ Symmetrische Matrizen $A = A^T$

Es gilt: Alle EW sind reell, $A = A^T \Rightarrow$ ist immer diagonalisierbar; T kann orthogonal gewählt werden.

$A \in \mathbb{C} ; A = (\bar{A})^T = A^H \leftarrow$ Hermitsch.

T unitär : $T^{-1} = T^H$

2.5. Numerik der Eigenwertprobleme

Aufgabe

Technik

a) Bestimmung einzelner EV und EW, vor allem des absolut grössten oder des absolut kleinsten

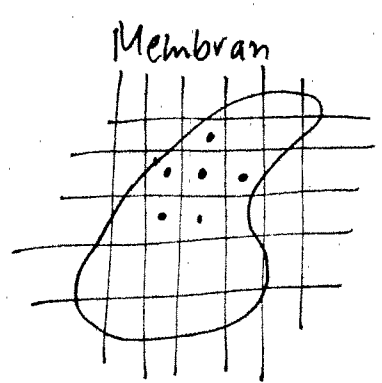
Vektoriteration mit Startvektor \underline{x}_0
 $\underline{x}_{k+1} = A \underline{x}_k$

b) Bestimmung aller EW und EV

Matrixzerlegung
 $T^{-1}AT = D$

a) Vektoriteration

Typisches Problem: Schwingung eines Kontinuums (Hochspannungsleitungen, Membran) in diskreter Approximation.



n Felder

$\Delta u = \lambda u$
 $u = 0$ auf Rand

Diskrete Approximation.

Sei $\underline{x} \in \mathbb{R}^n$, x gibt Auslenkung der Mitte des Feldes an

$$A \underline{x} = \lambda \underline{x}$$

$$A \in \mathbb{R}^{n \times n}$$

Spezielle Situation:

- A, x sind gross, d.h. n gross
- A ist dünnbesetzte Matrix
- $A = A^T$
- Nur wenige EW am unteren Spektrum gesucht.

▷ Algorithmus: Wähle (fast) beliebigen Startvektor \underline{x}_0
Iterationen $\underline{x}_{k+1} = A \underline{x}_k$; $k = 0, 1, 2, \dots$

Resultat: Folge $\{\underline{x}_k\}$ konvergiert nach EV zum grössten EW.

Kurze Wiederholung:

$$\text{EW-Probleme: } A \underline{x} = \lambda \underline{x}$$

Schlüsselbegriffe:

EW, EV, algebraische Vielfachheit, Diagonalform, ER, ...

Numerik der EW-Probleme

! Nicht das charakteristische Polynom bestimmen!

Grund: λ_k sind schlecht bestimmt durch die Koeffi. eines Polynoms. (siehe Bsp.)

Aufgabenstellung (mündlich wiederholen)bzw.
Fragestellunga) Vektoriteration

Bsp.: schwingende Membran (Zeichnung verteilen)

kontinuierlich:

diskret:

$$\begin{cases} \Delta u = \lambda u & \text{in } G \\ u = 0 & \text{auf } \partial G \end{cases}$$

$$\begin{cases} A \underline{x} = \lambda \underline{x} & \text{in } G \\ \underline{x} = 0 & \text{auf } \partial G \end{cases}$$

$\underline{x} \equiv$ Vektor der Auslenkungen der Gitterpunkte. wobei $\lambda = \omega^2$

↑
Eigenfrequenz

Algorithmus: (direkte Vektoriteration)

▷ Startin: $\underline{x}_{k+1} = A \underline{x}_k$; $k=0, 1, 2, \dots$

mit (fast) beliebigem Startvektor \underline{x}_0

▷ Normierung: $\underline{x}_{k+1} = \frac{A \underline{x}_k}{\|A \underline{x}_k\|}$

▷ Rayleigh-Quotient: $\lambda_{\max} = \lim_{k \rightarrow \infty} \frac{\underline{x}_k^T A \underline{x}_k}{\underline{x}_k^T \underline{x}_k}$ (Behav's Serie)



▷ Satz: Ann: λ_1 ein einfacher EW von $A \in \mathbb{R}^{n \times n}$ ($A = A^T$)
 und $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$; desweiteren
 stehe $\underline{x}_0 \in \mathbb{R}^n$ nicht senkrecht auf dem ER
 von λ_1 .

Folge: $\underline{y}_k := \frac{\underline{x}_k}{\|\underline{x}_k\|}$ mit $\underline{x}_{k+1} = A \underline{x}_k$ konvergiert
 gegen einen normierten EV von A zum EW λ_1 .

▷ Beweis:

$\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$ bilden eine Orthonormalbasis von
 Eigenvektoren von A mit $A \underline{v}_i = \lambda_i \underline{v}_i$.

Dann gilt $\underline{x}_0 = \sum_{i=1}^n \alpha_i \underline{v}_i$ mit $\alpha_1 = \langle \underline{x}_0, \underline{v}_1 \rangle \neq 0$

Folglich ist:

$$\underline{x}_k = A^k \underline{x}_0 = A^k \sum_{i=1}^n \alpha_i \underline{v}_i = \sum_{i=1}^n \alpha_i \lambda_i^k \underline{v}_i$$

$$= \alpha_1 \lambda_1^k \left[\underline{v}_1 + \sum_{i=2}^n \frac{\alpha_i}{\alpha_1} \left(\frac{\lambda_i}{\lambda_1} \right)^k \underline{v}_i \right]$$

Da $|\lambda_i| < |\lambda_1| \quad \forall i=2, \dots, n$,
 $=: \underline{z}_k$

gilt

$\lim_{k \rightarrow \infty} \underline{z}_k = \underline{v}_1$ und deshalb

$$\underline{y}_k = \frac{\underline{x}_k}{\|\underline{x}_k\|} = \frac{\underline{z}_k}{\|\underline{z}_k\|} \xrightarrow{k \rightarrow \infty} \underline{v}_1$$

□



Wir haben gezeigt die Folge $y_k = \frac{x_k}{\|x_k\|}$ konvergiert
in Richtung nach gegen EV zum grössten EW. (3)

▷ Bemerkungen:

- erhalten nur den EV zum grössten EW (N1)
- Berechnung ist "günstig"; A muss nicht gespeichert werden
- falls $\alpha_1 \cong 0$, konvergiert y_k dennoch gegen v_1
(Rundungsfehler)
- Konvergenzgeschwindigkeit hängt vom Quotienten $\left| \frac{\lambda_2}{\lambda_1} \right|$ ab
 \Rightarrow Langsame Konvergenz wenn $|\lambda_2|$ und $|\lambda_1|$ dicht zusammen liegen (N2)
- Bestimmung von λ_{\min} : Inverse Vektariteration

Allgemein:

Angenommen, wir hätten einen Schätzer $\mu \approx \lambda_i$ eines beliebigen
EW der Matrix A, so dass

$$|\lambda_i - \mu| < |\lambda_j - \mu| \quad \forall j \neq i.$$

Dann ist $(\lambda_i - \mu)^{-1}$ der betragsgrösste EW der
Matrix $(A - \mu I)^{-1}$.

Iterationsvorschrift

$$\Rightarrow (A - \mu I) x_{k+1} = x_k \quad \text{für } k = 0, 1, \dots$$

▷ Bei jedem Iterationsschritt wird das LGS gelöst, jedoch nur
für verschiedene rechte Seiten x_k .

$\Rightarrow (A - \mu I)$ muss nur einmal zerlegt werden!



Zusammenfassung:

vielleicht weglassen

Matrix	EW	EV	Grund
A	λ_i	\underline{v}_i	$A = TDT^{-1}$
$A^{-1} (\mu=0)$	λ_i^{-1}	\underline{v}_i	$A^{-1} = TD^{-1}T^{-1}$
$(A - \mu I)$	$(\lambda_i - \mu)$	\underline{v}_i	$(A - \mu I) = T(D - \mu I) T^{-1}$

Konvergenzfaktor: (Linear)

$\max_{j \neq i} \left| \frac{\lambda_i - \mu}{\lambda_j - \mu} \right| < 1$ für besonders gute Schätzer $\ll 1$

Simultane Vektoriteration (Unterraum-Iteration)

Ziel: p kleinste EW + dazugehörige EV bestimmen, alle paarweise orthonormal sein sollten.

Algorithmus:

$\triangleright p$ ^{orthogonale} Startvektoren: $X_0 = (x_0^{(1)}, \dots, x_0^{(p)}) \in \mathbb{R}^{n+p}$

\triangleright Iteration

Matlab: $[Q_k R_k^T] = qr(X_k); X_k = [x_k^{(1)} \dots]$

$X_{k+1} = AX_k$

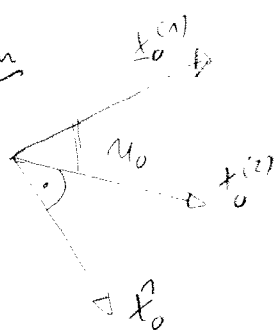
$\Rightarrow Q_k = [q_k^{(1)}, \dots, q_k^{(p)}]$
 $\langle q_k^{(i)}, q_k^{(j)} \rangle = \delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

- Vektoren $\underline{x}_k^{(i)}, i=1, \dots, p$ werden nicht mehr orthonormal sein

- Nach jedem Iterationsschritt eine orthog. Basis $\underline{x}_k^{(1)}, \underline{x}_k^{(2)}, \dots, \underline{x}_k^{(p)}$ im

Unterraum $U_k \in \mathbb{R}^b$ konstruieren mit Hilfe der QR-Zerlegung

Warum



\rightarrow strebe beide zu \underline{x}_1



b) QR-Algorithmus

Gebräuchlichste Methode zur Approximation aller EW und EV einer allgemeinen Matrix $A \in \mathbb{R}^{n \times n}$

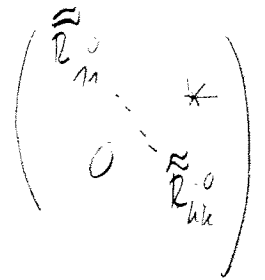
Sei $A_0 := A$

$$\begin{cases} \text{zerlege } A_0 = Q_0 R_0 = \boxed{Q_0} \begin{matrix} \triangleright \\ R_0 \end{matrix} ; Q_0^T Q_0 = I \\ \text{bilde } A_1 = R_0 Q_0 \end{cases}$$

Wegen $R_0 = Q_0^{-1} A_0 \Rightarrow A_1 = Q_0^{-1} A_0 Q_0 =$

$\Rightarrow A_1$ ist ^{orthogonal} "ähnlich" zu A_0

$\Rightarrow A_1$ hat dieselben EW wie A_0



Schursche Normalform (Satz von Schur erwähnen)

Algorithmus:

Gegeben $A_0 = A \in \mathbb{R}^{n \times n}$

$$A_k = Q_k R_k ; A_{k+1} = R_k Q_k , k = 0, 1, \dots$$

Satz: (wird nicht bewiesen)

Falls A_0 lauter EW von verschiedenem Betrag hat, konvergiert die Matrixfolge A_k gegen eine obere Dreiecksmatrix \tilde{R} .

Bemerkung:

- (i) EW von A erscheinen als Ein- und Zweierblöcke auf der Diagonalen von \tilde{R} .
- (ii) Bei betragsgleichen EW, entstehen "Kästchen" auf der Diagonalen, die nicht konvergieren.
- (iii) EV: $\underline{x}^{(i)} = Q \underline{y}^{(i)}$, $\underline{y}^{(i)}$ sind EV von A_k und $Q := Q_1 \cdot \dots \cdot Q_n$



function QR_ZerlegungFuerEW

mit Matlab zeige:

7

MATLAB
Beispiel

```

clc;
A = [0,0,4,-5,2; -2,1,8,-10,4;-2,-2,14,-15,6;...
     -2,-2,14,-17,8;-2,-2,15,-20,10];

```

```

eig(A)
n = 40
for i=1:n
    [Q,R] = qr(A);
    Aneu = R*Q;
    A = Aneu;
end;
A
return;

```

```

ans =

    2.0000
    1.0000
    3.0000
    1.0000 + 1.0000i
    1.0000 - 1.0000i

```

n =

10

A =

2.0010	5.9357	8.3191	40.1465	-11.7702
-0.0001	3.0005	1.5153	9.5280	-1.3649
-0.0000	0.0000	1.4992	-4.2026	1.0293
0.0000	-0.0003	0.2850	0.4908	0.1333
0.0002	-0.0012	-0.0485	-0.0496	1.0084

n =

40

A

2.0000	5.9386	28.0212	32.0702	-2.3863
-0.0000	3.0000	6.1357	7.5277	0.8402
0.0000	-0.0000	-0.5000	-3.7528	0.0000
0.0000	-0.0000	0.8660	2.5000	-0.0000
0.0000	-0.0000	0.0000	0.0000	1.0000

```
>> eig([-0.5 -3.7528; 0.8660, 2.500])
```

ans =

```

    1.0000 + 1.0000i
    1.0000 - 1.0000i

```

▷ EV $\hat{v}^{(3)}$ zum EW $\hat{\lambda}_3 \approx 3,0$ von A_{40} :

$$(A_{40} - \hat{\lambda}_3 I) \hat{v}^{(3)} = 0$$

$$\hat{x}^{(3)} = Q \hat{v}^{(3)} \quad \text{mit } Q = Q_1 \cdot Q_2 \cdot \dots \cdot Q_k$$



2.6. Singulärwertzerlegung (SVD)

07.04.2008

G. Golub, Stanford

$$A \in \mathbb{R}^{m \times n}; \text{Rang } p \leq l := \min(m, n)$$

$$(1) \underline{x} \in \mathbb{R}^n \mapsto \underline{y} = A \underline{x}$$

Neu Koordinaten: $\underline{\xi} \in \mathbb{R}^n$, $V \in \mathbb{R}^{n \times n}$ orthogonal $V^T V = I_{n \times n}$

$$\underline{x} = V \underline{\xi} \quad (2a)$$

Bildraum: $\underline{\eta} \in \mathbb{R}^m$: $U \in \mathbb{R}^{m \times m}$ ortho. $U^T U = I_{m \times m}$

$$(2b) \quad \underline{y} = U \underline{\eta}$$

$$U \underline{\eta} = \underline{y} \stackrel{(1)}{=} A \underline{x} \stackrel{(2a)}{=} A V \underline{\xi}$$

$$\underline{\eta} = \underbrace{U^T A V}_{=: S} \underline{\xi}$$

Ziel: Wähle U und V , sodass S "einfach" wird.

Satz: Zu jeder Matrix $A \in \mathbb{R}^{m \times n}$ gibt es orthogonale Matrizen U, V , sodass

$$A = U S V^T;$$

dabei ist S diagonal mit nicht negativen Diagonalelementen

$$s_1 \geq s_2 \geq s_3 \geq \dots \geq s_l \geq 0$$

s_1, s_2, \dots, s_l heißen Singulärwerte.

Bsp: $m \geq n$, $l = n$

$$\begin{array}{c} m \\ \left[\begin{array}{c} n \\ A \end{array} \right] = \begin{array}{c} m \\ \left[\begin{array}{c} m \\ U \end{array} \right] \end{array} \begin{array}{c} S \\ \left[\begin{array}{c} s_1 \quad 0 \\ s_2 \quad \cdot \\ \vdots \quad \cdot \\ s_n \end{array} \right] \end{array} \begin{array}{c} n \\ \left[\begin{array}{c} V^T \end{array} \right] \end{array} \\ \underbrace{\hspace{10em}}_{U^T U = I} \end{array}$$

Bem:

(i) Guter, stabiler Algorithmus zur Erzeugung von SVD bekannt. Aufwand ca. $10 m^3$ Operationen.

→ nicht anfällig für Rundungsfehler

(ii) MATLAB: ~~svd~~ $[U, S, V] = \text{svd}(A)$

(iii) Zusammenhang mit symmetrischen EWP

$$A^T A = \underbrace{V S^T S V^T}_{\text{diag}(s_1^2, s_2^2, \dots, s_n^2)}$$

$$\Rightarrow s_k^2 = \lambda_k \quad ; \quad \lambda_k \text{ EW von } A^T A$$

$$s_k = \sqrt{\lambda_k}$$

▷ Anwendungen:

(i) Rang einer Matrix $\text{rank}(A) = p < \min(m, n)$

$$\Rightarrow S = \text{diag}(s_1, s_2, \dots, s_p, 0, \dots, 0)$$

Praktisch werden die s_k mit $k > p$ zwar klein, aber nicht 0.

$$\frac{s_k}{s_1} \approx \text{eps}$$

Maschinengenauigkeit

"Numerischer Rang" hängt von der Rechengenauigkeit ab.

$$s_{\text{rech}} = \max \{k\}$$

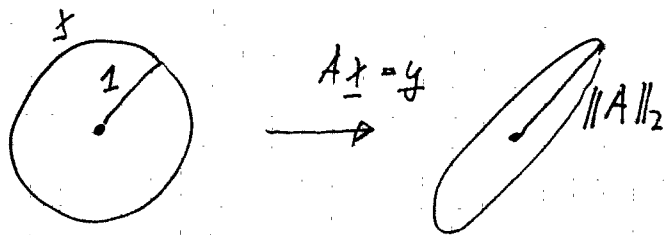
$$s_k > \text{eps} s_1$$

(ii) Euklidische Norm und Kondition

$$A \in \mathbb{R}^{n \times n}; A = USV^T$$

$$S = \text{diag}(s_1, \dots, s_n)$$

$\|A\|_2$ - maximale Länge des Bildes eines Einheitsvektors.



$$\Rightarrow s_1 = \|A\|_2$$

$$\|A^{-1}\|_2 = \frac{1}{s_n}$$

$$A^{-1} = VS^{-1}U^T$$



$$\text{diag}\left(\frac{1}{s_1}, \frac{1}{s_2}, \dots, \frac{1}{s_n}\right)$$

$$\text{cond}_2(A) = \frac{s_1}{s_n}$$

(iii) Lineare Abbildung:

Kern, Bild ansehen.

Sei $A \in \mathbb{R}^{m \times n}$, $m \geq n$

$\text{rank}(A) = s < n$

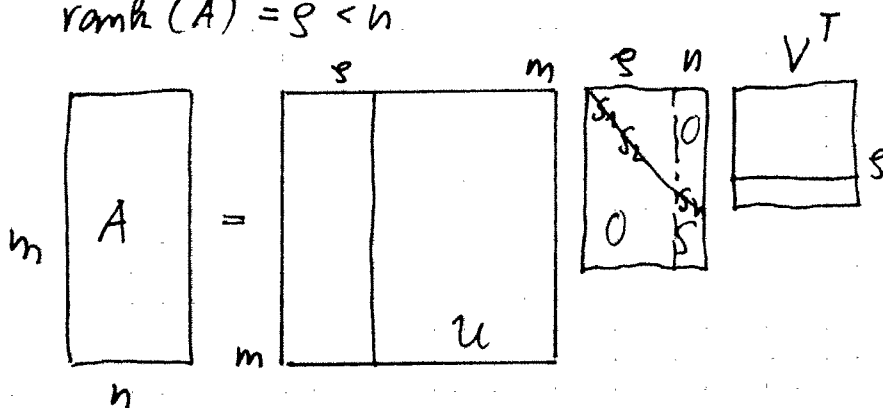
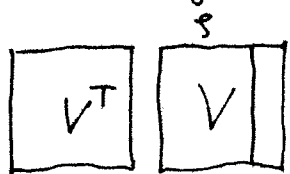


Abbildung: $\underline{x} \mapsto \underline{y} = A \underline{x}$



$\underbrace{\quad}_{n-s}$ orthogonale Basis von K_A

Kern: $K_A = \{ \underline{x} \in \mathbb{R}^n \mid A \underline{x} = 0 \}$

$\dim K_A = n - s$ orthogonale Basis von K_A
= Spaltenvektoren $s+1, s+2, \dots, n$
von V .

Bild: (Range)

$R_A := \{ \underline{y} \in \mathbb{R}^m \mid \exists \underline{x} \in \mathbb{R}^n \text{ mit } A \underline{x} = \underline{y} \}$

$\dim R_A = s$; ortho. Basis: Spalten $1, 2, \dots, s$ von U

(iv) Ausgleichsproblem mit nicht vollem Rang

Sei $A \in \mathbb{R}^{m \times n}$; $m \geq n$, $\text{rank}(A) = p < n$

Residuum: $\underline{r} = A\underline{c} - \underline{f}$

Gesucht \underline{c} , sodass $\|\underline{r}\|_2 \stackrel{!}{=} \min$

Sei SVD bekannt: $A = USV^T$

$$\begin{array}{c} \uparrow \\ = \text{diag}(s_1, \dots, s_p, 0, \dots, 0) \end{array}$$

Transformiertes Residuum:

$$\tilde{\underline{r}} = U^T \underline{r} = \underbrace{S(V^T \underline{c})}_{\underline{Y}} - \underbrace{U^T \underline{f}}_{\underline{g}} \quad \underline{Y} = V^T \underline{c}$$

Gesucht: $\underline{Y} = V^T \underline{c}$

Komponentenweise:

$$\tilde{r}_1 - s_1 y_1 - g_1 = 0 \quad y_1 = 0/s_1$$

$$\tilde{r}_2 = s_2 y_2 - g_2 = 0$$

$$\tilde{r}_p = s_p y_p - g_p = 0$$

$$\tilde{r}_{p+1} = 0 y_{p+1} - g_{p+1}$$

\vdots

$$\tilde{r}_n = 0 y_n - g_n$$

$$\tilde{r}_{n+1} = 0 y_{n+1} - g_{n+1}$$

Lösung:

$$y_k = \begin{cases} g_k/s_k & k=1, 2, \dots, p \\ \text{beliebig} & k=p+1, \dots, n \end{cases}$$

\Rightarrow es gibt ∞ -viele Lösungen.

$$\underline{c} = V\underline{y}$$

Auswahl ausgezeichnete Lösung:

Wähle y_{s+1}, \dots, y_n

sodass $\| \underline{s} \|_2 \stackrel{!}{=} \min$

Orthogonal.

$$\Rightarrow \| \underline{s} \|_2 = \| \underline{y} \|_2$$

\Rightarrow Minimum-Norm-Lösung

$$0 = y_{s+1} = y_{s+2} = \dots = y_n$$

3. Iterative Methoden für LGS

(1) Problem: $A\underline{x} = \underline{b}$

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

Speziell : n sehr gross $\cdot 10^6$

A dünnbesetzt (sparse)

Gesucht: \underline{x}

A sparse $\Rightarrow A^{-1}$ vollbesetzt

Nie A^{-1} berechnen !!

Benutze nur die Abbildung $\underline{x} \mapsto A\underline{x}$

Ziel: Verfahren die nur die Abbildung benutzen
(Operatorprinzip)

3.1. Jacobi - Iteration

Sei $A = C - B$

Aus (1) folgt

$$C\underline{x} - B\underline{x} = \underline{b}$$

$$C\underline{x} = B\underline{x} + \underline{b}$$

$$\underline{x} = C^{-1}(B\underline{x} + \underline{b})$$

Fixpunktform:

Wähle Startvektor x_0

$$\underline{x}_{k+1} = C^{-1} (B \underline{x}_k + \underline{b})$$

▷ Voraussetzungen:

C^{-1} regulär

C sodass, $C \underline{x}_{k+2} = B \underline{x}_k + \underline{b}$ (3) leicht lösbar.

▷ Fehlertheorie:

(2) - (3)

$$C(\underline{x}_{k+2} - \underline{x}) = B(\underline{x}_k - \underline{x})$$

▷ Def: Fehler: $\underline{e}_k = \underline{x}_k - \underline{x}$

$$\Rightarrow C \underline{e}_k = B \underline{e}_k$$

$$\underline{e}_{k+1} = C^{-1} B \underline{e}_k$$

Normen: $\|\underline{e}_{k+1}\| \leq \|C^{-1} B\| \|\underline{e}_k\|$

Satz: Hinreichend für die Konvergenz ist, dass

$$\|C^{-1} B\| < 1$$

für mindestens eine Norm

(i) Jacobi-Iteration

$$A = (a_{ij})$$

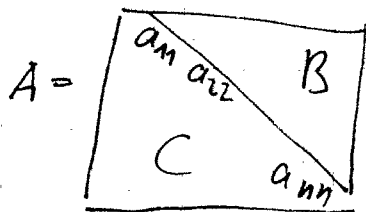
Sei A "diagonal dominant", d.h.

$$|a_{ii}| > \sum_{k \neq j} |a_{jk}| \quad (4)$$

Wähle $C = (C_{ij}) = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$

(4) $\Rightarrow \|C^{-1}B\|_{\infty} < 1$

(ii) Gauss-Seidel



konvergiert leicht schneller als Jacobi.

3.2. Konjugierte Gradienten (cg)

M. Hestenes, E. Stiefel ~ 51, ~ 52

Sei $A \in \mathbb{R}^{n \times n}$, symmetrisch $A^T = A$

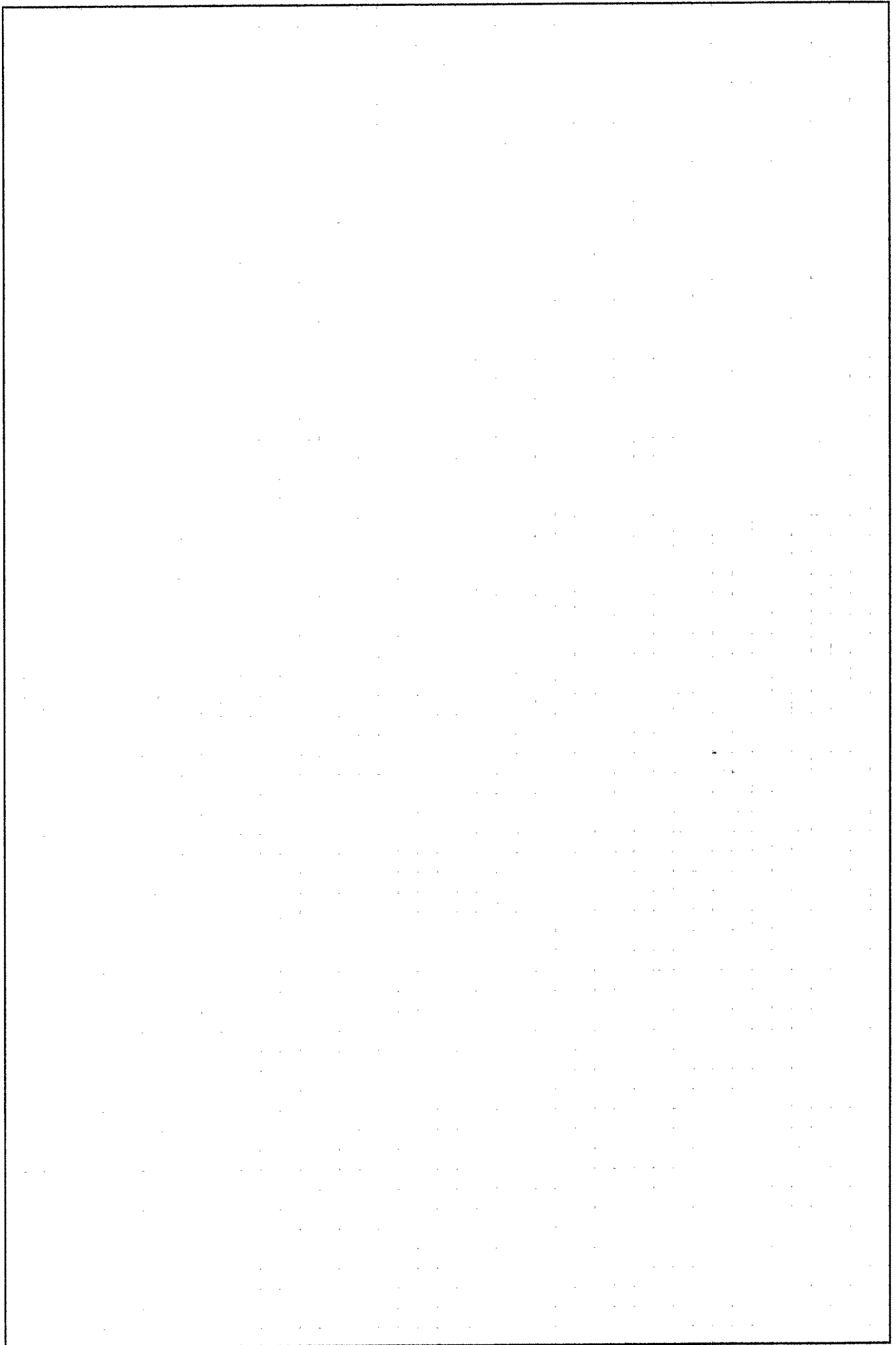
und A positivdefinit. d.h. $Q(x) := x^T A x > 0 \quad \forall x \neq 0$

Bem: (1) A hat lauter, positive EW

Beweis: λ EW, $v \in EV \Rightarrow Av = \lambda v$

$$0 < Q(v) = v^T A v = v^T v \lambda = \underbrace{v^T v}_{\|v\|_2^2} \lambda$$

$$\lambda > 0$$



Wiederholung:

Iterative Methoden für LGS

additive Zerlegung: $A = C - B$

Algorithmus: Wähle Startvektor \underline{x}_0

$$C \underline{x}_{k+1} = B \underline{x}_k + \underline{b} \quad ; \quad k = 0, 1, 2, \dots$$

$$\left(\underline{x}_{k+1} = C^{-1} B \underline{x}_k + C^{-1} \underline{b} \right)$$

↓
- wie Inverse berechnen

Bsp: für additive Zerlegung: $A = J - G$

$$\Rightarrow \underline{x}_{k+1} = G \underline{x}_k + \underline{b} \quad (\text{Serie})$$

Fehlertheorie: $\| \underline{e}_{k+1} \| \leq \| \underbrace{G}_{< 1} \| \cdot \| \underline{e}_k \|$

Bsp:

$$\begin{pmatrix} 4 & 1 & 0 \\ 1 & 4 & 1 \\ 0 & 1 & 4 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 5 \\ 6 \\ 5 \end{pmatrix}$$

$$4x_1 + x_2 = 5$$

$$\Rightarrow x_1^{(k+1)} = (-x_2^{(k)} + 5) \cdot \frac{1}{4}$$

$$x_2^{(k+1)} = (-x_1^{(k)} - x_3^{(k)} + 6) \cdot \frac{1}{4}$$

$$x_3^{(k+1)} = (-x_2^{(k)} + 5) \cdot \frac{1}{4}$$

Jacobi-Verfahren

besser

Gesamtschrittverfahren



▷ Matrix Darstellung

$$\underline{x}^{(k+1)} = \underbrace{\frac{1}{4} \begin{pmatrix} 0 & -1 & 0 \\ -1 & 0 & -1 \\ 0 & -1 & 0 \end{pmatrix}}_{\|G\| < 1} \underline{x}^{(k)} + \frac{1}{4} \begin{pmatrix} 5 \\ 6 \\ 5 \end{pmatrix}$$

▷ Gauss-Seidel-Verfahren (EinzelSchrittverfahren)

$$x_1^{(k+1)} = (-x_2^{(k)} + 5) \frac{1}{4}$$

$$x_2^{(k+1)} = (-x_1^{(k+1)} - x_3^{(k)} + 6) \frac{1}{4}$$

$$x_3^{(k+1)} = (-x_2^{(k+1)} + 5) \frac{1}{4}$$

$$x_3^{(k+1)} = \left[(+x_1^{(k+1)} + x_3^{(k)} + 6) \frac{1}{4} + 5 \right] \frac{1}{4}$$



3.2. Das Verfahren der konjugierten Gradienten

(3)

▷ Motivation:

Unter Berücksichtigung der Spezialstruktur der Matrix (Besetzungsstruktur: sparsity pattern) das LGS iterativ und möglichst "direkt" zu einer bestimmten Genauigkeit zu lösen.

(Bsp. Serie 5: Schwingende Membran: Laplace Gld. führte auf eine Block-Tridiagonalmatrix)

▷ Ziel: Löse LGS $A\underline{x} + \underline{b} = 0$ mit $A \in \mathbb{R}^{n \times n}$
 $A^T = A$ und positiv definit: $\underline{x}^T A \underline{x} > 0$

▷ Bem:

- (i) Alle EW von A sind positiv
- (ii) Wir suchen ein Minimum der Energie (Energienorm) (Bild)

▷ Satz: Die Lösung \underline{x} von $A\underline{x} + \underline{b} = 0$ ($A^T = A, \underline{x}^T A \underline{x} > 0$, $A \in \mathbb{R}^{n \times n}$) ist das Minimum der quadratischen Funktion

$$F(\underline{u}) := \frac{1}{2} \underline{u}^T A \underline{u} + \underline{u}^T \underline{b} \quad (*)$$

▷ Beweis:

$$\text{grad}_{\underline{u}} F(\underline{u}) = \left(\frac{\partial F}{\partial u_1}, \frac{\partial F}{\partial u_2}, \dots, \frac{\partial F}{\partial u_n} \right)^T = A \underline{u} + \underline{b} := \underline{v}$$

\underline{v} = Residuumvektor zum Vektor \underline{u} .

$\text{grad}_{\underline{u}} F(\underline{u}) = 0$ ist notwendige Bedingung für ein Extremum



Minimaleigenschaft:

$$\underline{u}, \underline{x} \in \mathbb{R}^n, \underline{x} \neq 0$$

$$\begin{aligned} F(\underline{u} + \underline{x}) &\stackrel{(*)}{=} \frac{1}{2} (\underline{u} + \underline{x})^T A (\underline{u} + \underline{x}) + (\underline{u} + \underline{x})^T \underline{b} \\ &= \frac{1}{2} \left[\underline{u}^T A \underline{u} + \underbrace{\underline{u}^T A \underline{x} + \underline{x}^T A \underline{u}}_{\underline{u}^T A^T \underline{x} = (A^T \underline{x})^T \underline{u} = \underline{x}^T A \underline{u}} + \underline{x}^T A \underline{x} \right] + \underline{u}^T \underline{b} + \underline{x}^T \underline{b} \\ &= F(\underline{u}) + \underbrace{\underline{x}^T (A \underline{u} + \underline{b})}_{=0 \text{ Extremum}} + \underbrace{\frac{1}{2} \underline{x}^T A \underline{x}}_{> 0 \text{ positiv def.}} \\ &> F(\underline{u}) \quad \square \end{aligned}$$

▷ Prinzip um Minimum zu finden

Zu einem gegebenen Näherungsvektor (Startvektor) und einem gegebenen, "geeignet festzulegenden Richtungsvektor" $\underline{p} \neq 0$ die Funktion $F(\underline{u})$ in dieser Richtung zu minimieren.

▷ 2D - Minimierung (vergleiche Bilder)

Geg.: Startpunkt \underline{x}_0
Suchrichtung \underline{p}_0

$$F(\underline{x}_0 + s \underline{p}_0) \stackrel{!}{=} \min$$

Ges.: Parameter $s = s_0 \in \mathbb{R}$, so dass

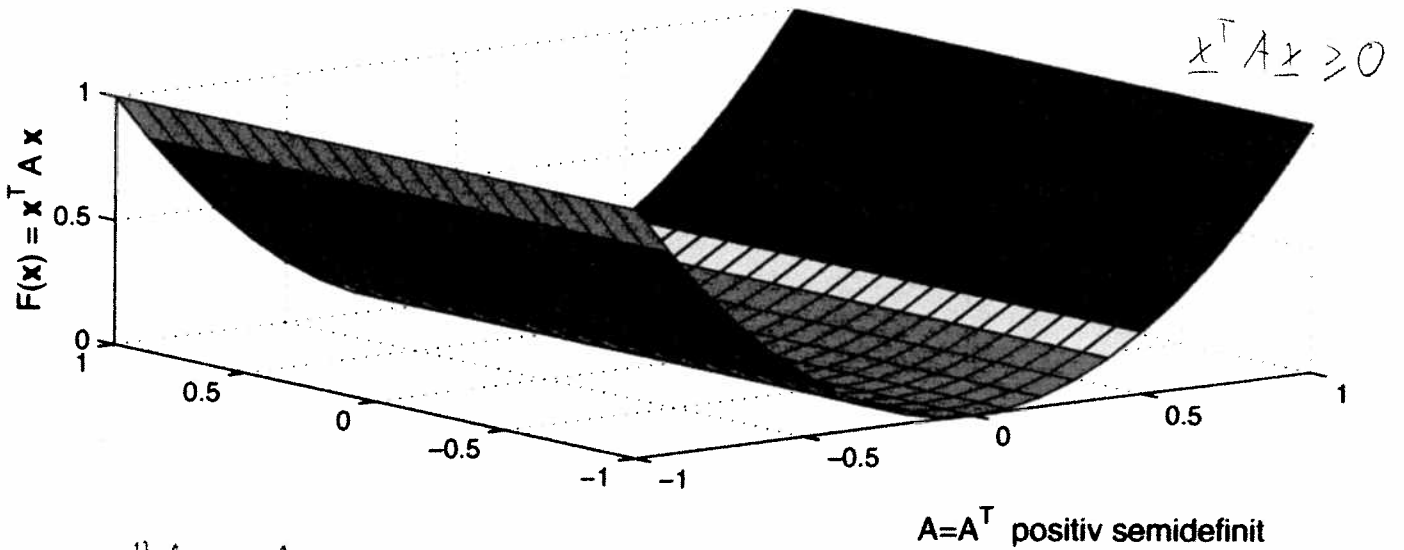
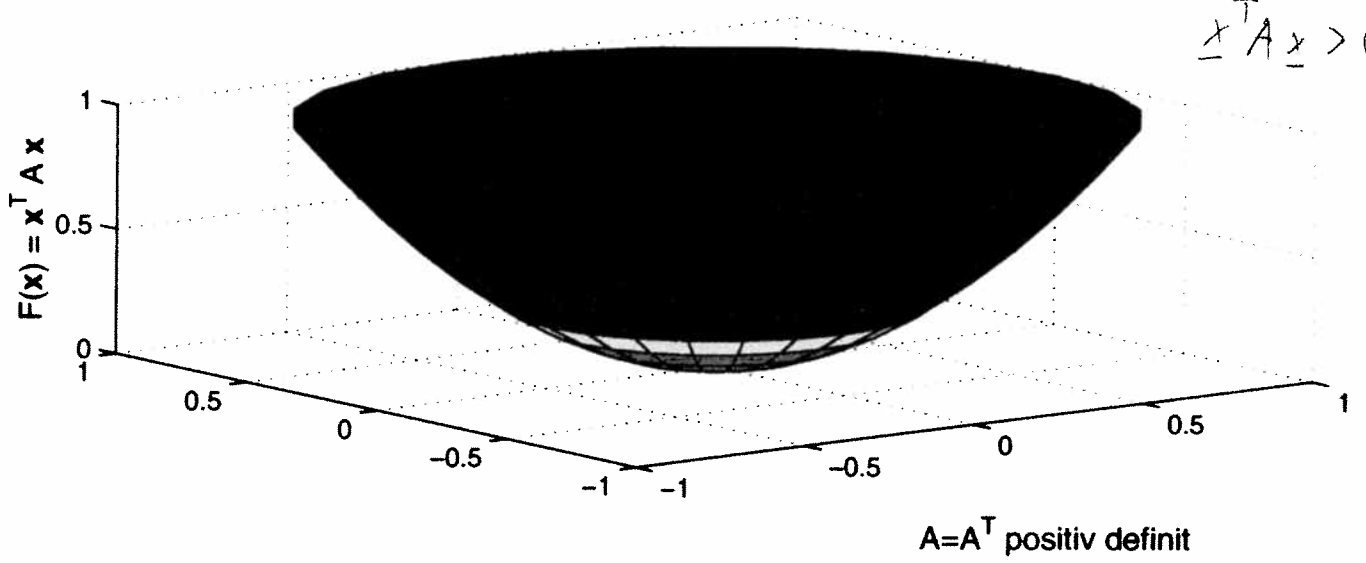
$$\begin{aligned} F(\underline{x}_0 + s \underline{p}_0) &= \frac{1}{2} (\underline{x}_0 + s \underline{p}_0)^T A (\underline{x}_0 + s \underline{p}_0) + (\underline{x}_0 + s \underline{p}_0)^T \underline{b} \\ &= F(\underline{x}_0) + \frac{1}{2} s^2 \underline{p}_0^T A \underline{p}_0 + s \underbrace{\underline{p}_0^T (A \underline{x}_0 + \underline{b})}_{= r_0} = F^*(s) \end{aligned}$$



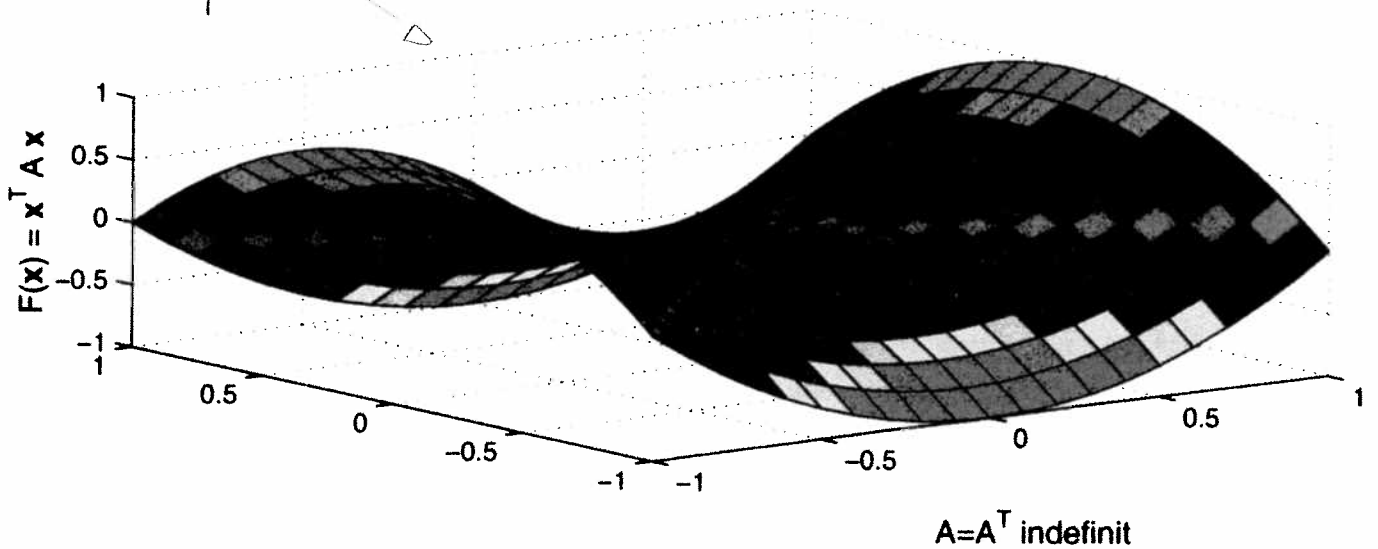
nicht aus ^{wie}VP Paraboloid sollte aber Ellipsoid sein

Bilder ①

Quadratische Formen in $x = (x_1, x_2)^T$



Sattelpunkt





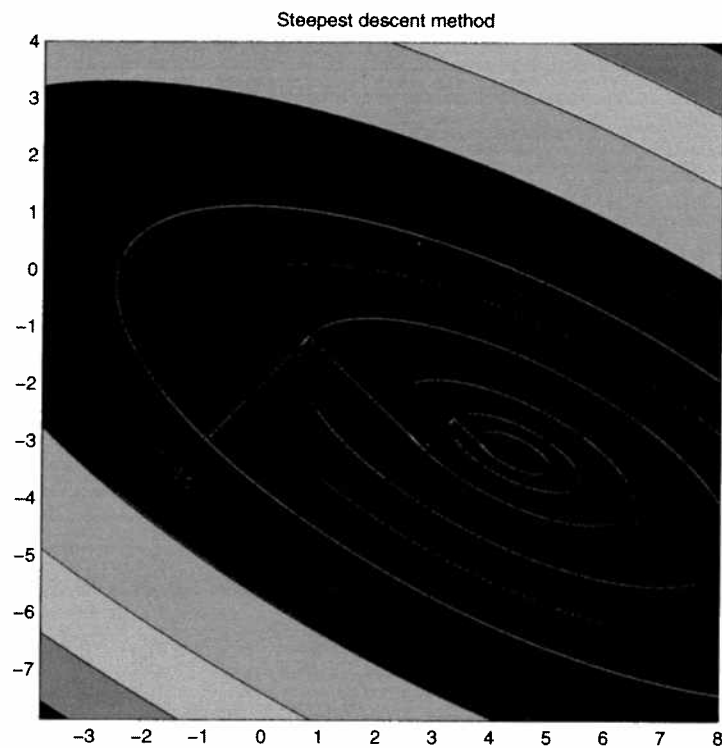


Figure 1: The steepest descent method for $N = 2$.

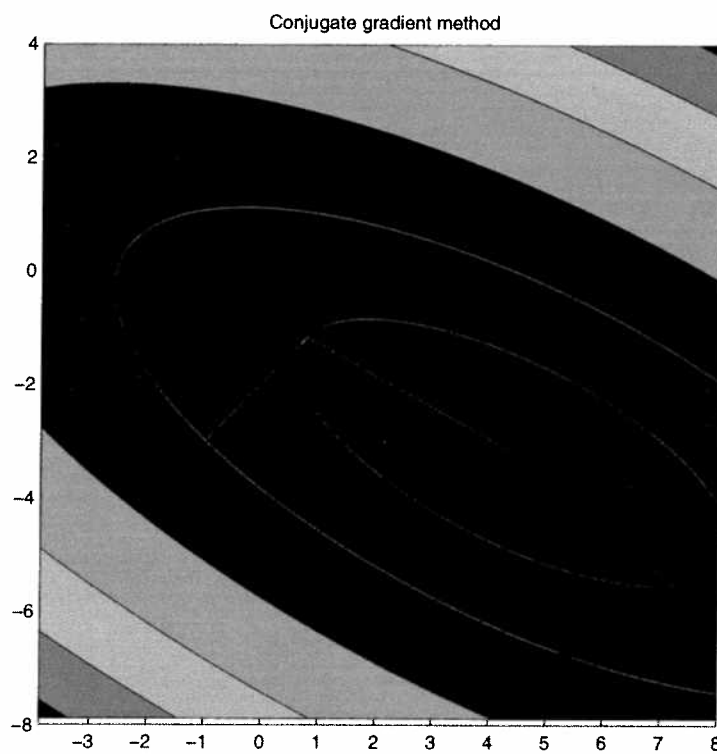


Figure 2: Conjugate directions — the CG method for $N = 2$.



$$\frac{d}{ds} F^*(s) = s p_0^T A p_0 + p_0^T r_0 \stackrel{!}{=} 0$$

(5)

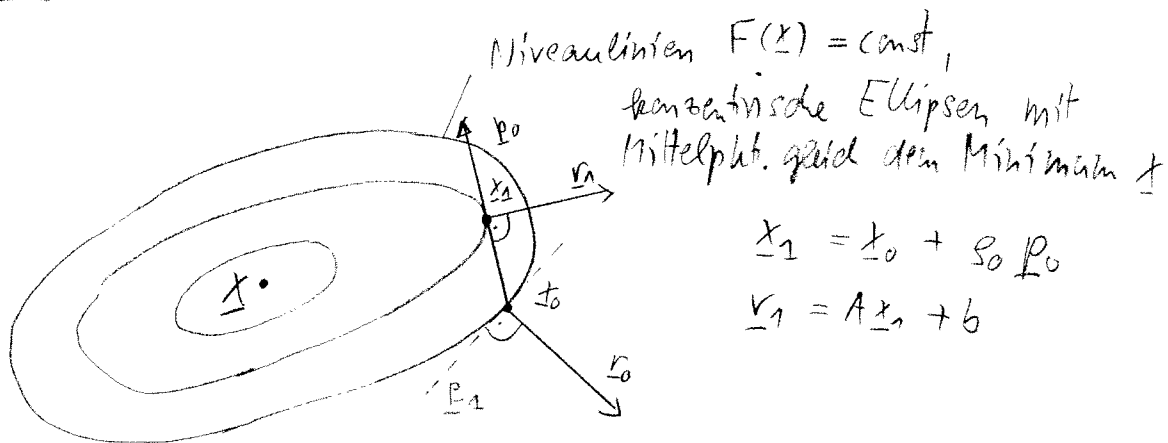
\Rightarrow

$$s_0 = - \frac{p_0^T r_0}{p_0^T A p_0}$$

 \leftarrow Somit haben wir die Länge s festgelegt.
 $\rightsquigarrow x_1 = x_0 + s_0 p_0$

$$\frac{d^2}{ds^2} F^*(s) = p_0^T A p_0 > 0 \quad \text{"Minimum"}$$

▷ Geometrische Interpretation eines Iterationsschrittes



$$x_{k+1} = x_k + s_0 p_k$$

$$r_k = A x_k + b$$

p ist Tangente an die Niveaulinie durch x_k

▷ Wie wählt man den Richtungsvektor p_k

a) $p_k := -r_k = -\text{grad } F(x_k)$

- naheliegend, denn steilster Abstieg, jedoch nicht immer optimal. (siehe Bild) \rightarrow Langsame Konvergenz
- langgestreckte Ellipsen (entsprechend $\text{cond}(A)$ groß) viele Schritte, obwohl in jedem Iterationsschritt die Richtung gewählt wird, welche die stärkste Abnahme von $F(x)$ garantiert.

\Rightarrow Mehr Information benötigt:

zwei (lin. unabh.) Suchrichtungen



b) Konjugierte Gradienten

6

▷ Geometrisch:

Die Richtung f , welche vom Punkt x_k den Mittelpunkt \underline{x} der Ellipse trifft, ^{ist} mit der Tangentialrichtung im Punkt x_k im Sinn der Kegelschnittgleichung konjugiert.

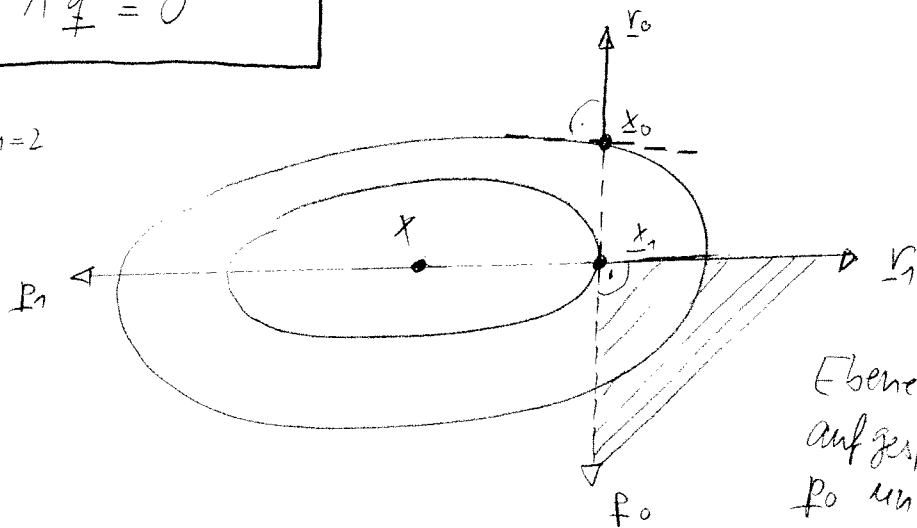
▷ Def:

Zwei Vektoren (Richtungen) $f, g \in \mathbb{R}^n$ heißen konjugiert (A -orthogonal), falls für A (mit $x^T A x > 0$ und $A^T = A$) gilt

$$f^T A g = 0 \quad (*)$$

Zeichnung

(2D) $n=2$



Ebene $E \in \mathbb{R}^n$
aufgespannt von
 f_0 und v_1

$$p_1 \in \text{span}(f_0, v_1)$$

▷ Ziel:

Im k -ten Iterationsschritt das Minimum von $F(x_k)$ bezüglich der Ebene E zu ermitteln.

\Rightarrow Richtungsvektor f_{k+1} muss konjugiert sein zu f_k bezüglich der Schnitteellipse und somit auch bezüglich des Ellipsoids $F(x_k)$

Aufgabe

bei gegebenem f_k finde eine Tangentialrichtung f_{k+1} an Ebene E (Niveaufläche $f_{k+1}^T A f_k = \text{const}$)



▷ Frage: Wie sieht die $(k+1)$ -te Suchrichtung aus?

▷ Ansatz:

$$p_1 = -r_1 + \alpha_1 p_0; \alpha_1 = ?$$

$$A = A^T$$

$$p_0, p_1 \text{ konjugiert} \Rightarrow 0 = p_1^T A p_0 \stackrel{\downarrow}{=} p_0^T A p_1$$

$$\stackrel{\downarrow}{=} p_0^T A (-r_1 + \alpha_1 p_0) = \alpha_1 p_0^T A p_0 - p_0^T A r_1 = 0$$

$$\Rightarrow \alpha_1 = \frac{p_0^T A r_1}{p_0^T A p_0} \quad (\text{damit haben wir auch } p_1 \text{ bestimmt})$$

Gesucht ist aber $x \approx x_2 = \underbrace{x_0 + p_0 p_0}_{x_1} + p_1 p_1$

aus $\frac{d}{ds} F^*(s) = 0$ und p_0, p_1 konjugiert

$$p_1 = -\frac{p_1^T r_1}{p_1^T A p_1}$$

falls $x_0 = 0$

▷ cg-Algorithmus

• Start: Wahl von $x_0 (=0)$; $r_0 = A x_0 + b$; $p_{-1} \stackrel{\downarrow}{=} 0$

• Iteration: $k = 0, 1, \dots$

Sonst definiere wir ihn
uns als $p_{-1} = 0$

$$(1) \quad \alpha_k = \begin{cases} \frac{p_{k-1}^T A r_k}{p_{k-1}^T A p_{k-1}} & \text{falls } k > 0 \\ 0 & \text{falls } k = 0 \end{cases}$$

falls $k > 0$

falls $k = 0$

dann ist gerade $\alpha_0 = 0$

$$(2) \quad p_k = \alpha_k p_{k-1} - r_k$$

$$(3) \quad p_k = -\frac{p_k^T r_k}{p_k^T A p_k}$$



$$(4) \quad \underline{x}_{k+1} = \underline{x}_k + \mathcal{S}_k \underline{p}_k$$

$$(5) \quad \underline{r}_{k+1} = A \underline{x}_{k+1} + \underline{b} = \underline{r}_k + \mathcal{S}_k A \underline{p}_k$$

▷ Satz: Orthogonalitätsrelationen

$$\bullet \quad \begin{array}{l} \underline{r}_{k+1} \perp \underline{p}_k \\ \underline{r}_{k+1} \perp \underline{p}_{k-1} \\ \underline{r}_{k+1} \perp \underline{r}_k \end{array} \quad \begin{array}{l} (\text{kann man mittels vollständige}) \\ (\text{vielleicht in Serie}) \end{array}$$

⇒ Suchrichtungen $\underline{p}_0, \underline{p}_1$ paarweise konjugiert

⇒ Residuen $\underline{r}_0, \underline{r}_1$ ($\nabla F(\underline{x}_0), \nabla F(\underline{x}_1)$)
paarweise orthogonal.

▷ Satz: Die Residuenvektoren \underline{r}_k bilden ein Orthogonalsystem, und die Richtungsvektoren \underline{p}_k sind paarweise konjugiert. Für $l \geq 2$ gelten

$$\underline{r}_{k+1}^T \underline{r}_k = 0 \quad k = 0, 1, \dots, l-2$$

$$\underline{r}_{k+1}^T \underline{p}_k = 0 \quad k = 1, 2, \dots, l-1$$

$$\underline{p}_k^T A \underline{p}_k = 0 \quad k = 1, 2, \dots, l-1$$

⇒ $\underline{r}_0, \underline{r}_1, \dots, \underline{r}_{n-1} \in \mathbb{R}^n$ bilden ein Orthogonalsystem in \mathbb{R}^n
d.h. es kann höchstens n von Null verschiedene
Vektoren erhalten

⇒ $\underline{r}_n = 0 \quad \Rightarrow \quad \boxed{x_k = x \quad \text{exakte Lösung}}$



c) Verbanditionierung

(9)

▷ Def: $\underline{x} \in \mathbb{R}^n$, $A \in \mathbb{R}^{n \times n}$, $\underline{x}^T A \underline{x} > 0$

A-Norm von \underline{x} : $\|\underline{x}\|_A := \sqrt{\underline{x}^T A \underline{x}}$ "Energienorm"
(siehe S.10)

▷ Satz: Konvergenzabschätzung des cg-Verfahrens (kein Beweis)

$$\underline{e}_k = \underline{x}_k - \underline{x}$$

$$\frac{\|\underline{e}_k\|_A}{\|\underline{e}_0\|_A} \leq 2 \left(\frac{\sqrt{\text{cond}_2(A)} - 1}{\sqrt{\text{cond}_2(A)} + 1} \right)^k \quad (*)$$

Annahme: $\|\underline{e}_k\|_A / \|\underline{e}_0\|_A \leq \varepsilon$

$$\stackrel{(*)}{\Rightarrow} \boxed{k \leq \frac{1}{2} \sqrt{\text{cond}_2(A)} \ln\left(\frac{2}{\varepsilon}\right) + 1}$$

↳ $\text{cond}(A)$ entscheidend für Konvergenzgröße! ▽

▷ Skalierung (Verbanditionierung)

neues LGS:

$$\tilde{A} \tilde{x} + \tilde{b} = 0 \quad (1)$$

$$\tilde{A} = C^{-T} A C^{-1}$$

$$\tilde{x} = C x$$

$$\tilde{b} = C^{-T} b$$

C regulär und "günstig"

▷ Ziel: $\text{cond}_2(\tilde{A}) \approx 1$

$$\text{cond} \equiv \kappa$$

Damit $\kappa_2(\tilde{A}) = \kappa_2(C^{-T} A C^{-1}) \ll \kappa_2(A)$ ist,

erhalten wir aus der Ähnlichkeitstransformation von \tilde{A} :



$$C^{-1} A C = C^{-1} C^T A \underbrace{C^{-1} C}_I = \underbrace{(C^T C)^{-1}}_{:=M} A$$

$M = C^T C$ ist symmetrisch und positiv definit
 ↑
 Vorconditionierungsmatrix

Wegen Ähnlichkeit:

$$k_2(\tilde{A}) = \frac{\lambda_{\max}(M^{-1}A)}{\lambda_{\min}(M^{-1}A)}$$

▷ cg-Algorithmus mit Vorconditionierung
 (siehe Skript)

Einführung einer zusätzlichen Vektorsequenz $\underline{z}_k := \underline{r}_k$
 $\underline{z}_k = M^{-1} \underline{r}_k$

▷ Nachtrag zur Energienorm:

$$\begin{aligned} \underbrace{\| \underline{x}_k - \underline{x} \|_A^2}_{\underline{e}_k} &= (\underline{x}_k - \underline{x})^T A (\underline{x}_k - \underline{x}) \\ &= \underline{x}_k^T A \underline{x}_k - 2 \underline{x}_k^T A \underline{x} + \underline{x}^T A \underline{x} \\ &= \underline{x}_k^T A \underline{x}_k + 2 \underline{x}_k^T \underline{b} + (A^{-1} \underline{b})^T \underline{b} \\ &= 2 F(\underline{x}_k) + (A^{-1} \underline{b})^T \underline{b} \end{aligned}$$

d.h. die Iterierte \underline{x}_k minimiert den Fehler \underline{e}_k in der Energienorm!



4. Gewöhnliche Differentialgleichungen (ODE)

①

4.1. Problemstellung und Bsp

(i) Radioaktiver Zerfall

$m(t) :=$ Menge radioaktiven Materials zur Zeit t

$$\boxed{\frac{dm(t)}{dt} = -\lambda m(t)} \quad ; t \in \mathbb{R}$$

↑
Proportionalitätsfaktor

Lösungen: $m(t) = A e^{-\lambda t}$ (allgemein)

$m(t_0) = m_0$ (vorhandene Menge des Materials zum Zeitpunkt t_0)

$m(t) = m_0 e^{-\lambda(t-t_0)}$ (speziell)

▷ Def: Anfangswertaufgabe $y: t \in \mathbb{R} \mapsto y(t) \in \mathbb{R}$

Gesucht ist eine Funktion $y(t)$, so dass gilt

$$\begin{cases} \frac{dy(t)}{dt} = f(t, y(t)), & (\text{AWP}) \\ y(t_0) = y_0 \end{cases}$$

$= y'(t)$ DGL 1. Ordnung

(ii) Erzwungene Schwingung eines stark gedämpften Schwingers

$m \frac{d^2 y}{dt^2} + \gamma \frac{dy}{dt} + c(y - y_0) = P \cos(\omega t)$

↑ Masse ↑ Dämpfungskraft ↑ Rückstellkraft ↑ Ruhelage ↑ Störkraft

↑ Kreisfrequenz



→ DGL 2ter Ordnung

(2)

Spezielle Werte für m, γ, c, y_0, P und w :

$$\Rightarrow \left\{ \begin{array}{l} \frac{d^2 y}{dt^2} + 200 \frac{dy}{dt} + 156.25 y = 80 \cos(t) + 156.25 \\ \text{(AWP) Anfangszustand (Bedingung)} \\ y(0) = 5 \\ y'(0) = -100 \end{array} \right.$$

Lösungsmethode: Variablensubstitution um die DGL höherer Ordnung auf System 1. Ord. zurückzuführen.

$$\begin{array}{l} y_1 := y \\ y_2 := y' \end{array} \quad \text{mit } \underline{y}(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \end{pmatrix} \in \mathbb{R}^2$$

$$\Rightarrow \left\{ \begin{array}{l} y_1'(t) = y_2(t) \\ \frac{dy_2}{dt} = -156.25 y_1(t) - 200 y_2(t) + 80 \cos(t) + 156.25 \\ y_1(0) = 5 \\ y_2'(0) = -100 \end{array} \right.$$

Allgemein wird eine DGL n -ter Ordnung in ein System erster Ordnung umgeschrieben.



(iii) Autonome DGL

$$y' = f(y)$$

t kommt nicht vor, d.h. $\frac{\partial f}{\partial t} = 0$

▷ Satz: Für $y'(t) = f(y)$ gilt:

Mit $y(t)$ ist auch $y(t-t_0)$ eine Lösung.

▷ Zurückführen ein allgemeines AWP auf ein autonomes

$$\begin{cases} y'(t) = f(t, y(t)) ; y \in \mathbb{R}^n & \text{"Nicht-autonom"} \\ y(t_0) = y_0 \end{cases}$$

Definitionen: $\underline{Y} := \begin{pmatrix} t \\ y \end{pmatrix} \in \mathbb{R}^{n+1} ; \underline{F}(\underline{Y}) := \begin{pmatrix} 1 \\ f(t, y) \end{pmatrix}$

$$\begin{cases} \frac{d}{dt} \underline{Y}(\tau) = \underline{F}(\underline{Y}(\tau)) & \text{"neues autonomes AWP"} \\ \underline{Y}(0) = \begin{bmatrix} t_0 \\ y_0 \end{bmatrix} \end{cases} \text{ mit unabhängiger Variable: } \tau = t - t_0$$

(iv) Existenzsatz für Differentialgleichungen

▷ Satz: (Existenz) $[a,b] \times \mathbb{R}^n$
Es sei $f(t, y) : D_f \subset \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$ eine stetige Fkt.

Dann besitzt die AWAufgabe (AWP) eine Lösung, die bis an den Rand von D_f fortgesetzt werden kann.

▷ Satz: (Eindeutigkeit)

Es sei $f(t, y) : D_f \subset \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$ stetig in D_f und $f(t, y)$ erfülle eine Lipschitzbedingung bezüg y , d.h.



es gibt eine Konstante $L > 0$, so dass:

(4)

$$\|f(t, y) - f(t, \tilde{y})\| \leq L \|y - \tilde{y}\|$$

$\forall (t, y) \in D_f$ und $(t, \tilde{y}) \in D_f$. Dann hat das AWP genau eine Lösung.

Bsp:

$$\begin{cases} y'(t) = -\sqrt{y} \\ y(0) = 0 \end{cases} \text{ erfüllt nicht Lipschitzbed.}$$

$$\Rightarrow \left. \begin{matrix} y_1(t) = 0 \\ y_2(t) = t^2 \end{matrix} \right\} \text{ keine Eindeutigkeit!}$$

4.2. Taylor- und Runge-Kutta-Verfahren

a) Das Euler-Verfahren

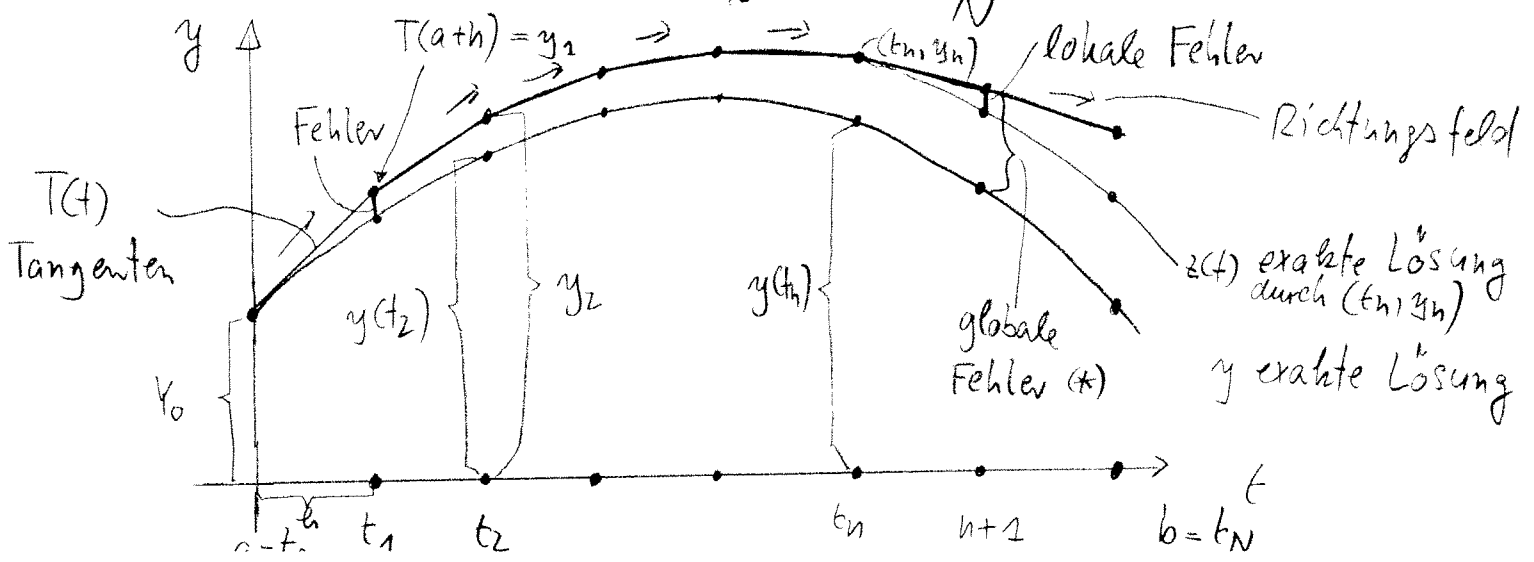
$$\begin{cases} y'(t) = f(t, y(t)) ; t \in [a, b] \\ y(t_0) = y_0 \end{cases}$$

Discretisation:

$$t_i = t_0 + i \cdot h \quad ; i = 1, \dots, N$$

Schrittweite $h = \frac{b-a}{N} = \text{const}$

(*) gibt die Fehlerordnung an





- Die exakte Lösung $y(t_i)$ wird durch numerische Werte y_i approximiert.

(5)

- Richtungsfeld tangential zur Lösungskurve

$$y(t_2) = y(a+h) \sim y_1 = y_0 + h y_0' = y_0 + h f(t_0, y_0)$$

$y(t_2)$ ist nicht bekannt, also beim nächsten Schritt muss man von (t_1, y_1) ausgehen.

$$y(t_2) \sim y_2 := y_1 + h f(t_1, y_1)$$

⇒ Euler-Verfahren für Vektor-DGL

geg: y_0 (AW)

Für $i = 0, 1, \dots, N$ berechne man

$$\underline{y}_{i+1} = \underline{y}_i + h \underline{f}(t_i, \underline{y}_i)$$

$$t_{i+1} = t_i + h$$



Numerik der gewöhnlichen Diff-Glch

28.04.2008

①

$$y'(t) = F(t, y(t))$$

↑ gegeben, "rechte Seite"

$$y(t=0) = y_0 \quad \text{gegeben}$$

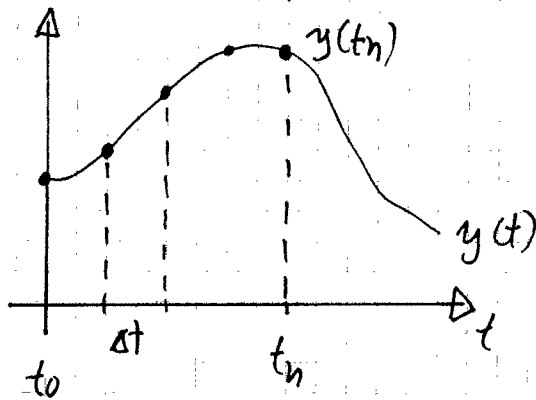
$$y(t) = ?$$

Idee (Numerik): Diskretisierung $t \in [0, T_{\text{end}}]$

$$t_0 = 0; \quad t_{n+1} = t_n + \Delta t_n$$

Hier $\Delta t_n = \Delta t$ ($\forall n$, ^{Schrittweiten} konstante Schrittweite)

$y(t)$ wird approximiert durch $\{y(t_n)\}_{n \in \{0, 1, 2, \dots\}}$

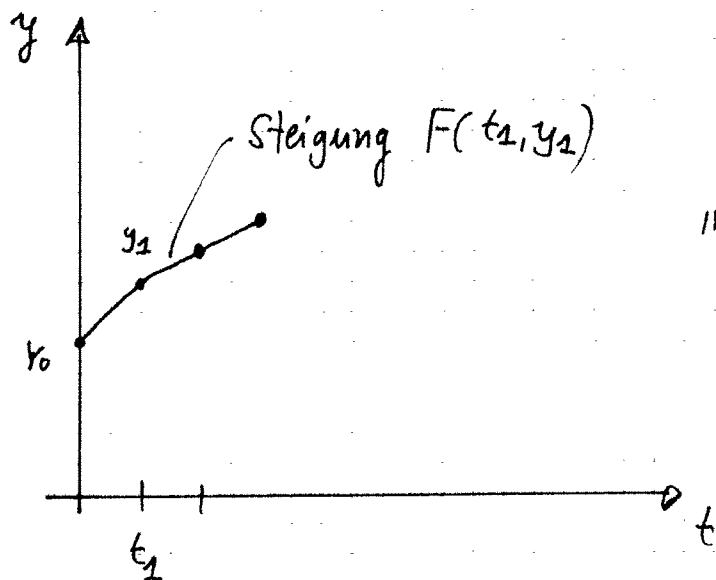


$$y(t_n) \approx y_n := y(t_n)$$

(1) Euler Methode

$$\underbrace{\frac{y_{n+1} - y_n}{\Delta t}}_{y'(t_n)} = F(t_n, y_n)$$

$$\leadsto \boxed{y_{n+1} = y_n + \Delta t F(t_n, y_n)}$$



"expliziter Euler"

$$(2) \quad \frac{y_{n+1} - y_n}{\Delta t} = F(t_{n+1}, y_{n+1})$$

"impliziter Euler"

$$\boxed{y_{n+1} = y_n + \Delta t F(t_{n+1}, y_{n+1})}$$

Bestimmungsgleichung für y_{n+1} . Nicht linear!

(Newton-Verfahren)

$$(3) \quad \frac{y_{n+1} - y_n}{\Delta t} = F(t_{n+1/2}, y_{n+1/2})$$

"Mittelpunktsregel"

$$\parallel \frac{1}{2}(y_n + y_{n+1})$$

$$\boxed{y_{n+1} = y_n + \Delta t F(t_{n+1/2}, \frac{1}{2}(y_n + y_{n+1}))} \quad \text{"implizit"}$$

$$(4) \quad \left\{ \begin{array}{l} y_{n+1/2}^{(*)} = y_n + \frac{\Delta t}{2} F(t_n, y_n) \\ y_{n+1} = y_n + \Delta t F(t_{n+1/2}, y_{n+1/2}^{(*)}) \end{array} \right.$$

"explizite
Mittelpunktsregel"

$$(5) F(t_{n+1/2}, y_{n+1/2}) \hat{=} F_{n+1/2} = \frac{1}{2} [F(t_n, y_n) + F(t_{n+1}, y_{n+1})]$$

$$y_{n+1} = y_n + \Delta t \frac{1}{2} [F(t_n, y_n) + F(t_{n+1}, y_{n+1})]$$

Trapezregel!

Struktur: Runge - Kutta - Methoden

$$k_i = F(t_n + c_i \Delta t, y_n + \Delta t \sum_{j=1}^s a_{ij} k_j)$$

$$y_{n+1} = y_n + \Delta t \sum_{i=1}^s b_i k_i$$

Jedes Verfahren ist durch die Werte a_{ij} , b_i , c_i gegeben.

Butcher - Tabelle

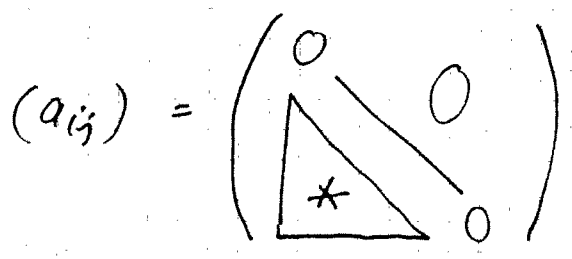
c_1	a_{11}	a_{12}	
\vdots	a_{21}	a_{22}	\dots
\vdots	\vdots	\vdots	\dots
c_s	a_{s1}	\dots	
	b_1	\dots	b_s

z. Bsp.

0	0	0
$\frac{1}{2}$	$\frac{1}{2}$	0
	0	1

(4)

"s" ist die Anzahl F-Auswertungen, Stufen.



\Leftrightarrow Verfahren explizit

Bsp:

0	0	0	0	0
$\frac{1}{2}$	$\frac{1}{2}$	0	0	0
$\frac{1}{2}$	0	$\frac{1}{2}$	0	0
1	0	0	1	0
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

Kürstufige RK-Verfahren

▷ Konvergenz (globale Fehler)

$$g_n := y_n - y(t_n) \quad (\text{am Ende } t_n = n\Delta t = T_{\text{end}})$$

\uparrow \uparrow
approx. exakt

▷ Def: Das Verfahren konvergiert falls

$$\lim_{\Delta t \rightarrow 0} \|g_n\| = 0$$

(Achtung: $n = \frac{T_{\text{end}}}{\Delta t} \rightarrow \infty$)

Konvergenzordnung $p \iff \|g_n\| = \mathcal{O}(\Delta t^p) \quad \Delta t \rightarrow 0$

▷ Satz: Konsistenz + Stabilität \implies Konvergenz

- Konsistenz gibt an wie gut ein einzelner Schritt ist.
- Stabilität gibt an ob ein Fehler für viele Schritte klein bleibt.

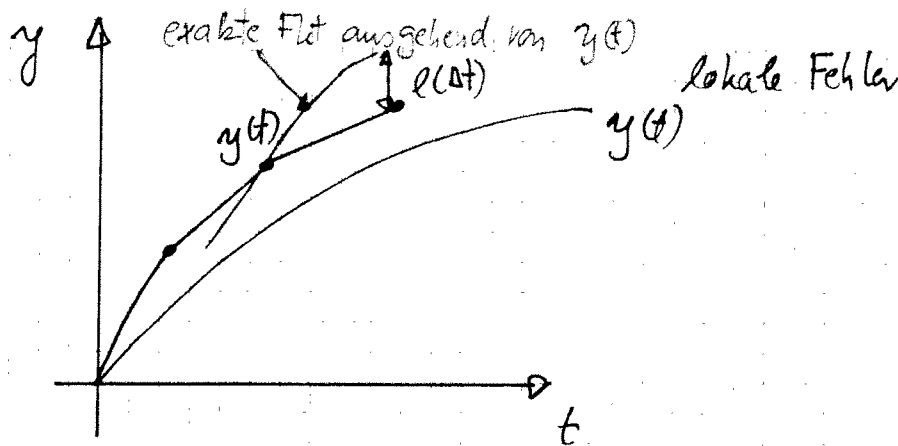
▷ Konsistenz: Abstraktes Verfahren:

$$y_{n+1} = y_n + \Delta t \underbrace{\Phi(t_n, y_n, \Delta t)}_{\text{Einschritt-Methode}}$$

Φ heisst Verfahrensvorschrift.

▷ Def: (Lokaler Fehler) $y(t)$ ist exakt

$$l(t, \Delta t) = y(t) + \Delta t \underbrace{\Phi(t, y(t), \Delta t)}_{\text{Ein Schritt ausgehend von } y(t)} - \underbrace{y(t + \Delta t)}_{\text{exakt}}$$



▷ lokaler Fehler

$$l(t, \Delta t) = y(t) + \Delta t \Phi(t, y(t), \Delta t) - y(t + \Delta t)$$

$$\|l(t, \Delta t)\| = O(\Delta t^{p+2}) \iff p \text{ Konsistenzordnung}$$

Bsp:

$$y_{n+1} = y_n + \Delta t F(y_{n+1/2}) \quad ; \quad y_{n+1} = y_n + \frac{\Delta t}{2} F(y_n)$$

angenommen $y(t)$ ist exakte Lsg. $y'(t) = F(y(t))$

$$l(t, \Delta t) = y(t) + \Delta t F(y(t) + \frac{\Delta t}{2} F(y(t))) - y(t + \Delta t)$$

$\underbrace{y(t)}_{\cong y_n} \quad \underbrace{F(y(t))}_{\cong y_n} \quad \underbrace{\frac{\Delta t}{2} F(y(t))}_{\cong y_n}$

Taylor
 1. Ordnung

$$l(t, \Delta t) = y(t) + \Delta t \left[F(y(t)) + F'(y(t)) \cdot \frac{\Delta t}{2} F(y(t)) + O(\Delta t^2) \right] - y(t + \Delta t)$$

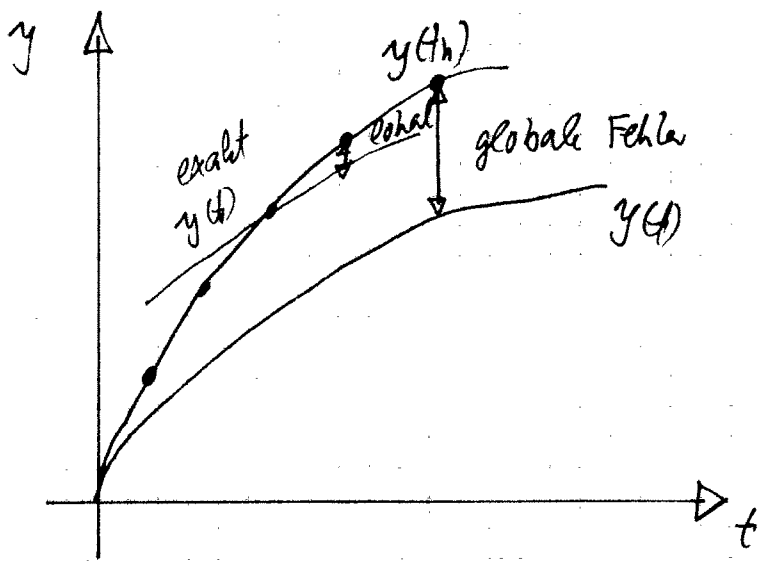
$$y(t + \Delta t) = y(t) + \underbrace{y'(t)}_{F(y(t))} \Delta t + \underbrace{y''(t)}_{F'(y(t)) \cdot F(y(t))} \frac{\Delta t^2}{2} + O(\Delta t^3)$$

$$\leadsto l(t, \Delta t) = O(\Delta t^3)$$

Konsistenzordnung : $p = 2$; Euler : $p = 1$

Runge-Kutta : $p = 4$

Falls die Konsistenzordnung p ist und das Verfahren stabil, dann ist die Konvergenzordnung auch p .



Der lokale Fehler summiert sich auf.
Deshalb $l = O(\Delta t^{p+1})$
aber $g = O(\Delta t^p)$

▷ Stabilität:

Wir betrachten die lineare Test Gleichung.

$$y' = \lambda y \quad (\text{z. Bsp. durch Linearisierung})$$

$$y' = F(y) = \underbrace{F'(0)}_{\hat{=} \lambda} y \quad ; \quad y \ll 1 \quad \lambda \in \mathbb{C} \text{ erlaubt}$$

▷ im Systemfall ist λ Eigenwert der Jacobi $F'(y_0)$

Lsg: $y(t) \sim e^{\lambda t}$

Weiter: Wir wählen ein Verfahren auf $y' = \lambda y$ an.

$$y_{n+1} = R(\lambda \Delta t) y_n$$

$$z := \lambda \Delta t$$

$R(z)$ heisst Verstärkungsfaktor

R verstärkt oder dämpft auch den Fehler \leadsto Stabilität

Def: (Stabilität)

Falls, für eine abfallende Lsg. d.h. $\operatorname{Re}(\lambda) < 0$, ~~es~~ der Verstärkungsfaktor $R(z)$ die Bedingung $\|R(z)\| \leq 1$ erfüllt, so ist das Verfahren (linear) stabil.

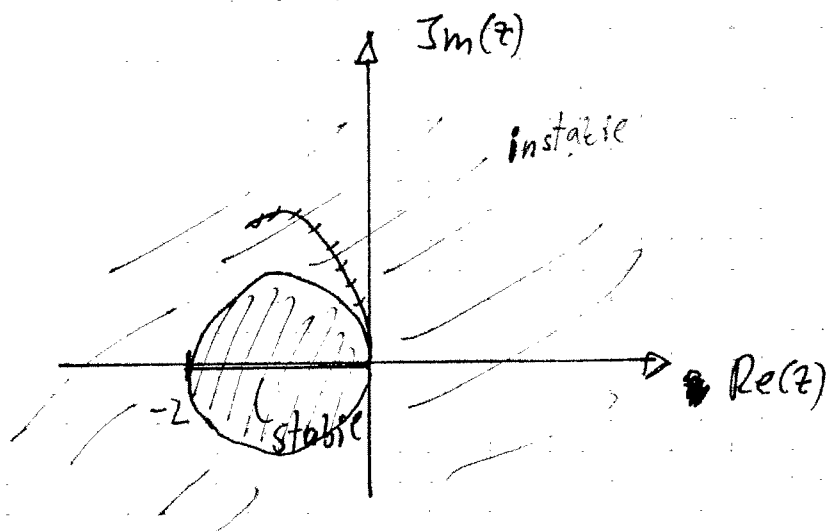
Bsp: Euler (explizit)

$$\begin{aligned} y_{n+1} &= y_n + \Delta t F(y_n) \quad ; \quad y' = \underbrace{\lambda y}_{F(y)} \\ &= y_n + \Delta t \lambda y_n \\ &= \underbrace{(1 + \Delta t \lambda)}_{R(\lambda \Delta t)} y_n \end{aligned}$$

$$\leadsto R(z) = 1 + z ;$$

Bedingung: Für $\lambda < 0$ muss $|R(z)| \leq 1$

$$|1+z| \leq 1 \quad ; \quad -2 \leq z \leq 0$$



Wähle : $\lambda \Delta t < 0$ ✓

$$\lambda \Delta t > -2 \quad \leadsto \quad \Delta t < \frac{2}{|\lambda|}$$

☛ Ist nur stabil für hinreichend kleine Zeitschritte.

▷ Def: (Stabilitätsgebiet)

$$S := \left\{ z \in \mathbb{C} \mid |R(z)| \leq 1 \right\}$$

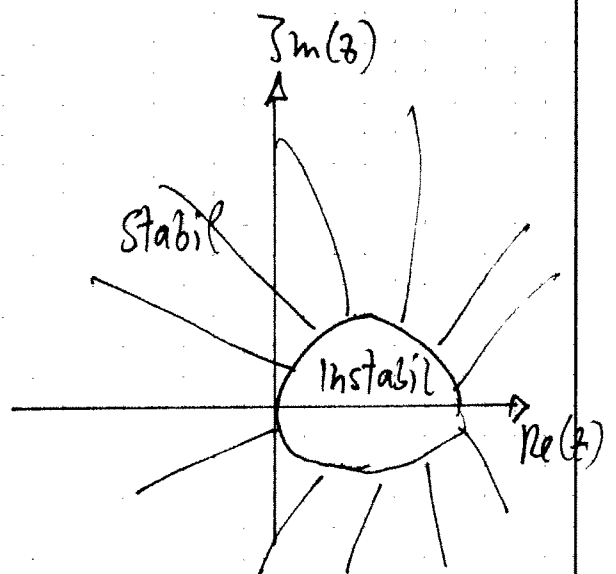
Bsp: impliziter Euler "

$$y_{n+1} = y_n + \Delta t \lambda y_{n+1}$$

$$y_{n+1} = \frac{1}{1 - \Delta t \lambda} y_n$$

$$\leadsto R(z) = \frac{1}{1-z}$$

Hier keine Zeitschrittbegrenzung.



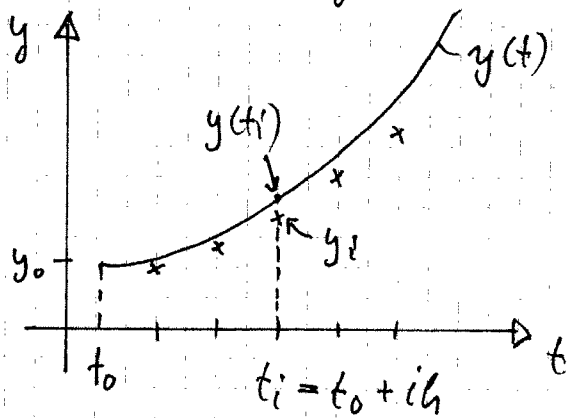
AWA gew. Diffgl.

$$y' = f(t, y(t)) \quad (1)$$

$$\text{Anfangsbedingung } y(t_0) = y_0 \quad (2)$$

Ges: $y(t)$ mit

Diskretisierung



$y_i \sim y(t_i)$
 num. appr. exakt

h Schrittweite

$$\text{lokaler Fehler: } l(t_0, h) = y(t_0 + h) - y_h = \mathcal{O}(h^{p+1})$$

$p = \text{Ordnung}$

Konsistenz: $p \geq 1$

4.4. Steife Diff. gl.

Testgleichung: $y' = qy \quad q \in \mathbb{C}$

Einschrittverfahren:

$$y_{n+1} = \underbrace{V(qh)}_{\text{Verstärkungsfaktor}} y_n \quad ; \quad qh \doteq z$$

$n = 0, 1, 2, \dots, \infty$

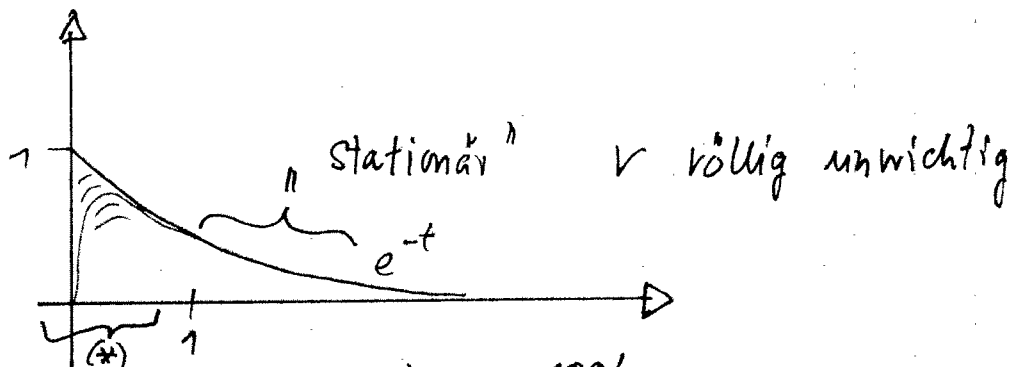
Stabilitätsgebiet:

$$S = \{z \in \mathbb{C} \mid |V(z)| \leq 1\}$$
$$= \{z \in \mathbb{C} \mid \{y_n\}_{n=0,1,\dots} \text{ beschränkt}\}$$

Warum?

a) Einführungsbsp:

$$\begin{cases} y'' + 101y' + 100y = 0 \\ y(0) = 0 \\ y'(0) = 99 \end{cases}$$



Exakt: $y(t) = e^{-t} - e^{-100t}$

(*) transiente Phase.

umschreiben in ein System

$$\begin{aligned}
 y &= z \\
 z &= y'
 \end{aligned}
 \quad
 \begin{pmatrix} y \\ z \end{pmatrix}' = \begin{pmatrix} 0 & 1 \\ -101 & 100 \end{pmatrix} \begin{pmatrix} y \\ z \end{pmatrix}$$

EW von $\begin{pmatrix} 0 & 1 \\ 101 & 100 \end{pmatrix}$

$$q_1 = -100; \quad q_2 = -1$$

Diagonalisierbar:

$$\begin{aligned}
 u' &= -u && : h, \Phi \\
 v' &= -100v && : h, \Phi
 \end{aligned}$$

Num. Lösung mit Euler:

$$y' = qy; \quad q_1 = -100; \quad q_2 = -1$$

$$y_{(h)} = (1 + hq)^n$$

Stabilität:

$$\left. \begin{aligned} |1 + hq_2| &\leq 1 \\ |1 + hq_1| &\leq 1 \end{aligned} \right\} \Rightarrow 0 < h \leq \frac{2}{-q_1}$$

$$y(0,2) = e^{-0,2} - e^{-20} = \underbrace{0,81873\dots}_u - \underbrace{2 \cdot 10^{-9}}_v$$

- (ii) Für $t \geq 0,2$ möchte man mit grösserem Schritt rechnen, angepasst an $-q_2 h_2 = 0,1 \rightarrow h_2^* = 0,1$. Sobald man das macht wird $1 + q_1 h^* = -9$. \Rightarrow Der unwesentliche Beitrag zur Lösung erzwingt

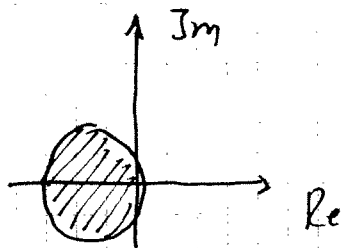
viel zu kleinen Schritt. \Rightarrow Steifheit.

\Rightarrow Stabilitätsgebiet.

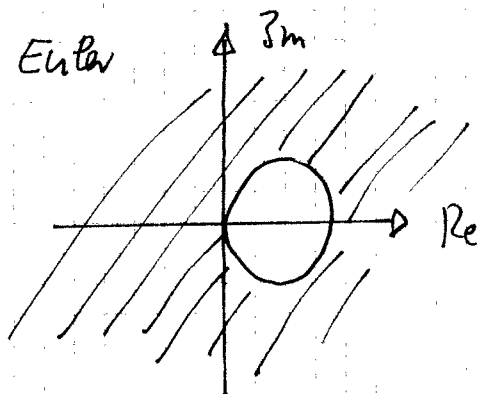
b) θ -Verfahren

$$y_{n+2} = y_n + h \left[(1-\theta) f(t_n, y_n) + \theta f(t_{n+1}, y_{n+1}) \right]$$

$\theta = 0$, Euler



$\theta = 1$, imp. Euler



Genauigkeit:

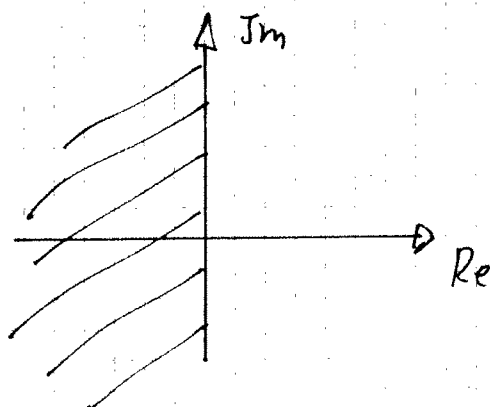
$$(1 + h\eta) \sim e^{h\eta}$$

anwenden auf $\theta = \frac{1}{2}$ (Trapezregel)

$$y_{n+2} = y_n + \frac{h}{2} \left[\underbrace{f(t_n, y_n)}_{\eta y_n} + \underbrace{f(t_{n+1}, y_{n+1})}_{\eta y_{n+1}} \right]$$

auf $y' = \eta y$

$$y_{n+1} = \underbrace{\frac{1 + z/2}{1 - z/2}}_{V(z)} y_n$$



c) Auswirkung des Stützgebietes auf die Wahl der Schrittweite

$$y' = Ay$$

$A \in \mathbb{R}^{m \times m}$ diagonalisierbar

$$\text{d.h. } y' = \underbrace{T^{-1}}_J \underbrace{A}_\Lambda \underbrace{T}_z y \quad (1)$$

$\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$

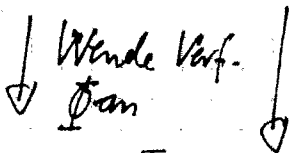
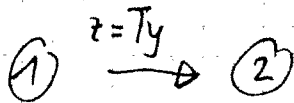
zur Analyse umtransformieren

$$z = Ty$$

$$z' = Ty' = \underbrace{TT^{-1}}_J \underbrace{A}_\Lambda \underbrace{T}_z y \quad (2)$$

$$z' = \Lambda z$$

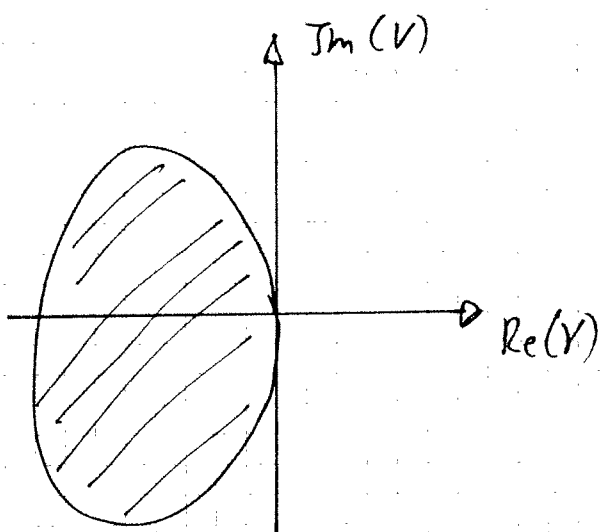
$$\left. \begin{aligned} z^1' &= \lambda_1 z^1 \\ z^2' &= \lambda_2 z^2 \end{aligned} \right\} \text{entkoppelt}$$



$$\{y_n\} \xrightarrow{\tilde{z}_n = Ty_n} \{z_n\}$$

Bei praktisch allen Verfahren gilt $\{\tilde{z}_n\} = \{z_n\}$

Typisches Stabgebiet eines expliziten Verfahrens



In der transienten Phase müssen alle Komponenten "genau" berechnet werden. Später ist $z^{(1)}$ unwichtig. Trotzdem muss $q_2 h \in S$ liegen.

Steife Probleme: \Rightarrow verwende implizite Verfahren.

A-Stabilität:

$$C^- = \{z \in \mathbb{C} \mid \operatorname{Re}(z) \leq 0\} \subset S$$

Imp. Euler: A-stabil; $p = 1$

Trapezregel: - " - ; $p = 2$

d) Rückwärtsdifferenzierungsformeln (BDF)

Wähle Interpolationspolynom $P(t)$ durch die P_k te.

$$(t_i, y_i) ; i = -k+1, \dots, 1$$

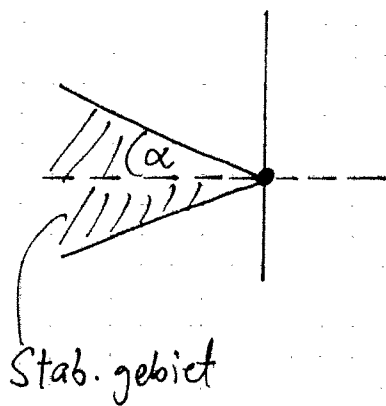
$$P'(t_{n+1}) = f(t_{n+1}, P(t_{n+1}))$$

\Rightarrow BDFarmeln

$h = 1$: impliziter Euler

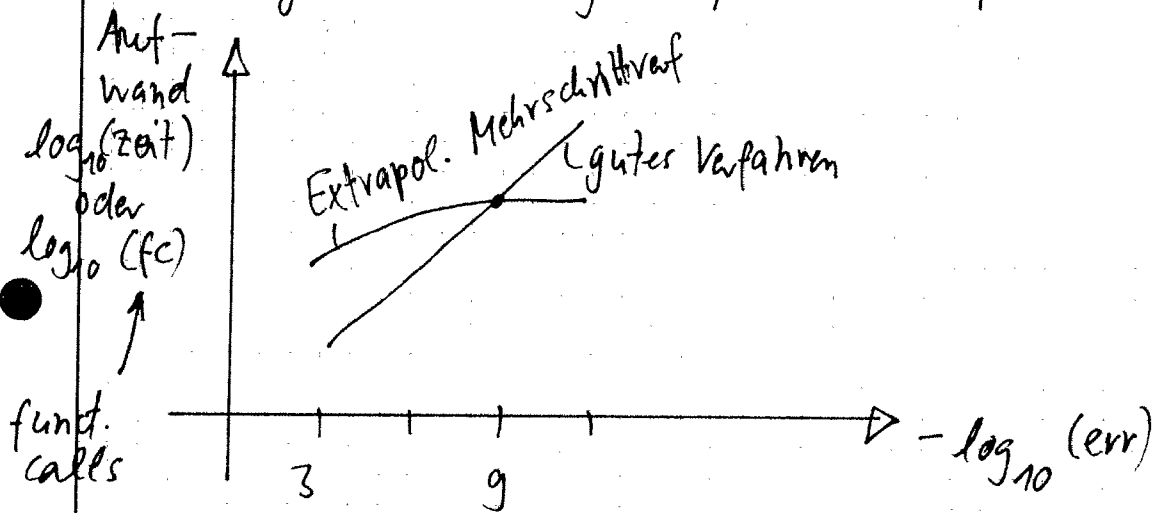
$$\sum \alpha_{ki} y_{n+i} = f(t_{n+1}, y_{n+1})$$

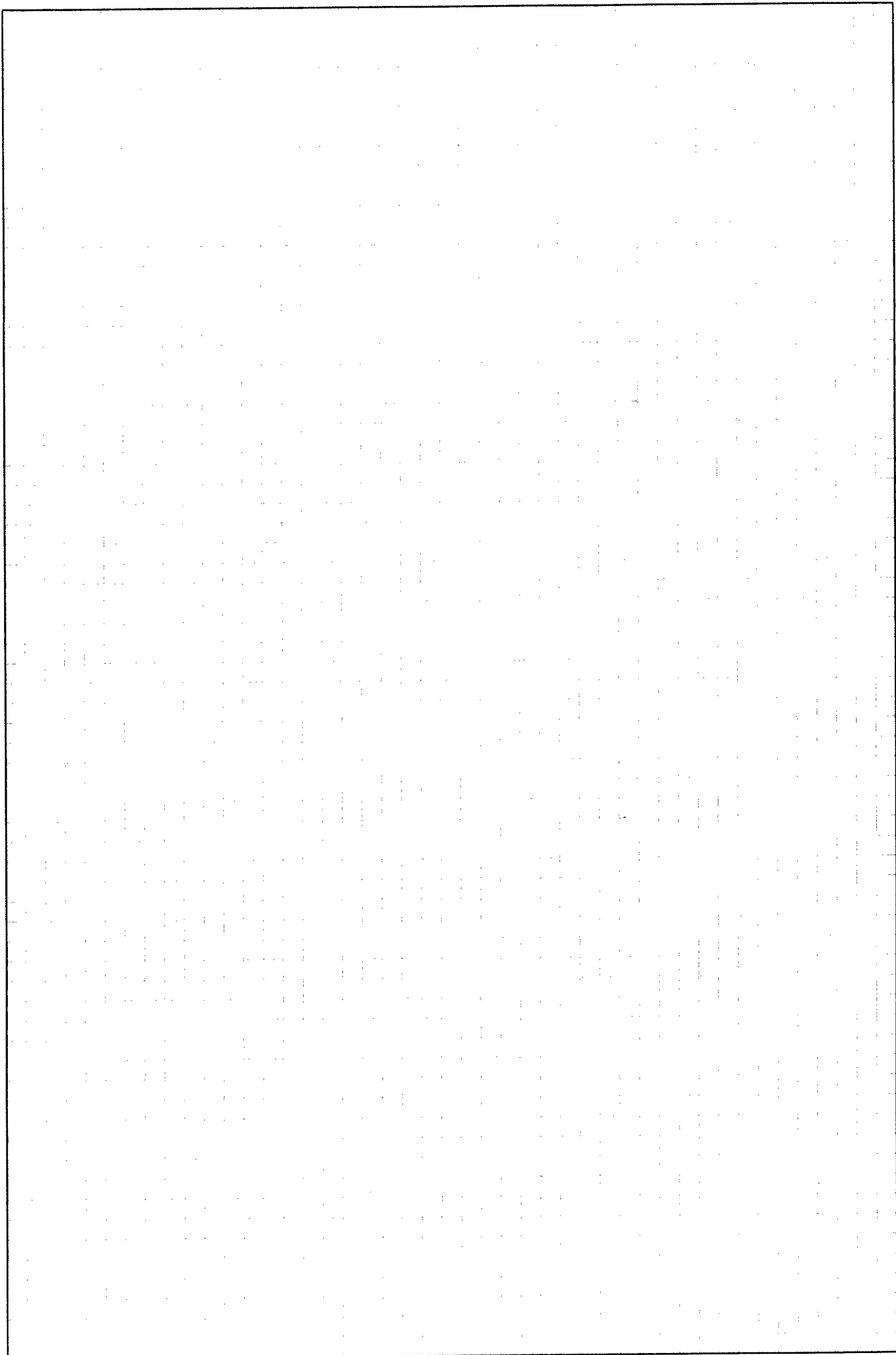
$A(\alpha)$ - stabil



e) Zur Praxis

Diagram: Genauigkeit, versus Aufwand





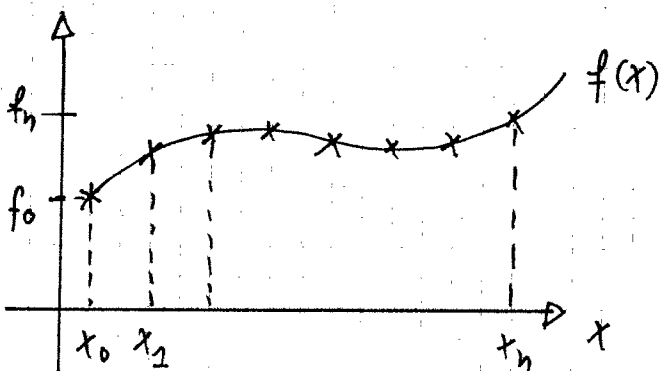
6. Interpolation

19.05.2008

Motivation : Vor allem zum Entwickeln von Verfahren,
z. Bsp. Rückwärtsdifferenzierungs Formeln.

6.1. Problemstellung

Rekonstruktion, Approximation einer Funktion aus
Funktionswerten.



$f(x_i) \hat{=} f_i$
diskrete Werte von
 $f(x)$

Gegeben : $n+1$ Stützstellen x_i
 $n+1$ Funktionswerte, f_i ; $i = 0, 1, \dots, n$

Gesucht : Approximation $p(x)$ für $f(x)$
mit : $\boxed{p(x_i) = f_i}$, $i = 0, 1, 2, \dots, n$
heißt : Interpolationsbedingungen

Es gibt ∞ -viele Lösungen. Wähle Funktionsklasse aus.
z. Bsp. trigonometrisches Polynom

Einfachster Fall :

$p(x)$ Polynom vom Grad $\leq n$

$n+1$ Parameter

$$\boxed{p(x) := a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n} \quad (2)$$

Gibt es eine eindeutige Lösung?

6.2. Lagrangepolynom

Wie soll man es nicht machen.

Interpolationsbedingung:

$$\Rightarrow \begin{pmatrix} 1 & x_0 & x_0^2 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^n \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ \vdots \\ a_n \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_n \end{pmatrix}$$

$\hat{=}$ Van der Monde matrix V

Bem: • V ist oft schlecht konditioniert.

• Gaußelim. $\frac{n^3}{3}$ -Operationen. (zu teuer)

Gute Methode: Basisfunktionen von Lagrange. $B_k(x)$

$$(3) \quad B_k(x_l) = \begin{cases} 1 & \text{für } k=l \\ 0 & \text{sonst} \end{cases}$$

$$(4) \quad B_k(x) = \prod_{\substack{l=0 \\ l \neq k}}^n (x - x_l) \bigg/ \prod_{\substack{l=0 \\ l \neq k}}^n (x_k - x_l)$$

Grad n

$$\text{Dann erfüllt } p(x) = \sum_{k=0}^n f_k B_k(x) \quad (5)$$

dann erfüllt die Interpol.-beding. (1)

Bem: Grad $\leq n$
Existenz gesichert.

Eindeutigkeit? Gl. system regulär.

Satz: Für $n+1$ verschiedene Stützstellen $x_k, k=0, 1, \dots, n$ und beliebige Stützwerte $f_k, k=0, 1, \dots, n$ gibt es genau ein Interpol. polynom. vom Grad $\leq n$ mit

$$p(x_k) = f_k$$

6.3. Die baryzentrische Formel

Effiziente Auswertung von p an gegebener "neuer" Stelle x .

$$L(x) := \prod_{i=0}^n (x - x_i) \quad ; \quad \text{Produkt aller linear Faktoren}$$

$$w_k := \frac{1}{\prod_{\substack{i=0 \\ i \neq k}}^n (x_k - x_i)} \quad ; \quad k = 0, 1, \dots, n$$

"Stützhoeff"

$$(4) \Rightarrow B_k(x) = L(x) \cdot \underbrace{\frac{w_k}{(x - x_k)}}_{:= c_k(x)} = L(x) \cdot c_k(x)$$

$$(5) \Rightarrow p(x) = L(x) \sum_{k=0}^n f_k \cdot c_k(x) \quad (6)$$

(6) gilt für beliebige Wahl von $\{f_k\}$

z.B. $f_k = 1, k=0, 1, \dots, n \Rightarrow p(x) = 1$

in (6) \Rightarrow

$$1 = L(x) \sum_{k=0}^n c_k(x) \rightarrow L(x) = \frac{1}{\sum_{k=0}^n c_k(x)}$$

Algorithmas zum Auswerten von $p(x)$

(i) Berechne Stützkoeff. $w_k = \frac{1}{\prod_{\substack{i=0 \\ i \neq k}}^n (x_k - x_i)}$; $k = 0, 1, \dots, n$

Aufwand : $O(n^2)$ (unabhängig von f_k und x)

(ii) Für Neustellen x

$$c_k := \frac{w_k}{x - x_k} ; k = 0, 1, \dots, n$$

$$p(x) = \frac{\sum_{k=0}^n c_k f_k}{\sum_{k=0}^n c_k}$$

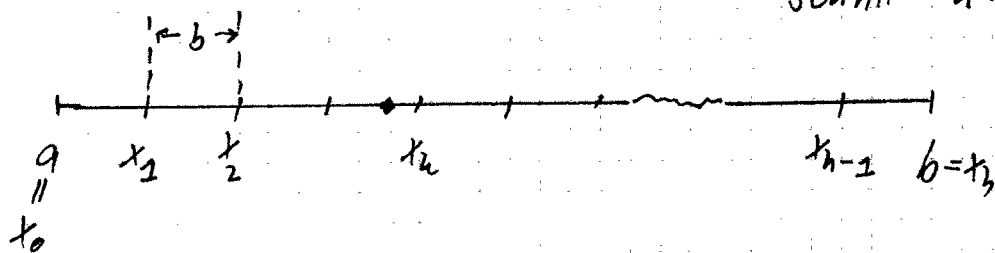
Aufwand : $O(n)$ ($2n + 3$ Punktoper.)

Die Stützkoeffizienten

Zwei wichtige Fälle

a) Äquidistante Stützstellen

Schritt $h > 0$



$$w_k = \frac{1}{h^n h (h-1)(h-2) \cdot 2 \cdot 1 \cdot (-1)^{n-k} \cdot 1 \cdot 2 \cdot 3 \cdot \dots \cdot (n-k)}$$

$$w_k^* = (-h)^n n! w_k = (-1)^k h \binom{n}{k} ; k = 0, 1, \dots, n$$

Binomialkoeff. mit altern. Vorzeichen.

Bsp: $n=3$

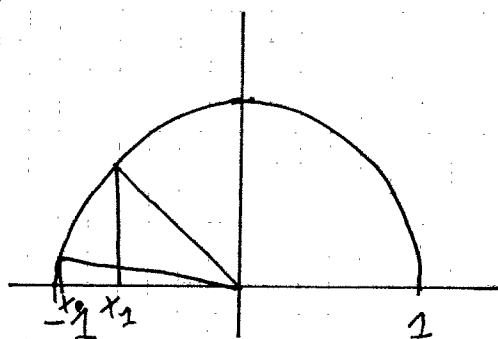
$$w_k^* = \{1, -3, 3, -1\}$$

b) Tschebyschev - Punkte

$$x_k := \frac{a+b}{2} + \frac{a-b}{2} \cos\left(\frac{(k+\frac{1}{2})\pi}{n+1}\right)$$

O. b. d. A. $a=1, b=-1$

$$x_k = \cos\left(\underbrace{\frac{(k+\frac{1}{2})\pi}{n+1}}_{\varphi_k}\right)$$



$$\frac{1}{w_k} = \prod_{\substack{i=0 \\ i \neq k}}^n (x_k - x_i) = \lim_{x \rightarrow x_k} \frac{\prod_{i=0}^n (x - x_i)}{x - x_k} \cdot \frac{x - x_k}{x - x_k}$$

$$= \lim_{x \rightarrow x_k} \frac{L(x) - L(x_k)}{x - x_k}$$

$$w_k = \frac{1}{L'(x_k)} \quad (7)$$

Berechne: $L(x)$ für T-Punkte

Setze: $x = \cos(\varphi)$

Beh: $L(x) = c \cdot \cos((n+1) \cdot \varphi) \quad (8)$

\hookrightarrow const.

denn: Muss gelten: $L(x_k) = 0$

$$L(x_k) = c \cdot \cos((n+1) \varphi_k) = c \cdot \cos\left(\left(k+\frac{1}{2}\right)\pi\right) = 0$$

Differenziere (8) nach φ

$$(8) \Rightarrow -L'(x) \sin(\varphi) = (-1)^k \sin((n+2)\varphi) \cdot (n+1)$$

$$w_k = \frac{1}{L'(x_k)} = (-1)^k (-1) \sin\left(\frac{(k+1/2)\pi}{n+1}\right)$$

6.4. Problematisch bei der Polynominterpolation

Satz: Sei f auf dem Intervall $J := [a, b]$,
 $n+1$ stetig diffbar.
und $p(x)$ das zu $x_0 < x_1 < x_2 < \dots < x_n$
gehörige Interpolationspolynom vom Grad $\leq n$
mit $p(x_k) = f(x_k) := f_k$

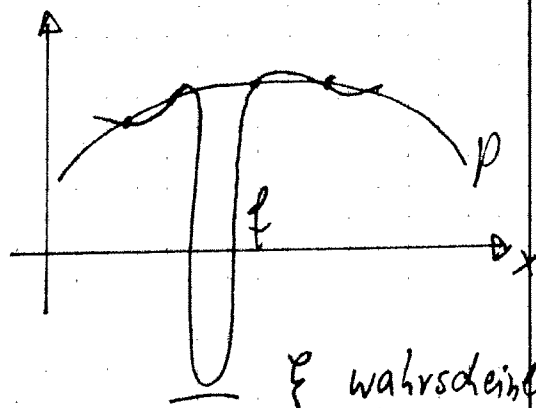
Dann existiert zu einer beliebigen Neustelle

$$x \in J \text{ ein } \xi_x \in J = [a, b] \text{ mit}$$
$$f(x) - p(x) = \frac{1}{(n+1)!} L(x) f^{(n+1)}(\xi_x)$$

Folgerung: Grosse Interpolationsfehler können 2 Ursachen haben.

$$(i) \left| f^{(n+1)}(\xi_x) \right| \text{ gross}$$

Es liegt an f .



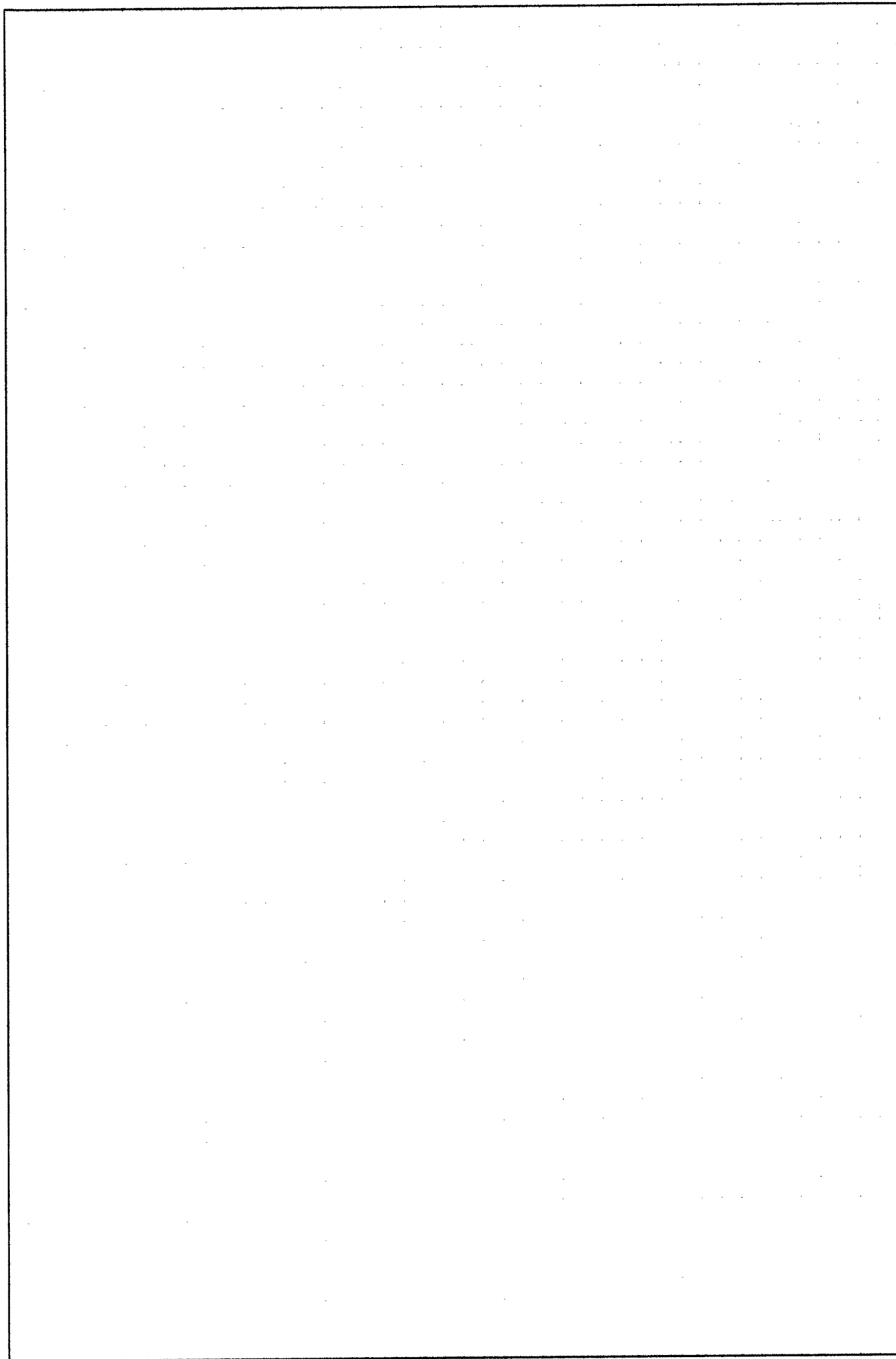
ξ_x wahrscheinlich hier

Abhilfe : Lege Stützstellen dichter.

(ii) $|L(x)|$ gross

Viele äquidistante Stützstellen sind ungünstig

Besser verdichte Stützstellen am Intervallrand z.B. : T-Punkte



① Fixpunktiteration

Gleichung, Form, $x_{k+1} = g(x_k)$ für $k = 0, 1, 2, \dots$

Konvergenz, abstoßend, anziehend $|g'(ξ)| < 1$
(linear)

Banachscher Fixpunktsatz,

Nullstellensuche:

(Skalar + N Dimensionen)

$f(x) = 0$, Newton-, Sekantenverfahren, Quasi-Newton

$$x_{k+1} = x_k - \frac{f(x_k)}{\underbrace{f'(x_k)}_{\neq 0}} \quad \text{„linearisierung“}$$

Konvergenz: Ordnung, Abbruchkriterium (absoluten, relativen Fehler)
Effizienz, Jacobi-Matrix

② Lineare Gleichungssysteme

regulär, singular

Gauß, Faktorisierung, LU-Zerlegung (LR)

Permutationsmatrix

Normen (Höldersehe, p-Norm); Matrixnorm: $\|A\| := \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$

Kondition, Einfluss der Fehler

QR-Zerlegung (Ausgleichsrechnung), Givensrotation.
(Least squares)



③ Eigenwertprobleme

Definitionen und Numerik (schwingende Membran) $\begin{cases} \Delta u = \lambda u \\ u = 0 \end{cases}$
Vektiteration, Rayleigh-Quotient
(Inverse) und QR für alle EW und EV

④ Singularwertzerlegung

für Rang einer Matrix, lineare Abb., Ausgleichsproblem

⑤ Iterative Methoden für LGS

Jacobi-Iteration, CG: Minimierung einer quadratischen Funktion

$A = A^T$ positiv definit

A-orthogonal (konjugiert), A-Norm
Algorithmus

Konditionierung

$$\|x\|_A = \sqrt{x^T A x}$$

⑥ Numerik der ODEs

DGL 1. und 2. Ordnung, AWP, Autonome DGL, Systeme

$$\begin{cases} y'(t) = f(t, y(t)) \\ y(t_0) = y_0 \end{cases}$$

Euler, (implizit), Mittelpunktsregel, ...

RK-Verfahren:

$P(i)$

$$P(i) \leq i-1$$

$$\begin{cases} k_i = f\left(t_n + c_i \Delta t, y_n + \Delta t \sum_{j=1}^s a_{ij} k_j\right) \\ y_{n+1} = y_n + \Delta t \sum_{i=1}^s b_i k_i \end{cases}$$



globaler, lokaler Fehler, Konsistenz

(3)

Konsistenz + Stabilität \implies Konvergenz

$\| \ell(t, \Delta t) \| = \mathcal{O}(\Delta t^{p+1}) \iff p$ Konsistenz (Konvergenz) Ordnung

(Taylor)

Stabilität: $y' = \lambda y$, $\lambda \in \mathbb{C}$, $y(t) \sim e^{\lambda t}$

$$y_{n+1} = R(\lambda \Delta t) y_n$$

Stife Diffgl. $A(\omega)$, A , Stabilität

⑦ Interpolation

Bedingung: $p(x_i) = f_i$

Lagrangepolynom; baryzentrische Formel

Stützkoeffizienten, Tschebyscheff-Punkte

