

Inhalt:

0. Gleitkommaarithmetik, etc

1. a) Eliminationsverfahren zur Lsg. (in Gl.)

1 b) Lineare Ausgleichsrechnung und SVD (Singularwertzerlegung)

2 a) Interpolation

b) numerische Quadratur

3. Iterationsverfahren

a) nichtlineare Gleichungssysteme

b) lineare GL

4. Eigenwertproblem

0. Grundlagen, ZahlendarstellungSei $r \in \mathbb{R}$, „ p -adische Darstellung“

$$r = \epsilon \sum_{k=-\infty}^n a_k b^k \quad b > 1, b \in \mathbb{N}, k \in \mathbb{Z}, 0 \leq a_k < b$$

 $\epsilon = \pm 1$, „Vorzeichen“ $a_k < b-1$ für unendlich viele k
 $a_k \neq 0 \Leftrightarrow r \neq 0$, b -Basis, Bsp $b=10$ Decimal
Gleitkommazahlen (normiert)

$$x = \epsilon \cdot f \cdot b^e \quad \begin{matrix} \text{Exponent} \\ e \in \mathbb{Z} \end{matrix}, b^{-1} \leq f < 1$$

Mantisse

$$\text{Bsp.: } x = 3,1415 \quad x = +0,31415 \cdot 10^{-1} \quad (e=-1)$$

$$f = 0, d_1, \dots, d_m = \sum_{k=1}^m d_k b^{-k}, d_1 \neq 0 \quad (1 \leq d_k < b) \quad k=1, \dots, m$$

m-Stellenzahl: $r \leq e \leq R$

$$IM(b, m, (r, R)) = \mathbb{N}(b, m)$$

Single Precision: $b=2, m=24$ Double Precision: $b=2, m=53, r=-1024, R=1023$ kleinste darstellbare Zahl $|x_{\min}| = 0,1 \cdot b^r = b^{r-1} (\neq 0) > \text{UNDERFLOW}$ größte darstellbare Zahl $|x_{\max}| = 0, (b-1) \dots (b-1) \cdot b^R = (1 - b^{-m}) b^R < \text{OVERFLOW}$

Rundungen

$f\ell: \mathbb{R} \rightarrow M(b, m)$
 $x \in \mathbb{R}, f\ell(x) = \begin{cases} \sum_{k=1}^m d_k b^{-k}, & \text{falls } d_{m+1} < \frac{b}{2} \\ \sum_{k=1}^m d_k b^{-k} + b^{-m}, & \text{falls } d_{m+1} \geq \frac{b}{2} \end{cases}$
 „Standardrundung“

Bsp: $b=10, m=3, x_1=0,25367, f\ell(x_1)=0,254$
 $x_2=0,2535, f\ell(x_2)=0,254$
 $x_3=0,2533, f\ell(x_3)=0,253$

absoluter Fehler: $|x - f\ell(x)|$

relative Fehler: $\frac{|x - f\ell(x)|}{|x|} \leq \frac{\frac{b}{2} b^{-(m+1)} \cdot b^e}{b^{-1} \cdot b^e} = \frac{b^{-(m+1)}}{2} = \text{eps}$ „Maschinen-Genauigkeit“

$$\text{eps} = \inf \{ \delta > 0 : f\ell(1+\delta) > 1 \}$$

Gleitkommarithmetik

arithm. $\cdot (+, -, \cdot, \div)$

Seien $x, y \in M(b, m)$

$$x \oplus y = f\ell(x \diamond y) \in M(b, m)$$

Es gilt $x \oplus y = (x \diamond y)(1 + \varepsilon)$, mit $|\varepsilon| \leq \text{eps}$

Bsp.: $M(10, 3), x = 6580 = 0,659 \cdot 10^4$

$$y = 1 = 0,1 \cdot 10^1, z = 4 = 0,4 \cdot 10^1$$

$$x + y + z = 6590 + 1 + 4 = 6595 = 0,6595 \cdot 10^4$$

$$f\ell(x + y + z) = 0,660 \cdot 10^4 = 6600$$

$$x \oplus y = f\ell(0,6591 \cdot 10^4) = 0,659 \cdot 10^4$$

$$(x \oplus y) \oplus z = f\ell(0,6594) \cdot 10^4 = 0,659 \cdot 10^4$$

$$\tilde{y} \oplus z = f\ell(0,1 + 0,4) \cdot 10^1 = 0,5 \cdot 10^1 = 5$$

$$x \oplus (y \oplus z) = f\ell(6590 + 5) = f\ell(6595) \cdot 10^4 = 0,660 \cdot 10^4$$

$\Rightarrow (x \oplus y) \oplus z \neq x \oplus (y \oplus z)$ nicht assoziativ

Grundregel 1: Summiere in Reihenfolge auf steigender Beträge!

Bsp: $x - y = 0,73563 - 0,73441 = 0,00122 = 0,122 \cdot 10^{-2} (M(10,3))$

$$f\ell(x \ominus f\ell(y)) = 0,736 \ominus 0,734 = 0,002 = 0,2 \cdot 10^{-2}$$

$$\text{Relativer Fehler: } \frac{0,122 \cdot 10^{-2} - 0,2 \cdot 10^{-2}}{0,122 \cdot 10^{-2}} \approx 0,64 = 64\%$$

Grundregel 2: Vermeide Auslöseeffekte bei Subtraktion bei nahezu gleich großen Zahlen!

Auslösung

1.4. Landau-Symbol

Def. $f, g : D \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$

$$x_0 \in \bar{D} \subset \bar{\mathbb{R}}$$

$$1) f(x) = O(g(x)), \text{ falls } \left| \frac{f(x)}{g(x)} \right| \leq C \quad \forall x \in U_g(x_0)$$

$$2) f(x) = o(g(x)), \text{ falls } \lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = 0 \quad (x \neq x_0)$$

$$\text{Bsp.: } f(x) = \frac{x^2 + 1}{x^3 + 1} \quad (x \neq -1),$$

$$1) f(x) = O\left(\frac{1}{x}\right), x \rightarrow 0 \text{ (d.h. } x_0 = 0)$$

$$2) f(x) = o\left(\frac{1}{x^2}\right)$$

1.5. Fehlerfortpflanzung

X, Y norm. Räume z.B. $X = \mathbb{R}^n$

$f: X \rightarrow Y$ ($f: D \rightarrow Y, D \subset X$ Gebiet)

$$x \mapsto f(x) = y, \tilde{x} = x - \Delta x, \tilde{y} = y - \Delta y$$

$$S_x = \frac{\|\Delta x\|}{\|x\|} \left(= \frac{\|\Delta x\|}{\|x\|}, X, Y = \mathbb{R} \right) \quad \delta_y = \frac{\|\Delta y\|}{\|y\|} \quad \left. \right\} \text{relativer Fehler}$$

(abs. Fehler $\|\Delta x\|$)

Problem ist gut konditioniert, falls (Problem: $x \mapsto y = f(x)$)

$$\text{Fr: } \frac{\delta_y}{\delta_x} \text{ „klein“} (\delta_x \rightarrow 0) \quad \tilde{y} = f(\tilde{x}) = y + \Delta y$$

Sei $f: \mathbb{R}^n \rightarrow \mathbb{R}, x \in \mathbb{R}^n, x \mapsto f(x) \in \mathbb{R}, f \in C^k(\mathbb{R}^n), k = 1, 2, 3$

$$\|x\|^2 = \|x\|_2^2 := \sum_{i=1}^n |x_i|^2, \quad x \in \mathbb{R}^n, \tilde{x} \in U_g(x), \quad x = (x_1, x_n)^T$$

$$f(\tilde{x}) = f(x) + f'(x)(\tilde{x} - x) + O(\|\tilde{x} - x\|)$$

$$= f(x) + f'(x)(\tilde{x} - x) + \frac{1}{2} \underbrace{(\tilde{x} - x)^T}_{H_f(x)} f''(x) (\tilde{x} - x) + O(\|\tilde{x} - x\|^2) \quad (\text{bzw. } O(\|\tilde{x} - x\|^3))$$

„Taylorentwicklung“

$$\text{Jacobi-Matrix: } f'(x) = (\nabla f(x))^T = \left(\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_n}(x) \right)$$

$f': \mathbb{R}^n \rightarrow \mathbb{R}$ linear

$$f(\tilde{x}) = f(x) + \sum_{j=1}^n \frac{\partial f}{\partial x_j} (\tilde{x}_j - x_j) + O(\|\tilde{x} - x\|)$$

lineare Näherung

$$\begin{aligned}
 \left| \frac{\Delta y}{y} \right| &= \left| \frac{f(\tilde{x}) - f(x)}{f(x)} \right| = \sum_{j=1}^n \underbrace{\left| \frac{\partial f(x)}{\partial x_j} - \frac{(\tilde{x}_j - x_j)}{f(x)} \right|}_{q_j} \\
 &\leq \sum_{j=1}^n \left| \frac{\partial f(x)}{\partial x_j} - x_j \right| \cdot \underbrace{\left| \frac{(\tilde{x}_j - x_j)}{x_j} \right|}_{\delta x_j} + O(||\tilde{x} - x||^2) \\
 &\leq \sum_{j=1}^n q_j \delta x_j (+ \dots)
 \end{aligned}$$

$q_j, \delta x_j$, "Fehlerverstärkungsfaktoren"

Bsp: 1) $y = f(x_1, x_2) := x_1 \cdot x_2$

$$\frac{\partial f}{\partial x_1}(x_1, x_2) = x_2, \quad \frac{\partial f}{\partial x_2}(x_1, x_2) = x_1$$

$$q_1 = \frac{\partial f}{\partial x_1} \frac{x_1}{f(x)} = x_2 \cdot \frac{x_1}{x_1 x_2} = 1 = q_2$$

$$\delta y = \left| \frac{\Delta y}{y} \right| = q_1 \delta x_1 + q_2 \delta x_2 = \delta x_1 + \delta x_2$$

2) $y = f(x_1, x_2) = x_1 + x_2$

$$\frac{\partial f}{\partial x_1}(x_1, x_2) = 1 = \frac{\partial f}{\partial x_2}(x_1, x_2) \Rightarrow q_1 = 1 \cdot \left| \frac{x_1}{x_1 + x_2} \right| + 1 \cdot \left| \frac{x_2}{x_1 + x_2} \right| = q_2$$

$$\delta y = \left| \frac{x_1}{x_1 + x_2} \right| \delta x_1 + \left| \frac{x_2}{x_1 + x_2} \right| \delta x_2, \text{ für } x_1 \approx -x_2 \text{ ist } q_1, q_2 \gg 1$$

"Auslösung"

3) Gesucht: kleinere Nullstelle y_2 von

$$y^2 - 2x_1 y + x_2 = 0$$

$$y_2 = x_1 - \sqrt{x_1^2 - x_2} = f(x_1, x_2)$$

$$\frac{\partial f}{\partial x_1}(x_1, x_2) = 1 - \frac{1}{2} \frac{2x_1}{\sqrt{x_1^2 - x_2}} \Rightarrow q_1 = \left(1 - \frac{x_1}{\sqrt{x_1^2 - x_2}} \right) \frac{|x_1|}{|x_1 - \sqrt{x_1^2 - x_2}|}$$

$$\frac{\partial f}{\partial x_2}(x_1, x_2) = +\frac{1}{2} \frac{2x_2}{\sqrt{x_1^2 - x_2}} \Rightarrow q_2 = \frac{1}{2\sqrt{x_1^2 - x_2}} \cdot \left| \frac{x_2}{x_1 - \sqrt{x_1^2 - x_2}} \right| \left(= \frac{1}{2} - \frac{1}{2} q_1 \right)$$

Sei $M = (1, 0, 5)$ ($m=5$)

$$x_1 = 0,60002 \cdot 10^1, x_2 = 0,1 \cdot 10^{-1}$$

$$q_1 \approx 1, q_2 \approx 1 \quad \text{"gut konditioniert"}$$

Lösung: $y = 0,83336 \cdot 10^{-3}$

Algorithmus 1 zur Berechnung von $y = f(x_1, x_2)$

NumI, VL

$$a := x_1 + x_2;$$

$$b := a - x_2;$$

$$c := \sqrt{b};$$

$$y := x_1 - c; (-x_1 - \sqrt{x_1^2 - x_2^2})$$

$$\Rightarrow y = 0,9 \cdot 10^{-3} \quad \text{S}y \text{ groß!}$$

Algorithmus 2

$$(y = f(x_1, x_2) = \frac{x_2}{x_1 + \sqrt{x_1^2 - x_2^2}})$$

$$a := x_1 + x_2;$$

$$b := a - x_2;$$

$$c := \sqrt{b};$$

$$d := x_1 + c;$$

$$y = \frac{x_2}{d}$$

$$\Rightarrow y = 0,83333 \cdot 10^{-3} \quad \text{S}y \text{ "klein"}$$

Kondition - Eigenschaft des Problems

Stabilität - Eigenschaft des Lösungsalgorithmus

Hauptsache
trennen

Vorwärtsanalyse

x -Daten (ideal)

\tilde{x} -gestörte Daten (real)

Problem gesucht: $y = f(x)$ ideal

berechnet $\tilde{y} = \tilde{f}(\tilde{x})$ real

z.B. absoluter Fehler

$$\|\Delta y\| = \|f(x) - \tilde{f}(\tilde{x})\|$$

$$\|f(x) - \tilde{f}(\tilde{x})\| = \|f(x) - f(\tilde{x})\| + \|f(\tilde{x}) - \tilde{f}(\tilde{x})\| + \|\tilde{f}(\tilde{x}) - \tilde{f}(\tilde{x})\|$$

"Datenfehler" "Verfahrenfehler" "akkumulierter Rundungsfehler"

In der Regel pessimistisch. Gibt nach Rückwärtsanalyse

1.6. Komplexität von Algorithmen

Wichtig sind:

- 1) Aufbau: Abfolge der Operationen sollen eindeutig definiert sein
- 2) Bestimmtheit: Jede mögliche Situation sollte erfasst sein
- 3) Endlichkeit: Der Algo sollte nach „endlicher“ Zeit terminieren

Komplexität: $n \in \mathbb{N}$ Daten, $T: \mathbb{N} \rightarrow \mathbb{N}, n \mapsto T(n)$ Anzahl der Operationen,
Maß für den Arbeitsaufwand (\approx Komplexität)

Anforderungen: 1) ^{höhe} Genauigkeit (Stabilität...)
2) geringe Komplexität

2. Direkte Methoden zur Lsg. lin. GL

Problem: Gegeben: $\begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} = b \in \mathbb{R}^n, A = (a_{ij})_{i,j=1}^n$

Gesucht: $(x_i)_{i=1}^n = x \in \mathbb{R}^n$ mit $Ax = b$

Voraussetzung: A regulär, d.h. A nicht singulär

$$\Leftrightarrow \det A \neq 0 \Leftrightarrow \exists A^{-1} \in \mathbb{R}^{n \times n}$$

2.1.1. Dreiecksmatrix

Def. 2.1.1: $R \in \mathbb{R}^{n \times n}$ heißt „obere Δ -Matrix“, falls $r_{i,j} = 0$ für alle $i > j$.

$$\begin{pmatrix} r_{11} & r_{12} & \cdots & r_{1n} \\ 0 & r_{22} & \cdots & r_{2n} \\ 0 & 0 & \ddots & \cdots \\ 0 & 0 & \cdots & r_{nn} \end{pmatrix} - (\Delta)$$

R ist regulär $\Leftrightarrow \det R = \prod_{i=1}^n r_{i,i} \neq 0 \Leftrightarrow r_{i,i} \neq 0 \forall i = 1, \dots, n$

$$r_{1,1}x_1 + r_{1,2}x_2 + \dots + r_{1,n}x_n = b_1$$

$$\dots \quad r_{2,2}x_2 + \dots = b_2$$

$$r_{n-1,n-1}x_{n-1} + r_{n,n}x_n = b_{n-1}$$

$$r_{n,n}x_n = b_n \Rightarrow x_n = \frac{b_n}{r_{n,n}}$$

$$\Rightarrow x_{n-1} = \frac{1}{r_{n-1,n-1}}(b_{n-1} - r_{n-1,n}x_n)$$

$$\Rightarrow \text{(vollst. Induktion)} x_j = \frac{1}{r_{jj}} \left(b_j - \sum_{k=j+1}^n r_{j,k}x_k \right) \quad j = n-1, \dots, 1, \quad x_n = \frac{b_n}{r_{n,n}} \text{ „Rückwärts-Substitution“}$$

$$(j \text{ Zeile } r_{j,j}x_j + r_{j,j+1}x_{j+1} + \dots + r_{j,n}x_n = b_j)$$

Komplexität numbers of nonzeros

$$R \in \mathbb{R}^{n \times n} \quad \# R = nnz_R = O(n^2)$$

$$b \in \mathbb{R}^n \quad \# b = nnzb = O(n)$$

Für $j = n-1 \rightarrow 1$
in rückwärts raus

in Zeile j : $(n-j)$ Multiplikationen (+1 Division) $(n-j+1)$ Add.

$$= T(n) = \sum_{j=n-1}^1 (n-j) + n \left[\sum_{j=n-1}^{n-1} (n-j) = \sum_{l=1}^{n-1} l \right]$$

$$= n + \sum_{l=1}^{n-1} l = n + \frac{n(n-1)}{2} = O(n^2)$$

Substitution $l = n-j$
Für $j = n-1 \Rightarrow l = n-n+1 = 1$
 $j=1 = l = n-1$

2.1.2 Gauß-Elimination

Ziel: Transformation auf obere Δ -Gestalt mit Elementarumformungen

1) Vertauschen von Zeilen

2) Subtraktion des Vielfachen einer Zeile (k -te Zeile) von
einer anderen Zeile (i -te Zeile)

Für $k = 1, \dots, n-1$:

1) Vertausche Zeile 1 mit Zeile i . [Pivotisierung]

$$a_{k,j} \leftarrow \arg \max \{ |a_{i,k}| : i \geq k \}$$

→ Vertausche

$$2) \text{ für } i=k+1, \dots, n \quad a_{i,k} \leftarrow \frac{a_{i,k}}{a_{k,k}}$$

neue
Zeile
entsteht
durch
Multiplizieren
der Spalte
mit
zugehöriger
Zeile

$j = 1, \dots, n$
(oder $j = k+1, \dots, n$)

$$\text{Für } j=k+1, \dots, n \quad a_{i,j} \leftarrow a_{i,j} - l_{i,k} a_{k,j}$$

Beispiel

$$\begin{array}{c}
 \begin{array}{ccccc}
 A & & b & & \text{Pivot} \\
 \left[\begin{array}{ccc|c} 0 & 1 & -3 & 3 \\ 1 & 1 & 3 & -4 \\ 1 & -1 & 3 & 5 \end{array} \right] & \xrightarrow{\quad} & \left[\begin{array}{ccc|c} 1 & 1 & 3 & -4 \\ 0 & 1 & -3 & 3 \\ 1 & -1 & 3 & 5 \end{array} \right] & \xrightarrow{\text{III} \leftrightarrow \text{I}} & \left[\begin{array}{ccc|c} 1 & 1 & 3 & -4 \\ 0 & 1 & -3 & 3 \\ 0 & -2 & 0 & 9 \end{array} \right] \\
 \text{ ohne Pivot} & \xrightarrow{\quad} & \left[\begin{array}{ccc|c} 1 & 1 & 3 & -4 \\ 0 & 1 & -3 & 3 \\ 0 & -2 & 0 & 9 \end{array} \right] & \xrightarrow{\Delta \text{ Zeilen} \Rightarrow \text{Backward-Sub}} & -6x_3 = 15 \Rightarrow x_3 = -\frac{5}{2} \\
 & & \xrightarrow{\text{III} - (-2) \cdot \text{II}} & & x_2 - 3 \cdot \left(-\frac{5}{2}\right) = 3 \\
 & & \xrightarrow{l_{3,2} = -2} & & \Rightarrow x_2 = -\frac{9}{2} \\
 & & & & x_1 - \frac{9}{2} + (3) \cdot \left(-\frac{5}{2}\right) = -4 \\
 & & & & \Rightarrow x_1 = 8
 \end{array} \\
 \rightarrow R = \left(\begin{array}{ccc} 1 & 1 & 3 \\ 0 & 1 & -3 \\ 0 & 0 & -6 \end{array} \right), L = \underbrace{\left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ l_{ij} & 1 \end{array} \right)}_{P} \rightsquigarrow L = \left(\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & -2 & 0 \end{array} \right) \\
 L \cdot R = \left(\begin{array}{ccc} 1 & 1 & 3 \\ 0 & 1 & -3 \\ 1 & -1 & 3 \end{array} \right), L \cdot R = \underbrace{\left(\begin{array}{ccc} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{array} \right)}_P \wedge PLR = A
 \end{array}$$

2.1.3 LR-Zerlegung

Erweiterte Koeffizientenmatrix
 $(A, b) \in \mathbb{R}^{n \times (n+r)}$

Permutationsmatrizen

$$P = P_{k,r} = \left(\begin{array}{cccccc} 1 & & & & & & \\ & 1 & & & & & \\ & & 0 & 1 & & & \\ & & & 1 & & & \\ & & & & 0 & 1 & \\ & & & & & 1 & \\ & & & & & & 1 \end{array} \right) \quad \begin{array}{l} \leftarrow k\text{-te Zeile} \\ \leftarrow r\text{-te Zeile} \\ \leftarrow r\text{-te Spalte} \end{array}$$

$\hat{A} = PA$ Multiplikation von Links mit der Matrix P vertauscht die Zeilen k mit r.

- Proposition
- 1) $P = P_{k,k}$ ($k \neq r$) ist regulär
 - 2) $P^2 = I$, d.h. $P^{-1} = P$,
 - 3) $P^T = P$

Satz 2.1.3 (L-R-Zerlegung)

a) Sei $A \in \mathbb{R}^{n \times n}$, $\det A \neq 0$, dann $\exists P = P_{k_1, k_1} \dots P_{k_d, k_d}$ ($d \leq n$) Permutation und $L = \begin{pmatrix} 1 & 0 \\ l_{ij} & 1 \end{pmatrix}$ und R obere Δ -Matrix mit
(konstruktiv mittel Gauß-Elimination)

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

$$l^{(k)} = (0, \underset{k-\text{te}}{\overset{1}{\dots}}, 0, l_{k+1, k}, \dots, l_{n, k})^T \in \mathbb{R}^n, e_k = (0, \underset{k-\text{te Zeile}}{\overset{1}{\dots}}, 0, 0)^T \in \mathbb{R}^n$$

$$G_k = \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 0 & -l_{k+1, k} \\ -l_{k+2, k} & 1 \\ \vdots & \vdots \\ -l_{n, k} & 1 \end{pmatrix} = I - l^{(k)} e_k^+$$

Spaltenvektor

kanonische Einheitsfaktor

$\Rightarrow \det G_k = 1$, d.h. G_k ist regulär und es gilt:

$$G_k^{-1} = I + l^{(k)} e_k^T = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

$$\langle e_k, l^k \rangle = 0$$

$$G_k^{-1} G_k = (I + l^{(k)} e_k^T)(I - l^k e_k^+) = I + l^k e_k^+ - l^k e_k^+ - l^k e_k^+ l^k e_k^T$$

Beobachtung $(A, b)^{(k+1)} = G_k \cdot P_{k, r} (A, b)^k$ $r \geq k$
nach k -tem Eliminationsschritt

b) $\tilde{A} = PA = LR$, Diese Zerlegung ist eindeutig!

$$\mathcal{L} = \{L \mid \begin{pmatrix} 1 & 0 \\ l_{ij} & 1 \end{pmatrix} \in \mathbb{R}^{n \times n} : l_{ij} \in \mathbb{R}\}, \mathcal{R} \in \begin{pmatrix} \mathbb{R}^m & \mathbb{R}^{m \times n} \\ 0 & \mathbb{R}^{m \times n} \end{pmatrix} | r_{ij} \in \mathbb{R}\}$$

Beweis: $L_1, L_2 \in \mathcal{L} \Rightarrow L_1 \cdot L_2 \in \mathcal{L}$, $R_1, R_2 \in \mathcal{R} \Rightarrow R_1 \cdot R_2 \in \mathcal{R}$

a) $l^{(k)} = (0, \underset{k-\text{tes Zeil}}{\overset{1}{\dots}}, 0, l_{k+1, k}, \dots, l_{n, k})^T, e_k = (0, \underset{k-\text{te}}{\overset{1}{\dots}}, 0, -z_{k+1, k}, \dots, -z_{n, k})^T \in \mathbb{R}^n$

$$G_k = \begin{pmatrix} 1 & & & & 0 \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ 0 & -l_{k+1, k} & & & 0 \\ & & 1 & & \\ & & & -l_{n, k} & 0 \\ & & & & 1 \end{pmatrix} = I - l^{(k)} \cdot e_k^T \in \mathcal{L}$$

$$\Rightarrow G_k^{-1} = I + l^{(k)} e_k^+$$

$P_k = P_{r, r}, r > k$ mit k -tem Eliminationsschritt

$$G_k P_k A^{(k)} = A^{(k+1)}$$

PinO + Schritt

Es gilt $P_k \cdot e_j = e_j$, $k > j$

$$\underbrace{P_k G_j P_k}_0 = P_k (I - z_j e_j^T) P_k = (P_k P_k) - (P_k z_j)(P_k e_j)^T = I - P_k z_j e_j^T = \tilde{G}_j$$

nach dem k -ten Schritt

$$R \circ R := A^n = G_{n-1} P_{n-1} A^{(n-1)} = G_{n-1} B_{n-1} G_{n-2} P_{n-2} \cdots G_1 P_1 A \quad \textcircled{*}$$

$$(A, b) \overset{\text{vergleich}}{\sim} (A, b)_j$$

$$\textcircled{*} = (G_{n-1} P_{n-1} G_{n-2}) \overset{I}{\underset{P_{n-1} P_{n-1}}{\sim}} P_{n-2}$$

$$= G_{n-1} \tilde{G}_{n-2} P_{n-1} P_{n-2} G_{n-3} P_{n-2} P_{n-3} B_{n-3}$$

$$= G_{n-1} \tilde{G}_{n-2} \tilde{G}_{n-3} P_{n-1} P_{n-2} P_{n-3} G_{n-4} P_{n-3} P_{n-4} \cdots$$

$$= G_{n-1} \tilde{G}_{n-2} \tilde{G}_{n-3} \cdots P_{n-1} - P_n \circ A \quad \text{mit } \tilde{G}_j = I - P_{n+1} \cdot P_{j+1} e_j^T \\ \text{und } \tilde{G} \in \mathcal{L} \quad = I - z_j e_j^T$$

$$\Rightarrow PA = \tilde{G}^{-1} \cdot R = \underbrace{\tilde{G}_{n-2}^{-1} \cdot \tilde{G}_{n-1}^{-1}}_{\in \mathcal{L}} R = LR$$

b) Eindeutigkeit: Seien $L_1, L_2 \in \mathcal{L}, R_1, R_2 \in \mathbb{R}$

$$\text{mit } PA = L_1 R_1 = L_2 R_2 \quad | L_2^{-1}(\text{)})$$

$$L_2^{-1} L_1 R_1 = L_2^{-1} L_2 R_2$$

$$L_2^{-1} L_1 R_1 = I R_2 \quad | (\text{)}) R_1^{-1}$$

$$L_2^{-1} L_1 R_1 R_1^{-1} = I R_2 R_1^{-1}$$

$$L \in \mathcal{X} \Rightarrow L^{-1} \in \mathcal{L}$$

$$\mathcal{L} \ni L_2^{-1} L_1 = R_2 R_1^{-1} \in \mathcal{R}$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \Rightarrow \boxed{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}} \xrightarrow{\quad} I$$

$$\Rightarrow I = I$$

$$\Rightarrow L_2 = L_1 \quad \text{und } R_1 = R_2$$

Bem. Sei $A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$ ($= P \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$) Dann ex. keine LR-Zerlegung
 (ohne Permutation)
 (ohne Pivotisierung)

Satz 2.1.4 Falls alle Hauptunterdeterminanten

$$0 \neq \det A_k = \begin{pmatrix} a_{11} & \dots & a_{1k} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nk} \end{pmatrix} \in \mathbb{R}^{k \times k} \quad \forall k=1, \dots, n, \text{ dann existieren}$$

$L \in \mathbb{Z}, R \in \mathbb{R}$ mit $A = LR$

Beweis Übung ** (unnötig)

2.1.5 Komplexität des Gauß-AG.

Satz 2.1.5. Der Eliminationsschritt benötigt für $n=2$

$$\frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6} \quad \begin{array}{l} \text{Multiplikation D.h. T(n)} \\ \text{Divisionen} \end{array}$$

Beweis: Im k -ten Schritt benötigt man $(n-k)$ Divisionen

und $(n-k+1)(n-k)$ Multiplikation

$$\Rightarrow \text{Anz. der } (\times) \text{ (Schritt?)} \sum_{k=1}^{n-1} (n-k)(n-k+1) = \sum_{\ell=1}^{n-1} \ell(\ell+2) \quad \textcircled{D}$$

$\ell = n-k, k=1, \ell=n-1, k=n-1, \ell=1$

$$\textcircled{D} = \sum_{\ell=1}^{n-1} \ell^2 + 2 \sum_{\ell=1}^{n-1} \ell = \frac{n(n-1)(2n-1)}{6} + 2 \frac{(n-1)n}{2} = \frac{n^3}{3} + \frac{n^2}{2} - \frac{5n}{6} = O(n^3)$$

Aufwand für $A \cdot B \rightarrow A \cdot B \in \mathbb{R}^{n \times n} \rightarrow O(n^3)$

aber Strassen (1969) - $O(n^{\log_2(7)}) = O(n^{2.8})$

2.1.5 Vorteile der LR-Zerlegung

1. Berechnung von Determinanten

$$PA = LR, \det(PA) = \det(LR)$$

$$\text{d} \quad \det(P) \det(A) = \det(L) \det(R) = 1 \cdot \prod_{i=1}^n r_{i,i}$$

$$\det(P) = (-1)^{\text{# Vertauschungen}} = \sum_{\sigma \in S} \text{sgn } \sigma$$

$$\Rightarrow \det(A) = \det(r_{11} \dots r_{n,n}) \rightarrow \text{Aufwand } T(n) = O(n^3) \quad (\text{LR-Zerlegung})$$

In Gegensatz zu Entwicklungssatz $T(n) = O(n!)$

2. Mehrere Rechte Seiten $Ax_i = b_1, \dots, b_p$

$$\underbrace{LRx}_{=: z} = PAx = Pb \quad \text{Alg: 1) } C = PB \quad \begin{array}{l} \text{Permutation} \\ 2) \text{ L\"ose } Lz = C; \text{ Vorw\"arts-substitution } O(n^2) \end{array}$$

$$3) Rx = z \quad \text{R\"uckw\"arts-substitution } O(n^2)$$

3) Berechnung von $A^{-1} = (x_1, \dots, x_n)$

$$Ax_i = e_i, i=1, \dots, n \quad O(n^2) = O(n^3)$$

4) Nachiteration

Sei $x^* \in \mathbb{R}^n$ eine Nährungslösung von $Ax=b$

Residuum $r = b - Ax^* = A(x - x^*) \neq 0$, ^{setze} $x = x^* + s$ ^{nicht bekannt} $\Rightarrow s = x - x^*$

mit $As = r = A(x - x^*)$

d.h. Löse $As = r$

2.1.6 Bemerkung zu Pivot-Strategien

- Üblich Spaltenpivotisierung: maximales $|a_{j,k}|$ für $j = k, \dots, n$ im k -ten Schritt
- Vollständige oder Restmatrix-Pivotisierung:
Suche maximalen $|a_{j,k}|$ für $j = k, \dots, n$
für $k = k, \dots, n$
 - ↳ durchlaufen nicht nur die restliche Spalte, sondern gesamte restliche Zeilen & Spaltenvertauschung
 - ↳ Optimal, aber zuteuер!
- Äquilibrierung:

Def: $A = (a_{i,j}) \in \mathbb{R}^{n \times n}$ heißt zeilenweise äquilibriert, wenn

$$\sum_{j=1}^n |a_{i,j}| = 1 (\approx 1) \quad \forall i = 1, \dots, n$$

Seien $d_i = \sum_{j=1}^n |a_{i,j}|$, $i = 1, \dots, n$ \Leftarrow Summe der Einträge der Zeile i : d_i
und sei $D = \text{diag}(d_i) = \begin{pmatrix} d_1 & & \\ & \ddots & \\ & & d_n \end{pmatrix}$

Aus $Ax = b \Leftrightarrow \bar{A}x = \bar{b} = D^{-1}Ax = D^{-1}b$

d.h. $\bar{b} = D^{-1}b$, $\bar{A} = D^{-1}A = \text{diag}(d_i^{-1})A$ (leicht zu berechnen)

$\Rightarrow \bar{A}$ ist äquilibriert

→ Strategie: Äquilibriere A und spaltenweise Pivotisierung

Bsp: $Ax = \begin{pmatrix} 0,31 \cdot 10^{-3} & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -3 \\ -7 \end{pmatrix}$

ohne Pivotisierung: $b_{2,1} = \frac{1}{0,31 \cdot 10^{-3}}$

erster Eliminationsschritt $\Rightarrow \left[\begin{array}{cc|c} 0,31 \cdot 10^{-3} & 1 & -3 \\ 1 & -3225 & 9670 \end{array} \right] \Rightarrow x_1^* \approx -6,452 \rightsquigarrow \text{relativ schlecht}$
 $x_2^* \approx -2,998 \rightsquigarrow \text{relativ gut}$

Exakt: $x_1 = -4,00124$, $x_2 = -2,998759$ (ohne Runden)

Ursache: Auslöschung im Rückwärtseinsetzen: $x_2 \approx 3$

und $x_1^* = \frac{1}{0,00031} \underbrace{(-3 - x_2^*)}_{\approx 0}$

Mit Spaltenpivotisierung

nach 1. Eliminationsschritt:

$$\left[\begin{array}{cc|c} 1 & 1 & -7 \\ 0 & 0,9997 & -2,998 \end{array} \right] \quad \left. \begin{array}{l} \bar{x}_1 \approx -4,001 \\ \bar{x}_2 \approx -2,999 \end{array} \right\} \text{gut}$$

mit Äquilibrierung

Multiplizierte 1. Zeile mit 10^4

→ ähnliches Resultat wie mit Pivotisierung

2.1.7 Spezielle Matrizen - Bandmatrizen

Def $A = (a_{i,j}) \in \mathbb{R}^{n \times n}$ heißt (m, k) -Bandmatrix, falls

$a_{i,j} = 0$ für $|i-j| > m$ oder $|j-i| > k$

also $\begin{pmatrix} * & & & \\ \cancel{0 \dots} & * & & \\ & & \ddots & \\ & & & 0 \end{pmatrix}$ oder $\begin{pmatrix} * & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 0 \end{pmatrix}$

Triadiagonalmatrix: $(1, 1), m=k=1$

Hessenbergmatrix: $(1, n-1), m=1, k=n-1$

tridiagonalmatrix: $(b, a, c) = \begin{pmatrix} a_1 & c_1 & & 0 \\ b_2 & a_2 & & \\ & & \ddots & \\ 0 & & & b_n a_n \end{pmatrix}$

$$\begin{pmatrix} a_1 & c_1 & 0 \\ b_2 & a_2 & c_2 \\ 0 & b_3 & a_3 & c_3 \\ & & \ddots & \ddots & 0 \\ & & & & b_n & a_n \end{pmatrix} = \begin{pmatrix} a_1 & & & 0 \\ b_2 & a_2 & & \\ 0 & b_3 & a_3 & \\ & & \ddots & \ddots & 0 \\ & & & & b_n & a_n \end{pmatrix} \begin{pmatrix} 1 & & & 0 \\ & \ddots & & \\ 0 & & 1 & 0 \\ & & & \ddots & 0 \\ & & & & 1 & 0 \end{pmatrix}$$

Koeffizienten abgleich

$$\Rightarrow a_1 = x_1, c_1 = x_1 y_1 \Rightarrow y_1 = \frac{c_1}{a_1} = \frac{c_1}{a_1}$$

$$\Rightarrow a_i = b_i y_{i-1} + x_i \cdot 1 \Rightarrow x_i = a_i - b_i y_{i-1}$$

$$\text{und } c_i = x_i y_i \Rightarrow y_i = \frac{c_i}{x_i} \quad \text{Rekursiv}$$

$$T = \tilde{L} \cdot \tilde{R}$$

T-Triadiagonalmatrix heißt irreduzibel diagonaldominant

falls 1) $0 < |c_i| < |a_{i+1}|$

2) $|a_i| \geq |b_i| + |c_i|$ (bei „Strikt“ diagonaldominant)

3) $|a_n| > |b_n| > 0$ (minreichend)

Satz 2.1.9: Sei $A \in \mathbb{R}^{n \times n}$ ($\mathbb{C}^{n \times n}$) irreduzibel diagonal dominant NumI, VL

Dann ist $\det A \neq 0$ und es gelten

$$1) |x_i| < 1 \quad \forall i = 1, \dots, n-1$$

$$2) |a_{ii}| + |b_{ii}| > |d_{ii}| > |a_{ii}| - |b_{ii}| > 0$$

Beweis 1) Induktionsanfang: $i=1$

$$|x_1| = \left| \frac{c_1}{a_1} \right| < 1$$

IS $j-1 \Rightarrow j$:

$$|x_j| = \left| \frac{c_j}{\alpha_j} \right| = \frac{|c_j|}{|a_{j-1} - b_{j-1} x_{j-1}|} \leq \frac{|c_j|}{|a_{j-1} - x_{j-1} b_{j-1}|} \leq \frac{|c_j|}{|a_{j-1} - b_{j-1}|} \leq \frac{|c_j|}{|c_j|} = 1$$

$$2) |a_{ii}| + |b_{ii}| > |a_{ii} - x_{i-1} b_{i-1}| = |\alpha_i| > |a_{ii}| - |b_{ii}| \geq |c_{ii}| > 0$$

$$\Rightarrow \det T = \det \tilde{L} \cdot \det \tilde{R} = \alpha_1 \alpha_2 \cdots \alpha_n > 0, \text{ da } \alpha_i > 0 \quad \forall i = 1, \dots, n \quad \blacksquare$$

2.2. Symmetrische Matrizen & Cholesky-Zerlegung

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

A symmetrisch heißt positiv definit, falls

$$\forall x \in \mathbb{R}^n \text{ mit } x \neq 0 \quad x^T A x = \langle x, Ax \rangle > 0$$

$$\text{alternativ: } \exists \gamma > 0 \quad \forall x \in \mathbb{R}^n \quad x^T A x \geq \gamma \|x\|^2 = \gamma x^T x$$

kurz SPD (symm. pos. definit)

Theorem: A SPD \Leftrightarrow alle Eigenwerte $\lambda_1, \dots, \lambda_n$ von A sind > 0 , d.h. $\lambda_i > 0 \quad \forall i = 1, \dots, n$

\Leftrightarrow alle Hauptunterdeterminanten von A , also

$\det A_k > 0 \quad \forall k = 1, \dots, n$, sind positiv.

Satz 2.2.3 A ist SPD \Rightarrow es ex. genau eine linke ^{untere} Δ -Matrix $\tilde{L} = (\tilde{l}_{ij})$ ($i, j = 0, j > i$) mit $A = \tilde{L} \tilde{L}^T$ (\tilde{L}^T ist obere Δ -Matrix)

Beweis: \Rightarrow da $\det A_k > 0$ ex. L, R mit $A = L \cdot R$ und $\det A_k = \det L_k \det R_k$

$$\Rightarrow r_{k,k} > 0 \quad \forall k = 1, \dots, n \quad \text{für } \tilde{l}_{i,i} = \sqrt{a_{i,i}} \Rightarrow \tilde{r}_{ii} = \frac{r_{ii}}{\sqrt{r_{ii}}} = \prod_{i=1}^k r_{ii} > 0$$

$$\tilde{l}_{i,j} = \tilde{l}_{i,i}^{-1} \cdot r_j, \quad r_{ij} = \tilde{r}_{i,j} \cdot r_{i,j}$$

$$\Rightarrow A^T = A = \tilde{L} \cdot \tilde{R} = A^T = (\tilde{L} \tilde{R}^T)^T = \tilde{R}^T \tilde{L}^T$$

Aus der Eindeutigkeit der LR -Zerlegung $\Rightarrow \tilde{L} = \tilde{R}^T = A = \tilde{L} \tilde{L}^T$

mit der Eigenschaft $\det L > 0$

$$\Leftrightarrow \text{Sei } \tilde{\Sigma} \text{ regulär: } \Rightarrow \forall x \neq 0: x^T A x = x^T \tilde{\Sigma} \tilde{\Sigma}^T x \\ = (\tilde{\Sigma}^T x)^T (\tilde{\Sigma}^T x) \\ = \|\tilde{\Sigma}^T x\|^2 \geq 0 \quad \forall x \neq 0$$

2.2.2. Die Crout-Cholesky-Zerlegung

gesucht: x mit $AX=b$
symm. $A = A^T \in \mathbb{R}^n$ ($A^* = A^H = A$ hermitesch))

A positiv definit (SPD) $\Rightarrow \exists L \in \mathbb{R}^{n \times n}$ (Cholesky-Zerlegung) $(l_{ii})^2 = d_i$

Sei $A = L \cdot R = L \cdot \tilde{D} \cdot L^T$ die Crout-Cholesky-Zerlegung

$$l_{ii} = 1 \quad \forall i=1, \dots, n$$

$D = \text{diag}(d_i)$

Diagonalmatrix

$$\text{d.h. } (i \geq j) \quad a_{ij} = \sum_{k=1}^n l_{ik} \tilde{l}_{jk} = \sum_{k=1}^{j-1} l_{ik} d_k l_{jk} + l_{ij} d_j l_{jj} = \sum_{k=1}^{j-1} l_{ik} d_k l_{jk} + l_{ij} d_j l_{jj}$$

$$\tilde{l}_{ij} = d_j l_{ij} \Rightarrow l_{ij} = \frac{\tilde{l}_{ij}}{d_j}$$

$$\sum_{i=1}^j a_{ii} = l_{11} d_1 l_{11} = d_1$$

$$\sum_{i=2}^j a_{2,i} = l_{2,1} d_1 \Rightarrow \tilde{l}_{2,1} = a_{2,1} / l_{2,1} = \frac{a_{2,1}}{d_1}$$

$$\sum_{i=2}^j a_{i,1} = l_{i,1} d_1 \Rightarrow l_{i,1} = \frac{a_{i,1}}{d_1}, \quad i=2, \dots, n$$

$$\sum_{i=2}^j a_{2,2} = d_2 + l_{2,1} d_1 l_{2,1}$$

$$\Rightarrow d_2 = a_{2,2} - l_{2,1} d_1 l_{2,1}$$

$$\sum_{i=2}^j a_{i,2} = l_{i,2} d_2 \quad l_{i,2}$$

$$(i=j) \quad \text{Für } j=1, \dots, n, \quad d_j = a_{jj} - \sum_{k=1}^{j-1} l_{jk} \tilde{l}_{jk} = a_{jj} - \sum_{k=1}^{j-1} l_{jk}^2 d_k$$

Für $i=j+1, \dots, n$:

$$\begin{aligned} \tilde{l}_{ij} &= a_{ij} - \sum_{k=1}^{j-1} l_{ik} \tilde{l}_{jk} \\ l_{ij} &= \frac{\tilde{l}_{ij}}{d_j} \end{aligned} \quad \left. \begin{array}{l} \text{Alg(77)} \\ \text{Algorithmus 77} \end{array} \right\}$$

$$\text{Beispiel: } A = \begin{pmatrix} 2 & 6 & -2 \\ 6 & 21 & 0 \\ -2 & 0 & 16 \end{pmatrix} \Rightarrow d_1 = a_{11} = 2$$

$$\Rightarrow j=1, i=2 \quad \tilde{l}_{2,1} = a_{2,1} = 6, \quad l_{2,1} = \frac{\tilde{l}_{2,1}}{d_1} = \frac{6}{2} = 3$$

$$j=1, i=3 \quad \tilde{l}_{3,1} = a_{3,1} = -2, \quad l_{3,1} = \frac{\tilde{l}_{3,1}}{d_1} = \frac{-2}{2} = -1$$

$$j=2, \quad d_2 = a_{22} - l_{2,1} \tilde{l}_{2,1} = 21 - 18 \stackrel{d_1=2}{=} 3 = d_2$$

$$j=2, i=3, \quad \tilde{l}_{3,2} = a_{3,2} - l_{3,1} \tilde{l}_{2,1} = 0 - (-1) \cdot 6 = 6$$

$$l_{3,2} = \frac{\tilde{l}_{3,2}}{d_2} = \frac{6}{3} = 2$$

$$j=3, \quad d_3 = a_{33} - l_{3,1} \tilde{l}_{3,1} - l_{3,2} \tilde{l}_{3,2} = 16 - (-1)(-2) - 2 \cdot 6 = 2$$

$$L = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ -1 & 2 & 1 \end{pmatrix} \quad D = \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

$$A = \begin{pmatrix} 1 & 0 & 0 \\ 3 & 1 & 0 \\ -1 & 2 & 1 \end{pmatrix} \begin{pmatrix} 2 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 3 & -1 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix} \Rightarrow \begin{matrix} d_i > 0 \\ i=1, \dots, n \end{matrix} \quad A = \text{SPD}$$

Die klassische Cholesky-Zerlegung

$$A = L L^T \text{ ergibt sich aus } A = L D L^T \text{ mit } r_i = \sqrt{d_i} \quad (d_i > 0)$$

$$\text{d.h. } l_{i,j}^{\text{Cnd}} = l_{i,j} \sqrt{d_i}$$

berlinet
in Crout-Cholesky

Vorteile von Crout-Cholesky gegenüber Cholesky:

- 1) Cholesky erfordert SPD
- 2) Cholesky erfordert (benötigt) Wurzelfunktion

Arbeitsaufwand zur Cholesky-Zerlegung:

$$\text{Anzahl der arithmetischen Operationen } T(n) = \frac{1}{6} n^3 + \frac{1}{2} n^2 - \frac{2}{3} n$$

$$T(n) \Big|_{\text{Cholesky}} \sim \frac{1}{2} T(n) \Big|_{\text{Gauß}}$$

$$L D L^T x = b$$

$$\text{Löse } Lz = b, \text{ Vorwärts-Substitution} \quad z \quad Dy = z$$

$$y = D^{-1}z = \text{diag}(d_i^{-1})z$$

$$x = L^T y, \text{ Rück-Sub}$$

Bemerkung: Es wurde keine Pivotisierung durchgeführt

(Geringer Stabilität als pivot. Gauß)

Orthogonale Projektion (oft =) beste Approximation, QR-Zerlegung Num I, VL

Def (Skalarprodukt) V -reeller Vektorraum,

$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$ heißt Skalarprodukt, falls:

$$\langle x, y \rangle = \langle y, x \rangle \quad \forall x, y \in V$$

$$\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle \quad \forall x, y, z \in V, \forall \alpha, \beta \in \mathbb{R}$$

$$\langle x, x \rangle \geq 0, \langle x, x \rangle = 0 \Leftrightarrow x = 0 \quad \forall x \in V$$

Satz (Cauchy-Schwarz-Ungl.) $|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$

Def (Beste Approximation): Sei V normierter Vektorraum über \mathbb{R} und M ein Untervektorraum. y^* heißt Element bester Approximation an $x \in V$, falls $y^* \in M$ und $\forall y \in M$:

$$\underbrace{\|x - y\|}_{\text{Abstand } x \text{ zu } y} \geq \underbrace{\|x - y^*\|}_{\text{Abstand } x \text{ zu } y^*}$$

Def (orthogonale Projektion): Sei V ein reeller VR mit Skalarprodukt, M ein UVR, y^* heißt orthogonale Projektion von x auf M , falls $y^* \in M$ und $\forall y \in M: \langle y, x - y^* \rangle = 0$
verbindungsline "zw. x und seiner orth. Proj.

Satz Orthogonale Projektion ist eindeutig

Beweis: Seien y^* und \tilde{y} beide orth. Proj. von x auf M

Es gilt dann $\langle y^* - x, y \rangle = 0 \quad \forall y \in M$

$$\langle \tilde{y} - x, y \rangle = 0 \quad \forall y \in M$$

Daraus folgt: $\langle y^* - \tilde{y}, y \rangle = \langle y^* - x, y \rangle - \langle \tilde{y} - x, y \rangle = 0 - 0 = 0 \quad \forall y \in M$

Wähle nun $y = y^* - \tilde{y} \in M$, dann folgt

$$\langle y^* - \tilde{y}, y \rangle = 0 \quad \text{Nach Definitheit des SP folgt } y^* - \tilde{y} = 0, \text{ also } y^* - \tilde{y} = 0$$

Satz Sei $y^* \in M$ ein El. bester Approximation an $x \in V$. Dann ist y^* auch die orthogonale Projektion von x auf M .

Beweis: Sei $y^* \in M$ ein beliebiges Element bester Approximation.

Dann liegt $y^* + \alpha v$ auch in M , falls $v \in M$ und die Funktion

$$F: \mathbb{R} \rightarrow \mathbb{R}: \alpha \mapsto \underbrace{\|x - (y^* + \alpha v)\|}_M^2 \quad \text{nur bei } \alpha = 0 \text{ ein Minimum haben}$$

$$F(\alpha) = \|(x - y^*) - \alpha v\|^2 = \|x - y^*\|^2 + \alpha^2 \|v\|^2 - 2\alpha \langle x - y^*, v \rangle$$

$$\Rightarrow \text{Bei } \alpha = 0 \text{ muss } F' \text{ verschwinden}$$

$$\Rightarrow 0 = F'(0) = -2 \langle x - y^*, v \rangle$$

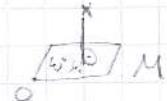
Da v beliebig gewählt folgt die Voraussetzung der orthg. Projektion
Korollar Das Element bester Approximation ist eindeutig

Beweis Übung

Satz Sei y^* die orthogonale Projektion von $x \in V$ auf M .

Dann ist y^* auch beste Approx.

Beweis Sei w irgendein Element aus M



$$\|x - w\|^2 = \|x - y^* + (y^* - w)\|^2$$

$$= \langle (x - y^*) + (y^* - w), (x - y^*) + (y^* - w) \rangle$$

$$\Rightarrow \langle x - y^*, x - y^* \rangle + \langle y^* - w, y^* - w \rangle + 2 \langle x - y^*, y^* - w \rangle$$

$$= \|x - y^*\|^2 + \underbrace{\|y^* - w\|^2}_{\geq 0} + \underbrace{2 \langle x - y^*, y^* - w \rangle}_{\in M} \geq \|x - y^*\|^2$$

VI
O

$$= 0, \text{ weil } (y^* - w) \perp (x - y^*) \\ \text{bzw. } y^* - w \perp \text{Projekt.}$$

Satz Ist $\dim(M) < \infty$ und a_1, \dots, a_m die lin. unabh. Spalten von A und $M = \text{span}\{a_1, \dots, a_m\}$ Dann ist $y^* = Ax$ genau dann das Element bester Approx. an x wenn

$$A^T A \alpha = A^T x \Leftrightarrow \underbrace{\|A\alpha - x\|}_{=y^*} \text{ minimal}$$

(Normalengleichung)

Beweis In einer UE würde gezeigt, dass $y^* = A(A^T A)^{-1} A^T x$

Bemerkung: $A^T A$ sym pos definit $\Rightarrow A^T A$ invertierbar

$\Rightarrow \alpha$ eindeutig bestimmt

Bemerkung: $A^T A$ zu berechnen ^{nur wenn} nicht immer sinnvoll!

Beispiel $A = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ \epsilon & \epsilon & \epsilon & \epsilon & \epsilon \end{pmatrix}$ hat Rang 5

$$A^T A = \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1+\epsilon^2 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1+\epsilon^2 \end{pmatrix}$$

Was wenn A nicht vollen Rang holt?

Wenn $\epsilon < \sqrt{\epsilon^*}$ ^{max rel. Fehler bei Rang}
 folgt $f(x^T A) = \begin{pmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{pmatrix}$
 hat Rang 1

Wollen α finden, so dass $\|A\alpha - y^*\|$ minimal

(gesucht ist Element bester Approximation $y^* = A\alpha$ im von den Spalten von A aufgespannten Untervektorraum)

$A\alpha$ eindeutig, aber α nicht eindeutig

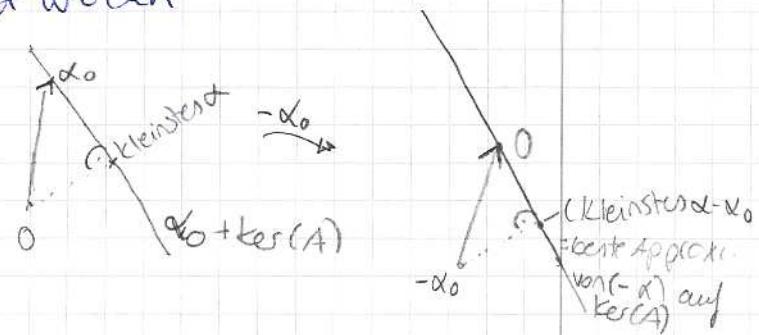
Sei α_0 sodass $A\alpha_0 = y^*$ E.b. $A \Rightarrow \alpha_0 + v$ auch gut,

wenn $v \in \text{Kern}(A)$ $A(\alpha_0 + v) = A\alpha_0 = y^*$

α kann aus affinem UVR $\alpha_0 + \text{ker}(A)$ gewählt werden.

Das kleinste $\alpha \in \alpha_0 + \text{Ker}(A)$ kann wieder durch

Beste Approximation bestimmt werden



QR-Zerlegung

Def: Die QR-Zerlegung einer $n \times m$ Matrix A ist eine Zerlegung von A in ein Produkt $A = Q \cdot R$, wobei Q eine orthogonale $n \times n$ -Matrix und R eine $n \times m$ obere Δ -Matrix ist

$$\begin{matrix} n \\ \vdots \\ m \end{matrix} \begin{matrix} A \end{matrix} = \begin{matrix} n \\ \vdots \\ n \end{matrix} \begin{matrix} Q \end{matrix} \begin{matrix} n \\ \vdots \\ m \\ 0 \end{matrix} \begin{matrix} R \end{matrix}$$

- Eine Matrix Q heißt orthogonal, wenn ihre Spalten jeweils Norm/Länge Eins haben und paarweise orthogonal/Senkrecht aufeinander stehen, also wenn $Q^T Q = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$. Es folgt, dass $Q^{-1} = Q^T$ und damit $Q^T Q = Q^{-1} Q = Q Q^{-1} = Q Q^T = (Q^T)^T Q^T$ und damit ist also auch Q^T orthogonal und deshalb können wir in der Definition einer orth. Matrix das Wort Spalten durch Zeilen ersetzen.

• Q lässt Abstände und Winkel konstant

$$\langle Qx, Qy \rangle = (Qx)^T (Qy) \Rightarrow x^T Q^T Q y = x^T y = \langle x, y \rangle$$

• Wir können die letzten $n-m$ Spalten von Q weglassen und die QR-Zerlegung als $A = QR$ mit $Q \in \mathbb{R}^{n \times n}$, $Q^T Q = \text{id}$ und $R \in \mathbb{R}^{m \times m}$ über Δ -Matrix definieren.

$$n \begin{bmatrix} A \\ m \end{bmatrix} = \begin{bmatrix} \square \\ m \end{bmatrix} n \begin{bmatrix} \square \\ m \end{bmatrix}^m$$

• Das Produkt orthogonaler Matrizen ist orthogonal

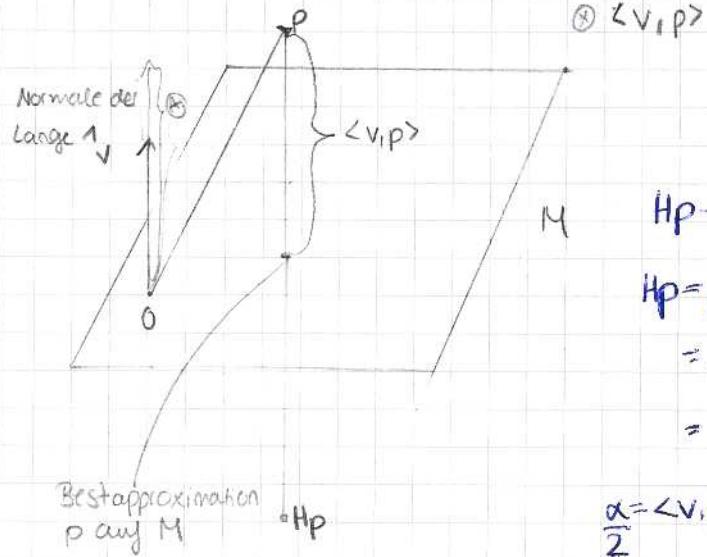
• Jede Matrix, deren assozierte lineare Abbildung das Skalarprodukt invariant lässt, ist orthogonal

⇒ Drehungen und Spiegelungen sind orthogonal
Spiegelungen (an Hyperebenen)

Def Eine Hyperebene im \mathbb{R}^n ist ein $(n-1)$ -dimensionaler möglicherweise affiner Unterraum

Def Eine Spiegelung f an einer Hyperebene M ist definiert durch die Eigenschaft, dass jeder Punkt p auf einen Punkt $f(p)$ abgebildet wird, der denselben Abstand zu M hat, so dass die Punkt-f(P) verbindende Gerade orthogonal auf M steht.

- Householder-Spiegelung/Transformation ist eine Spiegelung an einer Hyperebene, die durch 0 geht.
- orthogonales Komplement eines $(n-1)$ -dimensionalen Unterräums ist eindimensional und wird von den sogenannten Normalen aufgespannt d.h. wir können Householder-Spiegelung eindeutig durch diese Normale beschreiben.



$$H_p - p = \alpha \cdot v$$

$$H_p = p + \alpha v$$

$$= p - 2v <v, p>$$

$$= (I - 2vv^T)p$$

$$\alpha = \frac{<v, p>}{2} \text{ bitte überprüfen}$$

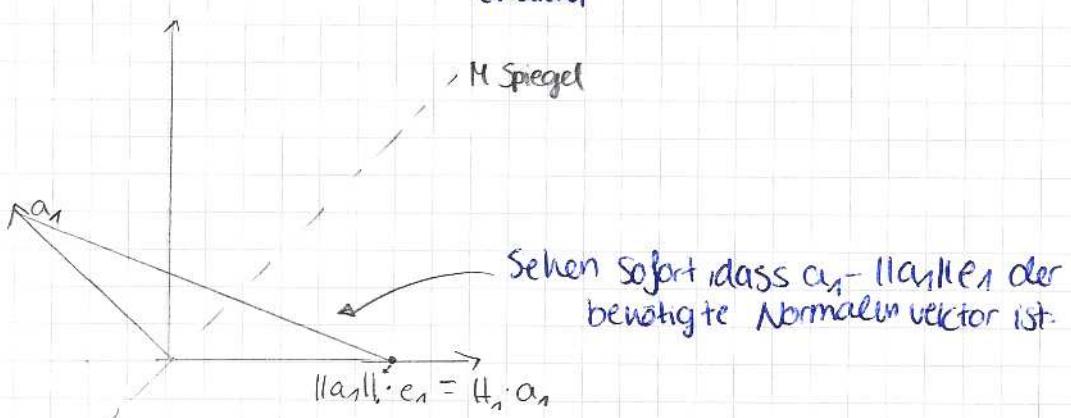
Satz Sei v eine Normale der Länge 1 auf M . Dann lässt sich die Spiegelung H an M schreiben als $H_p = (I - 2vv^T)p$

H ist sogar symmetrisch $H = H^T = H^{-1}$

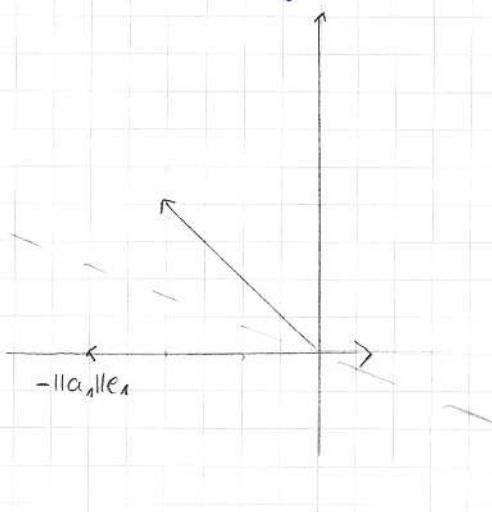
Ziel:

$$H_1 \begin{pmatrix} \square \\ a_1 \end{pmatrix} = \begin{pmatrix} \square \\ * \end{pmatrix} \text{ d.h. } H_1 \cdot a_1 = \begin{pmatrix} * \\ 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 1 \\ 0 \end{pmatrix} = \pm \|a_1\| \cdot e_1$$

weil H_1 ortho.
und damit längen-
erhaltend



Haben zwei Möglichkeiten, zu spiegeln. Die andere ist:



Numerisch stabiler ist aber die,

wodurch $(H_p - p)$ größer ist. Deshalb wählen wir:

$$w = a_1 + \text{sgn}(a_1) \|a_1\| e_1$$

Müssen nach normieren und erhalten

$$v = \frac{w}{\|w\|} \text{ als normierte Normale bzw. Einheitsnormale}$$

Next:

$$\begin{matrix} 1 & 0 \\ 0 & 0 \end{matrix} = \begin{matrix} x & -x \\ -x & 0 \end{matrix} = \begin{matrix} y & x \\ -x & 0 \end{matrix}$$

$$\tilde{H}_2 \quad R_1 \quad R_2$$

Haben die Bestimmung von H_2 also auf die von H_1 zurückgeführt

$$H_m \cdot H_2 \cdot H_1 \cdot A = R$$

$$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} = \begin{matrix} x & -x \\ -x & 0 \end{matrix} = \begin{matrix} y & x \\ -x & 0 \end{matrix}$$

$$\Rightarrow A = H_1 \cdots H_m R \Rightarrow Q = H_1 \cdots H_m$$

$$\begin{matrix} 0 & 0 \\ 0 & 0 \end{matrix} = \begin{matrix} x & -x \\ -x & 0 \end{matrix} = \begin{matrix} y & x \\ -x & 0 \end{matrix}$$

Auch möglich: Drehungen verwenden: Stichwort: "Givens-Rotation"

$$\begin{pmatrix} \cos \alpha & -\sin \alpha \\ \sin \alpha & \cos \alpha \end{pmatrix} \begin{pmatrix} * & * \\ * & * \end{pmatrix} = \begin{pmatrix} * & * \\ 0 & * \end{pmatrix}$$

23.4 Givens Rotation

\mathbb{R}^2 : orthogonale Matrizen

$$G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \text{ mit } c^2 + s^2 = 1$$

oder $c = \cos \varphi, s = \sin \varphi \quad \varphi \in [0, 2\pi]$

Grundaufgabe zu gegebenem $y = \begin{pmatrix} a \\ b \end{pmatrix}$, konstruiere G und $\begin{pmatrix} r \\ 0 \end{pmatrix}, r \in \mathbb{R}$ geeignet

$$\begin{pmatrix} r \\ 0 \end{pmatrix} = G \begin{pmatrix} a \\ b \end{pmatrix} \Rightarrow \begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix} = \begin{pmatrix} \sqrt{a^2 + b^2} \\ 0 \end{pmatrix} \Rightarrow c = \frac{a}{r}, s = \frac{b}{r}$$

$$\|Gy\|^2 = \|y\|^2 = \left\| \begin{pmatrix} r \\ 0 \end{pmatrix} \right\|^2 = r^2 = \left\| \begin{pmatrix} a \\ b \end{pmatrix} \right\|^2 = a^2 + b^2 \Rightarrow r = \sqrt{a^2 + b^2}$$

Probe: $\begin{pmatrix} \frac{a}{r} & \frac{b}{r} \\ -\frac{b}{r} & \frac{a}{r} \end{pmatrix} \cdot \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \frac{a^2 + b^2}{r} \\ -\frac{ab}{r} + \frac{ab}{r} \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$

$$G_{ijk} \begin{pmatrix} I & & 0 \\ 0 & c & s \\ 0 & -s & c \end{pmatrix} \xrightarrow{\text{j-Zeil}} \begin{pmatrix} a_{ik} \\ a_{ij} \\ a_{ik} \end{pmatrix} = \begin{pmatrix} r \\ \frac{a_{ij} + a_{ik}}{\sqrt{a_{ij}^2 + a_{ik}^2}} \\ 0 \end{pmatrix}$$

$$\rightsquigarrow G_1 \dots G_m A = R$$

Q

Anwendung z.B.: A Hessenberg-Matrix

Bemerkungen:

1) QR-Zerlegung eignet sich besonders für Matrizen, die „nahe zu singulär“ (schlecht konditioniert also die Matrizen) sind

$$\|Ax\| = \|Q Rx\| = \|Rx\|$$

1) ohne Pivotierung!

3) dünn besetzte Matrizen \rightsquigarrow Givens-Rotation

4) Komplexität: \approx doppelt soviel Operation wie LU-Zerl. (Gauß) für Householder

5) QR-Zerlegung fürsbesondere Householder tr. eignen sich zur stabilen Orthogonalisierung

6) ...

2.4. Fehleranalyse beim Lösen lin. Gl.

Sei V Vektorraum z.B. $V = \mathbb{R}^n$

Norm $\|\cdot\|: V \rightarrow \mathbb{R}$ mit

- 1) $\|u\| \geq 0, \forall u \in V$ und $\|u\| = 0 \Leftrightarrow u = 0$
- 2) $\|\alpha u\| = |\alpha| \|u\| \quad \forall u \in V \text{ und } \alpha \in K (\text{z.B. } \mathbb{R})$
- 3) $\|u+v\| \leq \|u\| + \|v\| \quad \forall u, v \in V$

Skalarprodukt $K = \mathbb{R} \quad \langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{R}$

mit 1) $\langle u, v \rangle = \langle v, u \rangle$

2) $\langle \alpha u + \beta v, w \rangle = \alpha \langle u, w \rangle + \beta \langle v, w \rangle$

3) $\langle u, u \rangle \geq 0$ f.a. $u \in V$ mit Gleichheit fallen $u=0$

$\|u\| := \sqrt{\langle u, u \rangle}$. Falls V vollständig ist. $(V, \langle \cdot, \cdot \rangle)$ ein Hilbertraum

Beispiel $V = \mathbb{R}^n$; $u = (u_1, \dots, u_n)^T \in V$

1) $\langle u, v \rangle = \sum_{i=1}^n u_i v_i$

2) Für $1 \leq p < \infty$ def.

$$\|u\|_p = \left(\sum_{i=1}^n |u_i|^p \right)^{1/p}, \quad p \text{-Norm, } l_p\text{-Norm}$$

3) $p=2$: $\|u\|_2^2 = \sum_{i=1}^n |u_i|^2 \quad \langle u, u \rangle \quad \text{Euklidische Norm}$

4) $p=\infty$ $\|u\|_\infty = \max_{1 \leq i \leq n} |u_i|$

Proposition: Es gilt $\forall u \in \mathbb{R}^n$

$$\|u\|_\infty \leq \|u\|_2 \leq \|u\|_1 \leq \sqrt{n} \|u\|_2 \leq n \|u\|_\infty$$

Ein Operator $A: V \rightarrow W$ (norm. Räume)

- linear: $\forall \alpha, \beta \in \mathbb{R}, u, v \in V$

$$A(\alpha u + \beta v) = \alpha A u + \beta A v$$

- beschränkt $\forall u \in V : \exists c > 0$

$$\|A u\|_W \leq c \|u\|_V$$

(A linear stetig \Leftrightarrow Lipschitz stetig \Leftrightarrow beschränkt)

$$\|A\| = \sup_{\|V\|_V=1} \|A V\|_W \quad (= \sup_{V \neq 0} \frac{\|A V\|_W}{\|V\|_V})$$

Satz 4.9 1) Sei $\mathcal{L}(V,W) = \{A: V \rightarrow W : \text{linz, beschränkt}\}$

dann $A \mapsto \|A\|$ eine Norm in $\mathcal{L}(V,W)$

2) Seien $A \in \mathcal{L}(V,W), B \in \mathcal{L}(W,Z) \Rightarrow BA = B \circ A \in \mathcal{L}(V,Z)$

und $\|BA\| \leq \|B\| \|A\|$

Bew 1) $u \in V$ bel. $0 \leq \|Au\|_W \leq c \|u\|$
 $\Rightarrow 0 \leq \|A\| \leq c$; $\|A0\| = \|0\| = 0$

Sei $\|A\| = 0$ d.h. $\sup_{\|u\|_V=1} \|Au\|_W = 0 \Rightarrow \|Au\|_W = 0 \forall u$

$\forall u \in V : Au = 0$ d.h. $A = 0$ d.h. $\|A\| = 0 \Leftrightarrow A = 0$
 $\|A\| = |\lambda| \|A\|$ k.l.s

$$\begin{aligned} \forall A, B \quad & \|A+B\| = \sup_{\|u\|_V=1} \|(A+B)u\|_W \\ & \leq \sup_{\|u\|_V=1} (\|Au\|_W + \|Bu\|_W) \\ & \leq \sup_{\|u\|_V=1} \|Au\|_W + \sup_{\|v\|_V=1} \|Bv\|_W = \|A\| + \|B\| \end{aligned}$$

2) Aus $\sup_{u \neq 0} \frac{\|Au\|_W}{\|u\|_V} = \|A\|$

$$\begin{aligned} \forall u \in V \quad & \|Au\|_W \leq \|A\| \|u\|_V \\ \Rightarrow \sup_{\|u\|_V=1} & \|Bu\|_Z = \sup_{\|u\|_V=1} \|B\| \|Au\|_W \\ & \leq \|B\| \sup_{\|u\|_V=1} \|Au\|_W = \|B\| \|A\| \end{aligned}$$

Sei $A \in \mathbb{R}^{n \times n}$ $\lambda_1, \dots, \lambda_n \in \mathbb{C}$ Eigenwerte von A

$\rho(A) = \max_{1 \leq i \leq n} |\lambda_i|$ „Spektralradius“

Sei $V = \mathbb{R}^n = W$, $A = (a_{ij})_{1 \leq i,j \leq n} \in \mathbb{R}^{n \times n}$

1) Spaltensummennorm

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

2) Zeilensummennorm

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

3) Spektralnorm

$$\|A\|_2 = \sqrt{\rho(A^T A)}$$

Satz 1) Sei $V = \mathbb{R}^n$ versehen mit der $\|\cdot\|_1$,

dann gilt für $A: V \rightarrow V$: $\|A\| = \|A\|_1$

2) Sei $V = \mathbb{R}^n$ versehen mit der $\|\cdot\|_\infty$

dann gilt für $A: V \rightarrow V$: $\|A\| = \|A\|_\infty$

3) " $\Rightarrow \|A\| = \|A\|_2$

4) $\|A\|_2 \leq \|A\|_1 \leq \|A\|_\infty$

$(a_{ij}) = A \in \mathbb{R}^{n \times m}$ $A: \mathbb{R}^m \rightarrow \mathbb{R}^n$ linear

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

② $\|A\|_2 = \sqrt{\lambda_{\text{max}}(A^T A)}$ Spektralnorm

größter Eigenwert,
sehr schwer zu bestimmen

Bew 1) Sei $\|u\|_1 = \|u\|_1 = \sum_{j=1}^m |u_j|$, $u \in \mathbb{R}^m$

$\lceil A: V \rightarrow W, V = \mathbb{R}^m, W = \mathbb{R}^n \rceil$

$(V, \|\cdot\|_V) = (\mathbb{R}^m, \|\cdot\|_1)$; $(W, \|\cdot\|_W) = (\mathbb{R}^n, \|\cdot\|_1)$

Operatornorm $\|A\| = \sup_{\substack{\|u\|_V=1}} \|Au\|_W$

$\Rightarrow \forall v \in V = \|Av\|_W \leq \|A\| \|v\|_V$

\leq " Sei $u \in V$ mit $\|u\|_1 = 1$

$$\begin{aligned}\|Au\|_1 &= \sum_{i=1}^n \left| \sum_{j=1}^m a_{ij} u_j \right| \leq \sum_{i=1}^n \sum_{j=1}^m |a_{ij}| |u_j| = \|u\|_1 \\ &= \sum_{j=1}^m |u_j| \left(\sum_{i=1}^n |a_{ij}| \right) \leq \max_{1 \leq i \leq n} \sum_{j=1}^m |a_{ij}| \underbrace{\sum_{j=1}^m |u_j|}_{\|u\|_1} \\ \|u\|_1 &= \max_{1 \leq i \leq n} \sum_{j=1}^m |a_{ij}| = \|A\|_1\end{aligned}$$

$\Rightarrow \|A\| \leq \|A\|_1$

operatornorm

Sei $j \max$, sodass $\sum_{i=1}^n |a_{ij}| = \max_{1 \leq j \leq m} \sum_{i=1}^n |a_{ij}|$

\geq Setze $j = j_{\max}$ und $u_j = e_j = (0 \ 0 \ 1 \ 0 \dots 0)^T \Rightarrow \|u_j\|_1 = 1$

$$\Rightarrow \|Au_j\|_1 = \|Ae_j\|_1 = \left\| (a_{i,j})_{i=1}^n \right\|_1 = \sum_{i=1}^n |a_{i,j}| = \|A\|_1 \leq \|A\|$$

2) $p=2$: $B = A^T A = B^T = B$ pos semi definit,

$$\text{denn } \langle v, Bv \rangle = v^T B v = v^T A^T A v = (Av)^T (Av) = \langle Av, Av \rangle = \|Av\|^2 \geq 0$$

Seien $\lambda_1, \dots, \lambda_m$ Eigenwerte und u_1, \dots, u_m , $\|u_j\|=1$ zugehörige EV.

$$Bu_j = \lambda_j u_j$$

Sei $u = \sum_{j=1}^m \alpha_j u_j \in \mathbb{R}^m$ mit $1 = \|u\|_2^2 = \sum_{j=1}^m |\alpha_j|^2$ bel

$$\|Au\|_2^2 = \langle Au, Au \rangle = u^T A^T A u$$

[u_j bilden ONB]
[weil B symm.]

$$= u^T B \sum \alpha_j u_j \stackrel{\text{linear}}{=} u^T \sum_{j=1}^m \lambda_j \alpha_j u_j$$

$$= \sum_{j,k=1}^m \underbrace{\alpha_k \lambda_{k,j} \alpha_j}_{= u^T} \underbrace{u_k^T u_j}_{= 1 \delta_{kj}} = \sum_{j=1}^m \lambda_j \alpha_j^2$$

$$\leq \max_{1 \leq j \leq m} |\lambda_j| \sum_{j=1}^m \alpha_j^2 = \sigma(B) = \sigma(A^T A)$$

$$\Rightarrow \sup_{\|u\|_2=1} \|Au\|_2 \leq \sqrt{\sigma(A^T A)} \quad \text{d.h. } \|A\| \leq \|A\|_2$$

Sei j , sodass $\lambda_j = \max_{1 \leq j \leq m} \lambda_j$. Wir nehmen $u = u_j$ EW zu λ_j

$$\|u\|_2 = \|u_j\|_2 = 1$$

$$\|Au\|_2^2 = \|Au_j\|_2^2 = \langle Au_j, Au_j \rangle = u_j^T A^T A u_j$$

$$= u_j^T B u_j = \lambda_j u_j^T u_j = \lambda_j = \sigma(B) = \sigma(A^T A) = \|A\|_2^2$$

$$\Rightarrow \|A\|_2^2 \leq \|A\|$$

3) Frobenius-Norm:

$$\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^m |a_{ij}|^2 \quad (\text{ist aber keine Operatormodulnorm}). \text{ Es gilt } \|A\|_2 \leq \|A\|_F$$

$$\text{Bsp: } A = \begin{pmatrix} 1 & -3 \\ 1 & 1 \end{pmatrix} \in \mathbb{R}^{2,2}$$

$$\Rightarrow 1) \|A\|_1 = \max \{ |1|, |-3|, |1| \} = 4$$

$$2) A^T A = \begin{pmatrix} 2 & -2 \\ -2 & 10 \end{pmatrix} =: B, \lambda_{1,2} = 6 \pm \sqrt{60}$$

$$\Rightarrow \|A\| = \sqrt{6+60} =$$

$$3) \|A\|_F^2 = 1^2 + 1^2 + (-3)^2 + 1^2 = 12 \Rightarrow \|A\|_F = \sqrt{12}$$

2.4.2. Kondition von Matrizen & Fehleranalyse

Def $A \in \mathbb{R}^{n \times n}$, dann heißt

$\|A\|_2$ Op-Norm in ℓ_2

$$\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2 \text{ die Kondition von } A$$

Sei x Lösung von $Ax=b$ und gestörtes Problem $A(x+\Delta x)=b+\Delta b$ A,b gegeben

$$\Rightarrow A\Delta x = \Delta b \text{ bzw. } \Delta x = A^{-1}\Delta b$$

$$\Rightarrow \frac{\|\Delta x\|}{\|x\|} = \frac{\|A^{-1}\Delta b\|}{\|x\|} \leq \frac{\|A^{-1}\|}{\|x\|} \frac{\|\Delta b\|}{\|b\|}$$

wenn A singulär, gibt es Probleme mit $\text{cond} = \infty$

rel. Fehler der Lsg.

$$= \frac{\|A^{-1}\|}{\|x\|} \frac{\|\Delta b\|}{\|b\|} \quad \leftarrow \text{rel. Fehler der Eingabe}$$

$$= \frac{\|A^{-1}\|}{\|x\|} \frac{\|A\|}{\|x\|} \frac{\|\Delta b\|}{\|b\|} \leq \frac{\|A^{-1}\|}{\|x\|} \|A\| \frac{\|x\|}{\|x\|} \frac{\|\Delta b\|}{\|b\|} = \text{cond}_2(A) \frac{\|\Delta b\|}{\|b\|}$$

„Kondition Abschätzung für Fehlersfortpflanzung bei Lösung von linearen Gleichungssystemen mit gestörten Eingabedaten“

Lemma 2.4.3 Sei $A \in \mathbb{R}^{n \times n}$ mit $\|A\| < 1$, dann existiert $(I+A)^{-1}$ und

$$\frac{1}{1+\|A\|} \leq \|(I+A)^{-1}\| \leq \frac{1}{1-\|A\|}$$

Bew. $(I+A)^{-1} = (I+(-A))^{-1} = \sum_{k=0}^{\infty} (-A)^k A^{-1}$ Neumannsche Reihe

Exkurs: $(I+(-A))x = b$, $x = -Ax + b = f(x)$ Benach FPS anwenden

$$\|f(x) - f(x')\| = \| -Ax + b - Ax' + b \| = \|A(x-x')\| \leq \|A\| \|x-x'\|$$

→ Kontraktion \Downarrow weil $\|A\| < 1$

$$\Rightarrow (I+A)^{-1}$$

Alternativ:

Sei $u \in \ker(I+A)$ d.h. $u+Au=0 \Leftrightarrow \|u\| = \|Au\| \geq \|u\|$ $\overset{\text{für } u \neq 0}{\rightarrow}$

$$\Rightarrow u=0 \Rightarrow \ker(I+A) = \{0\} \Rightarrow I+A \text{ regulär}$$

$$\|(I+A)x\| = \|x+Ax\| \geq \|x\| - \|Ax\| \geq \|x\| - \|A\|\|x\| = (1-\|A\|)\|x\|$$

$$\text{und } \|(I+A)x\| = \|x+Ax\| \leq (1+\|A\|)\|x\|$$

$$\|y\| = \|(I+A)x\| \geq (1-\|A\|)\|x\|$$

$$\|y\| \leq (1+\|A\|)\|x\|$$

$$\Rightarrow \frac{\|y\|}{(1+\|A\|)} = \|x\| = \|(I+A)^{-1}y\| \leq \frac{\|y\|}{(1-\|A\|)} \quad \forall y$$

$$\frac{1}{(1+\|A\|)} \leq \sup_{\|y\|=1} \|(I+A)^{-1}y\| \leq \frac{1}{1-\|A\|}, \frac{1}{1-\|A\|} \leq \|(I+A)^{-1}\| \leq \frac{1}{1-\|A\|} \quad \square$$

$\text{cond}_2(A) = \|A\| \|A^{-1}\|$, $\|A\|$ Spektralnorm \triangleq Op.-Norm in $\ell_2, \|\cdot\|_2$ NumI, VL

Lemma 2.4.13 $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$, $\|A\| < 1 \Rightarrow \|(I + A)^{-1}\| \leq \frac{1}{\|A\|}$

Satz 2.4.14

Sei $A \in \mathbb{R}^{n \times n}$ regulär. Sei $\Delta A \in \mathbb{R}^{n \times n}$ mit $\|\Delta A\| \leq \frac{1}{\|A^{-1}\|}$

$b, \Delta b \in \mathbb{R}^n$ gegeben

bwz. $\|\Delta A\| \|A^{-1}\| \leq 1$

$x \in \mathbb{R}^n$ Lsg von $Ax = b$ (d.h. $x = A^{-1}b$) (1)

und $(A + \Delta A)(x + \Delta x) = b + \Delta b \rightarrow x + \Delta x \in \mathbb{R}^n$ Lsg der gestörten (2) Gleichung.

Dann gilt $\|\Delta x\| \leq \frac{\text{cond}_2(A)}{1 - \text{cond}_2(A)\frac{\|A\|}{\|A^{-1}\|}} \frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|}$

$$\begin{aligned}
 \text{Beweis (2)} \quad M := (\Delta A + A) \Rightarrow \Delta x &= M^{-1}b + M^{-1}\Delta b - x \quad [\text{da } M = A + \Delta A = A(I + A^{-1}\Delta A)] \\
 &= M^{-1}b + M^{-1}\Delta b - A^{-1}b \quad \begin{matrix} \text{reg.} \\ \text{reg.} \\ \|A^{-1}\Delta A\| \leq \|A^{-1}\| \|\Delta A\| \end{matrix} \\
 &\quad - M^{-1}\Delta b + M^{-1}(A - M)A^{-1}b \\
 &= M^{-1}\Delta b + M^{-1}(A - (A + \Delta A))A^{-1}b \Rightarrow M^{-1}(I + \Delta A A^{-1})A^{-1}b \\
 &= M^{-1}(\Delta b - \Delta A A^{-1}b)
 \end{aligned}$$

$$\|\Delta x\| \leq \|M\|^{-1}(\|\Delta b\| + \|\Delta A A^{-1}b\|)$$

$$\begin{aligned}
 \Leftrightarrow \frac{\|\Delta x\|}{\|x\|} &\leq M^{-1} \left(\frac{\|\Delta b\|}{\|x\|} + \frac{\|\Delta A A^{-1}b\|}{\|x\|} \right) \\
 &\leq \|M^{-1}\| \|A\| \left(\frac{\|\Delta b\|}{\|A\| \|A^{-1}b\|} + \frac{\|\Delta A\| \|A^{-1}\| \cdot \|b\|}{\|A\| \|A^{-1}b\|} \right) \\
 &\stackrel{\|\Delta A\| \leq \frac{1}{\|A^{-1}\|}}{\leq} \|M^{-1}\| \|A\| \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right) \leq \frac{1}{\|A\|} \text{cond}_2(A) \left(\frac{\|\Delta b\|}{\|b\|} + \frac{\|\Delta A\|}{\|A\|} \right) \stackrel{\|\Delta A\| \leq \frac{1}{\|A^{-1}\|}}{\leq} \frac{1}{\|A\| - \frac{1}{\|A\|}} \quad (\square)
 \end{aligned}$$

$$\|M^{-1}\|: M^{-1} = (I + A^{-1}\Delta A)^{-1} A^{-1}$$

$$\begin{aligned}
 \|M^{-1}\| &\leq \|(I + A^{-1}\Delta A)^{-1}\| \|A^{-1}\| \\
 &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \|\Delta A\|} \leq \frac{\|A^{-1}\|}{1 - \frac{1}{\|A\|} \|\Delta A\|} = \frac{\|A^{-1}\|}{\frac{\|A\| - \|\Delta A\|}{\|A\|}} = \frac{\|A^{-1}\|}{\frac{\|A\| - \text{cond}_2(A)}{\|A\|}} = \frac{\|A^{-1}\|}{\text{cond}_2(A)}
 \end{aligned}$$

$$\text{Bsp } A = \begin{pmatrix} 3 & 1,001 \\ 6 & 1,997 \end{pmatrix}, b = \begin{pmatrix} 1,999 \\ 4,003 \end{pmatrix}$$

$$\tilde{A} = A + \Delta A = \begin{pmatrix} 3 & 1 \\ 6 & 1,997 \end{pmatrix},$$

$$\tilde{b} = b + \Delta b = \begin{pmatrix} 2,002 \\ 4 \end{pmatrix}$$

\Rightarrow mit Satz 2.4.19: $\|A\| \leq \|A\|_{\infty} = 6 + 1,997 = 7,997$

$$A^{-1} = -\frac{1}{6,015} \begin{pmatrix} 1,997 & -1,001 \\ -6 & 3 \end{pmatrix} \quad \|A^{-1}\| \leq \|A\|_{\infty} = 600$$

$$\Rightarrow \text{cond}_2(A) \leq \text{cond}_{\infty}(A) = 4798,2, \|\Delta b\|_{\infty} = 0,003$$

$$\Rightarrow \frac{\|Ax\|_{\infty}}{\|x\|_{\infty}} \leq 10,4898 \quad \text{d.h. (exakter) Rechnung} \Rightarrow \frac{\|Ax\|}{\|x\|_{\infty}} \approx 2,333$$

Abschätzung war korrekt, aber sehr pessimistische

Bemerkung

1) $\|A^{-1}\|$ ist aufwändig zu berechnen.

Aus L-R-Zerlegung (z.B. Cholesky) können $\|L^{-1}\|, \|R^{-1}\|$ leichter abgeschätzt werden.

2) QR-Zerlegung. Da $\|A\| = \|QR\| \leq \underbrace{\|Q\|}_{=1} \cdot \|R\| = \|R\|$

$$\text{und } \|A^{-1}\| = \|R^{-1}Q^T\| = \|R^{-1}\|$$

$\Rightarrow \text{cond}_2(A) = \text{cond}_2(R)$ Stabilität der QR-Zerlegung

3) Die EW λ_{\max} und λ_{\min} $\Rightarrow \text{cond}_2(R) \approx \frac{1}{\lambda_{\min}}$ (gute Schätzung)

2.5. Lineare Regression (Ausgleichsrechnung)

- Methode der kleinsten Fehlerquadrate (Gauß) (Least-Squares)

Zu ~~den Stützstellen~~ $x_i \in \mathbb{R}^n, i=1, \dots, m$ existieren Stützwerte $y_i \in \mathbb{R}$ (mit $m > n$)

mit $y_i \approx f(x_i)$; Angenommen f : $y = (y_1, \dots, y_m)^T \in \mathbb{R}^m$ gegeben.

ges $x \in \mathbb{R}^n \Rightarrow y = f(x)$

Bsp: $f(x) = Ax, x \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n} \quad \boxed{m \times n \times b}$

d.h. Gleichungssystem ist überbestimmt

\Rightarrow Gl.-syst. $A x - b = 0$ i.A. nicht lösbar

d.h. $0 \neq r := Ax - b$ (Residuum)

Optimierungsaufgabe $\|r\| \rightarrow \text{Min}$ d.h. gesucht $x \in \mathbb{R}^n$ mit

$$x = \underset{x \in \mathbb{R}^n}{\text{Argmin}} \|r(x)\|, x \in \mathbb{R}^n \text{ d.h. } \|r(x)\| \leq \|r(v)\| \forall v \in \mathbb{R}^n \iff \|r(x)\|^2 \leq \|r(v)\|^2$$

$$\|r\| = \|r\|_2$$

$$\Rightarrow x := \underset{v \in \mathbb{R}^n}{\operatorname{argmin}} \|r(v)\|_2^2 = \langle Av - b, Av - b \rangle = f(v)$$

$$x \text{ lässt } \nabla f(x) = \mathbf{0}; \quad f(x) = \|r(x)\|^2$$

notwendig

$$\|Ax - b\|_2^2 \rightarrow \min$$

$$A = (a_{ij}), v = \begin{pmatrix} v_1 \\ \vdots \\ v_m \end{pmatrix}$$

$$Ax - b \stackrel{!}{=} 0?$$

$$\|Ax - b\| \rightarrow \min$$

Optimierungsproblem:

$$x \rightarrow J(x) = \|Ax - b\|^2 = \langle Ax - b, Ax - b \rangle \quad \text{in } x \in \mathbb{R}^n$$

$$J: \mathbb{R}^n \rightarrow \mathbb{R}, J \in C^1(\mathbb{R}^n)$$

Sei $x \in \mathbb{R}^n, x = \operatorname{argmin} \{J(x) : x \in \mathbb{R}^n\}$ eine Minimalstelle

\Rightarrow (Ana II) $J'(x): \mathbb{R}^n \rightarrow \mathbb{R}, J'(x) = \vec{0}$ (notwendige Bedingung)

$$J'(x) = \left(\frac{\partial J}{\partial x_1}(x), \dots, \frac{\partial J}{\partial x_n}(x) \right), \quad \nabla J(x) = J'(x)^T$$

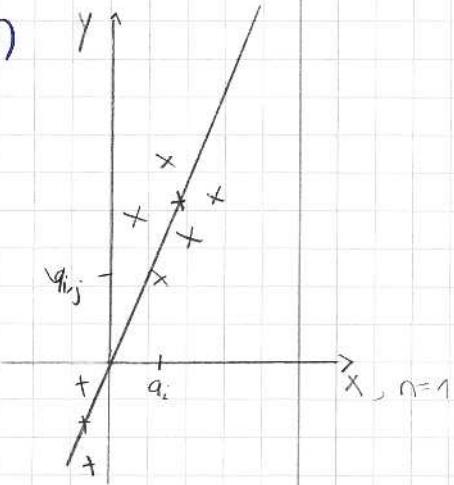
$$v \mapsto J(x+v) = J(x) + J'(x)v + O(|v|^2)$$

$$\begin{aligned} J(x+v) &= \langle A(x+v) - b, A(x+v) - b \rangle = \langle A(x+v), A(x+v) \rangle - \langle b, A(x+v) \rangle \\ &\quad - \langle A(x+v), b \rangle + \langle b, b \rangle \\ &= \langle A(x+v), A(x+v) \rangle - 2\langle b, A(x+v) \rangle + \langle b, b \rangle \\ &= \underbrace{\langle Ax, Ax \rangle}_{= J(x)} - 2\langle b, Ax \rangle + \underbrace{\langle b, b \rangle}_{= J'(x)} + 2\langle A^T A x, v \rangle - 2\langle A^T b, v \rangle + \|v\|^2 \\ &= \langle Ax - b, Ax - b \rangle + 2(A^T A x - A^T b)^T v + O(\|v\|^2) \end{aligned}$$

$$\text{mit } J'(x) = 2(A^T A x - A^T b)^T$$

$$x = \operatorname{argmin}_{v \in \mathbb{R}^n} J(v) \Rightarrow 0 = J'(x) = 2(A^T A x - A^T b), \text{ d.h. } \boxed{A^T A x = A^T b}$$

\rightsquigarrow Gaußsche Normalengl.



Satz 2.5.1 (Lineare Regression)

1) Das Optimierungsproblem $\|Ax - b\| \rightarrow \min$ besitzt eine Lösung

$x \in \mathbb{R}^n$ und $x \in \mathbb{R}^n$ löst die Gaußsche Normalengleichung (GNG)

$$A^T A x = A^T b$$

2) Falls Rang A maximal, d.h. $\text{Rang } A = n$ ($m \geq n$) dann ist GNG eindeutig lösbar, d.h. es existiert genau eine Minimalstelle x .

Bew. 1) Sei $\min J(x) = \|Ax - b\|^2$, d.h. x Min-Stelle

$$\Rightarrow \frac{1}{2} J'(x)^T = 0 = A^T A x - A^T b \quad (\text{GNG})$$

\Leftarrow Sei $\tilde{x} \in \mathbb{R}^n$ mit $A^T A \tilde{x} - A^T b = 0$ und $x \in \mathbb{R}^n$ bel.

$$\Rightarrow 2(x - \tilde{x})^T (A^T A \tilde{x} - b) = 0$$

$$J(x) = \|Ax - b\|^2 = \|A(x - \tilde{x}) + A\tilde{x} - b\|_A^2 = \langle \cdot, \cdot \rangle$$

$\left\langle A(x - \tilde{x}), A(\tilde{x}) \right\rangle$

$$= \langle A(x - \tilde{x}), A(x - \tilde{x}) \rangle + 2 \underbrace{\langle (x - \tilde{x})^T A^T A \tilde{x} \rangle}_{= 0, \text{ da } \tilde{x}} + \|A\tilde{x} - b\|^2$$

$$= \|A(x - \tilde{x})\|^2 + \|A\tilde{x} - b\|^2$$

$= 0, \text{ da } \tilde{x}$
Lösung der
GNG

$$\geq \|A\tilde{x} - b\|^2 = J(\tilde{x}) \Rightarrow \tilde{x} = \arg \min_{v \in \mathbb{R}^n} J(v)$$

2) $m \geq n$: $\text{Rang } A = n \Rightarrow \dim \text{Kern}(A) = 0$

$$\forall q \in \mathbb{R}^n: q^T A^T A q = \|Aq\|^2 > 0 \Leftrightarrow A = 0 \Leftrightarrow q = 0 \in \mathbb{R}^n$$

$\Rightarrow A^T A$ ist SPD $\Rightarrow A^T A$ ist regulär \blacksquare

Bemerkung:

Sei $A \in \mathbb{R}^{n \times m}$ $\text{cond}_2(A) \gg 1$, dann ist $\text{cond}_2(A^T A) \approx (\text{cond}_2(A))^2 \gg \gg 1$,

d.h. GNG ist oft extrem schlecht konditioniert.

Zur Lösung des Optimierungsproblems bietet sich das QR-Verfahren an.

Dann $J(x) = \|Ax - b\|^2 = \|Q(Rx - Q^T b)\|^2 = \|Rx - Q^T b\|^2$

$$= \|\tilde{R}x - \tilde{b}_1\|^2 + \|\tilde{b}_2\|^2$$

$\Rightarrow x \in \mathbb{R}^n$ löst $\tilde{R}x = \tilde{b}_1$ (Rücksubstitution)

$$\text{und } \min_{v \in \mathbb{R}^n} J(v) = \|\tilde{b}_2\|^2$$

($\text{cond}_2(R) = \text{cond}_2(A)$)

gut konditioniert, $Q^T b = \tilde{b} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{pmatrix}, \tilde{b}_1 \in \mathbb{R}^m$

$$R = \begin{bmatrix} \tilde{R} \\ 0 \end{bmatrix} = \tilde{R} \in \mathbb{R}^{m \times n}$$

von $Q^T A$
obere AM

$$Q^T b = \tilde{b} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{pmatrix}$$

$\tilde{b}_1 \in \mathbb{R}^m$

Beispiel: $A \cdot \begin{pmatrix} 3 & 7 \\ 0 & 12 \\ 4 & 1 \end{pmatrix} \in \mathbb{R}^{3 \times 2}$, $b = \begin{pmatrix} 10 \\ 1 \\ 5 \end{pmatrix} \in \mathbb{R}^3$

QR mit Givens-Rot

$$G_{1,3} = \begin{pmatrix} 3/15 & 0 & 4/15 \\ 0 & 1 & 0 \\ -4/15 & 0 & 3/15 \end{pmatrix}, \quad G_{2,3} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 12/13 & -5/13 \\ 0 & 5/13 & 12/13 \end{pmatrix}$$

$$G_{2,3} G_{1,3} A = \begin{pmatrix} 5 & 5 \\ 0 & 13 \\ 0 & 0 \end{pmatrix}$$

$$Q^T = G_{2,3} G_{1,3} \quad b = \begin{pmatrix} 10 \\ 37/13 \\ -55/13 \end{pmatrix} = \begin{pmatrix} \tilde{b}_1 \\ \tilde{b}_2 \\ \tilde{b}_3 \end{pmatrix} \Rightarrow \tilde{b}_1 = \begin{pmatrix} 10 \\ 37/13 \end{pmatrix}$$

$$\|b_2\| = \left(\frac{55}{13} \right)$$

(ose)

$$R \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 5 & 5 \\ 0 & 13 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 10 \\ 37/13 \end{pmatrix} \Rightarrow x = \begin{pmatrix} 301/169 \\ 37/169 \end{pmatrix}$$

2.6. Singulärwertzerlegung und Pseudoinverse

Satz 2.6.1. (SVD, Singular value decomposition)

Sei $A \in \mathbb{R}^{m \times n}$ beliebig, dann ex. $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$

orthogonal und $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r \geq 0$, $r = \min(m, n)$ Singulärwerte (σ_i)

mit der Darstellung:

$$A = U \cdot \Sigma V^T \quad \text{bzw. } \Sigma = U^T A V = \text{diag}(\sigma_1, \dots, \sigma_r) = \begin{pmatrix} \sigma_1 & 0 & & \\ 0 & \ddots & & \\ & & \ddots & 0 \\ 0 & & & 0 \end{pmatrix}_{(m > n)}$$

$$\text{oder: } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0 \quad (r \leq \min(m, n))$$

$$\Rightarrow \Sigma = \begin{pmatrix} \sigma_1 & 0 & & \\ 0 & \sigma_2 & 0 & \\ & & \ddots & 0 \\ 0 & & & 0 \end{pmatrix}_{m \geq n}$$

Eigenwertzerlegung: $A = U \Sigma U^T$ d.h. $V = U$

falls $A = A^T$

bzw. $U \Sigma U^{-1}$

falls A SPD

Beweisidee: $B = A^T A$ ist symm. und pos. semidef., Die nichttriv. Eigenwerte von B sind die Singulärwerte.

SVD (2.6.1)

$$A = (a_{i,j}) \in \mathbb{R}^{m \times n} \Rightarrow \exists \quad \sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = 0$$

$U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ orthogonal

$$A = U \operatorname{diag}(\sigma_i) V^T, \operatorname{diag}(\sigma_i) = \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & 0 \\ & & \ddots & \\ 0 & & & 0 \end{bmatrix} = \Sigma$$

Beweis Sei OBdA $m \leq n$ und $B = A^T A \in \mathbb{R}^{n \times n}$

$$\text{Symmetrisch. } U^T B U = U^T A^T A U = (A U)^T A U = \|A U\|^2 \geq 0 \quad \forall U \in \mathbb{R}^n$$

$$\Rightarrow B v_k = \lambda_k v_k \quad \|v_k\| = 1 \Rightarrow \lambda_k \geq 0, k=1, \dots, n$$

$$\dots \lambda_{k-1} \geq \lambda_k \geq 0$$

und $V = (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$ orthogonal

$$\text{Wir definieren } \tilde{\sigma}_i = \sqrt{\lambda_i}, i=1, \dots, n \quad \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = 0 := \sigma_{n+1}$$

$$C = A A^T \in \mathbb{R}^{m \times m}$$

$$\text{Für } j=1, \dots, r \text{ Seien } u_j := \frac{1}{\tilde{\sigma}_j} A v_j \quad (r \leq \min\{m, n\})$$

$$\text{Für } j=r+1, \dots, m \text{ Sei } \{u_j : j=r+1, \dots, m\} \text{ eine ONB von } (\operatorname{span}\{u_i : i=1, \dots, r\})^\perp \\ = \{w \in \mathbb{R}^m : u_i^T w = 0 \forall i=1, \dots, r\}$$

$$\text{Zeigen } \langle u_i, u_j \rangle = u_i^T u_j = \delta_{i,j} \quad (i > r \text{ oder } j > r)$$

$$\text{Für } i, j = 1, \dots, r \Rightarrow u_i^T u_j = u_i^T A v_j \cdot \frac{1}{\tilde{\sigma}_j} = \frac{1}{\tilde{\sigma}_j} v_i^T \underbrace{A^T A}_{B} v_j \cdot \frac{1}{\tilde{\sigma}_j} \\ = \frac{1}{\tilde{\sigma}_i \tilde{\sigma}_j} \lambda_j \underbrace{v_i^T v_j}_{=\delta_{i,j}} = \frac{\lambda_j}{\tilde{\sigma}_j^2} \delta_{i,j}$$

$$U^T A V = \operatorname{diag}(\tilde{\sigma}_j) U^T U = (\tilde{\sigma}_j \delta_{i,j}) = \sum_{i=1, \dots, m} \tilde{\sigma}_i^2$$

Korollar: $A = U \Sigma V^T$ aus Satz 2.6.1 $\Rightarrow \operatorname{Rang} A = r$ ($\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = 0$)

$$1) \operatorname{Rang} A = r$$

$$2) \operatorname{Im}(A) = \operatorname{Bild}(A) = \operatorname{span}\{u_1, \dots, u_r\}$$

$$3) (\operatorname{Kern}(A))^{\perp} = \operatorname{span}\{v_1, \dots, v_r\}$$

$$\Leftrightarrow (w \in \operatorname{Kern} A \Leftrightarrow w^T v_j = 0, j=1, \dots, r)$$

Bemerkung: Rang r -Darstellung von A

$$A = \sum_{k=1}^r \tilde{\sigma}_k u_k \otimes v_k = \sum_{k=1}^r \tilde{\sigma}_k u_k v_k^T \quad (\text{Komplexität nnz} = \mathcal{O}(r + nr^2 + mr^2)) \\ = \sum_{k=1}^r u_k \tilde{v}_k^+ \quad (\tilde{v}_k = \tilde{\sigma}_k v_k)$$

Alternative Darstellung

$$A = \tilde{U} R \tilde{V}^T, \tilde{U} \in \mathbb{R}^{m \times m}, \tilde{V} \in \mathbb{R}^{n \times n} \text{ orth. } R \in \mathbb{R}^{m \times n} \text{ regulär}$$

$$\begin{matrix} \tilde{U} \\ \tilde{V} \end{matrix} \stackrel{\text{def}}{=} \begin{matrix} U \\ V \end{matrix} \tilde{G}^{-1}$$

$$A = \tilde{U} R \tilde{V}^T = \tilde{U} \tilde{Q} \tilde{V}^T \quad \tilde{U} = \tilde{U} G \quad \tilde{R} = G^{-1} R F^{-1}$$

$$\tilde{V} = \tilde{V} F$$

$$\text{span} \{ \tilde{u}_j \mid j=1, \dots, p \} = \text{span} \{ \tilde{u}_j \mid j=1, \dots, r \} = \text{Im}(A)$$

$$u, v \rightarrow \tau(v, \tilde{u}) = \sum_{j=1}^r u_j \tilde{v}_j = uv^T, \tau: \mathbb{R}^{m \times r} \times \mathbb{R}^{n \times r} \rightarrow \mathbb{R}^{m \times n} \text{ bilinear}$$

Korollar 2.6.3

- 1) Operator Norm $\|A\| = \|A\|_{op} = \sigma_1$
- 2) Frobenius Norm: $\|A\|_F^2 = \sum_{i,j} |a_{ij}|^2 = \sum_{k=1}^r \sigma_k^2$
- 3) Nuklear (Spur) Norm $\|A\|_1 = \sum_{k=1}^r \sigma_k$
- 4) Sei $A \in \mathbb{R}^{m \times n}$ ($m=n$) $\text{cond}_2 A = \frac{\sigma_1}{\sigma_n}$

Def. Verallgemeinerte Kondition

$$A \in \mathbb{R}^{m \times n}, \text{cond}_2 A = \frac{\sigma_1}{\sigma_r}$$

Satz 2.4 Satz von Schmidt & Mirsky, $A = U \Sigma V^T$

$$1) \underset{\text{Rang } B \leq S}{\arg \min} \|A - B\|_F = \sum_{k=1}^S \sigma_k u_k v_k^T, \min_{\text{rang } B=S} \|A - B\|_F^2 = \sum_{k=S+1}^r \sigma_k^2$$

$$A - B = \sum_{k=1}^S \sigma_k u_k v_k^T - \sum_{k=1}^S \sigma_k u_k v_k^T = \sum_{k=S+1}^r \sigma_k u_k v_k^T$$

Beste Approximation
 \Rightarrow SV Durch dann
 kleinste singuläre
 Werte weglassen
 also nur bis
 bestimmtes S
 gehen

$$2) \underset{\text{Rang } B=S}{\arg \min} \|A - B\|_{op} = \sum_{k=1}^S \sigma_k u_k v_k^T, \min_{\text{rang } B=S} \|A - B\|_{op} = \sigma_{S+1}$$

3) (analog zu 1.1.1.)

$$4) A, B \in \mathbb{R}^{m \times n}, \text{ mit Singularwerten } \sigma_{A,i}, \sigma_{B,i} \geq 0$$

$$\sum_{i=1}^{\min\{m,n\}} |\sigma_{A,i} - \sigma_{B,i}|^2 \leq \|A - B\|_F^2$$



Beispiel

$$A = \begin{pmatrix} 0,48 & -0,36 \\ 0,006 & 0,008 \\ 0,64 & 0,48 \end{pmatrix} \xrightarrow{\text{SVD}} \underbrace{\begin{pmatrix} 0,6 & 0 & 0,8 \\ 0 & 1 & 0 \\ -0,8 & 0,6 \end{pmatrix}}_{U} \underbrace{\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0,01 & 0 \\ 0 & 0 & 0 \end{pmatrix}}_{\Sigma} \underbrace{\begin{pmatrix} 0,8 & -0,6 \\ 0,6 & 0,8 \end{pmatrix}}_{V^T}$$

$\sigma_1 = 1$
 $\sigma_2 = 0,01$

Bemerkung Die Berechnung des SVD mit Eigenwertlösung zu Matrix $B = A^T A$ hat den Nachteil, dass $\text{cond}_2 B \approx (\text{cond}_2 A)^2$ ist schlecht konditioniert.

Stabile Alternative: Alg von Golub & Kahan (Matlab)

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

Rang $A = n$

$\Rightarrow \|Ax - b\|^2 \rightarrow \min$ Gauß-N.G

$$A^T A x = A^T b \text{ bzw. } x = (A^T A)^{-1} A^T b = A^+ b$$

$A^+ = (A^T A)^{-1} A^T$; Pseudo-Inverse (Moore-Penrose-Inverse)

Satz 4.6.5: Es gilt: $A = U \Sigma V^T$, $r = n$

$$A^+ = V \underbrace{\Sigma^+}_{\in \mathbb{R}^{n \times m}} U^T, \quad \Sigma^+ = \begin{pmatrix} \sigma_1^{-1} & 0 \\ 0 & \sigma_n^{-1} \end{pmatrix}$$

$$\begin{aligned} \text{Beweis: } (A^T A)^{-1} A^T &= (V \Sigma^+ U^T (U \Sigma V^T))^{-1} (U \Sigma V^T)^T \\ &= (V \Sigma^+ (U^T U \Sigma V^T))^{-1} (V \Sigma^+ U^T) \\ &= V (\Sigma^+ \Sigma)^{-1} V^T V \Sigma^+ U^T \\ &= V (\Sigma^+ \Sigma)^{-1} \Sigma U^T \\ &= V \Sigma^+ U^T = A^+ \quad \square \end{aligned}$$

Satz (Mirsky-Schmidt) $A \in \mathbb{R}^{m \times n}$, $A = U \Sigma V^T$, $\text{rank } A = r \leq \min\{m, n\}$ NumI, VL

$B = \arg \min \{ \|A - B\|_F : \text{rank } B \leq s\}$, B beste rang s Approximation von A

Dann ist $B = U \Sigma_S V^T = \sum_{k=1}^s \tilde{\sigma}_k \underbrace{U_k V_k^T}_{\in \mathbb{R}^{m \times n}}$

$$\|A - B\|_F^2 = \sum_{k=s+1}^r \tilde{\sigma}_k^2$$

$$\Sigma_S = \begin{pmatrix} \tilde{\sigma}_1 & & 0 \\ & \ddots & \\ 0 & & 0_{(n-s) \times (n-s)} \end{pmatrix}$$

Beweis $r = s$, $B = A$, $\text{rang } A = r = s$ ($\|A - B\|_F = 0$)

$$\begin{aligned} 1) \text{ Sei } B = U \Sigma_S V^T \Rightarrow \|A - B\|_F^2 &= \|U \Sigma V^T - U \Sigma_S V^T\|_F^2 \\ \text{rang } B \leq s &= \|U (\Sigma - \Sigma_S) V^T\|_F^2 = \|\Sigma - \Sigma_S\|_F^2 = \sum_{k=1}^r \tilde{\sigma}_k^2 - \sum_{k=1}^s \tilde{\sigma}_k^2 \\ &= \sum_{k=s+1}^r \tilde{\sigma}_k^2 \geq 0 \\ &= \sum_{k=1}^r \tilde{\sigma}_k^2 - \sum_{k=1}^s \tilde{\sigma}_k^2 = \|U \Sigma V^T\|_F^2 - \sum_{k=1}^s \tilde{\sigma}_k^2 = \|A\|_F^2 \sum_{k=1}^s \tilde{\sigma}_k^2 \quad \otimes \end{aligned}$$

$$\begin{aligned} \text{Lemma 1)} \text{ Seien } x_k \in \mathbb{R}^m, y_k \in \mathbb{R}^n \text{ bel. } k=1, \dots, s, \| \sum_{k=1}^s x_k y_k^T - A \|_F^2 \\ \geq \|A\|_F^2 - \sum_{k=1}^s \tilde{\sigma}_k^2 \end{aligned}$$

$$\otimes \Rightarrow \|A - C\|_F^2 \quad \forall C \in \mathbb{R}^{m \times n}, \text{rang } C \leq s$$

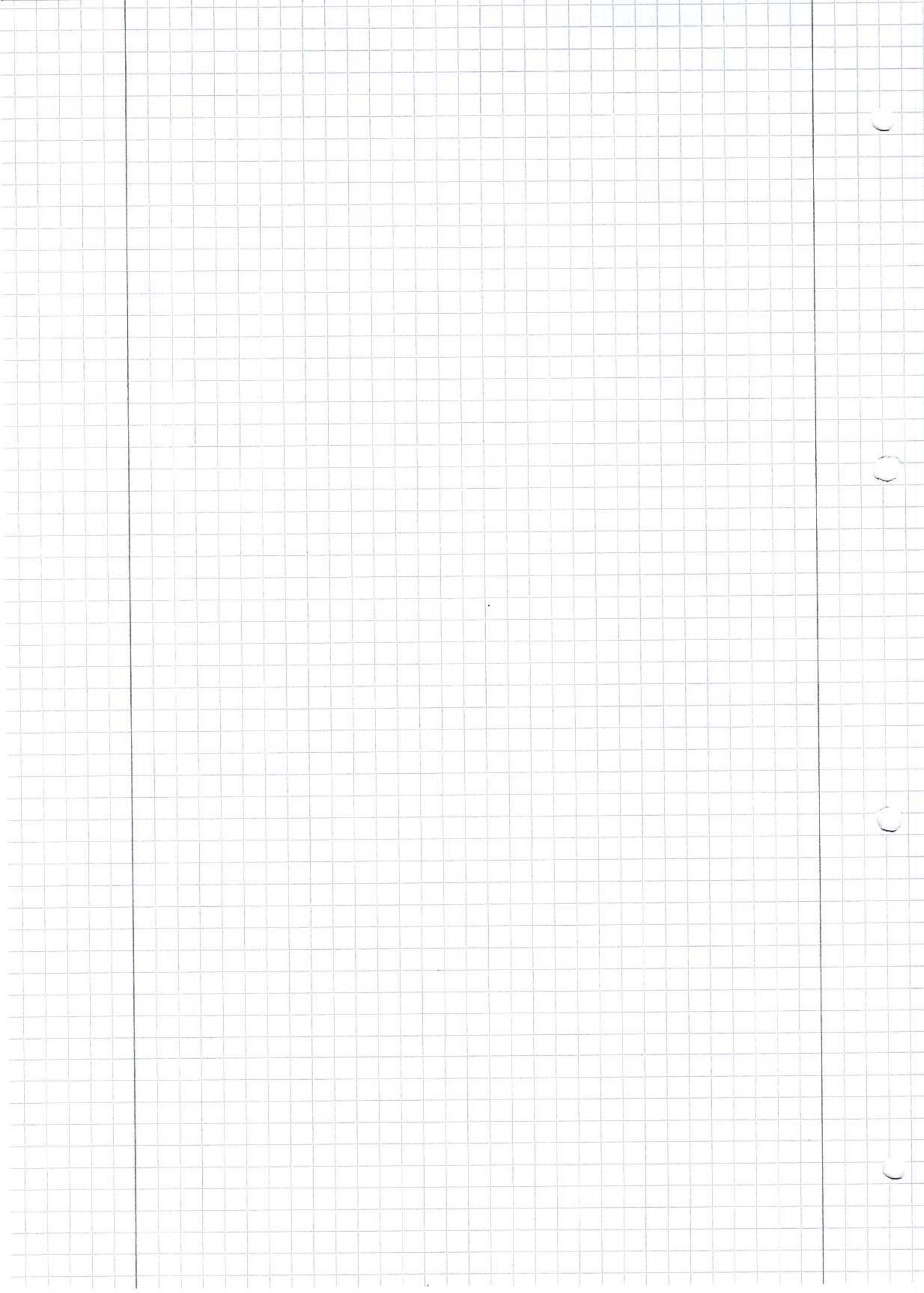
Beweis Lemma 1):

OBdA seien x_k orthogonal, d.h. $x_k^T x_j = \tilde{\sigma}_{kj}$

Andernfalls orthogonalisieren wir $\tilde{x}_k = \sum_{k'=1}^s x_{k'} S_{k'k}^{-1}$; $S \in \mathbb{R}^{s \times s}$ reg. $\Rightarrow x_k = \sum_{k'=1}^s \tilde{x}_{k'} S_{k'k}$

$$\sum_{k=1}^s x_k y_k^T = \sum_{k', k=1}^s \tilde{x}_{k'} S_{k'k}^{-1} y_k^T, \tilde{x}_k \text{ ON}$$

$\Rightarrow \Delta$ Beweis (ganz) siehe Bärwolff-Skript S.42
Satz 3.16



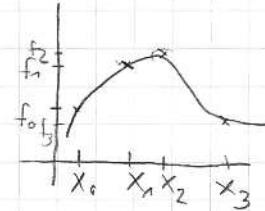
3. Interpolation

• kompl. Funktionen z.B. $f(x) = \sin(\ln(x^2+1))$

• (x_k, f_k) , $k = 0, 1, \dots, n$

• $[0, 10], x_k = \frac{k}{10} \cdot 10, k = 0, \dots, 10$

$$f_k = f(x_k), (x_k, f_k), k = 0, \dots, 10$$



3.1. Das allgemeine Interpolations-Problem

Def: 3.1. (lin. Funktional)

Sei \mathcal{F} ein lin. Raum (z.B. $C[a,b]$, $C^\infty(\Omega)$, $L_2(\Omega)$...).

Eine Abb. $\mu: \mathcal{F} \rightarrow \mathbb{R}$ heißt lineares Funktional, falls

$$\mu(\alpha_1 f_1 + \alpha_2 f_2) = \alpha_1 \mu(f_1) + \alpha_2 \mu(f_2), \quad \alpha_1, \alpha_2 \in \mathbb{R}, \quad f_1, f_2 \in \mathcal{F} \quad (3.1)$$

Beispiel • $\mathcal{F} = C[a,b], x_0 \in [a,b], \mu(f) = f(x_0)$

• $\mathcal{F} = C^1[a,b], x_0 \in [a,b], \mu(f) = f'(x_0)$

• $\mathcal{F} = R[a,b] \text{, } \mu(f) = \int_a^b f(x) dx$ Raum der Riemann-integrierbaren Fkt.

Allgemeines Interpolationsproblem

Sei \mathcal{F} ein Funktionenraum und G_n ein $(n+1)$ -dim

Teilraum von \mathcal{F} . Weiter seien lineare Funktionale

μ_0, \dots, μ_n auf G_n : $\left\{ \begin{array}{l} \text{zu } f \in \mathcal{F} \text{ finde } g_n \in G_n \text{ mit} \\ \mu_j(g_n) = \mu_j(f), \quad j = 0, \dots, n \end{array} \right. \quad (3.2)$

Lemma 3.1.1. Existenz und Eindeutigkeit

Die allgemeine Interpolationsaufgabe (3.2) hat genau dann für jeden $f \in \mathcal{F}$ eine eindeutige Lösung g_n , wenn

$$\det \left[(\mu_i | \varphi_j) \right]_{i,j=0}^n \neq 0 \quad \text{für eine (und damit jede) Basis } [\varphi_0, \dots, \varphi_n] \text{ von } G_n$$

Beweis Sei $g_n \in G_n$, dann existieren eind. bestimmte Koeffizienten $\{\alpha_0, \dots, \alpha_n\}$ mit $g_n = \sum_{j=0}^n \alpha_j \varphi_j$, also bedeutet (3.2)

$$\mu_i(g_n) = \sum_{j=0}^n \alpha_j \mu_i(\varphi_j) = (G\alpha)_i = \mu_i(f) = b_i, \quad G = (\mu_i(\varphi_j))_{i,j=0}^n$$

$Gx = b$ eindeutig lösbares LGS

$$\alpha = (\alpha_j)_{j=0}^n, \quad b = (b_j)_{j=0}^n$$

Bemerkung: $G_n = P_n \rightarrow$ Polynominterpolation

$G_{2n} = \{ F(x) = \frac{P(x)}{Q(x)}, P, Q \in P_n, Q \neq 0 \} \rightarrow$ rationale Interpolation

3.2 Polynominterpolation

3.2.1 Lagrange-Interpolation

$$(x_k, f_k), k=0, 1, \dots, n$$

Gesucht ist ein Polynom $P \in P_n$ vom Grad $p \leq n$, d.h.

$$P(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0, a_i \in \mathbb{R}, i=0, \dots, n$$

$$P(x_i) = f_i, i=0, \dots, n \quad (3.3)$$

Basis $\varphi_j = x^j, j=0, \dots, n$ \leftarrow Monom-Basis

Für (3.3.) erhalten wir $\mu_i(f) = f(x_i), i=0, \dots, n$

$$\mu_i(\varphi_j) = x_i^j \rightsquigarrow \{\mu_i(\varphi_j)\}_{i,j=0}^n \text{ ist die Vandermonde-Matrix}$$

$$\begin{pmatrix} 1 & x_0 & x_0^2 & x_0^3 & \dots & x_0^n \\ 1 & x_1 & x_1^2 & x_1^3 & \dots & x_1^n \\ 1 & x_2 & x_2^2 & x_2^3 & \dots & x_2^n \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & x_n^3 & \dots & x_n^n \end{pmatrix}$$

Definition 3.2.1.

$$L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{(x-x_j)}{(x_i-x_j)}, i=0, \dots, n \text{ Lagrange Basispolynome, (3.4)}$$

paarweise verschiedene x_0, \dots, x_n

Bemerkung $L_i^{(n)}(x_k) = S_{i,k} = \begin{cases} 1 & \text{für } i=k \\ 0 & \text{für } i \neq k \end{cases} \quad (3.5)$

$\rightsquigarrow L_i^{(n)}$ sind linear unabhängig

\rightsquigarrow Basis von P_n

Definition: Lagrange-Darstellung

Das Polynom $P(x) = \sum_{i=0}^n f_i L_i^{(n)}(x)$ wird Lagrange-Darstellung des Interpolationspolynoms zu den Stützstellen $(x_0, f_0), \dots, (x_n, f_n)$ genannt

$$P(x_k) = \sum_{i=0}^n f_i \underbrace{L_i^{(n)}(x_k)}_{S_{i,k}} = \sum_{i=0}^n f_i S_{i,k} = f_k, k=0, \dots, n$$

Eindeutigkeit: $P, Q \in P_n, P(x_k) = Q(x_k) = f_k \quad (k=0, \dots, n)$

$D(x) = P(x) - Q(x) \in P_n, D$ hat $n+1$ paarweise verschiedene NS (Nullstellen), $x_0, \dots, x_n \rightsquigarrow D \equiv 0$ Nullpolynom

Beispiel	x_k	f_k	$n=2$
	0	2	
	1	4	
	2	7	

$$L_0^{(2)}(x) = \frac{x-1}{0-1} \cdot \frac{x-2}{0-2} = \frac{1}{2}(x-1)(x-2)$$

$$L_1^{(2)}(x) = \frac{x-0}{1-0} \cdot \frac{x-2}{1-2} = -x(x-2)$$

$$L_2^{(2)}(x) = \frac{x-0}{2-0} \cdot \frac{x-1}{2-1} = \frac{1}{2}x(x-1)$$

$$\rightsquigarrow P_2(x) = \sum_{i=0}^2 f_i L_i^{(2)}(x) = 2 \cdot \frac{1}{2}(x-1)(x-2) + 4(-x(x-2)) + 7 \cdot \frac{1}{2}(x(x-1))$$

Problem (x_{n+1}, f_{n+1}) soll hinzugefügt werden \rightsquigarrow Basis-Polynome müssen neu berechnet werden

3.2.1. Monombasis

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

„schlecht konditioniert, großer Aufwand“

$$\begin{array}{c|ccccc} f_n & & & x & & \\ \hline f_0 & & & x & & \\ & 1 & & & & \\ & & x_0 & & x_0 & \\ & & & x_0 = 1 & x_0 = 2 & \end{array} + f'_0 \text{ ist gegeben, } P_2(x_0) = f_0 = 1$$

$$P_2'(x_0) = f'_0 = 0$$

$$P_2(x_n) = f_n = 1$$

$$P_2(x) = a_0 + a_1 x + a_2 x^2$$

$$P_2'(x) = a_1 + 2a_2 x$$

$$\rightarrow a_0 + a_1 + a_2 = 1$$

$$a_1 + 2a_2 = 0$$

$$a_0 + 2a_1 + 4a_2 = 1$$

$$\rightarrow \begin{pmatrix} 1 & 1 & 1 \\ 0 & 1 & 2 \\ 1 & 2 & 4 \end{pmatrix} \begin{pmatrix} a_0 \\ a_1 \\ a_2 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

3.2.2. Hermite-Interpolation

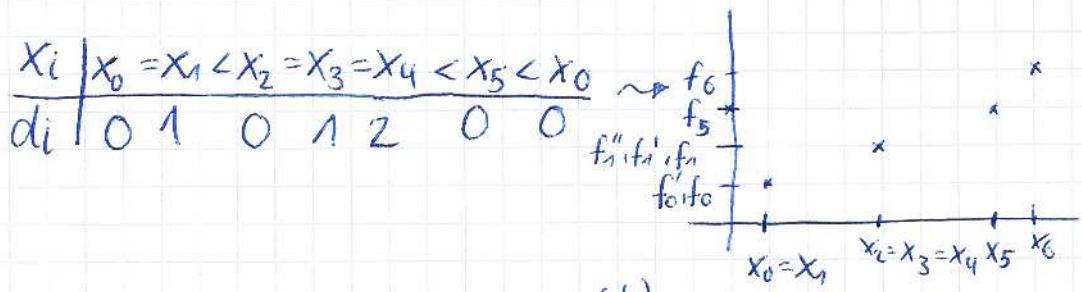
$$a = x_0 \leq x_1 \leq \dots \leq x_n = b$$

\triangleright Stützstellen können auch mehrfach auftreten

$$x_i \rightarrow f(x_i) = f_i, \quad f'(x_i) = f'_i, \quad f^{(k)}(x_i) = f_i^{(k)}$$

$\rightsquigarrow x_i$ soll in der Folge der Stützstellen $(k+1)$ -mal auftreten
gleich Knoten (Stützstellen) nummerieren wir mit der Vielfachheit

$$d_i = \max \{ j \in \mathbb{N} \mid x_i = x_{i-j} \}$$



$$\mu_i : C^n[a,b] \rightarrow \mathbb{R}, \mu_i(f) = f^{(di)}(x_i), i=0, \dots, n$$

Aufgabe der Hermite-Interpolation:

gegeben: μ_i ($i=0, \dots, n$), finde $P \in P_n$ mit $\mu_i(P) = \mu_i(f)$ ($i=0, \dots, n$) (3.8)

Die Lösung $P(f|x_0, \dots, x_n) \in P_n$ heißt Hermite-Interpolierende.

Satz zu jeder Funktion $f \in C^n[a,b]$ und jeder monotonen Folge $a = x_0 \leq x_1 \leq \dots \leq x_n = b$ von Stützstellen (Knoten)

gibt es genau ein Polynom $P \in P_n$, sodass $\mu_i(P) = \mu_i(f)$, $i=0, \dots, n$.

Beispiel: $x_0 = -1, f_0 = 1$. ? $P_2(x) = a_0 + a_1 x + a_2 x^2$

$x_1 = 0, f_1' = 0$	$\begin{array}{ c } \hline -1 \\ \hline \end{array}$	$P_2'(x) = a_1 + 2a_2 x$
$x_2 = 1, f_2 = 3$		

$$P_2(-1) = a_0 - a_1 + a_2 = 1$$

$$P_2'(0) = a_1 = 0$$

$$P_2(1) = a_0 + a_1 + a_2 = 3$$

$$\sim \begin{bmatrix} 1 & -1 & 1 \\ 0 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 3 \end{bmatrix}$$

3.2.3 Auswertung von Polynomen

Ziel $P = P(f|x_0, \dots, x_n)$ soll effizient berechnet werden

Schema von Aitken/Neville

Schema der dividierten Differenzen

x_k , $k=0, \dots, n$, paarweise verschieden: Basis von $P_n \rightsquigarrow$ Monombasis
Lagrange-Basis
Newton-Basis

Definition: Seien $x_0, x_1, \dots, x_n \in \mathbb{R}$ mit

$w_0 := 1, w_i(x) = \prod_{j=0}^{i-1} (x-x_j)$, $w_i \in P_i$, wir bezeichnen

$\{w_0, \dots, w_n\}$ als Newton-Basis des Polynomraumes P_n

$$\rightsquigarrow P(x) = c_0 w_0 + c_1 w_1 + \dots + c_n w_n$$

$$= c_0 + c_1(x-x_0) + c_2(x-x_0)(x-x_1) + \dots + c_n \underbrace{\prod_{j=0}^{n-1} (x-x_j)}_{w_n} \leq \sum_{i=0}^n c_i w_i(x)$$

$$(x_k, f_k), k=0, \dots, n \quad P(x_k) = f_k, k=0, \dots, n$$

$$P(x) = c_0 + c_1(x-x_0)$$

$$P(x_0) = c_0 \quad = f(x_0) = f_0$$

NumI, VL

$$P(x_1) = c_0 + c_1(x - x_0) \quad = f_1$$

$$P(x_2) = c_0 + c_1(x_1 - x_0) + c_2(x_2 - x_0)(x_2 - x_1) \quad = f_2$$

$$\vdots$$

$$P(x_n) = c_0 + c_1 w_1(x_n) + \dots + c_n w_n(x_n) = f_n$$

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 1 & (x_1 - x_0) & \dots & 0 \\ 1 & (x_2 - x_0) & (x_2 - x_1)(x_2 - x_0) & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & \dots & \dots & \dots & w_n(x_n) \end{pmatrix} \begin{pmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_n \end{pmatrix}$$

Basen von $P_n := \{1, x_1, x_1^2, \dots, x_1^n\}$, $P(x) = a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n$

$$Va = f, \quad a = \begin{pmatrix} a_0 \\ \vdots \\ a_n \end{pmatrix}$$

$$*) \left\{ L_i^{(n)}(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}, \quad i = 0, \dots, n \right\}, \quad P(x) = \sum_{k=0}^n b_k L_k^{(n)}(x) \rightsquigarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} b_0 \\ b_n \end{pmatrix} = \begin{pmatrix} f_0 \\ f_n \end{pmatrix}$$

$$b_k = f_k$$

$$*) \{ w_k(x), k = 0, \dots, n \}, \quad w_k(x) = \prod_{j=0}^{k-1} (x - x_j), \quad w_0(x) = 1,$$

$$P(x) = \sum_{j=0}^n c_j w_j(x) \rightarrow \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} c_0 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} f_0 \\ f_n \end{pmatrix}$$

$$P_n \rightarrow P_{n+1}, \quad w_{n+1} = \prod_{j=0}^n (x - x_j)$$

Lemma (Artken)

Für das Interpolationspolynom $P = P(f | x_0, \dots, x_n)$ gibt die

$$\text{Rekursionsformel } P(f | x_0, \dots, x_n)(x) = \underbrace{(x_0 - x)}_{x_0 - x_1} \underbrace{P(f | x_1, \dots, x_n)(x)}_{\Phi(x) \in P_n} - (x_n - x) P(f | x_0, \dots, x_{n-1})(x)$$

insbesondere gilt $P(f | x_k) = f_k$

$$\text{Beweis } \Phi(x_i) = \frac{(x_0 - x_i) f_i - (x_n - x_i) f_i}{x_0 - x_n} = f_i, \quad \Phi(x_0) = \dots = f_0$$

$$\Phi(x_n) = \dots = f_n$$

Algorithmus $P(f | x_i) = f_i, \quad i = 0, \dots, n$

$$x \text{ fest } P_{i,k} := P(f | x_{i-k+1}, \dots, x_i)(x), \quad i \geq k \quad \left. \begin{array}{l} P_{n,n} = (x - x_n) P_{n,n-1} - (x_n - x) P_{n-1,n-1} \\ P_{i,n} = (x - x_n) P_{i,n-1} - (x_n - x) P_{i-1,n-1} \end{array} \right\} (*)$$

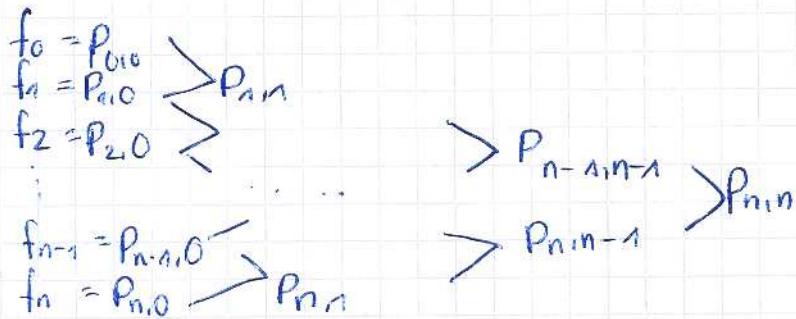
$$(*) \quad P_{n,n} = \frac{(x - x_n) P_{n,n-1} - (x_n - x) P_{n-1,n-1}}{x_0 - x_n}$$

$$\text{allgemeiner: } P_{i,k}, \quad i \geq k: \quad P_{i,k} = \frac{(x_{i-k} - x) P_{i,k-1} - (x_i - x) P_{i-1,k-1}}{x_{i-k} - x_i}$$

$$= P_{i,k-1} + \frac{x_i - x}{x_{i-k} - x_i} (P_{i,k-1} - P_{i-1,k-1})$$

$$1) P_{i,0} = f_i \quad i=0, \dots, n$$

$$2) P_{i,k} = P_{i,k-1} + \frac{x - x_i}{x_i - x_{i-k}} (P_{i,k-1} - P_{i-n,k})$$



Newton-Darstellung des Interpolationspolynoms und dividierte Differenzen

$$\underbrace{P(f|x_0, \dots, x_n)(x)}_{\in \mathbb{P}_n} = \underbrace{P(f|x_0, \dots, x_{n-1})(x)}_{\in \mathbb{P}_{n-1}} + \underbrace{Q_n(x)}_{\in \mathbb{P}_n},$$

$$Q_n(x_i) = P(f|x_0, \dots, x_n)(x_i) - P(f|x_0, \dots, x_{n-1})(x_i) = 0$$

$$\rightarrow Q_n(x) = a_n(x-x_0)(x-x_1) \cdots (x-x_{n-1}) \quad \text{d.h. vor } x^n$$

$$a_n = \frac{f_n - P(f|x_0, \dots, x_{n-1})(x_n)}{(x_n - x_0) \cdots (x_n - x_{n-1})}$$

→ führender Koeffizient von $P(f|x_0, \dots, x_n)$, d.h. vor x^n

Definition (n-te dividierte Differenz)

Der führende Koeffizient a_n des Interpolationspolynoms $P(f|x_0, \dots, x_n)(x)$ von f zu den Knoten $x_0 \leq x_1 \leq \dots \leq x_n$ heißt n -te dividierte Differenz von f in x_0, \dots, x_n und wird mit

$$[x_0 \dots x_n] f = a_n \quad ([x_0 \dots x_n] = a_n)$$

$$2. \text{ dividierte Differenz: } P(f|x_0, x_1, x_2) = a_2 x^2 + a_1 x + a_0$$

$$\text{Folgerung: } P(x) = \sum_{j=0}^n c_j w_j(x) = \sum_{j=0}^n a_n x^j \rightarrow c_n = a_n$$

Satz (Newton'sche Interpolation Sformel)

$$P(x) = \sum_{i=0}^n [x_0 \dots x_i] f w_i(x) \quad (= P(f|x_0 \dots x_n))$$

$$\text{Gilt } f \in C^{n+1}(\mathbb{R}) \rightarrow f(x) = P(x) + [x_0 \dots x_n] f \cdot w_{n+1}(x)$$

Beweis Induktion über $n \in \mathbb{N}$

Num I, VL

$$n=0 \rightsquigarrow P(x) = \sum_{i=0}^0 [x_0] f \omega_i = [x_0] f \omega_0 = [x_0] f = f_0$$

$$n>0 \quad P_{n-1} := P(f|x_0 - x_{n-1}) = \sum_{i=0}^{n-1} [x_0 - x_i] f \omega_i$$

$$P_n(x) = [x_0 - x_n] f x^n + a_{n-1} x^{n-1} + \dots + a_0 = [x_0 - x_n] f \omega_n(x)$$

$$+ \underbrace{Q_{n-1}(x)}_{\in P_{n-1}} = \textcircled{*}$$

$$\text{Es gilt } \omega_n(x_i) = 0, i=0, \dots, n-1 \rightsquigarrow Q_{n-1} = P_n - [x_0 - x_n] f \omega_n$$

verfülle die Interpolationsaufgabe
für x_0, \dots, x_{n-1}

$$\textcircled{*} = [x_0 - x_n] f \omega_n(x) + \sum_{i=0}^{n-1} [x_0 - x_i] f \omega_i(x) = \sum_{i=0}^n [x_0 - x_i] f \omega_i(x)$$

Lemma (Eigenschaften der dividiersten Differenzen, rekursive Berechnung)

(i) Für $x_i \neq x_k$ gilt $[x_0 - x_n] = \frac{[x_0 - \hat{x}_i - x_n] f - [x_0 - \hat{x}_k - x_n] f}{x_k - x_i}$

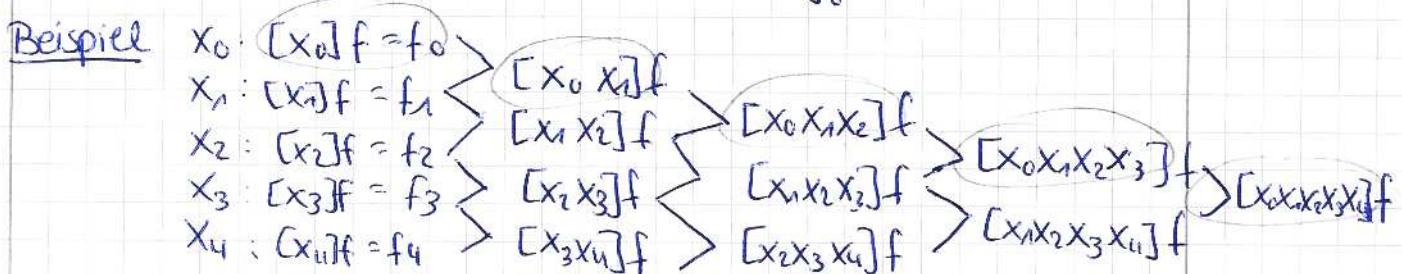
(ii) Für zusammenfallende Knoten

$$x_0 = x_1 = \dots = x_n \text{ gilt } \underbrace{[x_0 - x_0]}_{(n+1)\text{ mal}} = \frac{f^{(n)}(x_0)}{n!} \rightarrow [x_0, \dots, \hat{x}_i, \dots, x_n] f x^{n-1} + P_{n-2}, [x_0 - x_i - x_{i+1} - \dots - x_n] f$$

Beweis:

(i) $P(f|x_0 - x_n) = \frac{(x_i - x)}{x_i - x_0} P(f|x_0 - \hat{x}_i - x_n) - \frac{(x_k - x)}{x_k - x_0} P(f|x_0 - \hat{x}_k - x_n)$

(ii) $P(f|x_0 - x_i)(x) = \sum_{j=0}^n \frac{(x-x_0)^j}{j!} f^{(j)}(x_0) \rightsquigarrow \frac{f^{(j)}(x_0)}{j!} = [x_0 - x_0] f$



$$P_4(x) = [x_0] f + [x_0 x_1] f \omega_1(x) + [x_0 x_1 x_2] f \omega_2(x) + [x_0 x_1 x_2 x_3] f \omega_3(x) + [x_0 x_1 x_2 x_3 x_4] f \omega_4(x)$$

$$+ [x_0 x_1 x_2 x_3] f \omega_3(x) + [x_0 x_1 x_2 x_3 x_4] f \omega_4(x)$$

Beispiel

x_0	x	f
x_1	0	1
x_2	$\frac{3}{2}$	2
x_3	$\frac{5}{2}$	2
	$\frac{9}{2}$	3

$$\begin{array}{c|ccccc}
 & 0 & 1 & & & \\
 \begin{array}{c} 3/2 \\ 5/2 \\ 9/2 \end{array} & \left. \begin{array}{c} 2 \\ 2 \\ 3 \end{array} \right\} 0 & \left. \begin{array}{c} 2-1 \\ 3/2-0 \\ 5/2-0 \end{array} \right\} = \frac{2}{3} & \left. \begin{array}{c} 0-2 \\ 5/2-0 \\ -1/3 \end{array} \right\} = -\frac{4}{15} & \left. \begin{array}{c} 1 \\ 1/6 \\ 1/45 \end{array} \right\}
 \end{array}$$

$$\rightarrow P_3(x) = \underbrace{1}_{w_0(x)} + \underbrace{\frac{2}{3}(x-0)}_{w_1(x)} + \underbrace{\left(-\frac{4}{15}\right)(x-0)(x-\frac{3}{2})}_{w_2(x)} + \underbrace{\frac{1}{45}(x-0)(x-\frac{3}{2})(x-\frac{5}{2})}_{w_3(x)}$$

Polynom-Interpolation:

$I = [a, b]$, $a \leq x_0 \leq \dots \leq x_n \leq b$, x_i : $(n+1)$ -Knotenstellen pw.

verschieden. Gegeben f_k mit $f_k \approx f(x_k)$, $k=0, \dots, n$

Gesucht: $p_n \in \mathbb{P}_n$ ($\deg p_n \leq n$) mit $p_n(x_k) = f(x_k)$ $\forall k=0, \dots, n$

Lagrange-Basen $x \mapsto l_{k,j}$ $l_k(x_j) = \begin{cases} 1 & x_j = x_k \\ 0 & x_j \neq x_k \end{cases} \quad j, k = 0, \dots, n$

$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \frac{x - x_j}{x_k - x_j}, \quad p_n(x) = \sum_{k=0}^n f_k l_k(x)$$

$$\Rightarrow p_n(x_j) = \sum_{k=0}^n f_k l_k(x_j) = f_j \quad \forall j = 0, \dots, n \quad (\text{Lagrange-Darstellung})$$

Fehlerabschätzung zur Polynomapproximation

Satz Sei $f \in C^{n+1}(I)$, $x \in I$ bel.

Sei $p_n = P(f(x_0, \dots, x_n))$ das Interpolationspolynom zu $f_k := f(x_k)$ $k=0, \dots, n$, dann $\exists \xi \in \text{conv}\{x_0, \dots, x_n\} \subseteq I$

$$|f(x) - p_n(x)| \leq \frac{f^{(n+1)}(\xi)}{(n+1)!} \cdot \prod_{k=0}^n (x - x_k)$$

Beweis Sei $x' \in I$ bel. und fest und $\alpha \in \mathbb{R}$, sodass

$$F(x) = f(x) - p(x) - \alpha \prod_{k=0}^n (x - x_k) \quad F(x') = 0$$

d.h. $\alpha = \frac{f(x') - p_n(x')}{\prod_{k=0}^n (x - x_k)}$

$$\Rightarrow F \in C^{n+1}(I) \text{ und } F(x_k) = \underbrace{f(x_k) - p(x_k)}_{=0} - \alpha \prod_{j=0}^n (x_k - x_j) = 0$$

für $k = 0, \dots, n$, $F(x') = 0$

Falls $x' < x_k$ für ein k , ist die Aussage trivial ($0=0$)

$\Rightarrow F$ hat mindestens $n+2$ Nullstellen

Anwendung des Lemma von Rolle \Rightarrow es ex. eine Zwischenstelle

$$\xi \in \text{conv}\{x_0, \dots, x_n\} \text{ (konvexe Hülle)} \text{ mit } F^{(n+1)}(\xi) = 0 = f^{(n+1)}(\xi) + \underbrace{p_n^{(n+1)}(\xi) - \alpha(n+1)!}_{=0}$$

$$\Rightarrow \alpha = \frac{f^{(n+1)}(\xi)}{(n+1)!}$$

Grenzen der Polynominterpolation: Bsp: $I = [-5, 5]$, $f(x) = \frac{1}{1+x^2} \in C^\infty(\mathbb{R})$

$$x_k = -5 + \frac{10k}{n} \quad k=0, \dots, n \quad (\text{äquidistant})$$

$p_n(x)$ divergieren in der Nähe des Randes für $n \rightarrow \infty$

Folgerung: Äquidistante Knotenwahl ist für große n i.A. ungünstig.

Tschebychev-Knoten:

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



ohne Beweis: Seien x_k T-knoten und \tilde{x}_k beliebige Knoten $k=0, \dots, n$

Dann gilt: $\prod_{k=0}^n |x - x_k| \leq \prod_{k=0}^n |\tilde{x} - \tilde{x}_k|$. T-Knoten liefern quasi

eine optimale Knoten und es gilt $\|f - p_n\|_{\infty} \leq C_n \cdot \inf_{q \in \Pi_n} \|f - q\|_{\infty}$ mit

$$\inf_{q \in \Pi_n} \|f - q\|_{\infty} \text{ mit } C_n = 1 + \sum_{j=0}^n \prod_{\substack{k=0 \\ k \neq j}}^n \frac{(x_j - x_k)}{|x_j - x_k|}$$

3.5. Trigonometrische Interpolation

3.5. Diskrete Fourier-Transform.

$f: \mathbb{R} \rightarrow \mathbb{C}$ 1-periodisch (d.h. $f(x+k) = f(x) \forall x \in \mathbb{R}, k \in \mathbb{Z}$)

$$x \mapsto \psi_k(x) = c_k e^{ikx 2\pi}, k \in \mathbb{Z} \quad (\text{für } \varphi \in \mathbb{R} \text{ gilt } e^{i\varphi} = \cos \varphi + i \sin \varphi)$$

$$= c_k (e^{ix 2\pi})^k$$

Interpolationsaufgabe geg: x_0, \dots, x_{n-1} ($x_n = x_0 \in [0, 1]$)

$$\text{ges: } p_n(x) = \sum_{k=0}^{n-1} c_k \psi_k(x) = \sum_{k=0}^{n-1} c_k e^{2\pi i k x} \text{ mit } p_n(x_j) = f_j, j=0, \dots, n-1$$

$f_i = f(x_i)$ d.h. gerucht sind $x_k \in \mathbb{C}, k=0, \dots, n-1$

Seien $x_k = \frac{k}{n}, x_0, \dots, x_{n-1}$ lösen n Unbekannte, n -Gleichungen

$$\sum_{j=0}^{n-1} c_j e^{2\pi i j \frac{k}{n}} = f_k, \forall k=0, \dots, n-1 \quad (*)$$

Sei $w_n = e^{2\pi i \frac{k}{n}}$, für $k=0, \dots, n-1$, w_n^k n -te Einheitswurzel

$$w_n^0 = 1, w_n^k = 1$$

Lemma: Es gilt für $j=1, \dots, n-1$

$$\sum_{k=0}^{n-1} (w_n^j)^k = 0; \text{ Für } j=0: \sum_{k=0}^{n-1} (w_n^0)^k = \sum_{k=0}^{n-1} 1 = n$$

Beweis $j=0$ siehe oben

$j \neq 0$: Sei $q = w_n^j$ ($\Rightarrow |q|=1$)

$$\sum_{k=0}^{n-1} w_n^{jk} = \sum_{k=0}^{n-1} q^k = \frac{1-q^n}{1-q} = \frac{1-(w_n^j)^n}{1-w_n^j} = \frac{1-e^{2\pi i j \frac{n}{n}}}{1-w_n^j} = \frac{1-1}{1-w_n^j} = 0$$

geom. Reihe

Theorem: Die Lsg des Interpolations-Problems (*) lautet

$$c_j = \frac{1}{n} \sum_{k=0}^{n-1} f_k e^{-2\pi i j \frac{k}{n}} = (\underline{F} f)(j), \underline{F} = \frac{1}{n} (e^{-2\pi i \frac{k}{n}})_{j,k=0, \dots, n-1} \quad f = (f_k)$$

Bemerkung 1) Es gilt $\underline{F}^* = \underline{F}^H = (\bar{F}_{j,k})^T$ (Hermitsch) knj: $\underline{F}^H = \underline{E}^{-1} = (e^{2\pi i j \frac{k}{n}})_{j,k=0, \dots, n-1}$

$$2) p_n(x_k) = \sum_{j=0}^{n-1} c_j e^{2\pi i j \frac{k}{n}} = \sum_{j=0}^{n-1} \frac{1}{n} \sum_{k=0}^{n-1} f_k e^{2\pi i j \frac{k}{n}}$$

$$= \sum_{j=0}^{n-1} f_k \left(\frac{1}{n} \sum_{j=0}^{n-1} e^{2\pi i j \frac{(k-j)}{n}} \right) = \sum_{j=0}^{n-1} f_k \delta_{kj} = f_k \quad \forall k=0, \dots, n-1$$

Lemma 5.3.1

Numerische Mathematik I

FFT-Lemma

Die Fouriertransformation

$$DFT : (f_0, \dots, f_{n-1}) \mapsto (d_0, \dots, d_{n-1})$$

ist definiert durch

$$d_j = \frac{1}{n} \sum_{k=0}^{n-1} f_k e^{-2\pi i \frac{k}{n} j}$$

Es fällt auf, dass bei geradem j und geradem n der Term

$$e^{-2\pi i \frac{k}{n} j}$$

sich eine gerade Anzahl an Umdrehungen in der komplexen Ebene um den Nullpunkt dreht. Definiere $n = 2m$ und $j = 2l$. Genauer gesagt gilt dann

$$e^{-2\pi i \frac{k+m}{n} j} = e^{-2\pi i \frac{k+m}{2m} 2l} = e^{-2\pi i \frac{k+m}{m} l} = e^{-2\pi i \frac{k}{m} l} e^{-2\pi i l} = e^{-2\pi i \frac{k}{m} l} = e^{-2\pi i \frac{k}{n} j}.$$

Daraus ergibt sich das folgende Lemma:

Lemma 1. Sei $n = 2m$ und $DFT(f_0, \dots, f_{n-1}) = (d_0, \dots, d_{n-1})$. Dann gilt

$$d_{2l} = \frac{1}{m} \sum_{k=0}^{m-1} \frac{1}{2} (f_k + f_{k+m}) e^{-2\pi i \frac{lk}{m}} = DFT\left(\frac{1}{2}(f_0 + f_m), \dots, \frac{1}{2}(f_{m-1} + f_{2m-1})\right)$$

bzw.

$$d_{2l+1} = \frac{1}{m} \sum_{k=0}^{m-1} \frac{1}{2} (f_k - f_{k+m}) e^{-2\pi i \frac{k}{n}} e^{-2\pi i \frac{lk}{m}} = DFT\left(\frac{1}{2}(f_0 - f_m) e^{-2\pi i \frac{0}{n}}, \dots, \frac{1}{2}(f_{m-1} - f_{2m-1}) e^{-2\pi i \frac{m-1}{n}}\right)$$

Beachte, dass die DFTs auf der rechten Seite nur DFTs auf Vektoren von halber Länge sind. Wir können nun also eine DFT der Länge $2m$ berechnen, indem wir zwei DFTs der Länge m berechnen und dann geeignet zusammenfügen.



Schnelle FT (Fourier Transformation)

$$\omega_n = e^{-2\pi i \frac{k}{n}} = e^{-i \frac{\varphi}{n}}, \quad \varphi = \frac{2\pi k}{n}$$

$$f = (f_k)_{k=0}^{n-1}, \quad x_k = 2\pi \frac{k}{n}, \quad k = 0, \dots, n-1$$

$$d_j = (d_k)_{k=0}^{n-1}$$

$$\sum_{j=0}^{n-1} d_j e^{2\pi i \frac{jk}{n}} = f_k, \quad \forall k = 0, \dots, n-1$$

$$\sum_{j=0}^{n-1} d_j \omega_n^{-jk} = f_k, \quad f \text{ gegeben, dann}$$

$$\Rightarrow d_j = \frac{1}{n} \sum_{k=0}^{n-1} f_k \omega_n^{jk}, \quad d_j = \underline{F} f, \quad \underline{F} = \frac{1}{n} (w_n^{jk})_{j,k=0}^{n-1}$$

Lemma: Sei $n = 2m$, $\omega_m = e^{-2\pi i \frac{2}{n}} = \omega_n^2$ und $\ell = 0, \dots, m-1$

dann gelten $d_{2\ell} = \frac{1}{m} \sum_{k=0}^{m-1} (f_{k,\ell} + f_{k+m,\ell}) \omega_m^{\ell k}$ (1)

$$e^{2\pi i} = 1 \quad d_{2\ell+1} = \frac{1}{m} \sum_{k=0}^{m-1} (f_{k,\ell} - \omega_n^k f_{k+m,\ell}) \omega_m^{\ell k} \quad (2)$$

Beweis 1) $j = 2\ell$

$$(w_n^2)^{\ell k} = w_n^{2\ell(k+m)} = (w_n^2)^{\ell k} w_n^{2\ell m} \\ = w_m^{\ell k} w_m^{\ell m} = w_m^{\ell k}$$

$$w_n^2 = \omega_m, \quad w_m^{\ell m} = (\omega_m^{\ell})^m = 1, \quad \ell = 1$$

$$2) j = 2\ell + 1 \quad w_n^{(2\ell+1)(k+m)} = w_n^{(2\ell k + 2\ell m + k + m)} \\ = w_m^{\ell k} w_n^{\ell m} w_n^k w_n^m \\ = w_m^{\ell k} \cdot w_n^k (-1) \quad (\text{Dies lässt sich rekursiv fortsetzen})$$

$$w_n^m = e^{-2\pi i \frac{m}{2m}} = e^{-\pi i} = -1$$

Lemma: Sei $n = 2^P$, $n \in \mathbb{N}$ (2-er Potenz)

Dann gilt für den Aufwand (#arithmetisch. Operationen)

$$T(n) = T(2^P) \leq P 2^{P+1} = 2P 2^P = 2n \log_2 n$$

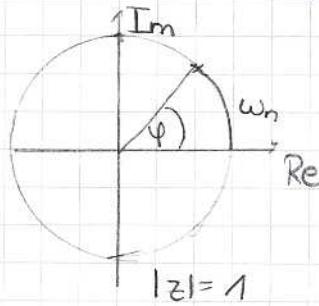
Beweis: In jedem Schritt (Lemma 3.5.2) gilt

$$T(n) \leq 2n + 2T\left(\frac{n}{2}\right)$$

Vollständige Induktion: $p = 1 \quad T(2) \leq 2 + 2 = 4 \leq 1 \cdot 2^2 = 4$

$$p \Rightarrow p+1$$

$$n = 2^{P+1} \quad T(2^{P+1}) \leq 2 \cdot 2^{P+1} + 2(2^P) 2^{P+1} = 2(1+2) 2^{P+1} \quad \blacksquare$$



4. Numerische Integration

4.1 Aufgabenstellung

Intervall $[a, b]$, Fkt $f: [a, b] \rightarrow \mathbb{R}$ unbekannt

$$I(f) = \int_a^b f(x) dx \approx I_h(f) \quad \text{Quadraturformel (*)}$$

$$\text{approximieren: } I_h(f) = \sum_{k=0}^n f(x_k) w_k$$

$x_k \in [a, b]$ Knotenstellen (pw. verschieden)

Gewichte: $w_k \in \mathbb{R}$ ((w_k, x_k) bestimmen die Quadraturformel)

Idee + Bsp: Sei $0 = t_0 < t_1 < t_2 < \dots < t_m = b$

$$\Rightarrow [a, b] = \bigcup_{j=0}^{m-1} [t_j, t_{j+1}], \text{ und } I(f) = \int_a^b f(x) dx = \sum_{j=0}^{m-1} \int_{t_j}^{t_{j+1}} f(x) dx$$

Zerlegung von $[a, b]$

Seien $t_{j+1} - t_j =: h$, äquidistant (Bsp.)

1) Rechteckregel: $I_h(f) = \sum_{j=0}^{m-1} h f(t_j)$

2) Trapezregel auf $[t_j, t_{j+1}]$

$$\tilde{I}_h(f) = h \frac{1}{2} (f(t_j) + f(t_{j+1}))$$

$$\text{zusammengesetzt } I_h(f) = \frac{h}{2} (f(a) + f(b)) + h \sum_{j=1}^{m-1} f(t_j)$$

3) Mittelpunktsregel

$$\int_{t_j}^{t_{j+1}} f(x) dx \approx \tilde{I}_h(f) = h f\left(\frac{t_j + t_{j+1}}{2}\right)$$

4.2 Newton-Cotes-Formeln

Intervall $[c, d] := [t_j, t_{j+1}]$ Teilintervall

Quadraturknoten $x_0, \dots, x_n \in [c, d]$ gewählt (Newton-Cotes) \Leftarrow äquidistant

Seien (Stützwerte) $f_k := f(x_k)$ und $x \mapsto p_n(x) = p_n(f)(x_0, \dots, x_n)(x) \in \mathcal{P}_n$

das zugehörige Interpolationspolynom

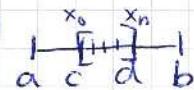
$$I(f) = \int_c^d f(x) dx \approx \int_c^d p_n(f)(x_0, \dots, x_n)(x) dx = \int_c^d p(x) dx = I_h(f)$$

Definition 4.1. Eine Q-Formel (*) heißt exakt von der Ordnung d ,

falls $\int_c^d p(x) dx = I_h(f) \forall p \in \mathbb{P}_d$ Polynom. vom Grad $\leq d$

Numerische Integration

d



$$\int_c^d f(t) dt = I(f) \quad h := [c, d]$$

$x_0 = c < x_1 < \dots < x_n = d$, Knoten, Fkt. weite $f(x_k) = f_k$

Z.B. Newton-Cotes Regeln, $x_k = kh + c = \frac{k(d-c)}{n} + c$

"äquidistante Knoten"

$$I_h(f) := \sum_{k=0}^n p_n(x_k) dx, \quad p_n(x) = P_n(x_0, \dots, x_n | f)(x) \in \Pi_n$$

$$l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^{n-1} \frac{x - x_j}{x_k - x_j}, \text{ Lagrange-Polynome}$$

$$\Rightarrow p_n(x) = \sum_{k=0}^n f(x_k) l_k(x) \Rightarrow I_h(f) = \int_c^d p_n(x) dx = \sum_{k=0}^n f(x_k) \underbrace{\int_c^d l_k(x) dx}_{=: w_k}$$

$$\underline{\text{Satz 4.2.}} \quad I_h(f) = \sum_{k=0}^n f(x_k) w_k, \text{ mit } w_k = \int_c^d \left(\prod_{\substack{j=0 \\ j \neq k}}^{n-1} \frac{x - x_j}{x_k - x_j} \right) dx$$

Bemerkung: Die Wahl der Knoten bestimmt die Quadraturgewichte

$$\underline{\text{Satz 4.1:}} \quad \text{Es gilt: } \left| \int_c^d f(x) dx - I_h(f) \right| \leq h \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} \max_{x \in [c, d]} \frac{n}{j=0} \frac{1}{(x-x_j)} \stackrel{n \leq h}{\approx}$$

für $f \in C^{n+1}([c, d])$ die Fehlerabschätzung

$$\leq \frac{h^{n+2}}{(n+1)!} \|f^{(n+1)}\|_\infty$$

$$\begin{aligned} \text{Beweis: } \left| \int_c^d (f - p_n)(x) dx \right| &\leq \int_c^d |f(x) - p_n(x)| dx \leq h \|f - p\|_\infty \\ &\leq h \underbrace{\frac{\|f^{(n+1)}(\xi)\|_\infty}{(n+1)!}}_{\text{Interpolations-}} \prod_{j=0}^n |x - x_j| \leq h \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} \prod_{j=0}^n |x - x_j| \\ &\leq \frac{h^{n+2}}{(n+1)!} \|f^{(n+1)}\|_\infty \end{aligned}$$

Bsp zu Newton-Cotes, wähle $c=0, d=1, ([0, 1]), h=1$

n	w_k	Fehlerabsch.
1	$\frac{1}{2} \quad \frac{1}{2}$	"Trapezregel"
2	$\frac{1}{6} \quad \frac{4}{6} \quad \frac{1}{6}$	"Simpsonregel"
3	$\frac{1}{18} \quad \frac{3}{18} \quad \frac{3}{18} \quad \frac{1}{18}$	$\frac{3}{8} \quad \frac{3}{8} \quad \frac{3}{8} \quad \frac{3}{8}$ - Regel
4	$\frac{7}{90} \quad \frac{32}{90} \quad \frac{12}{90} \quad \frac{32}{90} \quad \frac{7}{90}$	"Milne-Regel"

Bemerkung: Für N-C-Regeln der Ordnung $n \geq 7$ treten auch neg. Qu. gewichte ($w_k < 0$) auf!

⇒ Regeln sind weniger stabil und werden i.A. nicht verwendet

4.3. Romberg-Quadratur und Extrapolation

Zusammengenetzte Trapezregel $h := b-a$

$$\int_a^b f(x) dx = \sum_{i=1}^n \int_{x_{i-1}}^{x_i} f(x) dx \approx \frac{h}{2} (f(a) + f(b)) + h \sum_{i=1}^{n-1} f(x_i) =: T(h)$$

Satz 4.3.1

$$\left| \int_a^b f(x) dx - T(h) + \sum_{k=1}^m \gamma_{2k} h^{2k} \right| \leq C_{2m+2} |b-a| h^{m+2}$$

$$\text{mit } \gamma_{2k} = \frac{B_{2k}}{(2k)!} (f^{(2k+1)}(b) - f^{(2k+1)}(a))$$

B_{2k} „Benedelli-Zahlen“, $b_{2k} = \frac{B_{2k}}{(2k)!}$

$$g(t) = 1 + \frac{1}{2} t + \frac{t}{e^t - 1} = \sum_{k=1}^{\infty} b_{2k} t^{2k} \quad \text{„Potenzreihe“}$$

Beweis: „Euklidische Summenformel“ ("Iterativ")

Bemerkung: 1) Falls f periodisch ist ($m+T=b-a$)

und analytisch. Dann konvergiert die Trapezregel exponentiell.

2) $c_i := \frac{1}{T} \int_0^T f(x) e^{i \frac{2\pi}{T} x} dx \approx T(h)$ „Trapezregel $\hat{=} \text{diskreter FT}$

Trick „Extrapolation“

Es gilt $\int_a^b f(x) dx = \lim_{n \rightarrow 0} T(h)$

Sei $h := h_0 > h_1 > h_2$

Bsp. Romberg $h_k = h_0 2^k (h_0 \leq h_k)$

Wir interpolieren $T(h) \approx P_n(h^2)$, $P_n(x) = P_n(h_0, \dots, h_n) T(x)$, $x = h^2$

Definition 4.2 Sei $h \mapsto T(h)$, $h \in \mathbb{R}_+$ eine Funktion,

diese besitze eine „asymptotische Entwicklung“ für $h \rightarrow 0$ in h^p

der Ordnung $p \cdot m^p$, falls Konstanten $\tau_0, \dots, \tau_m \in \mathbb{R}$ ex. mit:

$$|T(h) - (\tau_{0,p} h^p + \tau_{1,p} h^{2p} + \dots + \tau_{m,p} h^{mp})| = O(h^{(m+1)p}), h \rightarrow 0$$

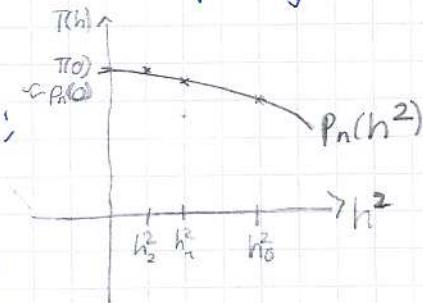
$$\text{d.h. } T(h) = \sum_{k=1}^m \tau_{k,p} h^{kp} + O(h^{(m+1)p}), h \rightarrow 0$$

Aus Satz 4.3.1 folgt, dass die Trapezregel $T(h)$ eine „asy. Entw.“

mit $p=2$ besitzt. ($\tau_0 := \int_a^b f(x) dx$)

Korollar (Extrapolation): $(P_n - \text{Extrapolationspolynom } (p=2))$

$$\left| \int_a^b f(x) dx - P_n(0) \right| \leq Ch^{(m+1)/2} \|f^{(m+1)/2}\|_\infty$$



4.4 Orthogonale Polynome und Gauß-Quadratur

4.4.1 Orthogonale Polynome

Lemma 4.5 Sei $x \mapsto w(x) > 0, x \in \Omega$ und f, g messbar, $d\mu = w(x)dx$
dann definiert $\int_{\Omega} f(x)\overline{g(x)} w(x) dx = \langle f, g \rangle_w$ ein inneres Produkt
im ~~Hilbertraum~~ Hilbertraum $L_2(\Omega, w) \stackrel{\text{a.e.}}{=} \{f : \|f\|_{w, \mu}^2 = \langle f, f \rangle_w\}$ a.e. almost everywhere

Lemma 4.6 Sei $\overline{\Pi}_j = \{p \text{ Polynom } p(x) = x^j + a_{j-1}x^{j-1} \dots a_0, a_j \in \mathbb{K}\} \subset \Pi_j$
Dann ex. genau ein $p \in \overline{\Pi}_j$ mit $\langle p_j, q \rangle_w = 0 \quad \forall q \in \Pi_{j-1} \quad (p_j \perp \Pi_{j-1})$

Beweis: per Induktion

$$\text{IA: } j=0, p_0 = 1$$

IS Seien $p_0, \dots, p_{j-1} \stackrel{\text{v.l., } i=0, \dots, j-1}{\text{gegeben}}$, $\Rightarrow \{x \mapsto x^i\}, p_{j-1}, \dots, p_0$ sind lin. unabh.
(Dimension(Π_j) = $j+1$); ergibt $\langle p_i, p_k \rangle_w = 0 \quad i \neq k, i, k \leq j$
 $\Rightarrow p_j(x) := x^j - \sum_{k=0}^{j-1} p_k(x) \underbrace{\langle p_k, x^j \rangle_w}_{\text{"Gram Schmidt"}}$
 $\Rightarrow p_j \in \Pi_j : \langle p_j, p_k \rangle_w = 0, \forall k=0, \dots, j-1$
 $\Rightarrow \langle p_j, q \rangle_w = 0, \forall q \in \Pi_{j-1} \quad \blacksquare$

Bemerkung Beweis ist konstruktiv. $p_j, j \in \mathbb{N}_0$, "Orthogonale Polynome"

Brgr. Gewicht w

Beispiel 1) $\Omega = [-1, 1], w=1$, Legendre-Polynome

$$p_j(x) = \frac{j!}{(2j)!} ((x^2 - 1)^j), j \in \mathbb{N}_0$$

2) Tschebyschev-Knoten $\Omega = (-1, 1)$

$$w(x) = \frac{1}{\sqrt{1-x^2}}$$

$$p_j(x) = T_j(x) = \cos(j \arccos x) \quad j \in \mathbb{N}_0$$

3) Laguerre-Polynome $\Omega = [0, \infty), w(t) = e^{-t}$

4) "Hermite-Polynome" $\Omega = \mathbb{R}, w(t) = e^{-t^2}$ (Stochastik / Statistik)

Lemma 4.8 Sei p_j gemäß oben $x \mapsto p_j(x), y \in [a, b] \subset \Omega, j \geq 1$

dann ex. genau j einfache Nst $x_i \in [a, b] : (a \leq x_1 < \dots < x_j \leq b)$

Beweis Sei $j \geq 1$. Seien $a \leq x_0 < \dots < x_k \leq b$ alle Nst von p_j .

$$\text{Da } \int_a^b p_j(x) w(x) dx = \langle p_j, 1 \rangle_w = \langle p_j, p_0 \rangle_w = 0 \Rightarrow k \geq 1 \quad (!!!)$$

Annahme $1 \leq k < j$

$$\text{Sei } q(x) = \prod_{i=1}^k (x - x_i) \Rightarrow q \in \Pi_k, (k < j) \quad \textcircled{*}$$

$$\Rightarrow x \mapsto p_j(x) q(x) = r(x) \stackrel{\substack{i=1 \\ i=0}}{\underset{\geq 0}{\prod}} (x - x_i)^2, r(x) \neq 0 \forall x \in [a, b]$$

d.h. Vorzeichen sign $p_j(x) q(x) = \text{const.}, \forall x \in [a, b]$

$$0 < \int_a^b p_j(x) q(x) w(x) dx = \langle p_j, q \rangle_w = 0 \quad \text{Q.E.D.}$$

$$\Rightarrow k = j \quad \blacksquare$$

Lemma 4.9. Sei $t_k \in [a, b], k = 1, \dots, j, a \leq t_1 < t_2 < \dots < t_j \leq b$

$$p_j \in \Pi_j \text{ (wie oben)} \quad \text{und } A = (a_{i,j})_{i,j=1}^n = (p_{i-1}(t_j))_{i,j=1}^n$$

$$\Rightarrow A \in \mathbb{R}^{(j-1) \times n}$$

$\Rightarrow A$ ist regulär

Beweis Annahme $\exists t_k$ mit A singulär d.h. $0 \notin \text{ker}(A)$

$$Au = 0 \text{ bzw. } \sum_{k=1}^j p_{k-1}(t_i) u_k = 0 \quad \forall i = 0, \dots, j-1$$

$$\Rightarrow q(x) = \sum_{k=1}^j p_{k-1}(x) u_k \in \Pi_{j-1}$$

$$\text{und } q(t_i) = \sum_{k=1}^j p_{k-1}(t_i) u_k = 0 \quad \forall i = 0, \dots, j-1$$

$$\Rightarrow q \in \Pi_{j-1} \text{ hat } j \text{ verschiedene NST.} \quad \Rightarrow q \equiv 0 \quad \text{Q.E.D.}$$

$$u_k = 0 \quad \forall k$$

4.4.2 Gauß-Quadratur

Sei $a \leq x_1 < \dots < x_j \leq b$ die Nullstellen von p_j , die Knotenpunkte

für die Gaußquadratur $\int_a^b f(x) w(x) dx = \sum_{k=1}^j p_j(x_k) v_k f$

Satz 4.10 Sei $v = (v_1, \dots, v_j)^T \in \mathbb{R}^j$ die Lösung von

$$Av = e_n \leftarrow P_0, P_0 \geq w, e_n = (1, 0, \dots, 0)^T \in \mathbb{R}^n$$

$$a_{i,j} = p_{j-1}(t_i), \text{ d.h. } \sum_{k=1}^j p_{k-1}(t_i) v_k = \delta_{0,i} \leftarrow P_0, P_0 \geq w, i = 0, \dots, j-1 \quad \text{(*)}$$

$$(i) w_k = v_k > 0$$

$$(ii) \text{ Sei } p \in \Pi_{2j-1} \Rightarrow \int_a^b p(x) w(x) dx = \sum_{k=1}^j p(x_k) v_k, \text{ Gauß-Quadrat.}$$

cst korrekt $\forall p \in \Pi_{2j-1}$ obwohl nur j Punkte gegeben sind

(iii) Es existieren keine j Punkte $0 \leq t_1 < \dots < t_j \leq b$ mit

$$\int_a^b p(x) w(x) dx = \sum_{k=1}^j p(x_k) w_k \quad \forall p \in \Pi_{2j}$$

$$\stackrel{\text{Lemma 4.9.}}{\Rightarrow} \exists ! v \in \mathbb{R}^j$$

\Rightarrow gilt als Gauß-Quadrat.

Gauß-Quadratur $\langle \cdot, \cdot \rangle_W, w(x) > 0$

NumI, VL

$$\langle f, g \rangle_W := \int\limits_a^b f(x)g(x)w(x)dx, \quad p_{i-1}, p_j \in \Pi_j^1 \text{ orthogonale Polynome}$$

zu p_j , Nullstellen $a \leq x_1 < \dots < x_n \leq b \Rightarrow$ Quadraturknoten der Gauß-Quadratur:

Seien $A = (a_{i,k}) = (p_{i-1}(x_k))_{i,k=1,\dots,n} \in \mathbb{R}^{n \times n}$

$\Rightarrow A$ regulär und $w = (w_1, \dots, w_j)^T \in \mathbb{R}^j$ Lösung von

$$Aw = \langle p_0, p_0 \rangle_W (1, 0, 0, \dots, 0)^T$$

w sind Quadraturgewichte (Behauptung)

j -Quadraturknoten, aber Approximationssordnung ist $\sum_{i=1}^{j-1} \alpha_i$

Beweis: ii) Sei $p \in \Pi_{2j-1}$ bel. $\Rightarrow \exists q, r \in \Pi_{j-1}$ mit $r(x) = \sum_{i=0}^{j-1} \alpha_i p_i(x)$

$$p(x) = p_j(x)q(x) + r(x) \quad \forall x \in \mathbb{R} \quad \alpha_i = \frac{\langle p_i, r \rangle_W}{\langle p_i, p_i \rangle_W}$$

$$\begin{aligned} a) \int\limits_a^b p(x)w(x)dx &= \int\limits_a^b q(x)p_j(x)w(x)dx + \int\limits_a^b r(x)w(x)dx \quad \langle p_i, p_i \rangle_W \\ &= \underbrace{\langle q, p_j \rangle_W}_{=0} + \sum_{i=0}^{j-1} \alpha_i \underbrace{\langle p_i, p_0 \rangle_W}_{=0} = 0 + \alpha_0 \underbrace{\langle p_0, p_0 \rangle_W}_{=1} = \alpha_0 \int\limits_a^b w(x)dx \end{aligned}$$

$$\begin{aligned} b) \sum_{k=1}^j p(x_k)w_k &= \sum_{k=1}^j \underbrace{q(x_k)p_j(x_k)}_{=Aw} w_k + \sum_{k=1}^j r(x_k)w_k \\ &= \sum_{i=0}^{j-1} \sum_{k=1}^j \alpha_i \underbrace{\langle p_i, p_j \rangle_W}_{=0} w_k = \sum_{i=0}^{j-1} \alpha_i \underbrace{\langle p_0, p_0 \rangle_W}_{=1} = \alpha_1 \underbrace{\langle p_0, p_0 \rangle_W}_{=1} \Rightarrow \text{ii)} \end{aligned}$$

i) Behauptung, $w_i > 0 \quad \forall i = 1, \dots, j$.

Sei $x \mapsto q(x) := \prod_{i=1}^j (x - x_i)^2 \geq 0, \quad q \in \Pi^{2j-2}$

$$0 < \int\limits_a^b q(x)w(x)dx = \sum_{i=1}^j q(x_i)w_i = \prod_{i=1}^j (x_i - x_i)^2 w_i \Rightarrow w_i > 0 \Rightarrow \text{i)}$$

iii) $\exists p \in \Pi_{2j}$ mit $\int\limits_a^b p(x)w(x)dx \neq \sum_{k=1}^j p(x_k)w_k$.

$$\text{Sei } p(x) = \prod_{i=1}^j (x - x_i)^2 \geq 0, \quad p \in \Pi_{2j} \Rightarrow 0 < \int\limits_a^b p(x)w(x)dx \quad \text{aber } 0 = \sum_{k=0}^j \underbrace{\int\limits_a^b p(x_k)w_k}_{=0}$$

Satz 4.11 $p_j, j \in \mathbb{N}, p_j(x_k) = 0$ usw. gewäß oben

$$\text{Sei } f \in C^0([a, b]) \quad \int\limits_a^b f(x)w(x)dx = \lim_{j \rightarrow \infty} \sum_{k=1}^j f(x_k)w_k$$

Beweis Approximationssatz von Weierstraß: $\forall \varepsilon > 0 \exists j \in \mathbb{N}$ mit $\|f - p_j\|_{\infty} < \varepsilon$

$$\begin{aligned} |I(f) - I_j(f)| &\leq |I(f) - I(p_j)| + |I(p_j) - I_j(f)| \leq \int\limits_a^b |f(x) - p_j(x)|w(x)dx + \sum_{k=1}^j |p_j(x_k) - f(x_k)|w_k \\ &\leq \|f - p_j\|_{\infty} \int\limits_a^b w(x)dx + \|f - p_j\|_{\infty} \sum_{k=1}^j w_k = 2\|f - p_j\|_{\infty} \underbrace{\langle p_0, p_0 \rangle_W}_{=1} = 2\langle f, p_j \rangle_W \leq 2\varepsilon \end{aligned}$$

$\xrightarrow{\varepsilon \rightarrow 0} 0 \quad \blacksquare$

Satz 4.12 $f \in C^{(2n)}([a,b])$, dann ex. $\{x_i\}_{i=1}^n$ mit

$$\left| \int_a^b f(x) w(x) dx - \sum_{k=1}^n f(x_k) w_k \right| \leq \frac{\|f^{(2n)}(\xi)\|}{(2n)!} \langle p_n, p_n \rangle_w$$

Beweis: Hermitte-Interpolation:

2n Knoten $a = x_0 < x_1 < \dots < x_n = b$
Funktion \uparrow Ableitung \downarrow

$$x \in [a, b]: f(x) = p(x) + [x, x_1, x_2, \dots, x_n] f \underbrace{\prod_{i=1}^j (x-x_i)^2}_{=(p_n(x))^2}$$

$$\begin{aligned} \left| \int_a^b f(x) w(x) dx - \sum_{k=1}^n f(x_k) w_k \right| &\leq \left| \int_a^b [x, x_1, x_2, \dots, x_n] f(p_n(x))^2 w(x) dx \right| \\ &\leq \|f^{(2n)}\|_\infty \int_a^b p_n^2(x) w(x) dx = \frac{\|f^{(2n)}\|}{(2n)!} \langle p_n, p_n \rangle_w \end{aligned}$$

Bemerkung: 1) Man kann zeigen, dass für analytische f die Gauß-Quadratur exponentiell konvergiert, d.h. $\exists \beta > 0$: Fehler $\leq C e^{-\beta j}$

Probleme: 1) uneigentliche Integrale bzw. unbeschränkte Integranden,
„graduierter Knotenwahl“ 2) Problem: hochoszillierende Funktion

5. Numerische Lösung nichtlinearer GLS

5.1. Aufgaben: Geg $D \subset \mathbb{R}^n$, $F: D \rightarrow \mathbb{R}^n$

$$x = (x_1, \dots, x_n)^T \in \mathbb{R}^n, F(x) = (F_1(x_1, \dots, x_n), \dots, F_n(x_1, \dots, x_n))^T \in \mathbb{R}^n$$

Gesucht: Nullstellen $x^* \in D$, d.h. $F(x^*) = 0 \in \mathbb{R}^n$

$$\begin{cases} F_1(x_1^*, \dots, x_n^*) = 0 \\ \vdots \\ F_n(x_1^*, \dots, x_n^*) = 0 \end{cases}$$

Beispiel: $n=2, D = \mathbb{R}^2$

$$\begin{cases} F_1(x_1, x_2) = x_1^2 - x_2 + 0,25 = 0 \\ F_2(x_1, x_2) = -x_1 + x_2^2 + 0,25 = 0 \end{cases} \quad x_1 = x_2 = 0,5, \quad x = \begin{pmatrix} 0,5 \\ 0,5 \end{pmatrix}$$

2) $F_1(x_1, x_2) = x_1 \sin x_1 - x_2 = 0$ hat abzählbar viele Lsg.

$$F_2(x_1, x_2) = x_2^2 - x_1 + 1 = 0$$

Sei $n=1$, $f \in C^0([a,b])$ und $f(a) \cdot f(b) < 0 \stackrel{\text{Zws}}{\Rightarrow} \exists x \in (a,b) : f(x) = 0, a_0 = a, b_0 = b$

$$x_0 = \frac{1}{2}(a_0, b_0) \text{ für } k=1, \dots, k_{\max}$$

Falls $f(a_k) \cdot f(x_0) < 0$, $a_{k+1} = a_k, b_{k+1} = x_0$

Andernfalls $b_{k+1} = b_k, a_{k+1} = x_0$

5.2 Fixpunktiteration gegeben

$$f(x^*) = y \Rightarrow F(x^*) = f(x^*) - g = 0, \quad F: D \rightarrow \mathbb{R}^n \quad (D \subset \mathbb{R}^n)$$

Gesucht x^* Nullst. von F

$$\phi(x^*) = x^* + F(x^*) = x^* \text{ „Fixpunktproblem“}$$

oder allgemeiner: Sei für $x \in D$

$$G(x) \in \mathbb{R}^{n \times n} \text{ regular}$$

$$F(x^*) = 0$$

$$\Leftrightarrow x^* - G(x^*) F(x^*) = x^*$$

$$\phi(x^*) :=$$

$$(1) \text{ Fixpunktgleichung } \phi(x^*) = x^*$$

Fixpunktiteration

$$\text{wähle } x_0 \in D, \quad x_1 = \phi(x_0), \quad x_2 = \phi(x_1), \dots, \quad x_{n+1} = \phi(x_n)$$

Hoffnung: $x_n \rightarrow x^*, n \rightarrow \infty?$

Def: 5.1: Sei $(X, \|\cdot\|)$ ein norm Raum (z.B. $X = \mathbb{R}^n, \|\cdot\| = \|\cdot\|_2$)

$D \subset X, \phi: D \rightarrow X$ heißt

1) Lipschitzstetig, falls $L > 0$ ex. mit $\|\phi(x) - \phi(y)\| \leq L \|x - y\| \quad \forall x, y \in D$ (\Rightarrow q.m st.)

2) Kontraktion, falls $0 < L < 1$

Sei $\phi: D \rightarrow \mathbb{R}^n, \phi \in C^1(D), x = (x_1, \dots, x_n)^T \in D$

$$(\text{Fréchet/totale}) \text{ Ableitung } \phi'(x) = \begin{pmatrix} \frac{\partial \phi_1}{\partial x_1} & \cdots & \frac{\partial \phi_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial \phi_n}{\partial x_1} & \cdots & \frac{\partial \phi_n}{\partial x_n} \end{pmatrix} \quad \left. \begin{array}{l} (x_1, \dots, x_n) \in \mathbb{R}^{n \times n} \\ \text{Jacobi-Matrix} \end{array} \right.$$

Satz 5.2 FP Satz von Banach

Sei $(X, \|\cdot\|)$ norm Raum, $D \subset X$ vollständig, $\phi: D \rightarrow D$ (Selbstabb) eine Kontraktion

\Rightarrow 1) es ex genau ein FP $x^* \in D, x^* = \phi(x^*)$

2) $\forall x_0 \in X, x_{n+1} := \phi(x_n)$ konvergiert $x_n \rightarrow x^*, n \rightarrow \infty$ / Kontraktionskonstante

3) a posterior- Abschätzung: $\|x_n - x^*\| \leq \frac{L}{1-L} \|x_n - x_{n-1}\|$

4) a priori- Abschätzung: $\|x_n - x^*\| \leq \frac{L^n}{1-L} \|x_1 - x_0\|$

5) $\|x^{n+1} - x^*\| \leq L \|x^n - x^*\|$

Korollar 5.3 Sei $\phi: D \rightarrow D$, $D \subset \mathbb{R}^n$ abgl., zshg. konvex, $\|\cdot\|$ Norm in \mathbb{R}^n

und $\|\phi'(x)\|$: Operatornorm von $\phi'(x) \in \mathbb{R}^{n \times n}$, $\phi \in C^1(D)$

Falls $\max_{x \in D} \|\phi'(x)\| < 1$, dann ist ϕ eine Kontraktion

mit $L := \max_{x \in D} \|\phi'(x)\|$, \Rightarrow Aussagen von FD von Banach

Beweis Seien $x, y \in D$ bel. $t \mapsto g(t) = x + t(y-x) = (1-t)x + ty$, $0 \leq t \leq 1$

$$\begin{aligned} \|\phi(x) - \phi(y)\| &= \left\| \int_0^1 \phi'(x + t(y-x))(y-x) dt \right\| \leq \int_0^1 \|\phi'(x + t(y-x))\| \|y-x\| dt \\ &\leq \underbrace{\max_{x \in D} \|\phi'(x)\|}_{=L} \|y-x\| \quad \blacksquare \end{aligned}$$

Definition $(X, \|\cdot\|)$ norm Raum, $(x_k)_{k \in \mathbb{N}}$, $x^*, x_k \in X$, $\lim_{k \rightarrow \infty} x_k = x^*$.

1) (x_k) heißt konvergent von Ordnung $p > 1$,

Falls $c > 0$, $k_0 \in \mathbb{N}$ ex. mit $\|x_{n+k} - x^*\| \leq c \|x_n - x^*\|^p$ ($p=2$, quadrat.)

2) „lineare Konvergenz“ falls $0 < L < 1$ ex mit Konvergenzfehler

$$\|x_{n+1} - x^*\| \leq L \|x_n - x^*\|$$

BSR FP Iteration ist linear konvergent

Newton Verfahren: $-G(x) := (F'(x))^{-1}$

$$\Rightarrow x_{n+1} := x_n + (F'(x_n))^{-1} F(x_n)$$

numerisch: $x_{n+1} = x_n + d_n$, $F'(x_n) dx = F(x_n)$ lin. Gl. lösen

$$F(x) = 0 \Leftrightarrow F(x) = F(x_n) + F'(x_n)(x - x_n) + O(\|x - x_n\|) \xrightarrow[\text{da linear}} \text{weglassen}$$

$$\Rightarrow -F(x_n) = F'(x_n)x_n - F'(x_n)x_n$$

$$F'(x_n)x_n = F'(x_n)x_n - F(x_n) \mid F'(x_n)^{-1}$$

$$x_{n+1} = x_n - F'(x_n)^{-1} F(x_n)$$

Satz 5.3 (Kantorowich-Newton)

Sei $D \subset \mathbb{R}$ offen und konvex

$F \in C^1(D)$ und $\sup_{x \in D} \|F'(x)\| = \beta$ und $\|F'(x) - F'(y)\| \leq L \|x-y\|$

und es ex. $x^* \in D$ mit $F(x^*) = 0$;

$0 < L$, (Hypothese konstante)

$$U_\omega(x^*) = U_\omega := \{x \in D : \|x - x^*\| \leq \omega\} \text{ und } \beta L \omega < 1$$

$$\Rightarrow \forall x \in U_\omega(x^*) \text{ konvergiert } x_n, (x_{n+1} = x_n - \frac{1}{2}(F'(x_n))^T F(x_n))$$

gegen x^* quadratisch (d.h. $p=2$) $\|x_{n+1} - x^*\| \leq C \|x_n - x^*\|^2$

$\phi: D \rightarrow \mathbb{R}^n, D \subset \mathbb{R}^n, x = (x_1, \dots, x_n)$

$\phi(x^*) = 0, x^* \in D$ gesucht

$$\text{N.V. } x^{n+1} = x^n + d^n$$

mit $\phi'(x_n)d_n = -\phi(x_n) \quad (x_n)_{n \in \mathbb{N}}, x^n \rightarrow x^*, n \rightarrow \infty?$

1) $\|\phi'(x)\|^{-1} \leq \beta, \forall x \in U_\omega(x^*) \cap D$ konvex

2) $\|\phi'(x) - \phi'(y)\| \leq L \|x-y\| \quad \forall x, y \in U_\omega(x^*) \cap D$ 3) $L \beta \omega < 1$
 $\Rightarrow x_0 \in U_\omega \Rightarrow x^n \in U_\omega(x^*), n \in \mathbb{N}, \|x^* - x^{n+1}\| \leq C \|x^n - x^*\|^2$
 $\rightarrow x^n \rightarrow x^*, n \rightarrow \infty, \text{quadr.}$

Beweis wir zeigen $\forall x, y \in U_\omega(x^*)$

$$\|\phi(x) - \phi(y) - \phi'(y)(x-y)\| \leq \frac{L}{2} \|x-y\|^2 \quad (* \text{ Vermutung})$$

Sei $t \mapsto \Psi(t) = \phi(y + t(x-y)), t \in [0,1]$



$$\Psi'(t) = \underbrace{\phi'(y+t(x-y))}_{\in \mathbb{R}^{n \times n}} \underbrace{(x-y)}_{\in \mathbb{R}^n}$$

$$\|\Psi(t) - \Psi(0)\| = \|\phi'(y+t(x-y)) - \phi'(y)(x-y)\|$$

$$\leq \|\phi'(y+t(x-y)) - \phi'(y)\| \|x-y\|$$

$$\leq L \|y + t(x-y) - y\| \|x-y\| \leq L t \|x-y\|^2$$

$$\|\phi(x) - \phi(y) - \phi(y)(x-y)\| = \|\Psi(1) - \Psi(0)\| \leq \int_0^1 \|\Psi'(t) - \Psi'(0)\| dt$$

$$\leq \frac{L}{2} \|x-y\|^2 \int_0^1 t dt \leq \frac{L}{2} \|x-y\|^2 \quad (\text{Beweis von } *) \blacksquare$$

Sei $x^n \in U_\omega(x^*)$:

$$x^{n+1} - x^* = x^n - (\phi'(x_n))^{-1}(\phi(x_n) - \phi(x^*)) - x^*$$

$$= (\phi'(x_n))^{-1} \left\{ \phi'(x_n)(x^n - x^*) \stackrel{=} 0 \right\} (\phi(x^n) - \phi(x^*))$$

$$\Rightarrow \|x^{n+1} - x^*\| = \|x^n - x^* - (\phi'(x_n))^{-1}(\phi(x^n) - \phi(x^*))\|$$

$$\leq \underbrace{\|\phi'(x_n)\|^{-1}}_{\leq \beta} \|\phi'(x_n)(x^n - x^*) - \phi(x^n) - \phi(x^*)\|$$

$$\stackrel{*}{=} \frac{\beta L}{2} \underbrace{\|x^n - x^*\|}_{\leq \omega} \|x^n - x^*\|$$

$$\leq \underbrace{\beta \omega}_{= q < 1} \|x^n - x^*\| = q \|x^n - x^*\|$$

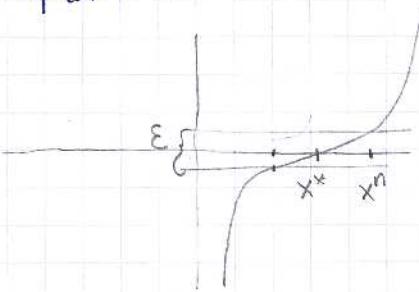
$\Rightarrow x^{n+1} \in U_\omega(x^*), x^n \rightarrow x^*, n \rightarrow \infty$ (Banach), mit $p=2$

Abbruchkriterium

$$\phi(x^*) = 0$$

Residuum

$$\|\phi(x^n)\| < \varepsilon$$



Alternative:

$$\left\| (\phi'(x^n))^{-1} \phi(x^n) \right\| < \varepsilon$$

d^n

a posteriori Fehler

Monotonietest (Konvergenztest)

geg. $0 < \theta < 1$: $\|\phi(x^{n+1})\| \leq \theta \|\phi(x^n)\|$ die Frage ist, ob das gut

(1) Idee $\|x^{n+1} - x^*\| \leq \theta \|x^n - x^*\|$, $\theta < 1$)

Alternativ: $\left\| (\phi'(x^{n+1}))^\top \phi(x^{n+1}) \right\| \leq \theta \left\| (\phi'(x^n))^\top \phi(x^n) \right\|$ $0 < \theta < 1$

Gedämpftes Newton Verfahren $x^{n+1} = x^n + \lambda_k^n d^n$, ($d^n = (\phi'(x^n))^{-1} \phi(x^n)$),

Für $k=0, \dots, K$

$$x^{n+1} = x^n + \lambda_{1/2}^n d^n$$

$$\text{mit } 0 < \lambda_k^n \leq 1$$

$$\lambda_0^n = 1, \lambda_1^n = \lambda_2^n = \dots = \lambda_K^n = 2^{-K}$$

↪ Monotonetest → bestanden $\rightarrow x^n \rightarrow x^{n+1}$ Iteration

$$\rightarrow \text{nicht bestanden} \Rightarrow \lambda_k^n \rightarrow \lambda_{k+1}^n = \frac{\lambda_k^n}{2}$$

Vereinfachte N-V

Für $n \gg 1$ ist die Berechnung von $\phi'(x^n)$ bzw. $(\phi'(x^n))^{-1} \in \mathbb{R}^{n \times n}$

sehr teuer.

Vereinfacht N-V: $\phi'(x^n) \approx \phi'(x^0)$, konvergiert unter obigen Vor.

$$\text{d.h. } x^{n+1} := x^n + (\phi'(x^0))^{-1} \phi(x^n), \text{ aber nun linear}$$

Praktisch macht man oft ein Update der Jacobi-Matrix $\phi'(x^n)$ nach k-reing- Iteration (z.B. $K=5, 10$ etc)

Quasi-Newton Verfahren

z.B. Sekantenverfahren ($n=1$), "Regula falsi", $n>1$, "Broyden Verfahren" (p. 1, 6)

Hier $1 < p < 2$. Oftmals Methoden der Wahl.

Nichtlineare Optimierung (ohne Nebenbed)

Problem: Kostenfkt $J: D \rightarrow \mathbb{R}$,

gesucht $x^* \in D$, mit $J(x^*) \leq J(x) \quad \forall x \in D$, $x^* = \arg \min [J](x); x \in D$

Sei $J \in C^k(D)$, $k=1, 2, \dots$, „stationärer Punkt“

Notwendige Bedingung $J'(x^*) = 0$

Gradienten Verfahren („steepest descent“)

$$x^{n+1} = x^n + \lambda^n d^n \text{ mit } d^n := -\text{grad } J(x^n) = -\nabla J(x^n)$$

Liniensuche für λ^n

$$(z.B. \lambda_{k+1}^n = \frac{1}{2} \lambda_k^n)$$

Monotonetest: z.B. $J(x^{n+1}) \leq \theta J(x^n) \quad \theta < 1$

\Rightarrow Konvergenz gegen stationäre Punkte, i.R. lokale Minima.

Alternativ: „Gauß-Newton-Verf.“

$x \mapsto \phi(x) = \nabla J(x) \cdot D \rightarrow \mathbb{R}^n$, zur Lösung $\nabla \phi(x^*) = \phi(x^*) = 0$

$$\phi'(x) = H_J(x) = J''(x), \text{Hessematrix von } J$$

6. Iterationsverfahren zur Lsg. linearer GL.

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

Ges: $x \in \mathbb{R}^n \quad Ax = b, Ax - b = 0$

Eliminationsmeth.: #arith. Op $O(n^3)$

Iterationsref: Matrix-Vektor multiplikation pro Iterationsschritt

$O(n^2)$ falls A voll besetzt

$O(n)$ falls A dünn besetzt d.h. #nnz A = $O(n)$

In der Praxis z.B. "FEM" usw. häufig $n > 10$

Finite-Elemente-Methode

Fixpunkt(FP)-Iteration: Sei $C (\approx A^{-1}) \in \mathbb{R}^{n \times n}$, „Vorkonditionieren“

$\phi(x) = x - C(Ax - b) = x$, „Fixpunkt-Problem“ $\phi: \mathbb{R}^n \rightarrow \mathbb{R}^n, D = \mathbb{R}^n$, vollst. $\|\cdot\| = \|\cdot\|_p$

FP-Iteration: $x^0 \in \mathbb{R}^n$ (bel. gewählt)

$$x^{n+1} = \phi(x^n) = x^n + C(b - Ax^n) = (I - CA)x^n + Cb$$

$M = (I - CA)$ Iterationsmatrix

$$\phi: \text{Kontraktion?} : \|(\phi(x) - \phi(x'))\| = \|(\mathbf{I} - CA)x + Cb - \{(\mathbf{I} - CA)x' + Cb\}\| \\ = \|(\mathbf{I} - CA)x\| \leq \|\mathbf{I} - CA\| \|x\|$$

$\|\mathbf{I} - CA\| < 1 \Leftrightarrow \phi$ Kontraktion

Satz 6.1. Sei $\varrho := \|\mathbf{I} - CA\|_{op} < 1$, dann konvergiert

$$\forall x^0 \in \mathbb{R}, x^{n+1} = x^n + C(b - Ax^n) \Rightarrow x \in \mathbb{R}^n \text{ (linear)}$$

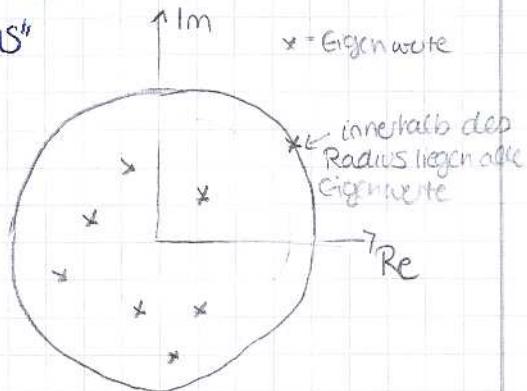
mit 1) $\|x^{n+1} - x\| \leq \varrho \|x^n - x\|$, ϱ -Fehlerverbrauchsfaktor

$$2) \|x^n - x\| \leq \frac{\varrho^n}{1-\varrho} \|x^0 - x\|, \quad (\|x^{n+1} - x\| \leq \varrho^{n+1} \|x^0 - x\|)$$

Beweis Banach-EPs

Def Seien $\lambda_i, i=1, \dots, n$ die Eigenwerte von A

$$\varrho(A) = \max_{1 \leq i \leq n} |\lambda_i| \quad \text{"Spektralradius"}$$



Satz 6.2. Sei $M = (\mathbf{I} - CA)$

$$x^{n+1} = Mx^n + Cb, \text{ dann konvergiert}$$

$$x^n \rightarrow x \Leftrightarrow \varrho(M) < 1$$

Beweis 1) Sei $\varrho(M) > 1$ $\forall x \neq x^0 \neq 0, x^0 \in \mathbb{R}$ bel. fest, wähle $v = \sum_{i=1}^n \alpha_i c_i$ Eigenwert

$$\Rightarrow \|x^{n+1} - x\| = \|Mx^n - Mv = (\mathbf{I} - CA)x^n - (\mathbf{I} - CA)v\| \\ = Mx^n + Cb - (Mx + Cb) = M(x^n - x) = M^{n+1}(x - x^0) = M^{n+1}|v|, \\ \Rightarrow \|x^n - x\| \leq \|M^n\| |v| \leq \left\| \sum_{k=1}^n \lambda_k^* \alpha_k c_k \right\| \leq \lambda_{\max}^n \max_{1 \leq k \leq n} |\alpha_k| \lambda_{\max} |c_k| = \alpha_{\lambda_{\max}} \lambda_{\max}^n |c_{\lambda_{\max}}| \\ = \lambda_{\max}^n \|x^0\| \xrightarrow{n \rightarrow \infty} \infty \\ \Rightarrow x^n \text{ divergiert}$$

$$\Leftarrow \varrho(M) < 1, \forall x_0 \in \mathbb{R}^n, \Rightarrow v = x^0 - x$$

$$\text{Es gilt } \|x^n - x\| = \|M^n v\|,$$

Sei $T \in \mathbb{R}^{n \times n}$ reg. mit $H = T^{-1}JT$ die Jordanzersetzung von H

$$\text{a.h. } J = \text{diag}(J_j) = J_e = \begin{pmatrix} \lambda_e^1 & & \\ & \ddots & \\ & & \lambda_e^m \end{pmatrix} = \lambda_e^1 + S_e, S_e = \begin{pmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix} \in \mathbb{R}^{m \times m}$$

$$H^n = (T^{-1}JT)^n = T^{-1}J^n T, \text{ es genügt}$$

$$\Rightarrow S_e^m = 0$$

$$J_e^n \rightarrow 0, n \rightarrow \infty \text{ zu zeigen } \forall \ell \quad n > m, m \text{ fest} \quad n \rightarrow \infty$$

$$\begin{aligned} J_e^n &= (\lambda_e I + S_p)^n = \sum_{k=0}^n \binom{n}{k} \lambda_e^{n-k} S_e^k \\ &= \sum_{k=0}^m \binom{n}{k} \lambda_e^{n-k} S_e^k \xrightarrow{n \rightarrow \infty} 0, \quad \lambda_e^{n-k} \xrightarrow{n \rightarrow \infty} 0 \text{ für } k=0, \dots, m \end{aligned}$$

Korollar 6.2 \forall induzierten Operatornormen $\| \cdot \|_{\text{op}}$

$$\text{gilt } S(M) \leq \| M \|_{\text{op}}$$

6.2.1 Richardson Iteration

$$C = \alpha I, \alpha \in \mathbb{R} \text{ geeignet}$$

$$R - \bar{E} \Rightarrow x^{n+1} = (I - \alpha A)x^n + \alpha b = x^n - \alpha Ax^n + \alpha b$$

konvergiert für $S(I - \alpha A) < 1$, z.B. für $\| I - \alpha A \| < 1 \quad \forall x^0 \in \mathbb{R}^n$
 $(A^{-1} \approx \alpha A^{-1})$

6.2.2 Gesamtschritt oder Jacobi-Vergleich

wir zerlegen $A = D + L + R$, $D = \text{diag}(a_{i,i})$, $L = \begin{pmatrix} 0 & & \\ a_{21} & 0 & \\ & \ddots & 0 \end{pmatrix}$, $R = \begin{pmatrix} 0 & * \\ 0 & 0 \end{pmatrix}$

Voraus: Sei $D \in \mathbb{R}^{n \times n}$ regulär, d.h. $a_{i,i} \neq 0 \quad \forall i = 1, \dots, n$

$$\text{Wähle } C = D^{-1} \quad (\approx A^{-1})$$

$$\Rightarrow x^{n+1} = x^n - D^{-1}(D + L + R)x^n + D^{-1}b = x^n - x^n - D^{-1}(L + R)x^n + D^{-1}b$$

$$\text{d.h. } x^n = (x_1^n \quad x_2^n \quad \dots \quad x_n^n)^T \in \mathbb{R}^n$$

$$x_k^{n+1} = a_{kk}^{-1} (b_k - \sum_{j=1, j \neq k}^n a_{kj} x_j^n) \quad k = 1, \dots, n$$

$$M = D^{-1}(L + R)$$

$$S(M) \leq 1? \quad \| D^{-1}(L + R) \|_1 = \max_{k \in \mathbb{N}} \sum_{j \neq k} \frac{1}{|a_{kj}|} |a_{kj}|$$

Def $A \in \mathbb{R}^{n \times n}$ heißt Strikt diagonal dominant

$$\max_{1 \leq k \leq n} \sum_{j=1}^n |a_{j,k}| \leq |a_{k,k}|$$

Sei $A \in \mathbb{R}^{n,n}$ regulär und sei $b \in \mathbb{R}^n$ genüge $AX=b$

Gesucht: $x \in \mathbb{R}^n$

$A = D + L + R$, D diagonal, R obere Δ -M., L untere Δ -M.

Mit Jakobi(Gesamtschritt)-Verfahren:

$$x_i^{k+1} := (a_{i,i})^{-1} (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^k - \sum_{j=i+1}^n a_{ij} x_j^k) \quad i=1, \dots, n \quad C = D^{-1}$$

Mit Gauß-Seidel (Einzelschritt)-Verfahren:

$$x_i^{k+1} := (a_{i,i})^{-1} (b_i - \sum_{j=1}^{i-1} a_{ij} x_j^k - \sum_{j=i+1}^n a_{ij} x_j^k) \quad i=1, \dots, n \quad C = (D+L)^{-1}$$

Def: $A = (a_{ij})_{i,j=1, \dots, n}$

A stark diagonaldominant, falls gilt

$$\sum_{j=1}^n |a_{ij}| < |a_{ii}| \quad \forall i=1, \dots, n$$

$$\sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| < |a_{jj}| \quad \forall j=1, \dots, n$$

6.4. A stark diagonaldominant, dann konvergieren

Jakobi und Gauß-Seidel für Wahl des Startvektors (linear)

6.5. A spd, dann konvergiert Gauß-Seidel ($\forall x_k$ linear)

Relaxationsverfahren

Iterationsvorschrift: $x^{k+1} = Gx^k + c$

Verfahren konv $\Leftrightarrow \sigma(G) < 1$

Um dies zu verbessern:

$\omega \in \mathbb{R}$ geeignet $x^{k+1} = \omega(Gx^k + c) + (1-\omega)x^k$

ω so gewählt, dass $\|(G + (1-\omega)I)\|$ minimiert wird

Angewandt auf Jakobi oder Gauß-Seidel \Rightarrow SOR-Verfahren

6.3. Gradienten und CG-Verfahren

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

dann definiert $\langle \cdot, \cdot \rangle_A = \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$

ein Skalarprodukt $\langle u, v \rangle_A = \langle u, Av \rangle = \langle Au, v \rangle = u^T A v$

mit der Norm $\|u\|_A = \sqrt{\langle u, u \rangle_A}$

Sei $u \in \mathbb{R}^n$, $u \mapsto J(u) := \frac{1}{2} \langle Au, u \rangle - \langle b, u \rangle$

und $x := \underset{u \in \mathbb{R}^n}{\operatorname{argmin}} J(u)$ erfüllt

$$J'(x) = 0 = (Ax)^T - b^T \Leftrightarrow Ax = b$$

Ritz-Galerkin-Verfahren

Sei $V_k = \operatorname{span}\{v_j, j=1, \dots, k\} \subset \mathbb{R}^n$. Wir suchen $x^k \in V_k$

$$x^k = \underset{u \in V_k}{\operatorname{argmin}} \{J(u) : u \in V_k\}$$

Gradientenverfahren (x^k) $\xrightarrow{k \in \mathbb{N}}, x^k \in \mathbb{R}^n$, Folge mit

$$J(x_k) \downarrow J(x) (\leq J(u) \forall u \in \mathbb{R}^n, \text{Min}) \xrightarrow{k \rightarrow \infty}$$

Idee: $x^{k+1} = x^k - \alpha^k \nabla J(x^k)$ ($\alpha^k \in \mathbb{R}_+$ geeignet)

$$= x^k + \alpha^k p_k$$

$$p_k = -\nabla J(x^k) = b - Ax^k =: r^k$$

$$x^k + \alpha^k p_k =: x + tp \quad (x = x^k, p = p^k \text{ fest})$$

$$\begin{aligned} J(x+tp) &= \tilde{J}(t) = \frac{1}{2} \langle A(x+tp), (x+tp) \rangle - \langle b, (x+tp) \rangle \\ &= \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle + t \langle Ax - b, p \rangle + \frac{1}{2} \langle Ap, p \rangle t^2 \rightarrow \min \end{aligned}$$

Führt zu $\frac{d}{dt} \tilde{J}(t) = 0 = \langle Ax - b, p \rangle + t \langle Ap, p \rangle$

$$\Leftrightarrow t = \frac{\langle b - Ax, p \rangle}{\langle Ap, p \rangle}$$

$$\alpha^k = \frac{\langle r^k, p^k \rangle}{\langle Ap^k, p^k \rangle}$$

Gradientenverfahren $p^k = b - Ax^k$

$$\alpha^k = \frac{\langle b - Ax^k, b - Ax^k \rangle}{\langle Ap^k, p^k \rangle} = \frac{\langle p^k, p^k \rangle}{\langle Ap^k, p^k \rangle}$$

$$x^{k+1} = x^k + \alpha^k p^k$$

(Vergleich Richardson-Vergleich)

Beobachtung: $x^k \in \text{span}\{b, Ar^0, A^2r^0, \dots, A^{k-1}r^0\} = K^k$
 $\text{span}\{x^j : j < k\} \oplus P_K$

Frage: existiert eine Linearkombination der Iterationen

d.h. $x_*^k \in K^k$, sodass das Funktional auf K^k minimiert?

Ja! \rightsquigarrow CG-Hestens-Stiefel

$$Ax = b, r = 0 = b - Ax$$

$$x_0 \in \mathbb{R}^n \text{ bel (z.B. } x_0 = 0) \quad x = x_0 + x^*$$

$$A(x_0 + x^*) = b \Leftrightarrow Ax^* = b - Ax_0 = \tilde{b}$$

$$\text{OBdA: Sei } x_0 = 0 \text{ gewählt} \Rightarrow r_0 = b$$

$$\begin{aligned} \text{Krylov-Raum } K_K \subseteq \mathbb{R}^n : K_K &= \text{span}\{K_0, b\} = \text{span}\{r_0, Ar_0, \dots, A^{K-1}r_0\} \\ &= \text{span}\{b, Ar_0, \dots, A^{K-1}b\} \\ &= \text{span}\{A^j b \mid j = 0 \dots K-1\} \end{aligned}$$

$$A \text{ ist spd: } J: \mathbb{R}^n \rightarrow \mathbb{R}: J(u) := \frac{1}{2} \langle u, Au \rangle - \langle b, u \rangle = \frac{1}{2} u^T Au - b^T u$$

$$x = A^{-1}b \Leftrightarrow x = \arg \min \{J(u) \mid u \in \mathbb{R}^n\}$$

$$J \text{ stkt konvex} \Rightarrow \exists ! x$$

Galerkin-Ritz-Vergfahren

$$\text{Sei } U_K \subseteq \mathbb{R}^n, \text{ z.B. } U_K := K_K = \text{span}\{p_1, \dots, p_K\}$$

$\exists ! x^k \in U_K \quad x^k = \arg \min \{J(u) \mid u \in U_K\}$ restriktionsoptimierungsproblem
 $U_K (= K_K)$ ist zulässige Menge (lineare Teilraum, konvex)

$$x^k \in U_K \text{ innerer Pkt } \exists \gamma \in C^1$$

\Rightarrow Notwendige Bed alle Richtungsableitungen von J um $x_k \in U_K$

$$\text{sind null, d.h. } 0 = J'(x_k)(u) \quad \forall u \in U_K \quad (\|u\|=1)$$

$$= (\nabla J(x_k))^T u$$

$$= \langle Ax_k - b, u \rangle \quad \forall u \in U_K \quad (\text{Schwache Formulierung})$$

$$\text{Ansatz: } U_K \ni x_K = \sum_{j=1}^K c_j p_j \Rightarrow \sum \langle p_i, Ap_j \rangle c_j = \langle p_i, b \rangle \quad \forall i = 1, \dots, K$$

$$\Rightarrow Ac = \underline{b} \in \mathbb{R}^K \quad \text{Galerkin-Gleichung}$$

$$\left\{ \begin{array}{l} \langle Ax_K, u \rangle \\ = \langle b, u \rangle \end{array} \right. \quad \forall u$$

Lemma Galerkin-Orthogonalität

$$x_k = \arg \{ J(u) \mid u \in U_k \} \Leftrightarrow r_k = b - Ax_k \perp U_k$$

$$(0 = \langle r_k, u \rangle = \langle b - Ax_k, u \rangle) \quad \forall u \in U_k$$

C-G-Verfahren

$$U_k = K_k = \text{span}\{r_0, Ar_0, A^2r_0\} = \text{span}\{p_{n-1}, p_k\}$$

mit A -orthogonalem p_k , d.h. $\langle p_j, p_k \rangle_A = \langle p_j, Ap_k \rangle$

$$\text{Sei } k < n, U_k = K_k, r_k = \underbrace{b - Ax_k}_{\in X_{k+1}} \quad (x_k = K_k = \sum_{j=1}^k \langle r_k, p_j \rangle_A p_j)$$

$$\notin K_k$$

$$\text{Gram-Schmidt} \quad p_{k+1} = r_k - \sum_{j=1}^k \frac{\langle r_k, p_j \rangle_A}{\langle p_j, p_j \rangle_A} p_j$$

$$= r_k - \frac{\langle r_k, p_k \rangle_A p_k}{\langle p_k, p_k \rangle_A} - \cancel{\frac{\langle r_k, p_{k-1} \rangle_A}{\langle p_{k-1}, p_{k-1} \rangle_A} p_{k-1}} = 0 \text{ wegen Galerkin-orthog.}$$

(wollen, dass x_k Galerkin-Lsg sind!!)

$$\text{Seien } x_k = \arg \min \{ J(u) \mid u \in K_k \}$$

$$\Rightarrow p_{k+1} = r_k + \beta_k p_k \quad \text{mit } \beta_k := -\frac{\langle r_k, A p_k \rangle}{\langle p_k, A p_k \rangle}$$

konjugester Gradient

Wir wählen die Iterationsfolge $x_i \in K_k$, sodass

$$x_k = \arg \min \{ J(u) \mid u \in K_k \}$$

$$K_{k+1} = \text{span}\{p_{n-1}, p_{k+1}\} = K_k \oplus \text{span}\{p_{k+1}\}$$

$$\Rightarrow x_{k+1} = x_k + \alpha_k p_{k+1} \quad \begin{matrix} \text{falls } v_k \text{ statt } p_k \text{ dann wie} \\ \text{Gradienten-verfahren} \end{matrix}$$

$$\Rightarrow x_{k+1} = \arg \min \{ J(u) \mid u \in K_{k+1} \} \Rightarrow \alpha_k = \frac{\langle p_{k+1}, p_{k+1} \rangle}{\langle A p_{k+1}, p_{k+1} \rangle}$$

$$r_k = b - Ax_k = A(x - x_k) = A(x - (x_{k-n} + \alpha_{k-n} p_k))$$

$$= r_{k-1} - \alpha_{k-n} A p_k$$

Lemma: Es gilt $\langle r_k, v_j \rangle = (r_k)^T v_j = \delta_{kj} \langle r_k, r_k \rangle$

Beweis: Galerkin-Orthogonalität. Sei $j < k$ (OBdA)

$$\Rightarrow r_j = b - Ax_j \in K_{j+1} \quad \text{u.G.O.} \Rightarrow r_k \perp K_k$$

$$\Rightarrow j < k \quad r_k \perp K_{j+1} \Rightarrow \langle r_k, r_j \rangle = 0 \blacksquare$$

Gradienten-update

Allgemeiner Ansatz

$$x_{k+1} = x_k + \alpha_k (r_k - \sum_{j=1}^k \beta_j p_j)$$

$$\Gamma K_{k+1} = X_k + \text{span}\{r_k\}$$

$$\forall k \quad \langle p_j, p_i \rangle_A = 0, i \neq j$$

$$(A x_{k+1} - b) \cdot u = (A(x_k + \alpha_k (r_k - \sum_{j=1}^k \beta_j p_j)) - b) u$$

Galerkin Gleichung mit $x_{k+1} = \arg \min$

$$\langle A(x_k + \alpha_k (r_k - \sum_{j=1}^k \beta_j p_j)) - b, p_i \rangle = 0 \quad \forall i = 1, \dots, k+1$$

$$\underbrace{\langle A(x_k - \alpha_k r_k - \alpha_k \beta_k p_k) - b, p_i \rangle}_{=0} - \sum_{j=1}^{k-1} \beta_j \underbrace{\langle A p_j, p_i \rangle}_{= S_{ij} \langle p_j, p_i \rangle} = 0 \quad \forall i \leq k-1$$

$$\Rightarrow \beta_j = 0 \quad r_j \in K$$

$$\Rightarrow x_{k+1} = x_k + \alpha_k (r_k + \beta_k p_k)$$

$$\text{Wegen } \alpha_k \langle r_k, p_{k+1} \rangle_A = \langle r_k, r_{k+1} - r_k \rangle$$

$$= \langle r_k, b - A x_{k+1} - r_k \rangle = \langle r_k, b - A(x_k + \alpha_k p_{k+1}) - (b - A x_k) \rangle$$

$$= \langle r_k, \alpha_k p_{k+1} \rangle_A$$

$$\beta_k = \langle r_k, p_{k+1} \rangle_A = \frac{\langle r_k, r_k \rangle}{\alpha_k} = \frac{\langle r_k, r_k \rangle}{\langle r_{k-1}, r_{k-1} \rangle}$$

$$\alpha_k = \frac{\langle r_{k-1}, r_{k-1} \rangle}{\langle A p_k, p_k \rangle}$$

Im Vergleich zum Gradienten 2-fache Matrix-Vektormultiplikation und doppeltes Speicher : ABER:

Theorem Es gibt Fehlerabschätzung

(wissen wir nicht schreibweise Fehler wie bei BFS)

$$\|x_k - x\|_A \leq 2 \left(\frac{1 - \sqrt{\lambda(A)}}{1 + \sqrt{\lambda(A)}} \right)^k \|x_0 - x\|_A \quad \begin{matrix} (\text{Iterationszahl}) \\ k \sim \sqrt{\lambda(A)} \end{matrix} \quad \gg$$

(Grad. Verfahren:

$$\|x^k - x\|_A \leq C \left(\frac{1 - \sqrt{\lambda(A)}}{1 + \sqrt{\lambda(A)}} \right)^k \|x_0 - x\|_A$$

(Iterationszahl)
 $k \sim \sqrt{\lambda(A)}$)

Das CG Verfahren ist das schnellste

Iterationsverfahren für Spd-Matrizen

In der Praxis „vorkonditioniertes CG Verfahren“

Nach n Schritten hat man die exakte Lsg.

7. Numerische Behandlung von Eigenwertproblemen

$$A \in \mathbb{C}^{n \times n}, \langle u, v \rangle := \bar{u}^T v$$

$(A^* \bar{A}) = A^H$, adjungierte Matrix zu A. $\langle A^H, u, v \rangle = \langle u, Av \rangle$

$$A = (a_{ij})_{i,j=1}^n, A^H = (\bar{a}_{ji})_{i,j=1}^n$$

$$(A \in \mathbb{R}^{n \times n} \Rightarrow A^H = A^T)$$

$$\begin{cases} z \in \mathbb{C} \\ z = a + ib, \operatorname{Re}(z) = a \\ \operatorname{Im}(z) = b \end{cases}$$

$$\bar{z} = a - ib$$

Unitäre Matrizen U: $U^{-1} = U^H, (U^H)^H = U$, U, U_1, U_2 ist unitär
bilden Gruppe

$$\text{Es gilt: } \|Ux\|_2 = \|x\|_2$$

(λ, v) ein Eigenpaar, falls i) $v \neq 0$ und ii) $Av = \lambda v$ $\forall \lambda \in \mathbb{C}, v \in \mathbb{C}^n$

Ähnliche Matrizen: A und B heißen ähnlich, falls eine reguläre Matrix $T \in \mathbb{C}^{n \times n}$ existiert mit $A = T^{-1}BT$ ($A \sim B$)

Falls $A \sim B$, haben A und B die gleichen EW $\lambda_1, \dots, \lambda_n \in \mathbb{C}$

\Rightarrow Falls U unitär existiert mit $A = U^HBU$, dann besitzen

A und B die gleichen EW.

$\{\lambda_1, \dots, \lambda_n\} = \sigma(A)$, Spektrum von A

Satz 7.1 (Schur-Normalform)

Zu $A \in \mathbb{C}^{n \times n}$ existiert $Q \in \mathbb{C}^{n \times n}$ unitär und eine rechte obere Δ -Matrix mit $A = Q^H R Q = Q^{-1} R Q$

$$A^H = (Q^H R Q)^H = Q^H R^H (Q^H)^H = Q^H R^H Q \xrightarrow{R^H = R = \operatorname{diag}(\lambda_i)}$$

Korollar: Falls $A^H = A$ (hermitisch)

(reell symm.)

Beweis: Es existiert ein Eigenpaar (λ_1, u_1) mit $\|u_1\|_2 = 1$

$\{v_2, \dots, v_n\}$ eine Orthonormalbasis ($\langle v_i, v_j \rangle = \delta_{ij}$) von

$\operatorname{span}\{u_1\}^\perp = \{v \in \mathbb{C}^n, \langle u_1, v \rangle = 0\} \Rightarrow \{u_1, v_2, \dots, v_n\}$ ist ONB von \mathbb{C}^n

Sei $P = (u_1, v_2, \dots, v_n) \Rightarrow P^H = P^{-1}$ ist unitär (Bew. Übung)

$$P^H A P = \begin{pmatrix} \lambda_1 & * & * \\ 0 & I & * \\ 0 & * & -* \end{pmatrix} = \begin{pmatrix} \lambda_1 & & \\ 0 & R_1 & \\ 0 & 0 & \ddots \end{pmatrix} \xrightarrow{\lambda_1, R_1, \dots, R_{n-1}} \begin{pmatrix} \lambda_1 & * & * \\ 0 & \lambda_2 & R_2 \\ 0 & 0 & \ddots \end{pmatrix} = R^H P^H A P P_1 P_2 \quad P_2 = \begin{pmatrix} I & 0 \\ 0 & P_2 \end{pmatrix}$$

$$P_n^H = P_2^H P^H A P_1 P_2 \cdots P_n = U A U^H = \begin{pmatrix} \lambda_1 & * \\ 0 & \lambda_n \end{pmatrix} \Rightarrow A = U^H R U$$

Satz 7.3 „Kreise des Gershgorin“

$$\text{Seien } r_i = \sum_{j=1}^n |a_{ij}|, c_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$$

$$\Rightarrow \sigma(A) \subset \bigcup_{i=1}^n \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}$$

$$c_i \cup \{z \in \mathbb{C} : |z - a_{ii}| \leq c_i\}$$



Beweis Sei $\lambda \in \mathbb{C}$ mit $|a_{ii} - \lambda| > r_i \forall i=1, \dots, n$

$$B = A - \lambda I \Rightarrow |b_{ii}| = |a_{ii} - \lambda| \quad \forall i=1, \dots, n \quad \forall i \neq j \quad b_{ij} = a_{ij}$$

$$\Rightarrow \sum_{\substack{j=1 \\ j \neq i}}^n |b_{ij}| = \sum_{j \neq i} |a_{ij}| = r_i < |b_{ii}| \quad \forall i=1, \dots, n$$

d.h. B ist diagonaldominant $\Rightarrow B$ ist regulär, d.h. λ ist kein EW von A

Satz 7.4 (Satz von Bauer & Fike)

Sei A diagonalisierbar, d.h. es ex. $T \in \mathbb{C}^{n \times n}$ reg.

mit $\text{diag}(\lambda_i)_{i=1, \dots, n}^{\text{Eigenwerte von } A} \Rightarrow D = T^{-1}AT$

$SA \in \mathbb{C}^{n \times n}$, λ ein EW von $B = A + SA$ und $1 \leq p \leq \infty$ dann gilt

$$\min_{1 \leq j \leq n} |\lambda_j - \lambda| \leq \|SA\|_p \text{ cond}_p(T)$$

Beweis $B = D - \lambda I + T^*SAT$

$$TBT^{-1} = TD T^{-1} + \lambda TT^{-1} + SA$$

$$\|A\|_p = \|A\|_{\text{op}} \text{ in } \ell_p$$

$= A + SA - \lambda I$ singular!

$$(D - \lambda I)^{-1}B = I + (D - \lambda I)^{-1}T^*SAT \text{ sing}$$

$$\Rightarrow g((D - \lambda I)^{-1}T^*SAT) > 1$$

$$\Rightarrow 1 < g(1) \leq \|D - \lambda I\|^{-1} \|T^*SAT\|_p$$

$$\leq \|D - \lambda I\|^{-1} \|T\|_p \|T^*\|_p \|SA\|_p \|T\|_p$$

$$\leq \frac{1}{\min_{1 \leq j \leq n} |\lambda_j - \lambda|} \|SA\|_p \text{ cond}(T) \blacksquare$$

Transformation auf Hessenberggestalt

"Hessenbergmatrix": $\begin{pmatrix} * & * & & \\ * & * & & \\ & & \ddots & \\ 0 & * & \cdots & * \end{pmatrix} = H$ (symm. Hessenb.m. sind tridiagonale)

Konstruiere Q unitär mit $Q^*AQ = H \Rightarrow \sigma(A) = \sigma(H)$

z.B. durch Givens-Rotation aber auch durch Householder

Iterative Näherungsverfahren zur Berechnung des EW Problems symm. Matrizen

$A \in \mathbb{R}^{n \times n}, A^T = A \quad (A^H = A)$

\Rightarrow dann ex. $Q \in \mathbb{R}^{n \times n}$ (Q unitär)

mit $A = QT \text{ diag}(\lambda_i) Q, \lambda_i \in \mathbb{R}, \lambda_1, \dots, \lambda_n \in \mathbb{R}$

Def 7.1 $\mu: \mathbb{R}^n \setminus \{0\} \rightarrow \mathbb{R}$ Rayleigh-Quotient von A

$$\mu(x) = \frac{\langle Ax, x \rangle}{\langle x, x \rangle}, \text{ Es gilt das folgende}$$

Satz 7.6 „Courant Fischer“

Es gilt das folgende Extremalprinzip

$$\lambda_{\max} = \max \left\{ \sum_{i=1}^n \mu(x) \mid x \in \mathbb{R}^n \right\}, \quad \lambda_{\min} = \min \left\{ \sum_{i=1}^n \mu(x) \mid x \in \mathbb{R}^n \right\} = \min \left\{ \langle Ax, x \rangle \mid x \in \mathbb{R}^n, \|x\|_2^2 = 1 \right\}$$

Satz 7.6 „Courant Fischer“

Eigenwertprobleme symmtr. Matrizen

$$A \in \mathbb{R}^{n \times n}, A^T = A \quad (A^H = A)$$

$$A = U^T \cdot \text{diag}(\lambda_i) \cdot U, \quad U = (u_1, \dots, u_n)$$

Satz 7.6 „Courant Fischer“

? Norm, wenn nur gesagt geschrieben wird

Sei $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$, EW mit $\mathbb{C}V u_1, \dots, u_n$ $\|u_i\| = 1$ $i=1, \dots, n$

$$1) \lambda_i = \max \left\{ \frac{\langle Ax, x \rangle}{\langle x, x \rangle} \mid 0 \neq x \in \text{span}\{u_{i+1}, \dots, u_n\} \right\} \Rightarrow \lambda_{\max} = \max_{x \in \mathbb{R}^n} \frac{\langle Ax, x \rangle}{\langle x, x \rangle} = \max \left\{ \frac{\langle Ax, x \rangle}{\langle x, x \rangle} \mid x \in \text{span}\{u_{i+1}, \dots, u_n\}, \|x\|=1 \right\}$$

$$2) \lambda_i = \min \left\{ \frac{\langle Ax, x \rangle}{\langle x, x \rangle} \mid 0 \neq x \in \text{span}\{u_i, u_{i+1}, \dots, u_n\} \right\}$$

Beweis $U \in \mathbb{R}^{n \times n}$ orthog. (unitär) $A = U^T \text{diag}(\lambda_i) U$

$$1) \lambda_i = \max \{ \lambda_k \mid 1 \leq k \leq i \} = \max \{ \langle Au_i, u_i \rangle \mid 1 \leq k \leq i, \|u_i\|=1 \}$$

$$\lambda_i = \langle Au_i, u_i \rangle \geq \langle Ax, x \rangle, \forall x \in \text{span}\{u_{i+1}, \dots, u_n\}, \|x\|=1$$

$$2) \text{ analog } \blacksquare$$

Potenzmethode (Poweriteration, von Mises) (Vektoriteration)

$$X_0 \in \mathbb{R}^n \text{ gewählt: } X_{k+1} = A \cdot X_k = (A^{k+1} X_0), \quad Y_{k+1} = \frac{X_{k+1}}{\|X_{k+1}\|}$$

$$\boxed{X_{k+1} = A X_k \Rightarrow X_{k+1} = \frac{X_{k+1}}{\|X_{k+1}\|}}$$

nur andere Schreibweise

$$(\|Y_k\| = 1, \forall k \in \mathbb{N})$$

Satz 7.7 λ_i EW von A, $|\lambda_1| \geq |\lambda_2| \geq \dots$ $\lambda_{\max, k} = \langle A Y_k, Y_k \rangle$

$u_i, x_0 \neq 0, X_0 \in \mathbb{R}^n$ bel., $X_{k+1} \in A X_k, Y_{k+1} = \frac{X_{k+1}}{\|X_{k+1}\|}$

„nicht ortho um Ev ansonsten erlösen“
 $\lim_{k \rightarrow \infty} Y_k = \pm u_i$, Eigenvektor zu $\lambda_i = \langle A u_i, u_i \rangle = \lambda_{\max}$

Beweis: Sei $x_0 = \sum_{j=1}^n \alpha_j u_j$ ($\|x_0\|=1$), da $0 \neq \langle x_0, u_1 \rangle \Rightarrow \alpha_1 \neq 0$

$$X_k = A X_{k-1} = A^k X_0 = A^k \left(\sum_{j=1}^n \alpha_j u_j \right) = \sum_{j=1}^n \alpha_j A^k u_j = \sum_{j=1}^n \alpha_j \lambda_j^k u_j$$

$$= \lambda_i^k (\alpha_1 u_i + \sum_{j \neq i} \left(\frac{\lambda_j}{\lambda_i} \right)^k \alpha_j u_j) \xrightarrow{\|X_k\| \rightarrow 0 \text{ für } k \rightarrow \infty, \text{ da } q = \left| \frac{\lambda_i}{\lambda_1} \right| < 1} u_i \text{ für } k \rightarrow \infty$$

zu beweisen: Konvergenzrate hängt von $\left| \frac{\lambda_2}{\lambda_1} \right| < 1$ ab.

Bemerkung: 1) für $|\lambda_1| \approx |\lambda_2|$ wird das Verfahren langsam

2) Berechnet immer nur den maximalen EW.

$|\lambda_i| \geq |\lambda_n|, \lambda_i \text{ EW von } A$

Dann gilt: $\left| \frac{1}{\lambda_n} \right| \leq \left| \frac{1}{\lambda_i} \right| \quad \left| \begin{array}{l} \text{EW von } A^{-1} \\ \text{mit den gleichen EV} \end{array} \right.$

Inverse Iteration ("Wielandt")

Löse $Ax_{k+1} = x_k$, d.h. $x_{k+1} = A^{-1}x_k, u_{k+1} = \frac{x_{k+1}}{\|x_{k+1}\|}$

Beschleunigen durch "Shift + Invert":

Sei $\mu \notin \sigma(A)$, d.h. $\min_{1 \leq j \leq n} |\mu - \lambda_j| > 0$ und $\lambda = \operatorname{argmin}_{1 \leq j \leq n} |\mu - \lambda_j|$

Sei $S = \frac{1}{\lambda - \mu}$, sei $u \in V$ zu $\lambda \in \sigma(A)$.

$\Rightarrow (A - \mu I)u = (\lambda - \mu)u$, d.h. $(\lambda - \mu) \in \sigma(A - \mu I)$ und der kleinste Eigenwert von $(A - \mu I)^{-1}$

$\Rightarrow (A - \mu I)^{-1}u = \frac{1}{\lambda - \mu}u = Su$, S größter? EW von $(A - \mu I)^{-1}$

Potenzmethode (Vektoriteration) mit $(A - \mu I)^{-1}$

d.h. $(A - \mu I)x_{k+1} = x_k$. (Problem: $\operatorname{cond}(A - \mu I) \gg 1$!)

Ausweg: Jakobi-Davidson)

bestes Algo um einen Eigenwert/Vektor? zu berechnen

Das QR-Verfahren von Francis:

$A_1 := A = QR_1$ Zerlegung (z.B. nach Householder)

$A_2 = R_1 Q_2 = Q_1^T A Q_1 \quad Q_k \in \mathbb{R}^{n \times n}$

$A_{kn} = R_k Q_k \quad \Rightarrow \bar{A}_k^T = A_k$

$\Rightarrow A_k \rightarrow \operatorname{diag}(\lambda_k), Q_k \cdots Q_1 \xrightarrow{k \rightarrow \infty} U$

$A^k (a_i^k)_{i,j=1}^n \quad \text{Satz 79 i)} \lim_{k \rightarrow \infty} Q_k = I$

ii) $\lim_{k \rightarrow \infty} R_k = \lim_{k \rightarrow \infty} A_k = \operatorname{diag}(\lambda_i)$

$|A| > \dots$

iii) $|a_j^{(k)}| = \left| \prod_{i=1}^k \frac{\lambda_i}{\lambda_j} \right| \quad \forall j > i$

Beweis per Induktion

Lanczos Verfahren

Krylov Raummethode: $\lambda_{\max}^{(k)} = \max \left\{ \frac{\langle Ax_i, x \rangle}{\langle x_i, x \rangle} \mid 0 \neq x \in \mathcal{K}^k \right\}$
 $\mathcal{K}^k = \text{span} \{ p_j : j = 1, \dots, k \}$ wie bei CG?

$$\langle Ap_i, p_j \rangle = \begin{pmatrix} * & 0 \\ * & * & * & * \\ 0 & * & * & * \\ * & * & * & * \end{pmatrix} \in \mathbb{R}^{k \times k}$$

3.6 Splines

Splines: Stückweise Polynome mit maximaler Glättigkeit. Trägt von Splines (schön?)

$k \geq 2$

$S_{k,\tau} = \{ f \in C^{k-1}([a,b]) \mid f|_{[t_j, t_{j+m}]} \in \Pi_{k-1}, j=0, \dots, l \}$ stückweise stetig für $k=2$
 Splinefkt der Ordnung k (Grad + 1)

Prop 3.6.1 1) dim $S_{k,\tau} = k+l$

2) $S_{k,\tau} = \text{span} \{ \Pi_{k-1}|_{[a,b]}, (x-t_j)_+^{k-1}, j=0, \dots, l-1 \}$

Beweis 2) Ein- und abw. $\Pi_k \in S_{k,\tau}$ $(x-t_j)_+^{k-1} \in C^{k-2}$

3.6.2. B-Splines

$$k=1, N_{j,1}(x) = \chi_{[t_j, t_{j+1}]}(x) = \begin{cases} 1, & x \in [t_j, t_{j+1}] \\ 0, & \text{sonst} \end{cases}$$

$$k=2, N_{j,2}(x) = \begin{cases} \frac{x-t_j}{t_{j+2}-t_j}, & x \in [t_j, t_{j+1}] \\ \frac{t_{j+2}-x}{t_{j+2}-t_{j+1}}, & x \in [t_{j+1}, t_{j+2}] \\ 0, & \text{sonst} \end{cases}$$

$$\text{Def: } N_{j,k}(x) = (t_{j+k} - t_j) \left[\frac{[t_j, t_{j+k}]}{(t_{j+k} - t_j)} (x-t_j)_+^{k-1} \right]$$

B-Splines der Ordnung k und Knoten j

Satz 3.6.3 Es gilt die Rekursionsformel

$$1) N_{j,k}(x) = \frac{x-t_j}{t_{k+j}-t_j} N_{j,k-1}(x) + \frac{t_{j+k}-x}{t_{k+j}-t_{j+1}} N_{j+1,k-1}(x)$$

$$2) N_{j,k}'(x) = (k-1) \left\{ \frac{N_{j,k-1}(x)}{t_{k+j}-t_j} - \frac{N_{j+1,k-1}(x)}{t_{k+j}-t_{j+1}} \right\}$$

Bew Skizze $(x-t_j)_+^{k-1} (x-t_{j+1})_+^{k-2} (x-t_{j+2})_+^{k-3}$

$[\] (x-t_j)_+^{k-1} = [\] \text{ II Leibniz Formel für Differenzen}$

$$(f(x)) = \begin{cases} f(x), & f(x) \geq 0 \\ 0, & \text{sonst} \end{cases}$$

- Satz 3.6.2
- $N_{j,k} \in \text{Span} \left\{ (T_j - x)_+^{k-1}, l=j, j+k \right\}$
 - $\text{Supp } N_{j,k} = \{x : N_{j,k}(x) \neq 0\} \subset [t_j, t_{j+k}]$
 - $N_{j,k}(x) \geq 0$
- Konvention: $\frac{0}{0} := 0$

Bew i) klar

$$\text{(ii)} \quad x < t_j, N_{j,k}(x) \subset \Pi_{k-1, l \neq j} U_l(x) \Rightarrow [t_{j+1}, t_{j+k}] (-x)_+^{k-1} = 0$$

$$x > t_{j+k}, N_{j,k}(x) = 0 \Rightarrow \text{(iii)}$$

iii) klar

Lemma 3.6.4 Marsden Identität

$$(x-y)^{k-1} = \sum_{i=1}^{k+l} \prod_{j=1}^{k-1} (t_{i+j}-y) N_{i,k}(x)$$

Beweis Induktion über k

i) $k=1$ klar

ii) Sei $k > 1$, $1 \leq j \leq k-1+l$ bewiesen

$$\begin{aligned} \Psi_{i,k} &= \prod_{j=1}^{k-1} (t_{i+j}-y) \\ \sum_{i=1}^{k+l} (\Psi_{i,k}(y) N_{i,k}(x)) &= \sum_{i=2}^{k+l} \underbrace{\left(\frac{x-t_i}{t_{i+k-1}-t_i} \Psi_{i,k}(y) + \frac{t_{i+k-1}-x}{t_{i+k-1}-t_i} \Psi_{i+1,k}(y) \right)}_{\text{Klammer}} N_{i,k}(x) \\ &= \sum_{i=2}^{k+l} \prod_{j=1}^{k-1} (t_{i+j}-y) \cdot \underbrace{\left[\frac{x-t_{i+k-1}}{t_{i+k-1}-t_i} (t_{i+k-1}-y) + \frac{t_{i+k-1}-x}{t_{i+k-1}-t_i} (t_{i+k-1}-y) \right]}_{\text{Klammer}} \\ &= (x-y) \sum_{i=2}^{k+l} N_{i,k-1}(x) = (x-y)(x-y)^{k-2} \end{aligned}$$

Erweiterte Kettenfolge

$$\begin{array}{c} \text{Hilf} \\ t_0 = t_0 \\ t_1 = t_1 \\ \vdots \\ t_k = t_k \\ t_{k+1} = t_{k+1} \\ \vdots \\ t_{k+l} = t_{k+l} \end{array}$$

t_i nicht notwendig gerade
verschieden $\Rightarrow t_i \Rightarrow t_i$

$$y \mapsto f(y) = (x-y)^{k-1}, f^{(l)}(y) = (k-l)(k-l-1)\dots(k-l-l+1) x^{k-l-1} = \sum_{i=1}^{k+l} \Psi_{i,k}^{(l)}(y) N_{i,k}(x)$$

$$\Rightarrow x^m \in \text{Span} \{N_{i,k} : l=0, \dots, k+l\} \quad \forall m \leq k-1$$

$$1 = \sum_{i=1}^m N_{i,k}(x), \text{"Zerlegung des Ein"}$$

Satz 3.6.5 1) $S_{k,l} = \text{span} \{N_{i,k} : i=1, \dots, k+l\}$

$$2) 1 = \sum_{i=1}^{k+l} N_{i,k}(x), \Pi_{k-1} \subset \text{span} \{N_{i,k}(x) : i=1, \dots, k+l\}$$

$$3) \text{ cont } \|f - v_k\|_{\infty} \leq C h^k \sup_{x \in [a, b]} |f^{(k)}(x)|$$

$v_k \in S_{k,T}$

$$h := \max \{t_{i+1} - t_i\}$$

De Boer Algorithmus (Carl de Boer)

$$\text{Sei } x \mapsto S(x) = \sum_{j=1}^{k+l} c_j N_{j,k}(x), x \in [a, b]$$

$$S(x) = \sum_{j=2}^{k+l} c_j^{[n]}(x) N_{j,k-n}(x), c_j^{[1]}(x) = \frac{x - t_j}{t_{j+k-1} - t_j} c_j + \frac{t_{j+k-1} - x}{t_{j+k-1} - t_j} c_{j-1}$$

$$S(x) = \sum_{j=1+n}^{k+l} c_j^{[n]}(x) N_{j,k-n}(x)$$

$$\text{mit } c_j^{(m)}(x) = \begin{cases} \frac{x - t_j}{t_{j+k-1} - t_j} c_j^{[n-1]}(x) \\ + \frac{t_{j+k-1} - x}{t_{j+k-1} - t_j} c_{j-1}^{[n-1]}(x) \end{cases}$$

$$= \sum_{j=k}^{l+k} c_j^{[k-1]} N_{j,1}(x) = c_m^{(m)}(x) \quad x \in [t_m, t_{m+n}]$$

Analog lassen sich damit auch die Ableitungen von S d.h. $S^{(m)}(x)$ berechnen

