

4.4 Rundungsfehleranalyse der Gauß-Elimination (LR-Zerlegung)

1
10.11.0

Absolutwertnotation

Für $A \in \mathbb{R}^{m \times n}$ ist

$$B = |A| \Rightarrow b_{ij} = |a_{ij}| \quad 1 \leq i \leq m, 1 \leq j \leq n$$

Die Abbildung $rd: \mathbb{R} \rightarrow \mathbb{F}$ erweitern wir entsprechend auf $\mathbb{R}^n, \mathbb{R}^{m \times n}$.

Dann gilt

$$rd(A) = A + A' \quad \text{mit} \quad |A'| \leq |A| \epsilon_{ps}$$

dies sind $m \cdot n$ Ungleichungen

Wdh. aus
Analyse von:
rd

$$rd(a_{ij}) = a_{ij}(1 + \epsilon_{ij}) = a_{ij} + \underbrace{a_{ij} \epsilon_{ij}}_{a'_{ij}}$$

$$|a'_{ij}| \leq |a_{ij}| \epsilon_{ps}$$

fl-Notation ~~→ wichtigste Not. früher einführen~~

Sei F eine Formel, dann bezeichne $fl(F)$ eine Berechnung der Formel F in \mathbb{R} .

Fließkommaarithmetik (dabei wird die Ausführungsreihenfolge meist offensichtlich sein, sonst ist sie anzugeben).
(links nach rechts, Klammern)

Somit ist zum Beispiel

$$fl(A+B) = (A+B) + H$$

$A, B \in \mathbb{F}^{n \times n}$ exakte Rechnung \mathbb{R}

mit $|H| \leq \epsilon_{ps} (|A+B|)$

(folgt aus

$$\begin{aligned} fl(a_{ij} + b_{ij}) &= a_{ij} \oplus b_{ij} \\ &= (a_{ij} + b_{ij})(1 + \epsilon_{ij}) \quad |\epsilon_{ij}| \leq \epsilon_{ps} \\ &= (a_{ij} + b_{ij}) + \epsilon_{ij} (a_{ij} + b_{ij}) \end{aligned}$$

Nachmal: |·| ist keine Norm sondern eine Matrix.

Man kann auch

$$fl(\sqrt{x^2})$$

oder $fl(\sin(x))$

schreiben.

Rückwärtsanalyse

Bisher haben ^{wir} die "Vorwärtsanalyse" der Rundungsfehler betrieben.
 z.B. gilt für das Skalarprodukt $x^T y$:

$$|fl(x^T y) - x^T y| \leq n \text{ eps} |x|^T |y| + O(\text{eps}^2) \quad \text{(absolut)}$$

oder

$$\frac{|fl(x^T y) - x^T y|}{|x^T y|} \leq n \text{ eps} \frac{|x|^T |y|}{|x^T y|} + O(\text{eps}^2) \quad \text{(relativ)}$$

Eine Alternative ist die sog. "Rückwärtsanalyse".

Dort versucht man das Fließkommaergebnis als exakter Ergebnis (ines modifizierten Ausdrucks) zu schreiben.

Beispiel 4.8 Betrachte Lösung des LGS $Ax = b$. Die GEM berechnet die numerische Lösung \hat{x} .

Vorwärtsanalyse: $\|\hat{x} - x\| \leq F(\text{eps}, n, A, b)$

Rückwärtsanalyse: $(A+E)\hat{x} = b$ mit $|E| \leq F'(\text{eps}, n, A)$

Mit dem Störungssatz **3.18** und $\|E\|_\infty = \| |E| \|_\infty$ folgt dann für d. Rundungsfehler:

$$\frac{\|\hat{x} - x\|_\infty}{\|x\|_\infty} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A)} \frac{\| |E| \|_\infty}{\|A\|_\infty}$$

Rundig: $rd(A) = A + \delta A$
 $|\delta A| \leq |A| \cdot \tau$

Störungssatz:

$$\frac{\|\delta A\|_\infty}{\|A\|_\infty} = \frac{\| |A| \|_\infty \tau}{\|A\|_\infty} \leq \frac{\|A\|_\infty \tau}{\|A\|_\infty} = \tau = \text{eps}$$

ist noch entsprechend abzuschätzen, ideal wäre $\| |E| \|_\infty \leq n^2 \text{ eps} \|A\|_\infty$

Somit ist ein direkter Vergleich mit der Konditionsanalyse möglich.

Rückwärtsanalyse des Skalarproduktes.

Hilfssatz 4.9 Es gilt $x, y \in \mathbb{F}^n$:

3
10.11.09

$$\hat{s} = fl(x^T y) = (x + f)^T y, \quad |f| \leq n \text{ eps } |x| + O(\text{eps}^2)$$

Beweis: Induktion über n .

$n=1$: $\hat{s}_1 = fl(x_1 y_1) = x_1 y_1 (1 + \delta_1) = (x_1 + \underbrace{x_1 \delta_1}_{=: f_1}) y_1$

also $|f_1| = |\delta_1| |x_1| \leq \text{eps } |x_1|$.

$n \geq 2$: $\hat{s}_n = fl(x^T y) = fl\left(\underbrace{\tilde{x}^T \tilde{y}}_{x^T y} + x_n y_n\right)$

$x = \begin{pmatrix} \tilde{x} \\ x_n \end{pmatrix}, y = \begin{pmatrix} \tilde{y} \\ y_n \end{pmatrix}$

$= \left(fl(\tilde{x}^T \tilde{y}) + fl(x_n y_n) \right) (1 + \epsilon_n)$

$= \left((\tilde{x} + \tilde{f})^T \tilde{y} + x_n y_n (1 + \delta_n) \right) (1 + \epsilon_n)$

$= (\tilde{x} + \tilde{f})^T \tilde{y} (1 + \epsilon_n) + x_n y_n (1 + (\delta_n + \epsilon_n) + \delta_n \epsilon_n)$

$= (\tilde{x} + \tilde{f} + \epsilon_n \tilde{x} + \epsilon_n \tilde{f})^T \tilde{y} + (x_n + (\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n) y_n$

$= \begin{bmatrix} \tilde{x} \\ x_n \end{bmatrix} + \begin{bmatrix} \tilde{f} + \epsilon_n \tilde{x} + \epsilon_n \tilde{f} \\ (\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n \end{bmatrix} \begin{bmatrix} \tilde{y} \\ y_n \end{bmatrix}$
 $\underbrace{\hspace{10em}}_{=: x} \quad \underbrace{\hspace{10em}}_{=: f} \quad \underbrace{\hspace{10em}}_{=: y}$

mit $|\tilde{f} + \epsilon_n \tilde{x} + \epsilon_n \tilde{f}| \leq (n-1) \text{eps } |\tilde{x}| + \text{eps } |\tilde{x}| + O(\text{eps}^2) \leq n \text{eps } |\tilde{x}| + O(\text{eps}^2)$

und $|(\delta_n + \epsilon_n)x_n + \delta_n \epsilon_n x_n| \leq 2 \text{eps } |x_n| + O(\text{eps}^2)$

Wg $n \geq 2$ gilt also $|f| \leq n \text{eps } |x| + O(\text{eps}^2)$ ▣

Bemerkung 4.90 $n \text{eps}$ in 4.9 ist der schlechteste Fall und sehr pessimistisch. Rundungsfehler in einer Operation hängt von den Argumenten ab und sie können sich auch wegheben (positiv/negativ!). Besser wäre eine statistische Betrachtung. ▣

Satz 4.19 (Lösen von Dreieckssystemen)

4
10.11.03

Es sind \hat{x} bzw. \hat{y} die numerischen Lösungen des unteren bzw. oberen Dreieckssysteme $Lx=b$ und $Ry=c$. Dann gilt

$$(L+F)\hat{x} = b \quad |F| \leq n \text{ eps } |L| + O(\text{eps}^2)$$

$$(R+G)\hat{y} = c \quad |G| \leq n \text{ eps } |R| + O(\text{eps}^2)$$

Beweis: (evtl als Übung). Induktion über n ,

$n=1$ $l_{11}x_1 = b_1 \rightarrow \hat{x}_1 = fl(b_1/l_{11}) = (b_1/l_{11})(1+\epsilon_1)$

$$\Leftrightarrow \frac{1}{1+\epsilon_1} \hat{x}_1 = \frac{b_1}{l_{11}} \Leftrightarrow \frac{l_{11}}{1+\epsilon_1} \hat{x}_1 = b_1$$

(1) $\frac{1}{1+\epsilon_1} = 1 + \delta_1$
 $\delta_1 \approx \frac{1}{1+\epsilon_1} - 1 = \frac{-\epsilon_1}{1+\epsilon_1} \approx -\epsilon_1$

$$\Rightarrow (1 - \epsilon_1 + \tilde{R}(\epsilon^2)) l_{11} \hat{x}_1 = b_1$$

$$\Leftrightarrow (l_{11} - \underbrace{\epsilon_1 l_{11}}_f + l_{11} \tilde{R}(\epsilon^2)) \hat{x}_1 = b_1$$

mit $|f| \leq \text{eps } |l_{11}| + O(\text{eps}^2)$

$n \geq 2$ $Lx=b \Leftrightarrow \begin{matrix} n-1 & 1 \\ L_1 & 0 \\ v^T & \alpha \end{matrix} \begin{pmatrix} x_1 \\ w \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$

Berechne Lösung \hat{x}_1 von $L_1 x_1 = b_1$ und rechne

$$\hat{w} = fl((\beta - v^T \hat{x}_1)/\alpha)$$

↙ Fehler in α

$$= fl((\beta - v^T \hat{x}_1)/\alpha) (1+\epsilon_n)$$

↘ Skalarprodukt

$$= ((\beta - fl(v^T \hat{x}_1))(1+\delta_n)/\alpha) (1+\epsilon_n) = ((\beta - fl(v^T \hat{x}_1))/\alpha) (1+\epsilon_n)(1+\delta_n)$$

H5.5.8

$$\stackrel{H5.5.8}{=} ((\beta - (v+f)^T \hat{x}_1)/\alpha) (1 + (\epsilon_n + \delta_n) + \epsilon_n \delta_n)$$

Satz 4.12 (Rückwärtsanalyse der LR-Zerlegung)

Sei $A \in \mathbb{F}^{n \times n}$. Es werde die LR-Zerlegung ohne Pivotierung berechnet. Dann gilt für die numerisch berechneten L^1, \hat{R}^1 :

$$L^1 \hat{R}^1 = A + H \quad \text{mit} \quad \|H\| \leq 3(n-1)\text{eps}(\|A\| + \|L^1\| \|\hat{R}^1\|) + O(\text{eps}^2)$$

Beweis: [Golub/Van Loan, TMM 3.3.1] Induktion über n .

$n=1$: Es ist $\hat{L}_{11}^1 = 1$ und $\hat{R}_{11}^1 = a_{11}$ und damit $H=0$.

$n \geq 2$: Schreibe

$$A = \begin{bmatrix} 1 & & & \\ \alpha & & & \\ & & & \\ v & & B & \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}$$

Dann macht die LR-Zerlegung:

a) $\hat{z} = fl(v/\alpha)$

b) $\hat{A}_1 = fl(B - \hat{z}w^T)$

c) Berechne LR-Zerlegung von \hat{A}_1

Fehler in \hat{z} :

(I₁) $\hat{z} = \frac{v}{\alpha} + f$ mit $|f| \leq \text{eps} \frac{|v|}{|\alpha|}$

Fehler in \hat{A}_1 :

$\hat{A}_1 = fl(B - \hat{z}w^T)$

← exakt gerundet

(I₂) Bistaff: $= B - fl(\hat{z}w^T) + G$

$|G| \leq \text{eps} |B - fl(\hat{z}w^T)|$

$= B - (\hat{z}w^T + \underbrace{G'}_=: F) + G$

$|G'| \leq \text{eps} |\hat{z}w^T| \leq \text{eps} |\hat{z}| |w|^T$

← ex. product! genau eine Operation pro matrix eintrag.

(II₂) $\hat{A}_1 = B - \hat{z}w^T + F$

mit $|F| = |G' + G| \leq |G'| + |G|$

$\leq \frac{\text{eps} |\hat{z}| |w|^T + \text{eps} |B - \hat{z}w^T - G'|}{|G'|}$

$\leq \text{eps} |\hat{z}| |w|^T + \text{eps} (|B| + |\hat{z}| |w|^T) + \text{eps}^2 |\hat{z}| |w|^T$

also $|F| \leq 2\text{eps} (|B| + |\hat{z}| |w|^T) + O(\text{eps}^2)$

Mit $\frac{1}{1 + (\epsilon_n + \delta_n) + \epsilon_n \delta_n} \approx 1 - (\epsilon_n + \delta_n) + R(\epsilon_n^2 + \delta_n^2 + \epsilon_n \delta_n)$ gilt

$$(1 - (\epsilon_n + \delta_n) + R(\dots)) \hat{w} = (\beta - (v + f)^T \hat{x}_1) / \alpha$$

$$\Leftrightarrow (v + f)^T \hat{x}_1 + (1 - (\epsilon_n + \delta_n) + R(\dots)) \alpha \hat{w} = \beta$$

Mit der Induktionsvoraussetzung $(L_1 + F_1) \hat{x}_1 = b_1$ folgt also für \hat{x} .

$$\begin{pmatrix} L_1 + F_1 & 0 \\ (v + f)^T & (1 - (\epsilon_n + \delta_n) + R(\dots)) \alpha \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{w} \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$$

$$\Leftrightarrow \left[\underbrace{\begin{pmatrix} L_1 & 0 \\ v^T & \alpha \end{pmatrix}}_{= L} + \underbrace{\begin{pmatrix} F_1 & 0 \\ f^T & (-(\epsilon_n + \delta_n) + R(\dots)) \alpha \end{pmatrix}}_{= F} \right] \begin{pmatrix} \hat{x}_1 \\ \hat{w} \end{pmatrix} = \begin{pmatrix} b_1 \\ \beta \end{pmatrix}$$

Wegen $|F_1| \leq (n-1) \epsilon \rho |L| + O(\epsilon \rho^2)$

$|f^T| \leq (n-1) \epsilon \rho |v|^T + O(\epsilon \rho^2)$

$|(-(\epsilon_n + \delta_n) + R(\dots)) \alpha| \leq 2 \epsilon \rho |\alpha| + O(\epsilon \rho^2)$

hier kommt das n-1 her!
nicht durch die Rekursion.
es geht also nicht besser

gilt für $n \geq 2$ dann

$$|F| \leq n \epsilon \rho |L| + O(\epsilon \rho^2)$$

(möglich wäre auch $\max(n-1, 2)$, ist aber nicht wirklich besser).

Analog für $Ry = c$.



Nun wird \hat{A}_1 LR-zersetzt und es gilt die Induktionsannahme. 7
10.11.09

$$(D_2) \quad L_1^{-1} \hat{R}_1^{-1} = \hat{A}_1 + H_1 \quad \text{mit} \quad |H_1| \leq 3(n-2) \cdot \text{eps} (|L_1^{-1}| + |R_1^{-1}|) + O(\text{eps}^2)$$

also ist die Blockform der LR-Zersetzung von A :

$$L^{-1} R^{-1} = \begin{bmatrix} 1 & 0 \\ \tilde{z} & L_1^{-1} \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ 0 & \hat{R}_1^{-1} \end{bmatrix} = \begin{bmatrix} \alpha & w^T \\ \tilde{z}\alpha & \tilde{z}w^T + \underbrace{L_1^{-1} \hat{R}_1^{-1}}_{= \hat{A}_1 + H_1} \end{bmatrix}$$

Darstellung
von
oben
in
Blockform

$$(D_1), (D_2), (D_3) \begin{bmatrix} \alpha & w^T \\ (\frac{\sigma}{\alpha} + f)\alpha & \tilde{z}w^T + \underbrace{(B - \tilde{z}w^T + F)}_{= \hat{A}_1} + H_1 \end{bmatrix}$$

$$= \begin{bmatrix} \alpha & w^T \\ 0 & B \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ \alpha f & H_1 + F \end{bmatrix}$$

$= A \qquad \qquad \qquad =: H$

Damit hat man schon auch die Form $L^{-1} A^{-1} = A + H$. Bleibt noch H abzuschätzen.

In H steckt $H_1 + F$, in $|H_1|$ steckt $|\hat{A}_1|$ also:

$$|\hat{A}_1| = |B - \tilde{z}w^T + F| \leq |B| + |\tilde{z}| |w^T| + |F|$$

Abh. für $|F| \rightarrow \leq |B| + |\tilde{z}| |w^T| + 2 \text{eps} (|B| + |\tilde{z}| |w^T|) + O(\text{eps}^2)$

$$\leq (1 + 2 \text{eps}) (|B| + |\tilde{z}| |w^T|) + O(\text{eps}^2)$$

$$|H_1 + F| \leq |H_1| + |F|$$

Indukt. Annahme $\rightarrow \leq 3(n-2) \text{eps} (|\hat{A}_1| + |L_1^{-1}| |R_1^{-1}|) + 2 \text{eps} (|B| + |\tilde{z}| |w^T|) + O(\text{eps}^2)$

Abh. für $|\hat{A}_1|$ aus $\rightarrow \leq 3(n-2) \text{eps} \left[(1 + 2 \text{eps}) (|B| + |\tilde{z}| |w^T|) + |L_1^{-1}| |R_1^{-1}| \right] + 2 \text{eps} (|B| + |\tilde{z}| |w^T|) + O(\text{eps}^2)$

$$\leq 3(n-1) \text{eps} \left[|B| + |\tilde{z}| |w^T| + |L_1^{-1}| |R_1^{-1}| \right] + O(\text{eps}^2)$$

Und damit

$$|H| = \begin{bmatrix} 0 & 0 \\ |\alpha| & |w| \end{bmatrix} \leq \begin{bmatrix} 0 & 0 \\ \epsilon |\alpha| & 3(n-1)\epsilon \left(|B| + |\hat{z}| |w|^T + |L_1^1| + |\hat{R}_1| \right) \end{bmatrix} + O(\epsilon^2)$$

(a₁)

19.11.09

$$\leq 3(n-1)\epsilon \begin{bmatrix} 0 & 0 \\ |w| & |B| + |\hat{z}| |w|^T + |L_1^1| + |\hat{R}_1| \end{bmatrix} + O(\epsilon^2)$$

$$\leq 3(n-1)\epsilon \left(\underbrace{\begin{bmatrix} |\alpha| & |w|^T \\ |w| & |B| \end{bmatrix}}_{|A|} + \underbrace{\begin{bmatrix} 1 & 0 \\ |\hat{z}| & |L_1^1| \end{bmatrix}}_{|L_1^1|} \underbrace{\begin{bmatrix} |w| & |w|^T \\ 0 & |\hat{R}_1| \end{bmatrix}}_{|\hat{R}_1|} \right)$$

() Terme sind
sinngefügig.
D.h. (1)

Nun sind noch die Dreieckssysteme anzulösen.

Folgerung

Nach Satz 4.13 ist \hat{x} exakte Lösung des modifizierten Systems

$$(A+E)\hat{x} = b.$$

Mit dem Störungssatz gilt dann in $\|\cdot\|_\infty$ -Norm (Beachte: $\|B\|_\infty = \|\|B\|_\infty\|$)

$$\frac{\|\hat{x} - x\|_\infty}{\|x\|_\infty} \leq$$

$$\text{Cond}(A) \cdot \left\{ \underbrace{3n \cdot \text{eps} \cdot \frac{\|A\|_\infty}{\|A\|_\infty}}_{=1} + 5n \cdot \text{eps} \cdot \frac{\|\hat{L}\| \|\hat{R}\|}{\|A\|_\infty} + O(\text{eps}^2) \right\}$$

Problem!

bis auf
Faktor \rightarrow

Vergleichbar
mit Rundungs-
fehler $\|rd(A)\|/\|A\|$

$$\Rightarrow \frac{\|SA\|_\infty}{\|A\|_\infty} \leq \text{eps} \frac{\|A\|_\infty}{\|A\|_\infty} = \text{eps}$$

O.K. takes Mal $rd(A) = A + SA$ mit $|SA| \leq |A| \cdot \text{eps}$

Nehmen an, dass $\|E\|_\infty / \|A\|_\infty \ll 1$ (brauchen ohnehin $\|E\| < \frac{1}{\|A^{-1}\|}$)

\rightarrow Nenner im Vorfaktor ist ≈ 1

- Erster Term vergleichbar mit dem aus der Konditionsanalyse
- Zweiter Term ist möglicherweise problematisch.

\hat{L} enthält Einträge der Form $\frac{\tilde{a}_{ij}}{\tilde{a}_{ii}}$, also $\frac{1}{\text{Pivot element}}$.

$|\text{Pivotelement}| \text{ klein} \Rightarrow |\hat{L}| \text{ groß} \Rightarrow \text{großer Rundungsfehler!}$

- Dies kann trotz guter Kondition von A passieren!

Beispiel: $A = \begin{bmatrix} \epsilon & 1 \\ 1 & 0 \end{bmatrix}$

\Rightarrow Gauß-Elimination (LR-Zerlegung) ist in dieser Form nicht numerisch stabil!

Satz 4.12

9
11.11.09

Seien \hat{L}^1 und \hat{R}^1 die numerisch berechnete LR-Zerlegung von $A \in \mathbb{F}^{n \times n}$ aus Satz 4.12. Sei weiter $\hat{y}^1 \in \mathbb{F}^n$ die numerische Lösung von $\hat{L}^1 y = b$ und schließlich $\hat{x}^1 \in \mathbb{F}^n$ die numerische Lösung von $\hat{R}^1 x = \hat{y}^1$. Dann gilt für \hat{x}^1 die Beziehung

$$(A+E)\hat{x}^1 = b$$

mit

$$|E| \leq n \text{ eps} (3|A| + 5|\hat{L}^1||\hat{R}^1|) + O(\text{eps}^2).$$

Beweis:

Wg Satz 4.19 gilt

$$(\hat{L}^1 + F)\hat{y}^1 = b \quad |F| \leq n \text{ eps} |\hat{L}^1| + O(\text{eps}^2)$$

$$(\hat{R}^1 + G)\hat{x}^1 = \hat{y}^1 \quad |G| \leq n \text{ eps} |\hat{R}^1| + O(\text{eps}^2)$$

Einsetzen liefert:

$$(\hat{L}^1 + F)\hat{y}^1 = (\hat{L}^1 + F)(\hat{R}^1 + G)\hat{x}^1 = (\underbrace{\hat{L}^1 \hat{R}^1}_{=A+H} + F\hat{R}^1 + \hat{L}^1 G + FG)\hat{x}^1 = b$$

Wg Satz 4.12 gilt

$$\hat{L}^1 \hat{R}^1 = A + H \quad |H| \leq 3(n-1) \text{ eps} (|A| + |\hat{L}^1||\hat{R}^1|) + O(\text{eps}^2)$$

also

$$(A+E)\hat{x}^1 = b \quad \text{mit} \quad E = H + F\hat{R}^1 + \hat{L}^1 G + FG$$

und

$$|E| \leq |H| + |F||\hat{R}^1| + |\hat{L}^1||G| + O(\text{eps}^2)$$

$$\leq \underbrace{3(n-1) \text{ eps} (|A| + |\hat{L}^1||\hat{R}^1|)}_{|H|} + \underbrace{n \text{ eps} |\hat{L}^1||\hat{R}^1|}_{|F|} + \underbrace{n \text{ eps} |\hat{L}^1||\hat{R}^1|}_{|G|} + O(\text{eps}^2)$$

$$\leq 3n \text{ eps} |A| + 5n \text{ eps} |\hat{L}^1||\hat{R}^1| + O(\text{eps}^2). \quad \square$$

$$\Leftrightarrow |e_{ij}| \leq n \text{ eps} \left(3|a_{ij}| + 5 \sum_{k=1}^n |e_{ik}| |\tau_{kj}| \right) + O(\text{eps}^2)$$

4.5 Pivotalisierung

Die Rundungsfehleranalyse in Satz 4.13 führt auf den unvorteilhaften Term $\|L\|_{\infty}$.

Mit der Wahl von r im Gauß-Algorithmus so dass

$$|a_{rk}^{(k)}| \geq |a_{ik}^{(k)}| \quad \forall k \leq i \leq n$$

passt schon,
aber $k < n$

gilt dann

$$|\hat{l}_{ij}| \leq 1 \text{ und damit } \|L\|_{\infty} \leq n.$$

Diese Wahl nennt man "Spaltenpivotalisierung" (oder maximales Spaltenpivot).

Beispiel 5.13 Aus [GrL] → Blatt.

$$\begin{bmatrix} -10^{-5} & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \quad A^{-1} = \frac{-1}{2+10^{-5}} \begin{bmatrix} 1 & -1 \\ -2 & -10^{-5} \end{bmatrix}$$

In exakter Arithmetik führt das Gaußsche Verfahren nach Elimination von a_{21} auf

$$\begin{bmatrix} -10^{-5} & 1 \\ 0 & 1+2 \cdot 10^5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \cdot 10^5 \end{bmatrix}$$

mit der Lösung

$$x_1 = -0.4999975, \quad x_2 = 0.999995$$

Nun führen wir das Verfahren in $\mathbb{F}(10, 4, 1)$ durch. Beim Multiplikator (Normierung)

$$q_{21} = (0.2 \cdot 10^1) \oslash (-0.1 \cdot 10^{-4}) = -0.2 \cdot 10^6$$

ergibt sich kein Rundungsfehler.

Für das neue a_{22} ergibt sich

$$\left[\begin{array}{cc|c} -0.1 \cdot 10^{-4} & 1 & -1 \\ 0.2 \cdot 10^6 & 0.2 \cdot 10^6 & 0.2 \cdot 10^6 \end{array} \right] \begin{array}{l} a_{22}^{(1)} = 0.1 \cdot 10^1 \ominus (-0.2 \cdot 10^6) \oslash (0.1 \cdot 10^1) \\ = 0.1 \cdot 10^1 \oplus 0.2 \cdot 10^6 = \boxed{0.2 \cdot 10^6} \end{array} \rightarrow \begin{array}{l} 0.2000 \cdot 10^6 \\ 0.000001 \cdot 10^6 \end{array}$$

$0.1 \cdot 10^{-5} \cdot 10^5 \cdot 10^4$
 \uparrow
 10^6

Hier wurde auf vier Stellen gerundet.

Damit ergibt sich (ohne Fehler)

$$b_2^{(1)} = \ominus (-0.2 \cdot 10^6) \oslash (0.1 \cdot 10^1) = 0.2 \cdot 10^6$$

und

$$x_2 = b_2^{(1)} \oslash a_{22}^{(1)} = 0.2 \cdot 10^6 \oslash 0.2 \cdot 10^6 = \boxed{1}, \text{ (statt } 0.999995)$$

$$x_1 = (0.1 \cdot 10^1 \ominus 0.1 \cdot 10^1 \oslash 1) \oslash (-0.1 \cdot 10^{-4}) = \boxed{0}$$

2
16.11.09

Es ist also *keine* Stelle im Ergebnis korrekt obwohl nur an einer *einzig*en Stelle (in der Berechnung von $a_{22}^{(1)}$) ein Rundungsfehler eingeführt wurde.

Darüberhinaus überprüfe man, dass für die Kondition von A gilt:

$$\kappa(A) = 3$$

Demnach ist das System gut konditioniert! Der Algorithmus, so wie er ist, ist numerisch nicht stabil.

Das Problem ist offensichtlich der große Multiplikator q_{21} der aus dem sehr kleinen a_{11} resultiert und der dafür sorgt, dass das ursprüngliche a_{22} in $a_{22}^{(1)}$ vollkommen ignoriert wird.

Im Prinzip haben wir in Fließkommaarithmetik das System

$$\begin{bmatrix} -10^{-5} & 1 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

exakt gelöst, was eine völlig andere Lösung hat als das ursprüngliche (11.1) (Rückwärtsanalyse).

Der große Multiplikator kann ganz einfach vermieden werden indem man eine Zeilenvertauschung durchführt, d. h. wir lösen

$$\begin{bmatrix} 2 & 1 \\ -10^{-5} & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}. \tag{11.2}$$

Nun erhält man

$$q_{21} = -0.1 \cdot 10^{-4} \oslash 0.2 \cdot 10^1 = -0.5 \cdot 10^{-5},$$

$$a_{22}^{(1)} = 0.1 \cdot 10^1 \ominus (-0.5 \cdot 10^{-5}) \oslash 0.1 \cdot 10^1 = 0.1 \cdot 10^1 \oplus 0.5 \cdot 10^{-5} = 0.1 \cdot 10^1,$$

$$b_2^{(1)} = 0.1 \cdot 10^1 \ominus 0.5 \cdot 10^{-5} \oslash 0 = 0.1 \cdot 10^1,$$

$$x_2 = 0.1 \cdot 10^1 \oslash 0.1 \cdot 10^1 = \boxed{1},$$

$$x_1 = (0 \ominus 0.1 \cdot 10^1 \oslash 0.1 \cdot 10^1) \oslash 0.2 \cdot 10^1 = \boxed{-0.5},$$

was in $\mathbb{F}(10, 4, 1)$ völlig in Ordnung ist. □

Spaltenpivotisierung ist nicht ausreichend wie folgendes Beispiel zeigt.

Fortb. von Beispiel 5.13

(ebenfalls aus [GO96]). Wir betrachten das 2×2 System

$$\begin{bmatrix} 10 & -10^6 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -10^6 \\ 0 \end{bmatrix}.$$

welches aus (11.1) durch Multiplikation der ersten Zeile mit -10^6 entsteht.

Die Spaltenpivotisierung erfordert keine Vertauschung. Allerdings entsteht für $a_{22}^{(1)} = 1 + 2 \cdot 10^5$ genau dasselbe Problem wie oben! □

Diese Probleme kann man durch eine Skalierung (Äquilibrierung). 3
16.11.09
des Gleichungssystems vermindern:

$$Ax = b \rightarrow D^{-1}Ax = D^{-1}b \quad \text{mit } d_{ii} = \sum_{j=1}^n |a_{ij}|$$

$$\Leftrightarrow \tilde{A}x = \tilde{b} \quad (\text{also } \|\tilde{A}\|_{\infty} = 1)$$

Rundungsfehleranalyse bei Pivotisierung ja

Analog zu Satz 4.13 zeigt man, dass für die Lösung \hat{x} bei Spaltenpivotisierung gilt: (Referenz GVL)

$$(A+E)\hat{x} = b \quad \text{mit } \|E\| \leq n \text{ eps} (3\|A\| + 5 P^T \|L\| \|R\|) + O(\text{eps}^2).$$

Nach Konstruktion ist $\|L\|_{\infty} \leq n$ und mit der Definition

$$\rho = \max_{i,j \in \mathbb{R}} \frac{|a_{ij}^{(k)}|}{\|A\|_{\infty}} \quad \text{„Wachstumsfaktor“}$$

Zeigt man [GVL] ...

$$\|E\|_{\infty} \leq \rho n^3 \|A\|_{\infty} \text{ eps} + O(\text{eps}^2),$$

In der Praxis ist $\rho \approx 10$

schlechtester Fall bei Spaltenpivotisierung ist $\rho = 2^{n-1}$,

Mit totaler Pivotisierung erreicht man

$$|a_{ij}^{(k)}| \leq k^{1/2} (2 \cdot 3^{1/2} \dots k^{1/(k-1)})^{1/2} \max |a_{ij}|$$

also deutlich kleineres Wachstum.

Totale Pivotsierung

4
16.11.09

Wähle $r, s \in \{1, \dots, n\}$ so dass

$$|a_{rs}^{(k)}| \geq |a_{ij}^{(k)}| \quad \forall k \leq i, j \leq n$$

und erreiche durch Zeilen und Spaltenvertauschung, dass

$$\tilde{a}_{kk}^{(k)} = a_{rs}^{(k)}$$

In Matrixform:

(sollten wir schon haben)

↙ Rechtsform, Formulation, Umbenennung der Variablen.

Schritt 1: $G_1 P_{r_1} A P_{s_1} P_{s_1}^T x = G_1 P_{r_1} b$

Schritt 2: $G_2 P_{r_2} G_1 P_{r_1} A P_{s_1} P_{s_2} P_{s_2}^T P_{s_1}^T x = G_2 P_{r_2} G_1 P_{r_1} b$

Schritt n und Umformulierung

$$\underbrace{G'_n \dots G'_1 P_{r_n} \dots P_{r_1}}_P A \underbrace{P_{s_1} \dots P_{s_n}}_{Q^T} \underbrace{P_{s_1}^T \dots P_{s_n}^T}_Q x = G'_n \dots G'_1 P_{r_n} \dots P_{r_1} b$$

$= R$

und damit

$$PAQ^T z = Pb$$

$$\boxed{PAQ^T = LR} \quad z = Qx$$

Lösen des LGS gelingt dann mit

$$LRz = Pb$$

$$b' = Pb$$

$$Ly = b'$$

$$Rz = y$$

$$x = Q^T z$$

(Q ist orthogonal)

Aufwand der Pivotisierung:

5
16.11.09

- $n^2/2$ Vergleiche bei Spaltenpivotisierung

- $n^3/3$ Vergleiche bei totaler Pivotisierung

Da Spaltenzugriffe teurer als eigentliche Rechnungen \Rightarrow Verdopplung des Zeitbedarfs bei totaler Pivotisierung.

Praktische Erfahrung zeigt keine Vorteile für Rundungsfehler bei totaler Pivotisierung

\Rightarrow Spaltenpivotisierung mit Zeilen skalierung ist effizient und numerisch stabil in der Praxis.

Symmetrisch positiv definite Matrizen

Satz 4.15 Eine symmetrisch positiv definite Matrix $A \in \mathbb{R}^{n \times n}$ ist stets ohne ^{stabil} Pivotisierung LR-zerlegbar. Für ^{die} Diagonalelemente der im Eliminationsprozess auftretenden Matrizen gilt

$$a_{ii}^{(k)} \geq \lambda_{\min}(A), \quad k \leq i \leq n.$$

Beweis: Betrachte einen Schritt in der LR-Zerlegung

(a).

$$A = \begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix} \begin{matrix} 1 \\ n-1 \end{matrix}$$

Elimination der Spalte v liefert

$$\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix} \underbrace{\begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix}}_A = \begin{bmatrix} \alpha & v^T \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix}$$

$$= \begin{bmatrix} \alpha & 0 \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix} \begin{bmatrix} 1 & v^T/\alpha \\ 0 & I \end{bmatrix}$$

Mit $\begin{bmatrix} 1 & v^T/\alpha \\ 0 & I \end{bmatrix}^{-1} = \begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}$ folgt:

$$\underbrace{\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix}}_{X^T} \underbrace{\begin{bmatrix} \alpha & v^T \\ v & B \end{bmatrix}}_A \underbrace{\begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}}_X = \underbrace{\begin{bmatrix} \alpha & 0 \\ 0 & B - \frac{1}{\alpha} v v^T \end{bmatrix}}_{\hat{A}}$$

X hat vollen Rang, A ist s.p.d. nach Vor. $\Rightarrow X^T A X = \hat{A}$ ist s.p.d.
 $B - \frac{1}{\alpha} v v^T$ ist Hauptuntermatrix von \hat{A} und somit auch sym. pos. definit.

b) Es gilt (Rayleigh-Quotient):

$$(Ax, x)_2 \geq \lambda_{\min}(A) (x, x)_2$$

und damit für $x = e^{(i)}$ (kartesischer Einheitsvektor)

$$a_{ii} = (Ae^{(i)}, e^{(i)})_2 \geq \lambda_{\min}(A) \underbrace{(e^{(i)}, e^{(i)})_2}_{=1} = \lambda_{\min}(A).$$

Für die Diagonalelemente von $B^{-\frac{1}{\alpha}} vv^T$ gilt mit $\tilde{e}^{(i)} = \begin{pmatrix} 0 \\ e^{(i)} \end{pmatrix}$

$$\left(B^{-\frac{1}{\alpha}} vv^T \right)_{ii} = \left(\tilde{A} \tilde{e}^{(i)}, \tilde{e}^{(i)} \right)_2 = \left(X^T A X \tilde{e}^{(i)}, \tilde{e}^{(i)} \right)_2$$

$$= \left(A X \tilde{e}^{(i)}, X \tilde{e}^{(i)} \right)_2$$

$$\geq \lambda_{\min}(A) \left(X \tilde{e}^{(i)}, X \tilde{e}^{(i)} \right)_2 = \lambda_{\min}(A) \left(1 + \frac{v_i^2}{\alpha} \right)$$

$$\underbrace{\begin{bmatrix} 1 & -v^T/\alpha \\ 0 & I \end{bmatrix}}_X \underbrace{\begin{bmatrix} a \\ e^{(i)} \end{bmatrix}}_{\tilde{e}^{(i)}} = \begin{bmatrix} -\frac{v_i}{\alpha} \\ e^{(i)} \end{bmatrix} \geq \lambda_{\min}(A).$$

□

Dies zeigt, dass die Pivotelemente bei der Gaußelimination echt nach unten beschränkt bleiben.

ü) Rundungsfehleranalyse: Kann man $|\hat{C}|$ $|\hat{R}|$ abschätzen?
 → Panachev Lemma 4.1.

Cholesky-Zerlegung

P
16.11.09

Die Pivotelemente die bei der LR-Zerlegung auftreten sind nach 4.15 stets positiv.

Mit $D = \text{diag}(R)$ gilt

$$A = L D \underbrace{D^{-1} R}_U = L D U$$

U ist obere Dreiecksmatrix mit $u_{ii} = 1$. Wegen der Symmetrie gilt $U = L^T$, also

$$A = L D L^T.$$

Da $d_{ii} > 0$ ist die Matrix " $D^{1/2}$ "

$$(D^{1/2})_{ii} = \sqrt{d_{ii}}$$

wohldefiniert und es gilt

$$A = \underbrace{L}_{\tilde{L}} \underbrace{D^{1/2} D^{1/2}}_D L^T = \tilde{L} \tilde{L}^T.$$

Dies ist die sog. Cholesky-Zerlegung einer symmetrisch positiv definiten Matrix.

Ausnutzung der Symmetrie erlaubt die Berechnung der Cholesky-Zerlegung in $\frac{n^3}{3} + O(n^2)$ Operationen (Faktor 2 schneller).

Diagonaldominante Matrizen

9
16.11.09

Definition 4.16 Eine Matrix $A \in \mathbb{R}^{n \times n}$ heißt diagonaldominant falls

$$\sum_{j=1, j \neq i}^n |a_{ij}| \leq |a_{ii}| \quad i=1, \dots, n.$$

Satz 4.17 Diagonaldominante, reguläre Matrizen erlauben eine LR-Zerlegung ohne Pivotisierung.

Beweis. $A = \begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix}$, ein Schritt der GEM liefert

$$\begin{bmatrix} 1 & 0 \\ -v/\alpha & I \end{bmatrix} \begin{bmatrix} \alpha & w^T \\ v & B \end{bmatrix} = \begin{bmatrix} \alpha & w^T \\ 0 & B - \frac{1}{\alpha} v w^T \end{bmatrix}$$

Da $\alpha \neq 0$ wg. Diagonaldominanz ist dies wohldefiniert.

Aus der Diagonaldominanz von A ergibt sich:

$$\text{(Zeile 1)} \quad \sum_{j=1}^{n-1} |w_j| \leq |\alpha| \Leftrightarrow \sum_{j=1}^{n-1} \frac{|w_j|}{|\alpha|} \leq 1$$

$$\text{(Zeile 2...n)} \quad |v_i| + \sum_{j=1, j \neq i}^{n-1} |b_{ij}| \leq |b_{ii}|$$

Für $B - \frac{1}{\alpha} v w^T$ rechnen wir dann nach:

$$\begin{aligned} \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij} - \frac{1}{\alpha} v_i w_j| + \left| \frac{v_i w_i}{\alpha} \right| &\leq \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| + |v_i| \sum_{\substack{j=1 \\ j \neq i}}^{n-1} \frac{|w_j|}{|\alpha|} + |v_i| \frac{|w_i|}{|\alpha|} \\ &= \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| + |v_i| \underbrace{\sum_{j=1}^{n-1} \frac{|w_j|}{|\alpha|}}_{\leq 1} \leq |v_i| + \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij}| \leq |b_{ii}| \end{aligned}$$

$$\Leftrightarrow \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |b_{ij} - \frac{1}{\alpha} v_i w_j| \leq |b_{ii}| - \left| \frac{v_i w_i}{\alpha} \right| \leq \left| b_{ii} - \frac{v_i w_i}{\alpha} \right| = \left| (B - \frac{1}{\alpha} v w^T) \right|$$

Somit ist $B - \frac{1}{\alpha} v w^T$ wieder diagonaldominant.

$$\begin{aligned} |x| &= |x - y + y| \leq |x - y| + |y| \\ &\Leftrightarrow |x| - |y| \leq |x - y| \end{aligned}$$

Rangbestimmung

10
16.11.09

Sei $A \in \mathbb{K}^{n \times n}$. Gilt nach k -Schritten der Gauß-Elimination

$$A^{(k)} = \begin{array}{|cc|} \hline \overbrace{\quad}^k & \overbrace{\quad}^{n-k} \\ \hline U_{ii}^{(k)} & * \\ \hline 0 & 0 \\ \hline \end{array} \begin{array}{l} k \\ n-k \end{array}$$

mit $U_{ii}^{(k)} \neq 0$ so ist $\text{Rang}(A) = k$. . . Somit

- Kann die GE zu Ende geführt werden so gilt $k=n$, mithin $\text{Rang}(A)=n$.

- Rangbestimmung erfordert totales Pivoting, da die erste Spalte der $n-k \times n-k$ Restmatrix 0 sein kann, was aber nicht $\text{Rang}(A)=k$ bedeutet.

Inversenberechnung

Zur Berechnung der Inversen geht man so vor:

- Berechne LR-Zerlegung von A . Aufwand $\frac{2}{3}n^3 + O(n^2)$
bei Spaltenpivoting

- Für $i=1, \dots, n$ löse $Ax^{(i)} = e^{(i)}$.

Aufwand $n \cdot 2n^2$ für Lösen der Dreieckssysteme.

- $A^{-1} = [x^{(1)}, \dots, x^{(n)}]$ besteht spaltenweise aus den $x^{(i)}$

- Gesamt Aufwand ist $\frac{8}{3}n^3 + O(n^2)$.

Tridiagonalsysteme

11
16.11.09

Definition 4.18

$A \in \mathbb{R}^{n \times n}$ heißt Tridiagonalmatrix falls

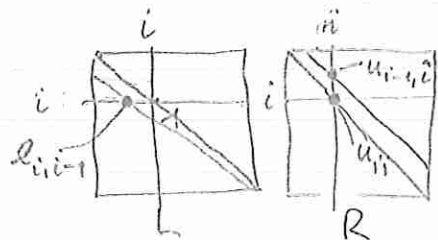
$$a_{ij} = 0 \text{ für } |i-j| > 1, \quad 1 \leq i \leq n$$

Eine Tridiagonalmatrix ist Spezialfall einer Bandmatrix:

$$a_{ij} = 0 \text{ für } j < i - m_l \text{ oder } j > i + m_r, \quad 1 \leq i \leq n$$

Die LR-Zerlegung einer Tridiagonalmatrix sei ohne Pivotisierung durchführbar. Dann gilt:

- L ist Tridiagonalmatrix
 - R ist Tridiagonalmatrix
- und damit:



zweite Spalte: $l_{i,i-1} \tau_{i-1,i} + \tau_{ii} = a_{ii} \Rightarrow \tau_{ii} = a_{ii} - l_{i,i-1} \tau_{i-1,i}$

$$l_{i,i-1} \tau_{i-1,i-1} = a_{i,i-1} \Rightarrow l_{i,i-1} = \frac{a_{i,i-1}}{u_{i-1,i-1}}$$

$$1 \cdot \tau_{i-1,i} = a_{i-1,i} \Rightarrow \tau_{i-1,i} = a_{i-1,i}$$

Nächstes Mal andere Reihenfolge!

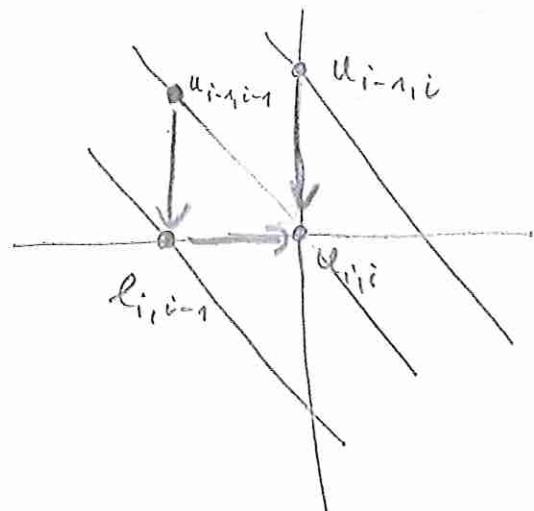
$$\tau_{11} = a_{11}, \quad \text{dabei}$$

for ($i = 2$ to n) do

$$\tau_{i-1,i} = a_{i-1,i}$$

$$l_{i,i-1} = \frac{a_{i,i-1}}{u_{i-1,i-1}}$$

$$\tau_{ii} = a_{ii} - l_{i,i-1} \tau_{i-1,i}$$



Nichtreguläre Systeme

12
17.11.09

Es sei nun $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ sowie $\text{Rang}(A)$ beliebig.
 $Ax = b$ hat dann genau eine, unendlich viele oder gar keine Lösung.

Einige Grundbegriffe aus der linearen Algebra: $m \times n$ $A: \mathbb{R}^n \rightarrow \mathbb{R}^m$
lineare Abb.

$$\text{Bild}(A) = \{y \in \mathbb{R}^m : y = Ax \text{ für ein } x \in \mathbb{R}^n\} \subseteq \mathbb{R}^m$$

$$\text{Kern}(A) = \{x \in \mathbb{R}^n : Ax = 0\} \subseteq \mathbb{R}^n$$

$$\text{Rang}(A) = \dim(\text{Bild}(A)) = \text{Rang}(A^T) = \dim(\text{Bild}(A^T))$$

(# l.u. Spalten = # l.u. Zeilen).

nicht trivial \rightarrow z.B. Beutelp. S. 101

Orthogonales Komplement:

$$m \times n \quad \text{Bild}(A)^\perp = \{y \in \mathbb{R}^m : (y, y')_2 = 0 \forall y' \in \text{Bild}(A)\}$$

$$(y, y')_2 = 0 \forall y' \in \text{Bild}(A)$$

$$\Leftrightarrow (y, Ax)_2 = 0 \forall x \in \mathbb{R}^n$$

$$\Leftrightarrow (A^T y)^T x = 0 \forall x \in \mathbb{R}^n \text{ geht für alle } x \text{ nur wenn } A^T y = 0$$

$$\Leftrightarrow y \in \text{Kern}(A^T)$$

und damit $\text{Bild}(A)^\perp = \text{Kern}(A^T)$.

Dies zeigt dann

$$\underbrace{\dim(\text{Bild}(A))}_{= \text{Rang}(A)} + \underbrace{\dim(\text{Kern}(A^T))}_{\text{Bild}(A)^\perp} = m$$

Zeilen

$$\underbrace{\dim(\text{Bild}(A^T))}_{= \text{Rang}(A)} + \dim(\text{Kern}(A)) = n$$

Den Lösungsbegriff für lineare Gleichungssysteme kann man auf die folgende Weise erweitern.

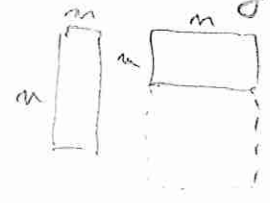
Satz 4.19 (Least Squares Lösung)

a) Es existiert ^{einz.} $\bar{x} \in \mathbb{R}^n$ so dass

$$\|A\bar{x} - b\|_2 = \min_{x \in \mathbb{R}^n} \|Ax - b\|_2.$$

b) Diese Bedingung ist äquivalent dazu dass \bar{x} Lösung von

$$A^T A \bar{x} = A^T b,$$



der sog. „Normalengleichung“.

c) Falls $\text{Rang}(A) = n$ (damit ist zwingend $m \geq n$) ist \bar{x} eindeutig bestimmt, sonst hat jede weitere Lösung die Form $\bar{x} + y$ mit $y \in \text{Kern}(A)$.

Beweis,

(b) \Rightarrow (a) \bar{x} sei Lösung der Normalengleichung. Für ein beliebiges $x \in \mathbb{R}^n$ gilt:

$$\begin{aligned} \|Ax - b\|_2^2 &= \|A(x - \bar{x} + \bar{x}) - b\|_2^2 = \underbrace{\|A\bar{x} - b\|_2^2}_{\in \text{Bild}(A)} + 2 \underbrace{\langle A\bar{x} - b, A(x - \bar{x}) \rangle_2}_{\substack{\in \text{Kern}(A) \perp \text{Bild}(A) \\ = \text{Bild}(A)^\perp}} + \underbrace{\|A(x - \bar{x})\|_2^2}_{\geq 0} \\ &= \|A\bar{x} - b\|_2^2 + \|A(x - \bar{x})\|_2^2 \\ &\geq \|A\bar{x} - b\|_2^2 \end{aligned}$$

da $A^T(A\bar{x} - b) = 0$

Damit erfüllt \bar{x} die Minimalitätsbedingung.

(a) \Rightarrow (b) Setze $F: \mathbb{R}^n \rightarrow \mathbb{R}$, $F(x) = \|Ax - b\|_2^2$. Notwendige Bedingung für ein Minimum ist $\nabla F(\bar{x}) = 0 \Leftrightarrow \frac{\partial F}{\partial x_k}(\bar{x}) = 0 \quad \forall k = 1, \dots, n$.

$$\begin{aligned} \frac{\partial F(\bar{x})}{\partial x_k} &= \frac{\partial}{\partial x_k} (Ax - b, Ax - b)_2 \Big|_{x=\bar{x}} = \frac{\partial}{\partial x_k} \left(\sum_{i=1}^m \left[\sum_{j=1}^n a_{ij} x_j - b_i \right]^2 \right) \Big|_{x=\bar{x}} \\ &= \left(\sum_{i=1}^m 2 \cdot \left[\sum_{j=1}^n a_{ij} x_j - b_i \right] a_{ik} \right) \Big|_{x=\bar{x}} = 2 \sum_{i=1}^m (A^T)_{ki} \left[\sum_{j=1}^n a_{ij} \bar{x}_j - b_i \right] \\ &= 2 (A^T (A\bar{x} - b))_k \quad \text{Nach Parametern.} \end{aligned}$$

Und damit $\nabla F(\bar{x}) = 2 A^T (A\bar{x} - b) \stackrel{!}{=} 0 \Leftrightarrow A^T A \bar{x} = A^T b$,
die Normalengleichung.

(c) Lösbarkeit der Normalgleichung.

14
19.11.09

$$\mathbb{R}^m = \text{Bild}(A) \oplus \text{Bild}(A)^\perp,$$

d.h. jedes $b \in \mathbb{R}^m$ besitzt eine eindeutige Zerlegung

$$b = s + r \text{ mit } s \in \text{Bild}(A), r \in \text{Bild}(A)^\perp = \text{Kern}(A^T)$$

Da $s \in \text{Bild}(A)$ gibt es $\bar{x} \in \mathbb{R}^n$ mit $A\bar{x} = s$. Für dieses gilt dann

$$A^T A \bar{x} = A^T s = A^T s + \underbrace{A^T r}_{=0 \text{ da } r \in \text{Kern}(A^T)} = A^T b,$$

also löst dieses \bar{x} auch die Normalgleichung.

Sei $\text{Rang}(A) = n$. und damit $m \geq n$ (wg. $\text{Rang}(A) \leq \min(m, n)$).

Betrachte $A^T A x = 0$ also $\text{Kern}(A^T A)$.

$$(*) A^T A x = 0 \Leftrightarrow A^T y = 0 \wedge y \in \text{Bild}(A)$$

Der alte, durchgestrichene Text bleibt gültig, der neue ist falsch! (26.6.2013)

~~Nun ist aber $\text{Kern}(A^T) \perp \text{Bild}(A)$ d.h. $\text{Kern}(A^T) \cap \text{Bild}(A) = \{0\}$.
d.h. (*) geht nur für $y=0$. Da $\text{Rang}(A)=n$ bedeutet dies $x=0$ also $A^T A$ regulär.~~

(Alternative: $A^T A$ ist symmetrisch und positiv definit wg. max. Rang).

Sei $\text{Rang}(A) < n$. Sei x_1 eine weitere Lösung der Normalgleichung

(also $A^T A x_1 = A^T b$). Dann gilt

$$b = \underbrace{A x_1}_{\in \text{Bild}(A)} + \underbrace{(b - A x_1)}_{\in \text{Kern}(A^T), \text{ da } A^T(b - A x_1) = A^T b - A^T A x_1 \stackrel{\downarrow}{=} 0}$$

Da die Zerlegung $\mathbb{R}^m = \text{Bild}(A) \oplus \text{Bild}(A)^\perp$ eindeutig ist muss

$$\hookrightarrow \text{d.h. } b = A x_1 + (b - A x_1) = \underbrace{A \bar{x}} + \underbrace{(b - A \bar{x})}$$

$$A x_1 = A \bar{x} \text{ sein}$$

und das heißt $x_1 = \bar{x} + y$ und

$$A x_1 = A \bar{x} - A y \stackrel{!}{=} A \bar{x} \Leftrightarrow \boxed{A y = 0} \text{ (was die Behauptung war) } \quad \square$$

Bemerkung: $A^T A$ ist symmetrisch und positiv semidefinit.

Lösung prinzipiell mit Cholesky-Zerlegung möglich.

QR-Zerlegung

1
12.6.13

Definition 4.20.

Sei $v \in \mathbb{R}^n$ gegeben, dann heißt die Matrix

$$Q_v = I - 2 \frac{v v^T}{v^T v} \quad \text{äußeres Produkt}$$

Householdermatrix oder Householderreflexion \square

Wir rechnen einige Eigenschaften dieser Matrix nach:

1) Q_v ist orthogonal, d.h. $Q_v^T Q_v = I$, $Q^T = Q$

$$Q^T = I - 2 \frac{v v^T}{v^T v} = I - \frac{2}{v^T v} (v v^T)^T = I - \frac{2}{v^T v} v v^T = Q$$

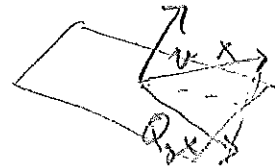
$$\begin{aligned} Q^T Q &= Q^2 = \left(I - 2 \frac{v v^T}{v^T v} \right) \left(I - 2 \frac{v v^T}{v^T v} \right) = I - 4 \frac{v v^T}{v^T v} + 4 \frac{(v v^T)(v v^T)}{v^T v v^T v} \\ &= I - 4 \frac{v v^T}{v^T v} + 4 \frac{v (v^T v) v^T}{(v^T v)(v^T v)} = I. \end{aligned}$$

2) $Q_v x = x - 2 \frac{v v^T}{v^T v} x = x - 2 \frac{v^T x}{v^T v} v$

$$\left((v v^T) x \right)_i = \sum_{j=1}^n v_i v_j x_j = (v^T x) v_i$$

3) $x = \alpha v \Rightarrow Q_v x = -x$ "Reflexion an der Ebene mit Normalen v "

$$Q_v x = \alpha v - 2 \alpha \frac{v^T v}{v^T v} v = -\alpha v = -x$$



4) $(x, v)_2 = v^T x = 0 \Rightarrow Q_v x = x$

$$Q_v x = x - 2 \frac{v^T x}{v^T v} v = x$$

5) zusammen: $x = \alpha v + w$ mit $(w, v)_2 = 0 \Rightarrow Q_v x = -\alpha v + w$

Q_v kann in ähnlicher Weise wie die Frobenius-Matrizen benutzt werden um eine Matrix auf obere Dreiecksform zu transformieren.

Problem. Wähle v so, dass

$$Q_v A = \begin{pmatrix} \alpha * & \dots & * \\ 0 & & 1 \\ \vdots & & \vdots \\ 0 & * & \dots & * \end{pmatrix}$$

Sei x die erste Spalte von A so bestimme v sodass

$$Q_v x = \beta e_1$$

$$\Leftrightarrow x - 2 \frac{v^T x}{v^T v} v = \beta e_1$$

Mit dem Ansatz $v = x + \alpha e_1$ erhalten wir

$$Q_v x = x - 2 \frac{v^T x}{v^T v} (x + \alpha e_1) = \left(1 - 2 \frac{v^T x}{v^T v}\right) x - 2\alpha \frac{v^T x}{v^T v} e_1$$

Um den ersten Term zu eliminieren:

$$1 - 2 \frac{v^T x}{v^T v} = 1 - 2 \frac{(x + \alpha e_1)^T x}{(x + \alpha e_1)^T (x + \alpha e_1)}$$

$$= 1 - 2 \frac{x^T x + \alpha x_1}{x^T x + 2\alpha x_1 + \alpha^2}$$

$$= \frac{x^T x + 2\alpha x_1 + \alpha^2 - 2x^T x - 2\alpha x_1}{x^T x + 2\alpha x_1 + \alpha^2} = \frac{-x^T x + \alpha^2}{x^T x + 2\alpha x_1 + \alpha^2} \stackrel{!}{=} 0$$

$$\Leftrightarrow \alpha^2 = x^T x = \|x\|_2^2 \quad \Leftrightarrow \boxed{\alpha = \pm \|x\|_2}$$

hier.

d.h. für $v = x \pm \|x\|_2 e_1$ gilt $Q_v x = \mp 2 \|x\|_2 e_1$

Welches Vorzeichen wählt man? Falls schon $x = \gamma e_1$ gilt, dann wäre $v = \gamma e_1 \pm |\gamma| e_1$, damit die Norm von v möglichst groß wird wählt man $+$ für $\gamma > 0$ und $-$ für $\gamma < 0$.

Satz 4.21 (QR-Zerlegung)

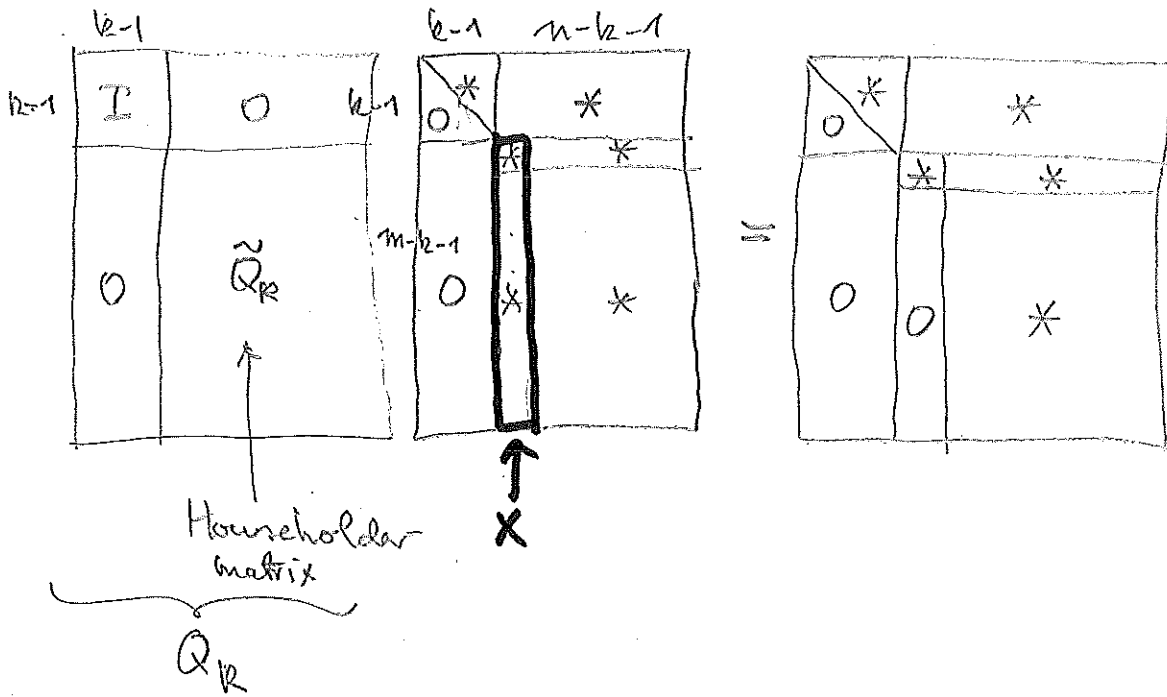
3
12.6.13

Zu jeder Matrix $A \in \mathbb{R}^{m \times n}$ mit $m \geq n$ und $\text{Rang}(A) = n$ existiert eine orthogonale Matrix $Q \in \mathbb{R}^{m \times m}$ und eine obere Dreiecksmatrix $R \in \mathbb{R}^{m \times n}$ so dass

$$A = QR$$

Die ersten n Spalten von Q bilden eine orthonormale Basis von A .

Beweis. Im Schritt $k = 1, \dots,$



- Q_k ist orthogonal, da \tilde{Q}_k Householder matrix. Somit gilt nach $n-1$ Schritten:

$$Q_{n-1} \dots Q_1 A = R$$

da $Q_k^{-1} = Q_k^T = Q_k$ gilt

$$A = \underbrace{Q_1 \dots Q_{n-1}}_{=: Q} R = QR$$

Q ist orthogonal, da $Q^T = (Q_1 \dots Q_{n-1})^T = Q_1^T \dots Q_{n-1}^T = Q_{n-1} \dots Q_1$
und $Q^T Q = Q_{n-1} \dots Q_1 Q_1 \dots Q_{n-1} = I$. □

Anwendung zur Lösung von $A^T A x = A^T b$.

4
12.6.13

$$A = QR \Rightarrow A^T A = (QR)^T QR = R^T \underbrace{Q^T Q}_{=I} R = R^T R$$

also $R^T R x = A^T b$
 $\underbrace{R^T R}_{=y}$

Somit 1) $R^T y = A^T b$

2) $R x = y$

$\Rightarrow \sim A^T A$ muss nicht explizit aufgestellt werden (n^3 Op, Fill-in).

- Vorteile bei Rundungsfehlerfortpflanzung

- Anwendung bei Berechnung von Eigenwerten (QR-Iteration)

Anwendung: Gaußsche Ausgleichsrechnung.

15
17.11.09

Gegeben:

(i) n Funktionen $u_1, \dots, u_n: \mathbb{R} \rightarrow \mathbb{R}$, sowie

(ii) m Datenpunkte $(x_i, y_i) \in \mathbb{R}^2$, $1 \leq i \leq m \geq n$

Gesucht: n Koeffizienten c_1, \dots, c_n so dass

↑ Bei $m = n$ ungl. Eindeutig.

$$u(x) = \sum_{j=1}^n c_j u_j(x)$$

und

$$(*) \quad \sum_{i=1}^m (u(x_i) - y_i)^2 \rightarrow \text{minimal.}$$

Dies lässt sich folgendermaßen formulieren:

$$c = (c_1, \dots, c_n)^T$$

$$y = (y_1, \dots, y_m)^T$$

$$a_{ij} = u_j(x_i)$$

Finde $c \in \mathbb{R}^n$ so dass $\|Ac - y\|_2^2$ minimal.

$$\text{Denn } \sum_{i=1}^m (u(x_i) - y_i)^2 = \sum_{i=1}^m \left(\underbrace{\sum_{j=1}^n c_j u_j(x_i) - y_i}_{(Ac - y)_i} \right)^2$$

$$= \|Ac - y\|_2^2.$$

Somit sind nach Satz ~~5.17~~^{4.19} die gesuchten c_i Lösung der Normalgleichung

$$A^T A c = A^T y.$$

Lösung: z.B. mit Cholesky-Zerl.

(ü): Zeige $\text{cond}_2(A^T A) = \text{cond}_2(A)^2$.

