

4. VL, 27.4.2009, G. Bärwolf

Einführung in die Numerische Mathematik

Theorem 1.27

Für die Berechnung von speziellen induzierten Matrixnormen gilt (A ($m \times n$)-Matrix, reell)

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}| \quad \text{Spaltensummennorm} \quad (1.25)$$

$$\|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| \quad \text{Zeilensummennorm} \quad (1.26)$$

$$\|A\|_2 = \sqrt{\lambda_{\max}} \quad \text{mit } \lambda_{\max} \text{ als größtes EW von } A^T A \quad (1.27)$$

Zum Beweis

(1.25) und (1.26) sollten als Übung nachgewiesen werden.

Für den Nachweis von (1.27) überlegt man, dass $A^T A$ ähnlich einer Diagonalmatrix mit den EW von $A^T A$ als Diagonalelementen ist. Die EW sind nichtnegativ und aus der Def. von $\|A\|_2$ folgt dann schließlich (1.27).

Definition 1.28

Unter dem Absolutbetrag einer Matrix $A \in \mathbb{C}^{m \times n}$ versteht man die Matrix

$$|A| = B \quad \text{mit } b_{ij} = |a_{ij}|,$$

also die Matrix mit den Absolutbeträgen ihrer Elemente. Gilt für $A, B \in \mathbb{R}^{m \times n}$ und $a_{ij} = b_{ij}$

Bemerkung 1.19

Für die "Beträge" von Matrizen gelten die Beziehungen

$$i) |A+B| \leq |A| + |B|$$

$$ii) |A \cdot B| \leq |A| |B|$$

$$iii) A \leq B, C \geq 0, D \geq 0 \rightarrow CAD \leq CBD$$

$$iv) \|A\|_p \leq \| |A| \|_p, \quad p \geq 1, p \in \mathbb{N}$$

$$v) \begin{cases} \|A\| = \||A|\| \text{ für } \|\cdot\|_1, \|\cdot\|_\infty, \|\cdot\|_F \\ |A| \leq |B| \rightarrow \|A\| \leq \|B\| \text{ für diese Normen} \end{cases}$$

Nachweis von i) - v) sind größtenteils trivial

Wir kehren zurück zur Fehlerfortpflanzung bei der Lösung linearer Gleichungssysteme. Seien A und \vec{b} durch $\Delta A = \varepsilon F$ bzw. $\vec{A}b = \varepsilon \vec{f}$ gestört, also statt $A\vec{x} = \vec{b}$ wird

$$(A + \varepsilon F)\vec{x}(\varepsilon) = \vec{b} + \varepsilon \vec{f} \quad (1.28)$$

betrachtet. F, \vec{f}, ε sind Matrizen und für $\varepsilon = 0$ ergibt sich das eigentliche Problem. A sei regulär. Der Satz über implizite Funktionen liefert dann die Existenz und Differenzierbarkeit der Abb. $\varepsilon \rightarrow \vec{x}(\varepsilon)$ und $\vec{x} = \vec{x}(0)$,

$$\dot{\vec{x}}(0) = A^{-1}(\vec{f} - F\vec{x})$$

Für die Taylorreihe folgt

$$\vec{x}(\varepsilon) = \vec{x} + \varepsilon \dot{\vec{x}}(0) + \mathcal{O}(\varepsilon^2),$$

also $\vec{x}(\varepsilon) - \vec{x} = \varepsilon \dot{\vec{x}}(0) + \mathcal{O}(\varepsilon^2) = \varepsilon A^{-1}(\vec{f} - F\vec{x}) + \mathcal{O}(\varepsilon^2)$ (3)

für den relativen Fehler gilt dann

$$\frac{\|\vec{x}(\varepsilon) - \vec{x}\|}{\|\vec{x}\|} \leq \varepsilon \frac{\|A^{-1}(\vec{f} - F\vec{x})\|}{\|\vec{x}\|} + \mathcal{O}(\varepsilon^2),$$

und im Fall einer Submultiplikation (normierter) Matrixnorm von A

$$\frac{\|\vec{x}(\varepsilon) - \vec{x}\|}{\|\vec{x}\|} \leq \varepsilon \|A^{-1}\| \left\{ \frac{\|\vec{f}\|}{\|\vec{x}\|} + \|F\| \right\} + \mathcal{O}(\varepsilon^2) \quad (1.29)$$

(1.29) sagt noch nichts über die Auswirkung der relativen Fehler

$$\frac{\|\Delta A\|}{\|A\|} \quad \text{bzw.} \quad \frac{\|\Delta \vec{b}\|}{\|\vec{b}\|}$$

aus, was aber von Interesse ist. Mit $\vec{b} = A\vec{x}$ folgt $\|\vec{b}\| \leq \|A\| \|\vec{x}\|$, also

$$\frac{1}{\|\vec{x}\|} \leq \frac{\|A\|}{\|\vec{b}\|}$$

und mit $\|F\| = \frac{\|F\|}{\|A\|} \|A\|$ folgt aus (1.29)

$$\begin{aligned} \frac{\|\vec{x}(\varepsilon) - \vec{x}\|}{\|\vec{x}\|} &\leq \|A\| \|A^{-1}\| \left\{ \varepsilon \frac{\|F\|}{\|A\|} + \varepsilon \frac{\|\vec{f}\|}{\|\vec{b}\|} \right\} + \mathcal{O}(\varepsilon^2) \\ &= \kappa(A) \left\{ \frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta \vec{b}\|}{\|\vec{b}\|} \right\} + \mathcal{O}(\varepsilon^2) \\ &\leq \kappa(A) \left\{ \frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta \vec{b}\|}{\|\vec{b}\|} \right\} \quad (1.30) \end{aligned}$$

(1.30) gilt für kleine ε . Im Folgenden soll $\textcircled{4}$
eine genauere Analyse vorgenommen werden

Hilfssatz 1.20

$F \in \mathbb{R}^{n \times n}$ und $\|\cdot\|$ submultiplikativ (normiert)
Für $\|F\| < 1$ ist $E - F$ nicht-singulär,
und es gilt

$$(E - F)^{-1} = \sum_{k=0}^{\infty} F^k$$

sowie

$$\|(E - F)^{-1}\| \leq \frac{1}{1 - \|F\|}$$

Beweis in Analogie zur geom. Reihe,
~~oder Verallgemeinerung derselben~~ \rightarrow

Hilfssatz 1.21

Es gelte $A\vec{x} = \vec{b}$ mit regulärer $A \in \mathbb{R}^{n \times n}$
und $\vec{b} \in \mathbb{R}^n, \vec{b} \neq \vec{0}$. $F = \Delta A$ und $\vec{f} = \Delta \vec{b}$ seien
Störungen mit relativem Fehler $\leq \delta$, d.h.

$$\frac{\|\Delta A\|}{\|A\|} \leq \delta, \quad \frac{\|\Delta \vec{b}\|}{\|\vec{b}\|} \leq \delta$$

Falls $\delta K(A) = r < 1$ gilt,

ist $A + \Delta A$ regulär und für die Lösung \vec{x} von
 $(A + \Delta A)\vec{x} = \vec{b} + \Delta \vec{b}$ gilt

$$\frac{\|\vec{x}\|}{\|\vec{x}\|} \leq \frac{1+r}{1-r}$$

Beweis Es gilt

5

$$\|A^{-1}F\| \leq \|A^{-1}\| \|F\| \leq \underbrace{\sigma \|A^{-1}\|}_{\kappa(A)} \|A\| = r < 1$$

$$\leadsto A + F = A(E + A^{-1}F) =: A(E - \tilde{F}) \quad \text{mit } \|\tilde{F}\| < 1$$

$$(E - \tilde{F})^{-1} \text{ ex. nach H.S. 1.20 } \leadsto$$

$$(A + F)^{-1} = (E - \tilde{F})^{-1} A^{-1} \text{ ex.}$$

Nach Def. von \tilde{x} folgt

$$(A + F)\tilde{x} = \vec{b} + \vec{f} = A\vec{x} + \vec{f} \quad | \cdot A^{-1}$$

$$\leadsto (E + A^{-1}F)\tilde{x} = \vec{x} + A^{-1}\vec{f}$$

$$\leadsto \|\tilde{x}\| \leq \underbrace{\|(E + A^{-1}F)^{-1}\|}_{\|(E - \tilde{F})^{-1}\|} \|\vec{x} + A^{-1}\vec{f}\|$$

$$, \|(E - \tilde{F})^{-1}\| \leq \frac{1}{1 - \|\tilde{F}\|} \stackrel{\leq r}{\leq} \frac{1}{1-r}$$

$$\leadsto \|\tilde{x}\| \leq \frac{1}{1-r} (\|\vec{x}\| + \|A^{-1}\| \|\vec{f}\|)$$

$$\leq \sigma \|\vec{b}\| = \sigma \|A\vec{x}\| \leq$$

$$\sigma \|A\| \|\vec{x}\|$$

$$\leq \frac{1}{1-r} (\|\vec{x}\| + \underbrace{\sigma \|A\| \|A^{-1}\|}_{r} \|\vec{x}\|) = \frac{1+r}{1-r} \|\vec{x}\|$$

Zur Abschätzung des rel. Fehlers
betrachten wir

$$\left. \begin{array}{l} A\tilde{x} + F\tilde{x} = \vec{b} + \vec{f} \\ A\vec{x} = \vec{b} \end{array} \right\} \Rightarrow A(\tilde{x} - \vec{x}) = \vec{f} - F\tilde{x}$$

$$\Rightarrow \tilde{x} - \vec{x} = A^{-1}\vec{f} - A^{-1}F\tilde{x} \quad (1.31)$$

Theorem 1.22

Unter den Voraussetzungen des H.S. 1.21 gilt

$$\frac{\|\vec{x} - \tilde{\vec{x}}\|}{\|\vec{x}\|} \leq \frac{2v}{1-v}, \quad v = \delta \kappa(A) \quad (1.32)$$

Beweis

Es ergibt sich aus (1.31)

$$\begin{aligned} \|\tilde{\vec{x}} - \vec{x}\| &\leq \|A^{-1} \vec{f}\| + \|A^{-1} F \tilde{\vec{x}}\| \\ &\leq \|A^{-1}\| \|\vec{f}\| + \|A^{-1} F\| \|\tilde{\vec{x}}\| \\ &\leq \delta \|\vec{b}\| \leq \delta \|A\| \|\vec{x}\| \\ &\leq \delta \kappa(A) \|\vec{x}\| + v \frac{1+v}{1-v} \|\vec{x}\| \\ &= v \left(1 + \frac{1+v}{1-v}\right) \|\vec{x}\| = \frac{v \cdot 2}{1-v} \|\vec{x}\| \end{aligned}$$

In den bisherigen Abschätzungen wurde der Fehler F der Matrix A in der Norm $\|F\|$ betrachtet und der Fehler von \vec{b} ebenso in der Vektornorm $\|\vec{f}\|$, d.h.

$$\|F\| \leq \delta \|A\| \quad \text{und} \quad \|\vec{f}\| \leq \delta \|\vec{b}\|$$

wurden vorausgesetzt. Betrachtet man nun die Komponenten, d.h. fordert man

$$|f_{ij}| \leq \delta |a_{ij}| \quad \text{für } i, j = 1, 2, \dots, n$$

dann lassen sich die Abschätzungen des H.S. 1.21 und des Theorems 1.22 verbessern.

Es gelten

Hilfsabsatz 1.23

(7)

Sei $A \in \mathbb{R}^{n \times n}$ regulär und $\vec{b} \in \mathbb{R}^n$, $\vec{b} \neq \vec{0}$
 Es gelte $A\vec{x} = \vec{b}$ und $(A + \Delta A)\vec{\tilde{x}} = \vec{b} + \Delta\vec{b}$
 und die Abschätzungen
 $|\Delta A| \leq \sigma |A|$, $|\Delta\vec{b}| \leq \sigma |\vec{b}|$.

Ist die Bedingung

$$\sigma \| |A^{-1}| |A| \|_M = r < 1$$

erfüllt, so ist $A + \Delta A$ regulär und es gilt

$$\frac{\|\vec{\tilde{x}}\|}{\|\vec{x}\|} \leq \frac{1+r}{1-r}$$

($\| \cdot \|_M$ ist eine der Matrixnormen $\| \cdot \|_{\infty}$, $\| \cdot \|_1$, $\| \cdot \|_F$ und $\| \cdot \|$ die damit verträgliche Normen).

Theorem 1.24

Unter den obigen Voraussetzungen gilt

$$\frac{\|\vec{\tilde{x}} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{2r}{1-r} \quad (1.33)$$

(für die obigen Normen)

Bemerkung 1.25

Die Aussage der Theoreme 1.22 und 1.24 findet man auch oft in der Form

$$\frac{\|\vec{\tilde{x}} - \vec{x}\|}{\|\vec{x}\|} \leq \frac{\kappa(A)}{1 - \|A^{-1}\| \|\Delta A\|} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta\vec{b}\|}{\|\vec{b}\|} \right)$$

Beispiel

$$(i) A \vec{x} = \vec{b} \iff \begin{bmatrix} 1 & 0 \\ 0 & 10^{-6} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 \\ 10^{-6} \end{bmatrix}, A^{-1} = \begin{bmatrix} 1 & 0 \\ 0 & 10^6 \end{bmatrix}$$

$$\|A\|_{\infty} = 1, \|A^{-1}\|_{\infty} = 10^6, \kappa_{\infty}(A) = 10^6$$

$$\Delta A = 0, \Delta \vec{b} = \begin{bmatrix} 10^{-7} \\ 0 \end{bmatrix} \Rightarrow \|\Delta \vec{b}\|_{\infty} = 10^{-7}, \|\vec{b}\|_{\infty} = 1$$

$$\rightarrow \frac{\|\Delta \vec{b}\|_{\infty}}{\|\vec{b}\|_{\infty}} \leq 10^{-7} = \delta$$

$$r = \delta \kappa_{\infty}(A) = 10^{-7} \cdot 10^6 = 0,1$$

Aus (1.32) ergibt

$$\frac{\|\Delta \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} \leq \frac{2 \cdot 0,1}{0,9} = \frac{2}{9} \quad \text{aber} \quad \|\Delta \vec{x}\|_{\infty} = |\Delta x_2| = 10^{-7}$$

$$\text{und} \quad \frac{\|\Delta \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} = 10^{-7}$$

→ Abschätzung wird zu pessimistisch

(ii) Anwendung des Theorems 1.24

$$|A| = A, |A^{-1}| = A^{-1}, |A| |A^{-1}| = E$$

$$\rightarrow \| |A| \cdot |A^{-1}| \|_{\infty} = 1 = \kappa_{\infty}(A)$$

$$\rightarrow \delta \cdot 1 = r < 1$$

$$\Delta A = 0, \Delta \vec{b} = \begin{bmatrix} 10^{-7} \\ 0 \end{bmatrix}$$

Abschätzung (1.32) bzw. (1.33) ergibt

$$\frac{\|\Delta \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} \leq \frac{2 \cdot 10^{-7}}{1 - 10^{-7}} \approx 2 \cdot 10^{-7}$$

also wird der tatsächliche rel. Fehler wesentlich realistischer abgeschätzt.

$$\frac{\|\Delta \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} = 10^{-7}$$

(9)

Schlussfolgerungen für die Bedingung, mit
Gleitpunkt-Zahlensystemen, wobei exakte
Bedingung vorausgesetzt wird

$$A \rightarrow \text{rd}(A) = A + \Delta$$

$$\vec{b} \rightarrow \text{rd}(\vec{b}) = \vec{b} + \vec{\delta}$$

Ergebnisse $\|\Delta\|_{\infty} \leq \text{eps} \|A\|_{\infty}, \|\vec{\delta}\|_{\infty} \leq \text{eps} \|\vec{b}\|_{\infty}$

$$\rightarrow (A + \Delta) \vec{x} = \vec{b} + \vec{\delta},$$

nach Theorem 1.22 folgt mit $\delta = \text{eps}$

$$\frac{\|\vec{x} - \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} \leq \frac{2 \text{eps} \kappa_{\infty}(A)}{1 - \text{eps} \kappa_{\infty}(A)}$$

angenommen $\text{eps} \kappa_{\infty}(A) \leq \frac{1}{2}$, dann gilt
bei exakter Rechnung

$$\frac{\|\vec{x} - \vec{x}\|_{\infty}}{\|\vec{x}\|_{\infty}} \leq \frac{2 \cdot \text{eps} \kappa_{\infty}(A)}{\frac{1}{2}} = 4 \text{eps} \kappa_{\infty}(A)$$

2. Lösung linearer Gleichungssysteme (durch sukzessive Lösung von Dreieckssystemen)

LR-Zerlegung

zu lösen ist $A\vec{x} = \vec{b}$

Dazu soll A als Produkt einer unteren Dreiecksmatrix L und einer oberen Dreiecksmatrix R geschrieben werden, d.h. man hat

$$A\vec{x} = \vec{b} \iff L \underbrace{R\vec{x}}_{\vec{y}} = \vec{b}$$

und löse zuerst

$$L\vec{y} = \vec{b}$$

und danach

$$R\vec{x} = \vec{y}$$

Realisierung mit dem Gaußsche Eliminationsverf.

Grundprinzip:

Rangschaltende Manipulationen der Matrix $[A|\vec{b}]$ durch linearkombinationen von Zeilen

$$L_{ij}(\lambda) = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & \ddots & \\ i & & & & \lambda & & \\ & & & & & \ddots & \\ & & & & & & 1 \end{pmatrix}$$

L_{ij} besteht aus der Einheitsmatrix vom Typ $n \times n$, wobei an der Position (i,j) die Zahl λ statt einer 0 eingefügt wurde

(Multiplikation $L_{ij}(\lambda) A$)

bewirkt die Addition des λ -fachen

der j -ten Zeile von A zur i -ten Zeile von A

d.h. durch geeignete Wahl von λ erzeugt man in

$$\tilde{A} = L_{ij}(\lambda)A$$

an der Position (ij) z.B. auch eine Null (A vom Typ $n \times n$).

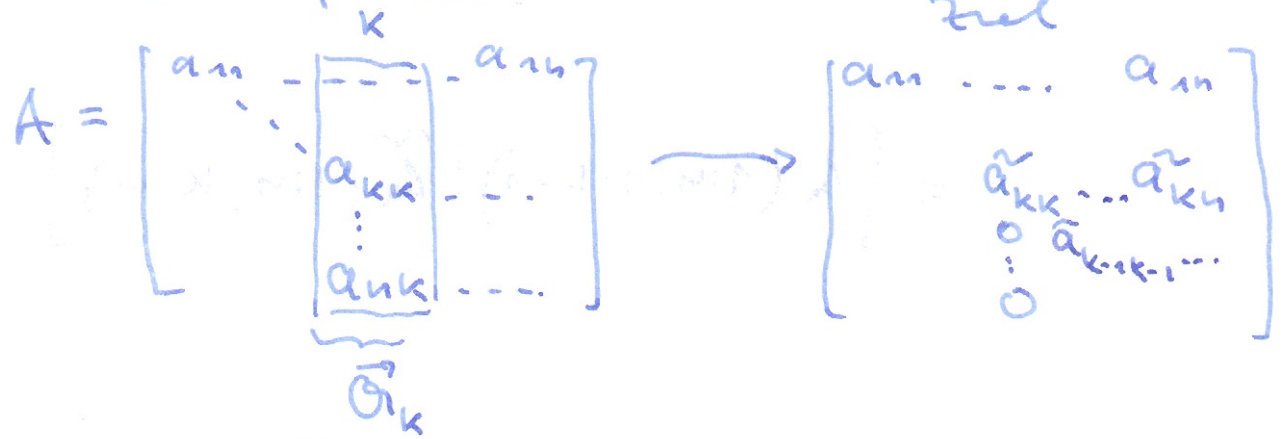
$L_{ij}(\lambda)$ hat den Rang n und die Determinante 1
 $\rightarrow \text{rg}(\tilde{A}) = \text{rg}(A)$

Durch mehrfache Multiplikation mit

$$L_{jk}, \quad j = k+1, \dots, n$$

erhält man bei geeigneter Wahl der λ unterhalb von \tilde{a}_{kk} Null-Einträge

Nun etwas präziser



Vor. $a_{kk} \neq 0$

Setzen $\vec{t} = \vec{t}^{(k)} = t^{(k)} \cdot \vec{e}_k = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ t_{k+k} \\ \vdots \\ t_{nk} \end{bmatrix}$ } k Komponenten

mit $t_{ik} = \begin{cases} 0 & i = 1, \dots, k \\ \frac{a_{ik}}{a_{kk}} & i = k+1, \dots, n \end{cases}$

\vec{e}_k sei der k -te Standardbasisvektor

Definition 2.1

$$M_k := \begin{bmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & t_{k+1,k} & & \\ & & \vdots & & \\ & & t_{n,k} & & \\ & & & & \ddots & \\ & & & & & 1 \end{bmatrix}$$

Spalte k

Frobenius-Matrix
"Gauß-Transform."

Man überlegt sich, dass M_k das Produkt der oben diskutierten Matrizen $L_{jk}(-t_j)$, $j = k+1, \dots, n$ ist.

Eigenwerte von M_k

$$M_k = E - \vec{t}^{(k)} \vec{e}_k^T,$$

$$M_k \vec{a}_k = \begin{bmatrix} a_{k1} \\ \vdots \\ a_{kk} \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

, a_{k1}, \dots, a_{kk}
bleiben bei der
Multiplikation
mit M_k unverändert

$$\text{rg}(M_k) = n,$$

$$\det(M_k) = 1$$

$$M_k^{-1} = E + \vec{t}^{(k)} \vec{e}_k^T, \text{ da}$$

$$M_k^{-1} M_k = E - \underbrace{\vec{t}^{(k)} \vec{e}_k^T \vec{t}^{(k)} \vec{e}_k^T}_{=0} = E$$

Wenn alles gut geht,

d.h. wenn jeweils $\tilde{a}_{k-1, k-1} \neq 0$ ist,

(15)

Dann erhält man nach der Multiplikation von A mit den Frobenius-Matrizen M_1, \dots, M_{n-1} , also

$$M_{n-1} M_{n-2} \dots M_1 A = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & \dots & r_{2n} \\ & & \ddots & \vdots \\ 0 & & & r_{nn} \end{bmatrix} := R$$

eine obere Dreiecksmatrix R , außerdem hat die Matrix

$$M_{n-1} M_{n-2} \dots M_1$$



die inverse Matrix

$$L = M_1^{-1} M_2^{-1} \dots M_{n-1}^{-1} = \begin{bmatrix} 1 & & & & \\ t_{21} & 1 & & & \\ t_{31} & t_{32} & \dots & & \\ \vdots & \vdots & \ddots & \ddots & \\ t_{n1} & t_{n2} & & & t_{nn}^{-1} \end{bmatrix}$$

so dass schließlich mit

$$A = LR$$

eine LR-Zerlegung vorliegt