

Numerische Verfahren

Jens-Peter M. Zemke
zemke@tu-harburg.de

Institut für Numerische Simulation
Technische Universität Hamburg-Harburg

29.04.2008



Lineare Systeme

- Zerlegung regulärer Matrizen
- Modifikationen des Gaußschen Verfahrens
- Störungen linearer Systeme
- Software für lineare Gleichungssysteme

Lineare Systeme

Wir betrachten in diesem Abschnitt das Problem, ein **lineares Gleichungssystem**

$$Ax = b$$

Lineare Systeme

Wir betrachten in diesem Abschnitt das Problem, ein **lineares Gleichungssystem**

$$Ax = b$$

mit einer gegebenen **regulären** Matrix

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

Lineare Systeme

Wir betrachten in diesem Abschnitt das Problem, ein **lineares Gleichungssystem**

$$Ax = b$$

mit einer gegebenen **regulären** Matrix

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

und einem gegebenen Vektor

$$b \in \mathbb{R}^n$$

zu lösen.

Zerlegung regulärer Matrizen

Wir betrachten das lineare Gleichungssystem

$$\mathbf{Ax} = \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{n \times n}, \quad \mathbf{b} \in \mathbb{R}^n. \quad (4.1)$$

Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ regulär.

Zerlegung regulärer Matrizen

Wir betrachten das lineare Gleichungssystem

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n. \quad (4.1)$$

Es sei $A \in \mathbb{R}^{n \times n}$ regulär. Dann existiert (vgl. LA Satz 4.42)

- ▶ eine Permutationsmatrix P ,
- ▶ eine normierte untere Dreiecksmatrix L ,
($\ell_{jj} = 1 \forall j = 1, \dots, n$) und
- ▶ eine obere Dreiecksmatrix R

mit

$$PA = LR. \quad (4.2)$$

Zerlegung regulärer Matrizen

Wir betrachten das lineare Gleichungssystem

$$Ax = b, \quad A \in \mathbb{R}^{n \times n}, \quad b \in \mathbb{R}^n. \quad (4.1)$$

Es sei $A \in \mathbb{R}^{n \times n}$ regulär. Dann existiert (vgl. LA Satz 4.42)

- ▶ eine Permutationsmatrix P ,
- ▶ eine normierte untere Dreiecksmatrix L ,
($\ell_{jj} = 1 \forall j = 1, \dots, n$) und
- ▶ eine obere Dreiecksmatrix R

mit

$$PA = LR. \quad (4.2)$$

Diese Zerlegung heißt die **LR-Zerlegung** (mit partieller Pivotisierung) der Matrix A .

Zerlegung regulärer Matrizen

Es gilt $P = E$ genau dann, wenn alle Hauptuntermatrizen

$$A(1:k, 1:k) := (a_{ij})_{i,j=1}^k, \quad k = 1, \dots, n,$$

von A regulär.

Zerlegung regulärer Matrizen

Es gilt $P = E$ genau dann, wenn alle Hauptuntermatrizen

$$A(1 : k, 1 : k) := (a_{ij})_{i,j=1}^k, \quad k = 1, \dots, n,$$

von A regulär.

Ist dies der Fall, so kann man die LR-Zerlegung von A mit dem **Gaußschen Eliminationsverfahren** bestimmen.

Zerlegung regulärer Matrizen

Es gilt $P = E$ genau dann, wenn alle Hauptuntermatrizen

$$A(1:k, 1:k) := (a_{ij})_{i,j=1}^k, \quad k = 1, \dots, n,$$

von A regulär.

Ist dies der Fall, so kann man die LR-Zerlegung von A mit dem **Gaußschen Eliminationsverfahren** bestimmen.

Dabei speichern wir die von Null verschiedenen Elemente von R in dem oberen Dreieck von A und die Elemente unterhalb der Diagonale von L in dem strikten unteren Dreieck von A ab.

Zerlegung regulärer Matrizen

Algorithmus 4.1: (LR-Zerlegung ohne Pivotsuche)

```
for i = 1:n-1
    for j = i+1:n
        % Bestimme die i-te Spalte von L
        a(j,i) = a(j,i)/a(i,i);
        for k = i+1:n
            % Datiere die j-te Zeile auf
            a(j,k) = a(j,k) - a(j,i)*a(i,k);
        end
    end
end
end
```

Ist die LR-Zerlegung von A bekannt, so kann man die Lösung des Gleichungssystems (4.1) schreiben als

$$\mathbf{Ax} = \mathbf{LRx} =: \mathbf{Ly} = \mathbf{b}, \quad \mathbf{Rx} = \mathbf{y},$$

Zerlegung regulärer Matrizen

... und durch **Vorwärtseinsetzen**

Algorithmus 4.2: (Vorwärtseinsetzen)

```
for j = 1:n
    y(j) = b(j);
    for k = 1:j-1
        y(j) = y(j) - a(j,k)*y(k);
    end
end
```

end

die Lösung y von $Ly = b$ bestimmen und durch **Rückwärtseinsetzen** ...

Zerlegung regulärer Matrizen

... und durch Rückwärtseinsetzen

Algorithmus 4.3: (Rückwärtseinsetzen)

```
for j = n:-1:1
    x(j) = y(j);
    for k = j+1:n
        x(j) = x(j) - a(j,k)*x(k);
    end
    x(j) = x(j)/a(j,j);
end
```

die Lösung \mathbf{x} von $\mathbf{R}\mathbf{x} = \mathbf{y}$, d.h. von $\mathbf{A}\mathbf{x} = \mathbf{b}$.

Zerlegung regulärer Matrizen

Als Aufwand der LR-Zerlegung erhält man

$$\begin{aligned}\sum_{i=1}^{n-1} \left(\sum_{j=i+1}^n (1 + \sum_{k=i+1}^n 2) \right) &= \sum_{i=1}^{n-1} \left((n-i) + 2(n-i)^2 \right) \\ &= \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n \\ &= \frac{2}{3}n^3 + O(n^2).\end{aligned}$$

flops (floating point operations).

Zerlegung regulärer Matrizen

Als Aufwand der LR-Zerlegung erhält man

$$\begin{aligned}\sum_{i=1}^{n-1} \left(\sum_{j=i+1}^n (1 + \sum_{k=i+1}^n 2) \right) &= \sum_{i=1}^{n-1} \left((n-i) + 2(n-i)^2 \right) \\ &= \frac{2}{3}n^3 - \frac{1}{2}n^2 - \frac{1}{6}n \\ &= \frac{2}{3}n^3 + O(n^2).\end{aligned}$$

flops (floating point operations).

Vorwärts- und Rückwärtseinsetzen erfordert jeweils $O(n^2)$ flops.

Zerlegung regulärer Matrizen

Existiert eine **singuläre Hauptuntermatrix** $A(1 : i, 1 : i)$ der regulären Matrix A , so bricht Algorithmus 4.1 ab, da das Pivotelement $a(i, i)$ bei der Elimination der i -ten Variable Null wird.

Zerlegung regulärer Matrizen

Existiert eine **singuläre Hauptuntermatrix** $A(1:i, 1:i)$ der regulären Matrix A , so bricht Algorithmus 4.1 ab, da das Pivotelement $a(i, i)$ bei der Elimination der i -ten Variable Null wird.

In diesem Fall gibt es ein $a_{ij} \neq 0, j > i$ (das im Verlaufe des Algorithmus erzeugt worden ist), und man kann die aktuelle **Zeile i mit der Zeile j vertauschen** und den Algorithmus fortsetzen.

Zerlegung regulärer Matrizen

Sammelt man diese Vertauschungen in der Permutationsmatrix P , so erhält man am Ende des so modifizierten Algorithmus 4.1 eine LR-Zerlegung der Matrix PA :

$$PA = LR.$$

Zerlegung regulärer Matrizen

Sammelt man diese Vertauschungen in der Permutationsmatrix P , so erhält man am Ende des so modifizierten Algorithmus 4.1 eine LR-Zerlegung der Matrix PA :

$$PA = LR.$$

Aus Stabilitätsgründen empfiehlt es sich, auch dann eine Vertauschung vorzunehmen, wenn a_{ii} zwar **von Null verschieden** ist, in der Restmatrix $A(i:n, i:n)$ aber Elemente vorhanden sind, deren Betrag **wesentlich größer** ist als der Betrag von a_{ii} .

Zerlegung regulärer Matrizen

Beispiel 4.4: Es sei

$$\mathbf{A} = \begin{pmatrix} 10^{-4} & 1 \\ 1 & 1 \end{pmatrix}.$$

Zerlegung regulärer Matrizen

Beispiel 4.4: Es sei

$$\mathbf{A} = \begin{pmatrix} 10^{-4} & 1 \\ 1 & 1 \end{pmatrix}.$$

Wir bezeichnen mit $\text{fl}(a)$ die Gleitpunktdarstellung von a . Dann liefert Algorithmus 4.1 bei dreistelliger Rechnung

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ \text{fl}(1/10^{-4}) & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 10^4 & 1 \end{pmatrix}$$

und

$$\mathbf{R} = \begin{pmatrix} 10^{-4} & 1 \\ 0 & \text{fl}(1 - 10^4) \end{pmatrix} = \begin{pmatrix} 10^{-4} & 1 \\ 0 & -10^4 \end{pmatrix},$$

Zerlegung regulärer Matrizen

Beispiel 4.4: Es sei

$$\mathbf{A} = \begin{pmatrix} 10^{-4} & 1 \\ 1 & 1 \end{pmatrix}.$$

Wir bezeichnen mit $\text{fl}(a)$ die Gleitpunktdarstellung von a . Dann liefert Algorithmus 4.1 bei dreistelliger Rechnung

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ \text{fl}(1/10^{-4}) & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 10^4 & 1 \end{pmatrix}$$

und

$$\mathbf{R} = \begin{pmatrix} 10^{-4} & 1 \\ 0 & \text{fl}(1 - 10^4) \end{pmatrix} = \begin{pmatrix} 10^{-4} & 1 \\ 0 & -10^4 \end{pmatrix},$$

und damit

$$\mathbf{LR} = \begin{pmatrix} 10^{-4} & 1 \\ 1 & 0 \end{pmatrix} \neq \mathbf{A}.$$

□

Zerlegung regulärer Matrizen

In vielen Fällen reicht es aus, vor dem Annullieren der Elemente der i -ten Spalte $j \in \{i, \dots, n\}$ zu bestimmen mit $|a_{ji}| \geq |a_{ki}|$ für alle $k = i, \dots, n$ und dann die i -te Zeile mit der j -ten Zeile zu vertauschen.

Zerlegung regulärer Matrizen

In vielen Fällen reicht es aus, vor dem Annullieren der Elemente der i -ten Spalte $j \in \{i, \dots, n\}$ zu bestimmen mit $|a_{ji}| \geq |a_{ki}|$ für alle $k = i, \dots, n$ und dann die i -te Zeile mit der j -ten Zeile zu vertauschen.

Dieses Vorgehen heißt **Spaltenpivotsuche** oder **partielle Pivotisierung**.

Zerlegung regulärer Matrizen

Algorithmus 4.5: (LR-Zerlegung mit Spaltenpivotsuche)

```
for i = 1:n-1
```

Wähle $j \geq i$ mit $|a_{ji}| \geq |a_{ki}|$ für alle $k \geq i$ und vertausche die i -te mit der j -ten Zeile

```
    for j = i+1:n
        a(j,i) = a(j,i)/a(i,i);
        for k = i+1:n
            a(j,k) = a(j,k) - a(j,i)*a(i,k);
        end
    end
end
end
```

Zerlegung regulärer Matrizen

Beispiel 4.6: Für Beispiel 4.4 erhält man mit Algorithmus 4.5

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ 10^{-4} & 1 \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} 1 & 1 \\ 0 & \text{fl}(1 - 10^{-4}) \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

Zerlegung regulärer Matrizen

Beispiel 4.6: Für Beispiel 4.4 erhält man mit Algorithmus 4.5

$$L = \begin{pmatrix} 1 & 0 \\ 10^{-4} & 1 \end{pmatrix}, \quad R = \begin{pmatrix} 1 & 1 \\ 0 & \text{fl}(1 - 10^{-4}) \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

und

$$LR = \begin{pmatrix} 1 & 1 \\ 10^{-4} & 1 + 10^{-4} \end{pmatrix}$$

ist eine gute Approximation von

$$PA = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} A = \begin{pmatrix} 1 & 1 \\ 10^{-4} & 1 \end{pmatrix}.$$

Zerlegung regulärer Matrizen

Nicht in allen Fällen führt die Spaltenpivotsuche zum Erfolg. Wendet man Algorithmus 4.5 bei dreistelliger Rechnung auf

$$\tilde{\mathbf{A}} = \begin{pmatrix} 1 & 10^4 \\ 1 & 1 \end{pmatrix}$$

an

Zerlegung regulärer Matrizen

Nicht in allen Fällen führt die Spaltenpivotsuche zum Erfolg. Wendet man Algorithmus 4.5 bei dreistelliger Rechnung auf

$$\tilde{\mathbf{A}} = \begin{pmatrix} 1 & 10^4 \\ 1 & 1 \end{pmatrix}$$

an, so erhält man

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} 1 & 10^4 \\ 0 & \text{fl}(1 - 10^4) \end{pmatrix} = \begin{pmatrix} 1 & 10^4 \\ 0 & -10^4 \end{pmatrix}$$

Zerlegung regulärer Matrizen

Nicht in allen Fällen führt die Spaltenpivotsuche zum Erfolg. Wendet man Algorithmus 4.5 bei dreistelliger Rechnung auf

$$\tilde{\mathbf{A}} = \begin{pmatrix} 1 & 10^4 \\ 1 & 1 \end{pmatrix}$$

an, so erhält man

$$\mathbf{L} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} 1 & 10^4 \\ 0 & \text{fl}(1 - 10^4) \end{pmatrix} = \begin{pmatrix} 1 & 10^4 \\ 0 & -10^4 \end{pmatrix}$$

und hiermit

$$\mathbf{LR} = \begin{pmatrix} 1 & 10^4 \\ 1 & 0 \end{pmatrix}.$$

□

Zerlegung regulärer Matrizen

Treten, wie in unserem Beispiel, in der Matrix Elemente sehr unterschiedlicher Größenordnung auf, so empfiehlt es sich, eine **vollständige Pivotsuche** auszuführen.

Zerlegung regulärer Matrizen

Treten, wie in unserem Beispiel, in der Matrix Elemente sehr unterschiedlicher Größenordnung auf, so empfiehlt es sich, eine **vollständige Pivotsuche** auszuführen.

In diesem Fall werden im i -ten Eliminationsschritt **ein Zeilenindex** $j \geq i$ und **ein Spaltenindex** $k \geq i$ bestimmt mit $|a_{jk}| \geq |a_{\ell m}|$ für alle $\ell, m \geq i$, und es wird vor der Elimination die i -te mit der j -ten Zeile und die i -te Spalte mit der k -ten Spalte getauscht.

Zerlegung regulärer Matrizen

Treten, wie in unserem Beispiel, in der Matrix Elemente sehr unterschiedlicher Größenordnung auf, so empfiehlt es sich, eine **vollständige Pivotsuche** auszuführen.

In diesem Fall werden im i -ten Eliminationsschritt **ein Zeilenindex** $j \geq i$ und **ein Spaltenindex** $k \geq i$ bestimmt mit $|a_{jk}| \geq |a_{\ell m}|$ für alle $\ell, m \geq i$, und es wird vor der Elimination die i -te mit der j -ten Zeile und die i -te Spalte mit der k -ten Spalte getauscht.

Man erhält dann eine Zerlegung

$$PAQ = LR,$$

wobei die Permutationsmatrix P die Zeilenvertauschungen in A vornimmt und die Permutationsmatrix Q die Spaltenvertauschungen.

Zerlegung regulärer Matrizen

Ist die reguläre Matrix A **symmetrisch** und existiert eine LR-Zerlegung, so kann man R als Produkt einer Diagonalmatrix D und einer normierten oberen Dreiecksmatrix \tilde{R} schreiben.

Zerlegung regulärer Matrizen

Ist die reguläre Matrix A **symmetrisch** und existiert eine LR-Zerlegung, so kann man R als Produkt einer Diagonalmatrix D und einer normierten oberen Dreiecksmatrix \tilde{R} schreiben. Hiermit gilt dann

$$A^T = \tilde{R}^T D L^T$$

mit einer normierten unteren Dreiecksmatrix \tilde{R}^T und einer oberen Dreiecksmatrix $D L^T$.

Zerlegung regulärer Matrizen

Ist die reguläre Matrix A **symmetrisch** und existiert eine LR-Zerlegung, so kann man R als Produkt einer Diagonalmatrix D und einer normierten oberen Dreiecksmatrix \tilde{R} schreiben. Hiermit gilt dann

$$A^T = \tilde{R}^T D L^T$$

mit einer normierten unteren Dreiecksmatrix \tilde{R}^T und einer oberen Dreiecksmatrix $D L^T$.

Da die LR-Zerlegung einer regulären Matrix eindeutig bestimmt ist, folgt $\tilde{R}^T = L$. Damit kann man A auch schreiben als

$$A = L D L^T$$

mit einer Diagonalmatrix D und L wie oben. Diese Zerlegung heißt die **LDLT-Zerlegung** von A .

Zerlegung regulärer Matrizen

Ist die symmetrische Matrix A zusätzlich **positiv definit**, so sind alle Hauptuntermatrizen von A ebenfalls positiv definit, also regulär und daher existiert in diesem Fall die LDLT-Zerlegung.

Zerlegung regulärer Matrizen

Ist die symmetrische Matrix A zusätzlich **positiv definit**, so sind alle Hauptuntermatrizen von A ebenfalls positiv definit, also regulär und daher existiert in diesem Fall die LDLT-Zerlegung.

Ferner sind die Diagonalelemente von $D = \text{diag}\{d_1, \dots, d_n\}$ positiv, und man kann mit

$$C = L \text{diag}\{\sqrt{d_1}, \dots, \sqrt{d_n}\}$$

die LDLT-Zerlegung von A schreiben als

$$A = CC^T. \quad (4.3)$$

Zerlegung regulärer Matrizen

Ist die symmetrische Matrix A zusätzlich **positiv definit**, so sind alle Hauptuntermatrizen von A ebenfalls positiv definit, also regulär und daher existiert in diesem Fall die LDLT-Zerlegung.

Ferner sind die Diagonalelemente von $D = \text{diag}\{d_1, \dots, d_n\}$ positiv, und man kann mit

$$C = L \text{diag}\{\sqrt{d_1}, \dots, \sqrt{d_n}\}$$

die LDLT-Zerlegung von A schreiben als

$$A = CC^T. \tag{4.3}$$

Diese Zerlegung heißt die **Cholesky-Zerlegung** von A .

Zerlegung regulärer Matrizen

Prinzipiell kann man die Cholesky-Zerlegung mit Hilfe des Gaußschen Eliminationsverfahrens bestimmen.

Zerlegung regulärer Matrizen

Prinzipiell kann man die Cholesky-Zerlegung mit Hilfe des Gaußschen Eliminationsverfahrens bestimmen.

Man benötigt dann $\frac{2}{3}n^3 + O(n^2)$ Operationen und n^2 Speicherplätze.

Zerlegung regulärer Matrizen

Prinzipiell kann man die Cholesky-Zerlegung mit Hilfe des Gaußschen Eliminationsverfahrens bestimmen.

Man benötigt dann $\frac{2}{3}n^3 + O(n^2)$ Operationen und n^2 Speicherplätze.

Durch direkten Vergleich der Elemente in (4.3) erhält man einen Algorithmus, der mit der **Hälfte** des Aufwandes auskommt.

Zerlegung regulärer Matrizen

Da C eine untere Dreiecksmatrix ist, gilt $c_{ij} = 0$ für $1 \leq i < j \leq n$, und daher ist

$$a_{ij} = \sum_{k=1}^n c_{ik}c_{jk} = \sum_{k=1}^i c_{ik}c_{jk}.$$

Zerlegung regulärer Matrizen

Da C eine untere Dreiecksmatrix ist, gilt $c_{ij} = 0$ für $1 \leq i < j \leq n$, und daher ist

$$a_{ij} = \sum_{k=1}^n c_{ik}c_{jk} = \sum_{k=1}^i c_{ik}c_{jk}.$$

Speziell für $i = j = 1$ bedeutet dies

$$a_{11} = c_{11}^2, \quad \text{d.h.} \quad c_{11} = \sqrt{a_{11}},$$

Zerlegung regulärer Matrizen

Da C eine untere Dreiecksmatrix ist, gilt $c_{ij} = 0$ für $1 \leq i < j \leq n$, und daher ist

$$a_{ij} = \sum_{k=1}^n c_{ik}c_{jk} = \sum_{k=1}^i c_{ik}c_{jk}.$$

Speziell für $i = j = 1$ bedeutet dies

$$a_{11} = c_{11}^2, \quad \text{d.h.} \quad c_{11} = \sqrt{a_{11}},$$

und für $i = 1$ und $j = 2, \dots, n$

$$a_{1j} = c_{11}c_{j1}, \quad \text{d.h.} \quad c_{j1} = a_{1j}/c_{11}.$$

Zerlegung regulärer Matrizen

Da C eine untere Dreiecksmatrix ist, gilt $c_{ij} = 0$ für $1 \leq i < j \leq n$, und daher ist

$$a_{ij} = \sum_{k=1}^n c_{ik}c_{jk} = \sum_{k=1}^i c_{ik}c_{jk}.$$

Speziell für $i = j = 1$ bedeutet dies

$$a_{11} = c_{11}^2, \quad \text{d.h.} \quad c_{11} = \sqrt{a_{11}},$$

und für $i = 1$ und $j = 2, \dots, n$

$$a_{1j} = c_{11}c_{j1}, \quad \text{d.h.} \quad c_{j1} = a_{1j}/c_{11}.$$

Damit ist die **erste Spalte** von C **bestimmt**.

Zerlegung regulärer Matrizen

Sind schon $c_{\mu\nu}$ für $\nu = 1, \dots, i-1$ und $\mu = \nu, \nu+1, \dots, n$ bestimmt, so erhält man aus

$$a_{ii} = c_{ii}^2 + \sum_{k=1}^{i-1} c_{ik}^2$$

Zerlegung regulärer Matrizen

Sind schon $c_{\mu\nu}$ für $\nu = 1, \dots, i-1$ und $\mu = \nu, \nu+1, \dots, n$ bestimmt, so erhält man aus

$$a_{ii} = c_{ii}^2 + \sum_{k=1}^{i-1} c_{ik}^2$$

das i -te **Diagonalelement**

$$c_{ii} = \sqrt{a_{ii} - \sum_{k=1}^{i-1} c_{ik}^2}$$

von C ,

Zerlegung regulärer Matrizen

und

$$a_{ij} = c_{ii}c_{ji} + \sum_{k=1}^{i-1} c_{ik}c_{jk}$$

Zerlegung regulärer Matrizen

und

$$a_{ij} = c_{ii}c_{ji} + \sum_{k=1}^{i-1} c_{ik}c_{jk}$$

liefert die i -te Spalte

$$c_{ji} = \frac{1}{c_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} c_{ik}c_{jk} \right), \quad j = i + 1, \dots, n.$$

Zerlegung regulärer Matrizen

und

$$a_{ij} = c_{ii}c_{ji} + \sum_{k=1}^{i-1} c_{ik}c_{jk}$$

liefert die i -te **Spalte**

$$c_{ji} = \frac{1}{c_{ii}} \left(a_{ij} - \sum_{k=1}^{i-1} c_{ik}c_{jk} \right), \quad j = i + 1, \dots, n.$$

Damit erhält man das folgende **Verfahren zur Bestimmung der Cholesky-Zerlegung**. Dabei überschreiben wir das untere Dreieck von A durch die wesentlichen Elemente von C .

Zerlegung regulärer Matrizen

Algorithmus 4.7: (Cholesky-Zerlegung)

```
for i = 1:n
    for k = 1:i-1
        a(i,i) = a(i,i) - a(i,k)*a(i,k);
    end
    a(i,i) = sqrt(a(i,i));
    for j = i+1:n
        for k = 1:i-1
            a(j,i) = a(j,i) - a(i,k)*a(j,k);
        end
        a(j,i) = a(j,i)/a(i,i);
    end
end
end
```

Zerlegung regulärer Matrizen

Der Aufwand des Cholesky-Verfahrens ist

$$\begin{aligned}\sum_{i=1}^n \left(\sum_{k=1}^{i-1} 2 + \sum_{j=i+1}^n (1 + \sum_{k=1}^{i-1} 2) \right) &= \sum_{i=1}^n \left(2(i-1) + (n-i)(2i-1) \right) \\ &= \frac{1}{3}n^3 + \frac{1}{2}n^2 - \frac{5}{6}n \\ &= \frac{1}{3}n^3 + O(n^2)\end{aligned}$$

flops und n Quadratwurzeln.

Zerlegung regulärer Matrizen

Der Aufwand des Cholesky-Verfahrens ist

$$\begin{aligned}\sum_{i=1}^n \left(\sum_{k=1}^{i-1} 2 + \sum_{j=i+1}^n (1 + \sum_{k=1}^{i-1} 2) \right) &= \sum_{i=1}^n \left(2(i-1) + (n-i)(2i-1) \right) \\ &= \frac{1}{3}n^3 + \frac{1}{2}n^2 - \frac{5}{6}n \\ &= \frac{1}{3}n^3 + O(n^2)\end{aligned}$$

flops und n Quadratwurzeln.

Da eine Quadratwurzel auf den meisten Rechnern mit IEEE 754-Arithmetik in etwa so teuer ist wie ein flop, ist der Aufwand also vergleichbar mit

$$\frac{1}{3}n^3 + \frac{1}{2}n^2 + \frac{1}{6}n.$$

Zerlegung regulärer Matrizen

Besitzt die symmetrische Matrix A eine LDLT-Zerlegung, so kann man diese mit einem ähnlichen Verfahren wie Algorithmus 4.7 bestimmen. Man beachte aber, dass eine Pivotsuche wie beim Gaußschen Eliminationsverfahren die Symmetrie der Matrix zerstört.

Zerlegung regulärer Matrizen

Besitzt die symmetrische Matrix A eine LDLT-Zerlegung, so kann man diese mit einem ähnlichen Verfahren wie Algorithmus 4.7 bestimmen. Man beachte aber, dass eine Pivotsuche wie beim Gaußschen Eliminationsverfahren die Symmetrie der Matrix zerstört.

Man muss also **gleichlautende Zeilen- und Spaltenvertauschungen** vornehmen. Auf diese Weise kann man nicht immer für eine reguläre Matrix ein von Null verschiedenes Pivotelement erzeugen, wie das Beispiel

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

zeigt.

Zerlegung regulärer Matrizen

Besitzt die symmetrische Matrix A eine LDLT-Zerlegung, so kann man diese mit einem ähnlichen Verfahren wie Algorithmus 4.7 bestimmen. Man beachte aber, dass eine Pivotsuche wie beim Gaußschen Eliminationsverfahren die Symmetrie der Matrix zerstört.

Man muss also **gleichlautende Zeilen- und Spaltenvertauschungen** vornehmen. Auf diese Weise kann man nicht immer für eine reguläre Matrix ein von Null verschiedenes Pivotelement erzeugen, wie das Beispiel

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

zeigt.

Auch wenn die LDLT-Zerlegung existiert, kann die Übertragung von Algorithmus 4.7 **instabil** sein.

Modifikationen des Gaußschen Verfahrens

In Algorithmus 4.1 haben wir die Elimination durchgeführt, indem wir im i -ten Schritt ein geeignetes Vielfaches der i -ten **Zeile** von den nachfolgenden Zeilen abgezogen haben. In Vektorschreibweise haben wir also

Algorithmus 4.8: (LR-Zerlegung; zeilenorientiert)

```
for i = 1:n-1
    for j = i+1:n
        a(j,i) = a(j,i)/a(i,i);
        a(j,i+1:n) = a(j,i+1:n) - a(j,i)*a(i,i+1:n);
    end
end
```

Modifikationen des Gaußschen Verfahrens

Vertauscht man die beiden inneren Schleifen, so erhält man eine **spaltenorientierte Version** des Gaußschen Eliminationsverfahrens.

Algorithmus 4.9: (LR-Zerlegung; spaltenorientiert)

```
for i = 1:n-1
    a(i+1:n,i) = a(i+1:n,i)/a(i,i);
    for k = i+1:n
        a(i+1:n,k) = a(i+1:n,k) - a(i,k)*a(i+1:n,i);
    end
end
```

Modifikationen des Gaußschen Verfahrens

Auch die i -Schleife kann man mit der j -Schleife und/oder der k -Schleife vertauschen. Man erhält insgesamt **6 Varianten** der Gauß-Elimination, die alle von der Zahl der Rechenoperationen her denselben Aufwand besitzen.

Modifikationen des Gaußschen Verfahrens

Auch die i -Schleife kann man mit der j -Schleife und/oder der k -Schleife vertauschen. Man erhält insgesamt **6 Varianten** der Gauß-Elimination, die alle von der Zahl der Rechenoperationen her denselben Aufwand besitzen.

Sie können sich aber auf **verschiedenen Plattformen** unter **verschiedenen Programmiersprachen** sehr unterschiedlich verhalten.

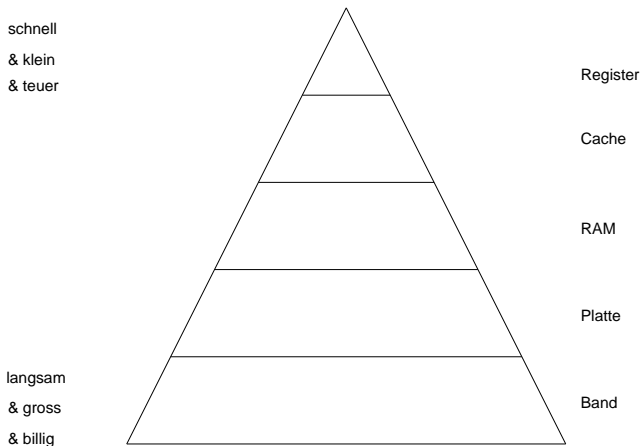
Modifikationen des Gaußschen Verfahrens

Auch die i -Schleife kann man mit der j -Schleife und/oder der k -Schleife vertauschen. Man erhält insgesamt **6 Varianten** der Gauß-Elimination, die alle von der Zahl der Rechenoperationen her denselben Aufwand besitzen.

Sie können sich aber auf **verschiedenen Plattformen** unter **verschiedenen Programmiersprachen** sehr unterschiedlich verhalten.

Der Grund hierfür ist, dass Speicher von Rechnern hierarchisch aufgebaut sind, beginnend mit sehr langsamen, sehr großen, sehr billigen Magnetband-Speichern, über schnellere, kleinere, teurere Plattenspeicher, den noch schnelleren, noch kleineren, noch teureren Hauptspeicher, den schnellen, kleinen und teuren Cache und die sehr schnellen, sehr kleinen und sehr teuren Register der CPU. Arithmetische und logische Operationen können nur in den Registern ausgeführt werden.

Hierarchischer Speicher



Modifikationen des Gaußschen Verfahrens

Daten können nur in den benachbarten Speicher bewegt werden, wobei der Datentransport zwischen den billigen Speicherarten sehr **langsam** geschieht, während der Transport zwischen den teuren Speichern an der Spitze der Hierarchie sehr **schnell** geschieht.

Modifikationen des Gaußschen Verfahrens

Daten können nur in den benachbarten Speicher bewegt werden, wobei der Datentransport zwischen den billigen Speicherarten sehr **langsam** geschieht, während der Transport zwischen den teuren Speichern an der Spitze der Hierarchie sehr **schnell** geschieht.

Insbesondere ist die Geschwindigkeit, mit der **arithmetische Operationen** ausgeführt werden, sehr viel höher als die Geschwindigkeit, mit der Daten transportiert werden.

Modifikationen des Gaußschen Verfahrens

Daten können nur in den benachbarten Speicher bewegt werden, wobei der Datentransport zwischen den billigen Speicherarten sehr **langsam** geschieht, während der Transport zwischen den teuren Speichern an der Spitze der Hierarchie sehr **schnell** geschieht.

Insbesondere ist die Geschwindigkeit, mit der **arithmetische Operationen** ausgeführt werden, sehr viel höher als die Geschwindigkeit, mit der Daten transportiert werden.

Faktoren zwischen 10 und 10000 (je nach Speicherebene) sind nicht ungewöhnlich. Algorithmen müssen also so gestaltet werden, dass der **Datentransport** zwischen verschiedenen Speicherebenen **möglichst klein** ist.

Modifikationen des Gaußschen Verfahrens

Die **Reihenfolge** der Schleifen im Gaußschen Eliminationsverfahren hat Einfluss auf die **Laufzeit eines Programms**.

Modifikationen des Gaußschen Verfahrens

Die **Reihenfolge** der Schleifen im Gaußschen Eliminationsverfahren hat Einfluss auf die **Laufzeit eines Programms**.

Eine Matrix wird in einem langen Array gespeichert. In der Sprache C geschieht das zeilenweise:

$$\begin{array}{ccccccccc} & a(1, 1) & \rightarrow & a(1, 2) & \rightarrow & a(1, 3) & \rightarrow & \dots & \rightarrow & a(1, n) \\ \rightarrow & a(2, 1) & \rightarrow & a(2, 2) & \rightarrow & a(2, 3) & \rightarrow & \dots & \rightarrow & a(2, n) \\ \rightarrow & a(3, 1) & \rightarrow & a(3, 2) & \rightarrow & a(3, 3) & \rightarrow & \dots & \rightarrow & a(3, n) \\ & \vdots & & \vdots & & \vdots & & & & \vdots \\ \rightarrow & a(n, 1) & \rightarrow & a(n, 2) & \rightarrow & a(n, 3) & \rightarrow & \dots & \rightarrow & a(n, n) \end{array}$$

Modifikationen des Gaußschen Verfahrens

Die **Reihenfolge** der Schleifen im Gaußschen Eliminationsverfahren hat Einfluss auf die **Laufzeit eines Programms**.

Eine Matrix wird in einem langen Array gespeichert. In der Sprache C geschieht das zeilenweise:

$$\begin{array}{ccccccc}
 & a(1, 1) & \rightarrow & a(1, 2) & \rightarrow & a(1, 3) & \rightarrow \dots \rightarrow a(1, n) \\
 \rightarrow & a(2, 1) & \rightarrow & a(2, 2) & \rightarrow & a(2, 3) & \rightarrow \dots \rightarrow a(2, n) \\
 \rightarrow & a(3, 1) & \rightarrow & a(3, 2) & \rightarrow & a(3, 3) & \rightarrow \dots \rightarrow a(3, n) \\
 & \vdots & & \vdots & & \vdots & & \vdots \\
 \rightarrow & a(n, 1) & \rightarrow & a(n, 2) & \rightarrow & a(n, 3) & \rightarrow \dots \rightarrow a(n, n)
 \end{array}$$

In der zeilenorientierten Version in Algorithmus 4.8 werden in der innersten Schleife Daten verwendet, die im Speicher **nahe beieinander** liegen.

Modifikationen des Gaußschen Verfahrens

Die **Reihenfolge** der Schleifen im Gaußschen Eliminationsverfahren hat Einfluss auf die **Laufzeit eines Programms**.

Eine Matrix wird in einem langen Array gespeichert. In der Sprache C geschieht das zeilenweise:

$$\begin{array}{ccccccc}
 & a(1, 1) & \rightarrow & a(1, 2) & \rightarrow & a(1, 3) & \rightarrow \dots \rightarrow a(1, n) \\
 \rightarrow & a(2, 1) & \rightarrow & a(2, 2) & \rightarrow & a(2, 3) & \rightarrow \dots \rightarrow a(2, n) \\
 \rightarrow & a(3, 1) & \rightarrow & a(3, 2) & \rightarrow & a(3, 3) & \rightarrow \dots \rightarrow a(3, n) \\
 & \vdots & & \vdots & & \vdots & & \vdots \\
 \rightarrow & a(n, 1) & \rightarrow & a(n, 2) & \rightarrow & a(n, 3) & \rightarrow \dots \rightarrow a(n, n)
 \end{array}$$

In der zeilenorientierten Version in Algorithmus 4.8 werden in der innersten Schleife Daten verwendet, die im Speicher **nahe beieinander** liegen.

Es sind daher nur wenige Datentransporte zwischen den verschiedenen Speicherebenen erforderlich.

Modifikationen des Gaußschen Verfahrens

Für die **spaltenorientierte Version** liegen die nacheinander benutzten Daten (für große n) sehr weit auseinander, und daher ist ein „**cache miss**“ (das benötigte Wort liegt nicht im Cache und es müssen zwei Blöcke zwischen dem Cache und dem Hauptspeicher ausgetauscht werden) oder sogar ein „**page fault**“ (das benötigte Wort liegt nicht einmal im Hauptspeicher und es müssen zwei Seiten zwischen dem Hauptspeicher und der Platte ausgetauscht werden) wahrscheinlicher.

Modifikationen des Gaußschen Verfahrens

Für die **spaltenorientierte Version** liegen die nacheinander benutzten Daten (für große n) sehr weit auseinander, und daher ist ein „**cache miss**“ (das benötigte Wort liegt nicht im Cache und es müssen zwei Blöcke zwischen dem Cache und dem Hauptspeicher ausgetauscht werden) oder sogar ein „**page fault**“ (das benötigte Wort liegt nicht einmal im Hauptspeicher und es müssen zwei Seiten zwischen dem Hauptspeicher und der Platte ausgetauscht werden) wahrscheinlicher.

In FORTRAN ist die Situation umgekehrt. Matrizen werden **spaltenweise** gespeichert, und für Algorithmus 4.9 ist die Datenlokalität hoch, während sie in Algorithmus 4.8 gering ist.

Modifikationen des Gaußschen Verfahrens

Um nicht für jeden neuen Rechner neue Software schreiben zu müssen, um die Modularität und Effizienz von Programmen zu erhöhen und um ihre Pflege zu erleichtern, wurden Standardoperationen der (numerischen) linearen Algebra, die „basic linear algebra subprograms“ (**BLAS**), definiert und es wurden Standardschnittstellen für ihren Aufruf festgelegt.

Modifikationen des Gaußschen Verfahrens

Um nicht für jeden neuen Rechner neue Software schreiben zu müssen, um die Modularität und Effizienz von Programmen zu erhöhen und um ihre Pflege zu erleichtern, wurden Standardoperationen der (numerischen) linearen Algebra, die „basic linear algebra subprograms“ (**BLAS**), definiert und es wurden Standardschnittstellen für ihren Aufruf festgelegt.

Diese werden (jedenfalls für Hochleistungsrechner) von den Herstellern auf **Hardwareebene realisiert** (nicht zuletzt, da die Hersteller wissen, dass die üblichen Benchmarktests Programme enthalten, die aus den BLAS-Routinen aufgebaut sind).

Modifikationen des Gaußschen Verfahrens

Die Benutzung der BLAS in eigenen Programmen bietet eine Reihe von Vorteilen:

- Die **Robustheit** von Berechnungen der linearen Algebra wird durch die BLAS erhöht, denn in Ihnen werden Details der Algorithmen und der Implementierung berücksichtigt, die bei der Programmierung eines Anwendungsproblems leicht übersehen werden, wie z.B. die Berücksichtigung von Overflow-Problemen.

Modifikationen des Gaußschen Verfahrens

Die Benutzung der BLAS in eigenen Programmen bietet eine Reihe von Vorteilen:

- Die **Robustheit** von Berechnungen der linearen Algebra wird durch die BLAS erhöht, denn in Ihnen werden Details der Algorithmen und der Implementierung berücksichtigt, die bei der Programmierung eines Anwendungsproblems leicht übersehen werden, wie z.B. die Berücksichtigung von Overflow-Problemen.
- Die **Portabilität** von Programmen wird erhöht, ohne auf Effizienz zu verzichten, denn es werden optimierte Versionen der BLAS auf Rechnern verwendet, für die diese existieren. Für alle anderen Plattformen existieren kompatible Standard-FORTRAN- oder C-Implementierungen.

Modifikationen des Gaußschen Verfahrens

Diese können als Public Domain Software bezogen werden von
<http://www.netlib.org/blas/>

bzw.

<http://www.netlib.org/clapack/cblas/>

Im **ATLAS-Projekt** („automatically tuned linear algebra software“) wurden und werden empirische Methoden entwickelt, um Bibliotheken hoher Performance zu erzeugen und zu pflegen, und so in der durch die Software diktierten Geschwindigkeit Schritt mit der Hardware-Entwicklung zu halten.

Modifikationen des Gaußschen Verfahrens

Es werden z.B. für den benutzten Rechner die Größen des Registers und Caches ermittelt und für die Level 2 und Level 3 BLAS die Blockgrößen angepasst. Die im ATLAS-Projekt entwickelten BLAS, für die es FORTRAN- und C-Interfaces gibt, können herunter geladen werden von

<http://www.netlib.org/atlas/>

Modifikationen des Gaußschen Verfahrens

Es werden z.B. für den benutzten Rechner die Größen des Registers und Caches ermittelt und für die Level 2 und Level 3 BLAS die Blockgrößen angepasst. Die im ATLAS-Projekt entwickelten BLAS, für die es FORTRAN- und C-Interfaces gibt, können herunter geladen werden von

<http://www.netlib.org/atlas/>

- Die **Lesbarkeit von Programmen** wird dadurch erhöht, dass mnemonische Namen für Standardoperationen verwendet werden und der Programmablauf nicht durch Codierungsdetails unterbrochen wird. Dies erleichtert auch die Dokumentation von Programmen.

Modifikationen des Gaußschen Verfahrens

Die erste Stufe der BLAS-Definitionen (**Level 1 BLAS** oder **BLAS1**) wurde 1979 durchgeführt und enthielt Vektor-Vektor-Operationen wie das Skalarprodukt oder die Summe eines Vektors und des Vielfachen eines weiteren Vektors.

Modifikationen des Gaußschen Verfahrens

Die erste Stufe der BLAS-Definitionen (**Level 1 BLAS** oder **BLAS1**) wurde 1979 durchgeführt und enthielt Vektor-Vektor-Operationen wie das Skalarprodukt oder die Summe eines Vektors und des Vielfachen eines weiteren Vektors.

Es wurde eine **Namenskonvention** eingeführt, die einen drei- bis fünfbuchstabigen Namen verwendet, dem ein Buchstabe zur Kennzeichnung des Datentyps vorangestellt wird (S, D, C oder Z).

Modifikationen des Gaußschen Verfahrens

Die erste Stufe der BLAS-Definitionen (**Level 1 BLAS** oder **BLAS1**) wurde 1979 durchgeführt und enthielt Vektor-Vektor-Operationen wie das Skalarprodukt oder die Summe eines Vektors und des Vielfachen eines weiteren Vektors.

Es wurde eine **Namenskonvention** eingeführt, die einen drei- bis fünfbuchstabigen Namen verwendet, dem ein Buchstabe zur Kennzeichnung des Datentyps vorangestellt wird (S, D, C oder Z).

Als Beispiele nennen wir nur die Function `_DOT`, für die durch

$$\text{ALPHA}=\text{DDOT}(\text{N}, \text{X}, 1, \text{Y}, 1)$$

das innere Produkt der doppelgenauen Vektoren x und y der Dimension n berechnet werden,

Modifikationen des Gaußschen Verfahrens

oder die Subroutine `_AXPY` („ $\alpha \times \text{plus } y$ “) mit dem Aufruf

```
CAXPY (N, ALPHA, X, 1, Y, 1)
```

durch die für die komplexen Vektoren x und y der Dimension n der komplexe Vektor $\alpha x + y$ berechnet wird und im Speicher von y abgelegt wird.

Modifikationen des Gaußschen Verfahrens

oder die Subroutine `_AXPY` („ $a \times$ plus y “) mit dem Aufruf

```
CAXPY (N, ALPHA, X, 1, Y, 1)
```

durch die für die komplexen Vektoren x und y der Dimension n der komplexe Vektor $\alpha x + y$ berechnet wird und im Speicher von y abgelegt wird.

Statt der Einsen in den Aufrufen kann man **Inkremente** angeben, so dass auch innere Produkte der Art

$$\sum_{j=0}^{n-1} a_{1+mj} a_{2+mj}$$

berechnet werden können, also bei spaltenweiser Speicherung einer (m, n) -Matrix A das innere Produkt der ersten und zweiten Zeile.

Modifikationen des Gaußschen Verfahrens

Beispiel 4.10: Als Beispiel betrachten wir die Bestimmung des inneren Produktes zweier Vektoren der Dimension $n = 10^7$.

Modifikationen des Gaußschen Verfahrens

Beispiel 4.10: Als Beispiel betrachten wir die Bestimmung des inneren Produktes zweier Vektoren der Dimension $n = 10^7$.

Wir verwenden (wie auch bei den folgenden Beispielen zur Performance der BLAS Routinen) eine HP 9000/785/C3000 Workstation mit einer CPU 2.0.PA8500 mit der Taktfrequenz 400 MHz und den FORTRAN 90 Compiler f90-HP.

Modifikationen des Gaußschen Verfahrens

Beispiel 4.10: Als Beispiel betrachten wir die Bestimmung des inneren Produktes zweier Vektoren der Dimension $n = 10^7$.

Wir verwenden (wie auch bei den folgenden Beispielen zur Performance der BLAS Routinen) eine HP 9000/785/C3000 Workstation mit einer CPU 2.0.PA8500 mit der Taktfrequenz 400 MHz und den FORTRAN 90 Compiler f90-HP.

Die folgende Tabelle enthält die Laufzeiten eines naiven Codes mit einer Laufanweisung, einer BLAS1 Implementierung in FORTRAN 77, die aus der netlib heruntergeladen werden kann, einer BLAS-Implementierung, die mit dem FORTRAN 90 Compiler von HP geliefert wird, der BLAS1-Routine aus dem ATLAS-Projekt und einer optimierten BLAS-Routine der veclib von HP.

Modifikationen des Gaußschen Verfahrens

Implementierung	CPU Zeit	Relation
naiv	0.92	7.9
BLAS (netlib)	0.52	4.5
BLAS (f90-HP)	0.29	2.4
BLAS (ATLAS)	0.23	2.0
veclib	0.12	1.0



Modifikationen des Gaußschen Verfahrens

Die BLAS1 haben zu **effizienten Implementierungen** von Algorithmen auf skalaren Maschinen geführt, auf Vektorrechnern oder Parallelrechnern werden weitere Standardoperationen benötigt. Man kann z.B. das Produkt einer Matrix A mit einem Vektor x codieren als

Modifikationen des Gaußschen Verfahrens

Die BLAS1 haben zu **effizienten Implementierungen** von Algorithmen auf skalaren Maschinen geführt, auf Vektorrechnern oder Parallelrechnern werden weitere Standardoperationen benötigt. Man kann z.B. das Produkt einer Matrix A mit einem Vektor x codieren als

Algorithmus 4.11: (Matrix-Vektor Produkt mit DAXPY)

```
Y = zeros(n, 1);  
for j = 1:n  
    DAXPY(n, X(j), A(:, j), 1, Y, 1)  
end
```

Modifikationen des Gaußschen Verfahrens

Die BLAS1 haben zu **effizienten Implementierungen** von Algorithmen auf skalaren Maschinen geführt, auf Vektorrechnern oder Parallelrechnern werden weitere Standardoperationen benötigt. Man kann z.B. das Produkt einer Matrix A mit einem Vektor x codieren als

Algorithmus 4.11: (Matrix-Vektor Produkt mit DAXPY)

```
Y = zeros(n, 1);  
for j = 1:n  
    DAXPY(n, X(j), A(:, j), 1, Y, 1)  
end
```

Dabei wird aber **nicht berücksichtigt**, dass der Ergebnisvektor y im Register gehalten werden könnte. BLAS1 hat also den Nachteil, dass zu wenige (nützliche) flops ausgeführt werden im Verhältnis zu (nutzlosen) Datenbewegungen.

Modifikationen des Gaußschen Verfahrens

Für die Ausführung eines `_AXPYs` sind z. B. $3n + 1$ Speicherzugriffe erforderlich (die Vektoren x und y und der Skalar α müssen gelesen werden, und der Ergebnisvektor y muss geschrieben werden) und es werden $2n$ flops ausgeführt. Das Verhältnis ist also $2/3$.

Modifikationen des Gaußschen Verfahrens

Für die Ausführung eines `_AXPYs` sind z. B. $3n + 1$ Speicherzugriffe erforderlich (die Vektoren x und y und der Skalar α müssen gelesen werden, und der Ergebnisvektor y muss geschrieben werden) und es werden $2n$ flops ausgeführt. Das Verhältnis ist also $2/3$.

Dies wurde 1988 verbessert mit der Definition den **Level 2 BLAS** oder **BLAS2**, die **Matrix-Vektor-Operationen** enthalten, wie z.B. die Subroutine `_GEMV`, durch die $y \leftarrow \alpha Ax + \beta y$ berechnet wird, das Produkt einer Dreiecksmatrix mit einem Vektor, Rang-1- oder Rang-2-Aufdatierungen von Matrizen wie $A \leftarrow \alpha xy^T + A$, oder Lösungen von Gleichungssystemen mit Dreiecksmatrizen.

Modifikationen des Gaußschen Verfahrens

Für die Ausführung eines `_AXPYs` sind z. B. $3n + 1$ Speicherzugriffe erforderlich (die Vektoren x und y und der Skalar α müssen gelesen werden, und der Ergebnisvektor y muss geschrieben werden) und es werden $2n$ flops ausgeführt. Das Verhältnis ist also $2/3$.

Dies wurde 1988 verbessert mit der Definition den **Level 2 BLAS** oder **BLAS2**, die **Matrix-Vektor-Operationen** enthalten, wie z.B. die Subroutine `_GEMV`, durch die $y \leftarrow \alpha Ax + \beta y$ berechnet wird, das Produkt einer Dreiecksmatrix mit einem Vektor, Rang-1- oder Rang-2-Aufdatierungen von Matrizen wie $A \leftarrow \alpha xy^T + A$, oder Lösungen von Gleichungssystemen mit Dreiecksmatrizen.

Für die Subroutine `_GEMV` sind im n -dimensionalen Fall $n^2 + 3n + 2$ Speicherzugriffe erforderlich, und es werden $2n^2$ flops ausgeführt. Das Verhältnis ist also 2.

Modifikationen des Gaußschen Verfahrens

Algorithmus 4.9 ist schon fast eine BLAS2 Implementierung der LR-Zerlegung. Man muss nur noch die Modifikationen der Spalten zu einer Rang-1-Modifikation des unteren $(n - i, n - i)$ -Blocks umschreiben.

Algorithmus 4.12: (LR-Zerlegung; BLAS2)

```
for i=1: n-1
    a(i+1:n, i)=a(i+1:n, i)/a(i, i);
    a(i+1:n, i+1:n)=a(i+1:n, i+1:n)-a(i+1:n, i)*a(i, i+1:n);
end
```


Modifikationen des Gaußschen Verfahrens

Beispiel 4.13: Als Beispiel zur Performance der BLAS2 betrachten wir die Bestimmung des Produktes einer $(10^4, 10^3)$ -Matrix mit einem Vektor. Hierfür erhält man die Laufzeiten

Implementierung	CPU Zeit	Relation
naiv	4.32	61.7
BLAS2 (netlib)	1.39	19.9
BLAS2 (f90-HP)	0.62	8.9
BLAS2 (ATLAS)	0.17	2.4
veclib	0.07	1.0



Modifikationen des Gaußschen Verfahrens

In die **Level 3 BLAS** oder **BLAS3** wurden schließlich **Matrix-Matrix-Operationen** wie

$$C \leftarrow \alpha AB + \beta C$$

aufgenommen.

Modifikationen des Gaußschen Verfahrens

In die **Level 3 BLAS** oder **BLAS3** wurden schließlich **Matrix-Matrix-Operationen** wie

$$C \leftarrow \alpha AB + \beta C$$

aufgenommen.

Um die **Wiederverwendung von Daten** in den Registern oder im Cache in möglichst hohem Maße zu erreichen, werden die beteiligten Matrizen in Blöcke unterteilt und die Operationen blockweise ausgeführt.

Modifikationen des Gaußschen Verfahrens

In die **Level 3 BLAS** oder **BLAS3** wurden schließlich **Matrix-Matrix-Operationen** wie

$$C \leftarrow \alpha AB + \beta C$$

aufgenommen.

Um die **Wiederverwendung von Daten** in den Registern oder im Cache in möglichst hohem Maße zu erreichen, werden die beteiligten Matrizen in Blöcke unterteilt und die Operationen blockweise ausgeführt.

Auf diese Weise erreicht man, dass für die obige Operation bei $4n^2 + 2$ Speicherzugriffen $2n^3 + O(n^2)$ flops ausgeführt werden, das Verhältnis von nützlicher Arbeit zu Speicherzugriffen also auf $n/2$ steigt.

Modifikationen des Gaußschen Verfahrens

Beispiel 4.14: Als Beispiel zur Performance der BLAS3 betrachten wir die Bestimmung des Produktes zweier $(10^3, 10^3)$ Matrizen. Hierfür erhält man die Laufzeiten

Implementierung	CPU Zeit	Relation
naiv	462.42	247.34
BLAS3 (netlib)	320.00	171.1
BLAS3 (f90-HP)	7.25	3.9
BLAS3 (ATLAS)	4.26	2.3
veclib	1.87	1.0



Störungen linearer Systeme

Wir wollen nun den Begriff der **Kondition** einer Matrix einführen, deren Wert — wie oben angedeutet — verantwortlich ist für die numerische Behandelbarkeit eines linearen Gleichungssystems.

Störungen linearer Systeme

Wir wollen nun den Begriff der **Kondition** einer Matrix einführen, deren Wert — wie oben angedeutet — verantwortlich ist für die numerische Behandelbarkeit eines linearen Gleichungssystems.

Wir betrachten dazu neben dem linearen Gleichungssystem

$$\mathbf{Ax} = \mathbf{b} \quad (4.4)$$

mit der regulären Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$

Störungen linearer Systeme

Wir wollen nun den Begriff der **Kondition** einer Matrix einführen, deren Wert — wie oben angedeutet — verantwortlich ist für die numerische Behandelbarkeit eines linearen Gleichungssystems.

Wir betrachten dazu neben dem linearen Gleichungssystem

$$Ax = b \quad (4.4)$$

mit der regulären Matrix $A \in \mathbb{R}^{n \times n}$ ein **gestörtes System**

$$(A + \Delta A)(x + \Delta x) = b + \Delta b \quad (4.5)$$

und fragen uns, wie die Lösung des ursprünglichen Systems auf diese Störungen reagiert.

Störungen linearer Systeme

Bemerkung 4.15: **Kleine Störungen** kann man bei der praktischen Lösung linearer Gleichungssysteme grundsätzlich **nicht ausschließen**.

Störungen linearer Systeme

Bemerkung 4.15: **Kleine Störungen** kann man bei der praktischen Lösung linearer Gleichungssysteme grundsätzlich **nicht ausschließen**.

Einerseits können die Eingangsdaten des Systems aus **Messungen** herrühren und somit a priori fehlerbehaftet sein. Andererseits wird man bei der Benutzung eines elektronischen Rechners durch die endliche Genauigkeit der Zahlendarstellungen auf dem Rechner **Eingabefehler** machen müssen.

Störungen linearer Systeme

Bemerkung 4.15: **Kleine Störungen** kann man bei der praktischen Lösung linearer Gleichungssysteme grundsätzlich **nicht ausschließen**.

Einerseits können die Eingangsdaten des Systems aus **Messungen** herrühren und somit a priori fehlerbehaftet sein. Andererseits wird man bei der Benutzung eines elektronischen Rechners durch die endliche Genauigkeit der Zahlendarstellungen auf dem Rechner **Eingabefehler** machen müssen.

Man muss also grundsätzlich davon ausgehen, dass man mit gestörten Systemen anstelle der wirklich zu lösenden Systeme rechnen muss. Allerdings kann man meistens **annehmen**, dass die **Störungen klein** sind.

Störungen linearer Systeme

Bevor wir die Wirkung von Störungen untersuchen können, benötigen wir noch einige Aussagen über sogenannte **Matrixnormen**, durch die wir die Größe einer Matrix messen.

Störungen linearer Systeme

Bevor wir die Wirkung von Störungen untersuchen können, benötigen wir noch einige Aussagen über sogenannte **Matrixnormen**, durch die wir die Größe einer Matrix messen.

Definition 4.16: Es sei $A \in \mathbb{C}^{(m \times n)}$ eine Matrix, $\|\cdot\|_n$ eine Vektornorm auf \mathbb{C}^n und $\|\cdot\|_m$ eine Vektornorm auf \mathbb{C}^m .

Störungen linearer Systeme

Bevor wir die Wirkung von Störungen untersuchen können, benötigen wir noch einige Aussagen über sogenannte **Matrixnormen**, durch die wir die Größe einer Matrix messen.

Definition 4.16: Es sei $A \in \mathbb{C}^{(m \times n)}$ eine Matrix, $\|\cdot\|_n$ eine Vektornorm auf \mathbb{C}^n und $\|\cdot\|_m$ eine Vektornorm auf \mathbb{C}^m .

Dann heißt

$$\|A\|_{m,n} := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_m}{\|\mathbf{x}\|_n}$$

die **Matrixnorm** von A , die den Normen $\|\cdot\|_n$ und $\|\cdot\|_m$ zugeordnet ist.

Störungen linearer Systeme

Bevor wir die Wirkung von Störungen untersuchen können, benötigen wir noch einige Aussagen über sogenannte **Matrixnormen**, durch die wir die Größe einer Matrix messen.

Definition 4.16: Es sei $A \in \mathbb{C}^{(m \times n)}$ eine Matrix, $\|\cdot\|_n$ eine Vektornorm auf \mathbb{C}^n und $\|\cdot\|_m$ eine Vektornorm auf \mathbb{C}^m .

Dann heißt

$$\|A\|_{m,n} := \max_{\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_m}{\|\mathbf{x}\|_n}$$

die **Matrixnorm** von A , die den Normen $\|\cdot\|_n$ und $\|\cdot\|_m$ zugeordnet ist.

Eine zugeordnete Matrixnorm „misst“ also die **maximale Verlängerung**, die ein Vektor x durch Anwendung von A erfahren kann.

Störungen linearer Systeme

Bemerkung 4.17: Wegen

$$\frac{\|A\mathbf{x}\|_m}{\|\mathbf{x}\|_n} = \left\| A \left(\frac{\mathbf{x}}{\|\mathbf{x}\|_n} \right) \right\|_m$$

Störungen linearer Systeme

Bemerkung 4.17: Wegen

$$\frac{\|\mathbf{Ax}\|_m}{\|\mathbf{x}\|_n} = \left\| \mathbf{A} \left(\frac{\mathbf{x}}{\|\mathbf{x}\|_n} \right) \right\|_m$$

genügt es, das Maximum über Vektoren der **Länge 1** zu erstrecken:

$$\|\mathbf{A}\|_{m,n} = \max\{\|\mathbf{Ay}\|_m : \|\mathbf{y}\|_n = 1\}.$$

Störungen linearer Systeme

Bemerkung 4.17: Wegen

$$\frac{\|\mathbf{Ax}\|_m}{\|\mathbf{x}\|_n} = \left\| \mathbf{A} \left(\frac{\mathbf{x}}{\|\mathbf{x}\|_n} \right) \right\|_m$$

genügt es, das Maximum über Vektoren der **Länge 1** zu erstrecken:

$$\|\mathbf{A}\|_{m,n} = \max\{\|\mathbf{Ay}\|_m : \|\mathbf{y}\|_n = 1\}.$$

Die **Existenz** des Maximums (daß es also einen Vektor \mathbf{x} mit $\|\mathbf{x}\|_n = 1$ und $\|\mathbf{Ax}\|_m = \|\mathbf{A}\|_{m,n}$ gibt) folgt aus der Stetigkeit der Abbildung $\mathbf{x} \mapsto \|\mathbf{Ax}\|_m$. □

Störungen linearer Systeme

Bemerkung 4.18: $\|\cdot\|_{m,n}$ ist eine Norm auf $\mathbb{C}^{(m \times n)}$, denn

$$\begin{aligned}\|\mathbf{A}\|_{m,n} = 0 &\Leftrightarrow \|\mathbf{Ax}\|_m = 0 \quad \forall \mathbf{x} \in \mathbb{C}^n \\ &\Leftrightarrow \mathbf{Ax} = \mathbf{0} \quad \forall \mathbf{x} \in \mathbb{C}^n \\ &\Leftrightarrow \mathbf{A} = \mathbf{O},\end{aligned}$$

$$\begin{aligned}\|\lambda \mathbf{A}\|_{m,n} &= \max\{\|\lambda \mathbf{Ax}\|_m : \|\mathbf{x}\|_n = 1\} \\ &= \max\{|\lambda| \cdot \|\mathbf{Ax}\|_m : \|\mathbf{x}\|_n = 1\} \\ &= |\lambda| \cdot \|\mathbf{A}\|_{m,n},\end{aligned}$$

$$\begin{aligned}\|\mathbf{A} + \mathbf{B}\|_{m,n} &= \max\{\|\mathbf{Ax} + \mathbf{Bx}\|_m : \|\mathbf{x}\|_n = 1\} \\ &\leq \max\{\|\mathbf{Ax}\|_m + \|\mathbf{Bx}\|_m : \|\mathbf{x}\|_n = 1\} \\ &\leq \max\{\|\mathbf{Ax}\|_m : \|\mathbf{x}\|_n = 1\} + \max\{\|\mathbf{Bx}\|_m : \|\mathbf{x}\|_n = 1\} \\ &= \|\mathbf{A}\|_{m,n} + \|\mathbf{B}\|_{m,n}.\end{aligned}$$

Störungen linearer Systeme

Diese ist zusätzlich **submultiplikativ** im folgenden Sinne:

$$\|\mathbf{Ax}\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{x}\|_n \quad \forall \mathbf{x} \in \mathbb{C}^n, \mathbf{A} \in \mathbb{C}^{(m \times n)}, \quad (4.6)$$

$$\|\mathbf{AB}\|_{m,p} \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \quad \forall \mathbf{A} \in \mathbb{C}^{(m \times n)}, \mathbf{B} \in \mathbb{C}^{(n,p)}. \quad (4.7)$$

Störungen linearer Systeme

Diese ist zusätzlich **submultiplikativ** im folgenden Sinne:

$$\|\mathbf{Ax}\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{x}\|_n \quad \forall \mathbf{x} \in \mathbb{C}^n, \mathbf{A} \in \mathbb{C}^{(m \times n)}, \quad (4.6)$$

$$\|\mathbf{AB}\|_{m,p} \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \quad \forall \mathbf{A} \in \mathbb{C}^{(m \times n)}, \mathbf{B} \in \mathbb{C}^{(n,p)}. \quad (4.7)$$

Die Ungleichung (4.6) folgt sofort aus der Definition 4.16; denn es ist

$$\|\mathbf{A}\|_{m,n} \geq \frac{\|\mathbf{Ax}\|_m}{\|\mathbf{x}\|_n}$$

für alle $\mathbf{x} \in \mathbb{C}^n$.

Störungen linearer Systeme

Diese ist zusätzlich **submultiplikativ** im folgenden Sinne:

$$\|\mathbf{Ax}\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{x}\|_n \quad \forall \mathbf{x} \in \mathbb{C}^n, \mathbf{A} \in \mathbb{C}^{(m \times n)}, \quad (4.6)$$

$$\|\mathbf{AB}\|_{m,p} \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \quad \forall \mathbf{A} \in \mathbb{C}^{(m \times n)}, \mathbf{B} \in \mathbb{C}^{(n,p)}. \quad (4.7)$$

Die Ungleichung (4.6) folgt sofort aus der Definition 4.16; denn es ist

$$\|\mathbf{A}\|_{m,n} \geq \frac{\|\mathbf{Ax}\|_m}{\|\mathbf{x}\|_n}$$

für alle $\mathbf{x} \in \mathbb{C}^n$.

Die Ungleichung (4.7) erschließt man mit (4.6) wie folgt: Für alle $\mathbf{x} \in \mathbb{C}^p$ ist

$$\|\mathbf{ABx}\|_m = \|\mathbf{A}(\mathbf{Bx})\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{Bx}\|_n \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \cdot \|\mathbf{x}\|_p,$$

Störungen linearer Systeme

Diese ist zusätzlich **submultiplikativ** im folgenden Sinne:

$$\|\mathbf{Ax}\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{x}\|_n \quad \forall \mathbf{x} \in \mathbb{C}^n, \mathbf{A} \in \mathbb{C}^{(m \times n)}, \quad (4.6)$$

$$\|\mathbf{AB}\|_{m,p} \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \quad \forall \mathbf{A} \in \mathbb{C}^{(m \times n)}, \mathbf{B} \in \mathbb{C}^{(n,p)}. \quad (4.7)$$

Die Ungleichung (4.6) folgt sofort aus der Definition 4.16; denn es ist

$$\|\mathbf{A}\|_{m,n} \geq \frac{\|\mathbf{Ax}\|_m}{\|\mathbf{x}\|_n}$$

für alle $\mathbf{x} \in \mathbb{C}^n$.

Die Ungleichung (4.7) erschließt man mit (4.6) wie folgt: Für alle $\mathbf{x} \in \mathbb{C}^p$ ist

$$\|\mathbf{ABx}\|_m = \|\mathbf{A}(\mathbf{Bx})\|_m \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{Bx}\|_n \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p} \cdot \|\mathbf{x}\|_p,$$

Also ist $\|\mathbf{AB}\|_{m,p} = \max\{\|\mathbf{ABx}\|_m : \|\mathbf{x}\|_p = 1\} \leq \|\mathbf{A}\|_{m,n} \cdot \|\mathbf{B}\|_{n,p}$. □

Störungen linearer Systeme

Bemerkung 4.19: Die Matrixnorm $\|A\|_{m,n}$ ist die **kleinste nichtnegative reelle Zahl** μ , mit der

$$\|Ax\|_m \leq \mu \cdot \|x\|_n \quad \forall x \in \mathbb{C}^n$$

ist. $\|A\|_{m,n}$ ist demnach die **maximale Verlängerung**, die ein x durch Abbildung mit A erfahren kann, wobei x selbst in der $\|\cdot\|_n$ -Norm und Ax in der Norm $\|\cdot\|_m$ des Bildraumes gemessen wird. □

Störungen linearer Systeme

Bemerkung 4.19: Die Matrixnorm $\|A\|_{m,n}$ ist die **kleinste nichtnegative reelle Zahl** μ , mit der

$$\|Ax\|_m \leq \mu \cdot \|x\|_n \quad \forall x \in \mathbb{C}^n$$

ist. $\|A\|_{m,n}$ ist demnach die **maximale Verlängerung**, die ein x durch Abbildung mit A erfahren kann, wobei x selbst in der $\|\cdot\|_n$ -Norm und Ax in der Norm $\|\cdot\|_m$ des Bildraumes gemessen wird. □

Wir betrachten nun **nur noch** den Fall, daß im Urbildraum und im Bildraum **dieselbe Norm** verwendet wird, auch wenn beide Räume verschiedene Dimension haben, und verwenden für die Matrixnorm dasselbe Symbol wie für die Vektornorm.

Störungen linearer Systeme

Bemerkung 4.19: Die Matrixnorm $\|A\|_{m,n}$ ist die **kleinste nichtnegative reelle Zahl** μ , mit der

$$\|Ax\|_m \leq \mu \cdot \|x\|_n \quad \forall x \in \mathbb{C}^n$$

ist. $\|A\|_{m,n}$ ist demnach die **maximale Verlängerung**, die ein x durch Abbildung mit A erfahren kann, wobei x selbst in der $\|\cdot\|_n$ -Norm und Ax in der Norm $\|\cdot\|_m$ des Bildraumes gemessen wird. □

Wir betrachten nun **nur noch** den Fall, daß im Urbildraum und im Bildraum **dieselbe Norm** verwendet wird, auch wenn beide Räume verschiedene Dimension haben, und verwenden für die Matrixnorm dasselbe Symbol wie für die Vektornorm.

Für $A \in \mathbb{R}^{(5,9)}$ bezeichnet also $\|A\|_\infty$ die Matrixnorm von A gemäß Definition 4.16, wenn sowohl im Urbildraum \mathbb{R}^9 als auch im Bildraum \mathbb{R}^5 die **Maximumnorm** verwendet wird.

Störungen linearer Systeme

Satz 4.20

Es sei $A \in \mathbb{C}^{(m \times n)}$ gegeben.

- ▶ Die **Zeilensummennorm**

$$\|A\|_{\infty} := \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}|$$

ist der Maximumnorm $\|x\|_{\infty} := \max_{i=1, \dots, n} |x_i|$ zugeordnet.

- ▶ Die **Spektralnorm**

$$\|A\|_2 := \max\{\sqrt{\lambda} : \lambda \text{ ist Eigenwert von } A^H A\}$$

ist der Euklidischen Norm zugeordnet.

Störungen linearer Systeme

Satz 4.20 (Fortsetzung)

- ▶ Die **Spaltensummennorm**

$$\|\mathbf{A}\|_1 := \max_{j=1,\dots,n} \sum_{i=1}^m |a_{ij}|$$

ist der Summennorm $\|\mathbf{x}\|_1 := \sum_{i=1}^n |x_i|$ zugeordnet.

Störungen linearer Systeme

Satz 4.20 (Fortsetzung)

- Die **Spaltensummennorm**

$$\|A\|_1 := \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|$$

ist der Summennorm $\|x\|_1 := \sum_{i=1}^n |x_i|$ zugeordnet.

Beweis.

Zeilensummennorm: Für alle $x \in \mathbb{C}^n$ gilt

$$\begin{aligned} \|Ax\|_\infty &= \max_{i=1, \dots, m} \left| \sum_{j=1}^n a_{ij} x_j \right| \leq \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}| \cdot |x_j| \\ &\leq \|x\|_\infty \cdot \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}|, \end{aligned}$$

Störungen linearer Systeme

Beweis.

und daher

$$\|\mathbf{A}\|_{\infty} \leq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|. \quad (4.8)$$

Störungen linearer Systeme

Beweis.

und daher

$$\|\mathbf{A}\|_{\infty} \leq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|. \quad (4.8)$$

Sei $k \in \{1, \dots, m\}$ mit

$$\sum_{j=1}^n |a_{ij}| \leq \sum_{j=1}^n |a_{kj}|$$

für alle i .

Störungen linearer Systeme

Beweis.

und daher

$$\|\mathbf{A}\|_{\infty} \leq \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}|. \quad (4.8)$$

Sei $k \in \{1, \dots, m\}$ mit

$$\sum_{j=1}^n |a_{ij}| \leq \sum_{j=1}^n |a_{kj}|$$

für alle i .

Wir definieren $x \in \mathbb{C}^n$ durch $x_j := 1$, falls $a_{kj} = 0$, und $x_j := \overline{a_{kj}}/|a_{kj}|$, sonst.

Störungen linearer Systeme

Beweis.

Dann gilt $\|\mathbf{x}\|_\infty = 1$ und

Störungen linearer Systeme

Beweis.

Dann gilt $\|\mathbf{x}\|_\infty = 1$ und

$$\begin{aligned}\|\mathbf{Ax}\|_\infty &= \max_{i=1,\dots,m} \left| \sum_{j=1}^n a_{ij}x_j \right| \geq \left| \sum_{j=1}^n a_{kj}x_j \right| \\ &= \left| \sum_{j=1}^n |a_{kj}| \right| = \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|,\end{aligned}$$

Störungen linearer Systeme

Beweis.

Dann gilt $\|\mathbf{x}\|_\infty = 1$ und

$$\begin{aligned}\|\mathbf{Ax}\|_\infty &= \max_{i=1,\dots,m} \left| \sum_{j=1}^n a_{ij}x_j \right| \geq \left| \sum_{j=1}^n a_{kj}x_j \right| \\ &= \left| \sum_{j=1}^n |a_{kj}| \right| = \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|,\end{aligned}$$

und daher

$$\begin{aligned}\|\mathbf{A}\|_\infty &= \max\{\|\mathbf{Ay}\|_\infty : \|\mathbf{y}\|_\infty = 1\} \\ &\geq \|\mathbf{Ax}\|_\infty \geq \max_{i=1,\dots,m} \sum_{j=1}^n |a_{ij}|,\end{aligned}$$

zusammen mit (4.8) also die Behauptung.

Störungen linearer Systeme

Beweis.

Spektralnorm: Es ist $A^H A \in \mathbb{C}^{(n \times n)}$ eine Hermitesche Matrix, und daher ist nach dem Rayleighschen Prinzip



Störungen linearer Systeme

Beweis.

Spektralnorm: Es ist $A^H A \in \mathbb{C}^{(n \times n)}$ eine Hermitesche Matrix, und daher ist nach dem Rayleighschen Prinzip

$$\max_{\mathbf{x} \neq \mathbf{0}} \frac{\|A\mathbf{x}\|_2^2}{\|\mathbf{x}\|_2^2} = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^H A^H A \mathbf{x}}{\mathbf{x}^H \mathbf{x}}$$

der maximale Eigenwert von $A^H A$.



Störungen linearer Systeme

Beweis.

Spektralnorm: Es ist $A^H A \in \mathbb{C}^{(n \times n)}$ eine Hermitesche Matrix, und daher ist nach dem Rayleighschen Prinzip

$$\max_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{Ax}\|_2^2}{\|\mathbf{x}\|_2^2} = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}^H \mathbf{A}^H \mathbf{Ax}}{\mathbf{x}^H \mathbf{x}}$$

der maximale Eigenwert von $A^H A$.

Spaltensummennorm: Zeigt man ähnlich wie bei der **Zeilensummennorm**.

Störungen linearer Systeme

Beispiel 4.21: Sei

$$\mathbf{A} = \begin{pmatrix} 1 & 0.1 & -0.1 \\ 0.1 & 2 & -0.4 \\ 0.2 & 0.4 & 3 \end{pmatrix}.$$

Störungen linearer Systeme

Beispiel 4.21: Sei

$$\mathbf{A} = \begin{pmatrix} 1 & 0.1 & -0.1 \\ 0.1 & 2 & -0.4 \\ 0.2 & 0.4 & 3 \end{pmatrix}.$$

Es ist

$$\|\mathbf{A}\|_{\infty} = \max\{1.2, 2.5, 3.6\} = 3.6$$

$$\|\mathbf{A}\|_1 = \max\{1.3, 2.5, 3.5\} = 3.5$$

Störungen linearer Systeme

Die Spektralnorm von A ist die Quadratwurzel aus dem maximalen Eigenwert von

$$A^T A = \begin{pmatrix} 1.05 & 0.38 & 0.46 \\ 0.38 & 4.17 & 0.39 \\ 0.46 & 0.39 & 9.17 \end{pmatrix}.$$

Störungen linearer Systeme

Die Spektralnorm von A ist die Quadratwurzel aus dem maximalen Eigenwert von

$$A^T A = \begin{pmatrix} 1.05 & 0.38 & 0.46 \\ 0.38 & 4.17 & 0.39 \\ 0.46 & 0.39 & 9.17 \end{pmatrix}.$$

Nach einiger Rechnung ergibt dies

$$\|A\|_2 \approx \sqrt{9.2294} \approx 3.04.$$



Störungen linearer Systeme

Die Spektralnorm von A ist die Quadratwurzel aus dem maximalen Eigenwert von

$$A^T A = \begin{pmatrix} 1.05 & 0.38 & 0.46 \\ 0.38 & 4.17 & 0.39 \\ 0.46 & 0.39 & 9.17 \end{pmatrix}.$$

Nach einiger Rechnung ergibt dies

$$\|A\|_2 \approx \sqrt{9.2294} \approx 3.04. \quad \square$$

Die Spektralnorm ist die genaueste (im Sinne der Euklidischen Geometrie des \mathbb{C}^n) und teuerste Norm (im Sinne der Berechenbarkeit) der drei vorgestellten Normen.

Störungen linearer Systeme

Die Spektralnorm von A ist die Quadratwurzel aus dem maximalen Eigenwert von

$$A^T A = \begin{pmatrix} 1.05 & 0.38 & 0.46 \\ 0.38 & 4.17 & 0.39 \\ 0.46 & 0.39 & 9.17 \end{pmatrix}.$$

Nach einiger Rechnung ergibt dies

$$\|A\|_2 \approx \sqrt{9.2294} \approx 3.04. \quad \square$$

Die Spektralnorm ist die genaueste (im Sinne der Euklidischen Geometrie des \mathbb{C}^n) und teuerste Norm (im Sinne der Berechenbarkeit) der drei vorgestellten Normen.

Wir werden noch eine bessere Deutung der Spektralnorm (mittels der sogenannten SVD) kennenlernen, aber oft begnügt man sich mit einer billigeren verwandten Norm, der **Schur-Norm**.

Störungen linearer Systeme

Satz 4.22

Für alle $\mathbf{A} \in \mathbb{C}^{(m \times n)}$ gilt

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F := \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

Störungen linearer Systeme

Satz 4.22

Für alle $\mathbf{A} \in \mathbb{C}^{(m \times n)}$ gilt

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_S := \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

Bemerkung 4.23: $\|\mathbf{A}\|_S$ heißt die **Schur-Norm** oder **Erhard-Schmidt-Norm** (oder — speziell in der anglo-amerikanischen Literatur — **Frobenius-Norm** der Matrix \mathbf{A}).

Störungen linearer Systeme

Satz 4.22

Für alle $\mathbf{A} \in \mathbb{C}^{(m \times n)}$ gilt

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_S := \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

Bemerkung 4.23: $\|\mathbf{A}\|_S$ heißt die **Schur-Norm** oder **Erhard-Schmidt-Norm** (oder — speziell in der anglo-amerikanischen Literatur — **Frobenius-Norm** der Matrix \mathbf{A}). Diese ist zwar eine Vektornorm auf dem $\mathbb{C}^{(m \times n)}$ (die Euklidische Norm auf dem $\mathbb{C}^{m \cdot n}$), aber keine einer Vektornorm zugeordnete Matrixnorm, denn für eine solche gilt im Fall $n = m$ stets

$$\|\mathbf{E}\| = \max\{\|\mathbf{E}\mathbf{x}\| : \|\mathbf{x}\| = 1\} = 1.$$

Es ist aber $\|\mathbf{E}\|_S = \sqrt{n}$. □

Störungen linearer Systeme

Beweis.

Wegen der Cauchy-Schwarzschen Ungleichung gilt für alle $\mathbf{x} \in \mathbb{C}^n$

$$\|\mathbf{Ax}\|_2^2 = \sum_{i=1}^m \left| \sum_{j=1}^n a_{ij}x_j \right|^2 \leq \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right) \sum_{j=1}^n |x_j|^2 = \|\mathbf{A}\|_F^2 \cdot \|\mathbf{x}\|_2^2,$$

und daher gilt $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_F$. □

Störungen linearer Systeme

Beweis.

Wegen der Cauchy-Schwarzschen Ungleichung gilt für alle $\mathbf{x} \in \mathbb{C}^n$

$$\|\mathbf{Ax}\|_2^2 = \sum_{i=1}^m \left| \sum_{j=1}^n a_{ij}x_j \right|^2 \leq \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2 \right) \sum_{j=1}^n |x_j|^2 = \|\mathbf{A}\|_S^2 \cdot \|\mathbf{x}\|_2^2,$$

und daher gilt $\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_S$. □

Beispiel 4.24: Für die Matrix \mathbf{A} aus Beispiel 4.21 erhält man

$$\|\mathbf{A}\|_2 \leq \|\mathbf{A}\|_S = \sqrt{14.39} \leq 3.80. \quad \square$$

Störungen linearer Systeme

Wir betrachten nun das gestörte System (4.5) und nehmen an, dass die **Störungen** der Matrixelemente so **klein** sind, dass auch die Matrix $A + \Delta A$ **regulär** ist (dass dies für hinreichend kleine Störungen bei regulärer Matrix A überhaupt stets möglich ist, wird aus dem Satz 4.25 unten folgen).

Störungen linearer Systeme

Wir betrachten nun das gestörte System (4.5) und nehmen an, dass die **Störungen** der Matrixelemente so **klein** sind, dass auch die Matrix $A + \Delta A$ **regulär** ist (dass dies für hinreichend kleine Störungen bei regulärer Matrix A überhaupt stets möglich ist, wird aus dem Satz 4.25 unten folgen).

Löst man mit dieser Annahme (4.5) nach Δx auf, so erhält man für den durch die Störungen ΔA und Δb hervorgerufenen **absoluten Fehler** wegen $Ax = b$

$$\begin{aligned}\Delta x &= (A + \Delta A)^{-1}(\Delta b - \Delta Ax) \\ &= (E + A^{-1}\Delta A)^{-1}A^{-1}(\Delta b - \Delta Ax).\end{aligned}$$

Störungen linearer Systeme

Wir betrachten nun das gestörte System (4.5) und nehmen an, dass die **Störungen** der Matrixelemente so **klein** sind, dass auch die Matrix $\mathbf{A} + \Delta\mathbf{A}$ **regulär** ist (dass dies für hinreichend kleine Störungen bei regulärer Matrix \mathbf{A} überhaupt stets möglich ist, wird aus dem Satz 4.25 unten folgen).

Löst man mit dieser Annahme (4.5) nach $\Delta\mathbf{x}$ auf, so erhält man für den durch die Störungen $\Delta\mathbf{A}$ und $\Delta\mathbf{b}$ hervorgerufenen **absoluten Fehler** wegen $\mathbf{A}\mathbf{x} = \mathbf{b}$

$$\begin{aligned}\Delta\mathbf{x} &= (\mathbf{A} + \Delta\mathbf{A})^{-1}(\Delta\mathbf{b} - \Delta\mathbf{A}\mathbf{x}) \\ &= (\mathbf{E} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\mathbf{A}^{-1}(\Delta\mathbf{b} - \Delta\mathbf{A}\mathbf{x}).\end{aligned}$$

Mit einer beliebigen Vektornorm $\|\cdot\|$ auf dem \mathbb{R}^n und der gleich bezeichneten zugeordneten Matrixnorm kann man also **abschätzen**:

$$\|\Delta\mathbf{x}\| \leq \|(\mathbf{E} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\| \cdot \|\mathbf{A}^{-1}\| (\|\Delta\mathbf{b}\| + \|\Delta\mathbf{A}\| \cdot \|\mathbf{x}\|).$$

Störungen linearer Systeme

Für $\mathbf{b} \neq \mathbf{0}$ und folglich $\mathbf{x} \neq \mathbf{0}$ erhält man daraus für den **relativen Fehler von \mathbf{x}** , also $\|\Delta\mathbf{x}\|/\|\mathbf{x}\|$, die Abschätzung

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|(\mathbf{E} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\| \cdot \|\mathbf{A}^{-1}\| \left(\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{x}\|} + \|\Delta\mathbf{A}\| \right). \quad (4.9)$$

Störungen linearer Systeme

Für $\mathbf{b} \neq \mathbf{0}$ und folglich $\mathbf{x} \neq \mathbf{0}$ erhält man daraus für den **relativen Fehler von \mathbf{x}** , also $\|\Delta\mathbf{x}\|/\|\mathbf{x}\|$, die Abschätzung

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \|(\mathbf{E} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\| \cdot \|\mathbf{A}^{-1}\| \left(\frac{\|\Delta\mathbf{b}\|}{\|\mathbf{x}\|} + \|\Delta\mathbf{A}\| \right). \quad (4.9)$$

Um in dieser Ungleichung $\|(\mathbf{E} + \mathbf{A}^{-1}\Delta\mathbf{A})^{-1}\|$ weiter abzuschätzen, benötigen wir das folgende **Störungslemma**, das gleichzeitig darüber Auskunft gibt, unter wie großen Störungen der Matrixelemente die **Existenz der Inversen** für die gestörte Matrix noch gesichert ist.

Störungen linearer Systeme

Satz 4.25 (Störungslemma)

Es sei $\mathbf{B} \in \mathbb{R}^{n \times n}$ und es gelte für eine beliebige einer Vektornorm zugeordneten Matrixnorm die Ungleichung $\|\mathbf{B}\| < 1$. Dann ist die Matrix $\mathbf{E} - \mathbf{B}$ regulär, und es gilt

$$\|(\mathbf{E} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}.$$

Störungen linearer Systeme

Satz 4.25 (Störungslemma)

Es sei $\mathbf{B} \in \mathbb{R}^{n \times n}$ und es gelte für eine beliebige einer Vektornorm zugeordneten Matrixnorm die Ungleichung $\|\mathbf{B}\| < 1$. Dann ist die Matrix $\mathbf{E} - \mathbf{B}$ regulär, und es gilt

$$\|(\mathbf{E} - \mathbf{B})^{-1}\| \leq \frac{1}{1 - \|\mathbf{B}\|}.$$

Beweis.

Für alle $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x} \neq \mathbf{0}$, gilt

$$\|(\mathbf{E} - \mathbf{B})\mathbf{x}\| \geq \|\mathbf{x}\| - \|\mathbf{B}\mathbf{x}\| \geq \|\mathbf{x}\| - \|\mathbf{B}\|\|\mathbf{x}\| = (1 - \|\mathbf{B}\|)\|\mathbf{x}\| > 0,$$

d.h. $(\mathbf{E} - \mathbf{B})\mathbf{x} = \mathbf{0}$ ist nur für $\mathbf{x} = \mathbf{0}$ lösbar und daher ist $\mathbf{E} - \mathbf{B}$ regulär.

Störungen linearer Systeme

Beweis.

Die Abschätzung der Norm der Inversen von $E - B$ erhält man so:

$$\begin{aligned} 1 &= \|(E - B)^{-1}(E - B)\| \\ &= \|(E - B)^{-1} - (E - B)^{-1}B\| \\ &\geq \|(E - B)^{-1}\| - \|(E - B)^{-1}B\| \\ &\geq \|(E - B)^{-1}\| - \|(E - B)^{-1}\| \cdot \|B\| \\ &= (1 - \|B\|) \cdot \|(E - B)^{-1}\|. \end{aligned}$$



Störungen linearer Systeme

Mit dem Störungslemma (mit $B = -A^{-1} \Delta A$) können wir nun den **relativen Fehler** des gestörten Problems aus (4.9) weiter abschätzen, wobei wir $\|A^{-1} \Delta A\| \leq \|A^{-1}\| \cdot \|\Delta A\|$ und $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$ beachten:

$$\begin{aligned} \frac{\|\Delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\Delta A\|} \left(\|A\| \frac{\|\Delta b\|}{\|b\|} + \|\Delta A\| \right) \\ &\leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \frac{\|\Delta A\|}{\|A\|}} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right). \end{aligned} \quad (4.10)$$

Störungen linearer Systeme

Mit dem Störungslemma (mit $B = -A^{-1} \Delta A$) können wir nun den **relativen Fehler** des gestörten Problems aus (4.9) weiter abschätzen, wobei wir $\|A^{-1} \Delta A\| \leq \|A^{-1}\| \cdot \|\Delta A\|$ und $\|b\| = \|Ax\| \leq \|A\| \cdot \|x\|$ beachten:

$$\begin{aligned} \frac{\|\Delta x\|}{\|x\|} &\leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\Delta A\|} \left(\|A\| \frac{\|\Delta b\|}{\|b\|} + \|\Delta A\| \right) \\ &\leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \frac{\|\Delta A\|}{\|A\|}} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta b\|}{\|b\|} \right). \end{aligned} \quad (4.10)$$

Aus der Abschätzung (4.10) liest man ab, dass für kleine Störungen der Matrixelemente (so dass der Nenner nicht wesentlich von 1 abweicht) der relative Fehler der rechten Seite und der Matrixelemente um den Faktor $\|A^{-1}\| \cdot \|A\|$ verstärkt wird. Diesen Verstärkungsfaktor nennen wir die **Kondition der Matrix A**.

Störungen linearer Systeme

Definition 4.26: Es sei $A \in \mathbb{R}^{n \times n}$ regulär und $\| \cdot \|_p$ eine einer (gleichbezeichneten) Vektornorm zugeordnete Matrixnorm auf $\mathbb{R}^{n \times n}$.

Störungen linearer Systeme

Definition 4.26: Es sei $A \in \mathbb{R}^{n \times n}$ regulär und $\|\cdot\|_p$ eine einer (gleichbezeichneten) Vektornorm zugeordnete Matrixnorm auf $\mathbb{R}^{n \times n}$.

Dann heißt

$$\kappa_p(A) := \|A^{-1}\|_p \cdot \|A\|_p$$

die **Kondition** der Matrix A (oder des linearen Gleichungssystems (4.4)) bezüglich der Norm $\|\cdot\|_p$.

Störungen linearer Systeme

Definition 4.26: Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ regulär und $\|\cdot\|_p$ eine einer (gleichbezeichneten) Vektornorm zugeordnete Matrixnorm auf $\mathbb{R}^{n \times n}$.

Dann heißt

$$\kappa_p(\mathbf{A}) := \|\mathbf{A}^{-1}\|_p \cdot \|\mathbf{A}\|_p$$

die **Kondition** der Matrix \mathbf{A} (oder des linearen Gleichungssystems (4.4)) bezüglich der Norm $\|\cdot\|_p$.

Wir fassen unsere Überlegungen zusammen:

Störungen linearer Systeme

Satz 4.27

Es seien $\mathbf{A}, \Delta\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{b}, \Delta\mathbf{b} \in \mathbb{R}^n$, $\mathbf{b} \neq \mathbf{0}$, gegeben, so dass \mathbf{A} regulär ist und $\|\mathbf{A}^{-1}\| \cdot \|\Delta\mathbf{A}\| < 1$ mit einer Vektornorm zugeordneten Matrixnorm $\|\cdot\|$ gilt. Dann existiert neben der Lösung des linearen Gleichungssystems (4.4) auch die Lösung $\mathbf{x} + \Delta\mathbf{x}$ des gestörten Systems (4.5) und es gilt mit der Kondition $\kappa(\mathbf{A}) := \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|$ die Abschätzung

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(\mathbf{A})}{1 - \kappa(\mathbf{A}) \cdot \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left(\frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} \right).$$

Störungen linearer Systeme

Satz 4.27

Es seien $\mathbf{A}, \Delta\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{b}, \Delta\mathbf{b} \in \mathbb{R}^n$, $\mathbf{b} \neq \mathbf{0}$, gegeben, so dass \mathbf{A} regulär ist und $\|\mathbf{A}^{-1}\| \cdot \|\Delta\mathbf{A}\| < 1$ mit einer Vektornorm zugeordneten Matrixnorm $\|\cdot\|$ gilt. Dann existiert neben der Lösung des linearen Gleichungssystems (4.4) auch die Lösung $\mathbf{x} + \Delta\mathbf{x}$ des gestörten Systems (4.5) und es gilt mit der Kondition $\kappa(\mathbf{A}) := \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\|$ die Abschätzung

$$\frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{\kappa(\mathbf{A})}{1 - \kappa(\mathbf{A}) \cdot \frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|}} \left(\frac{\|\Delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\Delta\mathbf{b}\|}{\|\mathbf{b}\|} \right).$$

Bemerkung 4.28: Für jede reguläre Matrix \mathbf{A} und jede Norm $\|\cdot\|$ gilt $\kappa(\mathbf{A}) \geq 1$, denn

$$1 = \|\mathbf{E}\| = \|\mathbf{A}\mathbf{A}^{-1}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1}\| = \kappa(\mathbf{A}).$$

□

Störungen linearer Systeme

Bemerkung 4.29: Werden die Rechnungen bei der Lösung eines linearen Gleichungssystems mit der Mantissenlänge ℓ ausgeführt, so haben die Daten von A und b bereits einen relativen Eingabefehler der Größe $5 \cdot 10^{-\ell}$.

Störungen linearer Systeme

Bemerkung 4.29: Werden die Rechnungen bei der Lösung eines linearen Gleichungssystems mit der Mantissenlänge ℓ ausgeführt, so haben die Daten von A und b bereits einen relativen Eingabefehler der Größe $5 \cdot 10^{-\ell}$.

Gilt $\kappa(A) = 10^\gamma$, so ist (abgesehen von Rundungsfehlern, die sich im numerischen Lösungsverfahren ergeben) für die numerische Lösung mit einem relativen Fehler der Größe $5 \cdot 10^{\gamma-\ell}$ zu rechnen.

Störungen linearer Systeme

Bemerkung 4.29: Werden die Rechnungen bei der Lösung eines linearen Gleichungssystems mit der Mantissenlänge ℓ ausgeführt, so haben die Daten von A und b bereits einen relativen Eingabefehler der Größe $5 \cdot 10^{-\ell}$.

Gilt $\kappa(A) = 10^\gamma$, so ist (abgesehen von Rundungsfehlern, die sich im numerischen Lösungsverfahren ergeben) für die numerische Lösung mit einem relativen Fehler der Größe $5 \cdot 10^{\gamma-\ell}$ zu rechnen.

Grob gesprochen **verliert man** also beim Lösen eines linearen Gleichungssystems γ **Stellen**, wenn die Koeffizientenmatrix eine Kondition der Größenordnung 10^γ besitzt. Dieser Verlust von Stellen ist nicht dem jeweilig verwendeten Algorithmus zuzuschreiben. Er ist **problemimmanent**. \square

Störungen linearer Systeme

Beispiel 4.30: Wir betrachten das lineare Gleichungssystem

$$\begin{pmatrix} 1 & 1 \\ 1 & 0.999 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 2 \\ 1.999 \end{pmatrix},$$

das offensichtlich die **Lösung** $\mathbf{x} = (1, 1)^T$ besitzt.

Störungen linearer Systeme

Beispiel 4.30: Wir betrachten das lineare Gleichungssystem

$$\begin{pmatrix} 1 & 1 \\ 1 & 0.999 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 2 \\ 1.999 \end{pmatrix},$$

das offensichtlich die Lösung $\mathbf{x} = (1, 1)^T$ besitzt.

Für den Vektor $\mathbf{x} + \Delta\mathbf{x} := (5, -3.002)^T$ gilt

$$A(\mathbf{x} + \Delta\mathbf{x}) = \begin{pmatrix} 1.998 \\ 2.001002 \end{pmatrix} =: \mathbf{b} + \Delta\mathbf{b}.$$

Störungen linearer Systeme

Beispiel 4.30: Wir betrachten das lineare Gleichungssystem

$$\begin{pmatrix} 1 & 1 \\ 1 & 0.999 \end{pmatrix} \mathbf{x} = \begin{pmatrix} 2 \\ 1.999 \end{pmatrix},$$

das offensichtlich die Lösung $\mathbf{x} = (1, 1)^T$ besitzt.

Für den Vektor $\mathbf{x} + \Delta\mathbf{x} := (5, -3.002)^T$ gilt

$$A(\mathbf{x} + \Delta\mathbf{x}) = \begin{pmatrix} 1.998 \\ 2.001002 \end{pmatrix} =: \mathbf{b} + \Delta\mathbf{b}.$$

Es ist also

$$\frac{\|\Delta\mathbf{b}\|_\infty}{\|\mathbf{b}\|_\infty} = 1.001 \cdot 10^{-3} \quad \text{und} \quad \frac{\|\Delta\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} = 4.002,$$

Störungen linearer Systeme

... und daher gilt für die Kondition

$$\kappa_{\infty}(\mathbf{A}) \geq \frac{4.002}{1.001} 10^3 = 3998.$$

Störungen linearer Systeme

... und daher gilt für die Kondition

$$\kappa_{\infty}(\mathbf{A}) \geq \frac{4.002}{1.001} 10^3 = 3998.$$

Tatsächlich gilt

$$\mathbf{A}^{-1} = \begin{pmatrix} -999 & 1000 \\ 1000 & -1000 \end{pmatrix}$$

und daher $\kappa_{\infty}(\mathbf{A}) = 4000$.

Störungen linearer Systeme

... und daher gilt für die Kondition

$$\kappa_{\infty}(\mathbf{A}) \geq \frac{4.002}{1.001} 10^3 = 3998.$$

Tatsächlich gilt

$$\mathbf{A}^{-1} = \begin{pmatrix} -999 & 1000 \\ 1000 & -1000 \end{pmatrix}$$

und daher $\kappa_{\infty}(\mathbf{A}) = 4000$.

Man sieht an diesem Beispiel, dass die Abschätzung in (4.10) scharf ist. \square

Störungen linearer Systeme

Der nächste Satz gibt eine geometrische Charakterisierung der Konditionszahl. Er sagt, dass der relative Abstand einer regulären Matrix zur nächsten singulären Matrix in der Euklidischen Norm gleich dem Reziproken der Kondition ist.

Störungen linearer Systeme

Der nächste Satz gibt eine geometrische Charakterisierung der Konditionszahl. Er sagt, dass der relative Abstand einer regulären Matrix zur nächsten singulären Matrix in der Euklidischen Norm gleich dem Reziproken der Kondition ist.

Satz 4.31

Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ regulär. Dann gilt

$$\min \left\{ \frac{\|\Delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} : \mathbf{A} + \Delta \mathbf{A} \text{ singular} \right\} = (\kappa_2(\mathbf{A}))^{-1}.$$

Störungen linearer Systeme

Der nächste Satz gibt eine geometrische Charakterisierung der Konditionszahl. Er sagt, dass der relative Abstand einer regulären Matrix zur nächsten singulären Matrix in der Euklidischen Norm gleich dem Reziproken der Kondition ist.

Satz 4.31

Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ regulär. Dann gilt

$$\min \left\{ \frac{\|\Delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} : \mathbf{A} + \Delta \mathbf{A} \text{ singular} \right\} = (\kappa_2(\mathbf{A}))^{-1}.$$

Beweis.

Siehe dazu das Skript „Grundlagen der Numerischen Mathematik“ von Heinrich Voß, Abschnitt 4.4. □

Software für lineare Gleichungssysteme

Sehr hochwertige Public Domain Software ist in den Bibliotheken [LAPACK](#) und [ScaLAPACK](#) erhältlich unter der Adresse

<http://www.netlib.org/lapack/>

bzw.

<http://www.netlib.org/scalapack/>

Software für lineare Gleichungssysteme

Die FORTRAN 77-Bibliothek LAPACK (und die Übertragungen in andere Sprachen: die C-Version CLAPACK, die C++ Version LAPACK++ und die FORTRAN 90 Version LAPACK90, die ebenfalls in der Netlib frei erhältlich sind) ist für PCs, Workstations, Vektorrechner oder Parallelrechner mit gemeinsamen Speicher geeignet, ScaLAPACK für Parallelrechner mit verteiltem Speicher oder vernetzte Workstations.

Software für lineare Gleichungssysteme

Die FORTRAN 77-Bibliothek LAPACK (und die Übertragungen in andere Sprachen: die C-Version CLAPACK, die C++ Version LAPACK++ und die FORTRAN 90 Version LAPACK90, die ebenfalls in der Netlib frei erhältlich sind) ist für PCs, Workstations, Vektorrechner oder Parallelrechner mit gemeinsamen Speicher geeignet, ScaLAPACK für Parallelrechner mit verteiltem Speicher oder vernetzte Workstations.

Beide Bibliotheken wurden unter Benutzung von BLAS3 geschrieben. Der Quellcode ist frei zugänglich und sehr gut dokumentiert. Die kommerziellen Bibliotheken IMSL oder NAG verwenden (zum Teil geringfügig modifizierte) LAPACK-Routinen. MATLAB verwendet ab Version 6.1 LAPACK-Routinen.