

Vorlesung 401-1652-10L : Numerische Mathematik I

Frühlingssemester 2014 : Prof. Ralf Hiptmair (SAM, D-MATH)

[Vorlesungsniederschrift, Rohformat, URL: <http://www.math.ethz.ch/education/bachelor/lectures/fs2014/math/nm>]

Kapitel 2 : Direkte Lösungsverfahren für lineare Gleichungssysteme (LGS)

LGS : $\underline{A}\underline{x} = \underline{b}$, $\underline{b} \in \mathbb{C}^n$, $\underline{A} \in \mathbb{C}^{n,n}$, A regulär (=invertierbar) $\Rightarrow [\underline{x} = \underline{A}^{-1}\underline{b}]$

Grundsatz : Berechne nie numerisch die Inverse einer Matrix

$$\begin{pmatrix} a_{11} & a_{12} & \cdots & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & \cdots & a_{2n} \\ \vdots & \vdots & \cdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \cdots & \vdots \\ a_{n1} & a_{n2} & \cdots & \cdots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \underline{x} \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \underline{b} \\ \vdots \\ b_n \end{pmatrix}$$

Lösungsalgorithmus : Gaußelimination

$$\begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \iff \begin{cases} x_1 + x_2 = 4 \\ 2x_1 + x_2 - x_3 = 1 \\ 3x_1 - x_2 - x_3 = -3 \end{cases}$$

 $\hat{=}$ Pivotzeile

$$\begin{pmatrix} 1 & 1 & 0 \\ 2 & 1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ 1 \\ -3 \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{1} & 1 & 0 \\ 0 & -1 & -1 \\ 3 & -1 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -3 \end{pmatrix} \rightarrow \begin{pmatrix} \mathbf{1} & 1 & 0 \\ 0 & -1 & -1 \\ 0 & -4 & -1 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ -15 \end{pmatrix}$$

$$\rightarrow \begin{pmatrix} 1 & 1 & 0 \\ 0 & \mathbf{-1} & -1 \\ 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} 4 \\ -7 \\ 13 \end{pmatrix}$$

Obere Dreiecksmatrix

gestaffeltes Gleichungssystem, "äquivalent" (= gleiche Lösungsmenge) zum Ausgangs-LGS

Dann Rücksubstitution $x_3 \rightarrow x_2 \rightarrow x_1$

```

1 function x = GENopivot(A,b)
2 % Raw Gaussian elimination without pivoting for
3 % linear system Ax=b: "Linear algebra version", UNSTABLE!
4 [m,n] = size(A); A = [A,b]; ← rechte Seite → letzte Spalte
5 % Forward elimination
6 for i=1:n-1
7 if (A(i,i) == 0.0), error('Zero pivot'); end
8   for j=i+1:n
9     fac = A(j,i)/A(i,i);
10    for k=i+1:n+1
11      A(j,k) = A(j,k) - fac*A(i,k);
12    end
13  end
14 end
15 % Backward substitution
16 x = [zeros(n-1,1);A(n,n+1)/A(n,n)];
17 for i=n-1:-1:1
18   x(i) = A(i,n+1);
19   for k=i+1:n
20     x(i) = x(i) - A(i,k)*x(k);
21   end
22   x(i) = x(i)/A(i,i);
23 end

```

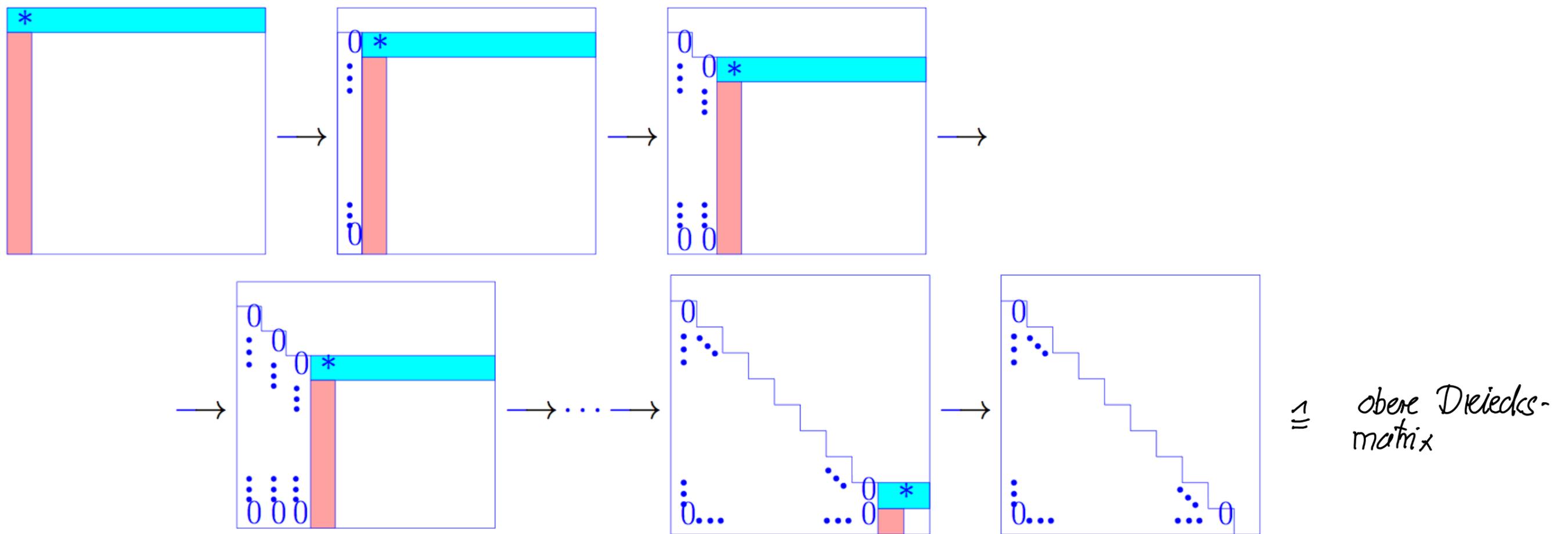
Numerische Sünde! wg. Rundungsfehlern

Anstatt $(x == 0)$ Frage ob

$$\text{abs}(x) \leq \text{eps} * \text{abs}(x_{\text{ref}})$$

← hier $A(1:n,1:n) \hat{=}$ obere Dreiecksmatrix (bei exakter Arithmetik)

Vorwärts elimination :



2.1. Gestaffelte LGS

$$\underline{L} \underline{x} = \begin{pmatrix} l_{11} & 0 & \cdots & \cdots & 0 \\ l_{21} & l_{22} & 0 & \cdots & \\ l_{31} & l_{32} & \cdots & \cdots & \vdots \\ \vdots & \vdots & \cdots & \cdots & 0 \\ l_{n1} & l_{n2} & \cdots & l_{nn-1} & l_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ \vdots \\ b_n \end{pmatrix}$$

untere Dreiecksmatrix, $l_{ii} \neq 0 \Rightarrow L$ regulär

Algorithmus 2.1 Vorwärtseinsetzen

Input: Reguläre untere Dreiecksmatrix $L \in \mathbb{R}^{n \times n}$, $\underline{b} \in \mathbb{R}^n$.

Output: Lösung \underline{x} von $L\underline{x} = \underline{b}$.

for $i = 1, 2, \dots, n$ **do** $x_1 \rightarrow x_2 \rightarrow \dots \rightarrow x_n$

$$x_i := \frac{1}{l_{ii}} \left(b_i - \sum_{k=1}^{i-1} l_{ik} x_k \right)$$

end for

$$U\underline{x} = \begin{pmatrix} u_{11} & u_{12} & \cdots & \cdots & u_{1n} \\ 0 & u_{22} & \cdots & \cdots & u_{2n} \\ 0 & 0 & \ddots & & \vdots \\ \vdots & \vdots & \cdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & u_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ \vdots \\ b_n \end{pmatrix}$$

$u_{ii} \neq 0$

Algorithmus 2.2 Rückwärtseinsetzen

Input: Reguläre obere Dreiecksmatrix $U \in \mathbb{R}^{n \times n}$, $\underline{b} \in \mathbb{R}^n$.

Output: Lösung \underline{x} von $U\underline{x} = \underline{b}$.

for $i = n, n-1, \dots, 1$ **do**

$$x_i := \frac{1}{u_{ii}} \left(b_i - \sum_{k=i+1}^n u_{ik} x_k \right)$$

end for

Rechenaufwand: # der für einen Algorithmus erforderlichen elementaren Operationen (Flops)

Alg. 2.2: \updownarrow R.A. = $\sum_{i=1}^n \underset{\substack{\uparrow \\ \text{DIV}}}{1} + \underset{\substack{\uparrow \\ \text{SUB}}}{1} + \underset{\substack{\uparrow \\ \text{MUL}}}{(n-i)} + \underset{\substack{\uparrow \\ \text{ADD}}}{(n-i-1)} = 2n + 2n^2 - n - n(n+1) = n^2$

[ertl. negative Terme = 0 setzen]

koreliert mit Rechenzeit

Einfluss von Rundungsfehlern

Lemma 2.4. Die von forwardsubst (oder Algorithmus 2.2 für U anstelle von L) in Gleitkommaarithmetik \mathbb{F} berechnete Lösung \hat{x} erfüllt

$$(L + \Delta L)\hat{x} = \underline{b}, \quad \text{für ein } \Delta L \text{ mit } |\Delta L| \leq \gamma_n(\mathbb{F})|L|,$$

wobei $\gamma_n(\mathbb{F}) := \frac{nu(\mathbb{F})}{1-nu(\mathbb{F})}$ (wie in Lemma 1.15).

↑
Komponentenweise Abschätzung

```

1 function x = forwardsubst(L,b)
2 % Solution of lower triangular linear system of equations
3 % by means of forward substitution
4 [n,m] = size(L); x = zeros(n,1);
5 for i=1:n
6     for j=1:i-1, x(i) = x(i) + L(i,j)*x(j); end
7     x(i) = (b(i) - x(i))/L(i,i);
8 end
    
```



numerisch stabil

$\Delta L \triangleq$ Rückwärtsfehler : Lemma 2.4 \Rightarrow

Beweis : (\rightarrow siehe Rundungsfehleranalyse für Euklidisches Skalarprodukt, Kap. 1, Seite 10)

$$\tilde{x}_i = (b_i - \sum_{j=1}^{i-1} l_{ij} (1 + \theta_j) \tilde{x}_j) (1 + \delta) / l_{ii} (1 + \delta), \quad |\theta_j| \leq \gamma_j(\mathbb{F}) \leq \gamma_n(\mathbb{F})$$

$$\Rightarrow \underbrace{l_{ii} (1 + \theta_2)}_{=: \tilde{l}_{ii}} \tilde{x}_i = b_i - \sum_{j=1}^{i-1} \underbrace{l_{ij} (1 + \theta_j)}_{=: \tilde{l}_{ij}} \tilde{x}_j \quad \frac{1}{(1+\delta)(1+\delta)} = 1 + \theta_2$$

$$\Rightarrow \tilde{L} \tilde{x} = \underline{b}, \quad \tilde{L} = L + \Delta L, \quad (\Delta L)_{ij} = l_{ij} \cdot \theta_j$$

□

2.2. LU-Zerlegung und Gaußelimination

Neue Perspektive : $A \underline{x} = \underline{b}$, $A \in \mathbb{C}^{n \times n}$ regulär, $\underline{b} \in \mathbb{C}^n$

Es sei bekannt : $A = L \cdot U$, $L, U \in \mathbb{C}^{n \times n}$

$$LU \underline{x} = \underline{b} \Rightarrow \underbrace{\underline{z} = L^{-1} \underline{b}}_{\text{gestaffelte Gleichungssysteme}}, \quad \underline{x} = U^{-1} \underline{z}$$

$L \triangleq$ untere Dreiecksmatrix } invertierbar
 $U \triangleq$ obere Dreiecksmatrix }

[MATLAB : $\underline{x} = U \setminus (L \setminus \underline{b})$]

Lösen von $A\underline{x} = \underline{b}$ in MATLAB: $\underline{x} = A \backslash \underline{b}$ 'backslash'

The diagram shows three matrices in large blue parentheses. The first is a yellow square labeled 'A'. This is followed by an equals sign. The second is a yellow square representing a lower triangular matrix L, with a blue diagonal line and a '1' in the top-left corner. The third is a yellow square representing an upper triangular matrix U, with a blue diagonal line and a '0' in the bottom-left corner. A small blue 'L' is written on the diagonal of the second matrix, and a small blue 'U' is written on the diagonal of the third matrix.

Definition 2.5 (LU-Zerlegung)

Sei $A \in \mathbb{R}^{n \times n}$. Dann besitzt A eine **LU-Zerlegung**, falls es eine obere Dreiecksmatrix U und eine untere Dreiecksmatrix L mit $l_{11} = \dots = l_{nn} = 1$ gibt, so dass $A = LU$.

Satz 2.6. (Eindeutigkeit der LU-Zerlegung)

Sei $A \in \mathbb{R}^{n \times n}$ regulär und habe eine LU-Zerlegung $LU = A$. Dann ist $(U)_{ii} \neq 0$ für $i = 1, \dots, n$, und die Zerlegung ist eindeutig.

Beweis: L, U sind regulär, $A = L_1 M_1 = L_2 M_2 \Rightarrow M_1 M_2^{-1} = L_1^{-1} L_2 = I$

Lemma Die invertierbaren oberen und unteren Dreiecksmatrizen bilden eine Gruppe mit der Matrixmultiplikation als Verknüpfung.

\uparrow obere Δ -Mat
 \uparrow untere Δ -Mat
 (mit Einheitsdiagonalen)

$\Rightarrow M_1 = M_2 \wedge L_1 = L_2$
 (wg. Eindeutigkeit der Inversen) □

Theorem 2.9. (Existenz einer LU-Zerlegung)

Sei $A \in \mathbb{R}^{n \times n}$ regulär. Dann besitzt A eine LU-Zerlegung genau dann, wenn alle **Hauptminoren** $A_k, k = 1, \dots, n$ regulär sind. Hierbei ist $A_k \in \mathbb{R}^{k \times k}$ gegeben durch $(A_k)_{ij} = (A)_{ij}, i, j = 1, \dots, k$. [MATLAB: $A_k = A(1:k, 1:k)$]



Beweis: Induktion : $n = 1$ ✓

$n-1 \rightarrow n$:

Ansatz
(Block-Matrix-Perspektive)

$$\left(\begin{array}{c|c} A_{n-1} & \underline{b} \\ \hline \underline{a}^T & \alpha \end{array} \right)$$

regulär

LU-Zerlegung von A_n

$$= \left(\begin{array}{c|c} L_{n-1} & 0 \\ \hline \underline{x}^T & 1 \end{array} \right) \left(\begin{array}{c|c} U_{n-1} & \underline{y} \\ \hline 0 & \xi \end{array} \right)$$

$$\Rightarrow \begin{aligned} L_{n-1} \underline{y} &= \underline{b} \Rightarrow \underline{y} = L_{n-1}^{-1} \underline{b} \\ \underline{x}^T U_{n-1} &= \underline{a}^T \Rightarrow \underline{x} = (U_{n-1})^{-T} \underline{a} \end{aligned}$$

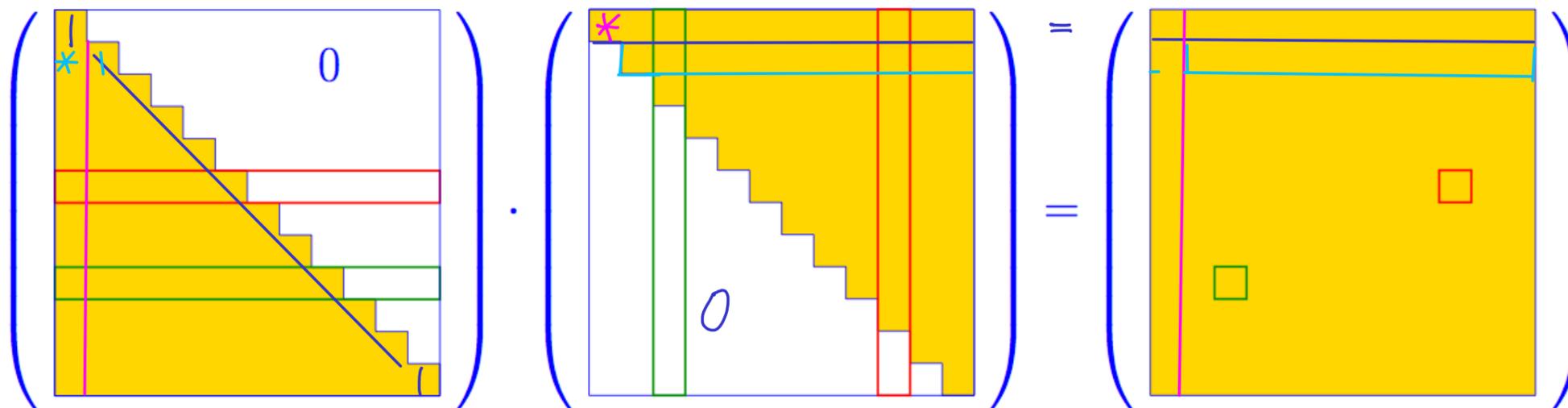
Induktionsannahme :

$$A_{n-1} = \underset{\substack{\uparrow \\ \text{regulär}}}{L_{n-1}} \cdot \underset{\uparrow}{U_{n-1}} \quad (\text{LU-Zerlegung}) \quad \xi = \alpha - \underline{x}^T \underline{y}$$

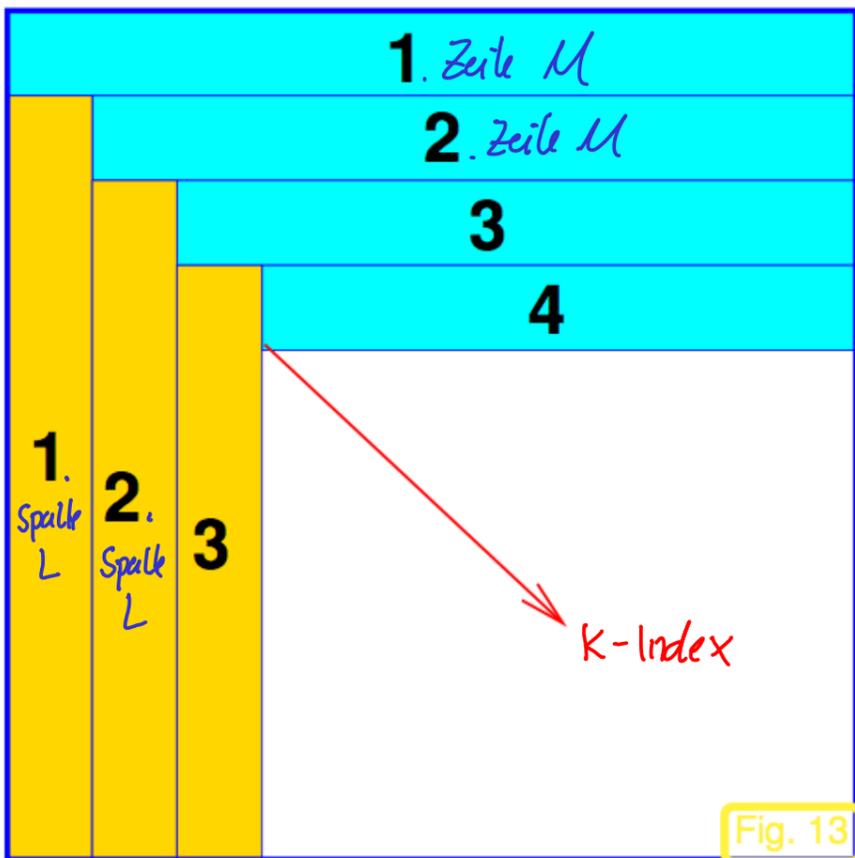
□

▷ Rekursiver Algorithmus zur Berechnung der LU-Zerlegung

Praxis : Parkettierungsalgorithmus



$$\mathbf{LU} = \mathbf{A} \Rightarrow a_{ik} = \sum_{j=1}^{\min\{i,k\}} l_{ij} u_{jk} = \begin{cases} \sum_{j=1}^{i-1} l_{ij} u_{jk} + 1 \cdot u_{ik} & , \text{ if } i \leq k , \\ \sum_{j=1}^{k-1} l_{ij} u_{jk} + l_{ik} u_{kk} & , \text{ if } i > k . \end{cases}$$



```

1 function [L,U] = lufak(A)
2 % Algorithm of Crout: LU-factorization of A : Instabil !
3 n = size(A,1); if (size(A,2) ~= n), error('n ~= m'); end
4 L = eye(n); U = zeros(n,n);
5 for k=1:n
6 % Compute row of U
7 for j=k:n, U(k,j) = A(k,j) - L(k,1:k-1)*U(1:k-1,j); end
8 % Compute column of L
9 for i=k+1:n, L(i,k) = (A(i,k) - L(i,1:k-1)*U(1:k-1,k)) /U(k,k); end
10 end
11

```

↓
= 0 ist möglich
↔ Nullpivot in GE

Zusammenhang mit GE :

```

1 function x = GENopivot(A,b)
2 % Raw Gaussian elimination without pivoting for
3 % linear system Ax=b: "Linear algebra version", UNSTABLE!
4 [m,n] = size(A); A = [A,b];
5 % Forward elimination
6 for i=1:n-1
7 if (A(i,i) == 0.0), error('Zero pivot'); end
8 for j=i+1:n
9 fac = A(j,i)/A(i,i);
10 for k=i+1:n+1
11 A(j,k) = A(j,k) - fac*A(i,k);
12 end
13 end
14 end
15 % Backward substitution
16 x = [zeros(n-1,1);A(n,n+1)/A(n,n)];
17 for i=n-1:-1:1
18 x(i) = A(i,n+1);
19 for k=i+1:n
20 x(i) = x(i) - A(i,k)*x(k);
21 end
22 x(i) = x(i)/A(i,i);
23 end

```

```

1 function [L,U] = gelu(A)
2 % LU-factorization based on Gaussian elimination
3 % No pivoting: UNSTABLE!
4 [m,n] = size(A); L = eye(n);
5 % Forward elimination
6 for i=1:n-1
7 if (A(i,i) == 0.0), error('Zero pivot'); end
8 for j=i+1:n
9 fac = A(j,i)/A(i,i); L(j,i) = fac;
10 for k=i+1:n
11 A(j,k) = A(j,k) - fac*A(i,k);
12 end
13 end
14 end
15 U = triu(A);

```

|| numerisch äquivalent

```

1 function [L,U] = lufak(A)
2 % Algorithm of Crout: LU-factorization of A
3 n = size(A,1); if (size(A,2) ~= n), error('n ~= m'); end
4 L = eye(n); U = zeros(n,n);
5 for k=1:n
6 % Compute row of U
7 for j=k:n, U(k,j) = A(k,j) - L(k,1:k-1)*U(1:k-1,j); end
8 % Compute column of L
9 for i=k+1:n, L(i,k) = (A(i,k) - L(i,1:k-1)*U(1:k-1,k)) /U(k,k); end
10 end
11

```

Rechenaufwand für 'Lufak':

$$\begin{aligned} \#flops &= \sum_{k=1}^n \left(\sum_{j=k}^n \underbrace{\{(k-1) + (k-1)\}}_{\substack{\uparrow \\ \text{MUL} \quad \text{ADD}}} + \sum_{i=k+1}^n \{1 + 2(k-1)\} \right) = \sum_{k=1}^n (n-k+1)2(k-1) + (n-k)(1+2(k-1)) = \\ &= \sum_{k=1}^n 2k^2 + \dots = 2 \cdot \frac{1}{6} (2n^3 + 3n^2 + n) + \dots \\ &= \frac{2}{3} n^3 + O(n^2) \quad \text{für } n \rightarrow \infty \\ &\quad \uparrow \\ &\quad 3 \text{ geschachtelte Schleifen} \end{aligned}$$

2.3. Rundungsfehleranalyse der Gausselimination

Algorithmus: Lufak + Vorwärtstechnung + Rücksubstitution ('backwardsubst')

Theorem 2.15. (Rückwärtsfehler beim Lösen eines LGS mit Gausselimination)

Sei die LU-Zerlegung mit `lufak` in Gleitkommaarithmetik \mathbb{F} durchführbar*. Dann erfüllt die in \mathbb{F} berechnete Lösung \hat{x} des linearen Gleichungssystems $Ax = b$

$$(A + \Delta A)\tilde{x} = b, \quad \text{für ein } \Delta A \text{ mit } |\Delta A| \leq (3\gamma_n + \gamma_n^2)|\hat{L}||\hat{U}|,$$

wobei \hat{L} und \hat{U} die von Algorithmus `lufak` in \mathbb{F} berechneten Faktoren der LU-Zerlegung sind und $\gamma_n(\mathbb{F}) := \frac{nu(\mathbb{F})}{1-nu(\mathbb{F})}$ (wie in Lemma 1.15).

* keine Division durch 0

▷ stabil, falls $|\hat{L}||\hat{U}| \approx |A|$

Beweis: ① LU-Zerlegung

```

1 function [L,U] = lufak(A)
2 % Algorithm of Crout: LU-factorization of A
3 n = size(A,1); if (size(A,2) ~= n), error('n ~= m'); end
4 L = eye(n); U = zeros(n,n);
5 for k=1:n
6     % Compute row of U
7     for j=k:n, U(k,j) = A(k,j) - L(k,1:k-1)*U(1:k-1,j); end → Skalarprodukt (*)
8     % Compute column of L
9     for i=k+1:n, L(i,k) = (A(i,k) - L(i,1:k-1)*U(1:k-1,k)) / U(k,k); end (**)
10 end [ (1+θ₂) U(k,k)L(i,k) = A(i,k) - L(i,1:k-1)*U(1:k-1,k) ]
11

```

Technik wie bei R.F.A. für Euklidisches Skalarprodukt

$$(*) \quad \underbrace{\left| 1 \cdot \tilde{u}_{kj} + \sum_{i=1}^{k-1} \tilde{l}_{ki} \tilde{u}_{ij} - a_{kj} \right|}_{\text{Rechenfehler im Skalarprodukt}} \stackrel{(*)}{\leq} \gamma_k(\mathbb{F}) \sum_{i=1}^k |\tilde{l}_{ki}| |\tilde{u}_{ij}|$$

Erinnerung: Abschnitt 1.2

$$s_n = \sum_{j=1}^n x_j y_j$$

$$\Rightarrow |\hat{s}_n - s_n| \leq \gamma_n(\mathbb{F}) \sum_{j=1}^n |x_j| |y_j| \quad (*)$$

$$(**) \quad \left| \sum_{j=1}^k \tilde{l}_{ij} \tilde{u}_{jk} - a_{ik} \right| \stackrel{(*)}{\leq} \gamma_k(\mathbb{F}) \sum_{j=1}^k |\tilde{l}_{ij}| |\tilde{u}_{jk}|$$

$$\Downarrow \quad \left| \tilde{L} \tilde{U} - A \right| \leq \gamma_n(\mathbb{F}) |\tilde{L}| |\tilde{U}| \Rightarrow \tilde{L} \tilde{U} = A + \Delta \bar{A}, \quad |\Delta \bar{A}| \leq \gamma_n(\mathbb{F}) |\tilde{U}| |\tilde{L}|$$

② Vorwärtsrechnung: $\tilde{L} \tilde{x} = \underline{b}$

```

1 function x = forwardsubst(L,b)
2 % Solution of lower triangular linear system of equations
3 % by means of forward substitution
4 [n,m] = size(L); x = zeros(n,1);
5 for i=1:n
6     for j=1:i-1, x(i) = x(i) + L(i,j)*x(j); end → Skalarprodukt
7     x(i) = (b(i) - x(i))/L(i,i);
8 end
    
```

Anwendung von (*) auf

$$\left| \sum_{j=1}^{i-1} \tilde{l}_{ij} \tilde{x}_j + 1 \cdot \tilde{x}_i - b_i \right| \leq \gamma_i(\mathbb{F}) \sum_{j=1}^i |\tilde{l}_{ij}| |\tilde{x}_j| \quad [\tilde{l}_{ii} = 1]$$

$$\tilde{x}_i = b_i - \sum_{j=1}^{i-1} \tilde{l}_{ij} \tilde{x}_j \quad \triangleright \quad \text{Aus der Rückwärtsfehleranalyse des Skalarprodukts}$$

↑
Skalarprodukt

$$\tilde{x}_i = b_i - \sum \hat{\tilde{l}}_{ij} \tilde{x}_j \quad \text{mit} \quad |\hat{\tilde{l}}_{ij} - \tilde{l}_{ij}| \leq \gamma_i(\mathbb{F}) |\tilde{l}_{ij}|$$

$$\triangleright \quad \text{Berechnete Lsg. } \tilde{\tilde{x}} \text{ erfüllt: } (\tilde{L} + \Delta L) \tilde{\tilde{x}} = \underline{b} \quad \text{mit} \quad |\Delta L| \leq \gamma_n(\mathbb{F}) |\tilde{L}|$$

③ Rückwärtsrechnung: $\tilde{U} \underline{x} = \tilde{\tilde{y}}$ → Berechnetes $\hat{\tilde{x}}$ erfüllt: $(\tilde{U} + \Delta U) \hat{\tilde{x}} = \tilde{\tilde{y}}$
mit $|\Delta U| \leq \gamma_n(\mathbb{F}) |\tilde{U}|$

$$\Rightarrow (\hat{L} + \Delta L) (\hat{U} + \Delta U) \hat{\tilde{x}} = \underline{b}$$

$$\Leftrightarrow \underbrace{(A + \Delta \bar{A} + \Delta L \cdot \tilde{U} + \hat{L} \cdot \Delta U + \Delta L \cdot \Delta U)}_{= \text{Störung } \Delta A} \hat{\tilde{x}} = \underline{b}$$

$$|\Delta A| \leq |\Delta \tilde{A}| + |\Delta L| |\tilde{M}| + |\hat{L}| |\Delta M| + |\Delta L| \cdot |\Delta M| \leq \gamma_n(\mathbb{F}) |\hat{L}| |\tilde{U}| + \gamma_n |\hat{L}| |\tilde{M}| + \gamma_n |\tilde{U}| |\hat{L}| + \gamma_n^2 |\tilde{M}| |\hat{L}| \quad \square$$

Normbasierte Formulierung:

Satz 2.15 \Rightarrow $\max_i \sum_j |(\Delta A)_{ij}| \leq (3\gamma_n + \gamma_n^2) \max_i \sum_j (|\hat{L}| |\tilde{M}|)_{ij}$

$$\|\Delta A\|_\infty \leq (3\gamma_n + \gamma_n^2) \|\hat{L}\|_\infty \|\tilde{M}\|_\infty$$

In $\|\cdot\|_2$: Verwende $\frac{1}{\sqrt{n}} \|M\|_\infty \leq \|M\|_2 \leq \sqrt{n} \|M\|_\infty$

$$\Rightarrow \|\Delta A\|_2 \leq \underbrace{n(3\gamma_n + \gamma_n^2)}_{\sim n^2 \mu(\mathbb{F})} \|\hat{L}\|_2 \|\tilde{M}\|_2$$

$$\begin{aligned} \|v\|_2^2 &= \sum_{i=1}^n |v_i|^2 \\ &\leq n \max_i |v_i|^2 = n \|v\|_\infty^2 \\ \|v\|_\infty^2 &= \max_i |v_i|^2 \\ &\leq \sum_{i=1}^n |v_i|^2 = \|v\|_2^2 \\ \|M\|_2 &= \sup_{z \neq 0} \frac{\|Mz\|_2}{\|z\|_2} \leq \\ &\leq \sup_{z \neq 0} \frac{\sqrt{n} \|Mz\|_\infty}{\|z\|_\infty} = \sqrt{n} \|M\|_\infty \end{aligned}$$

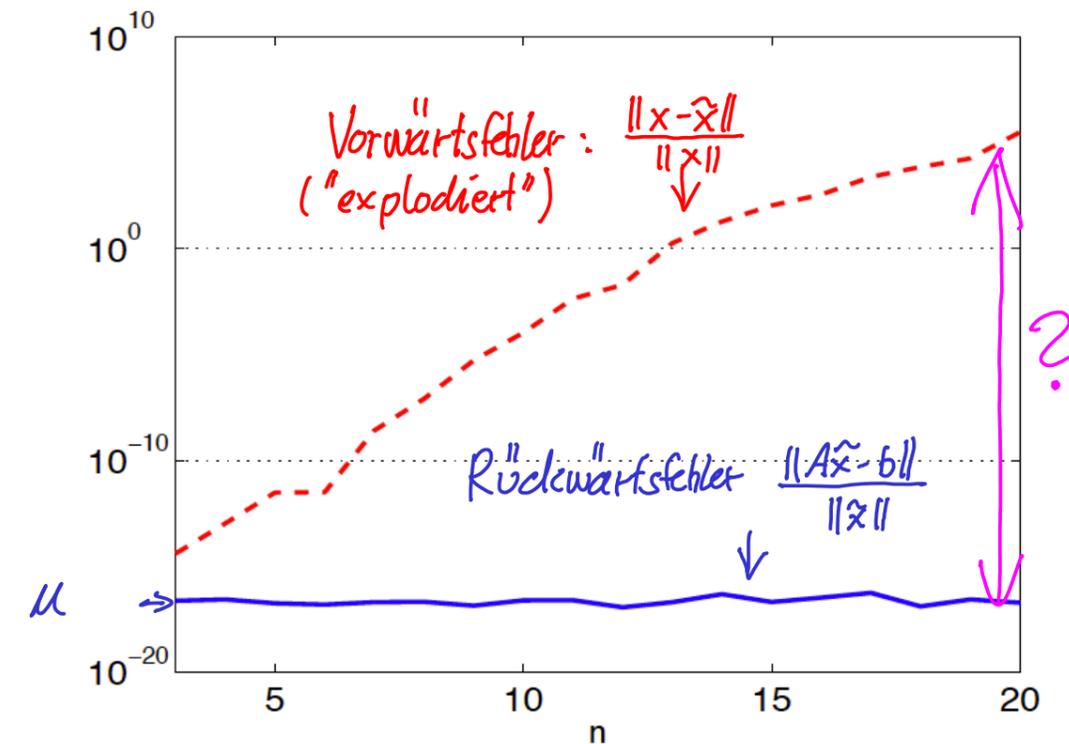
2.4. Sensitivität und Kondition

Bsp: Hilbertmatrix $A \in \mathbb{R}^{n,n}$, $a_{ij} := \frac{1}{i+j-1}$, $1 \leq i, j \leq n$

```

1- digits(100); rw = []; vw = [];
2- for n = 3:20,
3-   A = hilb(n); b = ones(n,1);
4-   x = A\b; % Solve LSE with std precision -> Gaußelimination
5-   rw = [rw;n,norm(A*x-b)/norm(x)];
6-   xe = vpa(A)\vpa(b); % high precision solve -> "exakte Lösung"
7-   vw = [vw;n,norm(double(x-xe))/norm(double(xe))]; -> relativer Fehler
8- end
9- figure('name','Gaussian elimination for Hilbert matrix');
10- plot(rw(:,1),rw(:,2),'b-',vw(:,1),vw(:,2),'r--');
11- xlabel('n','fontsize',14);

```



Warum betrachten wir $\frac{\|A\hat{x} - b\|}{\|\hat{x}\|}$?

Theorem 2.17. (Schranke für Rückwärtsfehler für Lösung eines LGS)

Sei $\hat{x} \neq 0$ die berechnete Lösung des linearen Gleichungssystems $Ax = b$. Dann gilt

$$\min \{ \|\Delta A\|_2 : (A + \Delta A)\hat{x} = b \} = \frac{\|r\|_2}{\|\hat{x}\|_2}, \quad \triangleleft \text{Normen verfügbarer Vektoren}$$

wobei $r = b - A\hat{x}$ das **Residuum** zu \hat{x} ist.

↑
"Rückwärtsfehler"

Beweis: (i) $\Delta A \cdot \hat{x} = r \Rightarrow \|\Delta A\|_2 \|\hat{x}\|_2 \geq \|\Delta A \cdot \hat{x}\|_2 = \|r\|_2 \Rightarrow \|\Delta A\|_2 \geq \frac{\|r\|_2}{\|\hat{x}\|_2}$

(ii) Mit $\Delta A = \frac{1}{\|\hat{x}\|_2^2} r \hat{x}^T \in \mathbb{C}^{n \times n}$: $(A + \Delta A)\hat{x} = (A + \frac{1}{\|\hat{x}\|_2^2} r \hat{x}^T)\hat{x} = A\hat{x} + r = b$

$$\|\Delta A\|_2 = \frac{1}{\|\hat{x}\|_2^2} \|r \hat{x}^T\|_2 = \frac{1}{\|\hat{x}\|_2^2} \|r\|_2 \|\hat{x}\|_2 = \frac{\|r\|_2}{\|\hat{x}\|_2}$$

□

2.4.1. Kondition einer Abbildung

$X, Y \cong$ normierte Vektorräume, **Problemabbildung**: $F : \mathcal{D}(F) \subset X \rightarrow Y$
↑ ↑
Datenraum Ergebnisraum

Wie wirkt sich eine (infinitesimal) kleine Störung von $x \in \mathcal{D}(F)$ auf $F(x)$ aus?

Def. [Absolute Kondition(szahl)]

Für eine Abbildung $F : \mathcal{D}(F) \subset X \rightarrow Y$ zwischen normierten Vektorräumen X und Y ist die **absolute Kondition(szahl)** von $F(x)$ für $x \in \mathcal{D}(F)$ gegeben durch

$$\text{cond}_{\text{abs}}(F(x)) = \limsup_{\delta \rightarrow 0} \left\{ \frac{\|F(x + \Delta x) - F(x)\|_Y}{\|\Delta x\|_X} : 0 < \|\Delta x\|_X < \delta, x + \Delta x \in \mathcal{D}(F) \right\}.$$

▷ cond_{abs} abhängig von $\|\cdot\|_X, \|\cdot\|_Y$!

→ misst Verstärkung kleiner absoluter Datenfehler

Differenzielle Konditionsanalyse: F 2x stetig diff.-bar

$$X = \mathbb{R}^m, Y = \mathbb{R}^n : F(\underline{x} + \Delta \underline{x}) = F(\underline{x}) + \mathcal{D}F(\underline{x}) \Delta \underline{x} + O(\|\Delta \underline{x}\|^2) \quad [\text{Mehrdim. Taylorformel}]$$

↑ Jacobimatrix $(\frac{\partial F_i}{\partial x_j})_{\substack{i=1, \dots, n \\ j=1, \dots, m}}$

$$\triangleright \text{cond}_{\text{abs}}(F(\underline{x})) = \lim_{\|\Delta \underline{x}\| \rightarrow 0} \frac{\|F(\underline{x}) + \mathcal{D}F(\underline{x}) \Delta \underline{x} - F(\underline{x})\| + O(\|\Delta \underline{x}\|^2)}{\|\Delta \underline{x}\|} \leq \|\mathcal{D}F(\underline{x})\|$$

Def. 2.19 [Relative Kondition(szahl)]

Für eine Abbildung $F : \mathcal{D}(F) \subset X \rightarrow Y$ zwischen normierten Vektorräumen X und Y ist die **relative Kondition(szahl)** von $F(x)$ für $x \in \mathcal{D}(F)$ mit $F(x) \neq 0$ gegeben durch

$$\text{cond}_{\text{rel}}(F(x)) = \lim_{\delta \rightarrow 0} \sup \left\{ \frac{\|F(x + \Delta x) - F(x)\|_Y}{\|F(x)\|_Y} \frac{\|x\|_X}{\|\Delta x\|_X} : 0 < \|\Delta x\|_X < \delta, x + \Delta x \in \mathcal{D}(F) \right\} \rightarrow \text{misst Verstärkung kleiner relativer Datenfehler}$$

Falls F 2x stetig diff.-bar :

$$\text{cond}_{\text{rel}}(F(x)) = \lim_{\Delta x \rightarrow 0} \frac{\|\mathcal{D}F(x) \Delta x\|_Y + O(\|\Delta x\|^2)}{\|\Delta x\|_X} \cdot \frac{\|x\|_X}{\|F(x)\|_Y} \leq \|\mathcal{D}F(x)\| \frac{\|x\|_X}{\|F(x)\|_Y}$$

Lemma. [Relative Konditionszahl differenzierbarer Funktionen]

Ist $F : \mathcal{D}(F) \subset X \rightarrow Y$ stetig differenzierbar in einem inneren Punkt $x \in \mathcal{D}(F)$ mit $F(x) \neq 0$, so gilt

$$\text{cond}_{\text{rel}}(F(x)) = \frac{\|\mathcal{D}F(x)\|_{\mathcal{L}(X,Y)}}{\|F(x)\|_Y} \|x\|_X \cdot$$

Bsp: Rel. Kondition der Subtraktion: $F(a,b) = a - b$, $X = \mathbb{R}^2$, $Y = \mathbb{R}$, Eukl. Norm

$$DF(a,b) = \left(\frac{\partial F}{\partial a}(a,b), \frac{\partial F}{\partial b}(a,b) \right) = (1, -1) \Rightarrow \|DF(a,b)\|_2 = \|(1, -1)^T\|_2 = \sqrt{2}$$

$$\text{cond}_{\text{rel}}(F(a,b)) = \sqrt{2} \frac{\sqrt{a^2 + b^2}}{|a - b|} \rightarrow \infty \text{ für } b \rightarrow a \text{ vgl. Auslöschung, Sect. 1.5}$$

Kondition - Vorwärtsfehler - Rückwärtsfehler

$$\frac{\|F(x+\Delta x) - F(x)\|}{\|F(x)\|} = \frac{\|\tilde{F}(x) - F(x)\|}{\|F(x)\|}$$

Für $\|\Delta x\| \approx 0$

$$\approx \text{cond}_{\text{rel}}(F(x)) \frac{\|\Delta x\|}{\|x\|}$$

$$\tilde{F}(x) = F(x + \Delta x)$$

↳ durch Rundungsfehler gestörte Abbildung, realisiert durch Algorithmus

↳ relativer Vorwärtsfehler unter Rundungsfehlerinfluss

↳ relativer Rückwärtsfehler

▷ Falls $\text{cond}_{\text{rel}}(F(x)) \gg 1 \Rightarrow$ winzige Rückwärtsfehler trotz riesiger Vorwärtsfehler
 ↳ numerisch gutes Ergebnis
 "numerisch harmlos"

2.4.2. Kondition einer Matrix

Bsp: Kondition Matrix x Vektor $F(x) = Ax$, $A \in \mathbb{R}^{n \times n}$ regulär, $X = \mathbb{R}^n$, $Y = \mathbb{R}^n$

$$DF(x) = A \stackrel{[\text{Def 2.19}]}{\Rightarrow} \text{cond}_{\text{rel}}(F(x)) = \|A\| \frac{\|x\|}{\|Ax\|} = \|A\| \frac{\|A^{-1}Ax\|}{\|Ax\|} \leq \|A\| \|A^{-1}\|, \|\cdot\| \stackrel{1}{=} \text{Vektor- oder assoz. Matrixnorm}$$

[$x \neq 0$]

Bsp: Matrixinversion: $F(X) = X^{-1}$, $X \in GL(n) := \{M \in \mathbb{R}^{n \times n}, M \text{ invertierbar}\} \subset \mathbb{R}^{n \times n}$ mit $\|\cdot\|$

Hilfsmittel: **Neumannsche Reihe** ("geometrische Reihe für Matrizen"): $M \in \mathbb{R}^{n \times n}$: $\|M\| < 1 \Rightarrow (I - M)^{-1} = \sum_{j=0}^{\infty} M^j$

$$(X + \Delta X)^{-1} - X^{-1} = X^{-1} [(I - \Delta X \cdot X^{-1})^{-1} - I] \stackrel{N.R.}{=} X^{-1} \sum_{j=1}^{\infty} (-\Delta X \cdot X^{-1})^j, \text{ falls } \|\Delta X \cdot X^{-1}\| < 1$$

$$= -X^{-1} \cdot \Delta X \cdot X^{-1} + O(\|\Delta X\|^2) \quad (*)$$

$$\text{cond}_{\text{rel}}(F(X)) = \lim_{\Delta X \rightarrow 0} \frac{\| -X^{-1} \cdot \Delta X \cdot X^{-1} \| + O(\|\Delta X\|^2)}{\|X^{-1}\|} \cdot \frac{\|X\|}{\|\Delta X\|} \leq \|X^{-1}\| \|X\|$$

\uparrow
 $\|\cdot\|$ submultiplikativ

Nichtasymptotische Analyse:

Proposition 2.21. [Relative Kondition der Matrixinversion]

Sei $A \in \mathbb{R}^{n \times n}$ regulär und $\Delta A \in \mathbb{R}^{n \times n}$ Störung von A mit $\|A^{-1} \Delta A\| < 1$ für eine submultiplikative Matrixnorm $\|\cdot\|$. Dann gilt

$$\underbrace{\frac{\|(A + \Delta A)^{-1} - A^{-1}\|}{\|A^{-1}\|}}_{\text{Rel. Fehler in } A^{-1}} \leq \|A\| \|A^{-1}\| \underbrace{\frac{\|\Delta A\|}{\|A\|}}_{\text{Rel. Fehler in } A} (1 + O(\|\Delta A\|)).$$

Beweis: (i) N.R. (wie oben): $A + \Delta A$ invertierbar: $A + \Delta A = A (I + A^{-1} \Delta A)$

$$(ii) \|(A + \Delta A)^{-1} - A^{-1}\| = \|-A^{-1} \cdot \Delta A (A + \Delta A)^{-1}\| \leq \|A^{-1} \cdot \Delta A\| \|(A + \Delta A)^{-1} - A^{-1} + A^{-1}\| \leq \|A^{-1} \Delta A\| (\|(A + \Delta A)^{-1} - A^{-1}\| + \|A^{-1}\|)$$

< 1

$\in GL(n)$ nach Vorausss.
 [wird gezeigt mit N.R.]

$$\Rightarrow \underbrace{(1 - \|A^{-1} \cdot \Delta A\|)}_{1 + O(\|\Delta A\|)} \|(A + \Delta A)^{-1} - A^{-1}\| \leq \|A^{-1}\|^2 \|\Delta A\| = \|A^{-1}\| (\|A\| \|A^{-1}\|) \frac{\|\Delta A\|}{\|A\|}$$

□

Def. [Matrixkondition]Für eine Matrixnorm $\|\cdot\|$ und eine invertierbare Matrix $A \in \mathbb{C}^{n \times n}$ heisst

$$\kappa(A) := \|A\| \|A^{-1}\|$$

die **Kondition** der Matrix A .→ Abhängig von der Matrixnorm $\|\cdot\|$ MATLAB: $\text{cond}(A)$

2.4.3 Kondition der Lösung von LGS

$$X = GL(n) \times \mathbb{R}^n \subset \mathbb{R}^{n \times n} \times \mathbb{R}^n, Y = \mathbb{R}^n : F(A, \underline{b}) = A^{-1}b$$

Sensitivität bzgl. Störung in \underline{b} : $\text{cond}_{\text{rel}}(F(\cdot, \underline{b})) = \|A^{-1}\| \|A\| = \kappa(A)$

Sensitivität (asymptotisch) bzgl. Störung in A : $\frac{\|(A + \Delta A)^{-1}b - A^{-1}b\|}{\|A^{-1}b\|} \frac{\|A\|}{\|\Delta A\|} = \frac{\|A^{-1} \cdot \Delta A \cdot A^{-1}b\| + O(\|\Delta A\|_2^2)}{\|A^{-1}b\|} \frac{\|A\|}{\|\Delta A\|}$
 $= \|A^{-1}\| \|A\| + O(\|\Delta A\|) = \kappa(A) + O(\|\Delta A\|)$

Nichtasymptotische Analyse:

Satz 2.22 [Sensitivität für lineare Gleichungssysteme]Sei $A \in \mathbb{C}^{n \times n}$ invertierbar und sei $\Delta A \in \mathbb{C}^{n \times n}$ so dass

$$\|A^{-1} \Delta A\| < 1$$

für eine submultiplikative Matrixnorm $\|\cdot\|$. Dann ist

$$(A + \Delta A)\hat{x} = \underline{b} + \Delta \underline{b}$$

eindeutig lösbar und es gilt mit $\underline{x} = A^{-1}\underline{b}$

$$\frac{\|\underline{x} - \hat{x}\|}{\|\underline{x}\|} \leq \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1} \Delta A\|} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta \underline{b}\|}{\|\underline{b}\|} \right).$$

Beweis: $\left. \begin{aligned} Ax &= \underline{b} \\ (A + \Delta A)\hat{x} &= \underline{b} + \Delta \underline{b} \end{aligned} \right\} \Rightarrow$

$$\underline{x} - \hat{x} = (A + \Delta A)^{-1} (\Delta A \underline{x} - \Delta \underline{b})$$

Mit Störungslemma: $\|(A + \Delta A)^{-1}\| \leq \|A^{-1}\| \frac{1}{1 - \|A^{-1} \Delta A\|}$

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

$$\Rightarrow \frac{\|\underline{x} - \hat{x}\|}{\|\underline{x}\|} = \frac{\|A^{-1}\| \|A\|}{1 - \|A^{-1} \Delta A\|} \left(\frac{\|\Delta A\|}{\|A\|} + \frac{\|\Delta \underline{b}\|}{\|A\| \|\underline{x}\|} \right)$$

□

Hilfsmittel :

Lemma [Störungslemma für lineare Abbildungen im $\mathbb{C}^{n \times n}$]

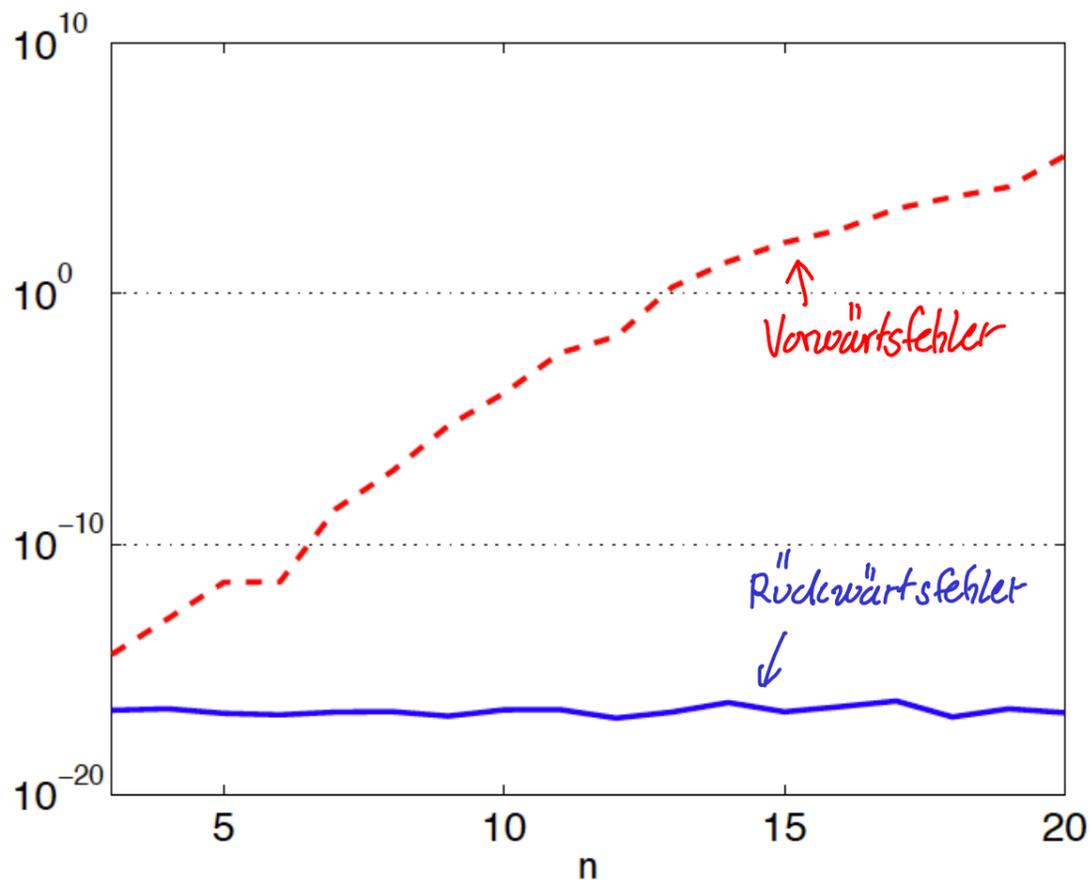
$$B \in \mathbb{C}^{n \times n}, \|B\| < 1 \Rightarrow I + B \text{ regulär}, \|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|}.$$

Beweis: (i) **N.R.** : $(I + B)^{-1} = \sum_{j=0}^{\infty} (-B)^j$, da $\|B\| < 1 \Rightarrow \|(I + B)^{-1}\| \stackrel{\Delta\text{-Ungl.}}{\leq} \sum_{j=0}^{\infty} \|B^j\| \stackrel{\text{Submult.}}{\leq} \sum_{j=0}^{\infty} \|B\|^j = \frac{1}{1 - \|B\|}$

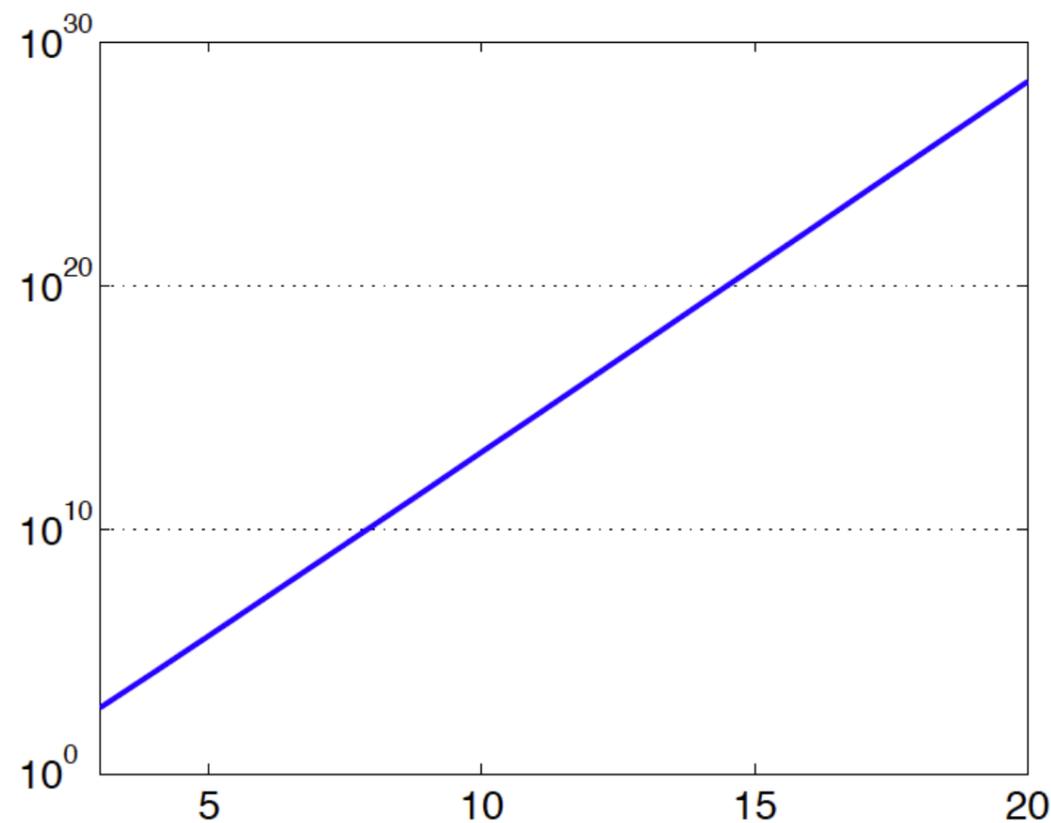
(ii) $v \in \mathbb{C}^n$: $\|(I + B)v\| \geq \|v\| - \|Bv\| \geq (1 - \|B\|)\|v\|$

$\|v\| = \|(I + B)(I + B)^{-1}v\| \stackrel{\downarrow}{\geq} \underbrace{(1 - \|B\|)}_{> 0} \|(I + B)^{-1}v\| \Rightarrow \frac{\|(I + B)^{-1}v\|}{\|v\|} \leq \frac{1}{1 - \|B\|} \quad \square$

Zum Beispiel Hilbertmatrix :



~



2.5. Gausselimination mit Pivotsuche

```

1 % Instability of Gaussian
2 % elimination w/o pivoting
3 A = [2^(-55), 1; 1, 1];
4 b = [1; 2];
5 x1 = A\b;
6 x2 = GENpivot(A,b);
7 [L,U] = lufak(A);
8 z = L\b; x3 = U\z;
9 format short; disp([x1,x2,x3]);

```

$$\triangle Ax = \underline{b}, \quad A = \begin{pmatrix} \varepsilon & 1 \\ 1 & 1 \end{pmatrix}, \quad \underline{b} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

$$\underline{x} = \begin{pmatrix} \frac{1}{1-\varepsilon} \\ \frac{1-2\varepsilon}{1-\varepsilon} \end{pmatrix} \approx \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \underline{\tilde{x}} \quad \text{für } \varepsilon \ll 1$$

$$\begin{pmatrix} 1 \\ 2 \end{pmatrix} - A \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} -\varepsilon \\ 0 \end{pmatrix} \stackrel{!}{=} \text{Residuum } \underline{r}$$

$$\frac{\|\underline{r}\|}{\|\underline{\tilde{x}}\|} \approx \varepsilon \rightarrow \text{numerisch stabile Lösung, Satz 2.17}$$

```

1 function [L,U] = lufak(A)
2 % Algorithm of Crout: LU-factorization of A
3 n = size(A,1); if (size(A,2) ~= n), error('n ~= m'); end
4 L = eye(n); U = zeros(n,n);
5 for k=1:n
6 % Compute row of U
7 for j=k:n, U(k,j) = A(k,j) - L(k,1:k-1)*U(1:k-1,j); end
8 % Compute column of L
9 for i=k+1:n, L(i,k) = (A(i,k) - L(i,1:k-1)*U(1:k-1,k)) / U(k,k); end
10 end
11

```

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

$$= \tilde{L} \cdot \tilde{U} \quad \tilde{L} = \begin{pmatrix} \varepsilon & 1 \\ 0 & 1-\varepsilon^{-1} \end{pmatrix}$$

$$\tilde{L} \tilde{U} \tilde{\underline{x}} = \underline{b} \Rightarrow \tilde{\underline{x}} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \rightarrow \text{numerisch instabil} \quad \left. \vphantom{\tilde{\underline{x}}} \right\} \rightarrow \text{riesiger Rückwärtsfehler}$$

$$\|\tilde{L}\| \|\tilde{U}\| = \begin{pmatrix} \varepsilon & 1 \\ 1 & 2\varepsilon^{-1} \end{pmatrix} \Rightarrow \|A\|, \text{ vgl. 2.14}$$

$$A \approx A_0 = \begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \Rightarrow \text{keine zulässige Eingabe für 'lufak'}$$

Faustregel: Numerische Instabilität schlägt zu bei Eingabedaten nahe bei unzulässigen Eingabedaten

⇒ Idee: Mache 'LUFak' A_0 -tauglich → (LA) Zeilenvertauschung

$$A \rightsquigarrow A_p = \begin{pmatrix} 1 & 1 \\ \varepsilon & 1 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & 0 \\ \varepsilon & 1 \end{pmatrix}}_{\tilde{L}} \begin{pmatrix} 1 & 1 \\ 0 & 1-\varepsilon \end{pmatrix} \Rightarrow \tilde{U} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

$$\tilde{L} \tilde{U} \hat{x} = \underline{b} \Rightarrow \hat{x} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

"Stabilisierte Gaußelimination"

```
1 function x = GELApivot(A,b)
2 % Raw Gaussian elimination with "theoretical" pivoting for
3 % linear system Ax=b: "Linear algebra version", UNSTABLE!
4 [m,n] = size(A); A = [A,b];
5 % Forward elimination
6 for i=1:n-1
7     l = find(A(i:n,i) ~= 0); % search nonzero pivot
8     A([i,l],:) = A([l,i],:); % swap rows → Zeilenvertauschung
9     for j=i+1:n
10        fac = A(j,i)/A(i,i);
11        for k=i+1:n+1
12            A(j,k) = A(j,k) - fac*A(i,k);
13        end
14    end
15 end
16 % Backward substitution
17 x = [zeros(n-1,1);A(n,n+1)/A(n,n)];
18 for i=n-1:-1:1
19     x(i) = A(i,n+1);
20     for k=i+1:n
21         x(i) = x(i) - A(i,k)*x(k);
22     end
23     x(i) = x(i)/A(i,i);
24 end
```

← Numerisch unzulässig !

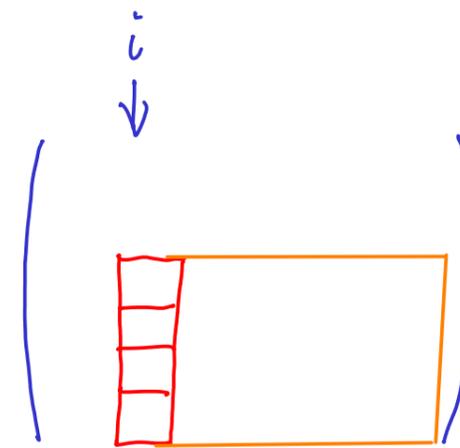
Numerisch:

Abfrage auf relative Grösse $\approx U(F)$

"Gaußelimination mit Spaltenpivotsuche"

```

1 function x = GEpivot(A,b)
2 % Raw Gaussian elimination without pivoting for
3 % linear system Ax=b: "Linear algebra version", UNSTABLE!
4 [m,n] = size(A); A = [A,b];
5 % Forward elimination
6 for i=1:n-1
7     % Column pivoting (Spaltenpivotsuche)
8     rowmax = max(abs(A(i:n,i:n)))';
9     [pvt,l] = max(abs(A(i:n,i))./rowmax);
10    if (abs(A(l,i)) < eps*max(rowmax)), disp('A nearly singular');
11    A([i,l],:) = A([l,i],:); % swap rows
12    for j=i+1:n
13        fac = A(j,i)/A(i,i);
14        for k=i+1:n+1
15            A(j,k) = A(j,k) - fac*A(i,k);
16        end
17    end
18 end
19 % Backward substitution
20 x = [zeros(n-1,1);A(n,n+1)/A(n,n)];
21 for i=n-1:-1:1
22     x(i) = A(i,n+1);
23     for k=i+1:n
24         x(i) = x(i) - A(i,k)*x(k);
25     end
26     x(i) = x(i)/A(i,i);
27 end
    
```



$l \in \{i, \dots, n\}$ so dass $\frac{|a_{eil}|}{\max_j |a_{ejl}|} \rightarrow \max$
 ↑
 relativ grösstes Pivotelement

→ wegen Skalierunginvarianz: gleiche numerische Lösung für $DA\underline{x} = D\underline{b}$ für beliebige reguläre Diagonalmatrix D .

'GEpivot' ist numerisch äquivalent zu

(i) Anwendung aller Zeilenvertauschungen auf A : $A_p = PA$
 ↳ Permutationsmatrix

(ii) LU-Zerlegung 'Lofak': $A_p = LM$

(iii) $x = M \setminus (L \setminus (P b))$

Satz 2.31 (Existenz der LU-Zerlegung mit Spaltenpivotsuche)

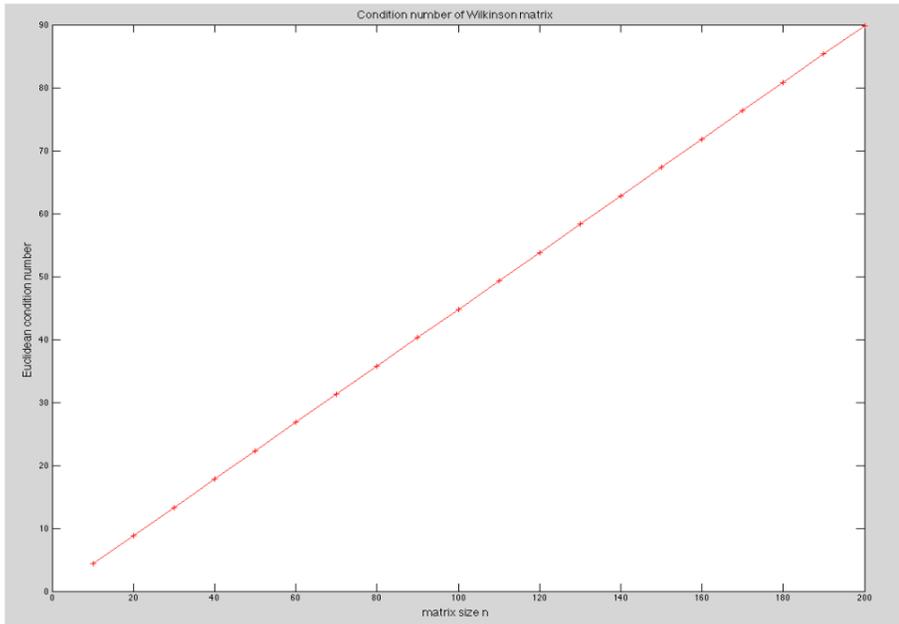
Zu jeder regulären Matrix $A \in \mathbb{C}^{n \times n}$ gibt es eine $n \times n$ Permutationsmatrix P so, dass PA eine LU-Zerlegung besitzt.

▷ GEpivot immer durchführbar (algebraisch)

Exakte LU-Zerlegung

$$A = LU, \quad l_{ij} = \begin{cases} 1 & , \text{if } i = j, \\ -1 & , \text{if } i > j, \\ 0 & \text{else} \end{cases}, \quad u_{ij} = \begin{cases} 1 & , \text{if } i = j, \\ 2^{i-1} & , \text{if } j = n, \\ 0 & \text{else.} \end{cases}$$

△ A gut konditioniert

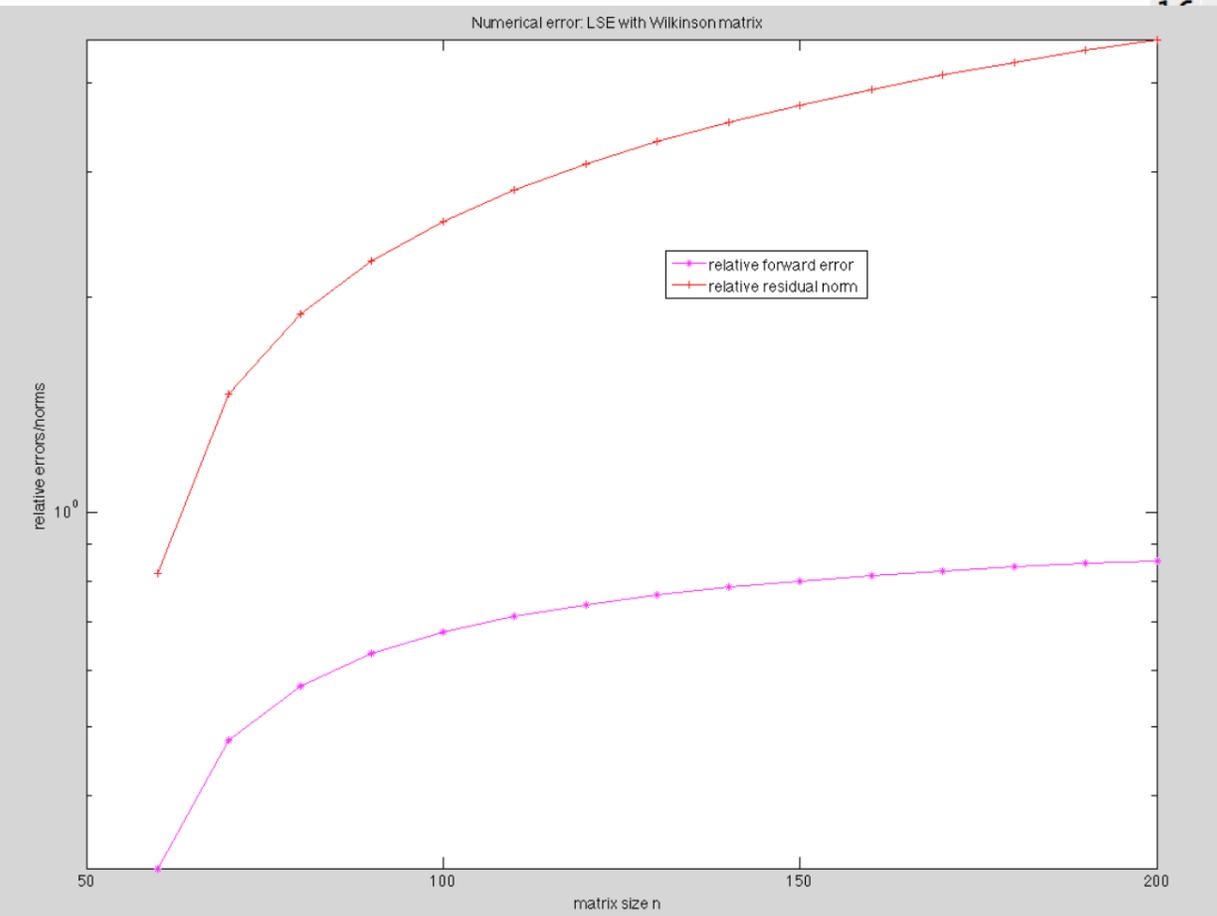


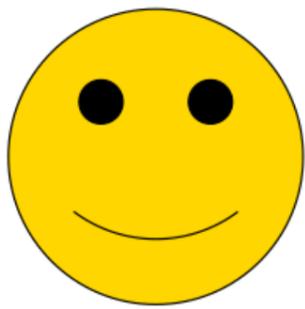
```

1 % Wilkinson's counterexample: instability of Gaussian elimination
2 % despite column pivoting
3 res = [];
4 for n=10:10:200
5     % Build Wilkinson matrix
6     A = [tril(-ones(n,n-1))+2*[eye(n-1);zeros(1,n-1)],ones(n,1)];
7     % imposed solution
8     x = ((-1).^(1:n))'; b = A*x;
9     % Numerical solution
10    xh = A\b; → G.F. mit Spaltenpivotssuche
11    % Residual
12    r = b-A*xh;
13    % Relative forward error
14    relerr = norm(xh-x)/norm(x); → Vorwärtsfehler
15    % Relative norm of residual
16    relresnorm = norm(r)/norm(xh); → Rückwärtsfehler
17    res = [res; n relerr relresnorm];
18 end
19 semilogy(res(:,1),res(:,2),'m-*',res(:,1),res(:,3),'r-+');
20 set(gca,'fontsize',14);
21 title('Numerical error: LSE with Wilkinson matrix');
22 xlabel('matrix size n','fontsize',14);
23 ylabel('relative errors/norms','fontsize',14);
24 legend('relative forward error','relative residual norm',...
25        'location','best');

```

△ Rückwärtsfehler / Vorwärtsfehler ≈ 1 (relativ)



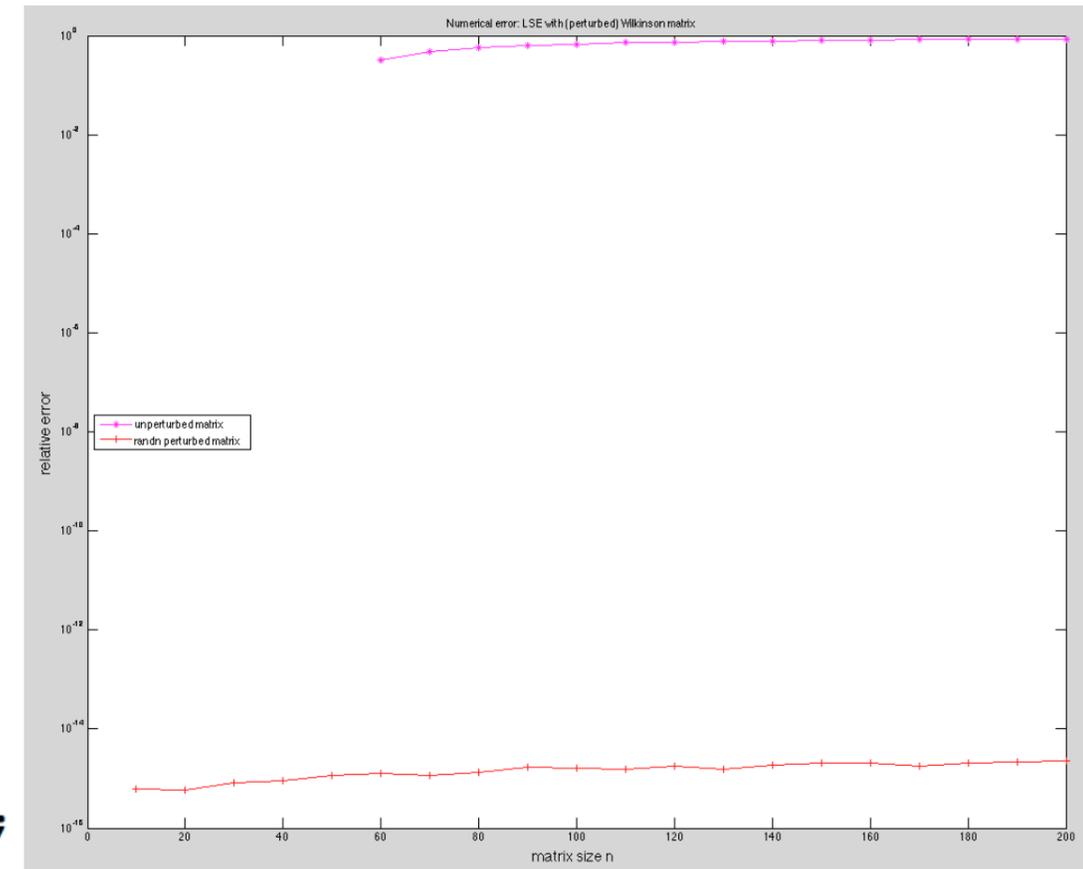


'GEpivot' funktioniert in der Praxis immer (= numerisch stabil)

```

1 % Curing Wilkinson's counterexample by random perturbation
2 % Theory: Spielman and Teng
3 res = [];
4 for n=10:10:200
5     % Build Wilkinson matrix
6     A = [tril(-ones(n,n-1))+2*[eye(n-1);zeros(1,n-1)],ones(n,1)];
7     % imposed solution
8     x = ((-1).^(1:n))';
9     relerr = norm(A\(A*x)-x)/norm(x);
10    % Randomly perturbed Wilkinson matrix by matrix with iid
11    % $N(0, {\mathrm{eps}})$ distributed entries
12    Ap = A + eps*randn(size(A));
13    relerrp = norm(Ap\(A*x)-x)/norm(x);
14    res = [res; n relerr relerrp];
15 end
16 semilogy(res(:,1),res(:,2),'m-*',res(:,1),res(:,3),'r-+');
17 title('Numerical error: LSE with (perturbed) Wilkinson matrix');
18 xlabel('matrix size n','fontsize',14);
19 ylabel('relative error','fontsize',14);
20 legend('unperturbed matrix','randn perturbed matrix','location','best');

```



Bem. (Algorithmische Relevanz der LU-Zerlegung)

ungeschickt !

```

1 % Setting: N >> 1, large matrix A ∈ ℝ^{n,n}
2 for j=1:N
3     x = A\b; → G.E.
4     b = some_function(x);
5 end

```

Rechenaufwand $O(Nn^3)$

schlau !

```

1 % Setting: N >> 1, large matrix A
2 [L,U] = lu(A);
3 for j=1:N
4     x = U\(L\b);
5     b = some_function(x);
6 end

```

Rechenaufwand $O(n^3 + Nn^2)$

MATLAB

$$[L,U] = lu(A)$$

↳ *permutierte* untere Δ -Matrix

$$[L,U,P] = lu(A)$$

↳ Permutationsmatrix

Bei mehrfachen rechten Seiten $\underline{b}_1, \dots, \underline{b}_m$: MATLAB: $X = A \setminus [b_1, b_2, \dots, b_m]$;

2.7. Cholesky-Zerlegung für SPD Matrizen

Def. 2.34. (Symmetrisch **positiv definite** Matrizen)

Sei $A \in \mathbb{R}^{n \times n}$. Eine Matrix $A \in \mathbb{R}^{n \times n}$ heißt

1. **positiv definit**, falls $\underline{x}^T A \underline{x} > 0 \quad \forall \underline{x} \in \mathbb{R}^n \setminus \{0\}$;

2. **positiv semi-definit**, falls $\underline{x}^T A \underline{x} \geq 0 \quad \forall \underline{x} \in \mathbb{R}^n$.

Eine symmetrische, positiv definite Matrix $A \in \mathbb{R}^{n \times n}$ heißt kurz **SPD**.

$$\rightarrow \begin{cases} \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R} \\ (\underline{x}, \underline{y}) \rightarrow \underline{x}^T A \underline{y} \end{cases} \text{ ist}$$

ein **Skalarprodukt**,
wenn $A = A^T$

Satz 2.35 (Eigenschaften von SPD Matrizen)

Sei $A \in \mathbb{R}^{n \times n}$ SPD. Dann gilt:

1. A ist regulär (invertierbar);

2. $a_{ii} > 0$ für alle $i \in \{1, \dots, n\}$;

3. $|a_{ij}| < \frac{1}{2}(a_{ii} + a_{jj})$ für $i \neq j$ und damit $\max_{ij} |a_{ij}| = \max_i a_{ii}$;

4. ist $X \in \mathbb{R}^{n \times n}$ invertierbar, so ist $X^T A X$ wieder SPD;

5. ist A eine Blockdiagonalmatrix $A = \text{diag}(A_1, A_2)$, so sind sowohl A_1 als auch A_2 SPD.

6. Alle EW von A sind positiv. $\Rightarrow \det(A) > 0$

zu 3.

$$\begin{pmatrix} * & * \\ & * \end{pmatrix} \begin{pmatrix} \square & \square \\ \square & \square \end{pmatrix} \begin{pmatrix} * \\ * \end{pmatrix}$$

$$\underline{x}^T (X^T A X) \underline{x} = (X \underline{x})^T A (X \underline{x})$$

Satz (LU-Zerlegung von SPD Matrizen)

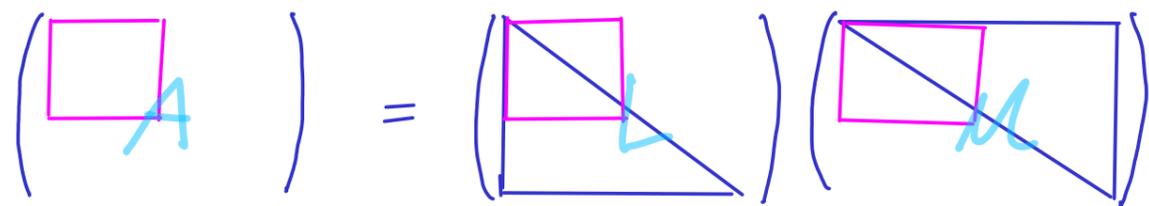
Jede SPD Matrix $A \in \mathbb{R}^{n \times n}$ besitzt eine LU-Zerlegung $A = LU$, wobei U positive Diagonalelemente hat.

Beweis:

$$(i) \begin{pmatrix} * \\ \vdots \\ * \\ 0 \end{pmatrix}^T \left(\begin{array}{c} \square \\ \square \\ \square \\ \square \end{array} \begin{array}{c} A \\ \\ \\ \end{array} \right) \begin{pmatrix} * \\ \vdots \\ * \\ 0 \end{pmatrix}$$

▷ Hauptminoren {SPD \Rightarrow regulär} : Satz 2.9.

(ii) $A = LM \Rightarrow 0 < \det(A) = \underbrace{\det(L)}_{=1} \det(M) = \prod_{i=1}^n (M)_{ii}$



i -te Hauptminoren von L/M liefern LU-Zerlegung des i -ten Hauptminors von $A \Rightarrow$ Induktion \square

Satz 2.36 (Cholesky-Zerlegung von SPD Matrizen)

Sei $A \in \mathbb{R}^{n \times n}$ SPD. Dann existiert eine eindeutige obere Dreiecksmatrix $U \in \mathbb{R}^{n \times n}$ mit positiven Diagonalelementen, so dass $A = U^T U$.

Beweis : $A = L\tilde{M} = A^T = \tilde{M}^T L^T = \hat{M}^T D L^T = \hat{M}^T (DL^T) \stackrel{!}{=} \text{wieder LU-Zerlegung}$

$[\tilde{M} = D\hat{M}, D = \text{diag}(\underbrace{m_{11}}_{>0}, \dots, \underbrace{m_{nn}}_{>0}) \Rightarrow (\hat{M})_{ii} = 1$

Eindeutigkeit LU-Zerl. $\Rightarrow L = \hat{M}^T \Rightarrow A = \hat{M}^T D \hat{M} = \underbrace{\hat{M}^T \sqrt{D}}_{U^T} \underbrace{\sqrt{D} \hat{M}}_U$ \square

Satz 2.39 (Gausselimination ohne Pivotsuche stabil für SPD Matrizen)

Sei $A \in \mathbb{R}^{n \times n}$ SPD und \hat{x} die mittels Cholesky-Zerlegung und Vorwärts-/Rückwärtseinsetzen berechnete Lösung von $Ax = b$. Dann gilt

$(A + \Delta A)\hat{x} = b, \quad \|\Delta A\|_2 \leq 4n^2(3n+1)u\|A\|_2.$

\uparrow
 \sim Rechenaufwand

Beweisskizze :

(i) Wie im Beweis zu Thm 2.19 :

$|\Delta A)_{ij}| \leq 4(3n+1)u(F) |M^T|_M$

\uparrow Cholesky-Faktor von A

$$(ii) \quad \Rightarrow \quad \| \Delta A \|_2 \leq 4n(3n+1) \nu(\mathbb{F}) \| |M|^T |M| \|_2 \leq 4n^2(3n+1) \nu(\mathbb{F}) \| A \|_2$$

$$\| |M|^T |M| \|_2 = \sup_{x \neq 0} \frac{x^T |M|^T |M| x}{\|x\|_2^2} = \| |M| \|_2^2 \leq \| |M| \|_F^2 = \| M \|_F^2 \leq n \| M \|_2^2 = n \| M^T M \|_2 = n \| A \|_2$$

$$[M = M^T \Rightarrow \| M \|_2 = \sup_{x \neq 0} \frac{x^T M x}{\|x\|_2^2} \text{ mit Spektralsatz}]$$

$$[\| M \|_F^2 = \sum_{i,j} (M)_{ij}^2 = \sum_j \sigma_j(M)^2 \leq n \sigma_{\max}(M)^2 = n \| M \|_2^2]$$

↑
Frobeniusnorm
(siehe Abschnitt 0.7)

↑
Singularwert
(siehe Abschnitt 0.10)

↑
Satz: $\| M \|_2 = \max_j \sigma_j(M)$

□

2.8. LU-Zerlegung

Def. 2.40 (Bandmatrix)

Eine Matrix $A \in \mathbb{R}^{n \times n}$ heisst **Bandmatrix** mit **Bandbreite** $p+q+1$, falls es $p, q \in \mathbb{N} \cup \{0\}$ gibt, so dass

$$a_{ij} = 0 \text{ f\u00fcr } j > i+p \text{ oder } i > j+q.$$

Die Zahl p heisst die **obere Bandbreite** und q die **untere Bandbreite**.

$p \stackrel{\text{!}}{=} \# \text{ nichtverschwindende Superdiagonalen}$

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1,p+1} & 0 & \dots & \dots & 0 \\ a_{21} & a_{22} & \dots & \dots & a_{2,p+2} & 0 & \dots & 0 \\ \vdots & \vdots & \dots & \dots & \dots & \dots & \dots & 0 \\ a_{q+1,1} & a_{q+1,2} & \dots & \dots & \dots & \dots & \dots & 0 \\ 0 & a_{q+2,2} & \dots & \dots & \dots & \dots & \dots & a_{n-p,n} \\ \vdots & 0 & \dots & \dots & \dots & \dots & \dots & \vdots \\ \vdots & \vdots & \dots & \dots & \dots & \dots & \dots & \vdots \\ 0 & 0 & \dots & 0 & a_{n,n-q} & \dots & a_{n,n-1} & a_{nn} \end{pmatrix}$$

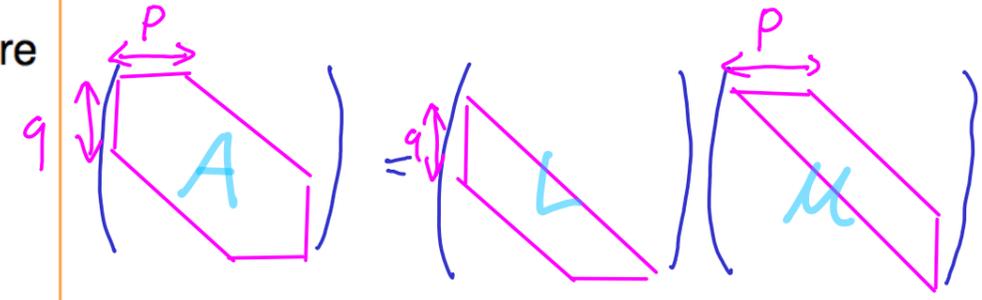
$q = \# \neq 0 \text{ Subdiagonalen}$

Satz 2.41 (Band-LU-Zerlegung)

Sei $A \in \mathbb{R}^{n \times n}$ eine Bandmatrix mit oberer Bandbreite p und unterer Bandbreite q und existiere die LU-Zerlegung $LU = A$. Dann haben L, U Bandstruktur, und es gilt:

$$l_{ij} = 0, \text{ falls } j > i \text{ oder } j < i - q,$$

$$u_{ij} = 0, \text{ falls } j < i \text{ oder } j > i + p.$$



Beweis: (Induktion, Blockmatrixperspektive, vgl. Beweis zu Satz 2.9)

$$\left(\begin{array}{c|c} A_{n-1} & \underline{b} \\ \hline \underline{a}^T & \alpha \end{array} \right) = \left(\begin{array}{c|c} L_{n-1} & \underline{0} \\ \hline \underline{x}^T & 1 \end{array} \right) \left(\begin{array}{c|c} U_{n-1} & \underline{y} \\ \hline \underline{0}^T & \xi \end{array} \right) \Rightarrow L_{n-1} \underline{x} = \underline{b} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ * \\ \vdots \\ * \end{pmatrix} \Rightarrow \underline{x} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ * \\ \vdots \\ * \end{pmatrix} \Rightarrow U \text{ Bandmatrix, obere Bandbreite } p$$

↑
untere Δ -Matrix

[U_{n-1} Bandmatrix nach Induktionsvor.]

Analog: $M_{n-1}^T \underline{x} = \underline{a} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ * \\ \vdots \\ * \end{pmatrix}$

↑
untere Δ -Matrix

□

```

1 function [L,U] = luband(A,p,q)
2 % Algorithm of Crout: LU-factorization of band matrix A
3 % p = upper bandwidth, q = lower bandwidth
4 n = size(A,1); if (size(A,2) ~= n), error('n ~= m'); end
5 L = eye(n); U = zeros(n,n);
6 for k=1:n
7     % Compute row of U (1+q nonzero entries only)
8     for j=k:min(n,k+q)
9         U(k,j) = A(k,j) - L(k,1:k-1)*U(1:k-1,j); end
10    % Compute column of L (p nonzero off-diagonal entries only)
11    for i=k+1:min(n,k+p)
12        L(i,k) = (A(i,k) - L(i,1:k-1)*U(1:k-1,k)) / U(k,k); end
13 end
    
```

⌊ Verkürzte Schleifen und Skalarprodukte
 Rechenaufwand: #flops = $O(npq)$

← } ineffizient implementiert

Hier: keine Pivot suche \Rightarrow i.A. instabil (aber bedenkenlos durchführbar für A SPD, Satz 2.39)
 \hookrightarrow kann Bandstruktur zerstören