

# Grundlagen

## Fehlerarten

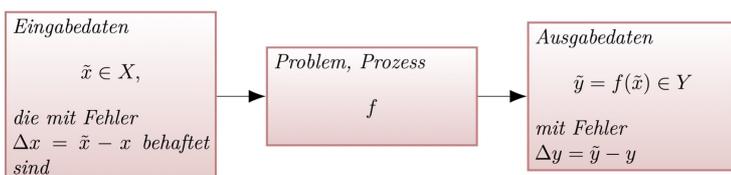
- Modellfehler: Idealisierungen / Vereinfachungen
- Datenfehler: Messfehler
- Verfahrensfehler: numerische Verfahren berechnen Lösungen nur näherungsweise
- Rundungsfehler: reelle Zahlen können mit Floats nur näherungsweise dargestellt werden



The numerical point of view goes back to the earliest mathematical writings. A tablet from the Yale Babylonian Collection (YBC 7289), gives a sexagesimal numerical approximation of the square root of 2, the length of the diagonal in a unit square.

## Konditionierung

beschreibt wie stark Fehler verstärkt werden resp. welche Genauigkeit man bestenfalls bei gestörten Eingabedaten erwarten kann



absolute Fehler:  $\|\vec{\Delta x}\|, \|\vec{\Delta y}\|$

relative Fehler:  $\delta_x = \frac{\|\vec{\Delta x}\|}{\|\vec{x}\|}, \delta_y = \frac{\|\vec{\Delta y}\|}{\|\vec{y}\|}$

absolute Kondition:  $\frac{\|\vec{\Delta x}\|}{\|\vec{\Delta y}\|}$

relative Kondition:  $\frac{\delta_y}{\delta_x}$

} je kleiner desto besser

falls  $f: X \rightarrow Y$  eine reelle Funktion einer reellen Variable ist, also  $X=Y=\mathbb{R}$  ist, gilt:

absolute Konditionszahl:  $K_{abs} = |f'(x)|$

relative Konditionszahl:  $K_{rel} = \left| f'(x) \frac{x}{f(x)} \right|$

} je kleiner desto besser

→ schlecht konditioniert:  $K_{rel} \gg 1$

## Auflösung

Verstärkung von Rundungsfehler bei Subtraktion zweier ungefähr gleich grosser Zahlen

→ schlecht konditionierte Operation

## Stabilität

ein Algorithmus ist stabil, falls die während der Laufzeit erzeugten Fehler in der Grössenordnung des durch die Kondition bedingten unvermeidbaren Fehlers bleiben

# digitale Zahlendarstellung

reelle Zahlen können i.A. nicht auf Computer dargestellt werden, da sie i.A. unendlich viele Nachkommastellen aufweisen

=> Rundungsfehler

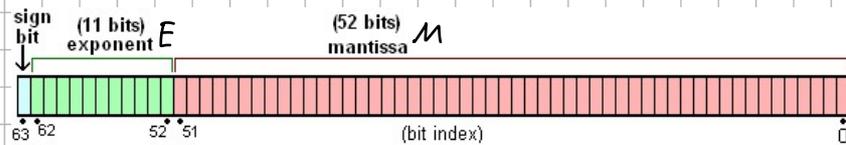
Maschinenzahlen: (endliche) Menge der auf Computer darstellbaren Zahlen

## normalisierte Gleitpunktdarstellung

$$x = f \cdot b^e \quad \text{wobei}$$

- $b$  die Basis des Zahlensystems ist:  $b \in \mathbb{N} \setminus \{1\}$
- der Exponent  $e$  eine ganze Zahl innerhalb gewisser fester Schranken ist:  $r \leq e \leq R$
- die Mantisse/Fraktion  $f$  eine feste Anzahl  $m$  von Stellen hat:
 
$$f = \pm 0.d_1 \dots d_m = \pm \left( \sum_{j=1}^m d_j \cdot b^{-j} \right), \quad d_j \in \{0, 1, \dots, b-1\}$$
- für  $x \neq 0$  stets  $d_1 \neq 0$  ist (Eindeutigkeit der Darstellung)

## Double/float64 nach IEEE



$$\begin{aligned}
 & (-1)^{\text{sign}} \cdot 1.M \cdot 2^{E - (1023)}_{10} \\
 &= (-1)^{\text{sign}} \cdot (1.b_{51}b_{50} \dots b_0)_2 \cdot 2^{(b_{10}b_{11} \dots b_{52})_2 - (1023)}_{10} \\
 &= (-1)^{\text{sign}} \cdot (1 + b_{51} \cdot 2^{-1} + b_{50} \cdot 2^{-2} + \dots + b_1 \cdot 2^{-51} + b_0 \cdot 2^{-52})_{10} \cdot 2^{(b_{52} \cdot 2^{10} + b_{51} \cdot 2^9 + \dots + b_{11} \cdot 2 + b_{10} \cdot 1 - 1023)}_{10}
 \end{aligned}$$

**Representation of Floating Point Numbers in Double Precision IEEE 754 Standard**

Value =  $N = (-1)^S \times 2^{E-1023} \times (1.M)$

$0 < E < 2047$   
 Actual exponent is:  
 $e = E - 1023$

exponent:  
 excess 1023  
 binary integer  
 added

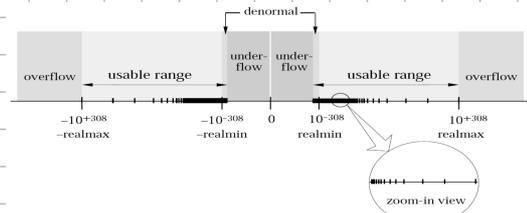
mantissa:  
 sign + magnitude, normalized  
 binary significand with  
 a hidden integer bit: 1.M

Magnitude of numbers that can be represented is in the range:  $2^{-1022} (1.0)$  to  $2^{1023} (2 - 2^{-52})$

Which is approximately:  $2.23 \times 10^{-308}$  to  $1.8 \times 10^{308}$

EECC250 - Shaaban

gewisse Bitmuster für  $\pm \infty$ , NaN, etc reserviert...



## Maschinengenauigkeit

$$\epsilon_{ps} = \frac{1}{2} \cdot 2^{1-m} \quad \text{mit } m = \text{Mantissenlänge (Anzahl Bits)}$$

≙ Auflösungsvermögen des Rechners

≙ kleinst mögliche auf dem Rechner darstellbare Zahl, die zu 1 addiert von der Rundung noch wahrgenommen wird

für Double gilt mit  $m=52$ :  $\epsilon_{ps} \approx 2.22 \cdot 10^{-16}$

-> falls  $K \geq \frac{1}{\epsilon_{ps}} \approx 10^{15}$  => schlecht konditioniert

```
import numpy as np
print(np.finfo(float).eps)
2.220446049250313e-16
```

## Addition von mehreren Zahlen

man sollte Summanden in (betragsmässig) aufsteigender Reihenfolge addieren, um ein genaueres Ergebnis zu erhalten

## Fehlerordnung

Fehler  $p$ -ter Ordnung  $\rightarrow$  wenn man Schrittweite halbiert, verringert sich der Fehler um  $(1/2)^p$

symbolisch:  $O(h^p)$  je grösser  $p$ , desto besser

## Ableitung numerisch berechnen

Approximation der ersten Ableitung:

Rückwärtsdifferenz:

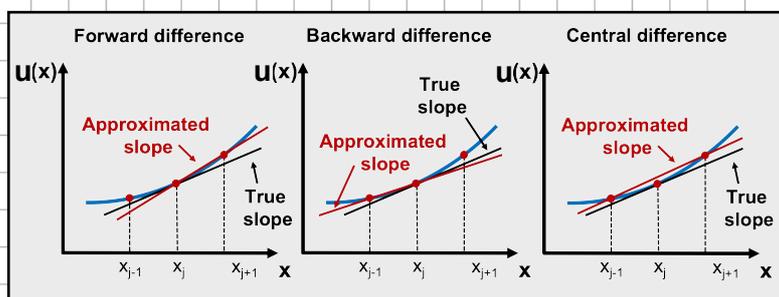
$$\frac{\partial}{\partial x} u(jh) \approx \frac{u_j - u_{j-1}}{h} \text{ mit } O(h)$$

Vorwärtsdifferenz:

$$\frac{\partial}{\partial x} u(jh) \approx \frac{u_{j+1} - u_j}{h} \text{ mit } O(h)$$

zentrale Differenz:

$$\frac{\partial}{\partial x} u(jh) \approx \frac{u_{j+1} - u_{j-1}}{2h} \text{ mit } O(h^2)$$



Approximation der zweiten Ableitung:

zentrale zweite Differenz:

$$\frac{\partial^2}{\partial x^2} u(jh) \approx \frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} \text{ mit } O(h^2)$$

# Lineare Gleichungssysteme



## Das lineare Gleichungssystem

$$\begin{cases} a_{1,1}x_1 + \dots + a_{1,n}x_n = b_1 \\ \vdots \\ a_{n,1}x_1 + \dots + a_{n,n}x_n = b_n \end{cases} \text{ bzw. } \underline{A} \cdot \vec{x} = \vec{b} \text{ mit } \vec{x} \in \mathbb{R}^n, \underline{A} \in \mathbb{R}^{n \times n}, \vec{b} \in \mathbb{R}^n$$

hat eine **eindeutige Lösung** falls:

$$\text{rang}(\underline{A}) = n \Leftrightarrow \det(\underline{A}) \neq 0 \Leftrightarrow \text{kern}(\underline{A}) = \{\vec{0}\} \Leftrightarrow \underline{A} \text{ invertierbar}$$

A heißt in diesem Fall **regulär**

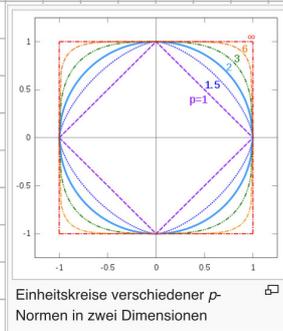
## Cramersche Regel

$$x_j = \frac{\det(\underline{A}_j)}{\det(\underline{A})}, \text{ wobei } \underline{A}_j \text{ aus } \underline{A} \text{ entsteht indem man die } j\text{-te Spalte durch } \vec{b} \text{ ersetzt}$$

**numerisch sehr aufwändig**

## Vektornorm (p-Norm)

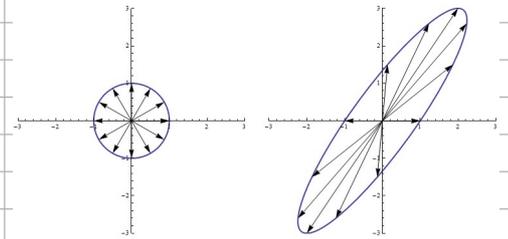
$$\|\vec{x}\|_p := \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$



## Matrixnorm (p-Norm)

$$\|\underline{B}\|_p := \sup_{\|\vec{x}\|_p=1} \|\underline{B}\vec{x}\|_p$$

gibt an wie stark die Matrix B einen Vektor maximal verformt



für  $\underline{B} \in \mathbb{R}^{m \times n}$  gilt:

$$\|\underline{B}\|_\infty = \max_{i=1, \dots, m} \sum_{k=1}^n |b_{i,k}|$$

$$\|\underline{B}\|_1 = \max_{k=1, \dots, n} \sum_{i=1}^m |b_{i,k}|$$

$$\underline{B} = \begin{bmatrix} b_{1,1} & \dots & b_{1,n} \\ b_{2,1} & \dots & b_{2,n} \\ \vdots & & \vdots \\ b_{m-1,1} & \dots & b_{m-1,n} \\ b_{m,1} & \dots & b_{m,n} \end{bmatrix}$$

für  $\underline{A} \in \mathbb{R}^{n \times n}$  gilt:

$$\|\underline{A}\|_\infty = \max_{i=1, \dots, n} \sum_{k=1}^n |a_{i,k}|$$

maximale absolute Zeilensumme

$$\|\underline{A}\|_1 = \max_{k=1, \dots, n} \sum_{i=1}^n |a_{i,k}|$$

maximale absolute Spaltensumme

$$\|\underline{A}\|_2 = \sqrt{\lambda_{\max}(\underline{A}^T \underline{A})}$$

Wurzel des größten Eigenwerts von  $\underline{A}^T \underline{A}$

## Konditionszahl einer Matrix

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

worst-case-wert für Fehlerverstärkungsfaktor

→ i.A. aufwändig zu berechnen, da Inverse aufwändig ist

## relativer Fehler der Lösung

$$\delta_{\vec{x}} = \frac{\|\Delta \vec{x}\|}{\|\vec{x}\|} = \frac{\|\vec{x} - \vec{\tilde{x}}\|}{\|\vec{x}\|}$$

$$\frac{\|\Delta \vec{x}\|}{\|\vec{x}\|} \leq \kappa(A) \frac{\|\Delta \vec{b}\|}{\|\vec{b}\|}$$

## Fehlerverstärkung von $\vec{b}$ zu $\vec{x}$

$$\text{Fehlerverstärkungsfaktor} := \frac{\delta_{\vec{x}}}{\delta_{\vec{b}}} = \frac{\frac{\|\Delta \vec{x}\|}{\|\vec{x}\|}}{\frac{\|\Delta \vec{b}\|}{\|\vec{b}\|}}$$

$$\frac{\delta_{\vec{x}}}{\delta_{\vec{b}}} \leq \kappa(A)$$

## Residuum

$$\vec{r} = \vec{b} - A\vec{x}$$

→ Mass für Genauigkeit einer Annäherung  $\vec{x}$  der Lösung des LGS  $A\vec{x} = \vec{b}$

→ Aussagekraft hängt von Kondition ab:

$$\frac{\|\vec{r}\|}{\kappa(A)\|\vec{b}\|} \leq \frac{\|\vec{x} - \vec{\tilde{x}}\|}{\|\vec{x}\|} \leq \frac{\kappa(A)\|\vec{r}\|}{\|\vec{b}\|}$$

falls schlecht konditioniert  $\Rightarrow \kappa(A)$  gross  
 $\Rightarrow$  Schranken für den relativen Fehler liegen weit auseinander

## Gauss-Elimination ohne Pivotisierung

gegeben: LGS  $A\vec{x} = \vec{b}$

gesucht:  $\vec{x}$

Vorgehen:

1. normal "Gaussian":  $(A|\vec{b}) \rightsquigarrow \dots \rightsquigarrow (R|\vec{c})$  mit rechter/oberer Dreiecksmatrix  $R$
2. Rückwärtseinsetzen:  $R\vec{x} = \vec{c}$  nach  $\vec{x}$  auflösen

## Condition Number



- If  $A$  is square and nonsingular, then  $\text{cond}(A) = \|A\| \cdot \|A^{-1}\|$
- If  $A$  is singular, then  $\text{cond}(A) = \infty$
- If  $A$  is nearly singular, then  $\text{cond}(A)$  is large.
- The condition number measures the ratio of maximum stretch to maximum shrinkage:

$$\|A\| \cdot \|A^{-1}\| = \left( \max_{\|x\|=1} \frac{\|Ax\|}{\|x\|} \right) \left( \min_{\|x\|=1} \frac{\|Ax\|}{\|x\|} \right)^{-1}$$

# LU-Zerlegung mit Skalierung und Pivotisierung

gegeben:  $A$

gesucht:  $\underline{P}, \underline{D}, \underline{L}, \underline{R}$  sodass  $\underline{L}\underline{R} = \underline{P}\underline{D}\underline{A}$

$$\underline{P}\underline{D}\underline{A} = \underline{L}\underline{R}$$
$$\begin{bmatrix} p_{11} & \dots & p_{1n} \\ \vdots & & \vdots \\ p_{n1} & \dots & p_{nn} \end{bmatrix} \begin{bmatrix} d_1 & & 0 \\ & d_2 & \\ 0 & & \ddots \\ & & & d_n \end{bmatrix} \begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} 1 & & & \\ * & 1 & & \\ \vdots & \vdots & \ddots & \\ * & * & \dots & 1 \end{bmatrix} \begin{bmatrix} * & * & \dots & * & * \\ & * & \dots & * & * \\ & & \ddots & & * \\ & & & 0 & * \\ & & & & * \end{bmatrix}$$

Vorgehen:

- bestimme Diagonalmatrix  $\underline{D}$ , welche auf der Diagonale die Zeilenskalierungsfaktoren enthält:

$$\underline{D} = \text{diag}(d_1, \dots, d_n) \quad \text{mit } d_i = \frac{1}{\sum_{k=1}^n |a_{i,k}|}$$

Skalierung verbessert die Kondition der Matrix

- LU-Zerlegung mit Spaltpivotisierung auf  $\underline{D}\underline{A}$  anwenden:

im  $j$ -ten Schritt schaut man ob unter dem Pivotelement irgendwo eine betragsmässig grössere Zahl vorhanden ist

→ falls dies der Fall ist wird diese Zeile mit der  $j$ -ten vertauscht  
→ dieselben Zeilen vertauscht man auch bei  $\underline{P}$  (diese speichert alle Vertauschungen)

dann die Elemente unterhalb Pivotelement eliminieren und die Faktoren in  $\underline{L}$  speichern  
→ dabei Folgendes beachten:

- nicht Zeilen mit Skalar multiplizieren!
- nur Vielfaches (Faktor  $a$ ) der Pivotzeile zu unteren Zeile addieren
- Faktor  $a$  in Spalte von  $\underline{L}$  notieren, deren Nummer mit der Nummer der Pivotzeile übereinstimmt

für  $\underline{L}$ : Zeile Faktor Spalte  
 $\underline{I} = a \cdot \underline{J}$

$\underline{P}$  (Permutationsmatrix):

Matrix mit Nullen und Einsen welche die Information über die Zeilenvertauschungen enthält

→ entsteht aus der Einheitsmatrix indem man alle Zeilenvertauschungen, welche man bei  $A$  vornimmt auch bei der Einheitsmatrix  $\underline{E}_n$  macht

$$\underline{P} = \underline{P}_{n-1, r_{n-1}} \underline{P}_{n-2, r_{n-2}} \dots \underline{P}_{2, r_2} \underline{P}_{1, r_1} \underline{E}_n$$

$$\underline{P}^{-1} = \underline{P}^T$$

Pivotisierung verbessert die Stabilität der LU-Zerlegung

## LGS mit LU-Zerlegung lösen

Gegeben:  $\underline{A}, \vec{b}$  sodass  $\underline{A}\vec{x} = \vec{b}$  und  $\underline{P}, \underline{D}, \underline{L}, \underline{R}$  sodass  $\underline{P}\underline{D}\underline{A} = \underline{L}\underline{R}$

Gesucht:  $\vec{x}$

Vorgehen:

1.  $\underline{L}\vec{y} = \underline{P}\underline{D}\vec{b}$  nach  $\vec{y}$  auflösen (Vorwärtseinsetzen)
2.  $\underline{R}\vec{x} = \vec{y}$  nach  $\vec{x}$  auflösen (Rückwärtseinsetzen)

## Inverse mit LU-Zerlegung

Gegeben:  $\underline{A}$

Gesucht:  $\underline{A}^{-1} := \begin{bmatrix} | & & | \\ \vec{x}_1 & \dots & \vec{x}_n \\ | & & | \end{bmatrix}$

Vorgehen:

$\underline{A}\vec{x}_i = \vec{e}_i$  für  $i=1, \dots, n$  lösen, wobei  $\vec{e}_i$  der  $i$ -te (kanonische) Einheitsvektor ist

## Nachiteration

wegen Rundungsfehler hat man nicht exakt  $\underline{L}$  und  $\underline{R}$  sondern Annäherungen  $\underline{L}$  und  $\underline{R}$   
beim Lösen von  $\underline{A}\vec{x} = \vec{b}$  erhält man somit auch nur eine Näherung  $\vec{x}_0$

es gibt also ein Residuum:  $\vec{r}_0 = \vec{b} - \underline{A}\vec{x}_0$

mithilfe von Nachiteration kann man das Residuum iterativ verkleinern:

Pseudocode:

Für  $k = 0, 1, 2, \dots, \text{nr\_of\_iter}$  berechne:

$$\underline{L}\vec{y}_k = \vec{r}_k; \quad // \text{Vorwärtseinsetzen}$$

$$\underline{R}\vec{\delta}_k = \vec{y}_k; \quad // \text{Rückwärtseinsetzen}$$

$$\vec{x}_{k+1} = \vec{x}_k + \vec{\delta}_k; \quad // \text{Lösung updaten/verbessern}$$

$$\vec{r}_{k+1} = \vec{b} - \underline{A}\vec{x}_{k+1}; \quad // \text{Residuum berechnen}$$

# Cholesky-Zerlegung, LDLT-Zerlegung

effizientere Alternative zu LU für sym. pos. def. Matrizen (spd-Matrizen)  
 → etwa doppelt so schnell um LGS zu lösen

falls  $A \in \mathbb{R}^{n \times n}$  symmetrisch ( $A^T = A$ ) und positiv definit ( $\vec{x}^T A \vec{x} > 0$  für alle  $\vec{x}$ ) ist:

- $A$  hat vollen Rang
- Inverse  $A^{-1}$  ist auch sym. pos. def.
- alle Eigenwerte von  $A > 0$
- Determinante  $\det(A) > 0$
- alle Diagonaleinträge von  $A > 0$
- grösster Eintrag von  $A$  liegt auf der Diagonalen
- $A$  besitzt (eindeutige) Cholesky-Zerlegung

gegeben: sym. pos. def. Matrix  $A$

gesucht:  $L, D$  sodass  $A = LDL^T$  mit  $\begin{cases} \text{normierter linker/unterer Dreiecksmatrix } L \\ \text{Diagonalmatrix } D \text{ mit } d_{ii} > 0 \forall i \end{cases}$

$$A = LDL^T \rightarrow \text{entspricht LU-Zerlegung für } B = DL^T$$

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \dots & a_{nn} \end{bmatrix} = \begin{bmatrix} 1 & & & \\ l_{21} & 1 & & \\ \vdots & \vdots & \ddots & \\ l_{n1} & \vdots & \dots & 1 \end{bmatrix} \begin{bmatrix} d_{11} & & & \\ & d_{22} & & \\ & & \ddots & \\ & & & d_{nn} \end{bmatrix} \begin{bmatrix} 1 & l_{21} & \dots & l_{n1} \\ & 1 & l_{32} & \dots & l_{n2} \\ & & \ddots & \ddots & \vdots \\ & & & 1 & l_{nn} \\ & & & & 1 \end{bmatrix}$$

oder: alles in einer Matrix speichern  
 → für in-place Implementierung

$$L_{ip} = \begin{bmatrix} d_{1,1} & & & \\ l_{2,1} & d_{2,2} & & \\ \vdots & \vdots & \ddots & \\ l_{n,1} & \vdots & \dots & d_{n,n} \\ l_{n1} & l_{n2} & \dots & l_{n,n} & d_{nn} \end{bmatrix}$$

alternative Form → Cholesky-Zerlegung  $A = CC^T$

LDLT kann einfach in Cholesky umgeformt werden:

$$A = CC^T \text{ wobei } C = LD^{1/2}, D^{1/2} = \text{diag}(\sqrt{d_{1,1}}, \dots, \sqrt{d_{n,n}}) \text{ Diagonalelemente von } C$$

umgekehrt kann man mit den Diagonalelementen von  $C$  wieder  $L$  erhalten:

$$L = CD^{-1/2}, D^{-1/2} = \text{diag}\left(\frac{1}{\sqrt{d_{1,1}}}, \dots, \frac{1}{\sqrt{d_{n,n}}}\right) \text{ Kehrwerte der Diagonalelemente von } C$$

## Pseudocode in-place LDLT-Zerlegung

Für  $k = 1, 2, \dots, n$ :  
 $d_{k,k} \leftarrow a_{k,k} - \sum_{j=1}^{k-1} l_{k,j}^2 d_{j,j}$   
 Für  $j = k+1, \dots, n$ :  
 $l_{j,k} \leftarrow (a_{j,k} - \sum_{i=1}^{k-1} l_{j,i} d_{i,i} l_{k,i}) / d_{k,k}$

## Pseudocode Cholesky-Zerlegung

Für  $k = 1, 2, \dots, n$ :  
 $c_{k,k} \leftarrow \text{sqr}t(a_{k,k} - \sum_{j=1}^{k-1} c_{k,j}^2)$   
 Für  $j = k+1, \dots, n$ :  
 $c_{j,k} \leftarrow (a_{j,k} - \sum_{i=1}^{k-1} c_{j,i} c_{k,i}) / c_{k,k}$

LGS  $A\vec{x} = \vec{b}$  mit LDLT:

1.  $L\vec{y} = \vec{b}$  nach  $\vec{y}$  auflösen
2.  $L^T\vec{x} = D^{-1}\vec{y}$  nach  $\vec{x}$  auflösen

LGS  $A\vec{x} = \vec{b}$  mit Cholesky:

1.  $C\vec{y} = \vec{b}$  nach  $\vec{y}$  auflösen
2.  $C^T\vec{x} = \vec{y}$  nach  $\vec{x}$  auflösen

# Lineare Ausgleichsrechnung

## Lineares Ausgleichsproblem

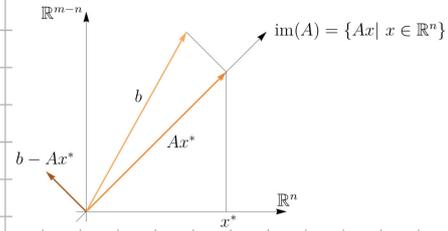
Gegeben: LGS  $A\vec{x} = \vec{b}$  mit  $A \in \mathbb{R}^{m \times n}$ ,  $m > n$  (Zeilen > Spalten)

$\Rightarrow$  überbestimmtes LGS ( $\hat{=}$  mehr Messungen als Modellparameter)  
 $\Rightarrow$  i.A. keine exakte Lösung

beste Lösung (im Sinne von kleinste Fehlerquadrate):

Normalgleichung:  $A^T A \vec{x}^* = A^T \vec{b}$

$\Rightarrow \vec{x}^*$  minimiert (quadratischen) Fehler:  $\|\vec{b} - A\vec{x}^*\|^2 = \min_{\vec{x} \in \mathbb{R}^n} \|\vec{b} - A\vec{x}\|^2$



## Normalgleichung mit LDLT lösen

1. berechne LDLT-Zerlegung von  $A^T A$ :  $L D L^T = A^T A$
2.  $L \vec{y} = A^T \vec{b}$  nach  $\vec{y}$  auflösen (Vorwärtseinsetzen)
3.  $L^T \vec{x} = D^{-1} \vec{y}$  nach  $\vec{x}$  auflösen (Rückwärtseinsetzen)

Nachteile:

- Berechnung von  $A^T A$  für grosse  $m$  aufwendig und teilweise ungenau
- Rundungsfehler in  $A^T A$  und  $A^T \vec{b}$  werden über Cholesky-Verfahren teilweise verstärkt
- kann im Extremfall instabil werden

$\rightarrow$  mit QR-Zerlegung versucht man diese Probleme zu umgehen

$\hookrightarrow$  stabiler und nur leicht höherer Rechenaufwand als LDLT

## spektrale Konditionszahl einer Matrix $A \in \mathbb{R}^{m \times n}$

$$K_2(A) = \frac{\sqrt{\lambda_{\max}(A^T A)}}{\sqrt{\lambda_{\min}(A^T A)}}$$

Wurzel des grössten Eigenwerts von  $A^T A$   
geteilt durch  
Wurzel des kleinsten Eigenwerts von  $A^T A$

falls Spalten von  $A$  linear abhängig ( $A$  singular)  $\Rightarrow K_2(A) = \infty$

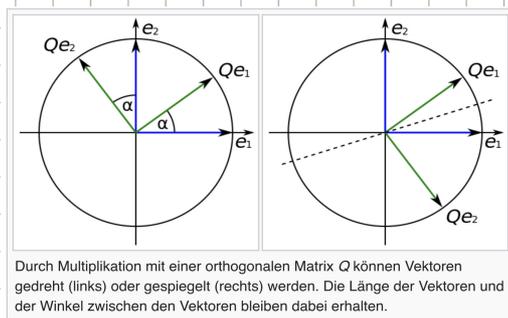
# orthogonale Matrizen

Definition: eine Matrix  $Q \in \mathbb{R}^{m \times m}$  ist orthogonal falls  $Q^T Q = E_m$

Inverse:  $Q^{-1} = Q^T$

Eigenschaften:

- $Q$  orthogonal  $\Rightarrow Q^T$  orthogonal
- $Q, \tilde{Q}$  orthogonal  $\Rightarrow Q\tilde{Q}$  orthogonal
- **Konditionszahl:  $\kappa_2(Q) = 1$**
- **$\kappa(A) = \kappa(AQ) = \kappa(QA)$**
- $\|Q\vec{x}\| = \|\vec{x}\|$  für beliebigen Vektor  $\vec{x}$
- $\|A\| = \|QA\| = \|AQ\|$  für beliebige Matrix  $A$

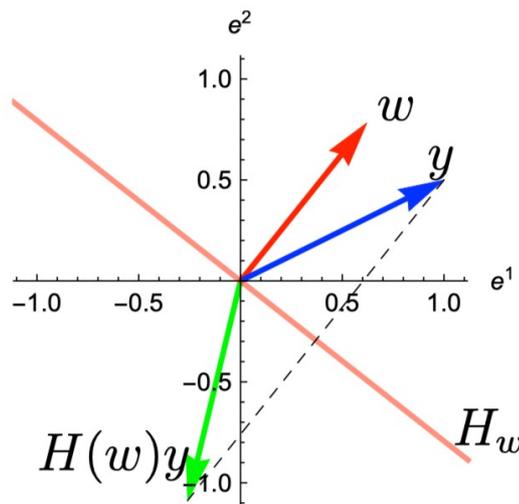


## Householder-Matrix

$$H(\vec{w}) = E_n - 2 \frac{\vec{w}\vec{w}^T}{\langle \vec{w}, \vec{w} \rangle}$$

$$= \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & -1 & \\ & & & \ddots \end{bmatrix} - \frac{2}{\langle \vec{w}, \vec{w} \rangle} \begin{bmatrix} w_1 w_1 & \dots & w_1 w_n \\ \vdots & & \vdots \\ w_n w_1 & \dots & w_n w_n \end{bmatrix}$$

$H(\vec{w})\vec{y}$  entspricht der Spiegelung von  $\vec{y}$  an der durch  $\vec{w}$  definierten Hyperebene  $H_w$  ( $\vec{w}$  steht senkrecht auf dieser Ebene)



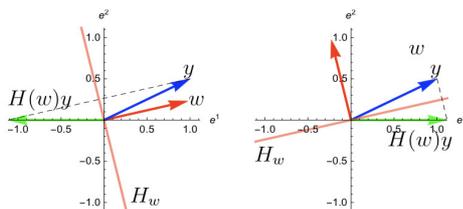
$H(\vec{w})$  ist **symmetrisch** ( $H^T = H$ ) und **orthogonal** ( $H^{-1} = H^T$ )

Bei der QR-Zerlegung wählt man  $\vec{w}$  in jedem Schritt so, dass  $\vec{y}$  durch die Spiegelung auf der  $e_1$ -Achse landet:

$$H(w)\vec{y} = \pm \|\vec{y}\|_2 \vec{e}_1$$

$$\Rightarrow \vec{w} = \vec{y} \pm \text{sign}(y_1) \|\vec{y}\|_2 \vec{e}_1$$

wobei  $\text{sign}(s) = \begin{cases} 1 & \text{für } s \geq 0 \\ -1 & \text{sonst} \end{cases}$



Um Auslöschung zu vermeiden wählt man:  $\vec{w} = \vec{y} + \text{sign}(y_1) \|\vec{y}\|_2 \vec{e}_1$

# QR-Zerlegung

$$\begin{aligned}
 1. \quad Q_1 A &= \begin{pmatrix} \left[ \begin{array}{c|c} & H(\vec{w}^{(1)}) \end{array} \right] & \begin{pmatrix} * & \dots & * \\ \color{red}{\vec{y}^{(1)}} & \vdots & \vdots \\ * & \dots & * \end{pmatrix} \end{pmatrix} = \begin{pmatrix} * & * & * & \dots & * \\ \color{red}{0} & | & * & \dots & * \\ \vdots & \color{red}{\vec{y}^{(2)}} & \vdots & \vdots & \vdots \\ \color{red}{0} & | & * & \dots & * \end{pmatrix} = \underline{A^{(2)}} \\
 2. \quad Q_2 Q_1 A &= \begin{pmatrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & \left[ \begin{array}{c|c} & H(\vec{w}^{(2)}) \end{array} \right] \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \dots & 0 \end{pmatrix} & \begin{pmatrix} * & * & * & \dots & * \\ 0 & \color{red}{\vec{y}^{(2)}} & * & \dots & * \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & * & \dots & * \end{pmatrix} \end{pmatrix} = \begin{pmatrix} * & * & \dots & * \\ 0 & \color{red}{*} & * & \dots & * \\ \vdots & \color{red}{0} & | & * & \dots & * \\ \vdots & \vdots & \color{red}{\vec{y}^{(3)}} & \vdots & \vdots \\ 0 & 0 & | & * & \dots & * \end{pmatrix} = \underline{A^{(3)}} \\
 3. \quad Q_3 Q_2 Q_1 A &= \begin{pmatrix} \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \left[ \begin{array}{c|c} & H(\vec{w}^{(3)}) \end{array} \right] \\ 0 & 0 & \dots & 0 \end{pmatrix} & \begin{pmatrix} * & * & \dots & * \\ 0 & * & * & \dots & * \\ \vdots & 0 & \color{red}{\vec{y}^{(3)}} & * & \dots & * \\ 0 & 0 & * & \dots & * \end{pmatrix} \end{pmatrix} = \begin{pmatrix} * & * & \dots & * \\ 0 & * & * & \dots & * \\ 0 & \color{red}{*} & * & \dots & * \\ \vdots & \color{red}{0} & | & * & \dots & * \\ \vdots & \vdots & \color{red}{\vec{y}^{(4)}} & \vdots & \vdots \\ 0 & 0 & | & * & \dots & * \end{pmatrix} = \underline{A^{(4)}} \\
 \text{usw...}
 \end{aligned}$$

im  $i$ -ten Schritt wahlt man jeweils:

$$\vec{w}^{(i)} = \vec{y}^{(i)} + \text{sign}(y_i^{(i)}) \|\vec{y}^{(i)}\| \vec{e}_i$$

wobei  $\text{sign}(s) = \begin{cases} 1 & \text{fur } s \geq 0 \\ -1 & \text{sonst} \end{cases}$

→ damit eliminiert man alle Eintrage von  $\vec{y}^{(i)}$  ausser dem obersten

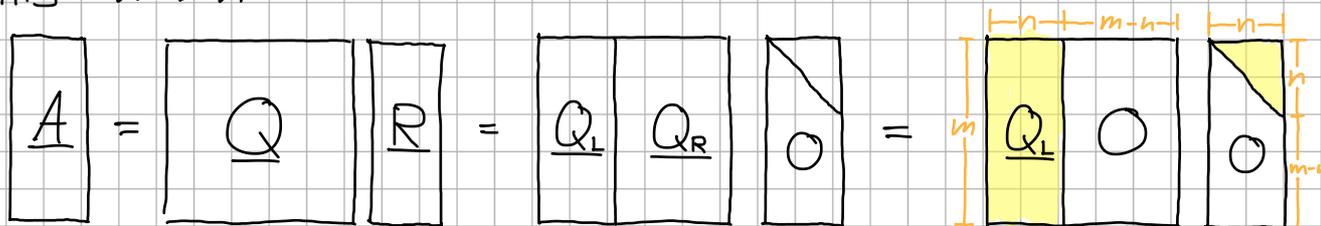
man erhalt somit eine rechte obere Dreiecksmatrix  $\underline{R} = \underline{Q}_n \dots \underline{Q}_1 A$

mit  $\underline{Q} := (\underline{Q}_n \dots \underline{Q}_1)^T = (\underline{Q}_n \dots \underline{Q}_1)^T = \underline{Q}_1^T \dots \underline{Q}_n^T$  folgt:

$$\underline{A} = \underline{Q} \underline{R} \quad \text{wobei} \quad \underline{Q} = \underline{Q}_1^T \dots \underline{Q}_n^T$$

\*: weil orthogonal

falls  $m > n$ :



⇒ man benotigt nur die ersten  $n$  Zeilen von  $\underline{R}$  und erste  $n$  Spalten von  $\underline{Q}$

⇒ reduzierte QR-Zerlegung:

$$\underline{A} = \underline{Q}_L \underline{R}_L$$

Normalengleichung mit QR von  $A$  (!) lösen

$$A^T A \vec{x} = A^T \vec{b} \Leftrightarrow (QR)^T QR \vec{x} = (QR)^T \vec{b} \Leftrightarrow R^T \underbrace{Q^T Q}_{=E} R \vec{x} = R^T Q^T \vec{b} \Leftrightarrow R^T R \vec{x} = R^T Q^T \vec{b}$$
$$\Leftrightarrow \boxed{R \vec{x} = Q^T \vec{b} = \vec{z}}$$

$$\begin{bmatrix} \square & \\ & \square \\ & & 0 \end{bmatrix} \begin{bmatrix} \vec{x} \\ \\ \end{bmatrix} = \begin{bmatrix} \square & \\ & \square \end{bmatrix} Q^T \begin{bmatrix} \vec{b} \\ \\ \end{bmatrix} = \begin{bmatrix} \vec{z} \\ \\ \end{bmatrix} \rightarrow \text{nach } \vec{x} \text{ auflösen (Rückwärtseinsetzen)}$$

**Vorteil:** i.A. schlecht konditionierte Matrix  $A^T A$  kommt nicht mehr vor, da die QR-Zerlegung von  $A$  direkt verwendet werden kann

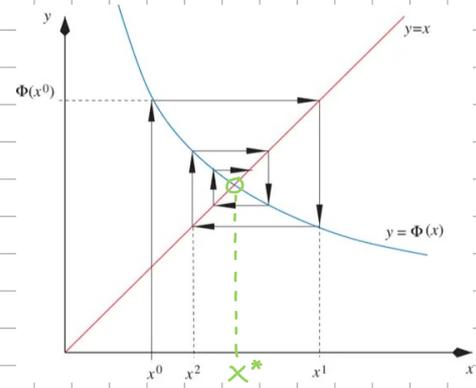
# Nichtlineare Gleichungssysteme ( $m=n$ resp $\# \text{Gleichungen} = \# \text{Unbekannte}$ )

## Fixpunktiteration

Fixpunktgleichung:  $\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} \Phi_1(x_1, \dots, x_n) \\ \vdots \\ \Phi_n(x_1, \dots, x_n) \end{pmatrix}$  bzw  $\vec{\Phi}(\vec{x}) = \vec{x}$

gesucht: Fixpunkt  $\vec{x}^*$ , sodass Fixpunktgleichung  $\vec{\Phi}(\vec{x}^*) = \vec{x}^*$  erfüllt ist

Iterationsvorschrift:  $\vec{x}_{k+1} = \vec{\Phi}(\vec{x}_k)$



## Voraussetzungen:

- es muss Fixpunktgleichung  $\vec{\Phi}(\vec{x}) = \vec{x}$  vorliegen  $\rightarrow$  falls man  $\vec{f}(\vec{x}) = \vec{0}_n$  hat, muss man  $\vec{x}$  auf beiden Seiten addieren:  $\vec{\Phi}(\vec{x}) := \vec{x} + \vec{f}(\vec{x}) = \vec{x}$
- $\vec{\Phi}: E \rightarrow E$  mit  $E \subseteq \mathbb{R}^n$  «Selbstabbildung»
- $\max \|D\Phi(\vec{x})\| < 1$  «Kontraktion»

$\Rightarrow$  Dann existiert genau ein Fixpunkt  $\vec{x}^* \in E$  und Fixpunktiteration konvergiert linear gegen  $\vec{x}^*$

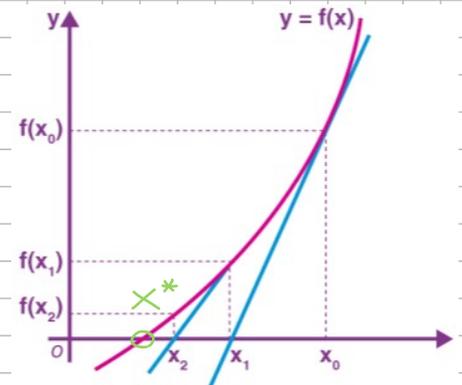
## Newtonverfahren

Nullstellengleichung:  $\begin{pmatrix} f_1(x_1, \dots, x_n) \\ \vdots \\ f_n(x_1, \dots, x_n) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$  bzw  $\vec{f}(\vec{x}) = \vec{0}_n$

gesucht: Lösung  $\vec{x}^*$ , sodass Nullstellengleichung  $\vec{f}(\vec{x}^*) = \vec{0}_n$  erfüllt ist

Iterationsvorschrift:

$Df(\vec{x}_k) \vec{\delta}_k = \vec{f}(\vec{x}_k)$  nach  $\vec{\delta}_k$  auflösen  
 $\vec{x}_{k+1} = \vec{x}_k - \vec{\delta}_k$



falls eine Lösung  $\vec{x}^*$  existiert und man gute Startwerte (idealerweise eine Näherungslösung für  $\vec{x}^*$ ) hat, konvergiert das Newtonverfahren unter gewissen Voraussetzungen quadratisch gegen  $\vec{x}^*$

mögliches Problem:

falls Jacobi in einer Nullstelle  $\vec{x}^*$  singular (nichtinvertierbar), weil bspw. die Graphen der Einzelnen Komponenten  $f_i$  die  $x_1 \dots x_n$ -Ebene bei  $\vec{x}^*$  nur berühren, aber keinen Nulldurchgang haben, kann  $\|\vec{\delta}_k\|$  in der Nähe von  $\vec{x}^*$  sehr gross werden und das Verfahren divergiert

## vereinfachtes Newtonverfahren

um Rechenaufwand zu reduzieren, benutzt man oft auch das vereinfachte Newtonverfahren, bei dem man die Jacobi, bzw. dessen LU-Zerlegung, nur einmal, am Anfang, für den Startwert berechnet:

Iterationsvorschrift:

$$\begin{aligned} \underline{Df}(\vec{x}_0) \vec{\delta}_k &= \vec{f}(\vec{x}_k) \text{ nach } \vec{\delta}_k \text{ auflösen} \quad (\text{am besten mit LU von } \underline{Df}(\vec{x}_0)) \\ \vec{x}_{k+1} &= \vec{x}_k - \vec{\delta}_k \end{aligned}$$

einzelnen Schritt weniger aufwändig, braucht jedoch normalerweise mehr Schritte

# Nichtlineare Ausgleichsrechnung ( $m > n$ resp. # Gleichungen $>$ # Unbekannte)

Motivation:

$m$  Messungen  $b_i$  zu Zeitpunkten  $t_i$ ,  $i \in \{1, \dots, m\}$ , für math. Modell mit  $n$  Parametern  $x_j$ ,  $j \in \{1, \dots, n\}$ :

$$\begin{cases} y(t_1, \vec{x}) = b_1 \\ \vdots \\ y(t_m, \vec{x}) = b_m \end{cases} \Leftrightarrow \begin{cases} y(t_1, \vec{x}) - b_1 = 0 \\ \vdots \\ y(t_m, \vec{x}) - b_m = 0 \end{cases} \quad \text{resp. } \vec{F}(\vec{x}) = \vec{0}_m$$

mit  $F_i(\vec{x}) := y(t_i, \vec{x}) - b_i$

Problemstellung:

gegeben:  $\vec{F}: D \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^m$  mit  $m > n \Rightarrow$  überbestimmtes nichtlineares GS

gesucht:  $\vec{x}^*$ , sodass  $\vec{F}(\vec{x}^*) = \vec{0}_m$

$\Rightarrow$  i.A. keine Lösung, da Bild im  $(\vec{F})$  ein nur  $n$ -dimensionales Objekt im  $m$ -dimensionalen  $\mathbb{R}^m$  ist

$\Rightarrow$  man **minimiert** stattdessen die **Residuenquadratsumme**  $\|\vec{F}(\vec{x})\|^2 = F_1(\vec{x})^2 + \dots + F_m(\vec{x})^2$   
resp die Hälfte davon  $\frac{1}{2} \|\vec{F}(\vec{x})\|^2 = \frac{1}{2} \vec{F}(\vec{x})^T \vec{F}(\vec{x}) =: \phi(\vec{x}) \rightarrow$  haben die gleichen Minima

$\phi$  hat in  $\vec{x}^*$  ein lokales Minimum genau dann wenn

- $\vec{\nabla} \phi(\vec{x}^*) = \vec{0}_n$
- $H\phi(\vec{x}^*)$  sym. pos. def.

Gradient von  $\phi$ :  $\vec{\nabla} \phi(\vec{x}) = \underline{DF}(\vec{x})^T \vec{F}(\vec{x})$

Hessesche Matrix von  $\phi$ :  $H\phi(\vec{x}) = \underline{DF}(\vec{x})^T \underline{DF}(\vec{x}) + \sum_{i=1}^m F_i(\vec{x}) \underline{HF}_i(\vec{x})$

numerisch aufwändig  
 $\hookrightarrow$  wird oft weggelassen

## Newton-Verfahren

Idee: Nullstelle von  $\vec{\nabla} \phi$  bestimmen  $\rightarrow$  normales Newtonverfahren auf  $\vec{\nabla} \phi(\vec{x})$  anwenden

Iterationsvorschrift:

$H\phi(\vec{x}_k) \vec{\delta}_k = \vec{\nabla} \phi(\vec{x}_k)$  nach  $\vec{\delta}_k$  auflösen

$\vec{x}_{k+1} = \vec{x}_k - \vec{\delta}_k$

# Gauss-Newton-Verfahren

Idee:  $\vec{F}(\vec{x})$  in jedem Schritt mit 1tem Taylorpolynom linearisieren und ein Minimum davon mit linearer Ausgleichsrechnung finden

$$\vec{F}(\vec{x}) \approx \vec{F}(\vec{x}_k) + \underline{DF}(\vec{x}_k) \underbrace{(\vec{x} - \vec{x}_k)}_{=\vec{\delta}} \stackrel{!}{=} \vec{0}_m$$
$$\Rightarrow \underbrace{\underline{DF}(\vec{x}_k)}_{m \times n} \underbrace{\vec{\delta}}_n = -\underbrace{\vec{F}(\vec{x}_k)}_m \rightarrow \text{überbestimmtes LGS, keine Lösung}$$
$$\Rightarrow \underbrace{\underline{DF}(\vec{x}_k)^T}_{n \times n} \underbrace{\underline{DF}(\vec{x}_k)}_n \vec{\delta} = -\underbrace{\underline{DF}(\vec{x}_k)^T}_{n} \underbrace{\vec{F}(\vec{x}_k)}_n \rightarrow \text{Normalgleichung}$$

## Pseudocode

Für  $k = 0, 1, 2, \dots$  :

$$\underline{J}_k = \underline{DF}(\vec{x}_k), \quad \vec{y}_k = \vec{F}(\vec{x}_k);$$

$$\underline{J}_k^T \underline{J}_k \vec{\delta}_k = -\underline{J}_k^T \vec{y}_k \quad \text{nach } \vec{\delta}_k \text{ auflösen};$$

$$\vec{x}_{k+1} = \vec{x}_k + \vec{\delta}_k;$$

$$\text{Falls } \|\underline{DF}(\vec{x}_k)^T \vec{F}(\vec{x}_k)\| < \text{tol}: \quad \text{// Norm des Gradienten von } \phi = \frac{1}{2} \|\vec{F}(\vec{x})\|^2$$

Abbruch; // falls keine Steigung mehr vorhanden  $\Rightarrow$  Minimum erreicht

um Stabilität zu verbessern wird Newtonschritt oft gedämpft  
 $\rightarrow$  mit Dämpfung verhindert man, dass man zu weit springt und damit ausreischen bergauf wandert

## Pseudocode für gedämpfte Variante

Für  $k = 0, 1, 2, \dots$  :

$$\underline{J}_k = \underline{DF}(\vec{x}_k), \quad \vec{y}_k = \vec{F}(\vec{x}_k);$$

$$\underline{J}_k^T \underline{J}_k \vec{\delta}_k = -\underline{J}_k^T \vec{y}_k \quad \text{nach } \vec{\delta}_k \text{ auflösen};$$

$$\mu = 1.0;$$

$$\text{Solange } \|\vec{F}(\vec{x}_k + \mu \vec{\delta}_k)\|^2 > \|\vec{F}(\vec{x}_k)\|^2 \quad \text{und } \mu > \mu_{\min}: \quad \text{// gehts bergauf?}$$

$$\mu = \mu/2; \quad \text{// dann lieber weniger weit springen}$$

$$\vec{x}_{k+1} = \vec{x}_k + \mu \vec{\delta}_k;$$

$$\text{Falls } \|\underline{DF}(\vec{x}_k)^T \vec{F}(\vec{x}_k)\| < \text{tol}: \quad \text{// Norm des Gradienten von } \phi = \frac{1}{2} \|\vec{F}(\vec{x})\|^2$$

Abbruch; // falls keine Steigung mehr vorhanden  $\Rightarrow$  Minimum erreicht

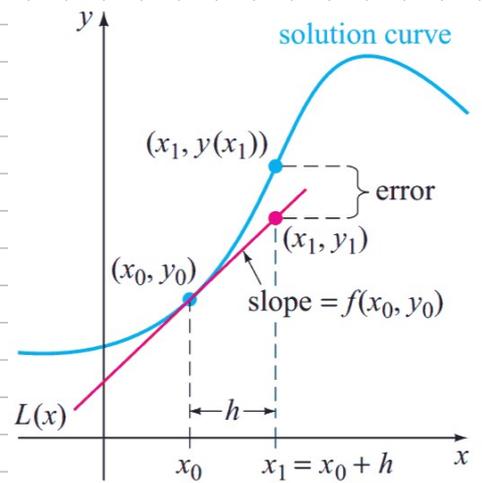
andere Variante um Stabilität zu verbessern: Levenberg-Marquardt-Verfahren  
 $\rightarrow$  Kombination aus Newton-Verfahren und Gradient-Descent:  
<https://www.youtube.com/watch?v=dPZo74SbkeQ>

# gewöhnliche DGL (ODE)

Anfangswertproblem:

$$y'(x) = f(x, y), \quad x \in [a, b] \quad \text{Steigung von } y \text{ bei } x$$

$$y(x_0) = y_0$$



Idee:

$[a, b]$  äquidistant unterteilen:

$$a = x_0 < x_1 < \dots < x_{n-1} < x_n = b$$

mit  $x_k = x_0 + k \cdot h$ ,  $h = \frac{b-a}{n}$

ausgehend von  $y_0 = y(x_0)$  mithilfe von  $y'(x) = f(x, y)$  ( $\rightarrow$  Richtungsfeld) Schritt für Schritt Näherungen  $y_k$  für  $y(x_k)$  bestimmen

## einstufige Verfahren

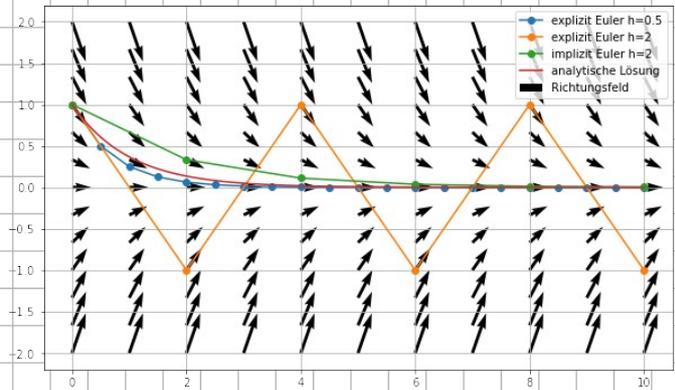
Euler explizit:

mit Steigung am aktuellen Punkt  $(x_k, y_k)$  einen Schritt machen:

$$f(x_k, y_k) \approx \frac{y_{k+1} - y_k}{h}$$

$$\Rightarrow y_{k+1} = y_k + h \cdot f(x_k, y_k) \quad \begin{array}{c} 0 \\ | \\ 1 \end{array}$$

- unterschätzt exakte Lösung
- oszilliert bei zu grosser schrittweite  $h$



Euler implizit:

mit Steigung am nächsten Punkt  $(x_{k+1}, y_{k+1})$  einen Schritt machen  
nicht direkt möglich, da  $y_{k+1}$  noch nicht bekannt  $\Rightarrow$  implizite Gleichung für  $y_{k+1}$

$$f(x_{k+1}, y_{k+1}) \approx \frac{y_{k+1} - y_k}{h}$$

$$\Rightarrow y_{k+1} = y_k + h \cdot f(x_{k+1}, y_{k+1}) \quad \begin{array}{c} 1 \\ | \\ 1 \end{array}$$

$\rightarrow$  mit Newtonverfahren nach  $y_{k+1}$  auflösen

- überschätzt exakte Lösung

## allgemeines einstufiges RK-Verfahren

$$k = f(x_{i-1} + c \cdot h, y_{i-1} + h \cdot a \cdot k) \quad \begin{array}{c|c} c & a \\ \hline & b \end{array}$$

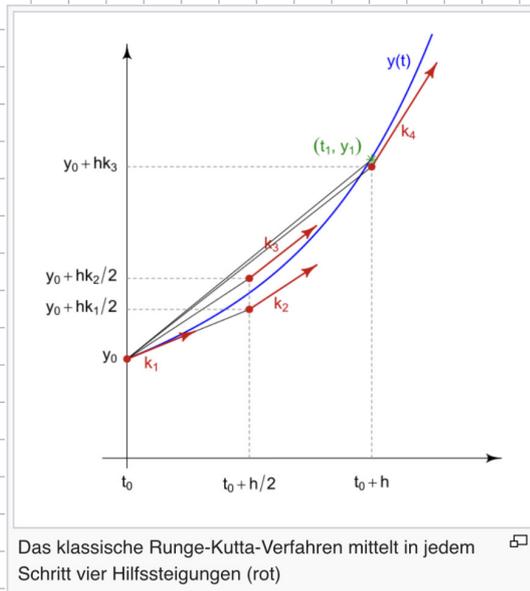
$$y_i = y_{i-1} + h \cdot b \cdot k$$

$\rightarrow$  explizit falls  $a=0$

## mehrstufige Verfahren

Verfahren nach Runge (RK2):

$$\begin{aligned} k_1 &= f(x_{i-1}, y_{i-1}) \\ k_2 &= f(x_{i-1} + 1/2 \cdot h, y_{i-1} + 1/2 \cdot h \cdot k_1) \\ y_i &= y_{i-1} + h \cdot k_2 \end{aligned} \quad \begin{array}{c|cc} 0 & 0 & 0 \\ 0.5 & 0.5 & 0 \\ \hline & 0 & 1 \end{array}$$



Verfahren nach Heun:

$$\begin{aligned} k_1 &= f(x_{i-1}, y_{i-1}) \\ k_2 &= f(x_{i-1} + h, y_{i-1} + h \cdot k_1) \\ y_i &= y_{i-1} + h \cdot (1/2 \cdot k_1 + 1/2 \cdot k_2) \end{aligned} \quad \begin{array}{c|cc} 0 & 0 & 0 \\ 1 & 1 & 0 \\ \hline & 0.5 & 0.5 \end{array}$$

Klassisches Runge-Kutta-Verfahren (RK4):

$$\begin{aligned} k_1 &= f(x_{i-1}, y_{i-1}) \\ k_2 &= f(x_{i-1} + 1/2 \cdot h, y_{i-1} + 1/2 \cdot h \cdot k_1) \\ k_3 &= f(x_{i-1} + 1/2 \cdot h, y_{i-1} + 1/2 \cdot h \cdot k_2) \\ k_4 &= f(x_{i-1} + h, y_{i-1} + h \cdot k_3) \\ y_i &= y_{i-1} + h \cdot (1/6 k_1 + 1/3 k_2 + 1/3 k_3 + 1/6 k_4) \end{aligned} \quad \begin{array}{c|cccc} 0 & 0 & 0 & 0 & 0 \\ 1/2 & 1/2 & 0 & 0 & 0 \\ 1/2 & 0 & 1/2 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ \hline & 1/6 & 1/3 & 1/3 & 1/6 \end{array}$$

## allgemeine s-stufige RK-Verfahren

$$\begin{aligned} k_i &= f(x_{i-1} + c_i h, y_{i-1} + h(a_{i1}k_1 + \dots + a_{is}k_s)) \\ &\vdots \\ k_s &= f(x_{i-1} + c_s h, y_{i-1} + h(a_{s1}k_1 + \dots + a_{ss}k_s)) \\ y_i &= y_{i-1} + h \cdot (b_1 k_1 + \dots + b_s k_s) \end{aligned} \quad \begin{array}{c|ccc} c_i & a_{i1} & \dots & a_{is} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \dots & a_{ss} \\ \hline & b_1 & \dots & b_s \end{array} \quad \text{resp.} \quad \vec{c} \mid \begin{array}{c} A \\ \hline \vec{b}^T \end{array}$$

(F)IRK-Verfahren (fully implicit):

mind. ein Eintrag von A oberhalb Hauptdiagonale  $\neq 0$

man muss in jedem Schritt Stufengleichung  $\vec{F}(\vec{k}) \stackrel{!}{=} \vec{0}_s$  lösen (z.B. mit Newton-Verfahren) um die Steigungen  $k_1, \dots, k_s$  zu berechnen:

$$\vec{F}(\vec{k}) = \begin{pmatrix} F_1(k_1, \dots, k_s) \\ \vdots \\ F_s(k_1, \dots, k_s) \end{pmatrix} := \begin{pmatrix} k_1 - f(x_{i-1} + c_1 h, y_{i-1} + h(a_{11}k_1 + \dots + a_{1s}k_s)) \\ \vdots \\ k_s - f(x_{i-1} + c_s h, y_{i-1} + h(a_{s1}k_1 + \dots + a_{ss}k_s)) \end{pmatrix} \stackrel{!}{=} \vec{0}_s$$

DIRK-Verfahren (diagonally implicit):

A ist linke (untere) Dreiecksmatrix (nur Nullen oberhalb Hauptdiagonale)

Vorteil gegenüber FIRK: Stufengleichungen  $\vec{F}(\vec{k}) \stackrel{!}{=} \vec{0}_s$  können sequenziell berechnet werden, da in j-ter Stufe  $k_1, \dots, k_{j-1}$  bereits bekannt sind

# SDIRK-Verfahren (singly diagonally implicit):

$A$  ist linke (untere) Dreiecksmatrix und alle Diagonalelemente sind gleich  
haben vor allem für grosse Systeme von gDgl'n einen Vorteil

# ERK-Verfahren (explicit)

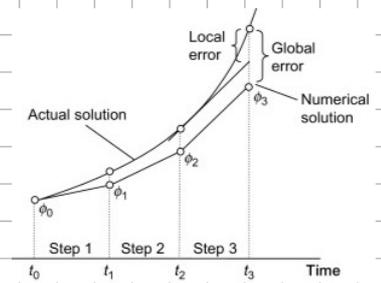
alle Einträge von  $A$  auf und oberhalb der Hauptdiagonale = 0

Vorteil: müssen keine Gleichungen gelöst werden, in der  $j$ ten Stufe steht explizite Formel für  $k_j$  → können direkt eins nach dem andern berechnet werden

# globaler Diskretisierungsfehler

globaler Fehler an Stelle  $x_k = a + kh$  ist definiert als:

$$e(x_k, h) := y(x_k) - y_k \quad (\text{exakter Wert} - \text{numerische Approximation})$$

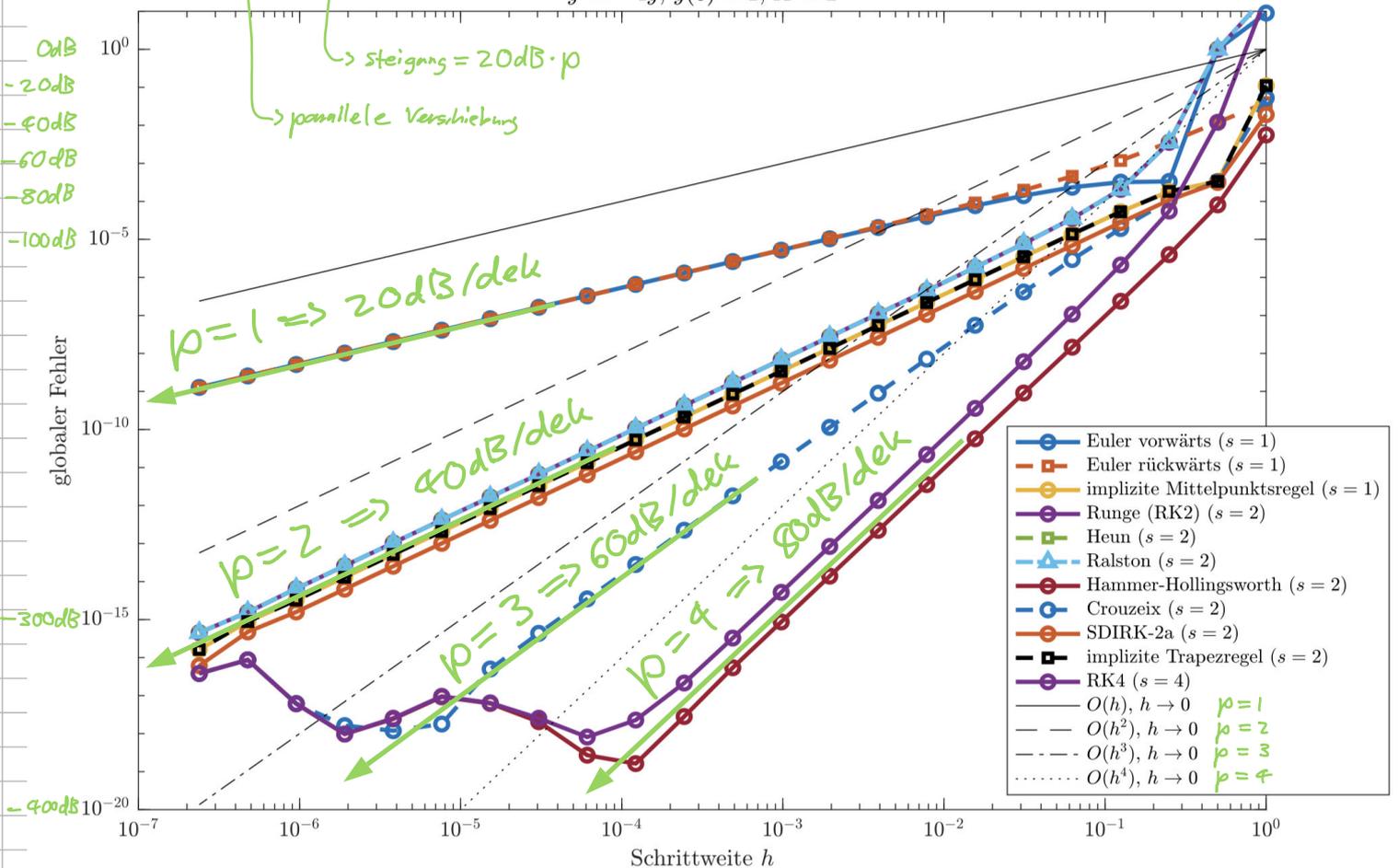


# Konvergenzordnung

Verfahren hat Konvergenzordnung  $p$  falls für irgendeine Konstante  $C > 0$

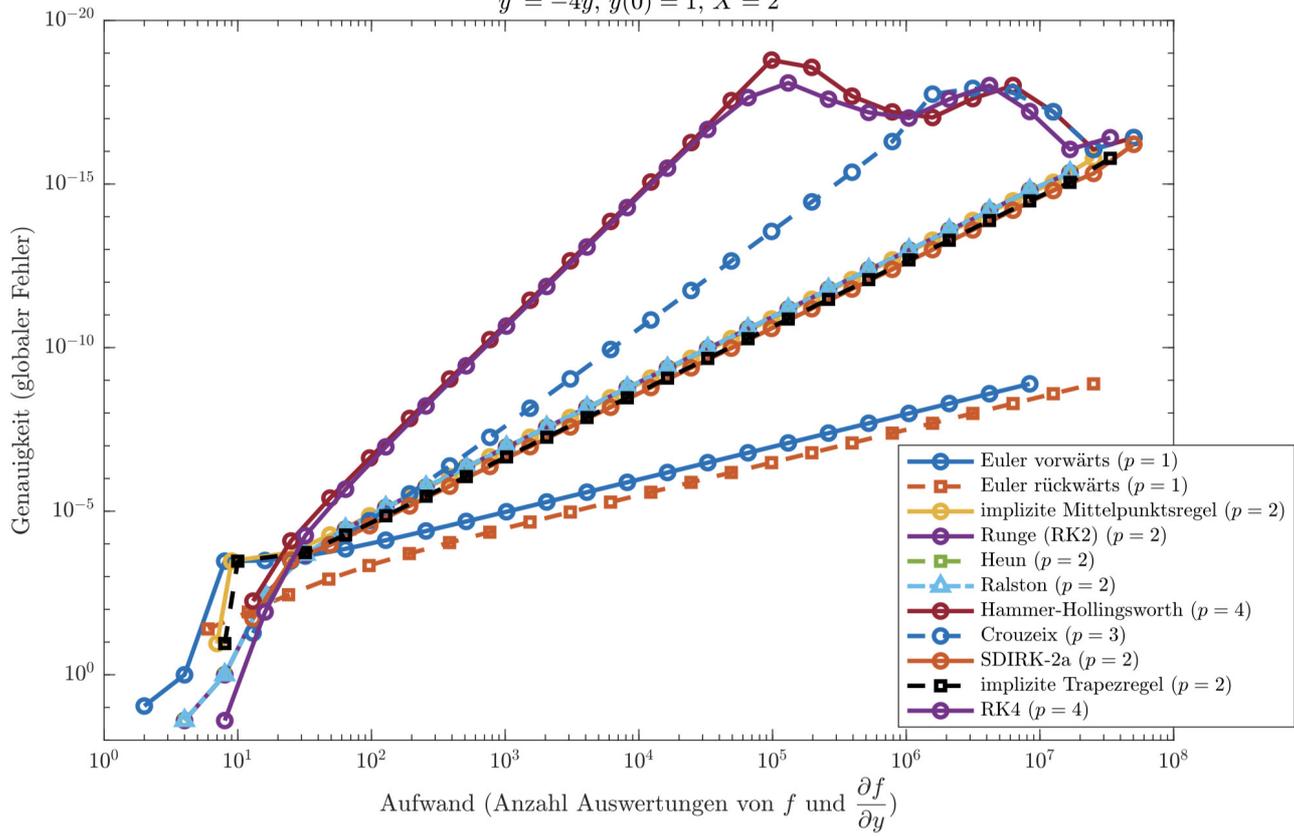
$$e(x_k, h) \leq C \cdot h^p, \quad h \rightarrow 0 \text{ für alle } x_k \text{ gilt}$$

$$y' = -4y, \quad y(0) = 1, \quad X = 2$$



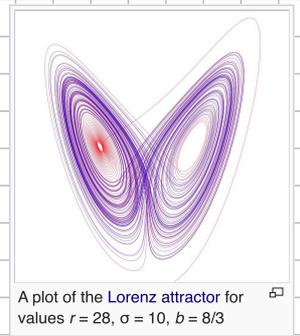
# Genauigkeits - Aufwand-Diagramm

$$y' = -4y, y(0) = 1, X = 2$$



⇒ bei gleichem Aufwand höhere Genauigkeit mit mehrstufigen RK-Verfahren

# Systeme von gewöhnlichen DGLn



Anfangswertproblem:

$$\vec{y}'(x) = \vec{f}(x, \vec{y}) \quad \text{resp} \quad \begin{pmatrix} \frac{d}{dx} y_1(x) \\ \vdots \\ \frac{d}{dx} y_n(x) \end{pmatrix} = \begin{pmatrix} f_1(x, y_1, \dots, y_n) \\ \vdots \\ f_n(x, y_1, \dots, y_n) \end{pmatrix}, \quad x \in [a, b]$$

$$\vec{y}(x_0) = \vec{y}_0$$

alle Verfahren für gDGL kann man auf Systeme von gDGLn anwenden, indem man  $y$  durch  $\vec{y}$ ,  $f$  durch  $\vec{f}$ ,  $k_i$  durch  $\vec{k}_i$ , etc. ersetzt

## allgemeine s-stufige RK-Verfahren

$$\begin{aligned} \vec{k}_i &= \vec{f}(x_{i-1} + c_i h, \vec{y}_{i-1} + h(a_{i1}\vec{k}_1 + \dots + a_{is}\vec{k}_s)) \\ &\vdots \\ \vec{k}_s &= \vec{f}(x_{i-1} + c_s h, \vec{y}_{i-1} + h(a_{s1}\vec{k}_1 + \dots + a_{ss}\vec{k}_s)) \\ \vec{y}_i &= \vec{y}_{i-1} + h \cdot (b_1\vec{k}_1 + \dots + b_s\vec{k}_s) \end{aligned}$$

$$\begin{array}{c|ccc} c_1 & a_{11} & \dots & a_{1s} \\ \vdots & \vdots & & \vdots \\ c_s & a_{s1} & \dots & a_{ss} \\ \hline & b_1 & \dots & b_s \end{array} \quad \text{resp} \quad \begin{array}{c} \vec{c} \\ \hline \vec{A} \\ \vec{b}^T \end{array}$$

(F)IRK-Verfahren (fully implicit):

man muss in jedem Schritt System von  $s \cdot n$  nichtlineare Gleichungen für  $s \cdot n$  Unbekannte lösen:

$$\vec{F}(\vec{k}) = \begin{pmatrix} \vec{F}_1(\vec{k}_1, \dots, \vec{k}_s) \\ \vdots \\ \vec{F}_s(\vec{k}_1, \dots, \vec{k}_s) \end{pmatrix} := \begin{pmatrix} \vec{k}_1 - \vec{f}(x_{i-1} + c_1 h, \vec{y}_{i-1} + h(a_{11}\vec{k}_1 + \dots + a_{1s}\vec{k}_s)) \\ \vdots \\ \vec{k}_s - \vec{f}(x_{i-1} + c_s h, \vec{y}_{i-1} + h(a_{s1}\vec{k}_1 + \dots + a_{ss}\vec{k}_s)) \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} \vec{0}_n \\ \vdots \\ \vec{0}_n \end{pmatrix} = \vec{0}_{sn}$$

ERK-Verfahren (explicit)

man muss keine Gleichungen lösen,  $\vec{k}_i$  können eins nach dem anderen direkt mit  $\vec{f}$  berechnet werden

## Stabilitätsanalyse

Stabilitätsfunktion  $R: \mathbb{C} \rightarrow \mathbb{C}$  eines allg.  $s$ -stufigen RK-Verfahrens mit Verfahrensparameter  $A \in \mathbb{R}^{s \times s}$  und  $\vec{b} \in \mathbb{R}^s$  ist gegeben durch:

$$R(z) = \frac{\det(E_s - zA + z\vec{1}_s \vec{b}^T)}{\det(E_s - zA)} = \frac{\det \left( \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} - z \begin{bmatrix} a_{11} & \dots & a_{1s} \\ \vdots & & \vdots \\ a_{s1} & \dots & a_{ss} \end{bmatrix} + z \begin{bmatrix} b_1 & \dots & b_s \\ \vdots & & \vdots \\ b_1 & \dots & b_s \end{bmatrix} \right)}{\det \left( \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} - z \begin{bmatrix} a_{11} & \dots & a_{1s} \\ \vdots & & \vdots \\ a_{s1} & \dots & a_{ss} \end{bmatrix} \right)} \in \mathbb{C}, \quad z \in \mathbb{C}$$

Stabilitätsgebiet  $S$  des Verfahrens:

$$S = \{z \in \mathbb{C} \text{ für die } |R(z)| < 1 \text{ ist}\} \subseteq \mathbb{C}$$

Verfahren ist im Punkt  $\vec{q}(x)$  des Richtungsfeldes stabil, falls für alle Eigenwerte  $\lambda_j$  der Jacobi  $\underline{J} = \partial_y \vec{f}(x, \vec{q}(x))$  das Produkt aus Schrittweite  $h$  und Eigenwert  $\lambda_j$  im Stabilitätsgebiet liegt:

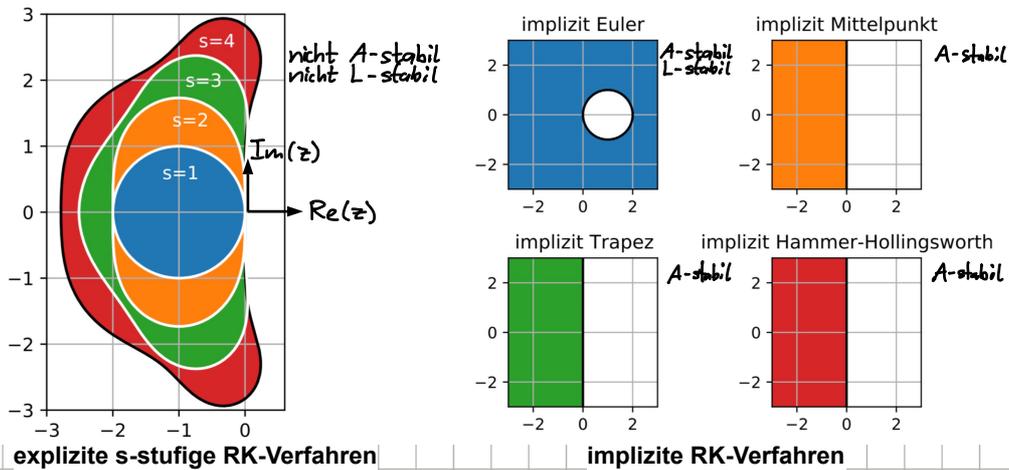
$$|R(h\lambda_j)| < 1 \text{ für alle Eigenwerte } \lambda_j \text{ von } \underline{J} \Rightarrow \text{Verfahren in } \vec{q}(x) \text{ stabil}$$

Verfahren ist

- A-stabil, falls „ $\text{Re}(z) < 0 \Rightarrow z \in S$ “ erfüllt ist  
 → komplette linke Halbebene im Stabilitätsgebiet
- L-stabil, falls A-stabil und zusätzlich  $\lim_{\text{Re}(z) \rightarrow -\infty} |R(z)| = 0$   
 → wenn man weit rechts geht, kommt man irgendwann ins Stabilitätsgebiet

für explizite Verfahren ist  $S$  immer beschränkt  $\Rightarrow$  nicht A-stabil

Stabilitätsbereiche  $S$  einiger RK-Verfahren:



### Energieerhaltung

- konservativ  $\Rightarrow$  Energie bleibt erhalten, falls dies beim physikalischen System auch so ist
- dissipativ  $\Rightarrow$  Energie geht verloren  $\rightarrow$  «numerische Reibung»
- anti-dissipativ  $\Rightarrow$  Energie nimmt zu  $\rightarrow$  unerwünscht!

### Übersicht

Runge-Kutta-Verfahren	s	p		Stabilität	Energie, $\lambda = -\nu, h > 0$
Euler vorwärts	1	1	explizit	–	anti-dissipativ $\rightarrow$ überschätzt exakte Lösung $\Rightarrow E \uparrow$
Euler rückwärts	1	1	implizit	L-stabil	dissipativ $\rightarrow$ unterschätzt exakte Lösung $\Rightarrow E \downarrow$
implizite Mittelpunktsregel	1	2	implizit	A-stabil	konservativ $\rightarrow$ Mischung aus beiden Euler $\Rightarrow E \text{ const.}$
Runge (RK2)	2	2	explizit	–	anti-dissipativ
Heun	2	2	explizit	–	anti-dissipativ
Ralston	2	2	explizit	–	anti-dissipativ
Hammer-Hollingsworth	2	4	implizit	A-stabil	konservativ
Crouzeix	2	3	implizit	L-stabil	dissipativ
SDIRK-2a	2	2	implizit	L-stabil	dissipativ
implizite Trapezregel	2	2	implizit	A-stabil	konservativ
RK4	4	4	explizit	–	anti-dissipativ für $h > 2\sqrt{2}$ konservativ für $h = 2\sqrt{2}$ dissipativ für $h < 2\sqrt{2}$

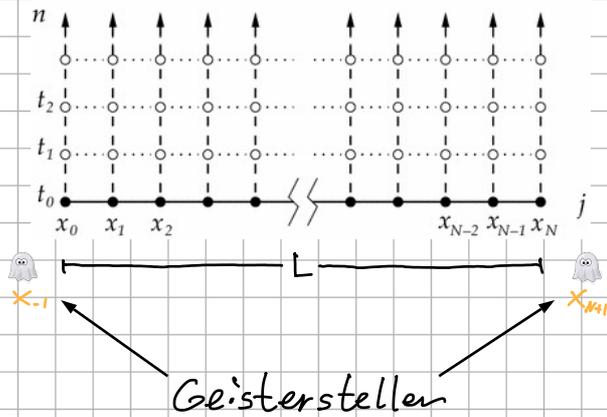
# partielle DGL (PDE)

## Problemstellung:

Wie entwickelt sich Temperaturverteilung  $u(t, x)$  in Stab der Länge  $L$ , wenn am Anfang ( $t=0$ ) die Temperatur gemäss  $\phi(x)$  verteilt ist

→ (vereinfachte) Wärmeleitungsgleichung:

$$\frac{\partial}{\partial t} u(t, x) = \frac{\partial^2}{\partial x^2} u(t, x), \text{ kurz: } \dot{u} = u''$$



## Diskretisierung:

Schrittweite in Zeit:  $t = n \Delta t$

Schrittweite im Ort:  $x = jh$ ,  $j = 0, \dots, N$  mit  $h = L/N$

→ manchmal zusätzlich noch Einführung von Geisterstellen  $x_{-1}, x_{N+1}$  notwendig

$u(t, x) = u(n \Delta t, jh) \approx U_{n,j}$  Temperatur bei Ort  $x_j$  zur Zeit  $t_n$

Anfangsbedingungen:  $\phi(x) \hat{=} \text{Temperaturverteilung des Stabes zur Zeit } t=0$

## Randbedingungen:

Dirichlet → Werte (durch Funktion) vorgegeben:  $u(t, 0) = L(t)$ ,  $u(t, L) = r(t)$

Neumann → (örtliche) Ableitung vorgegeben  $u'(t, 0) = g_l(t)$ ,  $u'(t, L) = g_r(t)$

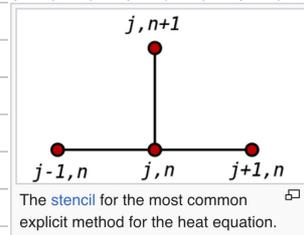
isotherme Randbedingungen: Dirichlet mit konstanten Werten:  $L(t), r(t) \text{ const.}$

adiabatische Randbedingungen: Neumann mit  $g_l(t) = g_r(t) = 0$

## FTCS (forward time, central space) explizites Verfahren

$$\dot{u} \approx \frac{U_{n+1,j} - U_{n,j}}{\Delta t}$$

$$u'' \approx \frac{U_{n,j+1} - 2U_{n,j} + U_{n,j-1}}{h^2}$$



$$\dot{u} = u'' \rightsquigarrow \frac{U_{n+1,j} - U_{n,j}}{\Delta t} = \frac{U_{n,j+1} - 2U_{n,j} + U_{n,j-1}}{h^2} \text{ Differenzengleichung}$$

nach  $U_{n+1,j}$  aufgelöst:

$$U_{n+1,j} = U_{n,j} + \frac{U_{n,j+1} - 2U_{n,j} + U_{n,j-1}}{h^2} \Delta t \text{ Temperatur bei Ort } x_j \text{ zur Zeit } t_{n+1}$$

→ für innere Punkte hat man also eine explizite Berechnungsvorschrift

→ äußeren beide Punkte → durch Randbedingungen gegeben

Voraussetzung für Stabilität:  $\frac{\Delta t}{h^2} \leq \frac{1}{2}$

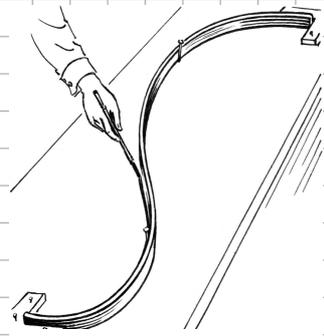
→ bei hoher Ortsauflösung ( $h$  klein,  $h^2$  sehr klein)  $\Rightarrow$  sehr viele Zeitschritte nötig  
→ oft nicht brauchbar (außer mit guter Grafikkarte)

Alternative: Crank-Nicolson-Verfahren → keine Stabilitätsbedingung



# Splines

gegeben:  $n+1$  Stützstellen  $(x_i, y_i), i=0, \dots, n$



Wortherkunft: Der Begriff stammt aus dem Schiffbau: eine lange dünne Latte (Straklatte, englisch spline), die an einzelnen Punkten durch Molche fixiert wird, biegt sich genau wie ein kubischer Spline. Dabei wird die Spannenergie minimal.

Spline  $k$ ten Grades:

Funktion welche stückweise aus Polynomen  $k$ ten Grades zusammengesetzt ist und durch vorgegebenen Stützstellen («Knoten») geht

$$S(x) = \begin{cases} S_1(x), & x \in [x_0, x_1] \\ S_2(x), & x \in [x_1, x_2] \\ \vdots \\ S_n(x), & x \in [x_{n-1}, x_n] \end{cases}$$

am Knoten müssen  $k-1$  Ableitungen übereinstimmen;

- bei  $k=1$  müssen die Knoten einfach (mit Massstab) verbunden werden
- bei  $k=3$  (kubische Splines) müssen an den Knoten auch die Steigung und Krümmung übereinstimmen

## stückweise lineare Interpolation ( $k=1$ )

$$S_i(x) = \frac{x_i - x}{x_i - x_{i-1}} y_{i-1} + \frac{x - x_{i-1}}{x_i - x_{i-1}} y_i \quad \text{wobei } i=1, 2, \dots, n$$

## kubische Splines ( $k=3$ )

$$S_i(x) = \frac{1}{6} \frac{(x_i - x)^3}{h_i} M_{i-1} + \frac{1}{6} \frac{(x - x_{i-1})^3}{h_i} M_i + \underbrace{\left( \frac{y_i - y_{i-1}}{h_i} - \frac{h_i}{6} (M_i - M_{i-1}) \right)}_{c_i} (x - x_{i-1}) + \underbrace{y_{i-1} - \frac{h_i^2}{6} M_{i-1}}_{d_i}$$

$$= \frac{1}{6} \frac{(x_i - x)^3}{h_i} M_{i-1} + \frac{1}{6} \frac{(x - x_{i-1})^3}{h_i} M_i + c_i (x - x_{i-1}) + d_i \quad \text{wobei } i=1, 2, \dots, n$$

mit  $c_i = \frac{y_i - y_{i-1}}{h_i} - \frac{h_i}{6} (M_i - M_{i-1}), d_i = y_{i-1} - \frac{h_i^2}{6} M_{i-1}, h_i = x_i - x_{i-1}$

für innere Momente  $M_i, i \in \{1, 2, \dots, n-1\}$ , müssen folgende  $n-1$  Gleichungen erfüllt sein:

$$\frac{h_i}{6} M_{i-1} + \frac{h_i + h_{i+1}}{3} M_i + \frac{h_{i+1}}{6} M_{i+1} = \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i}$$

aus Randbedingungen erhält man 2 weitere Gleichungen mit welchen man die äusseren beiden Momente  $M_0$  und  $M_n$  eliminieren kann

man erhält dann ein LGS für die inneren Momente:

$$A \vec{M} = \vec{b} \quad \text{mit } \vec{M} := (M_1, \dots, M_{n-1})^T \in \mathbb{R}^{n-1}, A \in \mathbb{R}^{(n-1) \times (n-1)}, \vec{b} \in \mathbb{R}^{n-1}$$

## Randbedingungen

Annahme:

$$a = x_0 < x_1 < \dots < x_{n-1} < x_n = b$$

mit  $x_i = a + i \cdot h$   
 $\rightarrow$  alle  $x_i$  denselben Abstand  $h$

$A \vec{M} = \vec{b}$  für verschiedene Randbedingungen

$\hookrightarrow$  nächste Seite

