

# §1 Quadratur

## §1.1. Motivation

Bsp  $\begin{cases} \dot{y} = f(t, y) \\ y(t_0) = y_0 \end{cases} \quad \left| \int_{t_0}^T dt \Rightarrow y(T) = y(t_0) + \int_{t_0}^T f(t, y(t)) dt \right.$

Energieverlust per Flächeneinheit:

$$Q = \int_0^{\infty} E(\lambda, T) d\lambda = \sigma T^4$$

Bsp Was ist die Periode eines Pendulums?

$$T = 4 \sqrt{\frac{l}{g}} K\left(\sin \frac{\alpha_0}{2}\right) \quad \text{wobei}$$

$$K(a) = \int_0^{\pi/2} \frac{1}{\sqrt{1 - a^2 \sin^2 s}} ds$$

Bsp Plank's Gesetz für die Strahlung eines schwarzen Körpers

$T$  = Temperatur in K

$\lambda$  = in  $\mu\text{m}$

$$E(\lambda, T) = \frac{c_1}{\lambda^5 \left( e^{\frac{c_2}{\lambda T}} - 1 \right)}$$

$\sigma$  = Stefan-Boltzmann-Konstante  
Ziel: berechne  $\sigma$ !

$$J = \int_a^b f(x) dx \approx Q(f, a, b) = \sum_{j=1}^n \omega_j f(x_j)$$

Gewichte Knoten

Ziel: Welche Knoten und Gewichte sollen wir wählen um den Fehler  $|J - Q|$  klein

mit möglichst wenige Funktionsauswertungen zu haben.

IDEA:  $f \approx$  einfache Funktion, dessen Integral leicht/analytisch berechenbar ist

z.B:  $f \approx$  Polynom  $\alpha_0 + \alpha_1 x + \dots + \alpha_n x^n$   
oder  $f \approx$  trigonometrischen Polynom.  
 $f \approx \alpha_0 + \alpha_1 e^{i1x} + \alpha_2 e^{i2x} + \dots$   
 $+ \alpha_{-1} e^{-i1x} + \alpha_{-2} e^{-i2x} + \dots$

Gegeben Knoten  $x_0, x_1, \dots, x_n$  kann man  $\alpha_0, \alpha_1, \dots, \alpha_n$  berechnen, so dass

$P_n(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_n x^n$  erfüllt

$P_n(x_j) = f(x_j)$  für  $j = 0, 1, 2, \dots, n$

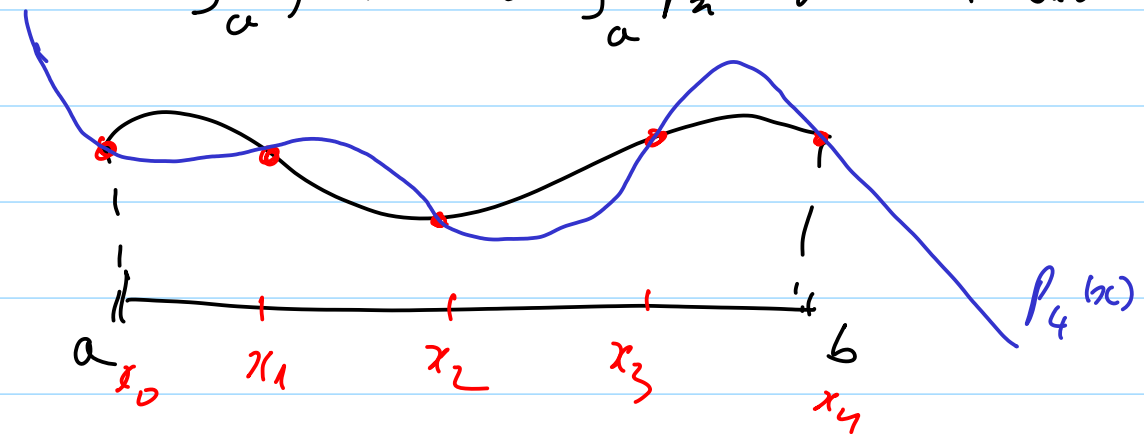
Im Prinzip, das ist das LGS



$$\begin{bmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & & x_1^n \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^n \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} f(x_0) \\ \vdots \\ f(x_n) \end{bmatrix}$$

Somit  $f(x) \approx P_n(x)$  und dann

$$\int_a^b f(x) dx \approx \int_a^b P_n(x) dx = \dots \text{ exakt}$$



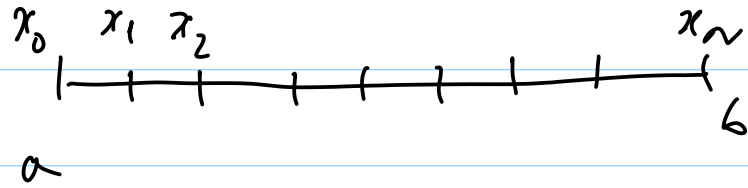
Man kann beweisen:

$$f \in C^n[a, b] \Rightarrow |J - Q| \leq \frac{1}{n!} (b-a)^{n+1} \max_{z \in [a, b]} |f^{(n)}(z)|$$

Glattheit ist wichtig

Länge des Intervalls ist wichtig!

IDEA: Zerlege  $\int_a^b f(x) dx = \sum_{k=0}^{N-1} \int_{x_k}^{x_{k+1}} f(x) dx$



$$x_k = x_0 + kh \quad \text{mit } h = \frac{b-a}{N} \text{ klein}$$

und wende die Quadraturformel auf jedes Intervall der Länge  $h$  (klein)

$\Rightarrow$  zusammengesetzte Quadraturformel.

$$\text{Fehler: } \left| \int_a^b f(x) dx - \sum_{k=0}^{N-1} Q(f, x_k, x_{k+1}) \right| \leq$$

$$\leq \sum_{k=0}^{N-1} \left| \int_{x_k}^{x_{k+1}} f(x) dx - Q(f, x_k, x_{k+1}) \right| \leq$$

$$\leq \sum_{k=0}^{N-1} \frac{1}{n!} \underbrace{(x_{k+1} - x_k)^{n+1}}_h \max_{z \in [x_k, x_{k+1}]} |f^{(n+1)}(z)| =$$

$$\sum_{k=0}^{N-1} \frac{h^{n+1}}{n!} \max_{z \in [x_k, x_{k+1}]} |f^{(n+1)}(z)| \leq C \cdot \frac{1}{n!} h^{n+1} \cdot N$$

$$h = \frac{b-a}{N} \Rightarrow N = \frac{b-a}{h}$$

$$\text{Somit: Fehler} \leq \frac{C}{n!} h^{n+1} \cdot \frac{(b-a)}{h} = \frac{(b-a)}{n!} C \underbrace{h^n}_{!}$$

!

Def Quadraturformel hat Ordnung  $n+1$   
wenn sie Polynome von Grad maximal  $n$   
exakt integriert.

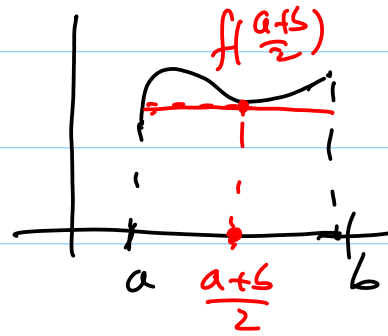
(das erste falsche Ergebniss:  $x^{n+1}$ )

$$\int_a^b p(x) dx = Q(p, a, b) \quad \text{für alle Polynome vom Grad } \leq n$$

und es gibt ein Polynom  $\tilde{p}$  vom Grad  $n+1$ :  
 $\int_a^b \tilde{p}(x) dx \neq Q(\tilde{p}, a, b).$

Bsp 1) Mittelpunktsregel

$$Q^M(f, a, b) = (b-a) f\left(\frac{a+b}{2}\right)$$



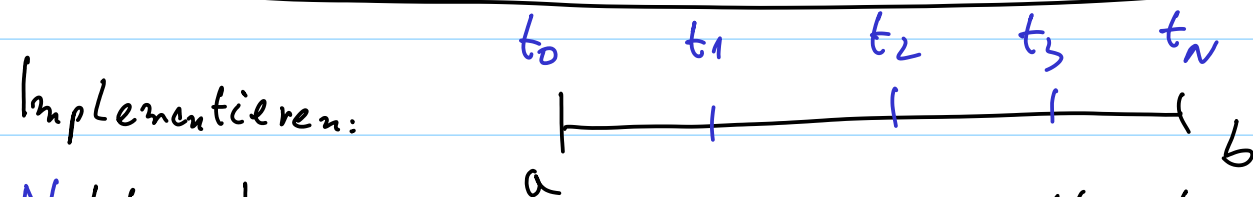
Bemerkung: 1) Polynome von Grad 0 und 1 werden exakt integriert (mittels MPR)

$\Rightarrow$  MPR hat Ordnung 2

Fehler:  $\frac{(b-a)^3}{24} f''(\xi)$  mit  $\xi \in [a, b]$

Bem 2) offene Quadraturformel:

Enden von  $[a, b]$  sind keine Knoten



$N$  Intervalle

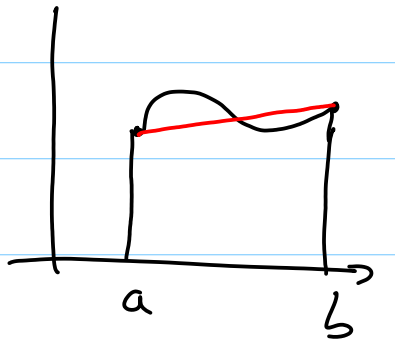
Argument:  $f, a, b, N;$

$$\int_a^b f(x) dx = \sum_{j=1}^N \int_{t_{j-1}}^{t_j} f(x) dx$$

Bsp 2) Trapezregel:

$$Q^T(f, a, b) = \frac{b-a}{2} f(a) + \frac{b-a}{2} f(b)$$

Gewichte  $\nearrow$  Knoten



$$= \sum_{j=1}^2 w_j f(x_j) \quad x_1 = a, x_2 = b, w_1 = w_2 = \frac{b-a}{2}$$

Bem 1)

Ordnung 2, Fehler:  $\frac{1}{12} (b-a)^3 f^{(2)}(\xi)$   
mit  $\xi \in [a, b]$

Bem 2) geschlossene Quadraturformel:

Enden von  $[a, b]$  sind Knoten.

$$\begin{aligned} \text{MPR: } \int_a^b f(x) dx &\approx \sum_{j=1}^N \underbrace{(t_j - t_{j-1})}_h f\left(\frac{t_{j-1} + t_j}{2}\right) \\ &= \frac{b-a}{N} \sum_{j=1}^N f(x_j) \quad \text{wobei } x_j = t_{j-1} + \frac{h}{2} \end{aligned}$$



Bsp3) Simpson Regel

$$Q^S(f, a, b) = \underbrace{\frac{b-a}{6}}_{w_1} \underbrace{f(a)}_{x_1} + \underbrace{\frac{b-a}{6} \cdot 4}_{w_2} \underbrace{f\left(\frac{a+b}{2}\right)}_{x_2} + \underbrace{\frac{b-a}{6}}_{w_3} \underbrace{f(b)}_{x_3}$$

$$= \sum_{j=1}^3 w_j f(x_j)$$

Bemerkung Ordnung 4; Fehler:  $\frac{1}{90} \left(\frac{b-a}{2}\right)^5 f^{(4)}(z)$   
mit  $z \in [a, b]$ .

Input:  $f, a, b, N$

$$h = \frac{b-a}{N}$$

$$\underline{t} = [t_0, t_1, \dots, t_N] \text{ mit } t_0 = a, t_N = b$$

$$(\text{linspace}(a, b, N+1))$$

$$\underline{x} = \underline{t}[:, -1] + \frac{h}{2}$$

$$\text{Output: } h \sum(f(\underline{x}))$$

$$\text{TR: } \int_a^b f(x) dx = \sum_{j=1}^N \int_{t_{j-1}}^{t_j} f(x) dx \approx \sum_{j=1}^N \left( \frac{h}{2} f(t_{j-1}) + \frac{h}{2} f(t_j) \right)$$

$$= h \sum_{j=1}^N \left( \frac{1}{2} f(t_{j-1}) + \frac{1}{2} f(t_j) \right) =$$

$$= \underbrace{\frac{h}{2} f(t_0) + \frac{h}{2} f(t_1)}_{j=1} + \underbrace{\frac{h}{2} f(t_1) + \frac{h}{2} f(t_2)}_{j=2} + \dots + \underbrace{\frac{h}{2} f(t_{N-1}) + \frac{h}{2} f(t_N)}_{j=N}$$

$$= \frac{h}{2} f(t_0) + h f(t_1) + h f(t_2) + \dots + h f(t_{N-1}) + \frac{h}{2} f(t_N)$$

$$= \frac{h}{2} f(t_0) + h \sum_{j=1}^{N-1} f(t_j) + \frac{h}{2} f(t_N)$$

$$\underline{x} = \underline{t}[:, -1]$$

$$\text{Output: } \text{sum}(h f(\underline{x})) + \frac{h}{2} f(a) + \frac{h}{2} f(b)$$

Simpson: ähnlich

$$\sum_{j=1}^N \left( \frac{h}{6} f(t_{j-1}) + \frac{h}{6} \cdot 4 f\left(\frac{t_{j-1} + t_j}{2}\right) + \frac{h}{6} f(t_j) \right)$$



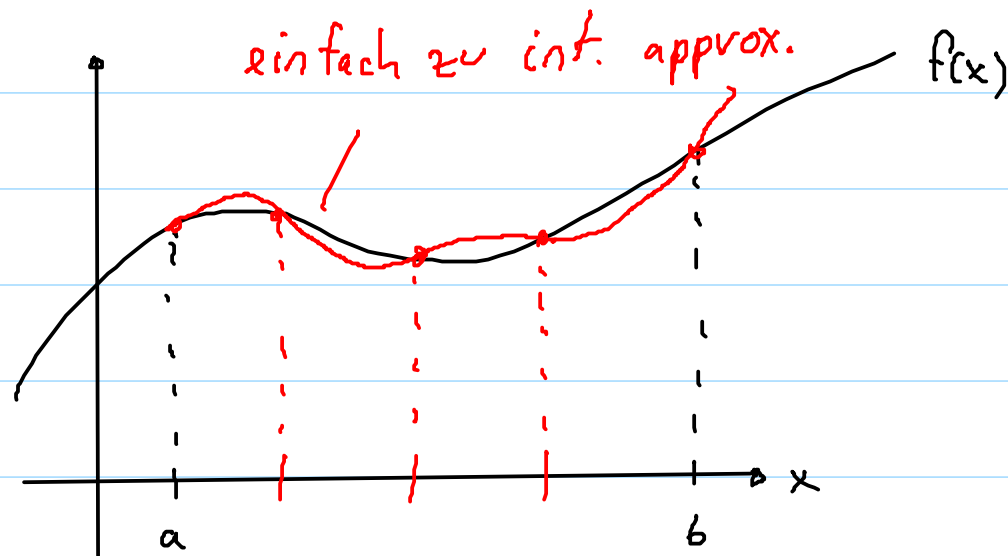
## Kurz Zusammenfassung

Ziel: Approx. von

$$Q(f) \approx \int_a^b f(x) dx = I(f)$$

Wozu: Oft leider nicht exakt berechenbar

Idee:



Einfache Fkt.: Polynome!

## Polynomiale Interpolation

Gez.  $n+1$  paarweise verschiedene  
Stützstellen/Knoten  $x_0, x_1, \dots, x_n$  und  
 dazugehörige Stützwerte  $y_0, y_1, \dots, y_n$   
 finde das Polynom  $n$ -ten Grades  
 , Interp.-Polynom (IP)

$$p_n(x) = \alpha_0 + \alpha_1 \cdot x + \dots + \alpha_n x^n \in \mathbb{P}_n$$

welches die Interpolationsbedingungen (IB)  
 erfüllt

$$p_n(x_j) = y_j \quad (j=0, \dots, n)$$

Die  $n+1$  Koeff.  $\alpha_0, \dots, \alpha_n$  ergeben sich  
 durch Lösen eines linearen Gleichungssystems (LGS)

$$\begin{pmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & & x_1^n \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^n \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix}$$

Das IP lässt sich auch direkt mit der sog. Lagrange'schen Interpolationsformel (LF)

$$p_n(x) = \sum_{j=0}^n y_j \cdot L_j^n(x)$$

wobei

$$L_j^n(x) = \prod_{\substack{i=0 \\ i \neq j}}^n \frac{x - x_i}{x_j - x_i} \quad (j=0, 1, \dots, n)$$

(LP)

die Lagrange-Polynome sind

Die LP haben folgende Eigenschaften

(LP1)  $L_j^n(x)$  sind Polynome von Grad  $n$

$$(LP2) \quad L_j^n(x_k) = \delta_{jk} = \begin{cases} 1 & j=k \\ 0 & \text{sonst} \end{cases}$$

Wegen (LP2) funktioniert

$$\begin{aligned} p_n(x_i) &= \sum_{j=0}^n y_j \cdot L_j^n(x_i) \\ &= 0 + \dots + 0 + y_i \cdot \underbrace{L_i^n(x_i)}_1 + 0 \dots \\ &= y_i \end{aligned}$$

IB erfüllt ✓.

Soweit waren die Stützpunkte beliebig.

Sei nun  $f: I=[a,b] \rightarrow \mathbb{R}$  und wir wollen mit einer IP approx.

$$p_n[f/x_0, \dots, x_n](x) \in \mathcal{P}_n$$

Dieses erfüllt die IB

$$p_n[f/x_0, \dots, x_n](x_j) = f(x_j) \quad (j=0, 1, \dots, n)$$

Numerische Integration = Quadratur

Ziel: Approx. von  $I(f) = \int_a^b f(x) dx$

Idee: Verwende polynomiale Interpolation um  $f(x)$  zu approx. und integriere

$$p[f/x_0, \dots, x_n]$$

(... Poly. sind so einfach!)

Def.: Eine endliche Rechenvorschrift der Form

$$Q(f) = \sum_{j=0}^n w_j \cdot f(x_j)$$

zur Approx. von  $I(f)$  nennt man Quadraturregel (QR) oder Quad.-formel

Die  $x_j \in [a,b]$  nennt man (Quadratur) Knoten oder Integrationsstützstelle

und die  $w_j$  (Quadratur) Gewichte

QR können nun ganz einfach hergeleitet werden.

Seien  $x_0, x_1, \dots, x_n \in I$  uns gegebene Knoten.

Dann ist das IP

$$p[f/x_0, \dots, x_n] = \sum_{j=0}^n f(x_j) \cdot \mathcal{L}_j^n(x)$$

Da  $p[f/x_0, \dots, x_n] \approx f(x)$

können wir versuchen

$$\begin{aligned}
 \int_a^b f(x) dx &\approx \int_a^b p[f|x_0, \dots, x_n] dx \\
 &= \int_a^b \sum_{j=0}^n f(x_j) \cdot L_j^{\wedge}(x) dx \\
 &= \sum_{j=0}^n \int_a^b f(x_j) \cdot L_j^{\wedge}(x) dx \\
 &= \sum_{j=0}^n f(x_j) \cdot \underbrace{\int_a^b L_j^{\wedge}(x) dx}_{\text{Gewichte!}} \\
 &= \sum_{j=0}^n f(x_j) \cdot w_j
 \end{aligned}$$

Also die Gewichte sind gegeben durch

$$w_j = \int_a^b L_j^{\wedge}(x) dx$$

$w_j$  hängen nicht von  $f$  ab!

Geg. Knoten: berechne  $w_j$  ein für alle Mal und tabelliere

Wichtige Bsp.

Bsp.: Mittelpunktsregel (MPR)

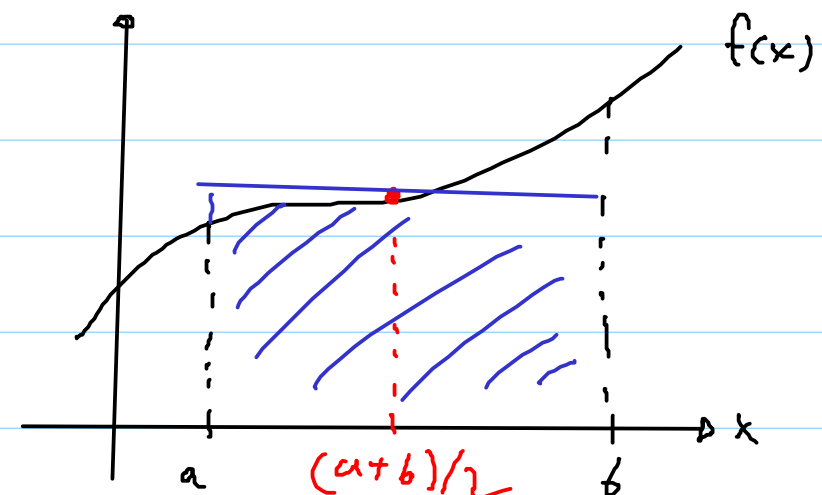
Knoten:  $x_0 = \frac{a+b}{2}$

LP :  $L_0^{\circ}(x) = 1$

Gewichte:  $w_0 = \int_a^b L_0^{\circ}(x) dx = b-a$

Damit

$$Q_0(f) = (b-a) \cdot f\left(\frac{a+b}{2}\right)$$



Bsp.: Trapezregel (TR)

Knoten:  $x_0 = a, x_1 = b$

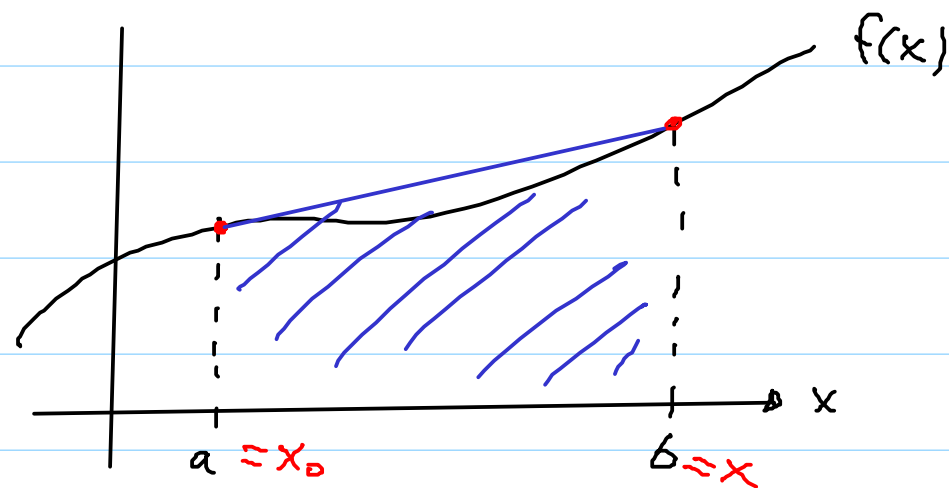
LP : ...

Gewichte:  $w_0 = \frac{b-a}{2}$

$w_1 = \frac{b-a}{2}$

Damit:

$$Q_1(f) = \frac{b-a}{2} (f(a) + f(b))$$



Bsp.: Simpson-Regel (SR)

Knoten:  $x_0 = a, x_1 = \frac{a+b}{2}, x_2 = b$

LP : ...

Gewichte:  $w_0 = \frac{b-a}{6}$

$w_1 = \frac{4(b-a)}{6}$

$w_2 = \frac{b-a}{6}$

Damit

$$Q_2(f) = \frac{b-a}{6} \left( f(a) + 4 \cdot f\left(\frac{a+b}{2}\right) + f(b) \right)$$

Die NQR, TR und SR sind Teil einer grossen Familie: Newton-Cotes  $QK_n$  (NC)

Bei NC wählt man die Knoten äquidistant verteilt im Integrations-Intervall

Bem.: (i) NPR, TR und SR sind sehr populär

(ii) NC QRn mit  $n > 6$  werden numerisch unbrauchbar (weil negative Gewicht auftauchen)

### Quadratfehler

Def.:  $E(f) = |Q(f) - I(f)|$  Quadratfehler (QF)

Def.: Eine QR hat Genauigkeitsgrad (GG)  $q \in \mathbb{N}$  falls sie alle Polynome bis zu und mit Grad  $q$  exakt integriert und  $q$  der GröÖte Wert mit dieser Eigenschaft (Manchmal Exaktheitsgrad)

Def.: Die Ordnung einer QR ist  $s = q+1$

Es genügt den GG für die Monome zu berechnen:

$$Q(x^k) = I(x^k) \quad k = 0, 1, \dots, q$$

$$Q(x^{q+1}) \neq I(x^{q+1})$$

Wegen der Linearität von  $I$  und  $Q$ , und die Monome bilden eine Basis von  $\mathcal{P}$

Weiter ist es bequem auf einem Referenzintervall (RI)  $[-1, 1]$  zu arbeiten

Jedes Intervall  $[a, b]$  lässt sich durch eine Substitution auf das RI transformieren



$$x = \frac{b-a}{2} t + \frac{a+b}{2} \quad \text{mit } t \in [-1, 1]$$

Transformiere

$$\begin{aligned} I(f) &= \int_a^b f(x) dx \\ &= \int_{-1}^1 f\left(\frac{b-a}{2} t + \frac{a+b}{2}\right) \frac{b-a}{2} dt \\ &= \frac{b-a}{2} \int_{-1}^1 \hat{f}(t) dt \end{aligned}$$

$$\approx \frac{b-a}{2} \sum_{j=0}^n \hat{w}_j \hat{f}(t_j)$$

Gewichte      Knoten auf RI

Übung: bestimme  $q$  (und  $s$ ) für  $\mathcal{NPR}$ ,  $\mathcal{TR}$  und  $\mathcal{SR}$

Für den QF lässt sich zeigen

$$E(f) \leq \frac{\|f^{(q+1)}\|_{\infty}}{(q+1)!} (b-a)^{q+2} \quad (\text{QF1})$$

$s$        $s+1$  - Ordnung

↑      ↑  
"Glattheit" von  $f$       Intervall-Breite

Zusammengesetzte QR<sub>n</sub>

Idee: Zerlege  $I = [a, b]$  in Teil-Intervalle und wende QR auf jedes an  
Mit gleich grossen Teil-Intervallen

$\tau I$ 

$$I_j = [x_{j-1}, x_j] \quad (j=1, \dots, N)$$

#  $\tau I$

und

$$x_j = a + h \cdot j \quad (j=0, \dots, N)$$

$$h = \frac{b-a}{N}$$

Fehler?

$$E^N(f) = |I(f) - Q_N^N(f)|$$

$$= \left| \sum_{j=1}^N I(f \text{ auf } I_j) - Q_N(f \text{ auf } I_j) \right|$$

 $\Delta$ -Ungl.

$$\leq \sum_{j=1}^N \underbrace{|I(f \text{ auf } I_j) - Q_N(f \text{ auf } I_j)|}_{E(f \text{ auf } I_j)}$$

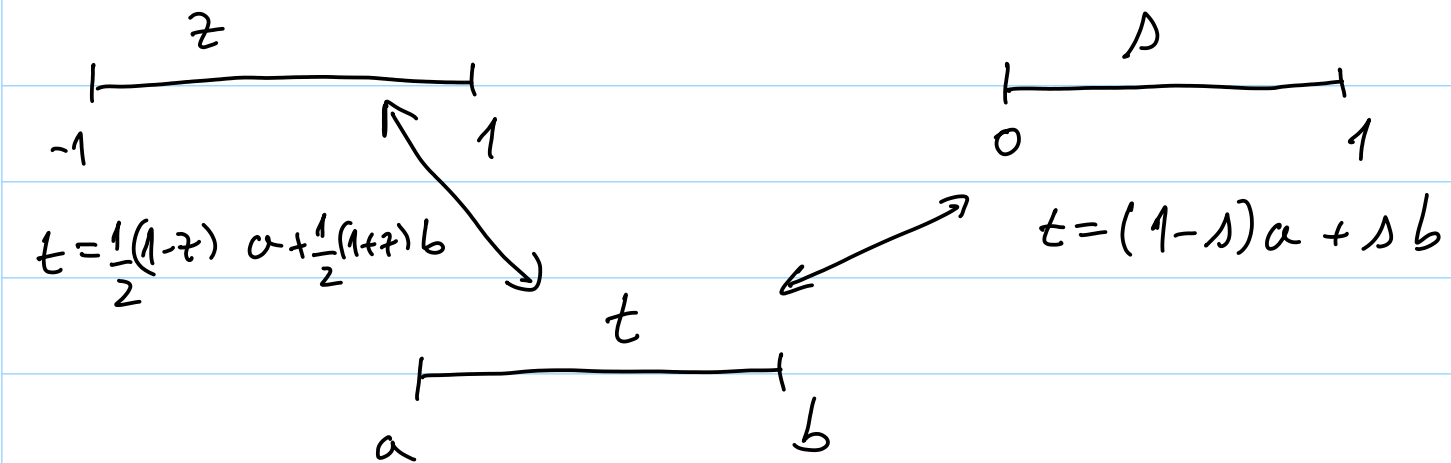
 $E(f \text{ auf } I_j)$ 

$$(QFA) \leq \sum_{j=1}^N \frac{\max_{x \in I_j} |f^{(q+1)}(x)|}{(q+1)!} \underbrace{(x_j - x_{j-1})}_{h}^{q+2}$$

$$\leq N \cdot \frac{\|f^{(q+1)}\|_{\infty}}{(q+1)!} h^{q+2}$$

$$= \frac{\|f^{(q+1)}\|_{\infty}}{(q+1)!} (b-a) \cdot h^{q+1}$$

## §1.2. Referenzintervalle und symmetrische QF



$$\int_a^b f(t) dt = \frac{b-a}{2} \int_{-1}^1 \hat{f}(z) dz \approx \frac{b-a}{2} \sum_{j=1}^n \hat{\omega}_j \hat{f}(\hat{c}_j)$$

mit  $\hat{f}(z) = f\left(\frac{1}{2}(1-z)a + \frac{1}{2}(1+z)b\right)$   $[-1, 1]$

Theorem Die Quadraturordnung einer symmetrischen QF ist gerade.

Beweis

Annahme, QF exakt für Polynome vom Grad  $(\max) 2m-2$ . Nehme  $f(x) = ax^{2m-1} + g(x)$   
 $g(x) = \text{Polynom vom Grad } 2m-2$ .

$$\int_{-1}^1 f(x) dx = a \underbrace{\int_{-1}^1 x^{2m-1} dx}_0 + \int_{-1}^1 g(x) dx = \sum_{k=1}^n \hat{\omega}_k g(\hat{c}_k)$$

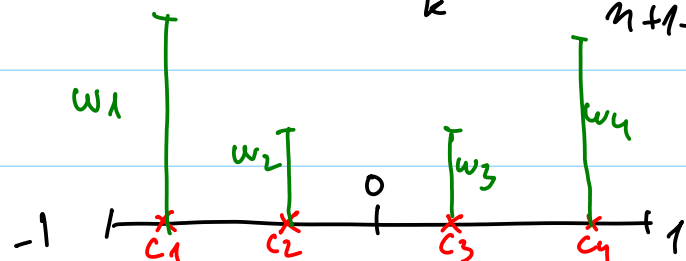
$$Q(f, -1, 1) = a \sum_{k=1}^n \hat{c}_k^{2m-1} \cdot \hat{\omega}_k + \sum_{k=1}^n \hat{\omega}_k g(\hat{c}_k)$$

Es reicht zu zeigen  $\sum_{k=1}^n \hat{\omega}_k \hat{c}_k^{2m-1} = 0$ .

$$\sum_{k=1}^n \hat{\omega}_k \hat{c}_k^{2m-1} = \sum_{k=1}^n \hat{\omega}_{n+1-k} (-\hat{c}_{n+1-k})^{2m-1} = - \sum_{k=1}^n \hat{\omega}_{n+1-k} \hat{c}_{n+1-k}^{2m-1} = - \sum_{j=n+1-k}^n \hat{\omega}_j \hat{c}_j^{2m-1} = 0$$

$j = n+1-k$

Definition QF auf  $[-1, 1]$  heißt symmetrisch falls  $c_k = -c_{n+1-k}$ ,  $w_k = w_{n+1-k}$



$$= - \sum_{j=1}^n \hat{\omega}_j \hat{c}_j^{2n-1} \Rightarrow \sum_{k=1}^n \hat{\omega}_k \hat{c}_k^{2n-1} = 0.$$

### §1.3 Fehler für Quadratur auf $[0,1]$

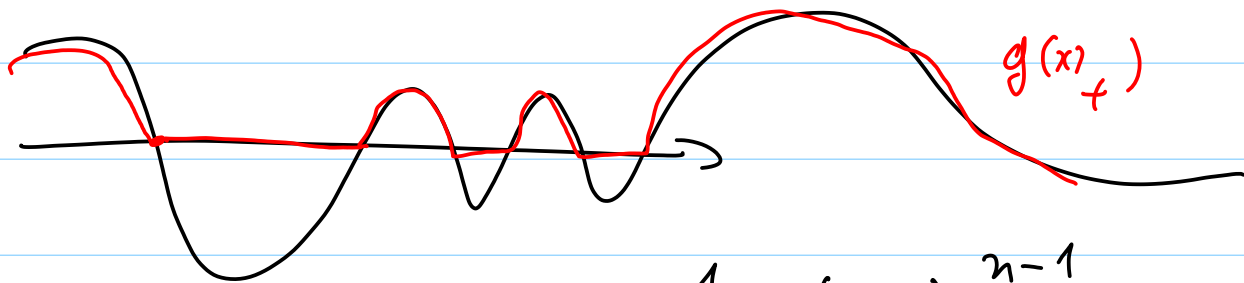
Fehler  $E(g) = \int_0^1 g(t) dt - \sum_{j=1}^n b_j g(c_j)$

$\downarrow$  Knoten  $[0,1]$   
 $\downarrow$  Gewichte  
 linear in  $g$ .

$$E(Ag + Bf) = A E(g) + B E(f)$$

$A, B \in \mathbb{R}, \quad g, f \text{ Funktionen.}$

das positive Teil  $g(x)_+ = \begin{cases} g(x) & \text{für } x \text{ so dass } g(x) > 0 \\ 0, & \text{sonst} \end{cases}$



Pearo-Kern:  $\alpha(z, t) = \frac{1}{(n-1)!} (t-z)_+^{n-1}$

$$K_n(z) = E(\alpha(z, \cdot)) = \int_0^1 \frac{1}{(n-1)!} (t-z)_+^{n-1} dt - \sum_{j=1}^n b_j \frac{(c_j-z)_+^{n-1}}{(n-1)!} = \frac{(1-z)^n}{n!} - \sum_{j=1}^n b_j \frac{(c_j-z)_+^{n-1}}{(n-1)!}$$

Theorem Sei  $Q$  eine QF mit  $Q$ -Ordnung  $n$  und sei  $g$   $n$ -mal stetig diff-bar.

Dann

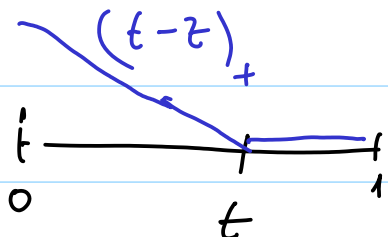
$$E(g) = \int_0^1 K_n(z) g^{(n)}(z) dz.$$

Beweis Taylor:

$$g(t) = g(0) + \dots + \frac{t^{n-1}}{(n-1)!} g^{(n-1)}(0) + \int_0^t \frac{(t-z)^{n-1}}{(n-1)!} g^{(n)}(z) dz$$

$q(t)$  Grad  $n-1$

$$\int_0^1 \frac{(t-z)_+^{n-1}}{(n-1)!} g^{(n)}(z) dz$$



$$E \text{ linear} \Rightarrow E(g) = E(g) + \int_0^1 \underbrace{E(\alpha(z, \cdot))}_{K_n(z)} g^{(n)}(z) dz$$

Wenden wir den Satz auf

$$g(t) = f(x_0 + th)$$

$$\int_{x_0}^{x_0+h} f(x) dx - h \sum_{j=1}^n b_j f(x_0 + c_j h) =$$

$$= h \int_0^1 g(t) dt - h \sum_{j=1}^n b_j g(c_j) = h E(g) =$$

$$= h h^n \int_0^1 k_n(z) f^{(n)}(x_0 + hz) dz = h \int_0^1 k_n(z) f^{(n)}(x_0 + hz) dz$$

Variablewechsel und n-te Ableitung

$$g(t) = f(x_0 + ht) \Rightarrow g'(t) = h f'(x_0 + ht) \Rightarrow \dots$$

$$g^{(n)}(t) = h^n f^{(n)}(x_0 + ht)$$

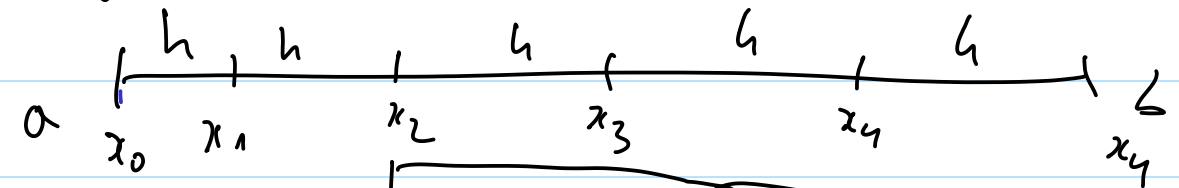
Zusammengesetzte Quadraturformel

$$|E(f)| = \left| \int_a^b f(x) dx - \sum_{k=1}^N \sum_{j=1}^n b_j f(x_{k-1} + c_j h) \right| \leq$$

$$\leq c \cdot h^n \max_{x \in [a, b]} |f^{(n)}(x)|$$

Glattheit von  $f$  wichtig

Ordnung der lokalen QF.

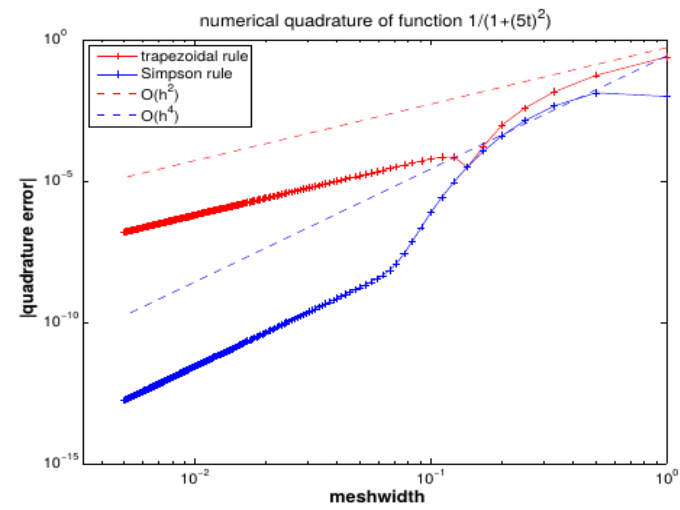


$$| \quad | \leq h^{n+1} \int_0^1 |k_n(z)| dz \cdot \max_{x \in [x_0, x_0+h]} |f^{(n)}(x)|$$

konstant

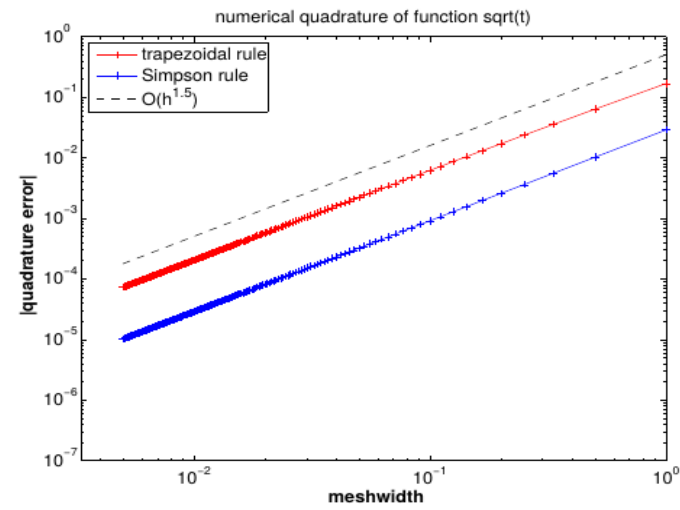
$$\text{MPR, TR} = \frac{1}{12}$$

$$\text{Simpson: } \frac{1}{2880}$$



Quadratur-Fehler für  $f_1(t) := \frac{1}{1+(5t)^2}$

Ordnung 2 für die Trapez-Regel und  
Ordnung 4 für Simpson-Regel



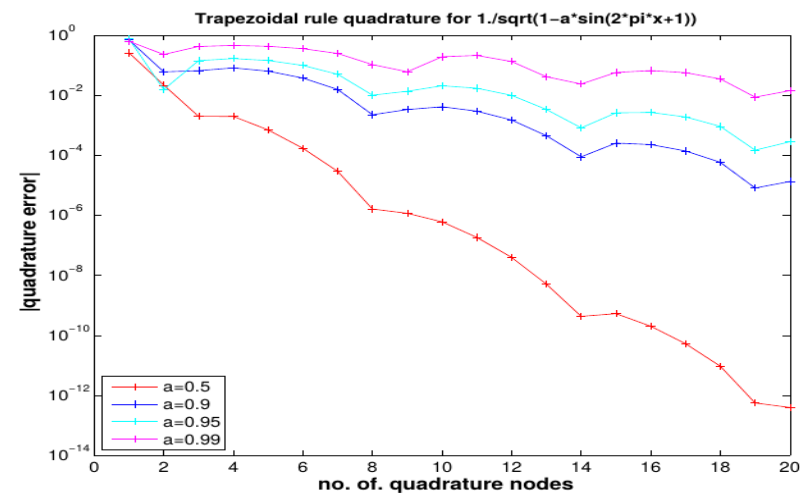
Quadratur-Fehler für  $f_2(t) := \sqrt{t}$

Unglattheit von  $f$  beschränkt die  
Konvergenzordnung

Die Trapez-Regel auf äquidistante Punkte hat Ordnung 2, aber angewendet auf den 1-periodischen glatten (analytischen) Integrand

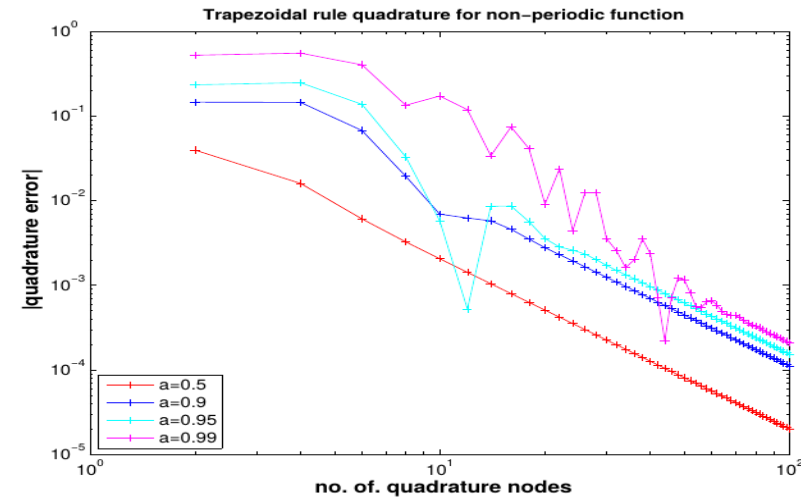
$$f(t) = \frac{1}{\sqrt{1 - a \sin(2\pi t - 1)}} \quad , \quad 0 < a < 1$$

erhalten wir folgende erstaunliche Ergebnisse (als “exakten Wert des Integrals” verwenden wir  $T_{500}$ ):



Quadratur-Fehler für  $T_n(f)$  auf  $[0, 1]$

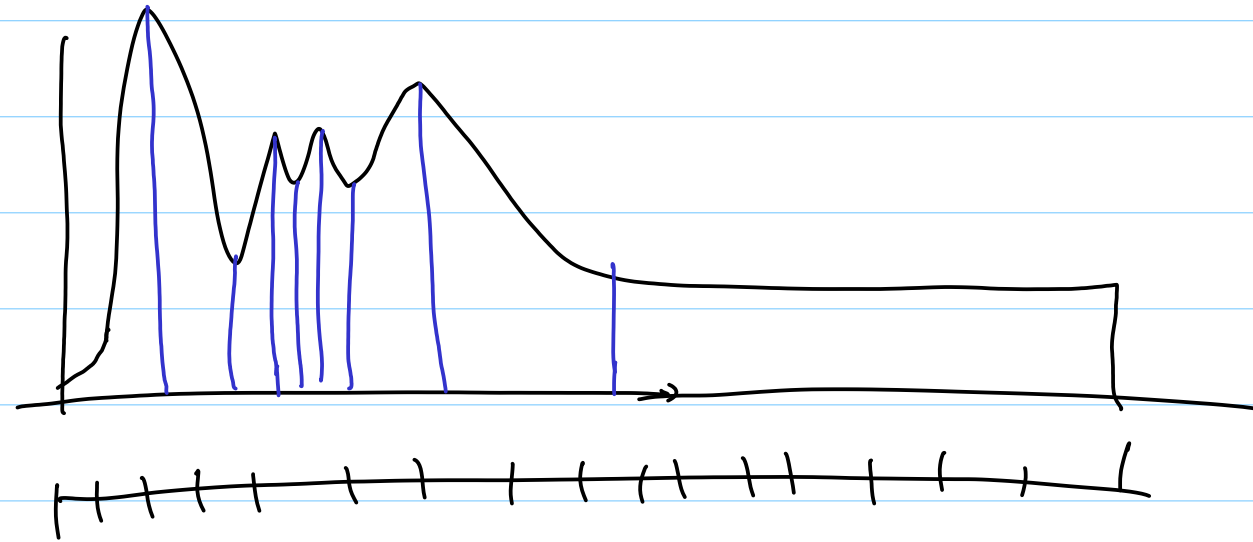
Exponentielle Konvergenz!



Quadratur-Fehler für  $T_n(f)$  auf  $[0, \frac{1}{2}]$

nur algebraische Konvergenz...

# §1.4. Adaptive Quadratur



Lokaler Fehler ~ Glattheit der Funktion.

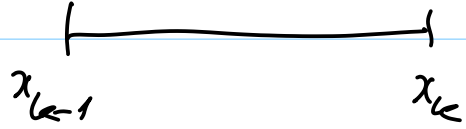
$$h^2 \max_{x \in [x_{k-1}, x_k]} |f^{(2)}(x)|$$
  
↓  
es lohnt  $h$  kleiner, dann wenn gross,  
sonst nicht.

Optimiere: Anzahl Funktionsauswertungen.

Intervalllänge      Gewichte      Knoten  $[0,1]$

$$\left| \int_a^b f(x) dx - \sum_{k=1}^N h_k \sum_{j=1}^n b_j f(x_{k-1} + c_j h_k) \right| \leq$$

$(b-a) \max_{k=1,2,\dots,N} \varepsilon_k$       wobei



lokale Fehler auf  $[x_{k-1}, x_k]$

$$\varepsilon_k = \left| \int_{x_{k-1}}^{x_k} f(x) dx - \sum_{j=1}^n b_j f(x_{k-1} + c_j h_k) \right|$$

IDEA: wähle  $h_k$  klein nur dort wo  $|f^{(2)}(x)|$  gross

Wo?

Wie schätze ich während der Rechnung  
den lokale Fehler, ohne  
weiter Informationen über  $f$ ?

$$\int_{x_{k-1}}^{x_k} f(x) dx = Q^T(f, x_{k-1}, x_k) + c h^3 \max_{z \in [x_{k-1}, x_k]} |f''(z)|$$

$$\int_{x_{k-1}}^{x_k} f(x) dx = Q^S(f, x_{k-1}, x_k) + c h^5 \max_{z \in [x_{k-1}, x_k]} |f^{IV}(z)|$$

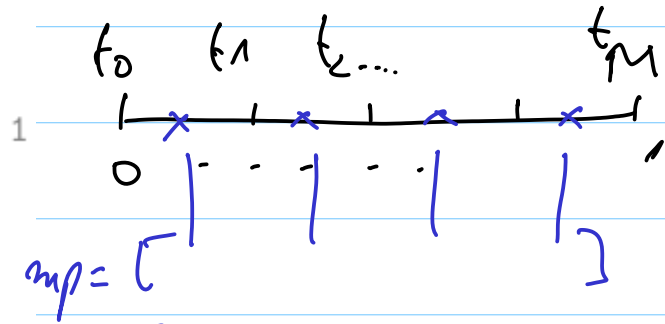
$$\left| \int_{x_{k-1}}^{x_k} f(x) dx - Q^T(f, x_{k-1}, x_k) \right| \approx$$

$$z^2 Q^S(f, x_{k-1}, x_k)$$

$|Q^S - Q^T|$  = Schätzung des lokalen Fehlers.

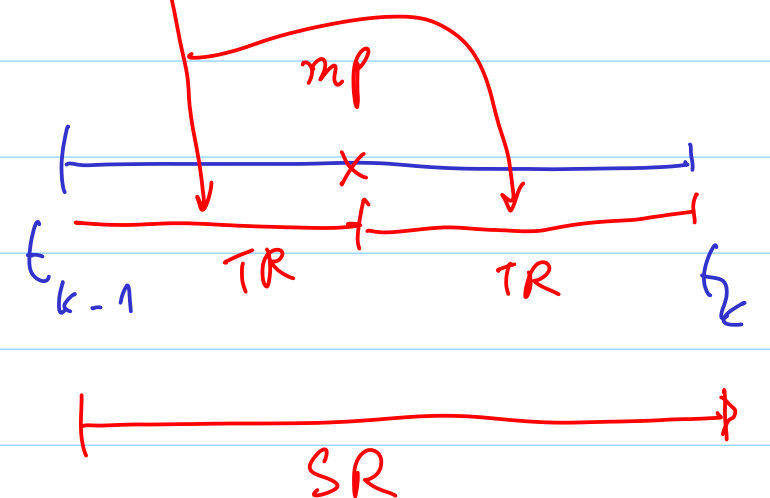
**IDEE:** verfeinere dort wo diese Schätzung gross ist!

```
if __name__ == '__main__':
    f = lambda x: exp(6*sin(2*pi*x))
    #f = lambda x: 1.0/(1e-4+x*x)
    M = arange(11.)/10 # 0, 0.1, ... 0.9, 1
    rtol = 1e-6; abstol = 1e-10
    I = adaptquad(f, M, rtol, abstol)
    exact, e = integrate.quad(f, M[0], M[-1])
    print 'adaptquad:', I, "exact:", exact
    print 'error:', abs(I-exact)
```



```
h = diff(M) # compute lengths of mesh intervals
mp = 0.5*( M[:-1]+M[1:] ) # compute midpoint positions
fx = f(M); fm = f(mp) # evaluate function at positions and
                        # midpoints
trp_loc = h*( fx[:-1]+2*fm+fx[1:] )/4 # local trapezoid rule
simp_loc = h*( fx[:-1]+4*fm+fx[1:] )/6 # local simpson rule
```

lokal!





```

I = sum(simp_loc) # use simpson rule value as
intermediate approximation for integral value
est_loc = abs(simp_loc - trp_loc) # difference of values obtained from
local composite trapezoidal rule and local simpson rule is used as an estimate
for the local quadrature error.
err_tot = sum(est_loc) # estimate for global error (sum
moduli of local error contributions)
# if estimated total error not below relative or absolute threshold, refine
mesh
if err_tot > rtol*abs(I) and err_tot > abstol:
    refcells = nonzero( est_loc > 0.9*sum(est_loc)/size(est_loc) )[0]
    I = adaptquad(f, sort(append(M, mp[refcells])), rtol, abstol) # add
midpoints of intervalls with large error contributions, recurse.
return I

```

$$\underset{\substack{\uparrow \\ \text{array.}}}{loc} > \frac{9}{10} \cdot \underset{\substack{\uparrow \\ \text{zahl}}}{tot} \cdot \frac{1}{\# \text{Intervalle}} \Rightarrow$$

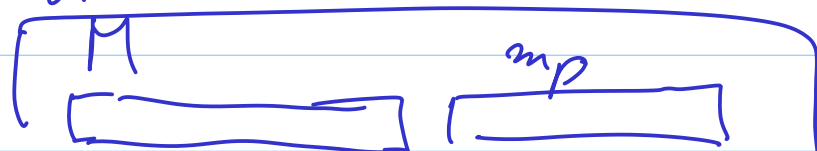
TT F F T F F T T  
 1 1 0 0 1 0 0 1 1  
 (true, false, ...)

nonzero  $\Rightarrow$  indices, wo das array keine 0 hat

$\hookrightarrow$  Intervalle, die zu verfeinern sind.

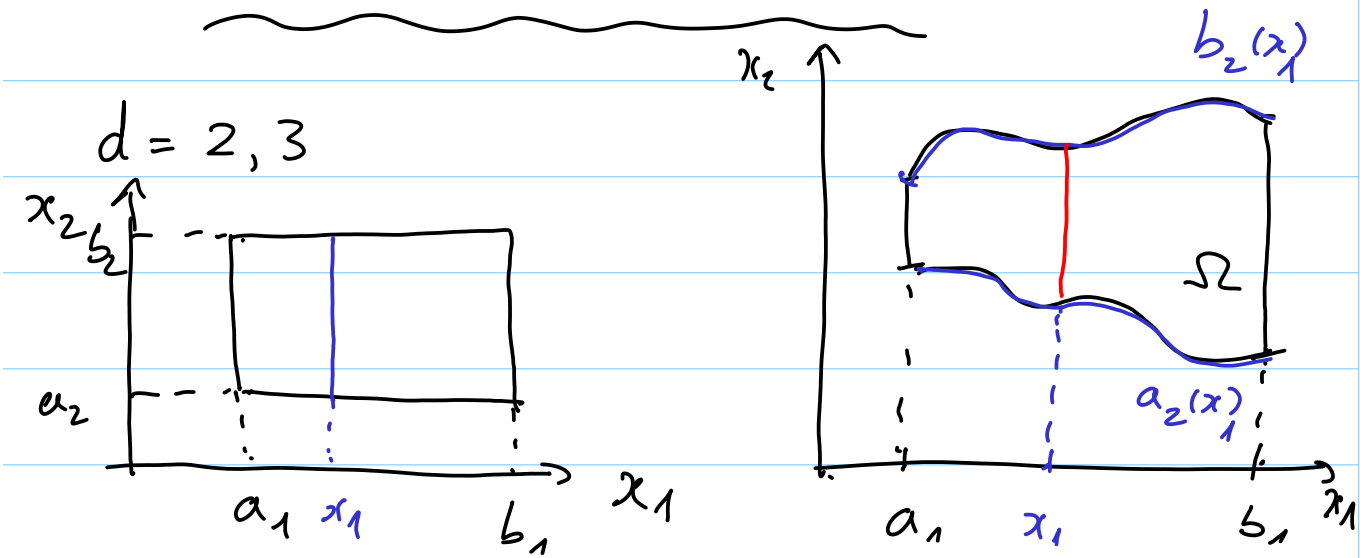
nehme die Mittelpunkte nur davon

append



, Sortiere  $\Rightarrow$  die neue Liste M.

## §1.5 Quadratur in $\mathbb{R}^d$



$$I = \int_{\Omega} f(x_2, x_1) dx_2 dx_1 = \int_{a_1}^{b_1} \int_{a_2(x_1)}^{b_2(x_1)} f(x_2, x_1) dx_2 dx_1$$

$$= \int_{a_1}^{b_1} F(x_1) dx_1$$

$$\approx \sum_{j_1=1}^{N_1} F(c_{j_1}^1) w_{j_1}^1$$

$\downarrow$  Knoten in  $Ox_1$ -Richtung  $\in [a_1, b_1]$

$$= \sum_{j_1=1}^{N_1} \int_{a_2(c_{j_1}^1)}^{b_2(c_{j_1}^1)} f(x_2, c_{j_1}^1) dx_2 \cdot w_{j_1}^1 \approx$$

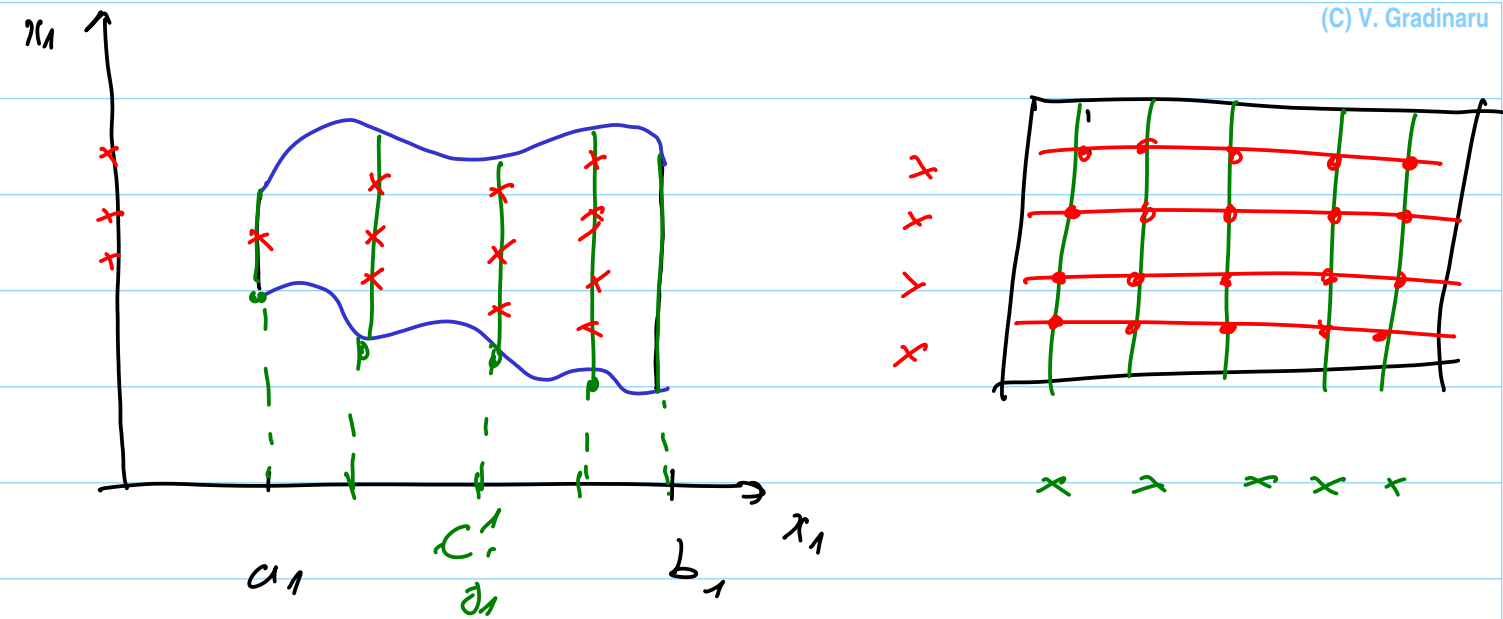
$$= \sum_{j_1=1}^{N_1} \sum_{j_2=1}^{N_2} f(c_{j_2}^2, c_{j_1}^1) w_{j_2}^2 w_{j_1}^1$$

$N_1 \cdot N_2$  Auswertungen von  $f$

$$\Omega = [0, 1]^d \subset \mathbb{R}^d$$

$$\int_{[0,1]^d} f(\underline{x}) d\underline{x} \approx \sum_{j_1=1}^{N_1} \sum_{j_2=1}^{N_2} \dots \sum_{j_d=1}^{N_d} f(c_{j_d}^d, c_{j_{d-1}}^{d-1}, \dots, c_{j_1}^1) \cdot w_{j_d}^d w_{j_{d-1}}^{d-1} \dots w_{j_1}^1$$

$N_1 \cdot N_2 \dots N_d$  Auswertungen von  $f$



\* Dünne Gitter / sparse grids

\* Monte Carlo

\* hochoszillatorische Funktionen.



# §1.6. Quadratur mit erhöhter Ordnung.

$\Delta$  Knoten auf Referenzintervall.

$\Delta$  Gewichte so bestimmt, dass Polynome höchstes Grades exakt mit der QF integriert werden.

Möchte Ordnung  $p = \Delta + m$

Jedes Polynom von Grad  $\leq \Delta + m - 1 = p - 1$  soll exakt integriert werden (mit QF)

IDEA: dividire  $f$  durch  $M(x) = (x-c_1)(x-c_2)\dots(x-c_n)$   
 $\text{Grad } M = \Delta$

Frage: Wenn wir frei sind sowohl Gewichte  $b_j$ , als auch Knoten  $c_j$ , wie hoch kann die Quadraturordnung werden?

$c_j$  äquidistant  $\Rightarrow \Delta$

2s Unbekannte, 2s Gleichungen.

$$\begin{cases} \int_0^1 1 dt = Q_\Delta(1, q, 1) \\ \int_0^1 t dt = Q_\Delta(t, q, 1) \\ \vdots \\ \int_0^1 t^{2\Delta-1} dt = Q_\Delta(t^{2\Delta-1}, q, 1) \end{cases}$$



$$f(x) = M(x)g(x) + r(x) \text{ mit } \text{Grad}(r) \leq \Delta - 1$$

$$\int_0^1 f(t) dt = \int_0^1 M(t)g(t) dt + \int_0^1 r(t) dt$$

||

$$\sum_{j=1}^{\Delta} b_j f(c_j) = \sum_{j=1}^{\Delta} b_j M(c_j)g(c_j) + \sum_{j=1}^{\Delta} b_j r(c_j)$$

$$\Rightarrow \int_0^1 M(t)g(t) dt = 0 \Rightarrow \langle M, g \rangle = 0 \text{ für alle } g$$

Polynome  $g$  mit  $\text{Grad}(g) \leq n-1$ .

$\langle M, g \rangle = \int_0^1 M(t)g(t)dt$  Skalarprodukt  
auf dem Raum der Polynome

Theorem Ordnung der QF ist  $n$   $\Leftrightarrow$   
 $\langle M, g \rangle = 0$  für alle Polynome  $g$  vom Grad  $\leq n-1$ .

$$P_n = \text{span}\{1, t, \dots, t^{n-1}\} : M \perp P_n$$

Theorem Ordnung einer QF ist höchstens  $n$ .

Beweis Annahme:  $p \geq 2n+1 \Rightarrow$

$$\int_0^1 M(t)g(t)dt = 0 \text{ für alle Polynome vom Grad } \leq p+1 \quad \Rightarrow$$

Nehme  $g = M$

$$\int_0^1 M(t)M(t)dt = 0 \Leftrightarrow \int_0^1 M(t)^2 dt = 0 \Rightarrow$$

$$\Rightarrow M(t) \equiv 0 \quad \rightarrow$$

$\Rightarrow$  wahr ist das Gegenteil der Annahme:  
 $p \leq 2n$ .

Orthogonale Polynome

$w: ]a, b[ \rightarrow \mathbb{R}$  Gewichtsfunktion  
stetig  $w(x) > 0$  für alle  $x \in ]a, b[$

$$\int_a^b |x|^k w(x) dx < \infty \text{ für } k=0, 1, 2, \dots$$

Betrachte den linearen Raum:

$$V = \left\{ f: ]a, b[ \rightarrow \mathbb{R}, f \text{ stetig}, \int_a^b |f(x)|^2 w(x) dx < \infty \right\}$$

Bem Alle Polynome liegen in  $V$   
auf  $V$ : Skalarprodukt

$$\langle f, g \rangle = \int_a^b f(x) g(x) w(x) dx.$$

T [Gram-Schmidt]:

Es existiert eine eindeutige Folge von Polynomen

$p_0, p_1, \dots$  mit  $p_k(x) = x^k + \text{Polynom vom Grad} \leq k-1$

so dass  $p_k \perp P_{k-1}$  bzgl.  $\langle \cdot, \cdot \rangle$ .

3-Term Rekurrenz:

$$p_{k+1}(x) = (x - \beta_{k+1}) p_k(x) - \gamma_{k+1}^2 p_{k-1}(x)$$

mit  $p_0(x) = 1$ ,  $p_{-1}(x) = 0$  und

$$\beta_{k+1} = \frac{\langle x p_k, p_k \rangle}{\langle p_k, p_k \rangle}, \quad \gamma_{k+1}^2 = \frac{\langle p_k, p_k \rangle}{\langle p_{k-1}, p_{k-1} \rangle}.$$

Bem  $c_1, c_2, \dots, c_s$  sind die Nullstellen von  $p_s$ .

Bsp  $w(x) \equiv 1 \Rightarrow$  orthogonale Polynome bzgl.  
 $\int_a^b f(x) g(x) dx \quad a=0, b=1$

Legendre Polynome

QF: Gauss-Quadratur

$\exists a, b \in ]-\infty, \infty[$ ,  $w(x) = e^{-x^2} \Rightarrow$  Hermite Polynome  
QF  $\Rightarrow$  Hermite-Quadratur.

Bsp Gauss-Quadratur:

1)  $n=1$  auf  $[0, 1]$ ,  $p_1(x) = x - \frac{1}{2}$ ;  $c_1 = \frac{1}{2}$ ,  $b_1 = 1$

auf  $[-1, 1]$   $c_j = \frac{1}{2}(1 + \gamma_j)$

$$c_1 = \frac{1}{2}(1 + \gamma_1) = 1$$

$$\frac{1}{2} = \frac{1}{2}(1 + \gamma_1) \Rightarrow$$

$$0 = \frac{1}{2}\gamma_1 \Rightarrow \gamma_1 = 0$$

Ordnung  $2n = 2 \cdot 1 = 2$

2)  $n=2$  auf  $[-1,1]$ ,  $p_2(x) = \frac{2}{2}x^2 - \frac{1}{2} = \frac{3}{2}(x^2 - \frac{1}{3})$

$$x_{1,2} = \pm \frac{1}{\sqrt{3}} ; c_1 = \frac{1}{2} - \frac{\sqrt{3}}{6}, c_2 = \frac{1}{2} + \frac{\sqrt{3}}{6}$$

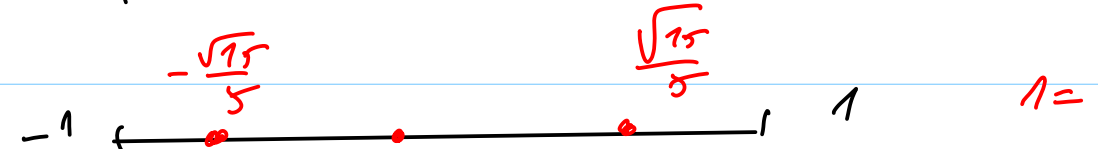
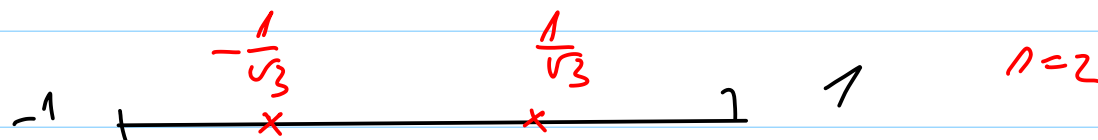
Ordnung  $2 \cdot n = 2 \cdot 2 = 4$ .

3)  $n=3$  auf  $[-1,1]$ ,  $p_3(x) = \frac{5}{2}x^3 - \frac{3}{2}x$

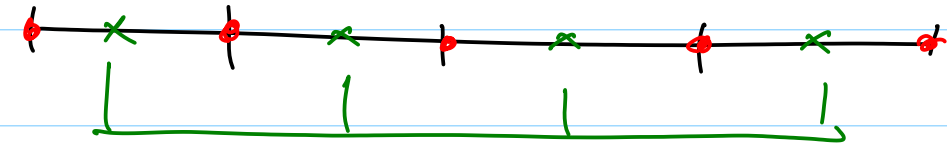
$$x_2=0, x_{1,3} = \pm \frac{\sqrt{15}}{5}, b_2 = \frac{8}{18}, b_1=b_3 = \frac{5}{18}$$

$$c_2 = \frac{1}{2}, c_{1,3} = \frac{1}{2} \pm \frac{\sqrt{15}}{10} \quad \text{Ordnung } 2 \cdot n = 6$$

Bez Gewichte der Gauss-QF sind positiv.



Bez Gauss-Knoten sind nicht verschachtelt  
 $\neq$  Trapez / Simpson  
 Nachterl.



Bez Knoten & Gewichte der Gauss-QF berechnet man via Eigenwertberechnung.

3-Term-Rekurrenz

$$\left\{ \begin{array}{l} p_k(x) = (a_k x + b_k) p_{k-1}(x) - c_k p_{k-2}(x) \quad \Rightarrow \\ p_{-1}(x) = 0, \quad p_0(x) = 1 \end{array} \right.$$

$$\left\{ \begin{array}{l} x p_{k-1}(x) = \boxed{\frac{c_k}{a_k}} p_{k-2}(x) - \boxed{\frac{b_k}{a_k}} p_{k-1}(x) + \boxed{\frac{1}{a_k}} p_k(x) \\ \text{für } k=1, 2, 3, \dots, n \end{array} \right.$$

$$\begin{aligned}
 x \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_{k-1} \\ \vdots \\ p_{n-2} \\ p_{n-1} \end{bmatrix} &= \begin{bmatrix} -\frac{b_1}{a_1} & \frac{1}{a_1} & & & \\ c_2/a_2 & -\frac{b_2}{a_2} & \frac{1}{a_2} & & \\ & c_k/a_k & -\frac{b_k}{a_k} & \frac{1}{a_k} & \\ & & & c_n/a_n & -\frac{b_n}{a_n} \\ & & & & 1/a_n \end{bmatrix} \begin{bmatrix} p_0 \\ p_1 \\ \vdots \\ p_{n-1} \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1/a_n \end{bmatrix} p_n \\
 \underbrace{\quad}_{\underline{f}(x)} & \quad \underbrace{\quad}_{\underline{A}} \quad ; \quad \underline{e}_n = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}
 \end{aligned}$$

$t_k \in W$  der Matrix  $\underline{A}$   
Eigenwerte.

$\underline{f}(t_k)$  ist  $\in V$  der Matrix  $\underline{A}$   
Eigenvektor.

Ben Es gibt ziemlich gute numerische  
Verfahren für die Berechnung  
von Approximationen  $a_n \in W, \in V$

$\Rightarrow$

$\text{eig}(\underline{A}) = \text{knoten}$

3-Term-Rekurrenz:  $x \underline{f}(x) = \underline{A} \underline{f}(x) + \frac{1}{a_n} p_n(x) \underline{e}_n$

Knoten  $\equiv$  Nullstellen von  $p_n(x) \Rightarrow p_n(t_k) = 0$   
 $\hookrightarrow t_1, t_2, \dots, t_n$  für  $k=1, 2, \dots, n$ .

$$\Rightarrow t_k \underline{f}(t_k) = \underline{A} \underline{f}(t_k) + 0 \Rightarrow$$

Gauss: Gewichte 2  $\underline{f}(t_0)$

Ben 1) Gauss nicht verschachtelt, nicht immer anwendbar

2)  $a, b$  sind keine Knoten  $\Rightarrow$  Gauss-QF ist offen.

3) Manchmal braucht man  $a$  oder  $b$  oder beide  
als Knoten, möchte trotzdem QF höchster  
Ordnung  $\Rightarrow 2n-1, 2n-2$

$\Rightarrow$  Selbe Prozedur aber mit 2 feste Knoten  $\Rightarrow$

$[-1, 1]$ ,  $w(x) = (1-x)(1+x)$ ,  $x_r = -1$  oder  $x_r = 1$

$\Rightarrow$  Radau-Quadratur.

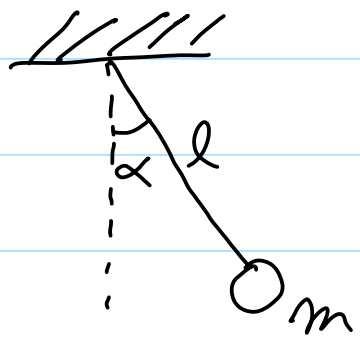
Beide  $-1, 1$  Knoten  $\Rightarrow$  Lobatto-Quadratur.



## §2. Einfache Verfahren für ODEs 1. Ordnung

ODE = ordinary differential equation  
gewöhnliche Differentialgleichungen

### §2.1 Linearisierung: global und lokal



$$m l \ddot{\alpha}(t) = -m g \sin \alpha(t)$$

$$(0) \begin{cases} \ddot{\alpha}(t) = -\frac{g}{l} \sin \alpha(t) \\ \alpha(0) = \alpha_0 \\ \dot{\alpha}(0) = \dot{\alpha}_0 \end{cases}$$

ODE 2. Ordnung da  $\ddot{\alpha}$   
autonom:  $t$  erscheint nicht explizit  
 $\ddot{y}(t) = f(y(t))$

$$\dot{\alpha}(t) = \frac{d}{dt} \alpha(t), \quad \ddot{\alpha}(t) = \frac{d^2}{dt^2} \alpha(t)$$

Die meiste Software ist für ODE 1. Ordnung

$$\dot{y}(t) = f(t, y(t))$$

"global": linearisiere die rechte Seite:

$$\sin \alpha = \alpha - \frac{1}{3!} \alpha^3 + \dots = \alpha + O(\alpha^3)$$

für kleines  $\alpha$ !

Neues Model.

$$(1) \begin{cases} \ddot{\beta}(t) = -\frac{g}{l} \beta(t) \\ \beta(0) = \alpha_0 \\ \dot{\beta}(0) = \dot{\alpha}_0 \end{cases} \quad \omega^2 = \frac{g}{l}$$

exakte Lösung

$$\beta(t) = \frac{\dot{\alpha}_0}{\omega} \sin(\omega t) + \alpha_0 \cos(\omega t)$$

Bsp Trick: Reduktion der Ordnung der ODE:

Notiere  $p(t) = \dot{\alpha}(t)$

$$\dot{p}(t) = \ddot{\alpha}(t) = -\frac{g}{l} \sin \alpha(t)$$

(0)  $\Leftrightarrow$  System ODEs 1. Ordnung:

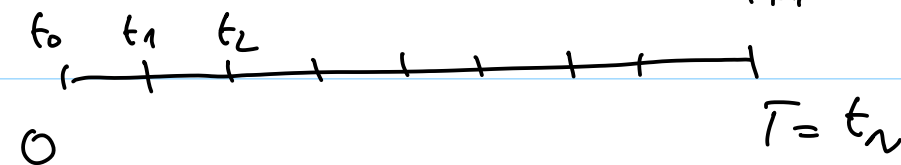
$$\begin{cases} \dot{\alpha} = p \\ \dot{p} = -\frac{g}{l} \sin \alpha \\ \alpha(0) = \alpha_0 \\ p(0) = \dot{\alpha}_0 \end{cases} \quad \underline{y} = \begin{bmatrix} \alpha \\ p \end{bmatrix}; \underline{y}(t) = \begin{bmatrix} \alpha(t) \\ p(t) \end{bmatrix} \in \mathbb{R}^2$$

$$\underline{f}(\underline{y}) = \begin{bmatrix} y_2 \\ -\frac{g}{l} \sin y_1 \end{bmatrix}$$

(0)  $\Leftrightarrow$  (2)  $\begin{cases} \dot{\underline{y}} = \underline{f}(\underline{y}) \\ \underline{y}(0) = \begin{bmatrix} \alpha_0 \\ \dot{\alpha}_0 \end{bmatrix} \end{cases}$

$\underline{y}: [0, T] \rightarrow \mathbb{R}^2$   
Unbekannte Funktion.  
ODEs 1. Ordnung  
autonom

Ziel: Möchte Approximation von  
 $\underline{y}(T) = \text{exakte Lösung zur Endzeit } T$   
Zeitgitter

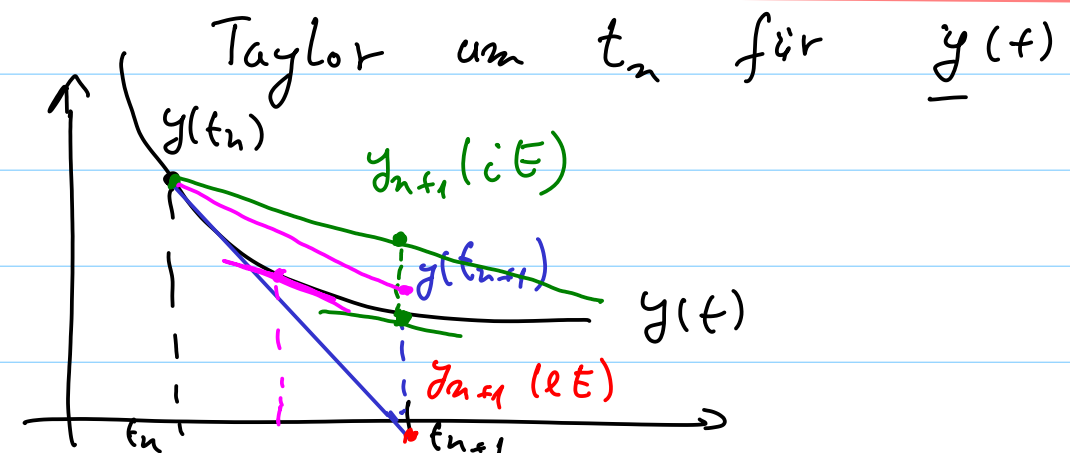
$$0 = t_0 < t_1 < t_2 < \dots < t_n < t_{n+1} < \dots < t_N = T$$


Zeitschritt  $h_n = t_{n+1} - t_n$

Wir werden Approximationen bauen

$$\underline{y}_n \approx \underline{y}(t_n)$$

IDEA: Lokal: Linearisiere die Lösung lokal:



$$\underline{y}(t_n+h) = \underline{y}(t_n) + h \dot{\underline{y}}(t_n) + O(h^2)$$

$\underline{y}_{n+1}$        $\underline{y}_n$        $f(t_n, \underline{y}_n)$

Methode:  $\underline{y}_{n+1} = \underline{y}_n + h f(t_n, \underline{y}_n)$  (eE)  
 expliziter Euler  
 mit lokalem Fehler  $O(h^2)$

Andere Herleitung

$$f(t_n, \underline{y}(t_n)) = \dot{\underline{y}}(t_n) = \lim_{h \rightarrow 0} \frac{\underline{y}(t_n+h) - \underline{y}(t_n)}{h} \approx$$

$$\Rightarrow f(t_n, \underline{y}(t_n)) = \frac{\underline{y}_{n+1} - \underline{y}_n}{h} \Rightarrow$$

$$\underline{y}_{n+1} = \underline{y}_n + h f(t_n, \underline{y}(t_n)) \quad (eE)$$

$$(eE): \begin{cases} \underline{y}_0 = \underline{y}(0) \\ \underline{y}_{n+1} = \underline{y}_n + h f(t_n, \underline{y}_n) \quad \text{für } n=0,1,2,\dots,N-1 \end{cases}$$

Taylor um  $t_{n+1}$  für  $\underline{y}(t)$ :

$$\underline{y}(t_{n+1}-h) = \underline{y}(t_{n+1}) - h \dot{\underline{y}}(t_{n+1}) + O(h^2)$$

$\underline{y}_n$        $\underline{y}_{n+1}$        $f(t_{n+1}, \underline{y}_{n+1})$

$$\Rightarrow \boxed{\underline{y}_n = \underline{y}_{n+1} - h f(t_{n+1}, \underline{y}_{n+1})} \quad (iE)$$

lokaler Fehler  $O(h^2)$       impliziter Euler

Ben  $f(t, y) = -\frac{y}{2} \ln y \Rightarrow \underline{y}_n = \underline{y}_{n+1} - h \left(-\frac{\underline{y}}{2}\right) \ln \underline{y}_{n+1}$   
 $\hookrightarrow$  algebraisches Problem  
 (Nullstelle!)

Neue Idee: linearisiere in der Mitte des Intervalls:  
Taylor um  $t^* = \frac{1}{2}(t_n + t_{n+1}) = t_n + \frac{1}{2}h$  für  $\underline{y}$ :

$$\underline{y}(t_{n+1}) = \underline{y}(t^*) + \frac{h}{2} \dot{\underline{y}}(t^*) + \frac{1}{2} \left(\frac{h}{2}\right)^2 \ddot{\underline{y}}(t^*) + O\left(\left(\frac{h}{2}\right)^3\right)$$

$$\underline{y}(t_n) = \underline{y}(t^*) - \frac{h}{2} \dot{\underline{y}}(t^*) + \frac{1}{2} \left(\frac{h}{2}\right)^2 \ddot{\underline{y}}(t^*) - O\left(\left(\frac{h}{2}\right)^3\right)$$

$$\Rightarrow \underline{y}(t_{n+1}) - \underline{y}(t_n) = h \dot{\underline{y}}(t^*) + O(h^3)$$

$$\parallel$$

$$\underline{f}(t^*, \underline{y}(t^*))$$

$$\Rightarrow \underline{y}(t_{n+1}) = \underline{y}(t_n) + h \underline{f}(t^*, \underline{y}(t^*)) + O(h^3)$$

$\underbrace{\quad}_{\underline{y}_{n+1}} \quad \underbrace{\quad}_{\underline{y}_n} \quad \underbrace{\quad}_{\quad}$

Ben  $\underline{y}(t^*)$  stört, brauche noch ein Trick!

$$1) \quad \underline{y}(t^*) \approx \frac{1}{2} (\underline{y}(t_n) + \underline{y}(t_{n+1})) \Rightarrow$$

$$\underline{y}_{n+1} = \underline{y}_n + h \underline{f}\left(t_n + \frac{h}{2}, \frac{1}{2}(\underline{y}_n + \underline{y}_{n+1})\right) \quad (\text{IMP})$$

implizite Mittelpunktsregel.  
lokaler Fehler  $O(h^3)$

$$2) \quad \underline{y}_{n+1} = \underline{y}_n + h \frac{1}{2} \left( \underline{f}(t_n, \underline{y}_n) + \underline{f}(t_{n+1}, \underline{y}_{n+1}) \right) \quad (\text{ETR})$$

implizite Trapezregel, lokaler Fehler  $O(h^3)$

Ben Es geht so nur wenn die exakte Lösung  $\underline{y}(t)$  genügend glatt ist, d.h.  $\ddot{\underline{y}}, \ddot{\underline{y}}, \dots$  existieren!

Bemerkung Addieren wir die 2 Gleichungen  $\Rightarrow$

$$\underline{y}(t_{n+1}) + \underline{y}(t_n) = 2\underline{y}(t^*) + \left(\frac{h}{2}\right)^2 \underline{\ddot{y}}(t^*) + O(h^4) \Leftrightarrow$$

$$\Rightarrow \underline{\ddot{y}}(t^*) = \frac{\underline{y}(t_{n+1}) - 2\underline{y}(t_n + \frac{h}{2}) + \underline{y}(t_n)}{\left(\frac{h}{2}\right)^2} + O(h^2)$$

$$\Rightarrow \underline{\ddot{y}}(t_n + \frac{h}{2}) \approx \frac{\underline{y}_{n+1} - 2\underline{y}_{n+\frac{1}{2}} + \underline{y}_n}{\left(\frac{h}{2}\right)^2} \quad \text{mit}$$

lokaler Fehler  $O(h^2)$

Bemerkung Selbe Rechnung um  $t_n \Rightarrow$

$$\underline{\ddot{y}}(t_n) \approx \frac{\underline{y}_{n+1} - 2\underline{y}_n + \underline{y}_{n-1}}{h^2} \quad \text{mit lokalem Fehler } O(h^2)$$

## § 2.2. Störmer-Verlet Verfahren

$$\begin{cases} \underline{\ddot{y}} = \underline{f}(t, \underline{y}) \\ \underline{y}(t_0) = \underline{y}_0 \\ \underline{\dot{y}}(t_0) = \underline{v}_0 \end{cases} \quad \begin{array}{ccccccc} & t_0 & & t_{k-1} & & t_k & & t_{k+1} & & t_N \\ & | & & | & & | & & | & & | \\ 0 & \text{-----} & & & & & & & & T \end{array}$$

$$\underline{f}(t_k, \underline{y}(t_k)) = \underline{\ddot{y}}(t_k) \approx \frac{\underline{y}_{k+1} - 2\underline{y}_k + \underline{y}_{k-1}}{h^2} \Rightarrow$$

$$\underline{y}_{k+1} = -\underline{y}_{k-1} + 2\underline{y}_k + h^2 \underline{f}(t_k, \underline{y}_k) \quad (\text{St.V.})$$

expliziter, 2-Schritt-Verfahren!

$\Rightarrow$  brauche Startwerte:  $\underline{y}(t_0) = \underline{y}_0$

$$\underline{y}_1 \approx \underline{y}(t_1)$$

eine Möglichkeit:  $\underline{y}_1 := \underline{y}_0 + h \underline{f}(t_0, \underline{y}_0)$  (eE) ☹

Taylor:  $\underline{y}_1 \approx \underline{y}(t_1) = \underbrace{\underline{y}(t_0)}_{\underline{y}_0} + h \underbrace{\dot{\underline{y}}(t_0)}_{\underline{v}_0} + \frac{h^2}{2} \underbrace{\ddot{\underline{y}}(t_0)}_{\underline{f}(t_0, \underline{y}_0)} + O(h^3)$

$$\underline{y}_1 := \underline{y}_0 + h \underline{v}_0 + \frac{h^2}{2} \underline{f}(t_0, \underline{y}_0)$$

für  $k=1, 2, 3, \dots$

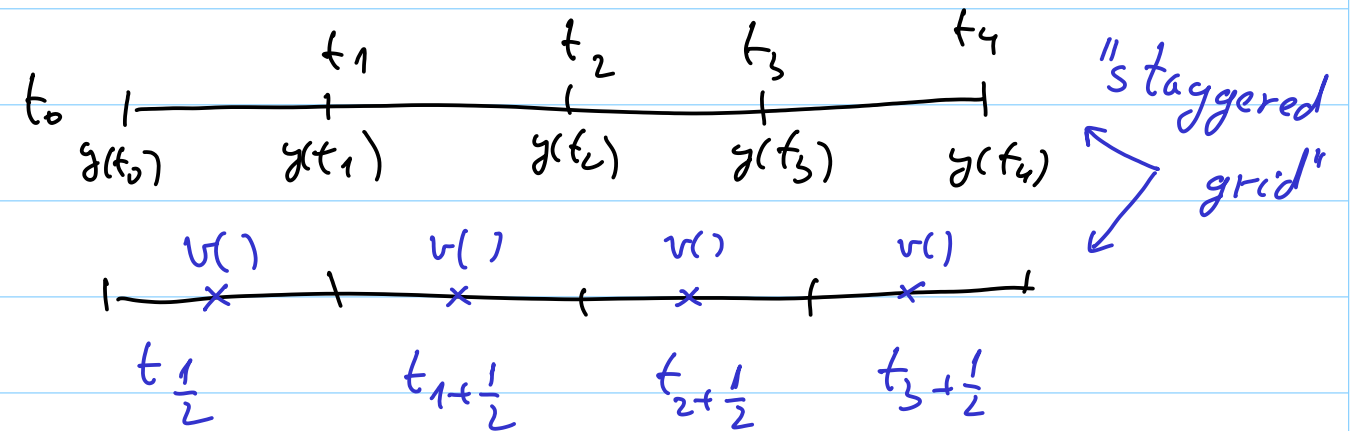
$$\underline{y}_{k+1} = -\underline{y}_{k-1} + 2\underline{y}_k + h^2 \underline{f}(t_k, \underline{y}_k)$$

lokal  $O(h^3)$ , global  $O(h^2)$

Notiere  $\underline{v}_{k+\frac{1}{2}} = \frac{\underline{y}_{k+1} - \underline{y}_k}{h}$  entspricht  $v(t_k + \frac{1}{2}h)$   
(gate Approximation an  $\dot{\underline{y}}(t_k + \frac{1}{2}h)$ )

in St.V)

$$\begin{cases} \underline{v}_{k+\frac{1}{2}} = \underline{v}_{k-\frac{1}{2}} + h \underline{f}(t_k, \underline{y}_k) \\ \underline{y}_{k+1} = \underline{y}_k + h \underline{v}_{k+\frac{1}{2}} \end{cases} \quad \text{"leap-frog"}$$



Besser noch: "velocity-Verlet"-Verfahren:

$$\begin{cases} \underline{y}_{k+1} = \underline{y}_k + h \underline{v}_k + \frac{h^2}{2} \underline{f}(t_k, \underline{y}_k) \\ \underline{v}_{k+1} = \underline{v}_k + h \frac{1}{2} (\underline{f}(t_k, \underline{y}_k) + \underline{f}(t_{k+1}, \underline{y}_{k+1})) \end{cases}$$

Begründung: Notiere  $\underline{v}_k = \frac{\underline{y}_{k+1} - \underline{y}_k}{2h}$  (StV)

$$\underline{v}_k = \frac{-2\underline{y}_{k-1} + 2\underline{y}_k + h^2 \underline{f}(t_k, \underline{y}_k)}{2h} = \frac{\underline{y}_k - \underline{y}_{k-1}}{h} + \frac{h}{2} \underline{f}(t_k, \underline{y}_k)$$

Somit

$$\begin{aligned} \underline{v}_{k+1} + \underline{v}_k &= \frac{\underline{y}_{k+1} - \underline{y}_k}{h} + \frac{h}{2} \underline{f}(t_{k+1}, \underline{y}_{k+1}) + \\ &+ \frac{\underline{y}_k - \underline{y}_{k-1}}{h} + \frac{h}{2} \underline{f}(t_k, \underline{y}_k) = \\ &= 2\underline{v}_k + \frac{h}{2} \left( \underline{f}(t_{k+1}, \underline{y}_{k+1}) + \underline{f}(t_k, \underline{y}_k) \right) \end{aligned}$$

$$\Rightarrow \underline{v}_{k+1} = \underline{v}_k + \frac{h}{2} \left( \underline{f}(t_{k+1}, \underline{y}_{k+1}) + \underline{f}(t_k, \underline{y}_k) \right)$$

Vorteile: explizit, genauer als (eE), (iE), erhalten die Energie!

## §2.3. Fehlerschätzung und Konvergenz

Theorem Taylor mit Rest als Integral:

$$\begin{aligned} f(x) &= f(a) + \frac{x-a}{1!} f'(a) + \dots + \frac{(x-a)^{n-1}}{(n-1)!} f^{(n-1)}(a) + \\ &+ \int_a^x \frac{(x-t)^{n-1}}{(n-1)!} f^{(n)}(t) dt \end{aligned}$$

Sei  $\underline{y}$  exakte Lösung von ODE

$$\dot{\underline{y}} = \underline{f}(t, \underline{y}) \quad \text{mit} \quad \underline{y}|_{t_0} = \underline{y}_0$$

$\underline{y}_n$  Approximation an  $\underline{y}(t_n)$ ,  $t_n = t_0 + nh$

durch (eE) definiert

$f: [t_0, T] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  stetig differenzierbar  
und Lipschitz:

$$\| \underline{f}(t, \underline{y}) - \underline{f}(t, \underline{z}) \| \leq L \| \underline{y} - \underline{z} \|$$

$\hookrightarrow$  konstante  $\in \mathbb{R}$ .

für alle  $\underline{y}, \underline{z}, t$ .

Theorem  $\| \underline{y}_n - \underline{y}(t_n) \| \leq M \cdot h$  für alle  $n$ , wobei

$$M = \frac{1}{L} \left( e^{L(T-t_0)} - 1 \right) \cdot \frac{1}{2} \max_{t \in [t_0, T]} \| \ddot{\underline{y}}(t) \|$$

Beweis 3 Schritte

1) lokaler Fehler: ein Schritt ( $e \in$ ) mit Start  $\underline{y}(t_n)$

$$\underline{y}(t_{n+1}) - \underline{y}_{n+1} = \underline{y}(t_{n+1}) - \left( \underline{y}(t_n) + h \underline{f}(t_n, \underline{y}(t_n)) \right) =$$

$\ddot{\underline{y}}(t_n) \leftarrow \text{ODE}$

$$= \underline{y}(t_{n+1}) - \underline{y}(t_n) - h \dot{\underline{y}}(t_n)$$

Satz von Taylor mit Rest als Integral  
( $a = t_n, x = t_n + h = t_{n+1}, f = \underline{y}$ )  $\Rightarrow$

$$\underline{y}(t_{n+1}) = \underline{y}(t_n) + h \dot{\underline{y}}(t_n) + \int_{t_n}^{t_{n+1}} (t_{n+1} - t) \ddot{\underline{y}}(t) dt$$

$$t = t_n + h\theta \quad dt = h d\theta$$

$$\Rightarrow \underline{y}(t_{n+1}) - \underline{y}(t_n) - h \dot{\underline{y}}(t_n) = h \int_0^1 (1-\theta) \ddot{\underline{y}}(t_n + h\theta) d\theta$$

$$\leq \max_{t \in [t_0, T]} \| \ddot{\underline{y}}(t) \|$$

Somit:

$$\| \underline{y}(t_{n+1}) - \underline{y}_{n+1} \| \leq \frac{1}{2} h^2 \cdot \max_{t \in [t_0, T]} \| \ddot{\underline{y}}(t) \| = C \cdot h^2$$



## 2) Fehlerfortpflanzung:

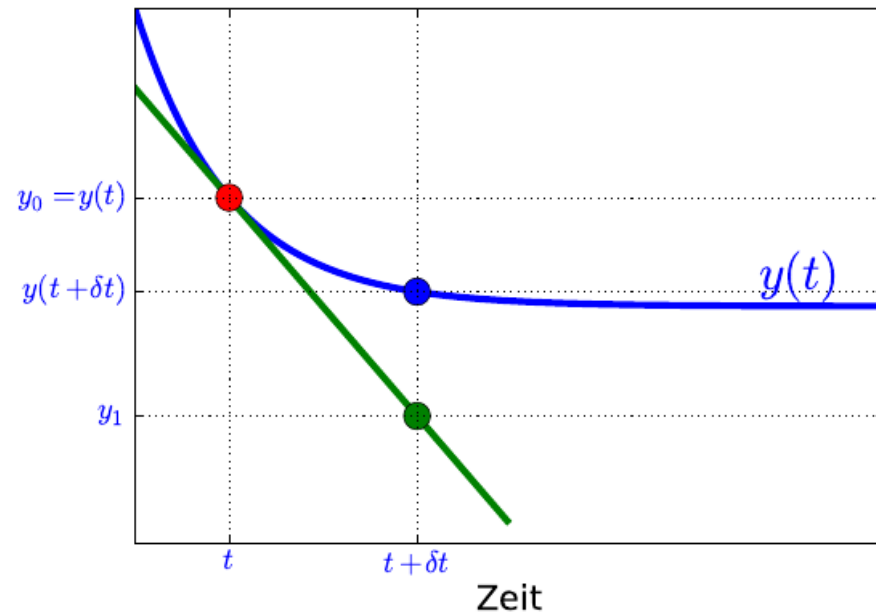
( $e \in E$ ) mit Startwert  $\underline{z}_n$ :  $\underline{z}_{n+1} = \underline{z}_n + h \underline{f}(t_n, \underline{z}_n)$

( $e \in E$ )  $\underline{w}_n$ :  $\underline{w}_{n+1} = \underline{w}_n + h \underline{f}(t_n, \underline{w}_n)$

$$\begin{aligned} \Rightarrow \|\underline{z}_{n+1} - \underline{w}_{n+1}\| &\leq \|\underline{z}_n - \underline{w}_n\| + h \|\underline{f}(t_n, \underline{z}_n) - \underline{f}(t_n, \underline{w}_n)\| \leq \\ &\leq \|\underline{z}_n - \underline{w}_n\| + h L \|\underline{z}_n - \underline{w}_n\| = \\ &\quad \underline{(1 + hL) \|\underline{z}_n - \underline{w}_n\|} \end{aligned}$$

## 3) Fehlerakkumulation

$$\begin{aligned} \|\underline{y}_n - \underline{y}(t_n)\| &\leq ch^2 + ch^2(1+hL) + ch^2(1+hL)^2 + \dots + \\ &\quad + ch^2(1+hL)^{n-1} = \\ &= ch^2 \frac{(1+hL)^n - 1}{1+hL - 1} = ch \frac{(1+hL)^n - 1}{L} \leq \\ &\leq c \frac{h}{L} (e^{nhL} - 1) = ch \frac{e^{(t_n - t_0)L} - 1}{L} \quad \underline{\text{g.g.d.}} \end{aligned}$$



$$\mathbf{y}_{n+1} = \mathbf{y}_n + h_n \mathbf{f}(t_n, \mathbf{y}_n), \quad n = 0, \dots, N-1.$$

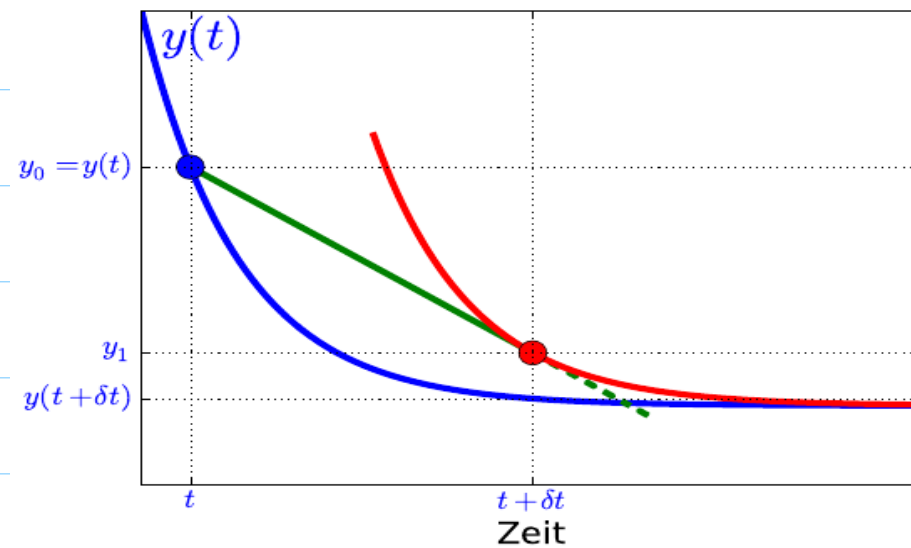
$$\mathbf{f}(t_n, \mathbf{y}(t_n)) = \dot{\mathbf{y}}(t_n) \approx \frac{\mathbf{y}(t_n + h_n) - \mathbf{y}(t_n)}{h_n}$$

$O(h^2)$  lokal

$\Downarrow$

$O(h)$

(CE)



$$\mathbf{y}_{n+1} = \mathbf{y}_n + h_n \mathbf{f}(t_{n+1}, \mathbf{y}_{n+1}), \quad n = 0, \dots, N-1$$

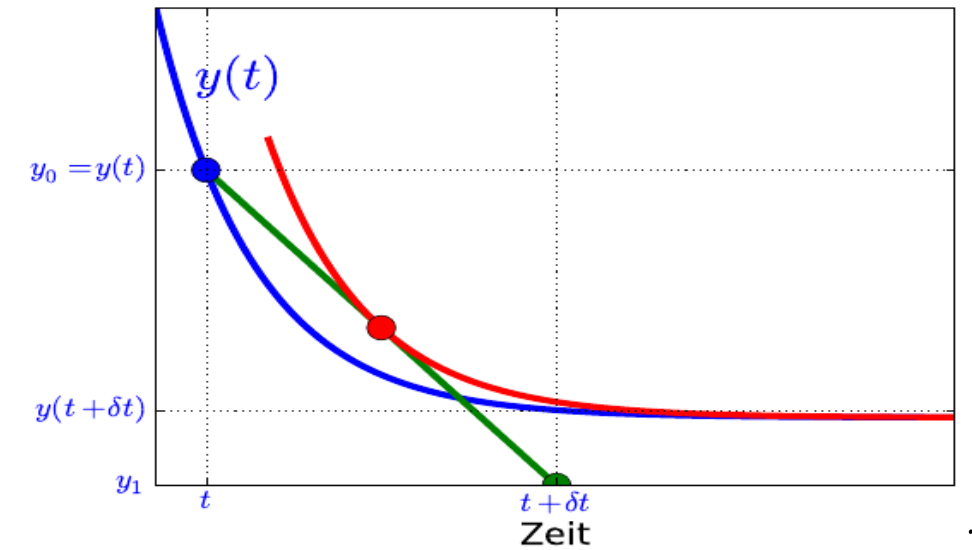
$$\mathbf{f}(t_{n+1}, \mathbf{y}(t_{n+1})) = \dot{\mathbf{y}}(t_{n+1}) \approx \frac{\mathbf{y}(t_{n+1} - h_n) - \mathbf{y}(t_{n+1})}{-h_n}$$

$O(h^2)$  lokal

$\Downarrow$

$O(h)$

(IE)



$$\mathbf{y}_{n+1} = \mathbf{y}_n + h_n \mathbf{f}\left(\frac{1}{2}(t_n + t_{n+1}), \frac{1}{2}(\mathbf{y}_n + \mathbf{y}_{n+1})\right).$$

$$\mathbf{f}\left(\frac{1}{2}(t_n + t_{n+1}), \frac{1}{2}(\mathbf{y}_n + \mathbf{y}_{n+1})\right) = \dot{\mathbf{y}}\left(\frac{1}{2}(t_n + t_{n+1})\right) \approx \frac{\mathbf{y}(t_n + h_n) - \mathbf{y}(t_n)}{h_n}.$$

$O(h^3)$  lokal

(MP)

$\Downarrow$

$O(h^2)$  global.

(iMP) lokaler Fehler:  $\bar{y} = \frac{1}{2}(y_n + y_{n+1})$   
 $t^* = \frac{1}{2}(t_n + t_{n+1})$

$$y(t_{n+1}) - y(t_n) - h f(t^*, \bar{y}) \stackrel{\text{Taylor}}{=} \\ = h \dot{y}(t^*) + o(h^3) - h f(t^*, \bar{y}) \stackrel{\downarrow \text{ODE}}{=}$$

$$= h f(t^*, y^*) - h f(t^*, \bar{y}) + o(h^3) = \\ = h (f(t^*, y^*) - f(t^*, \bar{y})) + o(h^3)$$

$$\|f(t^*, y^*) - f(t^*, \bar{y})\| \leq L \|y^* - \bar{y}\|$$

$$y^* - \bar{y} = \frac{1}{2} y(t^*) + \frac{1}{2} y(t^*) - \frac{1}{2} y(t_n) - \frac{1}{2} y(t_{n+1}) =$$

$$= \frac{1}{2} (y(t^*) - y(t_n)) + \frac{1}{2} (y(t^*) - y(t_{n+1}))$$

$$\stackrel{\text{Taylor}}{=} \left( \frac{h}{2} \dot{y}(t^*) + o(h^2) \right) - \left( \frac{h}{2} \dot{y}(t^*) + o(h^2) \right) = o(h^2)$$

Somit:

$$\|f(t^*, y^*) - f(t^*, \bar{y})\| \leq L \cdot C \cdot h^2 \Rightarrow$$

$$\|y(t_{k+1}) - y_{k+1}\| \leq h \cdot L \cdot C \cdot h^2 + o(h^3) = \underline{o(h^3)}.$$

(St.V.) Lokaler Fehler:

$$y(t_{n+1}) - (-y(t_{n-1}) + 2y(t_n) + h^2 f(t_n, y(t_n))) =$$

$$= y(t_{n+1}) + y(t_{n-1}) - 2y(t_n) - h^2 f(t_n, y(t_n)) =$$

$$\stackrel{\text{Taylor}}{=} \underbrace{h^2 \ddot{y}(t_n) - h^2 f(t_n, y(t_n))}_{\| \leftarrow \text{ODE} } + o(h^4) \Rightarrow$$

$\Rightarrow$  lokal  $o(h^4)$   $\nRightarrow$  global  $o(h^3)$

Startwert ist nur  $o(h^3)$   
 ausserdem Fehler fortpropagierung  $\Rightarrow$   $o(h^2)$ .

## § 2.4. Vorgehensweise bei Implementierung

Pendelgleichung:

$$\begin{cases} \ddot{\alpha} = -\frac{g}{l} \sin \alpha(t) & \text{Ziel: Approximation} \\ & \text{der Lösung \& Energien} \\ & \text{mittels} \\ \alpha(0) = \alpha_0 \\ \dot{\alpha}(0) = \dot{\alpha}_0 \end{cases} \quad (nA), (eE), (iE), (CHP), (STV)$$

1. Schritt: Umschreiben in ODE 1. Ordnung:

$$p = \dot{\alpha}, \quad \underline{y} = \begin{bmatrix} \alpha \\ p \end{bmatrix}, \quad \underline{f}(\underline{y}) = \begin{bmatrix} p \\ -\frac{g}{l} \sin \alpha \end{bmatrix}$$

2. Schritt: Idee der Lösung.

$$\underline{\dot{y}} = \underline{f}(\underline{y}) \quad y_n \approx y(t_n) \quad \text{für } t_n = 0 + n \cdot h$$

$$\omega^2 = g/l$$

$$(nA): \quad y_n = \frac{\dot{\alpha}}{\omega} \sin(\omega t_n) + \alpha_0 \cos(\omega t_n)$$

$$(eE) \quad \underline{y}_{n+1} = \underline{y}_n + h \underline{f}(\underline{y}_n)$$

$$(iE) \quad \underline{y}_{n+1} = \underline{y}_n + h \underline{f}(\underline{y}_{n+1})$$

$$\underline{y}_0 = [\alpha_0, \dot{\alpha}_0], \quad T, g, l \text{ gegeben}, \quad \omega^2 = \frac{g}{l}$$

↳ array

input variable

Rechte Seite:

```
def f(y):
    result = zeros(2)
    result[0] = y[1]
    result[1] = -g/l * sin(y[0])
    return result
```

(e)  $N$  Zeitschritte.

$$h = T/N$$

$$\begin{bmatrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \end{bmatrix}$$

$$y = ? \quad f\left(\frac{1}{2}(t_n + t_{n+1}), \frac{1}{2}(y_n + y_{n+1})\right) = \dot{y}\left(\frac{1}{2}(t_n + t_{n+1})\right) \approx \frac{y(t_n + h_n) - y(t_n)}{h_n}$$

$y[0]$  ...

für  $n=0, 1, \dots, N-1$

$$y[n+1] = y[n] + h \cdot f(y[n])$$

$$\text{plot}(y[0:], y[1:])$$

$g(t)$

(nA)

$$\frac{d}{dt} \left( \frac{\dot{z}}{\omega} \sin(\omega t) + \alpha_0 \cos(\omega t) \right) =$$
$$= \dot{z} \cos(\omega t) - \alpha_0 \omega \sin(\omega t) = dg(t)$$

$$y[0] = y_0$$

für  $n=0, 1, \dots, N-1$

$$y[n+1, 0] = g(y[n])$$

$$y[n+1, 1] = dg(y[n])$$

def.  $\text{potE}(y)$ :  $-y \cos(\alpha)$   
return  $(-g \cos y[0])$

(C) V. Gradinaru

$$\alpha = 1$$
$$l = 1$$

linE(y)

$$\text{return } \frac{1}{2} y[1]^2$$

def totE(y)

$$\text{return potE} + \text{linE}$$

(iE) für  $n=0, 1, 2, \dots, N-1$

$$\text{löse } \underline{y}_{n+1} = \underline{y}_n + h \underline{f}(\underline{y}_{n+1})$$

erste Möglichkeit: allgemein:

$$\underline{z} = \underline{a} + h \underline{f}(\underline{z}) \Leftrightarrow$$

$$\underline{z} - \underline{a} + h \underline{f}(\underline{z}) = 0 \Leftrightarrow \underline{F}(\underline{z}) = 0$$

mit  $\underline{F}: \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $\underline{F}(\underline{z}) = \underline{z} - \underline{a} - h \underline{f}(\underline{z})$

Finde Nullstelle ( $n$ ) von  $\underline{F}$

scipy.optimize.fsolve:

$\text{fsolve}(\underline{F}, \underline{z}_0)$

↳ Startwert; Hinweis wo

fsolve suchen soll

( $\underline{z}_0$  sollte nah an der Nullstelle sein)

$$\underline{z}_0 = \underline{y}_n$$

Startwert?



$\underline{z}_0 :=$  was (es) in diesem Schritt  
vorschlägt

$$\underline{z}_0 = \underline{y}_n + h \underline{f}(\underline{y}_n)$$

für  $n=0, 1, \dots, n-1$ :

$$\underline{z}_0 = \underline{y}_n + h \underline{f}(\underline{y}_n)$$

$$\underline{y}_{n+1} = \text{fsolve}(\underline{F}, \underline{z}_0)$$

zweite Möglichkeit: nutze dass für den Pendel  
 $f$  einfach ist  $\Rightarrow$

schreibe auf Papier das alg. Problem  
komponentenweise:

$$\begin{cases} z_1 - a_1 - h z_2 = 0 \\ z_2 - a_2 + h \frac{g}{l} \sin z_1 = 0 \Rightarrow z_2 = a_2 - h \frac{g}{l} \sin z_1 \end{cases}$$

$$\Rightarrow z_1 - a_1 - h a_2 + h^2 \frac{g}{l} \sin z_1 = 0$$

$$G(z) = z - a_1 - h a_2 + h^2 \frac{g}{l} \sin z$$

$$\underline{y}_{[n+1, 1]} = \text{fsolve}(G, \underline{z}_0)$$

↪ einfacher  
billiger!

### §3 Strukturhaltung

#### §3.1. Invariante und Hamilton Systeme

Bsp

autonome Lotka-Volterra:  $d=2$

$$\begin{cases} \dot{u} = (\alpha - \beta v)u \\ \dot{v} = (\delta u - \gamma)v \end{cases} \quad \begin{array}{l} \text{Konstanten} \\ \alpha, \beta, \gamma, \delta > 0 \\ u, v: \mathbb{R} \rightarrow \mathbb{R} \\ \text{Unbekannt} \end{array}$$

$$\left. \begin{aligned} \dot{u} &= (\alpha - \beta v)u = \left(\frac{\alpha}{v} - \beta\right)uv \\ \dot{v} &= (\delta u - \gamma)v = \left(\delta - \frac{\gamma}{u}\right)uv \end{aligned} \right\} \Rightarrow$$

$$\Rightarrow \left(\delta - \frac{\gamma}{u}\right)\dot{u} = \left(\frac{\alpha}{v} - \beta\right)\dot{v} \Rightarrow$$

$$\frac{d}{dt} \underbrace{\left(\delta u - \gamma \log u - \alpha \log v + \beta v\right)}_{I(u(t), v(t))} = 0 \quad \text{für alle } t \Rightarrow$$

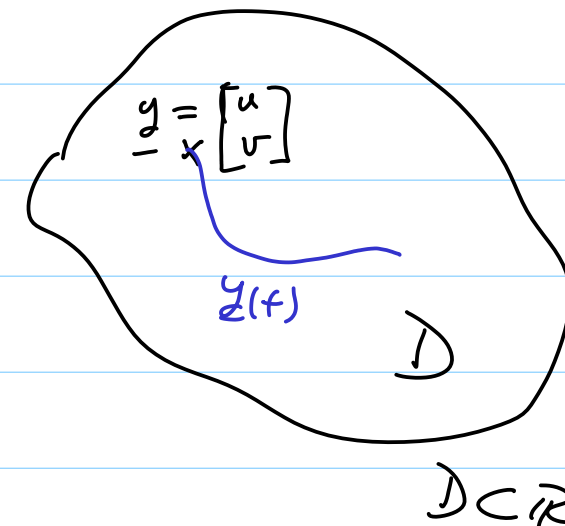
$$\frac{d}{dt} (I(u(t), v(t))) = 0 \quad \text{für alle } t \Rightarrow$$

$$\Rightarrow I(u(t), v(t)) = \text{konstant}$$

Def  $I$  heißt erstes Integral / Invariante der ODE  $\underline{\dot{y}} = \underline{f}(t, \underline{y})$  wenn

$I(\underline{y}(t)) = \text{konstant}$  für jede Lösung  $\underline{y} = \underline{y}(t)$  der ODE.

$$I: \underset{\substack{\mathbb{D} \subset \mathbb{R}^d \\ \underline{y}(t)}}{\mathbb{D} \subset \mathbb{R}^d} \rightarrow \mathbb{R}$$



$$\underline{y}: [0, T] \rightarrow \mathbb{R}^2$$

$$\underline{y}(t) \in \mathbb{D} \subset \mathbb{R}^2$$

Erinnerung  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(\underline{x}) \in \mathbb{R}$ .

partielle Ableitungen  $\frac{\partial}{\partial x_1} f(\underline{x}) =$  Ableitung von  $f$  nach der Variable  $x_1$

$$\frac{\partial}{\partial x_2}, \dots, \frac{\partial}{\partial x_d} \quad \underline{\text{grad}} f(\underline{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(\underline{x}) \\ \vdots \\ \frac{\partial f}{\partial x_d}(\underline{x}) \end{bmatrix}$$

Bsp  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $f(\underline{x}) = \|\underline{x}\|_2 = \left( \sum_{j=1}^d x_j^2 \right)^{1/2}$

Was ist  $\underline{\text{grad}} f(\underline{x}) = ?$

$$\frac{\partial f}{\partial x_1}(\underline{x}) = \frac{1}{2} \left( \sum_{j=1}^d x_j^2 \right)^{-\frac{1}{2}} \cdot 2x_1 = \frac{x_1}{\|\underline{x}\|_2} \Rightarrow$$

$$\Rightarrow \underline{\text{grad}} f(\underline{x}) = \frac{1}{\|\underline{x}\|_2} \underline{x}$$

Theorem  $I$  ist Invariante für  $\dot{\underline{y}} = \underline{f}(t, \underline{y})$   
 $\Leftrightarrow$

$$\underline{\text{grad}} I(\underline{y}) \cdot \underline{f}(t, \underline{y}) = 0 \text{ für alle}$$

$(t, \underline{y}) \in [0, T] \times D$

wo es eine (glatte) Lösung gibt

Beweis

Annahme,  $\underline{y} := \underline{y}(t)$  für jedes  $t$ , Lösung der ODE  
 $\downarrow$   
 $\Leftarrow: 0 = \underline{\text{grad}} I(\underline{y}(t)) \cdot \underline{f}(t, \underline{y}(t)) \stackrel{\downarrow}{=}$

$$= \underline{\text{grad}} I(\underline{y}(t)) \cdot \dot{\underline{y}}(t) = \frac{d}{dt} I(\underline{y}(t)) \Rightarrow$$

$\uparrow$   
Kettenregel!

$\Rightarrow I(\underline{y}(t)) = \text{konstant} \Rightarrow I$  ist Invariante der ODE



$\Rightarrow$ : Nehme  $\underline{z}_0 \in \mathcal{D}$ ,  $t_0 \in [0, T]$  beliebig.  
verwende die ODE  $\begin{cases} \dot{\underline{z}} = f(t, \underline{z}) \\ \underline{z}(t_0) = \underline{z}_0 \end{cases} \Rightarrow \underline{z}(t) =$

$I$  Invariante

$\downarrow$   
 $\Rightarrow \text{grad } I(\underline{z}(t)) \cdot f(t, \underline{z}(t)) = 0$  für alle  $t$

nehme  $t = t_0 \Rightarrow \text{grad } I(\underline{z}_0) \cdot f(t_0, \underline{z}_0) = 0$ .

Da  $\underline{z}_0, t$  beliebig  $\Rightarrow \text{grad } I(\underline{z}) \cdot f(t, \underline{z}) = 0$

für alle  $(t, \underline{z}) \in [t_0, T] \times \mathcal{D}$ .  
qed.

Notation:  $H: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $H(\underline{p}, \underline{q})$

$$\left[ \frac{\partial}{\partial p_1} H(\underline{p}, \underline{q}), \dots, \frac{\partial}{\partial p_d} H(\underline{p}, \underline{q}) \right]^T = \frac{\partial H}{\partial \underline{p}} = \text{grad}_{\underline{p}} H(\underline{p}, \underline{q})$$

Def Hamiltonische Differentialgleichung

$$H: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R} \quad H(\underline{p}, \underline{q})$$

$$\begin{cases} \dot{\underline{p}}_j = - \frac{\partial H}{\partial q_j}(\underline{p}(t), \underline{q}(t)) \\ \dot{\underline{q}}_j = \frac{\partial H}{\partial p_j}(\underline{p}(t), \underline{q}(t)) \end{cases} \quad \text{für } j=1, 2, \dots, d \quad \underline{p}, \underline{q}: [0, T] \rightarrow \mathbb{R}^d$$

autonomes Hamilton-System mit der  
Hamilton-Funktion  $H$ .

Bsp 1. Pendel:  $p, q: \mathbb{R} \rightarrow \mathbb{R}$  ( $d=1$ ) ( $m=1$ )

$$H(p, q) = \frac{1}{2} p^2 - \frac{g}{l} \cos(q) = E_{\text{tot}} \cdot \frac{1}{m l^2}$$

(Pendel.  $q = \alpha = \text{Winkel}$ ,  $p = \dot{q}$ )

$$\begin{cases} \dot{q} = \frac{1}{m} p \\ \dot{p} = -G'(\|q\|_2) \cdot \frac{q}{\|q\|_2} \end{cases}$$

2. Konservatives Feld, d.h.  $f(\underline{x}) = -\underline{\text{grad}} U(\underline{x})$

( $\underline{x} \in \mathbb{R}^d$ )

Vergleiche mit Newton's Gleichung:

z.B.  $U(\underline{x}) = G(\|\underline{x}\|_2)$

$$m \ddot{r}(t) = f(r(t))$$

$$f(\underline{q}) = -\underline{\text{grad}} G(\|\underline{q}\|) = -G'(\|\underline{q}\|_2) \cdot \underline{\text{grad}} \|\underline{q}\|_2$$

$$H(\underline{p}, \underline{q}) = \frac{1}{2m} \|\underline{p}\|^2 + G(\|\underline{q}\|_2)$$

Bem  $H(\underline{p}, \underline{q})$  Invariante für das Ham. System.

$p = m\dot{r} \Rightarrow$  Bewegungsgleichungen  
vom Massenpunkt  $m$   
im konservativen Zentralfeld.

Beweis  $\underline{y} = \begin{bmatrix} \underline{p} \\ \underline{q} \end{bmatrix}$

$$\underline{\partial} = \begin{bmatrix} \underline{0} & \underline{I}_d \\ -\underline{I}_d & \underline{0} \end{bmatrix} \stackrel{d=2}{=} \left[ \begin{array}{cc|cc} 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ \hline -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \end{array} \right]$$

Hamilton' System:

$$\dot{\underline{y}} = \underline{\partial}^{-1} \text{grad } H(\underline{y}) =: \underline{f}(\underline{y})$$

$\underline{I}(\underline{y})$  ist Invariante für  $\dot{\underline{y}} = \underline{f}(\underline{y}) \Leftrightarrow$

$$\text{grad } \underline{I}(\underline{y}) \cdot \underline{f}(\underline{y}) = 0 \quad \text{für alle } \underline{y}$$

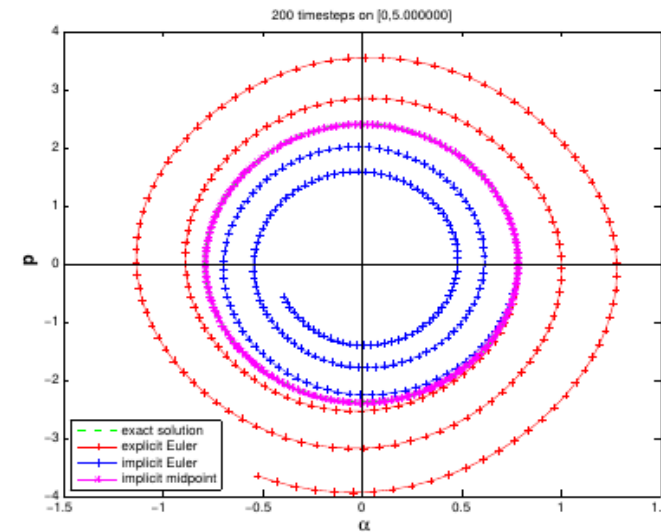
Überprüfen wir das für  $\underline{I} := H$  :

$$\text{grad } H(\underline{y}) \cdot \underline{f}(\underline{y}) = \text{grad } H(\underline{y})^T \underline{\partial}^{-1} \text{grad } H(\underline{y}) =$$

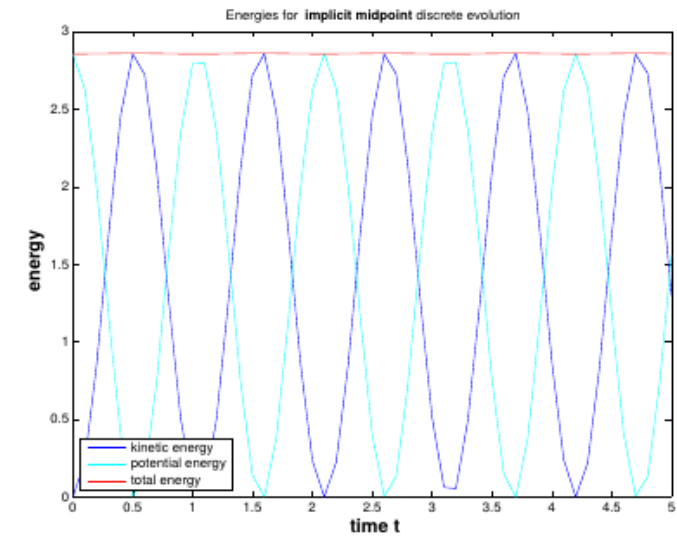
$$\text{grad } H(\underline{y}) = \begin{bmatrix} \underline{a} \\ \underline{b} \end{bmatrix}$$

$$= \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix} \begin{bmatrix} \underline{0} & -\underline{I} \\ \underline{I} & \underline{0} \end{bmatrix} \begin{bmatrix} \underline{a} \\ \underline{b} \end{bmatrix} = \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix} \begin{bmatrix} -\underline{b} \\ \underline{a} \end{bmatrix} = \underline{a} \underline{b} - \underline{b} \underline{a} = 0$$

$\underline{q+d}$



a) Vergleich e.E., i.E., i.M.



b) implizite Mittelpunktsregel: konstante Energie

Abb. 2.2.7. Strukturerhaltungseigenschaften der impliziten Mittelpunktsregel

Ben iMP, St-V erhalten die Energie;  $O(h^2)$   
Andere Methoden, die strukturerhaltend sind? höhere Ordnung?

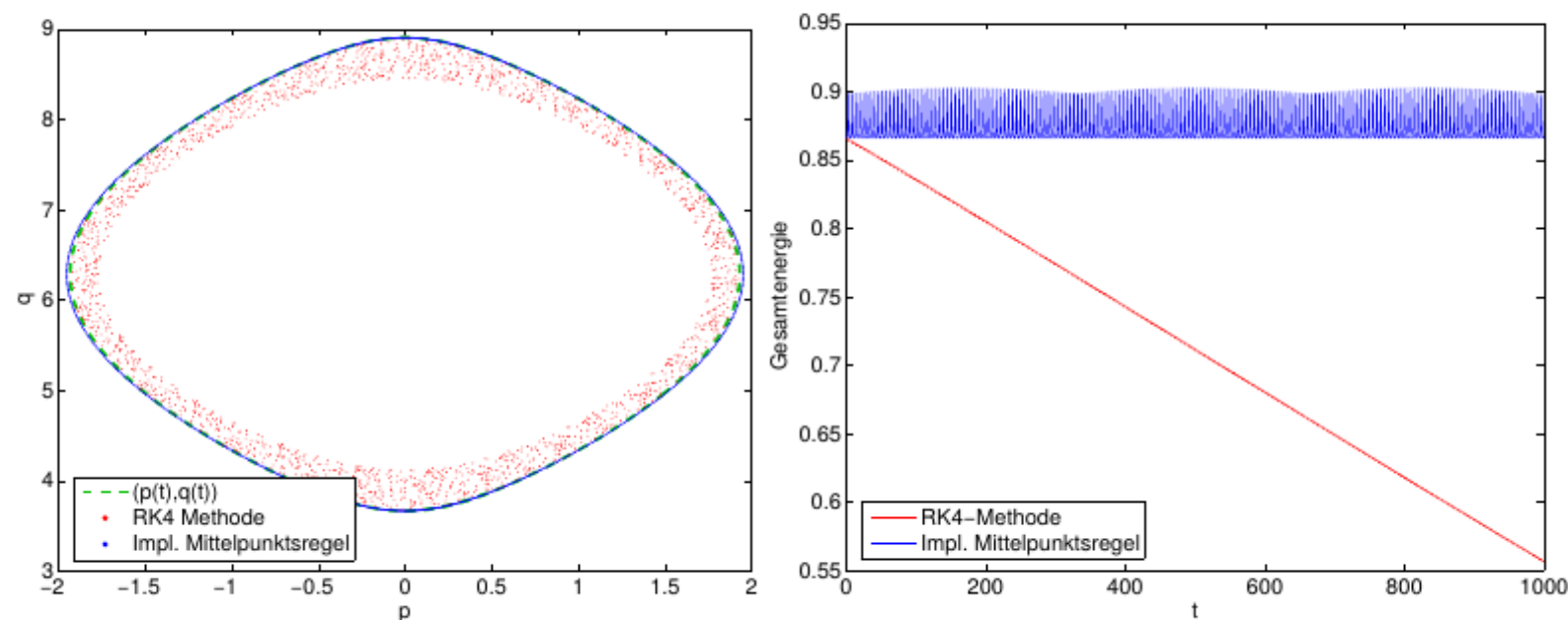
Solche weniger präzise Methoden sind für solch Dgl. (wo Strukturerhaltung wichtig ist) zu bevorzugen!

Theorem  $I: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $I(\underline{y}) = \frac{1}{2} \underline{y}^T \underline{B} \underline{y}$  mit

$\underline{B} \in \mathbb{R}^{d \times d}$ ,  $I$  Invariante für  $\dot{\underline{y}} = \underline{f}(\underline{y})$   
 $\underline{f}$  diff<sup>bar</sup>. Ser ( $\underline{y}_k$ ) aus iMP.

Dann  $I(\underline{y}_k) = I(\underline{y}_0)$  für alle  $k$ .

Wir vergleichen das klassische Runge-Kutta-Verfahren (2.5.34) der Ordnung 4 mit der impliziten Mittelpunktsregel bei konstanter Zeitschrittweite  $h = \frac{1}{2}$ :



### §3.2. Splitting Verfahren

autonome ODE:

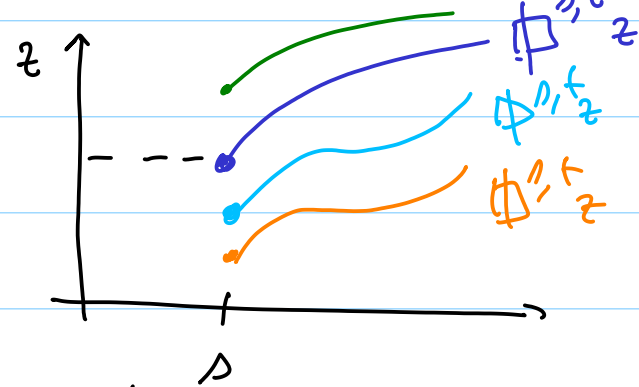
$$\dot{\underline{y}} = \underline{f}(\underline{y})$$

Bez Falls  $\dot{\underline{y}} = \underline{f}(t, \underline{y})$  nicht autonom ist, dann kann man sie autonomisieren, und zwar mit dem Trick:

$$\underline{z} = \begin{bmatrix} t \\ \underline{y} \end{bmatrix} \Rightarrow \begin{matrix} \dot{\underline{z}} = \begin{bmatrix} 1 \\ \underline{f}(\underline{z}) \end{bmatrix} \text{ ist autonom.} \\ \dot{t} = 1 \end{matrix}$$

$$\underline{z}(0) = \begin{bmatrix} t_0 \\ \underline{y}_0 \end{bmatrix}.$$

$$\text{ODE: } \begin{cases} \dot{y} = f(t, y) \\ y(1) = \bar{z} \end{cases} \quad \text{Lösung: } y(t) \text{ für } t > 1$$



$\Phi^{s,t} : D \rightarrow D$  zweiparametrische Familie  
von Abbildungen "Fluss der ODE"

$$\Phi^{t,t} = \text{Id} ; \quad \Phi^{s,t}_z = \Phi^{r,t} \circ \Phi^{s,r}_z$$

wenn ODE autonom ist:  $\Phi^{s,t} = \Phi^{0,t-s} = \Phi^{t-s}$

d.h. die Lösung der autonomen ODE

$$\dot{y} = \underline{f}(y)$$

ist translationsinvariant!

$$\dot{y} = \underline{f}(y)$$

$$\begin{aligned} u(t) &:= y(t+\tau) \Rightarrow \frac{d}{dt} u(t) = \dot{y}(t+\tau) \cdot 1 = \\ &= \underline{f}(y(t+\tau)) = \underline{f}(u) \end{aligned}$$

$$\Rightarrow \dot{u} = \underline{f}(u)$$

Ben Jede ODE kann man autonomisieren!

$$\dot{y} = f(t, y) ; \quad z = \begin{bmatrix} t \\ y \end{bmatrix} \Rightarrow$$

$$\dot{z} = \begin{bmatrix} 1 \\ y \end{bmatrix} = \begin{bmatrix} 1 \\ f(t) \end{bmatrix} \quad \text{autonom.}$$

(ODE):  $\underline{\dot{y}} = \underline{f}(y) = \underline{f}_a(y) + \underline{f}_b(y)$

Idee: wähle  $\underline{f}_a$ ,  $\underline{f}_b$  so dass

(a)  $\underline{\dot{y}} = \underline{f}_a(y)$  und

(b)  $\underline{\dot{y}} = \underline{f}_b(y)$

einfach/exakt lösbar sind.

Ben Für lineares Problem:

$$\underline{\dot{y}} = \underline{M} \underline{y} = (\underline{A} + \underline{B}) \underline{y} = \underline{A} \underline{y} + \underline{B} \underline{y}$$

mit  $\underline{A}, \underline{B}$  mit speziellen Eigenschaften....

exakte Lösung  $\underline{y}(t) = e^{\underline{M}t} \underline{y}_0 = e^{(\underline{A} + \underline{B})t} \underline{y}_0$

$$e^{\underline{M}t} = \sum_{n=0}^{\infty} \frac{1}{n!} (\underline{M}t)^n \quad (\text{nicht so rechnen})$$

$$e^{(\underline{A} + \underline{B})t} \neq e^{\underline{A}t} e^{\underline{B}t} \quad \text{z.B. wenn } \underline{A}\underline{B} \neq \underline{B}\underline{A}$$

Lösung von  $\begin{cases} \underline{\dot{y}} = \underline{A} \underline{y} \\ y(0) = y_0 \end{cases}$  ist  $e^{\underline{A}t}$

$\begin{cases} \underline{\dot{y}} = \underline{B} \underline{y} \\ y(0) = y_0 \end{cases}$  ist  $e^{\underline{B}t}$

$e^{\underline{B}t} (e^{\underline{A}t} \underline{y}_0) \cdot y(t) = e^{(\underline{A} + \underline{B})t} \underline{y}_0$

$e^{\underline{B}t} \underline{y}_0 \cdot e^{\underline{A}t} (e^{\underline{B}t} \underline{y}_0) \cdot e^{(\underline{A} + \underline{B})t} \underline{y}_0$

$$e^{(A+B)t} = \underline{I} + (\underline{A} + \underline{B})t + \frac{1}{2}(\underline{A}^2 + \underline{A}\underline{B} + \underline{B}\underline{A} + \underline{B}^2)t^2 + \dots$$

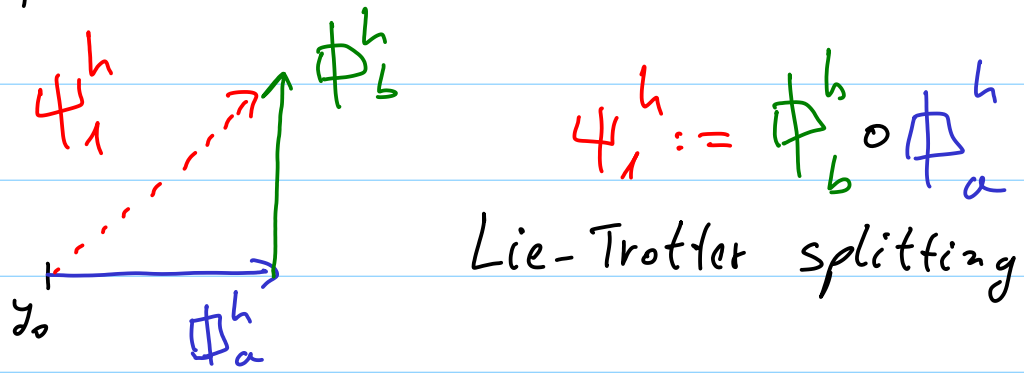
$$(e^{A^+})(e^{B^+}) = (\underline{I} + \underline{A}t + \frac{1}{2}\underline{A}^2t^2 + \dots)(\underline{I} + \underline{B}t + \frac{1}{2}\underline{B}^2t^2 + \dots) =$$
$$= \underline{I} + (\underline{B} + \underline{A})t + (\underline{A}\underline{B}t^2 + \frac{1}{2}\underline{A}^2t^2 + \frac{1}{2}\underline{B}^2t^2) + \dots$$

Der  $t^2$ -Termin ist in  $e^{(A+B)t}$ :  $\frac{1}{2}\underline{A}^2 + \frac{1}{2}(\underline{A}\underline{B} + \underline{B}\underline{A}) + \frac{1}{2}\underline{B}^2$

in  $e^{A^+}e^{B^+}$ :  $\frac{1}{2}\underline{A}^2 + \underline{A}\underline{B} + \frac{1}{2}\underline{B}^2$

ungleich wenn  $\underline{A}\underline{B} \neq \underline{B}\underline{A}$

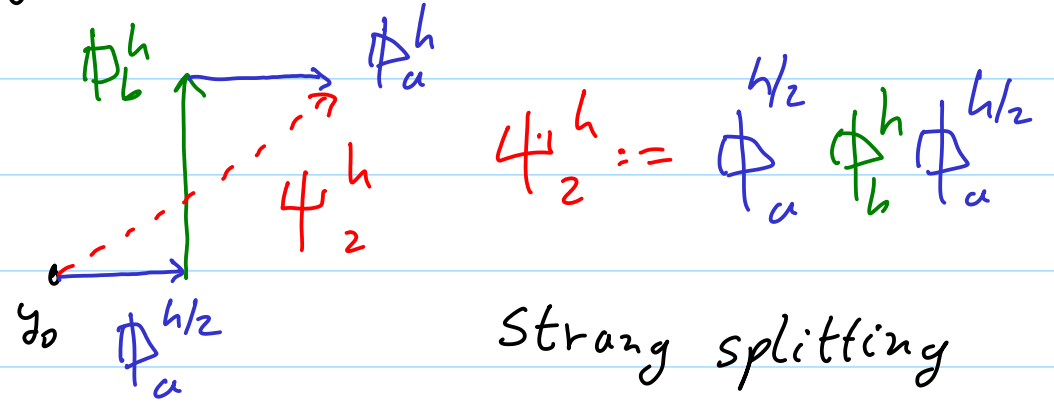
Man kann beweisen, dass dies eine Approximation liefert!



$$y_1 := \Phi_b^h \Phi_a^h y_0 \quad O(h) \text{ zur Endzeit } T$$

$$\Phi_a^h \Phi_b^h y_0 \quad O(h) \text{ zur Endzeit } T$$

Symmetrie schenkt uns ein Ordnung mehr:



Allgemein:  $\Psi^h = \prod_{i=1}^N \Phi_b^{b_i h} \Phi_a^{a_i h}$

mit  $\sum_{i=1}^N a_i = 1, \quad \sum_{i=1}^N b_i = 1$

Bsp  $N=1, a_1=1, b_1=1 \Rightarrow$  Lie-Trotter

$N=2, a_1=a_2=\frac{1}{2}, b_1=1, b_2=0 \Rightarrow$  Strang

Welche  $\underline{a}, \underline{b}$  wählen damit wir hohe Ordnung?  
 $f_a, f_b$  Erhaltungseigenschaften.

**Beispiel 2.4.1.** (Konvergenz einfacher Splitting-Verfahren)

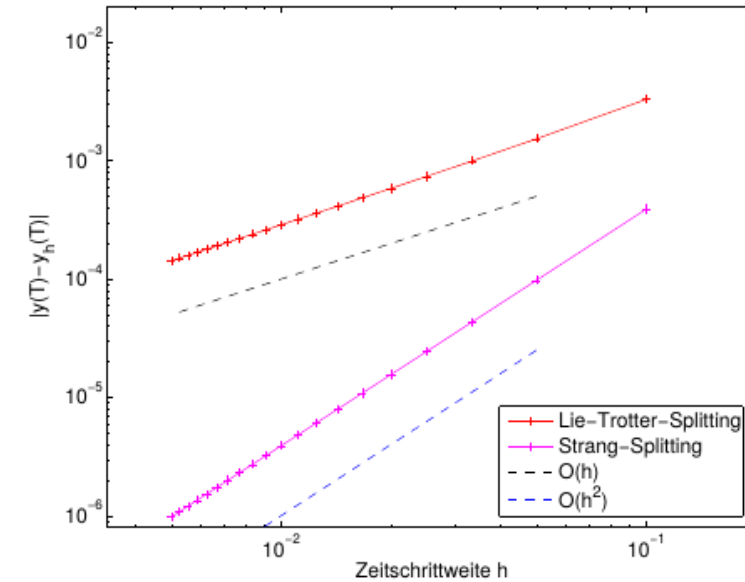
Sei

$$\dot{y} = \underbrace{\lambda y(1-y)}_{=:f_a(y)} + \underbrace{\sqrt{1-y^2}}_{=:f_b(y)}, \quad y(0) = 0.$$

Die Evolutionsoperatoren der zwei Teile sind analytisch bekannt:

$$\Phi_a^t y = \frac{1}{1 + (y^{-1} - 1)e^{-\lambda t}}, \quad \text{für } t > 0, y \in ]0, 1] \text{ und}$$

$$\Phi_b^t y = \begin{cases} \sin(t + \arcsin(y)), & \text{wenn } t + \arcsin(y) < \frac{\pi}{2}, \\ 1, & \text{sonst} \end{cases} \quad \text{für } t > 0, y \in [0, 1].$$

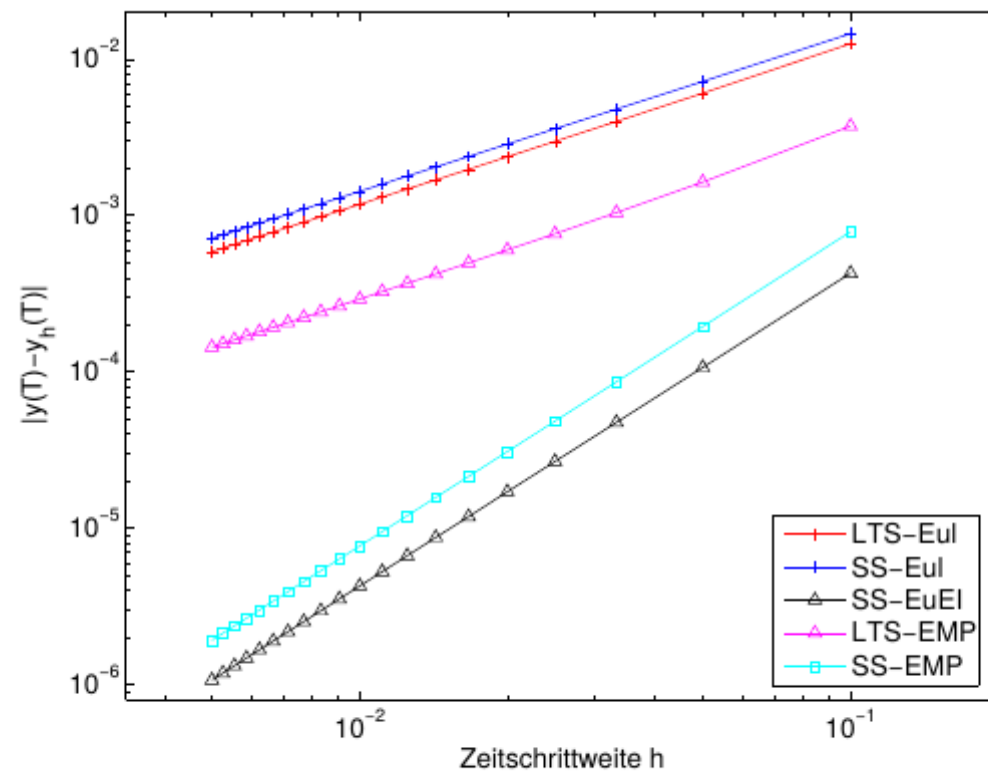


**Abb. 2.4.2.** Fehler zur Endzeit  $T = 1$ .

Numerisches Experiment:

$T = 1, \lambda = 1$ , Vergleich :  
 einer sehr genauen numerisch  
 Lösung.



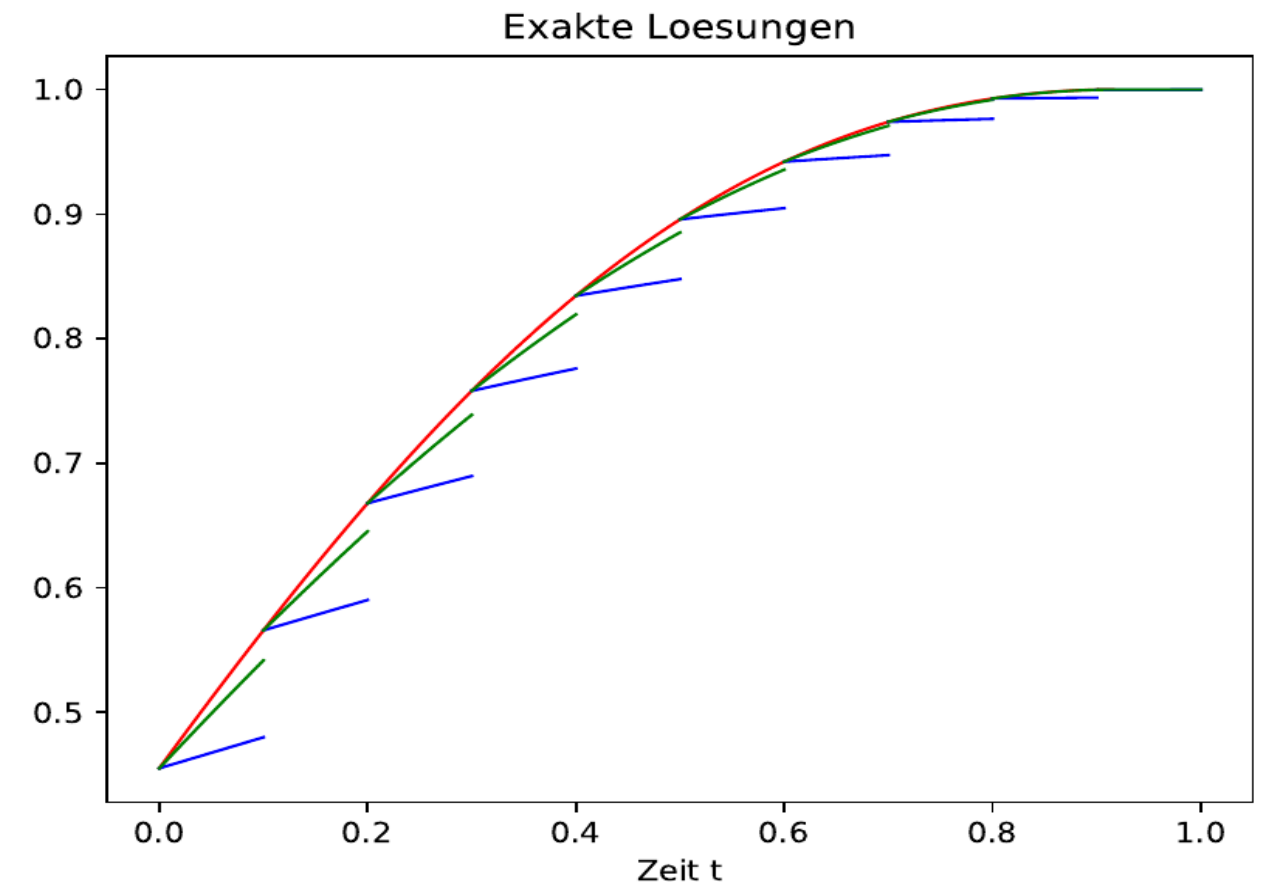
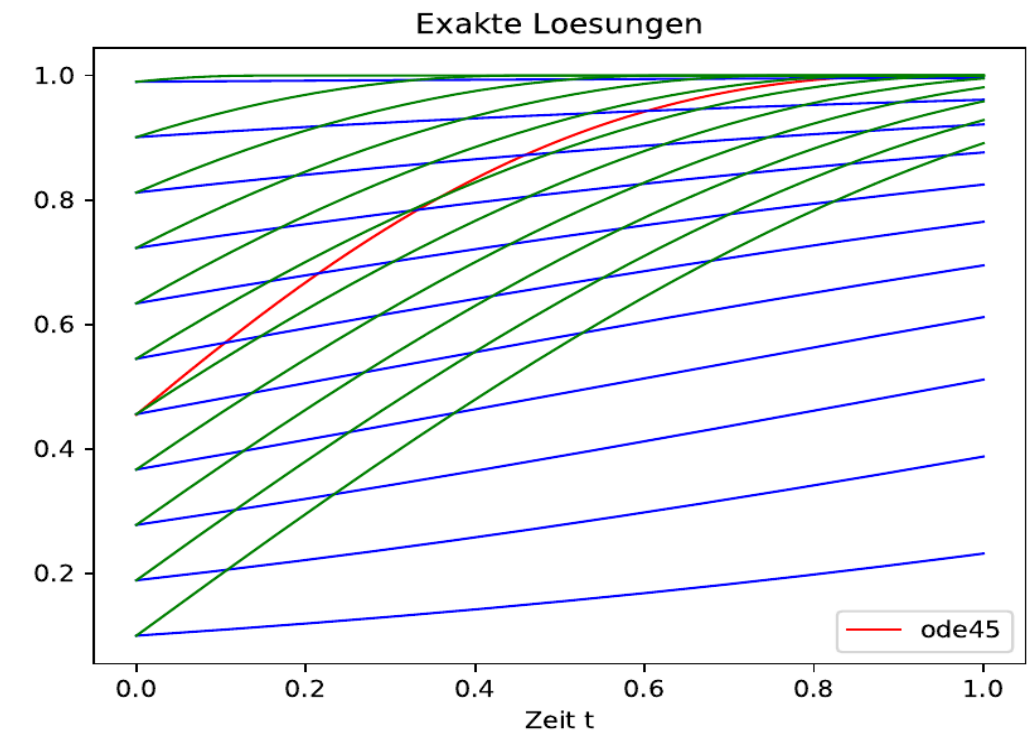
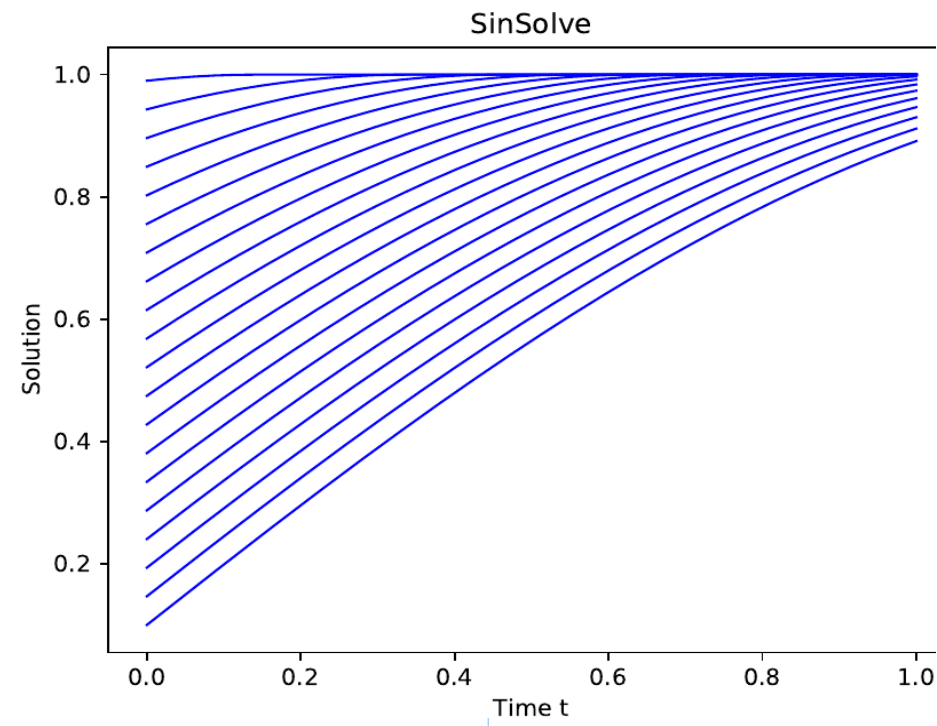
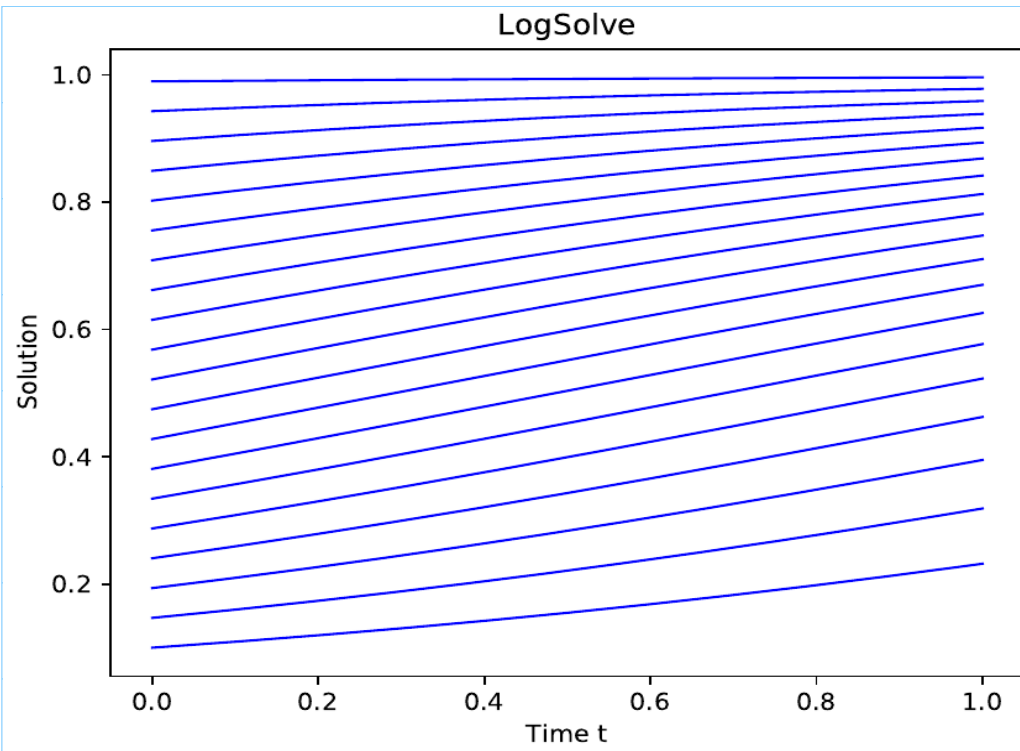


LTS-Eul	explizites Euler als $\Psi_{a,b}^h$ , $\Psi_{a,b}^h$ und Lie-Trotter-Splitting
SS-Eul	explizites Euler als $\Psi_{a,b}^h$ , $\Psi_{a,b}^h$ und Strang-Splitting
SS-EuEI	Strang-Splitting: explizites Euler als $\Psi_a^{h/2}$ , exaktes $\Phi_b^h$ und implizites Euler als $\Psi_a^{h/2}$
LTS-EMP	explizite Mittelpunkt-Regel als $\Psi_{a,b}^h$ , $\Psi_{a,b}^h$ und Lie-Trotter-Splitting
SS-EMP	explizite Mittelpunkt-Regel als $\Psi_{h,g}^h$ , $\Psi_{h,f}^h$ und Strang-Splitting

Abb. 2.4.5. Einfache Splitting-Verfahren

$$y_1 = y_0 + h f\left(t_0 + \frac{h}{2}, \underbrace{y_0 + \frac{h}{2} y_0}\right) \quad O(h^2)$$

ein halber Zeitschritt  $\in \in$  um  $y(\frac{h}{2})$  zu approximieren.



Bsp (Splitting-Verfahren für Newtonsche Gleichung)

$$\ddot{r} = a(r) \Leftrightarrow \dot{y} = \begin{bmatrix} \dot{r} \\ \dot{v} \end{bmatrix} = \begin{bmatrix} v \\ a(r) \end{bmatrix} =: F(y)$$

$$F(y) = \underbrace{\begin{bmatrix} 0 \\ a(r) \end{bmatrix}}_f + \underbrace{\begin{bmatrix} v \\ 0 \end{bmatrix}}_g \quad \text{Startwert } y_0 = \begin{bmatrix} r_0 \\ v_0 \end{bmatrix}$$

Die exakten Evolutionsoperatoren sind:

$$\begin{cases} \dot{r} = 0 \\ \dot{v} = a(r) \end{cases} \quad \begin{array}{l} \text{von } 0 \text{ zu } h: \\ \downarrow \\ \dot{v} = a(r_0) \Rightarrow v(h) = v_0 + h a(r_0) \end{array} \quad \begin{array}{l} r(h) = r(0) = r_0 \\ \downarrow \\ v(h) = v_0 + h a(r_0) \end{array}$$

$$\Phi_f^h \begin{bmatrix} r_0 \\ v_0 \end{bmatrix} = \begin{bmatrix} r_0 \\ v_0 + h a(r_0) \end{bmatrix}$$

$$\begin{cases} \dot{r} = v \\ \dot{v} = 0 \end{cases} \Rightarrow \begin{array}{l} \dot{r} = v_0 \Rightarrow r(h) = r_0 + h v_0 \\ v(h) = v_0 \end{array}$$

$$\Phi_g^h \begin{bmatrix} r_0 \\ v_0 \end{bmatrix} = \begin{bmatrix} r_0 + h v_0 \\ v_0 \end{bmatrix}$$

Kombinieren wir diese exakten Lösungen

(i) Lie-Trotter-Splittung:

$$\Psi^h \begin{bmatrix} r \\ v \end{bmatrix} = (\Phi_g^h \Phi_f^h) \begin{bmatrix} r \\ v \end{bmatrix} = \begin{bmatrix} r + h(v + h a(r)) \\ v + h a(r) \end{bmatrix}$$

symplektische Euler-Verfahren

(v) Strang-Splitting:

$$\Psi^h \begin{bmatrix} r \\ v \end{bmatrix} = \Phi_g^{h/2} \circ \Phi_f^h \circ \Phi_g^{h/2} \begin{bmatrix} r \\ v \end{bmatrix} =$$

$$= \begin{bmatrix} r + \frac{h}{2}v + \frac{1}{2}h^2 a\left(r + \frac{1}{2}hv\right) \\ v + h a\left(r + \frac{1}{2}hv\right) \end{bmatrix}$$

$\Leftrightarrow$

Notation

$$\begin{cases} r_{k+\frac{1}{2}} = r_k + \frac{1}{2}h v_k \\ v_{k+1} = v_k + h a\left(r_{k+\frac{1}{2}}\right) \\ r_{k+1} = r_k + \frac{1}{2}h v_{k+1} \end{cases}$$

$\equiv$  ein-Schritt-Formulierung des (SEV) !

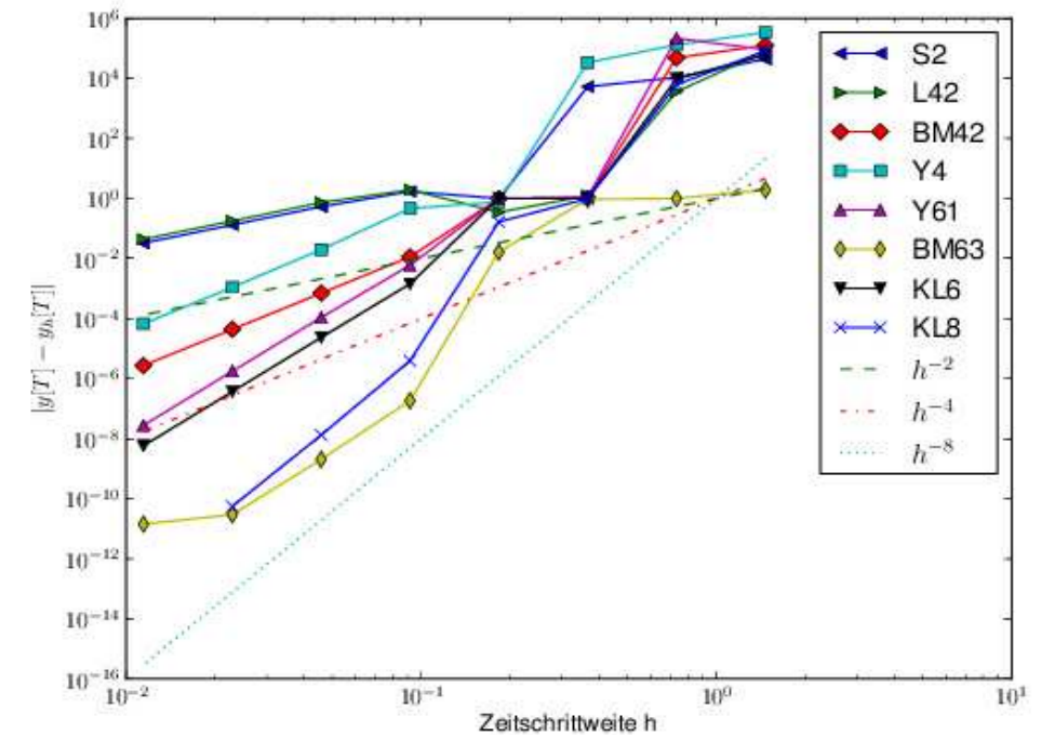
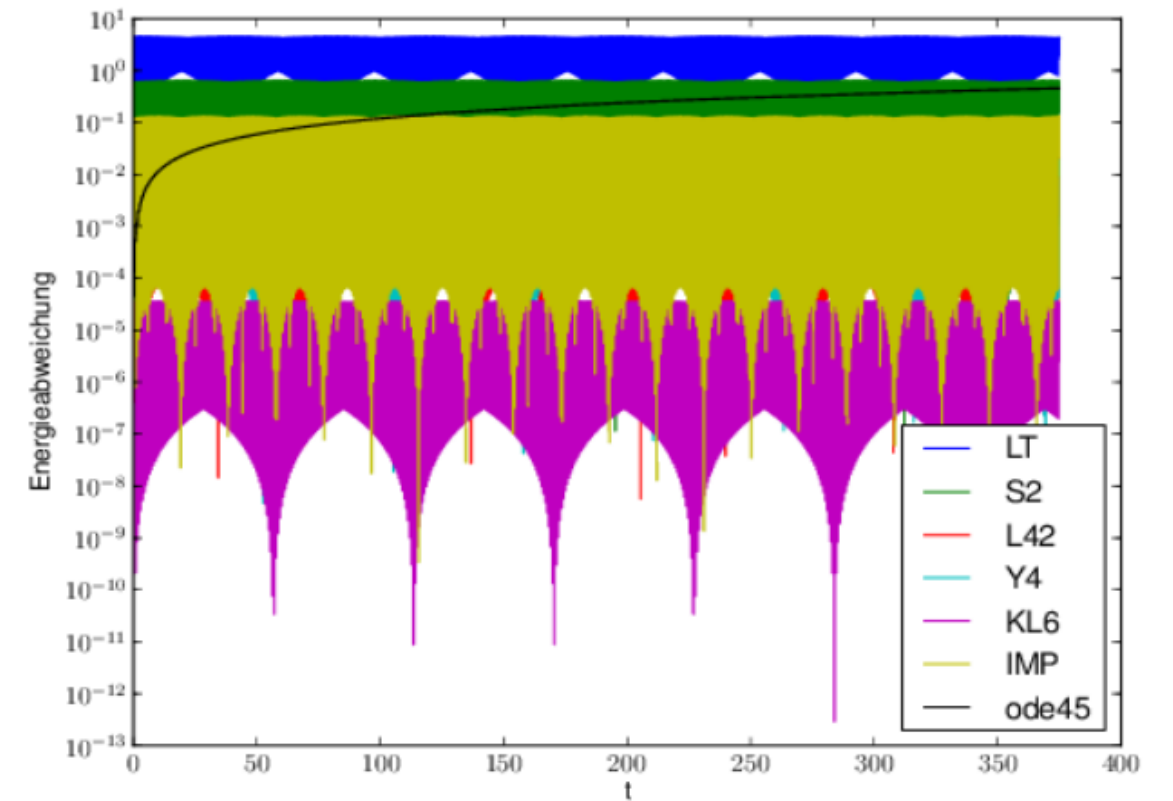
$$\begin{aligned} \Phi_g^h \begin{bmatrix} r_0 + \frac{h}{2}v_0 \\ v_0 \end{bmatrix} &= \begin{bmatrix} r_0 + \frac{h}{2}v_0 \\ v_0 + h a(r_0) \end{bmatrix} \\ \Phi_g^{h/2} \left( \underbrace{\begin{bmatrix} r_0 + \frac{h}{2}v_0 \\ v_0 \end{bmatrix}} \right) &= \Phi_g^{h/2} \begin{bmatrix} r_0 + \frac{h}{2}v_0 \\ v_0 + h a(r_0 + \frac{h}{2}v_0) \end{bmatrix} = \\ &= \begin{bmatrix} r_0 + \frac{h}{2}(v_0 + h a(r_0 + \frac{h}{2}v_0)) \\ v_0 + h a(r_0 + \frac{h}{2}v_0) \end{bmatrix} \end{aligned}$$

Ben Trick auch für separablen Hamilton-Systeme:  
 $H(\underline{p}, \underline{q}) = T(\underline{p}) + V(\underline{q})$

Lie-Trotter: 
$$\begin{cases} \underline{p}_1 = \underline{p}_0 - h \operatorname{grad} V(\underline{q}_0) \\ \underline{q}_1 = \underline{q}_0 + h \operatorname{grad} T(\underline{p}_1) \end{cases}$$

Symplektische Euler-Verfahren

Strang-Splitting  $\Rightarrow$  (St-V) !



Bsp  $H(p, q) = \frac{1}{2} p^2 - \frac{g}{l} \cos q + (-q) A \cos(\omega t)$

$$\begin{cases} \dot{p} = -\frac{\partial H}{\partial q} = -\frac{g}{l} \sin q + A \cos(\omega t) \\ \dot{q} = \frac{\partial H}{\partial p} = p \end{cases}$$

Nicht autonom!  $\Rightarrow$  autonomisieren

Unbekannte  $\dot{t} = 1$

$$\underline{u} = \begin{bmatrix} q \\ t \\ p \end{bmatrix} \Rightarrow \underline{\dot{u}} = \underline{f}(\underline{u}), \quad \underline{f}(\underline{u}) = \begin{bmatrix} p \\ 1 \\ -\frac{g}{l} \sin q + A \cos(\omega t) \end{bmatrix}$$

$$\underline{f}(\underline{u}) = \begin{bmatrix} p \\ 1 \\ 0 \end{bmatrix} + \underbrace{\begin{bmatrix} 0 \\ 0 \\ -\frac{g}{l} \sin q + A \cos(\omega t) \end{bmatrix}}_{(a)} = \underline{f}_a(\underline{u}) + \underline{f}_b(\underline{u})$$

(b)

$$(a): \begin{cases} \dot{q} = 0 \Rightarrow q(h) = q(0) = q_0 \\ \dot{t} = 0 \Rightarrow t(h) = t(0) = t_0 \\ \dot{p} = -\frac{g}{l} \sin q + A \cos(\omega t) \Rightarrow \end{cases}$$

$$\dot{p} = -\frac{g}{l} \sin q_0 + A \cos(\omega t_0) \Rightarrow$$

$$p(h) = p_0 - \left(\frac{g}{l} \sin q_0\right) h + A \cos(\omega t_0) h$$

$$(b) \begin{cases} \dot{q} = p & \Rightarrow q(h) = q_0 + p_0 h \\ \dot{t} = 1 & \Rightarrow t(h) = t_0 + h \\ \dot{p} = 0 & \Rightarrow p(h) = p_0 \end{cases}$$

$\tau = a_i h$ ,  $\tau = b_i h \Rightarrow$  kein Lieblings splitting!

Processing

$$\hat{\Psi}^h = \overset{\text{post-processor}}{\pi^h} \circ \Psi^h \circ \overset{\text{pre-processor}}{(\pi^h)^{-1}}$$

$$\begin{aligned} (\hat{\Psi}^h)^n &= \pi^h \circ \Psi^h \circ (\pi^h)^{-1} \circ \pi^h \circ \Psi^h \circ (\pi^h)^{-1} \circ \dots \circ \pi^h \circ \Psi^h \circ (\pi^h)^{-1} \\ &= \pi^h \circ (\Psi^h)^n \circ (\pi^h)^{-1} \end{aligned}$$

Processing bringt Vorteile falls:

+  $\hat{\Psi}^h$  genauer als  $\Psi^h$  ist

+  $\pi^h, (\pi^h)^{-1}$  günstig.

+ keine /wenige Ausgaben vor Endzeit

Bsp Strang-Splitting:

$$\Psi_2^h = \Phi_a^{h/2} \circ \Phi_b^h \circ \Phi_a^{h/2} \Rightarrow$$

$$\Phi_a^{h/2} \circ (\Phi_a^{h/2})^{-1}$$

$$\Rightarrow \Psi_2^h = \Phi_a^{h/2} \circ \Phi_b^h \circ \Phi_a^{h/2} \circ (\Phi_a^{h/2})^{-1}$$

Lie-Trotter

$\Rightarrow$  Strang-Splitting = Processing von Lie-Trotter!

Leicht gestörte Probleme

$$\dot{y} = f_a(y) + \varepsilon f_b(y) \quad \text{mit } \varepsilon \text{ klein}$$

optimierte Splitting-Verfahren:

$$O\left(\varepsilon h^{r_1} + \varepsilon^2 h^{r_2} + \varepsilon^3 h^{r_3} + \dots\right)$$

$$r_1 \geq r_2 + 1, \dots$$

$$\varepsilon h^4 + \varepsilon^2 h^2$$



## §4 Runge-Kutta - Verfahren

### §4.1. Grundidee

$$\begin{cases} \dot{\underline{y}} = \underline{f}(t, \underline{y}) \\ \underline{y}(t_0) = \underline{y}_0 \end{cases} \Rightarrow \underline{y}(t_1) = \underline{y}(t_0) + \int_{t_0}^{t_1} \underline{f}(t, \underline{y}(t)) dt$$

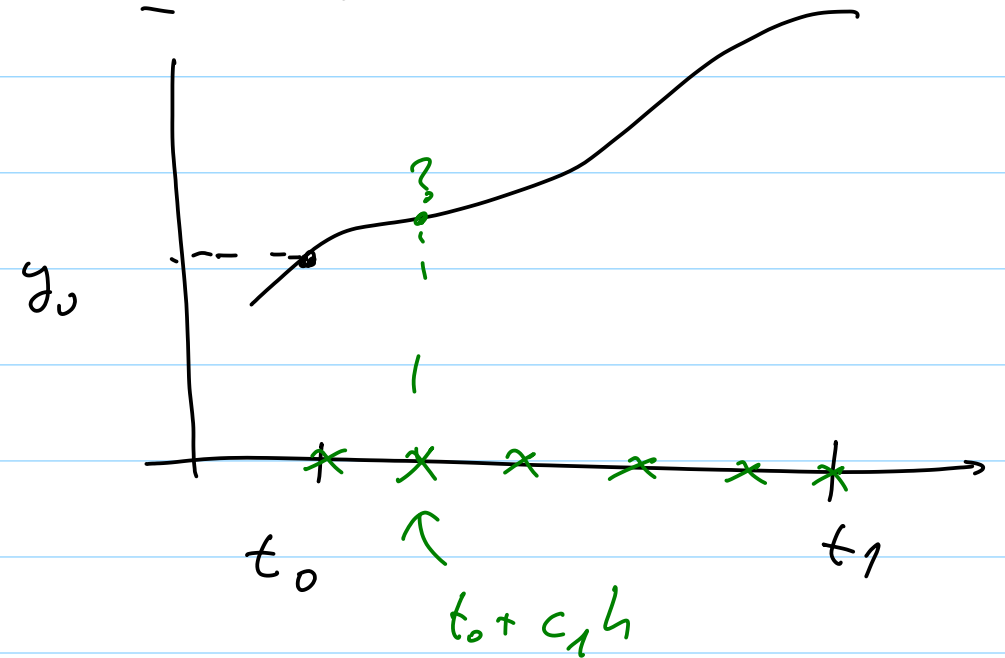
$h = t_1 - t_0$ , Referenzintervall  $[0, 1] \Rightarrow$

$$\underline{y}(t_1) = \underline{y}(t_0) + h \int_0^1 \underline{f}(t_0 + hz, \underline{y}(t_0 + hz)) dz$$

QF mit Gewichten  $b_i$ , Knoten  $c_i \Rightarrow$

$$\underline{y}(t_1) \approx \underline{y}(t_0) + h \sum_{i=1}^s b_i \underbrace{f(t_0 + hc_i, \underline{y}(t_0 + hc_i))}_{k_i}$$

$h$  erlaubt uns einen Fehler  $O(h^{p+1})$  für  $\underline{y}(t_1)$  zu bekommen, auch wenn für  $\underline{y}(t_0 + hc_i)$  Approx.  $O(h^p)$  verwende!



Bsp 1) QF = Trapezregel auf  $[0, 1]$ :

$$s=2, \quad c_1=0, \quad c_2=1, \quad b_1=b_2=\frac{1}{2}$$

$$\underline{y}_1 = \underline{y}_0 + h \left( \frac{1}{2} \underbrace{f(t_0 + h \cdot 0, \underline{y}(t_0 + h \cdot 0))}_{\substack{t_0 \\ \underline{y}(t_0) = \underline{y}_0}} + \frac{1}{2} \underbrace{f(t_0 + h \cdot 1, \underline{y}(t_0 + h \cdot 1))}_{\substack{t_1 \\ \underline{y}(t_1)}} \right)$$

Idee: verwende etwas billigeres für den inneren Term (unbekannt)

$$\underline{y}(t_0+h) \approx \underline{y}_0 + h \underbrace{f(t_0, \underline{y}_0)}_{\underline{k}_1} \quad (e \in E)$$

$$\begin{cases} \underline{k}_1 := \underline{f}(t_0, \underline{y}_0) \\ \underline{k}_2 := \underline{f}(t_0+h, \underline{y}_0 + h \underline{k}_1) \\ \underline{y}_1 = \underline{y}_0 + h \cdot \frac{1}{2} \underline{k}_1 + h \cdot \frac{1}{2} \underline{k}_2 \end{cases} \quad \text{explizite Trapezregel.}$$

Bsp 2) QF = Mittelpunktsregel

$$\underline{y}_1 = \underline{y}_0 + h \underline{f}\left(t_0 + h \frac{1}{2}, \underline{y}\left(t_0 + h \frac{1}{2}\right)\right) \quad (e \in E)$$

$$\underline{y}_0 + h \frac{1}{2} \underbrace{\underline{f}(t_0, \underline{y}_0)}_{\underline{k}_1}$$

$$\begin{cases} \underline{k}_1 := \underline{f}(t_0, \underline{y}_0) \\ \underline{k}_2 := \underline{f}\left(t_0 + \frac{h}{2}, \underline{y}_0 + \frac{h}{2} \underline{k}_1\right) \\ \underline{y}_1 = \underline{y}_0 + h \underline{k}_2 \end{cases} \quad \text{explizite Mittelpunktsregel.}$$

Def Runge-Kutta-Verfahren mit  $s$  Stufen.

Gegeben Butcher-Schema

$$\begin{array}{c|c} \begin{matrix} c_1 \\ c_2 \\ \vdots \\ c_s \end{matrix} & \underline{A} \in \mathbb{R}^{s \times s} \\ \hline \begin{matrix} 1 \end{matrix} & \begin{matrix} b_1 & b_2 & \dots & b_s \end{matrix} \end{array}$$

so dass  $b_1 + b_2 + \dots + b_s = 1$

$$\sum_{j=1}^s a_{ij} = c_i \quad \text{für } i=1,2,\dots,s$$

$$\begin{cases} \underline{k}_i = \underline{f}(t_0 + c_i \cdot h, \underline{y}_0 + h \sum_{j=1}^n a_{ij} \underline{k}_j) \\ \text{für } i=1, 2, \dots, p \end{cases}$$

Stufen bis  $j \leq i-1$  falls es explizit sein sollte!

Def Konsistenzordnung  $q$ :

lokale Fehler  $\|y(t_0+h) - y_1\| \leq c \cdot h^{q+1}$

$$\underline{y}_1 = \underline{y}_0 + h \sum_{i=1}^n b_i \underline{k}_i$$

ein  $n \times n$  nichtlineares  
algebraisches  
Gleichungssystem

Theorem

RK hat Konsistenzordnung  $q \Rightarrow$

QF hat Ordnung  $q$

(ist exakt für Polynome von Grad  $\max q-1$ )

$$\underline{A} = \begin{bmatrix} 0 & & 0 \\ & \ddots & \\ * & & 0 \\ & & & 0 \end{bmatrix}$$

$a_{ij} = 0$  für alle  $i \leq j \Rightarrow$   
 $\Rightarrow$  explizites Verfahren.

Beweis Nehme  $\begin{cases} \dot{y} = t^n \\ y(0) = 0 \end{cases} \Rightarrow y(t) = \frac{1}{n+1} t^{n+1}$

$$\underline{A} = \begin{bmatrix} & & 0 \\ * & \diagdown & \\ & & \end{bmatrix}$$

diagonal implizites Verfahren.

Fehler:  $|y(h) - y_1| = \left| \frac{1}{n+1} h^{n+1} - h \sum_{j=1}^n b_j (c_j h)^n \right| \leq c \cdot h^{q+1}$

$$\Leftrightarrow \left| \frac{1}{n+1} h^{n+1} - h^{n+1} \sum_{j=1}^n b_j c_j^n \right| \leq c \cdot h^{q+1} \quad | : h^{n+1} \Rightarrow$$

$$\left| \frac{1}{n+1} - \sum_{j=1}^n b_j c_j^n \right| \leq c \cdot h^{q-n}$$

$\downarrow$   $h \rightarrow 0$   
 $0$  solange  $q > n$

Für:  $n=0, 1, 2, \dots, q-1$ :

$$\frac{1}{n+1} = \sum_{j=1}^n b_j c_j^n$$

$\Leftrightarrow$  QF exakt für  $p_n(t) = t^n$  g.e.d.

Konsequenz RK  $s$  Stufen  $\Rightarrow$  Konsistenzordnung  $\leq 2s$

Bem 1)  $\sum_{j=1}^n b_j = 1 \Rightarrow$  mindestens Konsistenzordnung  $q=1$

2) RK hat mindestens Konsistenzordnung  $q=2$

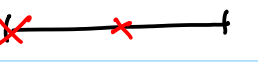
wenn  $\sum_{j=1}^n b_j c_j = \frac{1}{2}$

3)  $q=3$  brauchen wir

$$\sum_{j=1}^n b_j c_j^2 = \frac{1}{3}, \quad \sum_{j=1}^n b_j \sum_{n=1}^n a_{jn} c_n = \frac{1}{6}$$

Theorem RK explizit  $\Rightarrow q \leq s$

Gauss-Quadratur  $\Rightarrow q=2s$  

Radau-Quadratur  $\Rightarrow$  

Radau-Verfahren für ODE



Lobatto-Quadratur

$\hookrightarrow$  Lobatto-Verfahren für ODE

Theorem RK hat Konsistenzordnung  $q \Rightarrow$   
(globale) Konvergenzordnung  $q$

$$\|y(t_i) - y_i\| \leq c \cdot h^q \quad \text{für } i = 1, 2, \dots, N.$$

## §4.2. Kollokation

Def  $c_1, c_2, \dots, c_s \in [0, 1]$  verschieden.

Kollokationspolynom  $u(t)$  von Grad  $s$ :

$$\begin{cases} u(t_0) = y_0 \\ \dot{u}(t_0 + c_i h) = f(t_0 + c_i h, u(t_0 + c_i h)) \end{cases} \quad \text{für } i=1, 2, \dots, s$$

Bsp 1)  $s=1 \Rightarrow$  Polynom vom Grad 1:

$$u(t) = y_0 + (t - t_0)k$$

mit  $k$  bestimmt so, dass

$$\dot{u}(t_0 + c_1 h) = f(t_0 + c_1 h, u(t_0 + c_1 h))$$

$$c_1 = 0 \Rightarrow (E)$$

$$c_1 = 1 \Rightarrow (iE)$$

$$c_1 = \frac{1}{2} \Rightarrow (iMP)$$

$$2) \quad s=2 ; \quad c_1=0, c_2=1 \Rightarrow \text{implizite Trapezregel}$$

$$c_{1,2} = \frac{1}{2} \pm \frac{\sqrt{3}}{6} \Rightarrow \text{Gauss-Verfahren } O(h^4)$$

Theorem Die Kollokation mit Knoten  $c_1, \dots, c_s$

$$\begin{aligned} & \xLeftrightarrow{\quad} \text{\textit{s}-Stufiges RKV mit } a_{ij} = \int_0^{c_i} l_j(z) dz, \\ & b_i = \int_0^1 l_i(z) dz \end{aligned}$$

wobei

$$l_i(z) = \prod_{\substack{j=1 \\ j \neq i}}^s \frac{z - c_j}{c_i - c_j} \quad \text{Lagrange Polynom.}$$

$$l_i(c_j) = \begin{cases} 1, & i=j \\ 0, & i \neq j \end{cases}$$

Beweis Notiere  $k_i := \dot{u}(t_0 + c_i h)$ .

Eigenschaft von  $l_i \Rightarrow \dot{u}(t_0 + z h) = \sum_{j=1}^n k_j \cdot l_j(z) \int_0^{c_i}$

Polynom vom Grad  $n-1$

$$\Rightarrow \underline{u(t_0 + c_i h)} = y_0 + h \sum_{j=1}^n k_j \underbrace{\int_0^{c_i} l_j(z) dz}_{a_{ij}}$$

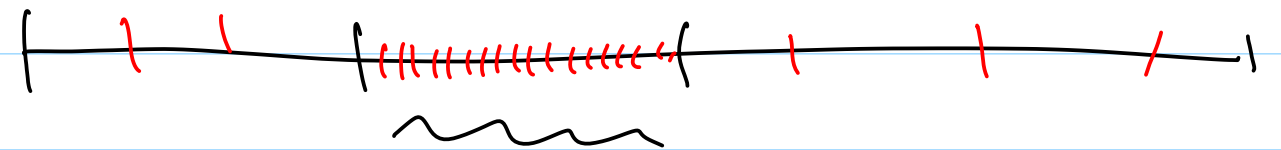
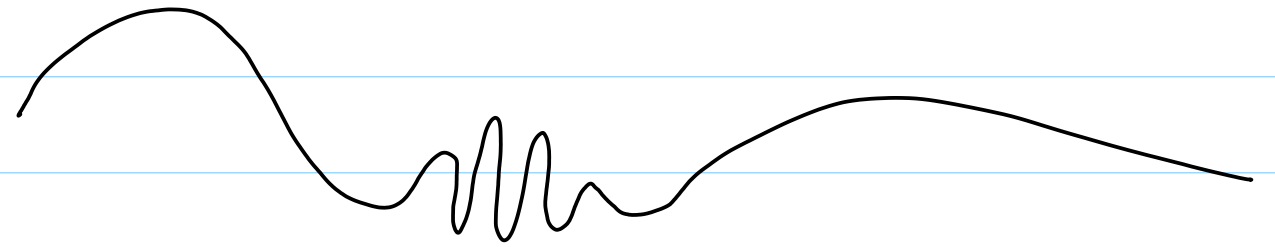
zwischen stellen.

$$\int_0^1 \Rightarrow y_1 = y_0 + h \sum_{j=1}^n k_j \underbrace{\int_0^1 l_j(z) dz}_{b_j}$$

Konsequenz

Kollokationsmethode hat dieselbe Ordnung wie die entsprechende QF.

### §4.3. Adaptivität



lokaler Fehler.

$$\Phi - \Psi_h \Rightarrow \text{Schätzung } \tilde{\Psi}_h - \Psi_h$$

$\tilde{\Psi}_h$  (genauere Methode,  $O(h^{p+2})$ )  
 $\Psi_h$  ( $O(h^{p+1})$ )

$$\left| \tilde{\Psi}_{y(t_k), t, t+h} - \Psi_{y(t_k), t, t+h} \right| = \text{est}_k$$

$$\underline{\text{est}_k} = c h^{p+1} = \text{tol} \Rightarrow h^* := h \sqrt[p+1]{\frac{\text{tol}}{\text{est}_k}}$$

SKRIPT!

# §4.4. Partitionierte RK-Verfahren

System ODE, so partitioniert:

$$\begin{cases} \dot{\underline{y}} = f(\underline{y}, \underline{z}) \\ \dot{\underline{z}} = g(\underline{y}, \underline{z}) \end{cases}$$

Idee: verwende 2 verschiedene RKV für  $\underline{y}$  und  $\underline{z}$ :

$$\begin{cases} \underline{k}_i = f\left(\underline{y}_0 + h \sum_{j=1}^s a_{ij} \underline{k}_j, \underline{z}_0 + h \sum_{j=1}^s \hat{a}_{ij} \underline{l}_j\right) \\ \underline{l}_i = g\left(\underline{y}_0 + h \sum_{j=1}^s a_{ij} \underline{k}_j, \underline{z}_0 + h \sum_{j=1}^s \hat{a}_{ij} \underline{l}_j\right) \end{cases}$$

für  $\underline{y}$ :  $\begin{array}{c|c} c & A \\ \hline 1 & b \end{array}$       für  $\underline{z}$ :  $\begin{array}{c|c} \hat{c} & \hat{A} \\ \hline 1 & \hat{b} \end{array}$

$$\underline{y}_1 = \underline{y}_0 + h \sum_{j=1}^s b_j \underline{k}_j$$

$$\underline{z}_1 = \underline{z}_0 + h \sum_{j=1}^s \hat{b}_j \underline{l}_j$$

Bsp 1)

(EE) mit  $b_1=1, a_{11}=1$  }  $\Rightarrow$  symplektische Euler-Verfahren  
(EE)  $b_1=1, \hat{a}_{11}=0$  } (für Newton-Gleichung)

Bsp 2)

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & \frac{1}{2} & \frac{1}{2} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

$$\begin{array}{c|cc} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

PRK  $\Rightarrow$

Störmer-Verlet!

Eine Verallgemeinerung vom Störmer-Verlet  
kommt hier natürlich heraus, aus  
3-stufiges Lobatto-Proc: (Ordnung  $h^4$ )

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{5}{24} & \frac{1}{3} & -\frac{1}{24} \\ 1 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \\ \hline & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array}$$

$$\begin{array}{c|ccc} 0 & \frac{1}{6} & -\frac{1}{6} & 0 \\ \frac{1}{2} & \frac{1}{6} & \frac{1}{3} & 0 \\ 1 & \frac{1}{6} & \frac{5}{6} & 0 \\ \hline & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array}$$

Ben Newton-Gleichung:

$$\ddot{\underline{y}} = \underline{g}(t, \underline{y}, \dot{\underline{y}})$$

Umschreiben in System 1. Ordnung:

$$\begin{cases} \dot{\underline{y}} = \underline{z} \\ \dot{\underline{z}} = \underline{g}(t, \underline{y}, \underline{z}) \end{cases}$$

PRK für  $\nearrow$  RK-Nyström-Verfahren (RKN)

$$\begin{cases} \underline{l}_i = \underline{g}(t, \underline{y}_0 + c_i h \dot{\underline{y}}_0 + h^2 \sum_{j=1}^n \bar{a}_{ij} \underline{l}_j, \dot{\underline{y}}_0 + h \sum_{j=1}^n \hat{a}_{ij} \underline{l}_j) \\ \underline{y}_1 = \underline{y}_0 + h \dot{\underline{y}}_0 + h^2 \sum_{i=1}^n \bar{b}_i \underline{l}_i \\ \dot{\underline{y}}_1 = \dot{\underline{y}}_0 + h \sum_{i=1}^n \hat{b}_i \underline{l}_i \end{cases}$$



mit

$$\bar{a}_{ij} = \sum_{k=1}^n a_{ik} \hat{a}_{kj}, \quad \bar{b}_i = \sum_{k=1}^n b_k \hat{a}_{ki}$$

Wenn  $g$  nicht noch von  $j$  abhängt, dann  
dann braucht man  $\hat{a}_{kj}$  gar nicht.

Bem PRK  $\equiv$  Splitting mit

$$\underline{u} = \begin{bmatrix} \underline{y} \\ \underline{z} \end{bmatrix}, \quad \underline{f}_a = \begin{bmatrix} f(\underline{u}) \\ 0 \end{bmatrix}, \quad \underline{f}_b = \begin{bmatrix} 0 \\ g(\underline{u}) \end{bmatrix}$$

$\Rightarrow$  einfachere Anwendung!

$$\begin{array}{ll} \text{BM}_{42} & \rightarrow O(h^4) \\ \text{BM}_{63} & \rightarrow O(h^6) \end{array} \quad \begin{array}{l} \text{symplektische} \\ \text{Verfahren.} \end{array}$$

## §5 Steife Differentialgleichungen

## §5.1. Einführung

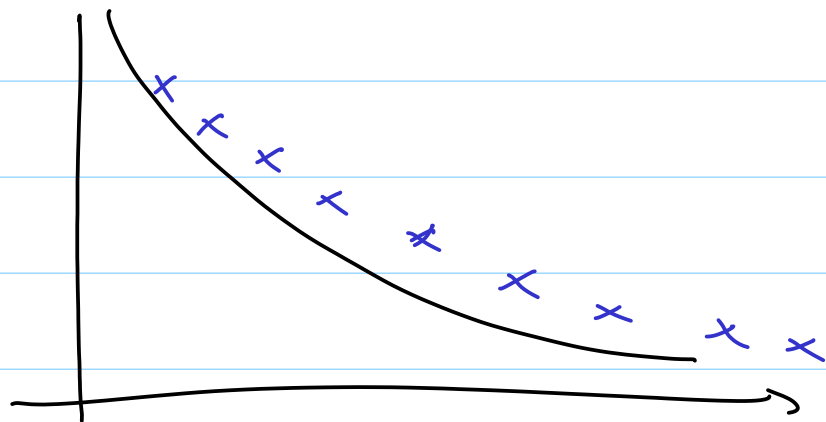
Modelproblem:  $\dot{y} = -\lambda y$  mit  $\lambda > 0$

$$y(t) = e^{-\lambda t} y_0 \rightarrow 0 \text{ für } t \rightarrow \infty$$

Der Abfall der Lösung ist schneller für grösseres  $\lambda$ .

Numerische Lösung: festes Zeitschritt  $h$

$$\left( y_{\lambda}^{kh} \right)_{k=1,2,3,\dots} \rightarrow 0 \text{ für } k \rightarrow \infty ?$$



$$(2E) \quad y_1 = y_0 + h f(y_0) = y_0 - h \lambda y_0 = (1 - h \lambda) y_0$$

$$y_2 = y_1 - h \lambda y_1 = (1 - h \lambda) y_1 = (1 - h \lambda)^2 y_0$$

...

$$y_k = (1 - h \lambda)^k y_0$$

$$\text{für } k \rightarrow \infty : \quad |1 - h \lambda| = 1 \Rightarrow y_k = \dots = y_0 \quad \text{☹️}$$

$$|1 - h \lambda| < 1 \Rightarrow y_k \rightarrow 0 \quad \text{😊}$$

$$|1 - h \lambda| > 1 \Rightarrow y_k \rightarrow \infty \quad \text{☹️}$$

$$\Rightarrow \text{Bedingung an } h: \quad |1 - h \lambda| < 1 \Rightarrow 0 < h < \frac{2}{\lambda} \quad \begin{matrix} \uparrow \\ \lambda > 0 \end{matrix}$$

Def ODE heisst steif, falls explizite Methoden einen Zeitschritt  $h$  sehr klein brauchen, kleiner als die Genauigkeit verlangt!

Bsp

$$\frac{d}{dt} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} -50 & 49 \\ 49 & -50 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$\dot{\underline{y}} = \underline{B} \underline{y} \quad \underline{B} \text{ symmetrisch: es gibt } \underline{S} \text{ so dass } (\underline{S} \underline{S}^T = \underline{I})$$

$$\underline{B} = \underline{S} \underline{D} \underline{S}^T$$

mit  $\underline{D} = \text{Diagonalmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}$

$$\dot{\underline{y}} = \underline{S} \underline{D} \underline{S}^T \underline{y} \Rightarrow \frac{d}{dt} (\underline{S}^T \underline{y}) = \underline{D} (\underline{S}^T \underline{y}) \quad \Rightarrow$$

$$\underline{z} = \underline{S}^T \underline{y}$$

$$\dot{\underline{z}} = \underline{D} \underline{z} \quad (\Rightarrow \begin{cases} \dot{z}_1 = \lambda_1 z_1 \\ \dot{z}_2 = \lambda_2 z_2 \\ \vdots \\ \dot{z}_d = \lambda_d z_d \end{cases})$$

$$\underline{D} = \begin{bmatrix} -1 & 0 \\ 0 & -99 \end{bmatrix}$$

$$\begin{cases} \dot{z}_1 = -z_1 \Rightarrow h < 2 \\ \dot{z}_2 = -99 z_2 \Rightarrow h < \frac{2}{99} \end{cases}$$

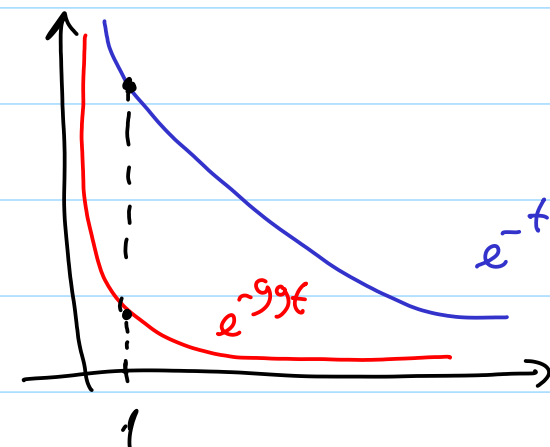
$$\begin{cases} y_1(t) = e^{-t} + e^{-99t} \\ y_2(t) = e^{-t} - e^{-99t} \end{cases}$$

Zur Zeit  $t=1$ 

$$y_1(1) = e^{-1} + \underbrace{e^{-99}}_{\text{sehr klein}}$$

sehr klein  $\Rightarrow$  für die exakte Lösung irrelevant!

aber ein explizites Verfahren will. ein sehr kleines  $h$ , diktiert genau von  $e^{-99t}$ !



$$(iE) \quad y_1 = y_0 + h f(y_1) = y_0 - \lambda h y_1 \Rightarrow$$

$$\Rightarrow (1 + \lambda h) y_1 = y_0 \Rightarrow y_1 = \frac{y_0}{1 + \lambda h}$$

$(h > 0, \lambda > 0)$

$$y_2 = \frac{1}{1 + \lambda h} y_1 = \frac{1}{(1 + \lambda h)^2} y_0$$

$\Rightarrow \dots$

$$y_k = \left( \frac{1}{1 + \lambda h} \right)^k y_0 \rightarrow 0 \quad \text{für } k \rightarrow \infty$$

$$\lambda > 0, h > 0 \Rightarrow 0 < \frac{1}{1 + \lambda h} < 1 \quad \text{egal was für } h!$$

Ben Implizite Verfahren stellen keine Bedingung an  $h$ .

Bsp explizite Trapezregel:

$$\begin{cases} k_1 = -\lambda y_0 \\ k_2 = -\lambda(y_0 + h k_1) \end{cases} \quad \dot{y} = -\lambda y$$

$$y_1 = \underbrace{\left[ 1 - \lambda h + \frac{1}{2}(\lambda h)^2 \right]}_{S(\lambda h)} y_0 = S(\lambda h) y_0$$

$$\Rightarrow \dots \quad y_k = S(\lambda h)^k y_0 \rightarrow 0 \quad \text{nur wenn}$$

$$|S(\lambda h)| < 1$$

mit s Stufen

explizite RK:  $S(z) = \text{Polynom vom Grad } s$

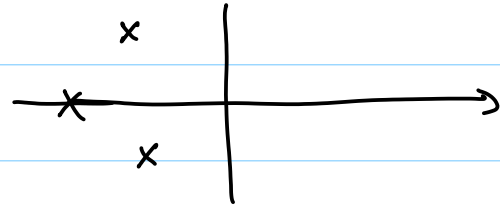
$$\text{implizite RK: } S(z) = \frac{p(z)}{q(z)} \quad \text{mit } p, q$$

Polynome vom Grad  $\leq s$ .

## § 5.2. Stabilität des RK-Verfahrens

Testproblem  $\dot{y} = \lambda y$ ,  $\lambda \in \mathbb{C}$ ,  $\operatorname{Re} \lambda < 0$

$$y(t) = e^{\lambda t} y_0$$



Bsp

Falls  $\lambda = si \in \mathbb{C} \Rightarrow y(t) = e^{sit} y_0 = y_0 (\cos t + i \sin t)$

Falls  $\lambda = -1 + si \in \mathbb{C} \Rightarrow y(t) = \underbrace{e^{-t}}_{\rightarrow 0 \text{ } t \rightarrow \infty} y_0 (\cos t + i \sin t)$

$$|y(t)| \rightarrow 0 \text{ für } t \rightarrow \infty$$

Numerische Verfahren:  $y_n = S(\lambda h)^n y_0$

Frage: wann  $|y_n| \rightarrow 0$  für  $n \rightarrow \infty$ ?

$S(z)$  = Stabilitätsfunktion ( $z = \lambda h$ )

$$(eE): S(z) = 1 - z$$

$$(eE): S(z) = \frac{1}{1+z}$$

$$(eTR): S(z) = 1 - z + \frac{1}{2} z^2$$

RK-Verfahren mit  $s$  Stufen:

$$y_1 = y_0 + h \sum_{i=1}^s b_i k_i$$

$$\begin{cases} k_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^s a_{ij} k_j) \\ i = 1, 2, \dots, s \end{cases}$$

$$\begin{array}{c|c} c & A \\ \hline & b \end{array}$$

Für das Testproblem  $\dot{y} = \lambda y \rightarrow$

$$\begin{cases} k_i = \lambda y_0 + \underbrace{(\lambda h)}_z \sum_{j=1}^s a_{ij} k_j \\ i = 1, 2, \dots, s \end{cases}$$

$$\underline{k} = \begin{bmatrix} k_1 \\ \vdots \\ k_s \end{bmatrix}$$

$$\text{LGS: } \underline{k} = \lambda y_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + z \underline{A} \underline{k} \quad (\Leftrightarrow)$$

$$\left( \underline{I} - z \underline{A} \right) \underline{k} = \lambda y_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

falls  $\underline{I} - z \underline{A}$  nicht invertierbar  $\rightsquigarrow$

$\underline{I} - z \underline{A}$  invertierbar  $\Rightarrow$

$$\underline{k} = \lambda y_0 \left( \underline{I} - z \underline{A} \right)^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \Rightarrow$$

$$y_1 = y_0 + h \lambda y_0 \sum_{i=1}^n b_i \left( \left( \underline{I} - z \underline{A} \right)^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \right)_i \Rightarrow$$

$$y_1 = y_0 \left( 1 + z \underline{b}^T \left( \underline{I} - z \underline{A} \right)^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \right) \Rightarrow$$

$$S(z) = 1 + z \underline{b}^T \left( \underline{I} - z \underline{A} \right)^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

Stabilitätsfunktion des  $n$ -stufiges RKV

Theorem Falls  $\underline{I} - z \underline{A}$  invertierbar ist, dann ist  $S$  eine (komplexwertige) rationale Funktion

$$S(z) = \frac{P(z)}{Q(z)} \quad \text{mit } P, Q \text{ Polynome}$$

von Grad  $\leq n$  und  $Q(z) = 0$  für  $z = \frac{1}{\mu}$  mit  $\mu$  Eigenwert von  $\underline{A}$ .

Falls RK-explizit, dann  $Q(z) \equiv 1$

Beweis

RK-Verfahren explizit:  $\underline{A} = \begin{bmatrix} 0 & \dots & 0 \\ * & \ddots & 0 \end{bmatrix}$   
 $n \times n$

$$\underline{A}^n = \underline{A} \cdot \underline{A} \cdot \dots \cdot \underline{A} = \begin{bmatrix} 0 & \dots & 0 \\ 0 & \dots & 0 \end{bmatrix}$$

$$(\underline{I} - z \underline{A})^{-1} = \underline{I} + z \underline{A} + (z \underline{A})^2 + \dots + (z \underline{A})^{n-1} +$$

$$+ \underbrace{(z \underline{A})^n}_{0} + \dots + \underbrace{(z \underline{A})^{n-1}}_{0}$$

$$= \underline{I} + z \underline{A} + z^2 \underline{A}^2 + \dots + z^{n-1} \underline{A}^{n-1}$$

$$= \text{Polynom vom Grad } n-1 \text{ (in } z \text{)}$$

Cramer'sche Regeln

$$\underline{v} = (\underline{I} - z \underline{A})^{-1} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \Rightarrow v_i(z) = \frac{p_n(z)}{\det(\underline{I} - z \underline{A})}$$

$$Q(z) = 0 \Leftrightarrow \det(\underline{I} - z \underline{A}) = 0 \Leftrightarrow z = \frac{1}{\mu}$$

mit  $\mu \in W$  von  $\underline{A}$ .

Konsequenz

1) Falls  $\frac{1}{\lambda h}$  nicht  $\in W$  von  $\underline{A}$ , dann ist RK:

$$y_n = S(\lambda h)^n y_0 \text{ mit } n=0,1,2,\dots$$

mit wohldefinierte Stabilitätsfunktion  $S(z)$

2)  $y_0 = 1 \Rightarrow$  exakte Lösung  $y(t) = e^{\lambda t}$

$$\text{num. Lösung } y_n = S(\lambda h)^n$$

RK mit Konvergenzordnung  $q$ :

$$\text{Fehler } |e^{n\lambda h} - S(\lambda h)^n| \leq O(|\lambda h|^{q+1})$$

$$|e^z - S(z)| \leq O(|z|^{q+1}) \Rightarrow$$

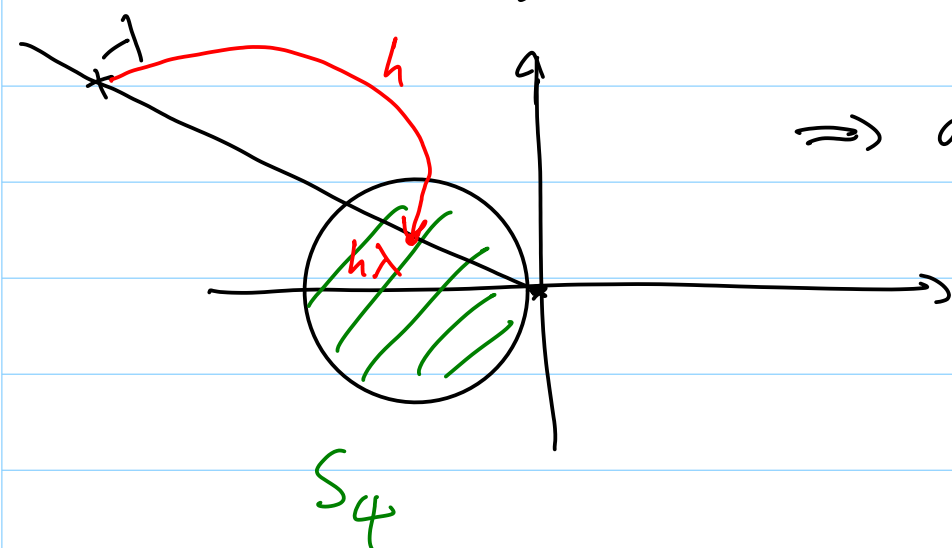
Stabilitätsfunktion  $\approx e^z$  und.

Taylorpolynome von  $e^z$  und  $S(z)$  um  $z=0$   
identisch bis zum Grad  $q$

Definition Stabilitätsgebiet des Verfahrens

$$S_4 = \{z \in \mathbb{C} \text{ so dass } |S(z)| < 1\}$$

$$y_n = S(z)^n y_0 \rightarrow 0 \text{ nur wenn } |S(z)| < 1$$

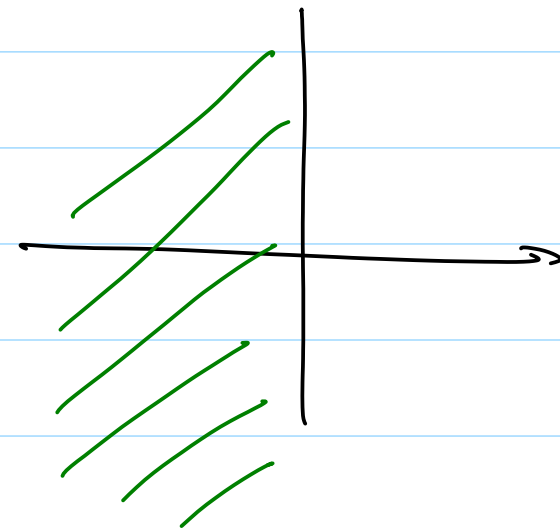


$\Rightarrow$  ablesen, wie klein  
muss  $h$  gewählt  
werden.

Ben RK explizit  $\Rightarrow S_4$  beschränkt.  $\Rightarrow$   
immer eine Schranke an  $h$ .

Ben

ode45 / ode45 verwenden explizite RK  
 $\Rightarrow S_4$  beschränkt  $\Rightarrow$  brauchen  
 $h$  klein!



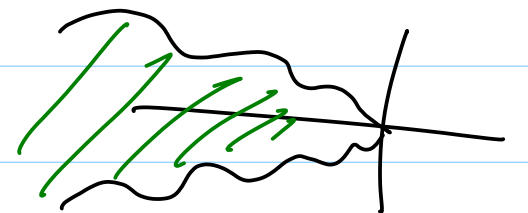
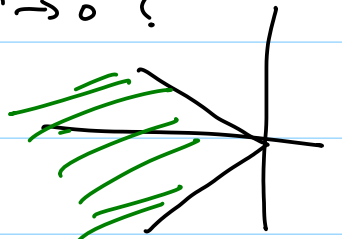
Def Ein Verfahren heisst  
A-stabil falls

$$\{z \in \mathbb{C} : \operatorname{Re} z < 0\} \subset S_4.$$

Def  $S(z) \approx e^z$ ;  $z \rightarrow -\infty$ ,  $S(-\infty) \rightarrow 0$ ?

Verfahren heisst L-stabil

falls  $\lim_{z \rightarrow -\infty} S(z) = 0$





$$\dot{y} = f(y) \quad f \text{ nicht linear}$$

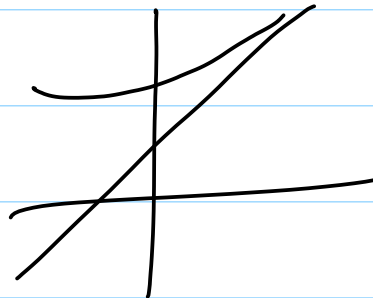
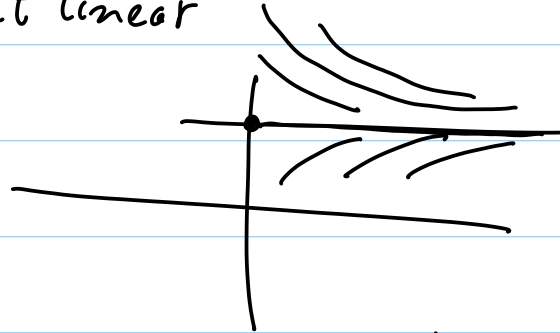
$$\underline{y}^* \in \mathbb{R}^d \text{ so dass}$$

$$f(\underline{y}^*) = 0 \quad \text{Stationär!}$$

linearisiere  
(Taylor um  $\underline{y}^*$  von  $f$ )

$$\dot{z} = \underline{D}f(\underline{y}^*) z.$$

Testproblem!



Solche Verfahren mit maximaler

Konvergenzordnung  $2n-1$ :

Radau-Quadratur / Radau-Verfahren.

Radau-Verfahren von Ordnung 3, 5  $\rightarrow$  Skript

Ben in jedem Zeitschritt:

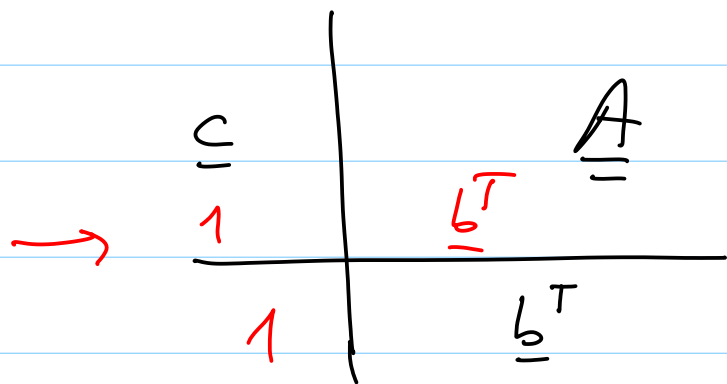
noch zu lösen: ein nicht-lineares  
algebraisches System mit  $d$ s Gleichungen

( $n$  Stufen:  $\underline{k}_1, \dots, \underline{k}_n \in \mathbb{R}^d$ )

$\hookrightarrow$  sehr teuer.



Ben RK L-stabil falls  $\underline{b}^T = a_{n,n} =$   
= letzte Zeile von  $\underline{A}$



$$c_s = 1$$



rechte Ecke des Referenzintervalls  
ist ein Quadraturknoten

$F(\underline{x}) = 0$ , solve mit Startwert

Idee: nicht-lineare alg. System  
linearisiere!

Ersetze:

$$k_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^s a_{ij} k_j)$$

falls autonom

durch:

$$\begin{cases} k_i = f(y_0) + h \underline{D} f(y_0) \left( \sum_{j=1}^s a_{ij} k_j \right) \end{cases}$$

$\mathbb{R}^d$        $\mathbb{R}^d$        $\mathbb{R}^{d \times d}$        $\in \mathbb{R}^d$

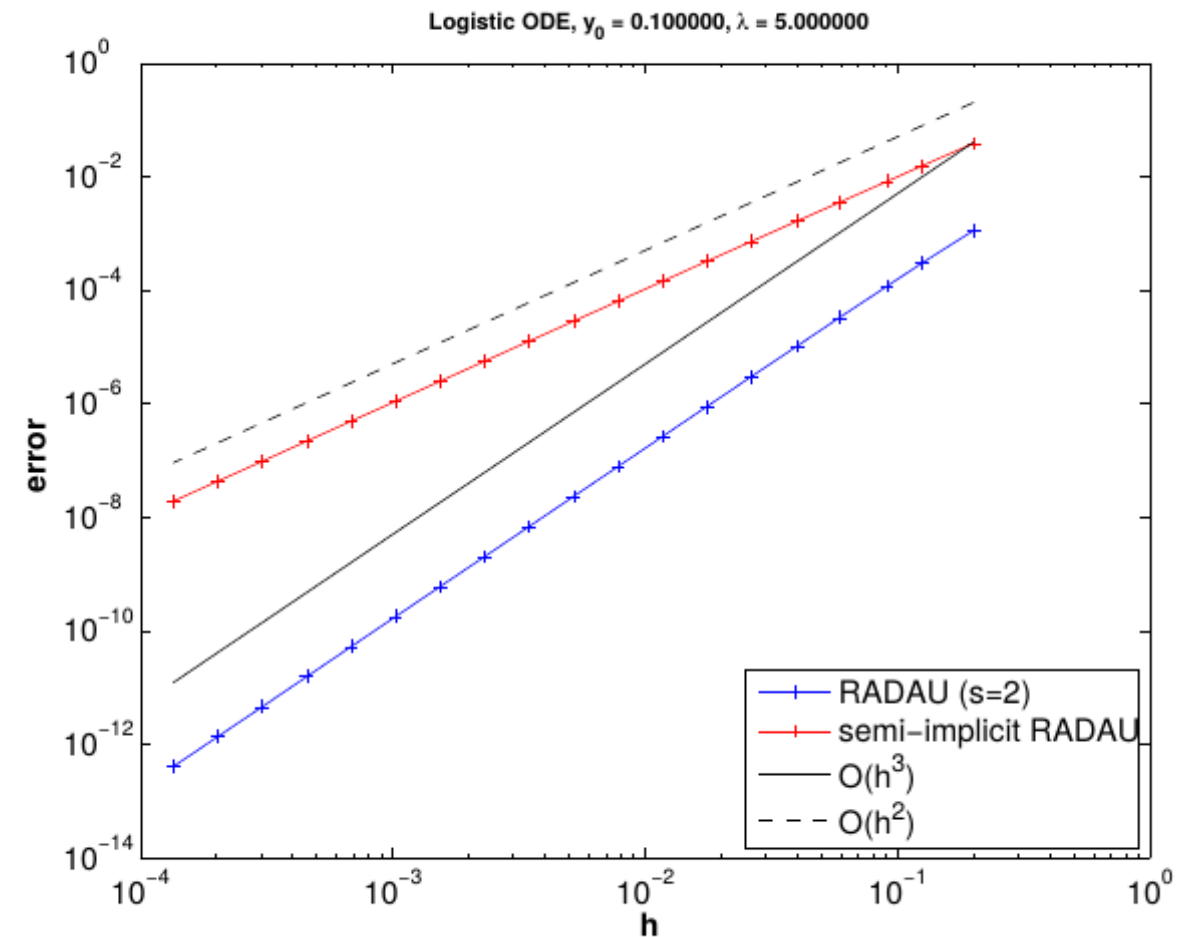
ein grosses lineares Gleichungssystem  $\mathbb{R}^{dxs}$



$$\underline{D} f(x) = \left[ \frac{\partial f_k}{\partial x_l}(x) \right]_{k,l=1,2,\dots,d}$$

Bsp  $\begin{cases} \dot{y} = \lambda y(1-y) \\ y(0) = 0.1 \end{cases} \quad \lambda = 5$

Skript: man verliert eine Konvergenzordnung!



Idee: Verwende ein Newton-Schritt für das nicht-lineare System  $\underline{k} = \underline{f}(\underline{y}_0 + h\underline{k})$

Bsp (iE)  $\underline{y}_1 = \underline{y}_0 + h \underline{f}(\underline{y}_1)$

$$\underline{F}(\underline{z}) = \underline{z} - \underline{y}_0 - h \underline{f}(\underline{z})$$

Löse  $\underline{F}(\underline{z}) = \underline{0}$  via Newton-Iteration (feuert  
mit Startwert  $\underline{z}_0 = \underline{y}_0$  (nur 1 Schritt))

$$\underline{D}\underline{F}(\underline{y}_0) \Rightarrow \underline{z}_1 = \underline{z}_0 - \underline{D}\underline{F}(\underline{z}_0)^{-1} \underline{F}(\underline{z}_0) \Rightarrow$$

$$\underline{y}_1 = \underline{y}_0 + [\underline{I} - h \underline{D}\underline{f}(\underline{z}_0)]^{-1} h \underline{f}(\underline{y}_0)$$

Brauche zusätzlich ein mehr guten Startwert für Newton  $\Rightarrow$  Ordnung retten!  
Das geht bei diagonal-implizite RK!

$$\underline{A} = \begin{bmatrix} \text{triangle} & 0 \end{bmatrix} \Rightarrow \text{gestaffeltes System}$$

$$\underline{k}_i = \underline{f}\left(\underline{y}_0 + h \sum_{j=1}^i a_{ij} \underline{k}_j\right)$$

$$\underline{F}(\underline{k}) = \underline{k} - \underline{f}\left(\underline{y}_0 + \underline{z} + h a_{ii} \underline{k}\right)$$

wobei  $\underline{z} = h \sum_{j=1}^{i-1} a_{ij} \underline{k}_j$

Newton-Schritt:

$$\underline{D}\underline{F}(\underline{k}) = \underline{I} - \underline{D}\underline{f}\left(\underline{y}_0 + \underline{z} + h a_{ii} \underline{k}\right) h a_{ii}$$

Startwert  $\underline{k}^{(0)}$  in Newton:

$$\underline{k}_i^{(0)} = \sum_{j=1}^{i-1} \frac{d_{ij}}{a_{ij}} \underline{k}_j$$

Spezielle Wahl von  $a_{ij}$ ,  $d_{ij}$  rettet die  
Konvergenzordnung.

linear-implizite RK Rosenbrock-Wanner-Methoden  
(ROW-Methoden)

$$\left( \underline{I} - h a_{ii} \underline{\partial} \right) \underline{k}_i = f \left( y_0 + h \sum_{j=0}^{i-1} (a_{ij} + d_{ij}) \underline{k}_j \right) - h \underline{\partial} \sum_{j=1}^{i-1} d_{ij} \underline{k}_j$$

mit

$$\underline{\partial} = \underline{D} f \left( y_0 + h \sum_{j=1}^{i-1} (a_{ij} + d_{ij}) \underline{k}_j \right)$$

wie bei vereinfachten Newton  
lassen wir das weg!

$\Rightarrow$  ROW2, ROW3  $\Rightarrow$  ode23s nicht in standard  
Rythm.

# §6 Nichtlineare algebraische Gleichungen

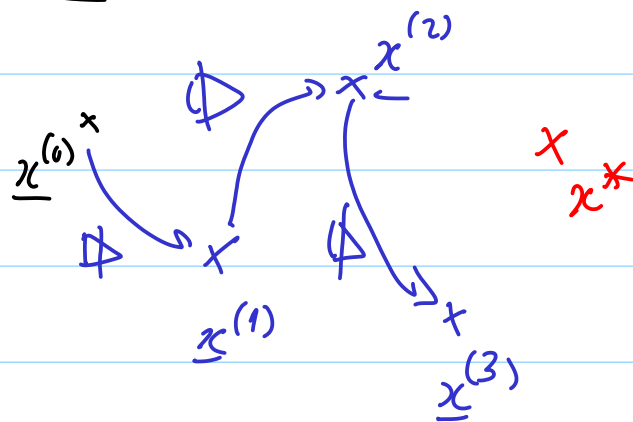
## §6.1. Einführung.

Finde  $\underline{x}^* \in \mathbb{R}^d$  so dass  $\underline{F}(\underline{x}^*) = \underline{0}$

Bsp  $d=1$ ,  $F(x) = x e^x - 1$  ( $\because x^* e^{x^*} - 1 = 0$ )

Iterativ:  $\underline{x}^{(0)}, \underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(k)} \rightarrow \underline{x}^*$

$\underline{x}^{(0)}$  gegeben  $\underline{x}^{(k+1)} = \Phi(\underline{x}^{(k)})$



Def  $\underline{x}^{(k+1)} = \Phi(\underline{x}^{(k)})$  heisst  
linear konvergent nach  $\underline{x}^*$

falls

es gibt  $L < 1$  so dass

$$\|\underline{x}^{(k+1)} - \underline{x}^*\| \leq L \|\underline{x}^{(k)} - \underline{x}^*\| \text{ für alle } k \in \mathbb{N}.$$

Bem

$$\begin{aligned} \|\underline{x}^{(k+1)} - \underline{x}^*\| &\leq L \|\underline{x}^{(k)} - \underline{x}^*\| \leq L^2 \|\underline{x}^{(k-1)} - \underline{x}^*\| \leq \dots \\ &\leq L^k \|\underline{x}^{(1)} - \underline{x}^*\| \leq L^{k+1} \|\underline{x}^{(0)} - \underline{x}^*\| \end{aligned}$$

$$0 < L < 1$$

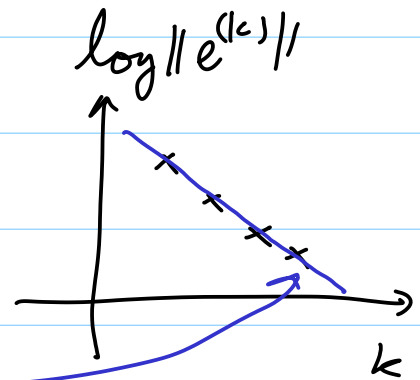
—————> Konvergenz.

Fehler  $\underline{e}_k = \underline{x}^{(k)} - \underline{x}^*$

$$\log \|\underline{e}_k\| \leq k \log L + \log \|\underline{e}^{(0)}\|$$

"linear"

Steigung!

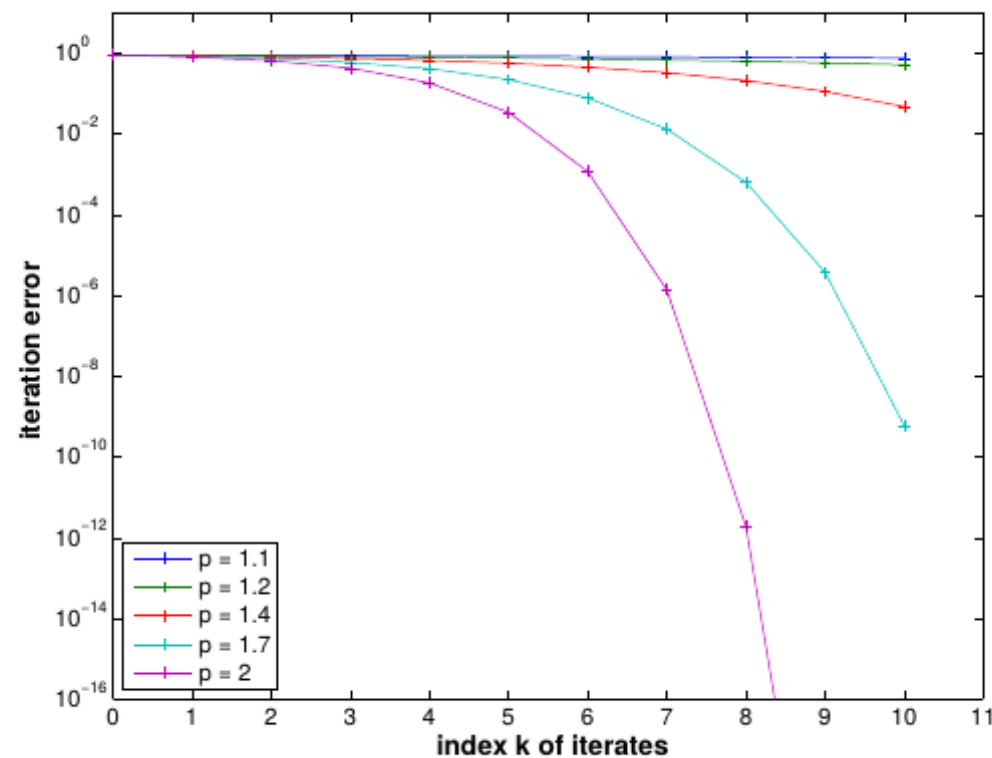


Def Konvergenz der Ordnung  $p$  des iterativen Verfahrens:

es gibt eine  $C > 0$  so dass:

$$\| \underline{x}^{(k+1)} - \underline{x}^* \| \leq C \| \underline{x}^{(k)} - \underline{x}^* \|^p$$

für alle  $k \in \mathbb{N}$ , mit  $C < 1$  für  $p = 1$ .



## §6.2. Fixpunktiteration

$$\underline{F}(\underline{x}^*) = 0 \Leftrightarrow \underline{x}^* = \underline{\Phi}(\underline{x}^*) \quad \text{Fixpunkt von } \underline{\Phi}$$

$$\underline{x}^{(k+1)} = \underline{\Phi}(\underline{x}^k)$$

$$\begin{array}{ccc} \text{falls } \underline{x}^{(k)} \text{ konvergiert} & \downarrow & k \rightarrow \infty \\ \underline{x}^* & & \underline{\Phi}(\underline{x}^*) \end{array}$$

Bsp 1)  $x e^x - 1 = 0 \Leftrightarrow x e^x = 1 \Leftrightarrow x = e^{-x}$

$$\underline{\Phi}_1(x) = e^{-x}$$

starte mit  $x^{(0)}$ ,  $x^{(k+1)} = \underline{\Phi}_1(x^{(k)})$

2)  $x e^x - 1 = 0 \Leftrightarrow x e^x - 1 + x = x \Leftrightarrow x(e^x + 1) = x + 1 \Leftrightarrow$

$$\Leftrightarrow x = \frac{x+1}{e^x + 1}$$

$$\underline{\Phi}_2(x) = \frac{x+1}{e^x + 1}$$

$$3) \quad x e^x - 1 = 0 \Leftrightarrow x e^x - 1 - x = -x \Leftrightarrow x = x + 1 - x e^x$$

$$\Phi(x) = x + 1 - x e^x.$$

$k$	$x^{(k+1)} := \phi_1(x^{(k)})$	$x^{(k+1)} := \phi_2(x^{(k)})$	$x^{(k+1)} := \phi_3(x^{(k)})$
0	0.5000000000000000	0.5000000000000000	0.5000000000000000
1	0.606530659712633	0.566311003197218	0.675639364649936
2	0.545239211892605	0.567143165034862	0.347812678511202
3	0.579703094878068	0.567143290409781	0.855321409174107
4	0.560064627938902	0.567143290409784	-0.156505955383169
5	0.571172148977215	0.567143290409784	0.977326422747719
6	0.564862946980323	0.567143290409784	-0.619764251895580
7	0.568438047570066	0.567143290409784	0.713713087416146
8	0.566409452746921	0.567143290409784	0.256626649129847
9	0.567559634262242	0.567143290409784	0.924920676910549
10	0.566907212935471	0.567143290409784	-0.407422405542253

↓  
lineare konvergenz

↓  
 $p=2$  (quadratische)  
konvergenz

↓  
keine  
konvergenz.

Bez Hinreichende Bedingung für lokale  
lineare konvergenz:

$U$  konvex,  $\Phi: U \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$  stetig differenzierbar

$$L = \sup_{x \in U} \|\underline{D}\Phi(x)\| < 1$$

Wenn  $\Phi(x^*) = x^*$  für  $x^* \in U$ , dann

konvergiert  $x^{(k+1)} = \Phi(x^{(k)})$  gegen  $x^*$

lokal mindestens linear.

→ wir müssen schon nah an  $x^*$  starten!

Bem  $\Phi: U \subset \mathbb{R} \rightarrow \mathbb{R}$   $\Phi$   $(m+1)$ -mal differenzierbar  
 $\Phi(x^*) = x^* \in U$ .

$$\text{Taylor: } \Phi(y) - \Phi(x) = \sum_{k=1}^m \frac{1}{k!} \Phi^{(k)}(x)(y-x)^k +$$

$$+ O(|y-x|^{m+1})$$

Theorem

Voraussetzung:  $\Phi^{(l)}(x^*) = 0$  für  $l=1,2,\dots,m \geq 1$

Dann konvergiert die Fixpunktiteration

$$x^{(k+1)} = \Phi(x^{(k)}) \text{ gegen } x^* \text{ lokal}$$

mit der Ordnung  $p \geq m+1$ .

Beweis Taylor für  $x=x^*$ ,  $y=x^{(k)} \Rightarrow$

$$\begin{aligned} x^{(k+1)} - x^* &= \Phi(x^{(k)}) - \Phi(x^*) = \\ &= \sum_{k=1}^m \frac{1}{k!} \underbrace{\Phi^{(k)}(x^*)}_{!!} (x^{(k)} - x^*)^k + O(|x^{(k)} - x^*|^{m+1}) \end{aligned}$$

$$\Rightarrow |x^{(k+1)} - x^*| \leq C \cdot |x^{(k)} - x^*|^{m+1}$$

Bsp  $\Phi_2(x) = \frac{x+1}{e^x+1} \Rightarrow \Phi_2'(x) = \frac{e^x+1 - (x+1)e^x}{(e^x+1)^2} = \frac{1-xe^x}{(e^x+1)^2}$

$x^*$  Nullstelle von  $1-xe^x$  gesucht  $\Leftrightarrow \underset{x^*}{1-x^*e^{x^*}} = 0 \Rightarrow$

$\Phi_2(x^*) = 0 \Rightarrow \Phi_2$  mindestens Konvergenzordnung?

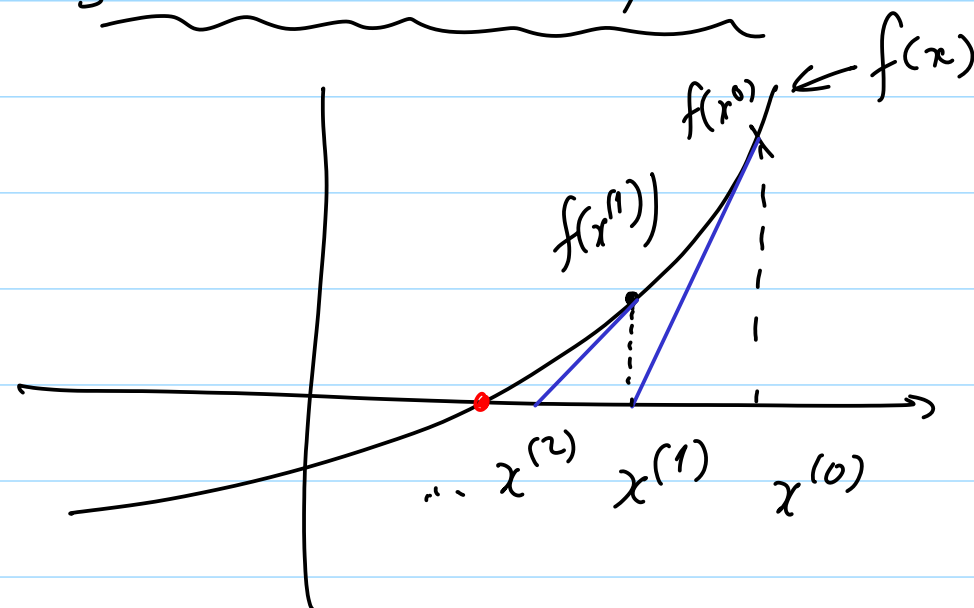
$$\Phi_2(x^*) = \frac{1-x^*e^{x^*}}{(e^{x^*}+1)^2} = \frac{0}{(e^{x^*}+1)^2}$$

(C) V. Gradinaru  
Für  $\Phi_1'(x^*) = -e^{-x^*} \in [-1, 0] \Rightarrow$  nur lineare Konv.

$$\Phi_3'(x^*) = -\frac{1}{x^*} < -1$$

### § 6.3. Newton Verfahren

Idee



approximiere  $f$  lokal durch eine lineare Funktion

$x^{(0)}$  gegeben

ersetze  $\underline{f(x)=0}$  durch  $\tilde{\underline{f(x)=0}}$



Newton:  $\tilde{f}(\underline{x}) = \underline{f}(\underline{x}^0) + \underline{\underline{Df}}(\underline{x}^{(0)}) (\underline{x} - \underline{x}^0)$

Löse  $\tilde{f}(\underline{x}) = 0 \Leftrightarrow \underline{f}(\underline{x}^0) + \underline{\underline{Df}}(\underline{x}^{(0)}) (\underline{x}^{(1)} - \underline{x}^{(0)}) = 0$

(LGS)  $\underline{\underline{Df}}(\underline{x}^{(0)}) \underline{x}^{(1)} = -\underline{f}(\underline{x}^0) + \underline{\underline{Df}}(\underline{x}^{(0)}) \underline{x}^{(0)}$

$\Rightarrow$  Newton Iteration:

gegeben  $\underline{x}^{(0)}$

für  $k = 0, 1, 2, \dots$

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \underline{\underline{Df}}(\underline{x}^{(k)})^{-1} \underline{f}(\underline{x}^{(k)})$$

Beim Berechne niemals die Inverse einer Matrix  
sondern löse nur LGS

Gauss-Elimination; LU-Zerlegung.

Cholesky-Zerlegung; QR-Zerlegung.

Theorem Newton-Verfahren konvergiert mit  
Ordnung  $p=2$  (lokal).

Beweis  
 $d=1$

$$x^{(k+1)} = x^{(k)} - f'(x^{(k)})^{-1} f(x^{(k)})$$

Fixpunkt iteration

$$x^* = \Phi(x^*); \Phi(x) = x - f'(x)^{-1} f(x)$$

Nullstellenproblem:  $f(x^*) = 0$

$$\Phi'(x) = 1 - \frac{f'(x)f'(x) - f(x)f''(x)}{(f'(x))^2} =$$

$$= \underbrace{1-1}_0 + \frac{f(x)f''(x)}{(f'(x))^2} \Rightarrow$$

$$\Phi'(x^*) = \frac{f(x^*)f''(x^*)}{(f'(x^*))^2} = 0 \quad (\text{nur wenn } f'(x^*) \neq 0)$$

Bem Was tun, falls  $f'(x)$  nicht bekannt?

1D: Sekantenverfahren: statt Tangente  
verwende die Sekante

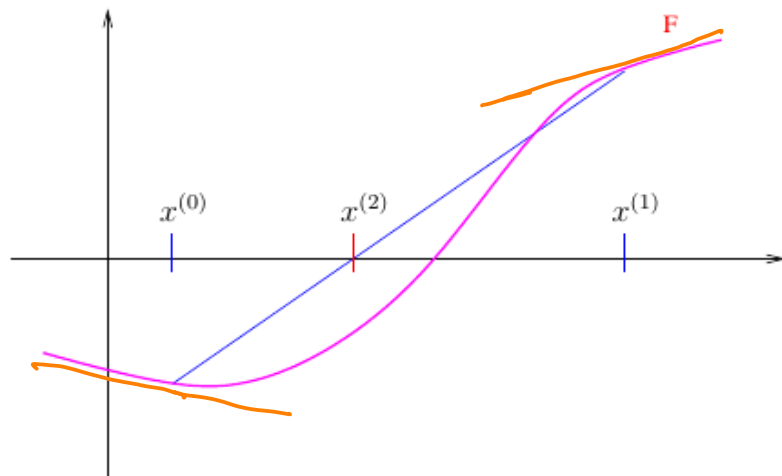
$$f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

Update:  $x^{(k+1)} = x^{(k)} - \lambda \quad (\lambda = x^{(k+1)} - x^{(k)})$

$$\lambda = \frac{x^{(k)} - x^{(k-1)}}{f(x^{(k)}) - f(x^{(k-1)})} f(x^{(k)})$$

Man braucht 2 Startwerte:  $x^{(0)}$  und  $x^{(1)}$

Achtung: so einfach nur in 1D.



Beispiel 3.6.3.  $F(x) = xe^x - 1$ ,  $x^{(0)} = 0$ ,  $x^{(1)} = 5$ .

$k$	$x^{(k)}$	$F(x^{(k)})$	$e^{(k)} := x^{(k)} - x^*$	$\frac{\log  e^{(k+1)}  - \log  e^{(k)} }{\log  e^{(k)}  - \log  e^{(k-1)} }$
2	0.00673794699909	-0.99321649977589	-0.56040534341070	
3	0.01342122983571	-0.98639742654892	-0.55372206057408	24.43308649757745
4	0.98017620833821	1.61209684919288	0.41303291792843	2.70802321457994
5	0.38040476787948	-0.44351476841567	-0.18673852253030	1.48753625853887
6	0.50981028847430	-0.15117846201565	-0.05733300193548	1.51452723840131
7	0.57673091089295	0.02670169957932	0.00958762048317	1.70075240166256
8	0.56668541543431	-0.00126473620459	-0.00045787497547	1.59458505614449
9	0.56713970649585	-0.00000990312376	-0.00000358391394	1.62641838319117
10	0.56714329175406	0.00000000371452	0.00000000134427	
11	0.56714329040978	-0.00000000000001	-0.00000000000000	

Wir beobachten eine fraktionale Konvergenzordnung  $p$ ! Man kann  $p = \frac{1}{2}(1 + \sqrt{5}) \approx 1.62$  beweisen. **Statt Konvergenzordnung 2, wie Newton!**

Beispiel 3.6.4. (Lokale Konvergenz des Sekantenverfahrens)

$$F(x) = \arctan(x)$$

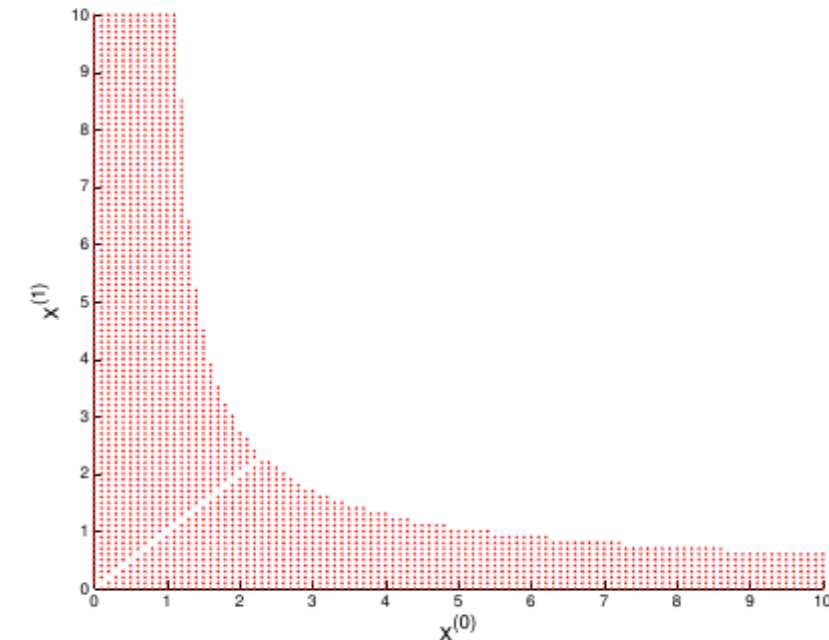


Abb. 3.6.5.  $\hat{=}$  Sekantenverfahren konvergiert für ein Paar  $(x^{(0)}, x^{(1)})$  als Startwert.

Bei Wahl vom Startwert sehr wichtig

Beispiel 3.8.1. (Lokale Konvergenz vom Newton-Verfahren)

$$F(x) = xe^x - 1 \implies F'(-1) = 0$$

$$x^{(0)} < -1 \implies x^{(k)} \rightarrow -\infty$$

$$x^{(0)} > -1 \implies x^{(k)} \rightarrow x^*$$

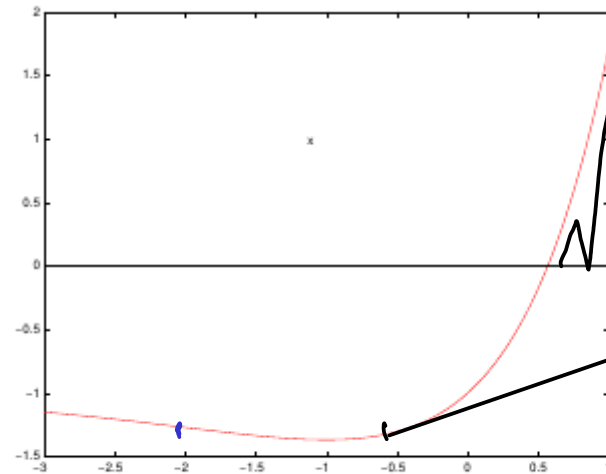
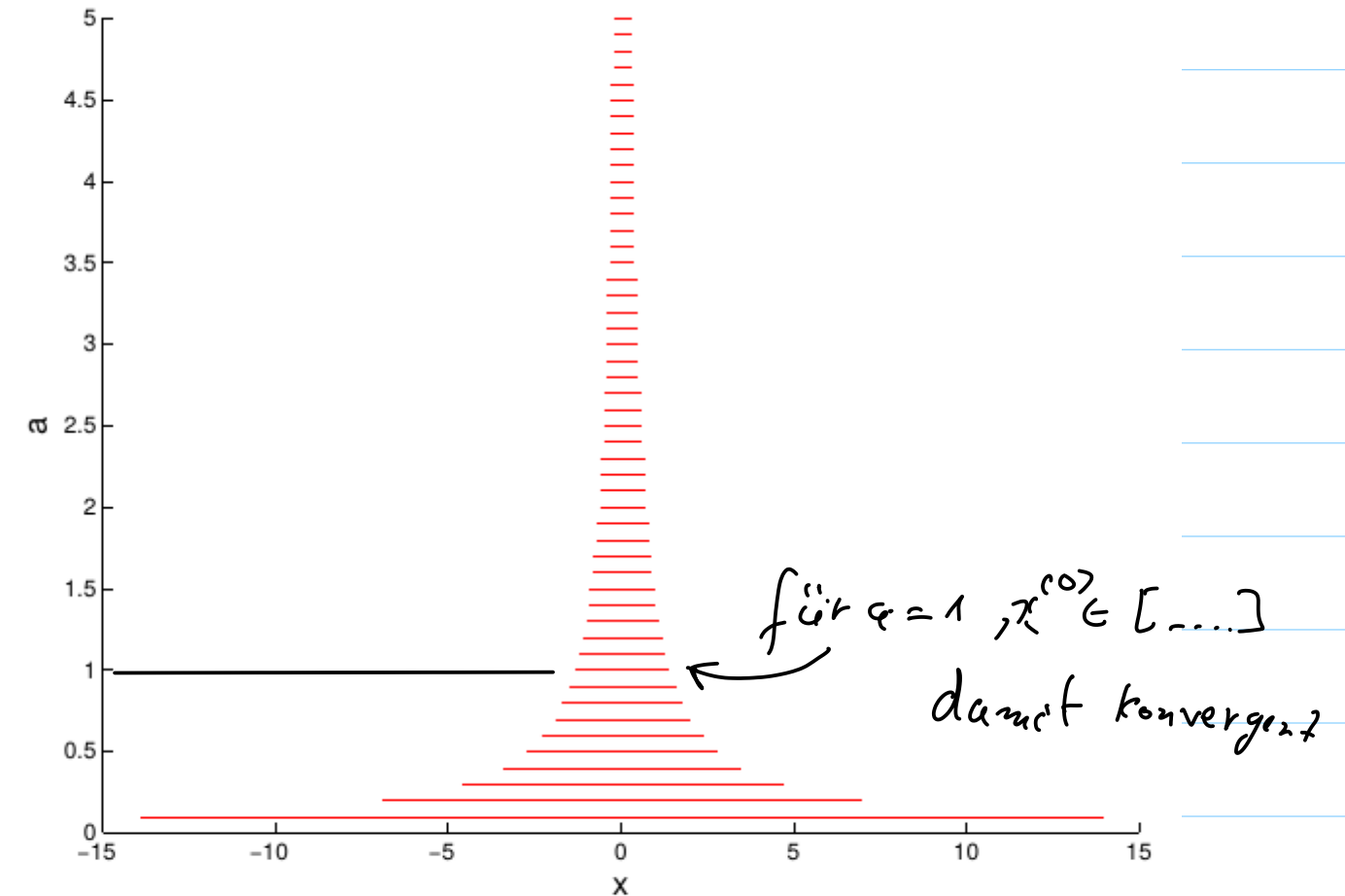


Abb. 3.8.2. Verhalten von  $F(x) = xe^x - 1$

$f(x) = \arctan(ax)$  mit  $a > 0$  Parameter



**Rote Zone** =  $\{x^{(0)} \in \mathbb{R}, x^{(k)} \rightarrow 0\}$

$x^{(0)}$  muss bereits sehr nah an der Lösung sein. ☹️

Bem Bei Divergenz: das Verfahren dämpfen:  
erlaube nur kurze Schritte in der  
Richtung der Tangente

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \lambda \underline{\Delta} \quad \text{mit } |\lambda| < 1$$

↳ Wahl von  $\lambda$  ist heuristisch.

Bem Um Rechenzeit zu sparen (in  $d$  Dimensionen)  
kann man auf das vereinfachte Newton-Verfahren  
ausweichen:

$$x^{(0)}; \quad \underline{J} = \underline{Df}(x^{(0)})$$

$$\underline{J} = \underline{L} \underline{U} \quad \text{kostet } \underline{O(d^3)}$$

verwende für die nächsten Schritte

diesen  $\underline{L}, \underline{U}$  ( $\Rightarrow Df(x^{(0)})$  auch in  
 $x^{(1)}, x^{(2)}, \dots$ )

für  $k = 0, 1, 2, \dots$

$$\text{löse } \underline{L} \underline{U} \underline{\Delta} = \underline{f}(x^{(k)}) \quad \text{nur } O(d^2)$$

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \underline{\Delta}$$

$$\underline{L} \underline{U} \underline{x} = \underline{b} \Leftrightarrow \underline{L} \underline{y} = \underline{b}$$

$\underline{y}$

$$\begin{bmatrix} x & * & \dots & x \\ 0 & x & \dots & x \\ 0 & 0 & x & \dots & x \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & x \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{d-1} \\ y_d \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{d-1} \\ b_d \end{bmatrix}$$

$\frac{d^2}{2}$

Rückwärts substitution:  $y_d = -\dots$   $O(d^2)$   
 $y_{d-1} = \dots$   
 $\dots$

Dasselbe für  $\underline{U} \underline{x} = \underline{y}$  mit  
Vorwärtssubstitution.

Typischerweise nur lineare Konvergenz.

Ben Was ist wenn wir für  $d > 1$  die Ableitung nicht kennen?

1D:  $f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}} = \underline{\underline{\partial_k}}$

dD:  $\underline{\underline{\partial_k}} (\underline{x^{(k)}} - \underline{x^{(k-1)}}) = \underline{f(x^{(k)})} - \underline{f(x^{(k-1)})}$   
 $\hookrightarrow d \times d$  Matrix  $\in \mathbb{R}^d$

Ben Viele mögliche  $\underline{\underline{\partial_k}}$ !

Broyden: baue  $\underline{\underline{\partial_k}}$  iterativ

$$\underline{\underline{\partial_k}} = \underline{\underline{\partial_{k-1}}} + \frac{1}{\|\underline{x^{(k)}} - \underline{x^{(k-1)}}\|_2^2} \underline{f(x^{(k)})} (\underline{x^{(k)}} - \underline{x^{(k-1)}})^T$$

Vorteil: in jedem Schritt nur ein Mal  $\underline{f}$  auszuwerten

Broyden - Verfahren:

$x^{(0)}$  gegeben

$$\underline{\underline{\partial_0}} = \underline{\underline{Df}}(x^{(0)})$$

für  $k=0, 1, 2, \dots$

löse

$$\underline{\underline{\partial_k}} \underline{\Delta} = \underline{f(x^{(k)})}$$

 teuer

$$\underline{x^{(k+1)}} = \underline{x^{(k)}} - \underline{\Delta}$$

$$\underline{\underline{\partial_{k+1}}} \approx \underline{\underline{\partial_k}} + \frac{1}{\|\underline{\Delta}\|_2^2} \underline{f(x^{(k+1)})} (-\underline{\Delta})^T$$

Man kann beweisen:

$$\lim_{k \rightarrow \infty} \frac{\|\underline{x^{(k+1)}} - \underline{x^*}\|}{\|\underline{x^{(k)}} - \underline{x^*}\|} = 0$$

d.h. Superlineare Konvergenz

Bessere Implementierung via  
Shermann-Morrisson-Formel.

$$\underline{\underline{\partial}}_{k+1} = \underline{\underline{\partial}}_k + \underbrace{\left( \frac{1}{\|\underline{\underline{\rho}}\|_2^2} \right)}_{\in \mathbb{R}} \cdot \underbrace{\begin{bmatrix} \underline{\underline{u}} \\ \underline{\underline{v}}^T \end{bmatrix}}_{\text{Rang 1-Matrix}}$$

$$\left( \underline{\underline{A}} + \underline{\underline{u}} \underline{\underline{v}}^T \right)^{-1} = \underline{\underline{A}}^{-1} - \frac{\underline{\underline{A}}^{-1} \underline{\underline{u}} \underline{\underline{v}}^T \underline{\underline{A}}^{-1}}{1 + \underline{\underline{v}}^T \underline{\underline{A}}^{-1} \underline{\underline{u}}}$$

(einfach überprüfen!)

Rang-1-update braucht nur  $\underline{\underline{A}}^{-1}$  (falls  $\underline{\underline{A}}$  invertierbar)

In jedem Schritt:

$$\underline{\underline{\partial}}_{k+1}^{-1} = \underline{\underline{\partial}}_k^{-1} + \frac{\underline{\underline{\partial}}_k^{-1} f(\underline{\underline{x}}^{(k+1)}) \underline{\underline{\rho}}^T \underline{\underline{\partial}}_k^{-1}}{\|\underline{\underline{\rho}}\|^2 - \underline{\underline{\rho}}^T \underline{\underline{\partial}}_k^{-1} f(\underline{\underline{x}}^{(k+1)})}$$

Iteration:  $\underline{\underline{\partial}}_0 = \underline{\underline{D}} f(\underline{\underline{x}}^{(0)})$

zerlege  $\underline{\underline{\partial}}_0 = \underline{\underline{L}} \underline{\underline{U}}$

Löse  $\underline{\underline{L}} \underline{\underline{U}} \underline{\underline{\rho}}^{(0)} = \underline{\underline{f}}(\underline{\underline{x}}^{(0)})$

$$\underline{\underline{x}}^{(1)} = \underline{\underline{x}}^{(0)} - \underline{\underline{\rho}}^{(0)}$$

$$\underline{\underline{\partial}}_1^{-1} = \underline{\underline{\partial}}_0^{-1} + \frac{\underline{\underline{\partial}}_0^{-1} f(\underline{\underline{x}}^{(1)}) \underline{\underline{\rho}}^{(0)T} \underline{\underline{\partial}}_0^{-1}}{\|\underline{\underline{\rho}}\|^2 - \underline{\underline{\rho}}^{(0)T} \underline{\underline{\partial}}_0^{-1} f(\underline{\underline{x}}^{(1)})}$$

$$\underline{\underline{x}}^{(2)} = \underline{\underline{x}}^{(1)} - \underline{\underline{\partial}}_1^{-1} f(\underline{\underline{x}}^{(1)})$$

$\underline{\underline{\rho}}^{(1)}$

$$\underline{\underline{\rho}}^{(1)} = \underline{\underline{\partial}}_1^{-1} f(\underline{\underline{x}}^{(1)})$$

$$= \underbrace{\underline{\underline{\partial_0^{-1} f(x^1)}}}_{\underline{w}} + \underbrace{\underline{\underline{\partial_0^{-1} f(x^{(1)})}}}_{\underline{z}} \underline{\rho^{(0)T}} \underbrace{\underline{\underline{\partial_0^{-1} f(x^{(1)})}}}_{\underline{z}}$$

$\underline{w} = \underline{\partial_0^{-1} f(x^{(1)})}$  berechnet als Lösung von  
LGS  $\underline{\partial_0} \underline{w} = \underline{f(x^1)}$

$$\underline{z} = \underline{\rho^{(0)T}} \underline{w} \in \mathbb{R}$$

$$\underline{\rho^{(1)}} = \underline{w} + \frac{\underline{w} \underline{z}}{\|\underline{\rho^{(0)}}\|^2 - \underline{z}} = \left(1 + \frac{\underline{z}}{\|\underline{\rho^{(0)}}\|^2 - \underline{z}}\right) \underline{w}$$

$$\underline{x^{(2)}} = \underline{x^{(1)}} - \underline{\rho^{(1)}}$$

Code 3.9.4: Broyden-Verfahren

```
from numpy.linalg import lu_solve, lu_factor, norm, solve
from numpy import dot, zeros
```

```
def fastbroyd(x0, F, J, tol=1e-12, maxit=20):
```

```
    x = x0.copy()
```

```
    lup = lu_factor(J)
```

```
    k = 0; s = lu_solve(lup, F(x))
```

```
    x -= s; f = F(x); sn = dot(s, s)
```

```
    dx = zeros((maxit, len(x)))
```

```
    dxn = zeros(maxit)
```

```
    dx[k] = s; dxn[k] = sn
```

```
    k += 1; tol *= tol
```

```
    while sn > tol and k < maxit:
```

```
        w = lu_solve(lup, f)
```

```
        for r in range(1, k):
```

```
            w += dx[r] * (dot(dx[r-1], w) / dxn[r-1])
```

```
        z = dot(s, w)
```

```
        s = (1 + z / (sn - z)) * w
```

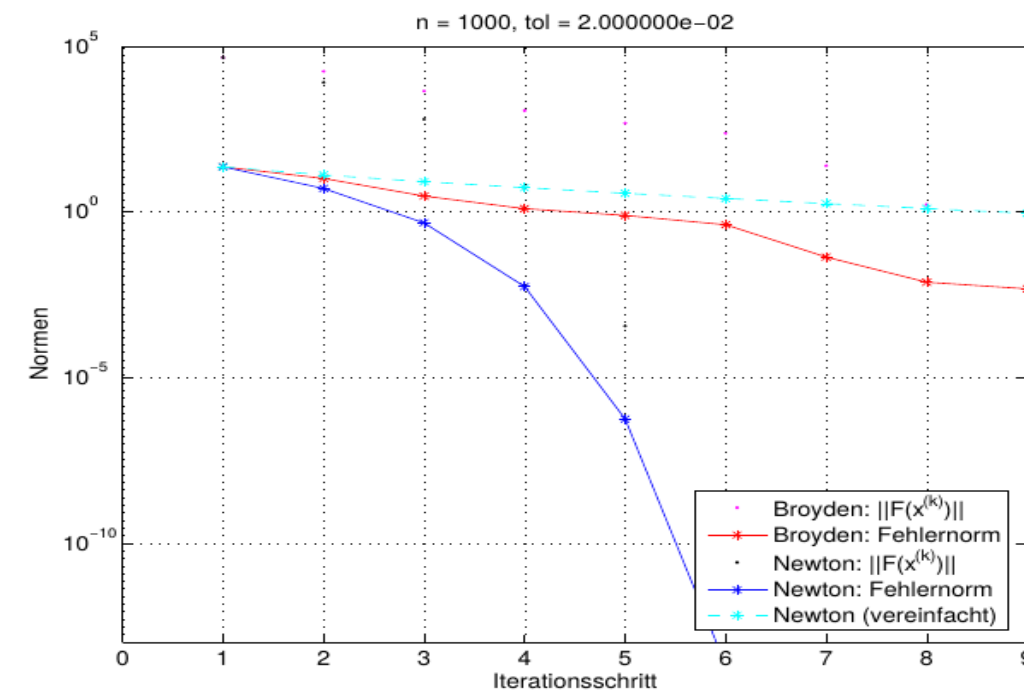
```
        sn = dot(s, s)
```

```
        dx[k] = s; dxn[k] = sn
```

```
        x -= s; f = F(x); k += 1
```

```
    return x, k
```

$f(x^{(0)})$ ;  $\partial = Df(x^{(0)})$   
"LU"  
wieder verwende! bildig!





Software: optimize.fsolve  
 ↳ auf eine andere Methode,

```
In [1]: import numpy as np
        from scipy import optimize
        import matplotlib.pyplot as plt

        optimize.bisect(lambda x: np.exp(x) - 2, -2, 2)
```

```
In [4]: x_root_guess = 2
        f = lambda x: np.exp(x) - 2
        fprime = lambda x: np.exp(x)
        print(optimize.newton(f, x_root_guess))
        print(optimize.newton(f, x_root_guess, fprime=fprime))
```

```
0.69314718056
0.69314718056
```

```
In [6]: def f(x):
        return [x[1] - x[0]**3 - 2 * x[0]**2 + 1, x[1] + x[0]**2 - 1]

        optimize.fsolve(f, [1, 1])
```

```
Out[6]: array([ 0.73205081,  0.46410162])
```

```
In [7]: import sympy
        x, y = sympy.symbols("x, y")
        f_mat = sympy.Matrix([y - x**3 - 2*x**2 + 1, y + x**2 - 1])
        f_mat.jacobian(sympy.Matrix([x, y]))
        def f_jacobian(x):
            return [[-3*x[0]**2-4*x[0], 1], [2*x[0], 1]]
        optimize.fsolve(f, [1, 1], fprime=f_jacobian)
```

```
Out[7]: array([ 0.73205081,  0.46410162])
```

```
In [9]: def f(r):
        return 2 * np.pi * r**2 + 2 / r
        r_min = optimize.brent(f, brack=(0.1, 4))
        r_min, f(r_min)
```

```
Out[9]: (0.54192607725571351, 5.5358104459320856)
```

```
In [10]: def f(X):
         x, y = X
         return (4 * np.sin(np.pi * x) + 6 * np.sin(np.pi * y)) + (x - 1)**2 + (y - 1)**2
```

```
In [11]: x_start = optimize.brute(f, (slice(-3, 5, 0.5), slice(-3, 5, 0.5)), finish=None)
         x_start, f(x_start)
```

```
Out[11]: (array([ 1.5,  1.5]), -9.5)
```

```
In [12]: x_opt = optimize.fmin_bfgs(f, x_start)
         x_opt, f(x_opt)
```

```
Optimization terminated successfully.
Current function value: -9.520229
Iterations: 4
Function evaluations: 28
Gradient evaluations: 7
```

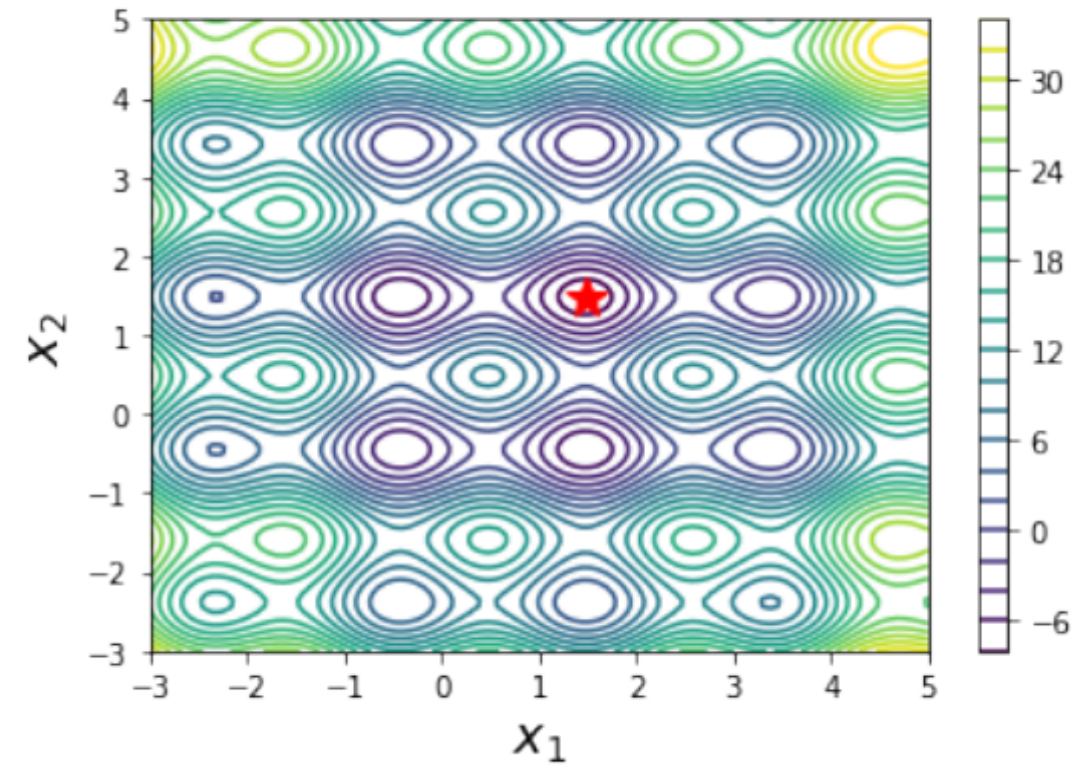
```
Out[12]: (array([ 1.47586906,  1.48365787]), -9.5202292730550155)
```



```

In [13]: def func_X_Y_to_XY(f, X, Y):
        """
        Wrapper for f(X, Y) -> f([X, Y])
        """
        s = np.shape(X)
        return f(np.vstack([X.ravel(), Y.ravel()])).reshape(*s)
fig, ax = plt.subplots(figsize=(6, 4))
x_ = y_ = np.linspace(-3, 5, 100)
X, Y = np.meshgrid(x_, y_)
c = ax.contour(X, Y, func_X_Y_to_XY(f, X, Y), 25)
ax.plot(x_opt[0], x_opt[1], 'r*', markersize=15)
ax.set_xlabel(r"$x_1$", fontsize=18)
ax.set_ylabel(r"$x_2$", fontsize=18)
plt.colorbar(c, ax=ax)
plt.show()

```



```

In [14]: result = optimize.minimize(f, x_start, method='BFGS')
        x_opt = result.x
        x_opt, f(x_opt)

```

```

Out[14]: (array([ 1.47586906,  1.48365787]), -9.5202292730550155)

```

## §7 Intermezzo LA

$$\underline{A} = \begin{bmatrix} | & | & \dots & | \\ \hline & & & \end{bmatrix}$$

$\uparrow$   $\underline{a}_i; \underline{A}_{:,i}$

$$\underline{A} \underline{x} = \sum_{i=1}^n x_i \underline{A}_{:,i}$$

$$\underline{q} = \begin{bmatrix} q_1 \\ \vdots \\ q_n \end{bmatrix}$$

$$\underline{q}^T = [q_1 \dots q_n], \quad \underline{q}^H = [\bar{q}_1 \dots \bar{q}_n]$$

$$\underline{Q} = [\underline{q}_1 \quad \underline{q}_2 \quad \dots \quad \underline{q}_n] \text{ orthogonal/unitär}$$

$$\underline{Q}^H \underline{Q} = \underline{I} \quad \underline{q}_1, \dots, \underline{q}_n \text{ orthonormal}$$

Unitäre Transformationen erhalten Winkel  
und Längen

$$\|\underline{Q} \underline{x}\|_2^2 = (\underline{Q} \underline{x})^H (\underline{Q} \underline{x}) = \underline{x}^H \underline{Q}^H \underline{Q} \underline{x} = \underline{x}^H \underline{x} = \|\underline{x}\|_2^2$$

Orthogonale Transformationen in  $\mathbb{R}^n$ :  
Rotationen, Spiegelungen, Permutation

Ben Gegeben  $\underline{a}_1, \underline{a}_2, \dots, \underline{a}_n$  lin. unabhängige  
Vektoren  $\Rightarrow$  ONB in  $\text{span}\{\underline{a}_1, \underline{a}_2, \dots, \underline{a}_n\}$   
mit Gram-Schmidt  
 $\hookrightarrow$  Wichtig: modifizierte Version.

### §7.1. Wichtigsten Matrixzerlegungen

1) LU- oder LR-Zerlegung: Gauss-Elimination

$$\underline{A} = \underline{L} \cdot \underline{R} = \begin{bmatrix} 1 & & 0 \\ * & \ddots & \\ * & & 1 \end{bmatrix} \cdot \begin{bmatrix} \diagup & * \\ & \diagdown \\ 0 & & \end{bmatrix}$$

Ben Falls  $\underline{A}$  symmetrisch  $\Rightarrow$

$$\underline{A} = \begin{bmatrix} \triangle & \square & \nabla \\ \underline{L} & \underline{D} & \underline{L}^T \end{bmatrix} = \underline{L} \underline{D} \underline{L}^T$$

A symmetrisch, pos. definit :

$$\underline{A} = \underline{L} \underline{\sqrt{D}} \underline{\sqrt{D}} \underline{L}^T = \underline{M} \underline{M}^T \text{ Cholesky-Zerlegung.}$$

Gauss-Elimination  $\underline{P} \underline{A} = \underline{L} \underline{R}$

2)  $\underline{Q} \underline{R}$ -Zerlegung.  $\underline{A} = \underline{Q} \underline{R}$   $\rightarrow$  obere Dreiecksmatrix  
 $\hookrightarrow$  orthogonale Matrix.

$$\underline{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad \underline{A} \underline{x} \in \mathbb{R}^m$$

ass. Norm für Matrizen :  $\|\underline{A}\| = \sup_{\underline{x} \neq 0} \frac{\|\underline{A} \underline{x}\|}{\|\underline{x}\|}$

A quadratisch

Def cond(A) =  $\|\underline{A}^{-1}\| \|\underline{A}\|$

$\hookrightarrow$  Indikator ob das Lösen von LGS

$\underline{A} \underline{x} = \underline{b}$  mit Gauss-Elimination anfällig an

Rundungsfehler ist oder nicht.

$$\underline{A} \underline{x} = \underline{b} \Leftrightarrow \underline{Q} \underline{R} \underline{x} = \underline{b} \Rightarrow \underline{Q}^H \mid \uparrow \text{keine Rundungsfehler-akumulation}$$

$\Rightarrow \underline{R} \underline{x} = \underline{Q}^H \underline{b}$  und Rückwärts substitution,  
 $\rightarrow$  viel stabilere Art LGS zu lösen  
 (kostet auch circa 3 Mal mehr als via LU).

3) Singularwertzerlegung (SVD) A beliebig

$$\underline{A} = \underline{U} \underline{\Sigma} \underline{V}^H$$

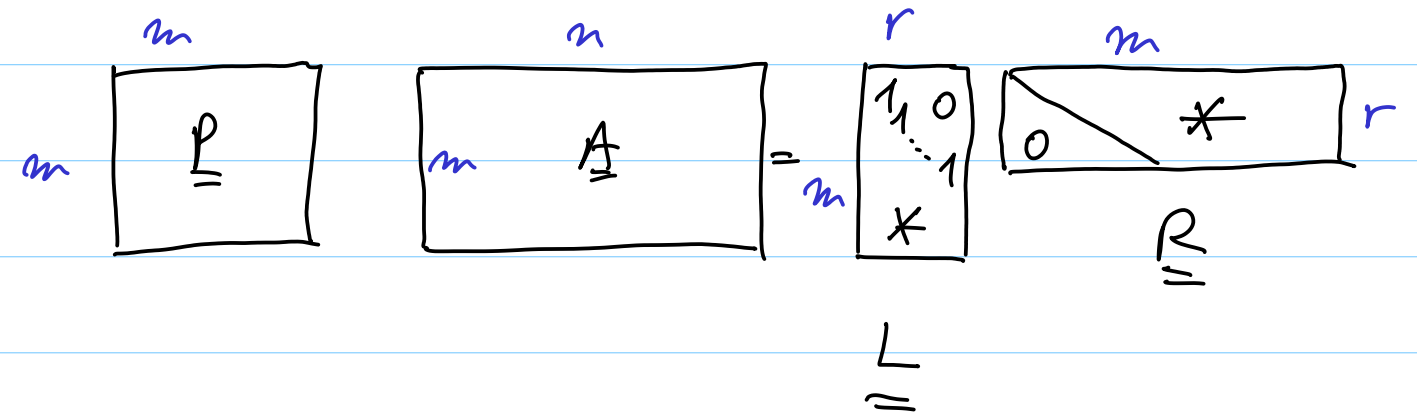
$\downarrow$  orthonormale Spalten  $\downarrow$  orth. Zeilen

$$\underline{\Sigma} = \begin{bmatrix} \sigma_1 & & & 0 \\ & \sigma_2 & & \\ & & \ddots & \\ 0 & & & \sigma_r & & 0 \\ & & & & \ddots & \\ & & & & & 0 \end{bmatrix}$$

$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$   
 mit  $r = \text{Rang}(\underline{A})$

4) Schur-Zerlegung.  $\underline{A}$  beliebig quadratisch  
 $\Rightarrow$  es gibt  $\underline{U}$  unitär, so dass

$$\underline{U}^H \underline{A} \underline{U} = \underline{T} \text{ obere Dreiecksmatrix.}$$



5) Polarzerlegung  $\underline{A} = \underline{Q} \underline{H}$   
 $\uparrow$  orthonormal  $\rightarrow$  symmetrisch, pos. definit

6)  $\underline{A}$   $n \times n$  Matrix,  $n$  lin. unabhängigen EV  $\Rightarrow$   
 $\underline{A} = \underline{S} \underline{\Lambda} \underline{S}^{-1}$   $\rightarrow$  diag(EV)

$$\underline{A} = \underline{S} \underline{\Lambda} \underline{S}^{-1}$$

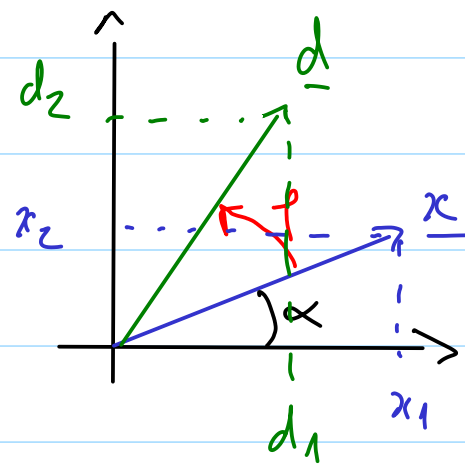
$\hookrightarrow$  EV als Spalten.

$\rightarrow$  orthogonal.

7)  $\underline{A}$   $n \times n$  symmetrisch  $\Rightarrow$   $\underline{A} = \underline{Q} \underline{\Lambda} \underline{Q}^H$

LU-Zerlegung: Zeilenstufenform. für  $m \times n$  Matrix  
 vom Rang  $r < n$ :

### § 7.2. Orthogonale Matrizen



$$\underline{D}(\varphi) = \begin{bmatrix} \cos \varphi & -\sin \varphi \\ \sin \varphi & \cos \varphi \end{bmatrix}$$

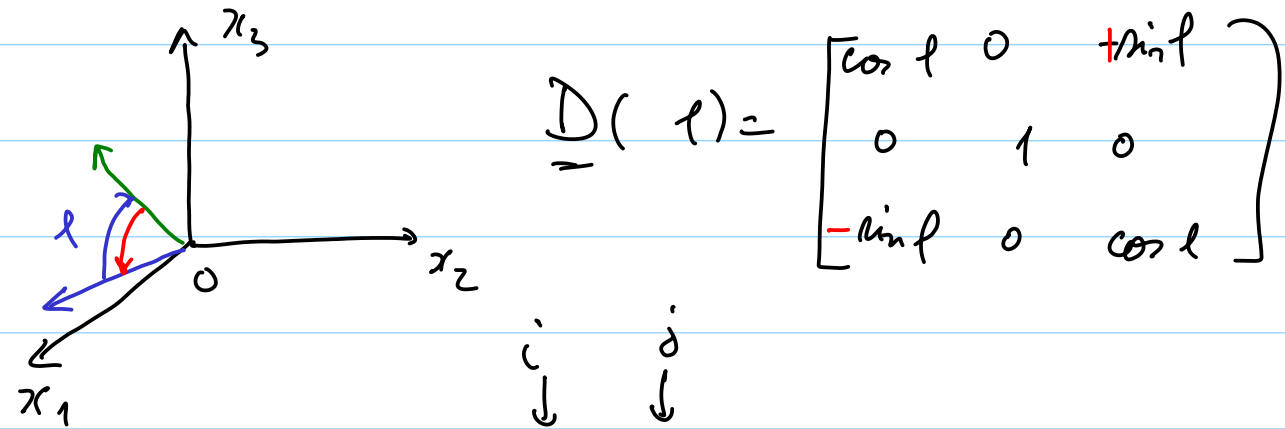
$$\underline{d} = \underline{D}(\varphi) \underline{x}$$

Nützlich  $\underline{D}(-\varphi) = \begin{bmatrix} \cos \varphi & \sin \varphi \\ -\sin \varphi & \cos \varphi \end{bmatrix} = \underline{D}(\varphi)^T$

$$G(\varphi) =$$

Givens-Rotation

Drehung in der  $x_1 O x_3$  - Ebene?



$$\underline{D}_{ij}(-\varphi) = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & \cos \varphi & \sin \varphi & \\ & -\sin \varphi & \cos \varphi & \\ & & & 1 \end{bmatrix} \begin{matrix} \leftarrow i \\ \leftarrow j \end{matrix}$$

$\downarrow$       $\downarrow$   
 $i$       $j$

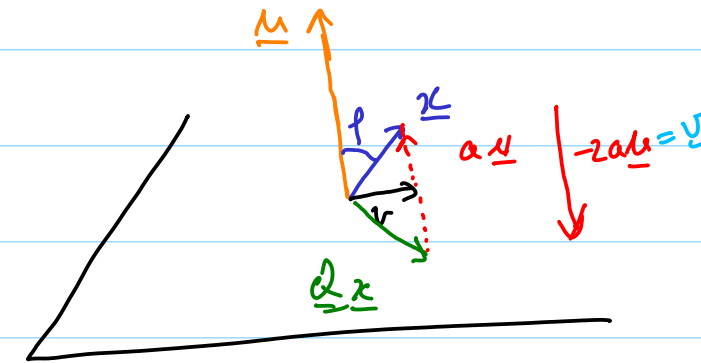
Permutationsmatrix.

$$\begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} \xrightarrow{\underline{P}} \begin{bmatrix} 1 & & & \\ & 1 & & \\ & 0 & & 1 \\ & 1 & 0 & \\ & & \ddots & \\ & & & 1 \end{bmatrix} \begin{matrix} \leftarrow \\ \leftarrow j' \end{matrix}$$

Bsp Spiegelung

$\underline{u} \perp \text{Ebene}$ ,  $\|\underline{u}\|_2 = 1$

$\varphi$  = Winkel zwischen  $\underline{x}$ ,  $\underline{u}$



$$\cos \varphi = \frac{a}{\|\underline{x}\|}$$

$$a = \|\underline{x}\| \cos \varphi$$

$$\underline{Q}\underline{x} = \underline{x} + (-2a\underline{u}) = \underline{x} - 2a\underline{u}$$

$$\underline{u}^T \underline{x} = \langle \underline{u}, \underline{x} \rangle = \|\underline{u}\| \cdot \|\underline{x}\| \cos \varphi = \|\underline{u}\| \cdot a = a$$

$$\underline{Q}\underline{x} = \underline{x} - 2(\underline{u}^T \underline{x}) \underline{u} = \underline{x} - 2\underline{u}(\underline{u}^T \underline{x}) = \underline{x} - 2\underline{u}\underline{u}^T \underline{x}$$

$$= \left( \underline{I} - 2 \underbrace{\underline{u}\underline{u}^T}_{\underline{P}_{\underline{u}}} \right) \underline{x}$$

$$\underline{Q} = \underline{I} - 2 \underline{P}_{\underline{u}}$$

Householdermatrix.

# § 7.3. QR-Zerlegung

(I)

Erste Möglichkeit: via (modifizierten) Gram-Schmidt.

$$\underline{A} = [\underline{v}_1 \ \underline{v}_2 \ \dots \ \underline{v}_n] \rightsquigarrow \underline{q}_1, \dots, \underline{q}_n \text{ ONB}$$

$$\underline{q}_1 := \frac{1}{\|\underline{v}_1\|} \underline{v}_1 \Rightarrow \underline{v}_1 = \|\underline{v}_1\| \underline{q}_1 = \langle \underline{q}_1, \underline{v}_1 \rangle \underline{q}_1$$

$$\underline{c}_2 = \underline{v}_2 - \rho_{\underline{q}_1} \underline{v}_2 = \underline{v}_2 - \langle \underline{q}_1, \underline{v}_2 \rangle \underline{q}_1$$

$$\underline{q}_2 := \frac{1}{\|\underline{c}_2\|} \underline{c}_2$$

$$\underline{c}_3 := \underline{v}_3 - \rho_{\substack{\underline{v}_2 \\ \text{span}\{\underline{q}_1, \underline{q}_2\}}} \underline{v}_3 = \underline{v}_3 - \langle \underline{q}_1, \underline{v}_3 \rangle \underline{q}_1 - \langle \underline{q}_2, \underline{v}_3 \rangle \underline{q}_2$$

$$\underline{q}_3 := \frac{1}{\|\underline{c}_3\|} \underline{c}_3$$

...

$$\underline{c}_k = \underline{v}_k - \rho_{\substack{\underline{v}_k \\ \text{span}\{\underline{q}_1, \dots, \underline{q}_{k-1}\}}} \underline{v}_k = \underline{v}_k - \langle \underline{q}_1, \underline{v}_k \rangle \underline{q}_1 - \dots - \langle \underline{q}_{k-1}, \underline{v}_k \rangle \underline{q}_{k-1}$$

$$\underline{q}_k = \frac{1}{\|\underline{c}_k\|} \underline{c}_k \dots$$

$$\underline{v}_1 = \langle \underline{q}_1, \underline{v}_1 \rangle \underline{q}_1 + 0 \underline{q}_2 + \dots + 0 \underline{q}_n$$

$$\underline{v}_2 = \langle \underline{q}_1, \underline{v}_2 \rangle \underline{q}_1 + \langle \underline{q}_2, \underline{v}_2 \rangle \underline{q}_2 + 0 \underline{q}_3 + \dots + 0 \underline{q}_n$$

$$\underline{v}_3 = \langle \underline{q}_1, \underline{v}_3 \rangle \underline{q}_1 + \langle \underline{q}_2, \underline{v}_3 \rangle \underline{q}_2 + \langle \underline{q}_3, \underline{v}_3 \rangle \underline{q}_3 + 0 \underline{q}_4 + \dots + 0 \underline{q}_n$$

...

$$\underline{v}_k = \langle \underline{q}_1, \underline{v}_k \rangle \underline{q}_1 + \dots + \langle \underline{q}_k, \underline{v}_k \rangle \underline{q}_k + 0 \underline{q}_{k+1} + \dots + 0 \underline{q}_n$$

Im  $k$ -ten Schritt:

$$\underline{v}_k = (q_1^H \underline{v}_k) \underline{q}_1 + (q_2^H \underline{v}_k) \underline{q}_2 + \dots + (q_k^H \underline{v}_k) \underline{q}_k$$

$\underline{v}_k$  = lineare Kombination von  $\underline{q}_1, \underline{q}_2, \dots, \underline{q}_k$   
mit koeff  $(q_1^H \underline{v}_k), \dots, (q_k^H \underline{v}_k) \in \mathbb{C}$

$$\underline{v}_k = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_k \end{bmatrix} \begin{bmatrix} (q_1^H \underline{v}_k) \\ \vdots \\ (q_k^H \underline{v}_k) \end{bmatrix}$$

das für jedes  $k=1, 2, \dots, n$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \dots & \underline{v}_k \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_k \end{bmatrix} \begin{bmatrix} q_1^H \underline{v}_1 & q_1^H \underline{v}_2 & \dots & q_1^H \underline{v}_k \\ \vdots & q_2^H \underline{v}_2 & \dots & q_2^H \underline{v}_k \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \vdots & \vdots & q_k^H \underline{v}_k \end{bmatrix}$$

}  $k$

←

$$\begin{bmatrix} \underline{v}_1 & \dots & \underline{v}_k \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_k \end{bmatrix} \begin{bmatrix} * & & \\ & * & \\ & & \ddots \\ 0 & & & * \end{bmatrix}$$

Spalten sind orthogonale Vektoren.

Nach  $n$  Schritten:

$$\begin{bmatrix} \underline{v}_1 & \dots & \underline{v}_n \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_n \end{bmatrix} \begin{bmatrix} * & & \\ & * & \\ & & \ddots \\ 0 & & & * \end{bmatrix}$$

$\begin{matrix} i \\ \downarrow \\ q_i^H \underline{v}_i \\ \leftarrow j \end{matrix}$



Ben Algorithmus entspricht zu Multiplikationen mit oberen Dreiecksmatrizen:

$$\underline{A} = \underbrace{\underline{R}_1 \underline{R}_2 \dots \underline{R}_n}_{\underline{R}^{-1}} = \underline{Q} \Rightarrow \underline{A} = \underline{Q} \underline{R}$$

$$r_{ij} = q_i^H v_j$$

$$\underline{R}_1 = \begin{bmatrix} \frac{1}{r_{11}} & -\frac{r_{12}}{r_{11}} & -\frac{r_{13}}{r_{11}} & \dots \\ & 1 & & \\ 0 & & 1 & \ddots \\ & & & 1 \end{bmatrix}$$

$$\underline{R}_2 = \begin{bmatrix} 1 & & & \\ & \frac{1}{r_{22}} & -\frac{r_{23}}{r_{22}} & \dots \\ & & 1 & \ddots \\ & & & 1 \end{bmatrix} \quad \underline{R}_3 = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \frac{1}{r_{33}} & \dots \\ & & & 1 \ddots \end{bmatrix}$$

Ben Gram-Schmidt produziert eine orthogonale Matrix Q indem man in jedem Schritt mit einer oberen Dreiecksmatrix multipliziert.

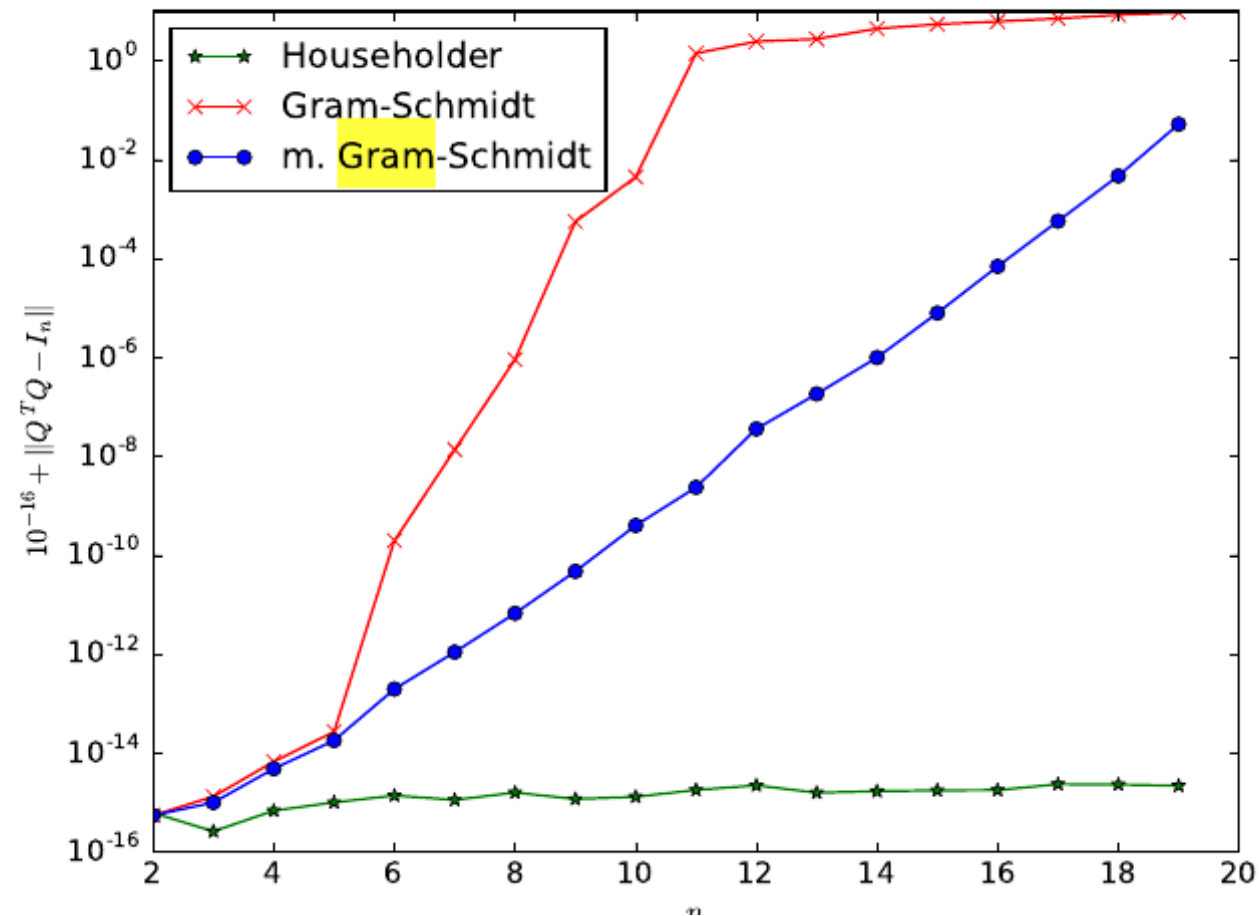
$\Rightarrow$  anfällig zu Rundungsfehler!

Ben Ein Vorteil von Gram-Schmidt ist, wenn man noch k Schritte aufhört, hat man bereits  $\underline{q}_1, \dots, \underline{q}_k$  orthogonal

$$\text{span} \{ \underline{q}_1, \dots, \underline{q}_k \} = \text{span} \{ \underline{v}_1, \dots, \underline{v}_k \}$$



$$\mathbf{A} = \begin{bmatrix} t_0^{n-1} & \dots & t_0^1 & 1 \\ t_1^{n-1} & \dots & t_1^1 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ t_{19}^{n-1} & \dots & t_{19}^1 & 1 \end{bmatrix}, \text{ mit } t_i = \frac{i}{(m-1)}.$$



classisches GS.:

$$P_{q_k} v_k = \begin{pmatrix} P_{q_{k-1}} & \dots & P_{q_2} & P_{q_1} & v_k \end{pmatrix} \rightarrow \text{alle gleichzeitig}$$

↓ nur auf  $q_k$  projiziert.

classisch.  $\underline{a}_1, \dots, \underline{a}_n$

für  $j=1, \dots, n$ :

$$\underline{v}_j = \underline{a}_j$$

für  $i=1, 2, \dots, j-1$ :

$$r_{ij} = \underline{q}_i^H \underline{a}_j$$

$$\underline{v}_j = \underline{v}_j - r_{ij} \underline{q}_i$$

$$r_{jj} = \|\underline{v}_j\|$$

$$\underline{q}_j = \frac{\underline{v}_j}{r_{jj}}$$

Sofort  $\underline{q}_1, \dots, \underline{q}_{k-1}$  verwenden.  
Projektion aller folgenden  
Spalten auf das neu gefundene  
 $\underline{q}_i$

für  $j=1, \dots, n$ :

$$\underline{v}_j = \underline{a}_j$$

für  $i=1, \dots, n$

$$r_{ji} = \|\underline{v}_i\|, \underline{q}_i = \frac{1}{r_{ii}} \underline{v}_i$$

für  $j=i+1, \dots, n$ :

$$r_{ij} = \underline{q}_i^H \underline{v}_j$$

$$\underline{v}_j = \underline{v}_j - r_{ij} \underline{q}_i$$

Fazit:

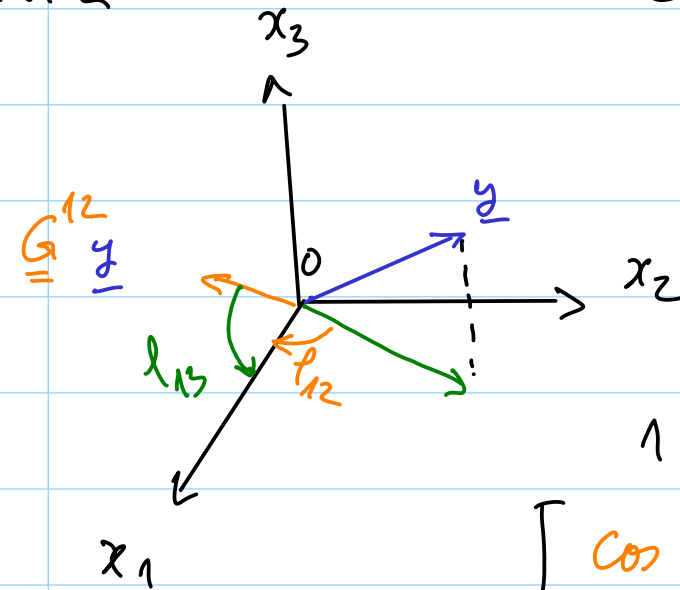
- 1) klassisches Gram-Schmidt <sup>ist</sup> nur für analytische Zwecke zu verwenden!
- 2) modifiziertes Gram-Schmidt <sup>ist</sup> nur für "kleine" Matrizen oder nur für einige Vektoren zu verwenden zu verwenden!
- 3) sonst II, III die folgen!

II Q R via Rotationen

$$\begin{bmatrix} \underline{A} \end{bmatrix} = \begin{bmatrix} \underline{Q} \end{bmatrix} \begin{bmatrix} \text{wavy}^* \\ 0 \end{bmatrix} = \underline{Q} \underline{R}$$

$$\begin{bmatrix} \underline{A} \end{bmatrix} = \begin{bmatrix} \underline{Q} \end{bmatrix} \begin{bmatrix} \text{wavy}^* \\ 0 \end{bmatrix} = \underline{Q} \underline{R}$$

Idee: Erzeuge die 0 unterhalb der Hauptdiagonale mittels orthogonale Matrizen  
→ Drehungen!



$$\underline{A} = \begin{bmatrix} y_1 & \tilde{z}_1 \\ y_2 & \tilde{z}_2 \\ y_3 & \tilde{z}_3 \end{bmatrix}$$

$$\underline{G}^{12}(t_{12}) = \begin{bmatrix} \cos t_{12} & \sin t_{12} & 0 \\ -\sin t_{12} & \cos t_{12} & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Givens-Drehung, die  $y_2$  zu 0 macht

$$\underline{G}^{12}(t_{12}) \underline{A} = \begin{bmatrix} * & * & \vdots \\ 0 & * & \vdots \\ y_3 & \tilde{z}_3 & \ddots \end{bmatrix}$$

$$\underline{G}^{13}(\varphi_{13}) = \begin{bmatrix} \cos \varphi_{13} & 0 & \sin \varphi_{13} \\ 0 & 1 & 0 \\ -\sin \varphi_{13} & 0 & \cos \varphi_{13} \end{bmatrix}$$

↪ Givens-Drehung, die  $z_3$  zu 0 macht

$$\underbrace{\underline{G}^{13}}_{\text{orthogonale Matrizen}} \underbrace{\underline{G}^{12}}_{\text{orthogonale Matrizen}} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} * \\ 0 \\ 0 \end{bmatrix}$$

orthogonale Matrizen

$$\underline{y} \in \mathbb{R}^n \quad \underbrace{\underline{G}^{1n} \dots \underline{G}^{13} \underline{G}^{12}}_{\underline{Q}^1} \underline{y} = \begin{bmatrix} * \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Dann kommt die zweite Spalte von  $\underline{A}$  dran:

$$\begin{bmatrix} * & z_1 & * \\ 0 & z_2 & * \\ 0 & z_3 & * \end{bmatrix}$$

$\underline{G}^{23}$  um  $z_3$  zu 0 zu machen

$$\underline{Q}^1 \underline{A} = \begin{bmatrix} * & x & x & \dots & x \\ 0 & x & x & \dots & x \\ 0 & x & x & \dots & x \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & x & x & \dots & x \end{bmatrix}$$

$$\underline{Q}^2 = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & \boxed{\text{Drehungen um}} \\ \vdots & \boxed{\begin{smallmatrix} x \\ 0 \\ \vdots \\ 0 \end{smallmatrix}} & \dots & \dots & \dots \\ 0 & \dots & \dots & \dots & \dots \end{bmatrix}$$

zu erzeugen:

$$\underbrace{\underline{Q}^{n-1} \dots \underline{Q}^2 \underline{Q}^1}_{\text{orthogonale Matrizen}} \underline{A} = \underline{R} \quad \Rightarrow$$

orthogonale Matrizen

$$(\underline{Q}^1)^T \dots (\underline{Q}^{n-1})^T$$

$$\underline{A} = \underbrace{(\underline{Q}^1)^T \dots (\underline{Q}^{n-1})^T}_{\text{orthogonale Matrix}} \underline{R} = \underline{Q} \underline{R}$$

$$G^{ij}(l) = \begin{bmatrix} 1 & & & & & \\ & 1 & & & & \\ & & \ddots & & & \\ & & & \cos l & \sin l & \\ & & & & 1 & \\ & & & & & \ddots \\ & & & & & & 1 \\ & & & & & & & \cos l & \sin l \\ & & & & & & & & 1 \\ & & & & & & & & & \ddots \\ & & & & & & & & & & 1 \end{bmatrix} \begin{matrix} \leftarrow i \\ \\ \\ \\ \\ \\ \\ \leftarrow j \\ \\ \end{matrix}$$

$$G^{ij}(l) \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ x_i \\ x_{i+1} \\ \vdots \\ x_{j-1} \\ x_j \\ x_{j+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ x_i \\ x_{i+1} \\ \vdots \\ x_{j-1} \\ x_j \\ x_{j+1} \\ \vdots \\ x_n \end{bmatrix} \begin{matrix} \uparrow \\ \\ \\ \text{red } r \\ \\ \\ \text{red } 0 \\ \\ \end{matrix}$$

$$r = \sqrt{x_i^2 + x_j^2}$$

$$\cos l = \frac{x_i}{r}$$

$$\sin l = \frac{x_j}{r}$$

Bsp  $\begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} \rightsquigarrow \begin{bmatrix} * \\ 0 \\ 0 \end{bmatrix}$   $r = \sqrt{4^2 + (-3)^2} = 5$   
 $\cos l = \frac{4}{5}, \sin l = -\frac{3}{5}$

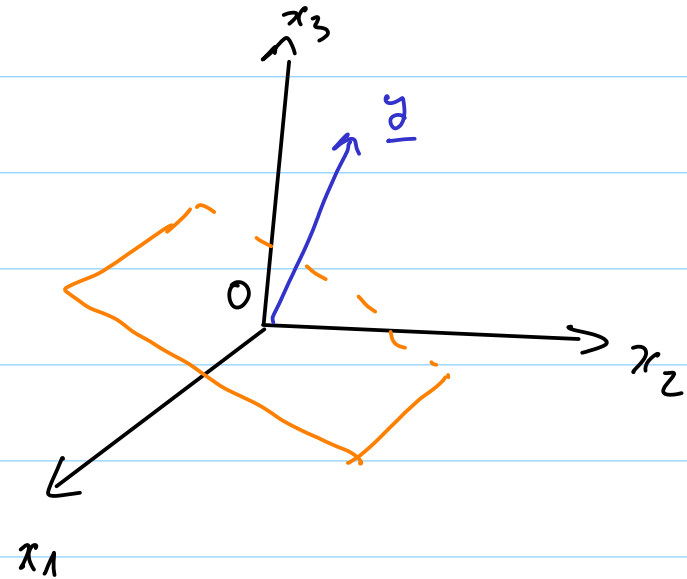
$$G^{12} \begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{4}{5} & -\frac{3}{5} & 0 \\ \frac{3}{5} & \frac{4}{5} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix}$$

$$G^{13} \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{5}{\sqrt{26}} & 0 & \frac{1}{\sqrt{26}} \\ 0 & 1 & 0 \\ -\frac{1}{\sqrt{26}} & 0 & \frac{5}{\sqrt{26}} \end{bmatrix} \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \sqrt{26} \\ 0 \\ 0 \end{bmatrix}$$

$$G = G^{13} G^{12}$$

Ben Diese Methode (Givens-Drehungen) ist günstig falls A sehr viele 0-Einträge hat. (dünnbesetzte Matrix).

### III QR-Zerlegung mit Spiegelungen (standard-Software).



$$\underline{\underline{Q}} = \underline{\underline{I}} - 2 \underline{\underline{u}} \underline{\underline{u}}^T$$

$$\underline{\underline{u}} \perp \text{Ebene}$$

$$\|\underline{\underline{u}}\|_2 = 1$$

$\underline{\underline{Q}}^1$  spiegelt  $\underline{\underline{y}}$  so dass  $\underline{\underline{Q}}^1 \underline{\underline{y}}$  auf  $\partial x_1$  liegt

$$\underline{\underline{Q}}^1 \underline{\underline{y}} = \begin{bmatrix} * \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \underline{\underline{Q}}^1 \underline{\underline{A}} = \begin{bmatrix} * & x & \dots & x \\ 0 & x & \dots & \\ \vdots & \vdots & & \\ 0 & x & \dots & x \end{bmatrix}$$

$\uparrow$   
 $\underline{\underline{z}}$

$\nwarrow$  Spiegelung in  $x_2 \vee x_3$

$\underline{\underline{Q}}^2$  spiegelt  $\underline{\underline{z}}$  so dass  $\underline{\underline{Q}}^2 \underline{\underline{z}}$  auf  $\partial x_2$  liegt

$$\underline{\underline{Q}}^2 \underline{\underline{Q}}^1 \underline{\underline{A}} = \begin{bmatrix} * & x & x & \dots & x \\ 0 & x & x & \dots & x \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & x & \dots & x \end{bmatrix} \quad \text{usw.}$$

Bsp  $\underline{\underline{x}} = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} \xrightarrow{\underline{\underline{Q}}} 3 \underline{\underline{e}}_1 = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} \quad \underline{\underline{e}}_1 = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

$$\|\underline{\underline{x}}\| = \sqrt{2^2 + 2^2 + 1^2} = 3$$

$$\underline{\underline{x}} = \|\underline{\underline{x}}\| \underline{\underline{e}}_1 + \underline{\underline{v}} \Rightarrow \underline{\underline{v}} = \underline{\underline{x}} - \|\underline{\underline{x}}\| \underline{\underline{e}}_1 = \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix}$$

$$\underline{\underline{u}} = \frac{1}{\|\underline{\underline{v}}\|} \underline{\underline{v}} = \frac{1}{\sqrt{6}} \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix} \Rightarrow$$

$$\Rightarrow \underline{\underline{Q}} = \underline{\underline{I}} - 2 \underline{\underline{u}} \underline{\underline{u}}^T = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} - \frac{2}{6} \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 2 & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} - \frac{2}{6} \begin{bmatrix} 1 & -2 & -1 \\ -2 & 4 & 2 \\ -1 & 2 & 1 \end{bmatrix} =$$

$$= \frac{1}{3} \begin{bmatrix} 2 & 2 & 1 \\ 2 & -1 & -2 \\ 1 & -2 & 2 \end{bmatrix} \Rightarrow \underline{\underline{Q}} \underline{\underline{x}} = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix}$$

$$\underline{Q}^{n-1} \underline{Q}^{n-2} \dots \underline{Q}^2 \underline{Q}^1 \underline{A} = \underline{R} \Rightarrow$$

$$\underline{A} = \underline{Q} \underline{R}, \quad \underline{Q} = \underline{Q}^1 \underline{Q}^2 \dots (\underline{Q}^{n-1})^T$$

Bsp

$$\underline{A} = \begin{bmatrix} 1 & 1 \\ 2 & 0 \\ 2 & 0 \end{bmatrix} \Rightarrow \underline{Q}^1 \underline{A} = \left[ \underline{Q}^1 \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \quad \underline{Q}^1 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right]$$

$$\underline{Q}^1 \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} \text{ do } \left\| \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \right\| = 3$$

$$\underline{v} = \underline{x} - 3 \underline{e}_1 = \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} - \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} = \begin{bmatrix} -2 \\ 2 \\ 2 \end{bmatrix}$$

$$\|\underline{v}\| = \sqrt{(-2)^2 + 2^2 + 2^2} = \sqrt{12}$$

$$\underline{u} = \frac{1}{\sqrt{12}} \begin{bmatrix} -2 \\ 2 \\ 2 \end{bmatrix} \Rightarrow \underline{Q}^1 = \underline{I} - 2 \underline{u} \underline{u}^T = \underline{I} - \frac{2}{12} \begin{bmatrix} -2 \\ 2 \\ 2 \end{bmatrix} \begin{bmatrix} -2 & 2 & 2 \end{bmatrix} =$$

$$= \frac{1}{6} \begin{bmatrix} 2 & 4 & 4 \\ 4 & 2 & -4 \\ 4 & -4 & 2 \end{bmatrix}$$

$$\underline{Q}^1 \underline{A} = \begin{bmatrix} 3 & \frac{1}{3} \\ 0 & \frac{2}{3} \\ 0 & \frac{2}{3} \end{bmatrix}$$

MSW.

Theorem  $\underline{A} \in \mathbb{R}^{m \times n}$ ,  $m \geq n$  voller Rang  
( $\text{rang}(\underline{A}) = n$ ).

Dann gibt es eine eindeutige Matrix  $\underline{Q} \in \mathbb{R}^{m \times m}$   
orthogonal, so dass

$$\underline{A} = \underline{Q} \begin{bmatrix} \underline{R} \\ \underline{0} \end{bmatrix}_m \quad \text{mit } \underline{R} \text{ obere Dreiecksmatrix}$$

und die Elemente auf der Diagonale  
von  $\underline{R}$  sind  $\geq 0$ .

Bem. Falls die Elemente auf der Diagonale von  
 $\underline{R} > 0$  sind, dann ist  $\underline{R}$  der Cholesky-Faktor  
von  $\underline{A}^T \underline{A}$

Beweis

$$\underline{A}^T \underline{A} = \begin{bmatrix} \underline{R}^T & \underline{0} \end{bmatrix} \underline{Q}^T \underline{Q} \begin{bmatrix} \underline{R} \\ \underline{0} \end{bmatrix} = \underline{R}^T \underline{R}$$

Theorem Falls  $\text{rang}(\underline{A}) = r < n$ , dann  
gibt es eine Permutation  $\underline{P}$

$$\underline{A} \underline{P} = \underline{Q} \begin{bmatrix} \underline{R}_{11} & \underline{R}_{12} \\ \underline{0} & \underline{0} \end{bmatrix} \begin{matrix} r \\ m-r \end{matrix}$$

mit  $\underline{R}_{11}$  obere Dreiecksmatrix mit  
Elementen auf der Diagonale  $\geq 0$ .

$\underline{R}_{11}$  ist eindeutig,  $\underline{R}_{12}$  nicht!

Beweis wähle  $r$  lin. unabhängige Spalten,  
permutiere sie nach vorne:

$$\underline{A} \underline{P} = \left[ \underline{A}_1 \mid \underline{A}_2 \right] \text{ und wende voriges}$$

Theorem (Givens-Drehung) an.

Ben QR-Zerlegung mit orthogonalen Transformationen darf nicht vorzeitig abgebrochen werden.

§7.4. Singularwertzerlegung: siehe LA, Skript

Bsp  $5 \geq 3 \geq 0.1 \geq 10^{-5} \geq 10^{-8} \geq 10^{-12} \geq 5 \cdot 10^{-16} \geq 2 \cdot 10^{-16} \geq 10^{-42} \geq 0 = 0$

$\sigma_1 \quad \sigma_2 \quad \sigma_3 \quad \sigma_4 \quad \sigma_5 \quad \sigma_6 \quad \sigma_7 \quad \sigma_8 \quad \sigma_9$

mathematische Rang:  $r=9$

numerische Rang: dort wo  $\sigma_p \gg \sigma_{p+1}$   $p=6$ , oder 7  
 oder auch  $p=4$   $p=5?$  } Anwendungsabhängig!

Ben Reduzierte SVD:  $\underline{A} = \begin{bmatrix} \underline{u}_1 & \dots & \underline{u}_r \end{bmatrix} \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_r \end{bmatrix} \begin{bmatrix} \underline{v}_1^T \\ \vdots \\ \underline{v}_r^T \end{bmatrix}$



$$\underline{u}_1, \underline{u}_2, \dots, \underline{u}_r, \underline{u}_{r+1}, \dots, \underline{u}_m \in \mathbb{R}^m \text{ ONB}$$

$$\underline{v}_1, \underline{v}_2, \dots, \underline{v}_r, \underline{v}_{r+1}, \dots, \underline{v}_n \in \mathbb{R}^n \text{ ONB}$$

$$\underline{u}_i^T \underline{A} \underline{v}_j = \begin{cases} \sigma_i, & \text{falls } i=j \leq r \\ 0, & \text{sonst} \end{cases} \quad (\Leftrightarrow)$$

$$\underline{U}^T \underline{A} \underline{V} = \underline{\Sigma} = \begin{array}{|c|c|} \hline \begin{array}{c} \sigma_1 \\ \vdots \\ \sigma_r \end{array} & \begin{array}{c} 0 \\ \vdots \\ 0 \end{array} \\ \hline \begin{array}{c} 0 \\ \vdots \\ 0 \end{array} & \begin{array}{c} 0 \\ \vdots \\ 0 \end{array} \\ \hline \end{array} \quad \text{ge.d.}$$

ONB Bild  $\underline{A}$

$$\underline{A} = \begin{array}{|c|c|} \hline \begin{array}{c} \text{ONB in } \text{Bild } \underline{A} \\ \text{ONB in } \text{Ker } \underline{A} \end{array} & \begin{array}{c} \text{ONB in } \text{Ker } \underline{A} \end{array} \\ \hline \end{array} \quad \underline{U}^T \quad \underline{\Sigma} \quad \underline{V}^T$$

$$\underline{U}_r = [\underline{u}_1 \dots \underline{u}_r], \quad \underline{V}_r = [\underline{v}_1 \dots \underline{v}_r]$$

$$\underline{v}_1, \dots, \underline{v}_r = \text{rechte singularvektoren} = \text{EV von } \underline{A}^T \underline{A}$$

$$\underline{u}_1, \dots, \underline{u}_r = \text{linke} \quad \text{---}, \quad \text{---} = \text{EV von } \underline{A} \underline{A}^T$$

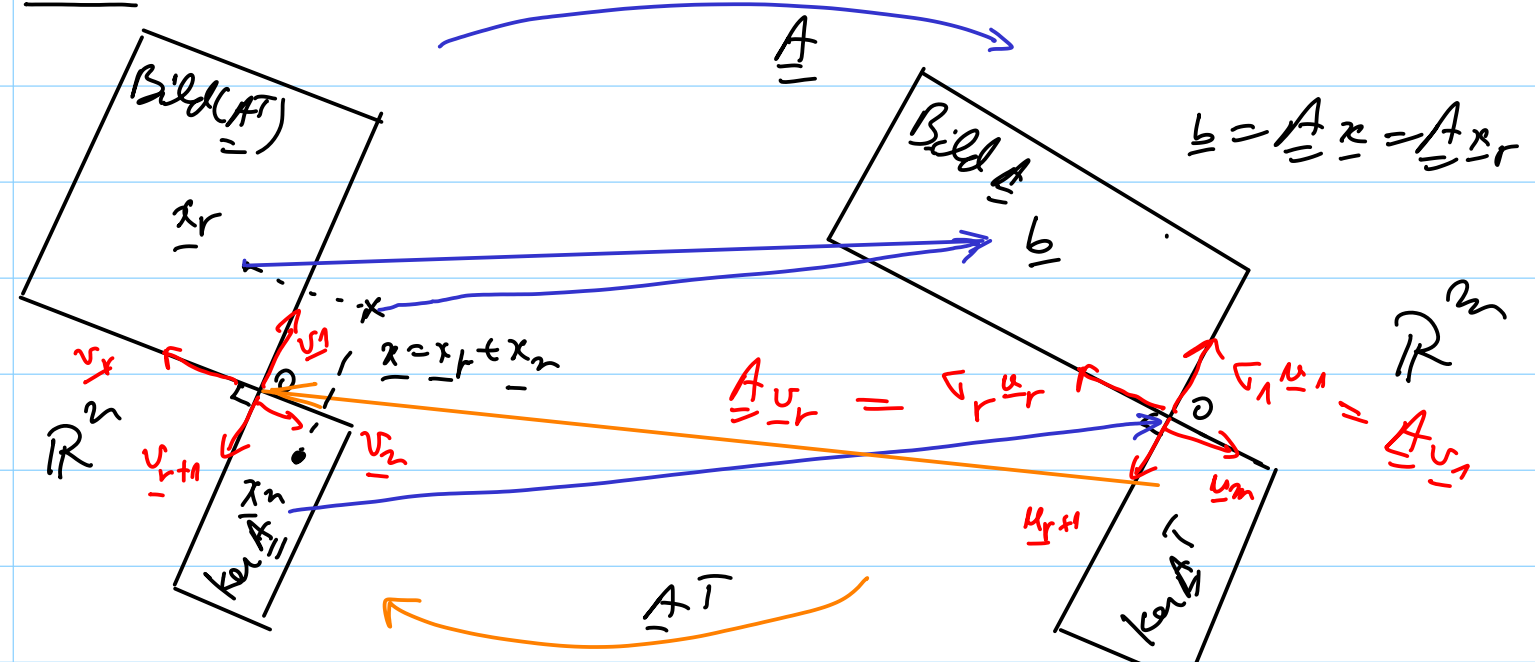
$$\underline{A} \underline{v}_1 = \sigma_1 \underline{u}_1, \dots, \underline{A} \underline{v}_r = \sigma_r \underline{u}_r$$

$$\underline{A} \underline{v}_{r+1} = 0, \dots, \underline{A} \underline{v}_n = 0$$

$$\text{Ker } \underline{A} = \text{span} \{ \underline{v}_{r+1}, \dots, \underline{v}_n \} \quad \text{ONB}$$

$$\text{Bild } \underline{A} = \text{span} \{ \underline{u}_1, \dots, \underline{u}_r \} \quad \text{ONB}$$

Beh Erinnern: Fundamentalsatz der LA:



# Anwendungen der SVD

1) Speicherplatzreduktion:

A braucht  $m \times n$  Plätze

SVD:  $r \cdot m + r \cdot n + r = r(m+n+1)$  Plätze!

$$= \max_{\|\underline{y}\|_2=1} \sqrt{\sigma_1^2 |y_1|^2 + \sigma_2^2 |y_2|^2 + \dots + \sigma_r^2 |y_r|^2} \leq$$

$$(\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0)$$

$$\leq \max_{\|\underline{y}\|_2=1} \sqrt{\sigma_1^2 \|\underline{y}\|_2^2} = \sigma_1 \text{ erreicht für } \underline{y} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}.$$

$$\Rightarrow \|\underline{A}\|_2 = \sigma_1$$

2) Berechnung der Norm einer Matrix

$$\|\underline{A}\|_2 = \max_{\|\underline{x}\|_2=1} \|\underline{A}\underline{x}\|_2 \quad \underline{U} \text{ orthogonal}$$

$$\|\underline{A}\underline{x}\|_2^2 = \|\underline{U} \underline{\Sigma} \underline{V}^T \underline{x}\|_2^2 \stackrel{\downarrow}{=} \|\underline{\Sigma} \underline{V}^T \underline{x}\|_2^2.$$

$$\|\underline{A}\|_2 = \max_{\|\underline{x}\|_2=1} \|\underline{A}\underline{x}\|_2 = \max_{\|\underline{x}\|_2=1} \|\underbrace{\underline{\Sigma} \underline{V}^T \underline{x}}_{\underline{y}}\|_2 =$$

$\underline{y} \Rightarrow \underline{x} = \underline{V} \underline{y}$

$$= \max_{\|\underline{V} \underline{y}\|_2=1} \|\underline{\Sigma} \underline{y}\|_2 = \max_{\|\underline{y}\|_2=1} \|\underline{\Sigma} \underline{y}\|_2 =$$

$$3) \underline{A} = \sigma_1 \underline{u}_1 \underline{v}_1^T + \sigma_2 \underline{u}_2 \underline{v}_2^T + \dots + \sigma_r \underline{u}_r \underline{v}_r^T$$

= Summe von Matrizen von Rang 1.

$$\{\underline{u}_1, \dots, \underline{u}_r\}, \{\underline{v}_1, \dots, \underline{v}_r\}$$

SVD  
 $\Rightarrow$  wird verwendet um eine Approximation von  $\underline{A}$  mit einer Matrix von niedrigem Rang zu bauen!

Daten mit Rausch / Fehlern:

$$\underline{A} = \sum_{j=1}^r \sigma_j \underline{u}_j \underline{v}_j^T \approx \sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^T \text{ mit } k \ll r.$$

$\neq$  Diagonalisierung verwenden wurde.

Falls  $\underline{A}$  symmetrisch:  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$

$$\underline{A} = \sum_{j=1}^n \lambda_j \underline{s}_j \underline{s}_j^T \quad \underline{A} = \underline{S} \underline{\Lambda} \underline{S}^{-1}$$

Aber:  $\rightarrow$  das geht so nur für quadratische Matrix  $\underline{A} \neq \text{SVD}$

$\rightarrow$  keine Ordnung für  $\lambda_j$ , falls  $\lambda_j$  komplex  $\neq \text{SVD}$

$\rightarrow$  keine Orthogonalität  $\underline{s}_j \neq \underline{s}_j^T$

Bsp

$$\underline{A} = \begin{bmatrix} 2 & 100 \\ 0 & 1 \end{bmatrix} \quad \lambda_1 = 2, \quad \lambda_2 = 1$$

$$\underline{A} = \underline{S} \underline{\Lambda} \underline{S}^{-1} = 2 \underbrace{\begin{bmatrix} 1 \\ 0 \end{bmatrix} \begin{bmatrix} 1 & 100 \end{bmatrix}}_{\lambda_1 \underline{s}_1 \underline{s}_1^T} + 1 \underbrace{\begin{bmatrix} 1 \\ -\frac{1}{100} \end{bmatrix} \begin{bmatrix} 0 & -100 \end{bmatrix}}_{\lambda_2 \underline{s}_2 \underline{s}_2^T}$$

Versuch der Approximation mit einer Rang 1 Matrix:  
mit EW, EV:

$$\text{Fehler} = \underline{A} - \lambda_1 \underline{s}_1 \underline{s}_1^T = \begin{bmatrix} 0 & -100 \\ 0 & 1 \end{bmatrix}$$

Fehler ist sehr gross

$$\underline{A} - \lambda_2 \underline{s}_2 \underline{s}_2^T = \begin{bmatrix} 2 & 200 \\ 0 & 0 \end{bmatrix}$$

also so gehtes NICHT ☹️

Mit SVD:

$$\underline{A} = 100,025 \begin{bmatrix} 0,9995 \\ 0,001 \end{bmatrix} \begin{bmatrix} 0,01999 & 0,9998 \end{bmatrix} + 0,02 \begin{bmatrix} 0,01 \\ 0,9995 \end{bmatrix} \begin{bmatrix} -0,9998 & 0,01999 \end{bmatrix}$$

$\sigma_1 \underline{u}_1 \underline{v}_1^T$

$\sigma_2$

$\underline{u}_2 \underline{v}_2^T$

$$\text{Fehler } \underline{A} - \sigma_1 \underline{u}_1 \underline{v}_1^T = \begin{bmatrix} 2 \cdot 10^{-4} & 0 \\ -2 \cdot 10^{-2} & 4 \cdot 10^{-4} \end{bmatrix}$$

kleiner Fehler ☺️

## Theorem [Eckart-Young]

Sei  $\underline{A} \in \mathbb{C}^{m \times n}$ . Für jedes  $k \leq \text{Rang}(\underline{A})$  gibt es eine abgebrochene SVD

$$\underline{A}_k = \sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^H$$

ist die beste Approximation von Rang  $k$  an  $\underline{A}$  im Sinne von:

$$\| \underline{A} - \underline{A}_k \| = \min_{\text{Rang}(\underline{X}) \leq k} \| \underline{A} - \underline{X} \| = \sigma_{k+1}$$

(Euklidische & Frobenius Norm).

4) PCA = principal component analysis.

Zwei Signale, gemessen an  $m=50$  Orten.

Messfehler, messen 10 Mal jedes Signal  $\Rightarrow n=20$  Messungen.

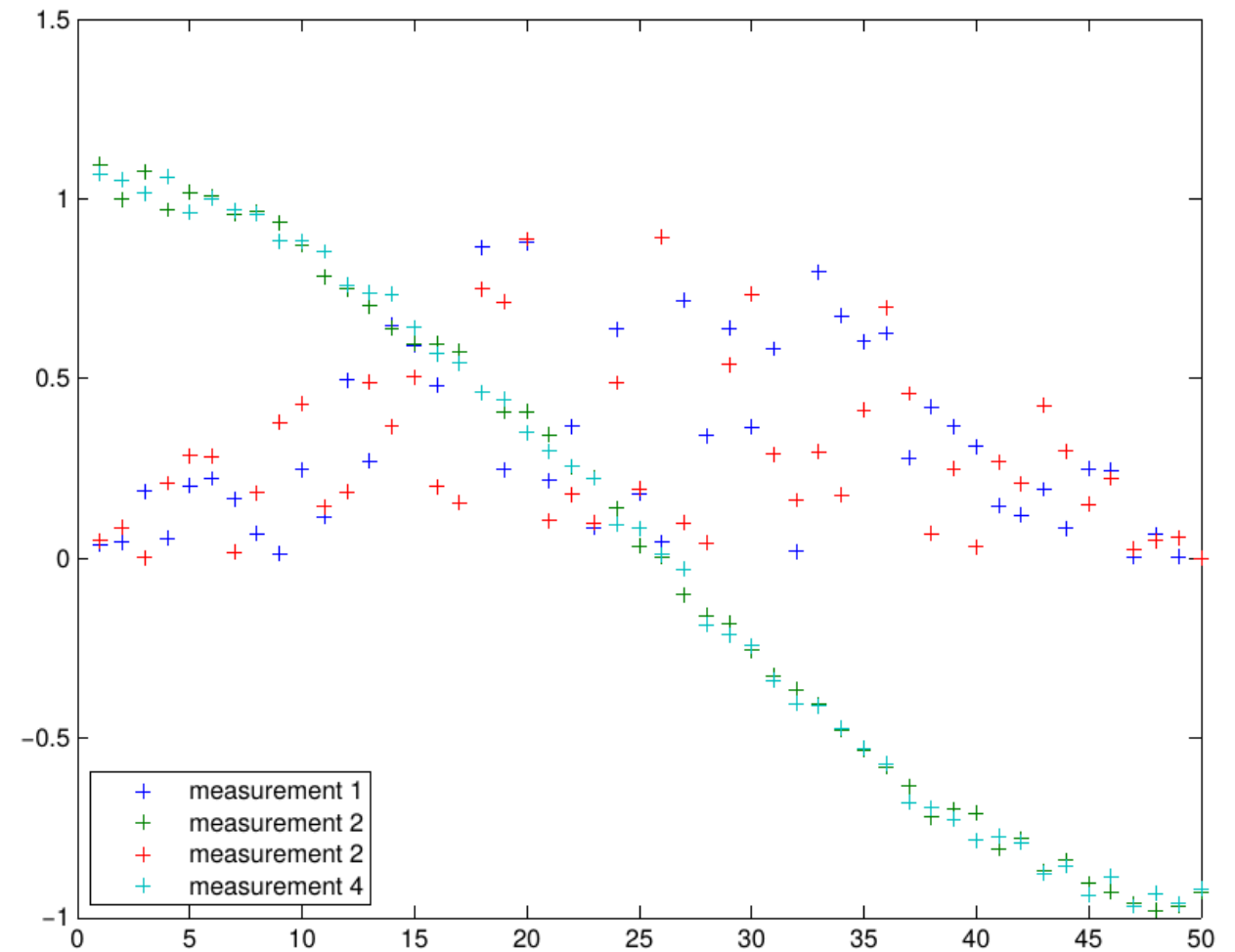
Daten in einer  $m \times n$  Matrix  $\underline{A}$  (durchgezogen)

Frage: wie viele Signale gemessen worden sind?

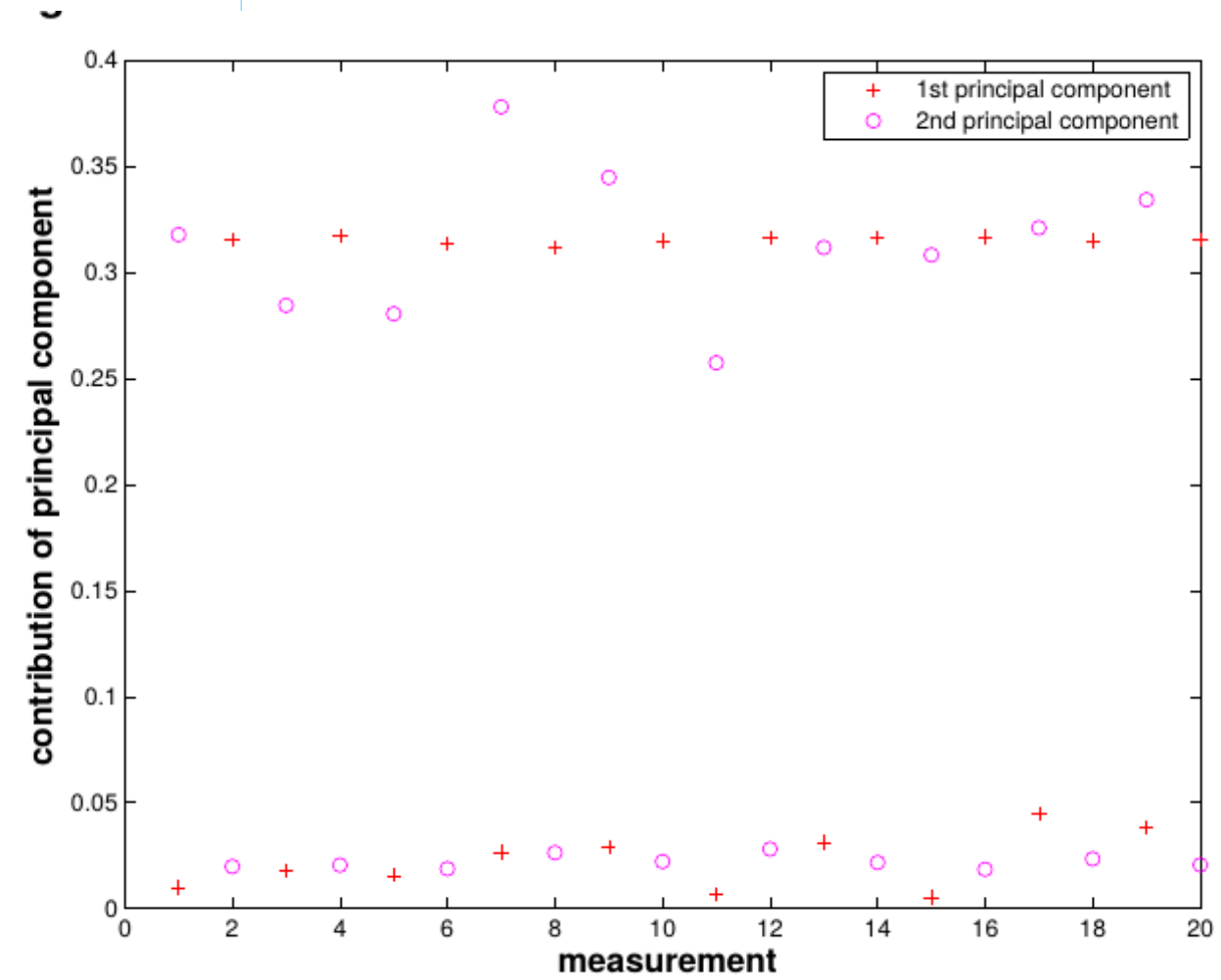
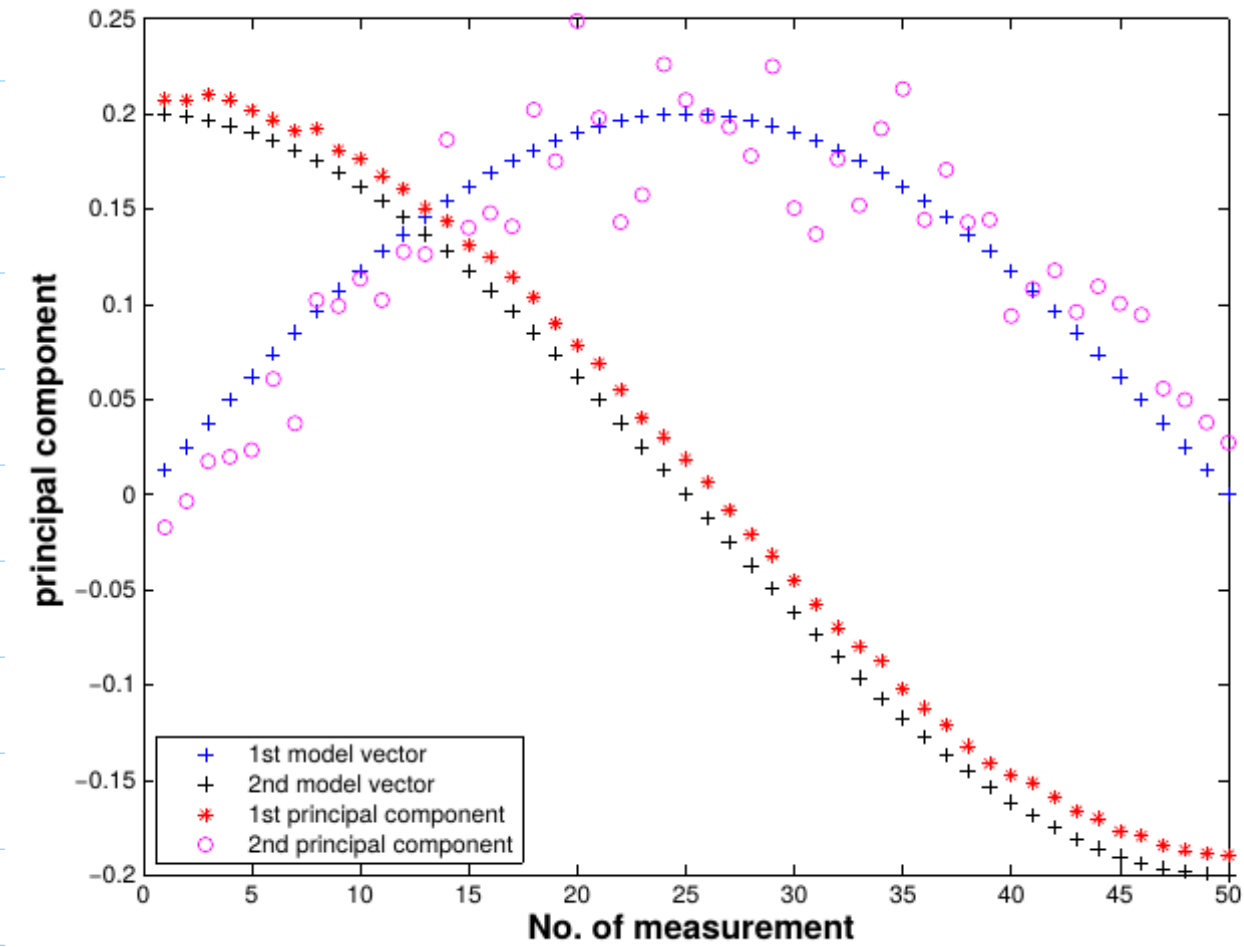
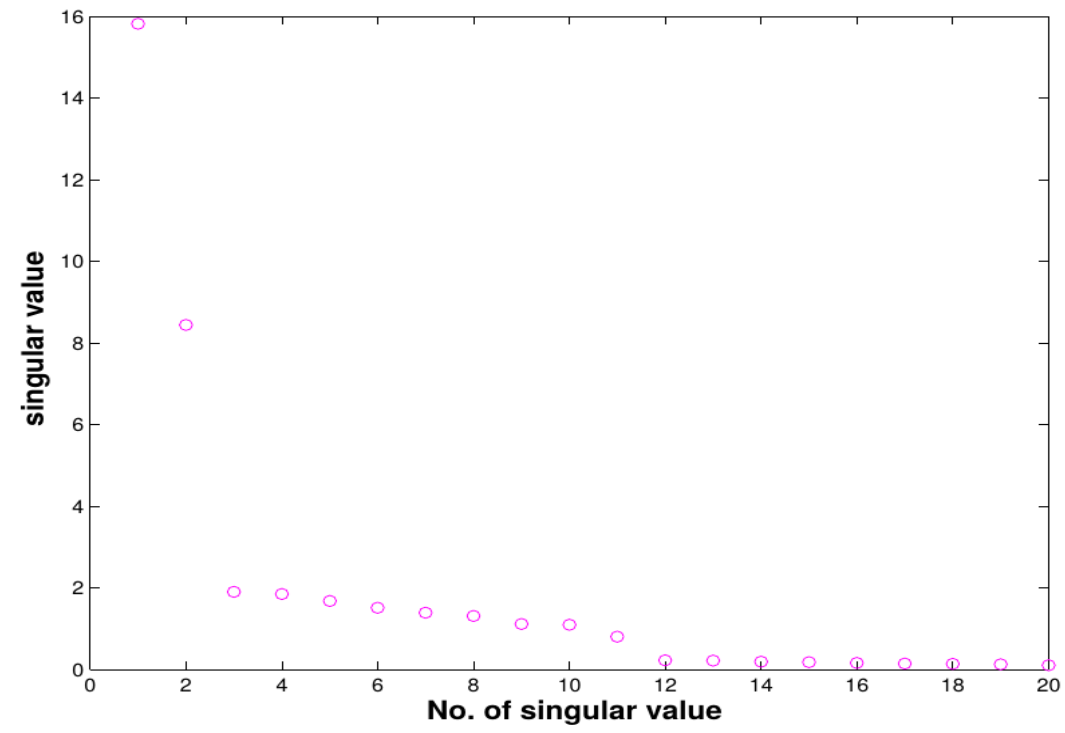
Wie sehen die echten Signale ungefähr aus?

In der Sprache der LA: Was ist  $\text{Rang}(\underline{A})$ ?

Wie viele unabhängige Vektoren / Messungen sind in  $\underline{A}$ ?



```
n = 10; m = 50
r = linspace(1,m,m)
x = sin(pi*r/m)
y = cos(pi*r/m)
A = zeros((2*n,m))
for k in range(n):
    A[2*k] = x*rand(m)
    A[2*k+1] = y+0.1*rand(m)
```



## §8 Ausgleichsrechnung

### §8.1. Motivation und Normalengleichung.

Bsp Modell  $f(t) = c \cdot e^{-at^2 + bt}$  mit physikalischen Parameter  $a, b, c$  die aus 200 Messungen zu bestimmen.

Lineares Modell:

$$\underline{y} = \underline{a}^T \underline{t} + c \cdot 1 \quad \text{mit } \underline{a} \in \mathbb{R}^n, c \in \mathbb{R} \text{ Parameter}$$

$m$  Messungen

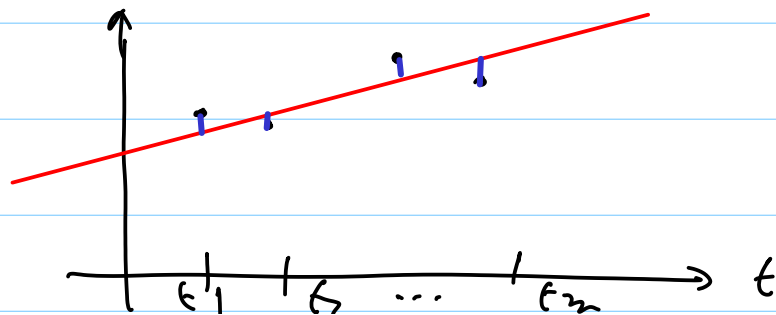
$$\underline{t}_1, \underline{t}_2, \dots, \underline{t}_m \in \mathbb{R}^n \quad \text{Messpunkte}$$

$$\downarrow \quad \downarrow \quad \dots \quad \downarrow$$

$$y_1 \quad y_2 \quad \dots \quad y_m$$

Gemessene Werte

Für  $n=1$



Vorschlag:  $\underline{a}, c$ , die  $\underline{q}, p$  das Minimum realisieren:

$$\min_{\substack{\underline{p} \in \mathbb{R}^n \\ q \in \mathbb{R}}} \sum_{i=1}^m |y_i - \underline{p}^T \underline{t}_i - 1 \cdot q|^2 = \min_{\underline{x} \in \mathbb{R}^{1+n}} \| \underline{A} \underline{x} - \underline{b} \|_2$$

$$\underline{x} = \begin{bmatrix} q \\ \underline{p} \end{bmatrix} \in \mathbb{R}^{1+n}, \quad \underline{A} = \begin{bmatrix} 1 & \underline{t}_1^T \\ 1 & \underline{t}_2^T \\ \vdots & \vdots \\ 1 & \underline{t}_m^T \end{bmatrix} \in \mathbb{R}^{m \times (1+n)}$$

$$\underline{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix}$$

Methode der kleinsten Quadrate  
(least squares Method).

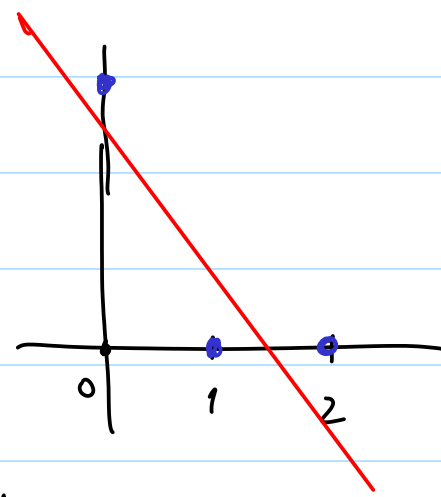
Bsp  $n=1$   $t_1=0, t_2=1, t_3=2$   
 $y_1=6, y_2=0, y_3=0$

Modell  $y = dt + c \cdot 1$

Messpunkte  $\begin{bmatrix} 0 \\ 1 \\ 2 \end{bmatrix}$ , Messwerte  $\begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$

$$\underline{A} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad \underline{b} = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$$

$$\min_{\underline{x} \in \mathbb{R}^2} \|\underline{A}\underline{x} - \underline{b}\|_2$$



alg. Problem:  $\begin{cases} b_1 = dt_1 + c_1 \\ b_2 = dt_2 + c_2 \\ b_3 = dt_3 + c_3 \end{cases} \Rightarrow \underline{A}\underline{x} = \underline{b}$   
 hat keine Lösung.

Möchte aber "Lösung" im Sinne.

$$\|\underline{A}\underline{x} - \underline{b}\|_2 \stackrel{!}{=} \text{Min!}$$

Residuum

I Algebraische Methode:

$$\underline{A}^T \mid \underline{A}\underline{x} = \underline{b} \Rightarrow \underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b}$$

quadratisch, symmetrisch.

Falls  $\text{Rang}(\underline{A}) = n \Rightarrow \underline{A}^T \underline{A}$  invertierbar  $\Rightarrow$

kann das LGS eindeutig lösen, am besten mit QR-Zerlegung!

weil stabiler als LU-Zerlegung!

$$\text{cond}(\underline{A}) \text{ gross} \Rightarrow \text{cond}(\underline{A}^T \underline{A}) = \text{cond}(\underline{A})^2$$

$$\text{cond}(\underline{A}) = 10^2 \rightarrow \text{cond}(\underline{A}^T \underline{A}) = 10^4 \rightarrow \text{LU!}$$

$$\underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b} \quad \text{Normalengleichung.}$$



Ben Das ist in der Tat, auch die Lösung von Min.-Problem:

$$\underline{A}: \mathbb{R}^n \rightarrow \mathbb{R}^m \quad ; \quad \underline{b} \in \mathbb{R}^m \Rightarrow \underline{b} = \underset{\substack{\uparrow \\ \text{Bild}(\underline{A})}}{\underline{p}} + \underset{\substack{\uparrow \\ \text{Ker}(\underline{A}^T)}}{\underline{e}}$$

$\psi$   
 $\underline{x} \mapsto \underline{A}\underline{x}$

es gibt eine Lösung  $\hat{\underline{x}}$  von  $\underline{A}\hat{\underline{x}} = \underline{p}$

$\underline{p} \perp \underline{e}$

Fehler  $\underline{b} - \underline{A}\hat{\underline{x}} = \underline{b} - \underline{p} = \underline{e} \in \text{Ker}(\underline{A}^T) \perp \text{Bild}(\underline{A}) \Rightarrow$

$$\|\underline{b} - \underline{A}\hat{\underline{x}}\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \|\underline{b} - \underline{A}\underline{x}\|_2^2 \Rightarrow \hat{\underline{x}} \text{ ist die Lösung des Min-Problems.}$$

$$\underline{A}^T \mid \Rightarrow \underline{A}^T \underline{b} = \underline{A}^T \underline{p} + \underbrace{\underline{A}^T \underline{e}}_{\substack{\perp \\ \underline{e} \in \text{Ker}(\underline{A}^T)}} = \underline{A}^T \underline{p} = \underline{A}^T \underline{A} \hat{\underline{x}}$$

## II Analytische Methode

$$f(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2, \quad f: \mathbb{R}^{1+n} \rightarrow [0, \infty)$$

Min=?

$$\Downarrow f(\underline{x}) = 0$$

$$f(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2 = (\underline{A}\underline{x} - \underline{b})^T (\underline{A}\underline{x} - \underline{b}) =$$

$$= \underbrace{\underline{x}^T \underline{A}^T \underline{A} \underline{x}}_{\underline{M} \text{ symmetrisch.}} - \underline{b}^T \underline{A} \underline{x} - \underline{x}^T \underline{A}^T \underline{b} + \underline{b}^T \underline{b} =$$

$$= \underline{x}^T \underline{M} \underline{x} - 2 \underline{x}^T \underbrace{\underline{A}^T \underline{b}}_{\underline{w}} + \underline{b}^T \underline{b} =$$

$$= \sum_{i=1}^n x_i \sum_{j=1}^n m_{ij} x_j - 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} =$$

$$= x_1 \sum_{j=1}^n m_{1j} x_j + \sum_{i=2}^n x_i \sum_{j=1}^n m_{ij} x_j - 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} =$$

$$\sum_{j=1}^n x_j \sum_{i=2}^n x_i m_{ij}$$

$$= x_1 m_{11} x_1 + x_1 \sum_{j=2}^n m_{1j} x_j + x_1 \sum_{i=2}^n x_i m_{i1} + \sum_{j=2}^n x_j \sum_{i=2}^n x_i m_{ij} -$$

$$- 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} \quad \Rightarrow \quad \downarrow M \text{ symmetrisch.}$$

$$= m_{11} x_1^2 + 2x_1 \sum_{i=2}^n x_i m_{i1} + \sum_{i,j=2}^n x_i m_{ij} x_j - 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b}$$

$$\frac{\partial}{\partial x_1} f(\underline{x}) = 2m_{11}x_1 + 2 \sum_{i=2}^n x_i m_{i1} - 2w_1 \quad \Rightarrow$$

$$Df(\underline{x}) = 2 \underline{M} \underline{x} - 2 \underline{A}^T \underline{b} =$$

$$= 2 \underline{A}^T \underline{A} \underline{x} - 2 \underline{A}^T \underline{b}.$$

kritische Punkte  $Df(\underline{x}) = 0 \Leftrightarrow$

$$\underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b} \quad \text{Normalengleichung.}$$

Minimum nur wenn  $D^2 f(\underline{x}) = 2 \underline{A}^T \underline{A}$

symmetrisch positiv definit

$\uparrow$   
nur wenn  $\text{Rang}(\underline{A}) = n$ .

Bem  $\text{cond}_2 \underline{M} = \|\underline{M}^{-1}\|_2 \|\underline{M}\|_2$

↳ wie sehr sich die Rundungsfehler sich auf das Lösen von  $\underline{M}\underline{x} = \underline{b}$  auswirken.

$\text{cond}(\underline{M}) \gg 1 \Rightarrow$  grosse Rundungsfehler (:(

$$\begin{aligned} \text{cond}_2(\underline{A}^T \underline{A}) &= \text{cond}_2(\underline{V} \underline{\Sigma}^T \underline{U}^T \underline{U} \underline{\Sigma} \underline{V}^T) = \text{cond}_2(\underline{V} \underline{\Sigma}^2 \underline{V}^T) \\ &= \text{cond}(\underline{\Sigma}^2) = \frac{\sigma_1^2}{\sigma_r^2} = \text{cond}(\underline{A})^2 \end{aligned}$$

wobei  $\sigma_1, \dots, \sigma_r$  die Singulärwerte von  $\underline{A}$  sind.

Bsp  $f \in L^2 = \{f: I \rightarrow \mathbb{R} ; \int_I |f(t)|^2 dt < \infty\}$

$f_n \in V_n, \dim V_n < \infty$

↙ Basis  $b_1, \dots, b_n$

$$f_n = \sum_{j=1}^n x_j b_j \Rightarrow f_n(t) = \sum_{j=1}^n x_j b_j(t)$$

verwende Messungen  $y_i \approx f(t_i)$

Aufgabe: gegeben  $(t_i, y_i)$  für  $i=1, 2, \dots, m$  Messungen,

finde die beste Approximation in  $V_n$

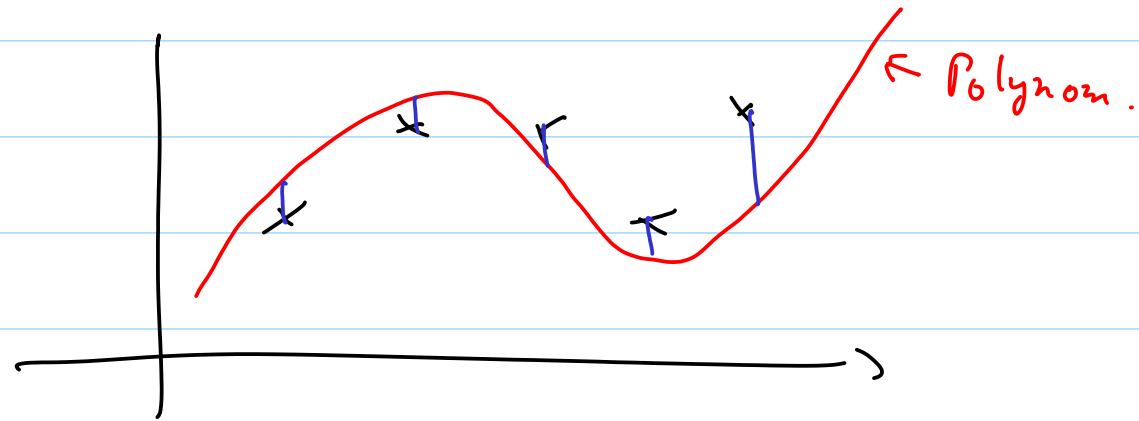
an  $f$ , d.h. finde  $x_1, x_2, \dots, x_n$  so dass

$$\sum_{i=1}^m |f_n(t_i) - y_i|^2 = \text{minimal} !$$

$$\begin{cases} \sum_{j=1}^n b_j(t_i) x_j = y_i \\ i = 1, 2, \dots, m \end{cases} \quad \underline{A} = \begin{bmatrix} b_j(t_i) \end{bmatrix}_{\substack{i=1, 2, \dots, m \\ j=1, 2, \dots, n}}$$

Ben Spezielle Wahl:  $V_n = P_n$  = Polynome vom Grad  
Maximal  $n-1$

Basis in  $V_n$ : Monome  $b_j(t) = t^{j-1}$   
 $\rightarrow$  polynomiales fit; polyfit.



## §8.2 Lösung mittels orthogonalen Transformationen

Fall 1  $\text{rang}(\underline{A}) = n$

$$\underline{A} = \underline{Q} \underline{R} = \underline{Q} \begin{bmatrix} \tilde{\underline{R}} \\ 0 \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix}$$

$\underline{Q}$  orthogonal

$$\min_{\underline{x} \in \mathbb{R}^n} \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \|\underline{Q}\underline{R}\underline{x} - \underline{I}\underline{b}\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underline{Q}\underline{R}\underline{x} - \underline{Q}\underline{Q}^H \underline{b}\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underline{Q}(\underline{R}\underline{x} - \underline{Q}^H \underline{b})\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underbrace{\underline{R}\underline{x} - \underline{Q}^H \underline{b}}_{\underline{\beta}}\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \left\| \begin{bmatrix} \tilde{\underline{R}} \\ 0 \end{bmatrix} \begin{bmatrix} \underline{x} \end{bmatrix} - \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \\ \beta_{n+1} \\ \vdots \\ \beta_m \end{bmatrix} \right\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \left\| \begin{bmatrix} \tilde{\underline{R}}\underline{x} \\ 0 \\ \vdots \\ 0 \end{bmatrix} - \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_n \\ \beta_{n+1} \\ \vdots \\ \beta_m \end{bmatrix} \right\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underline{\tilde{R}}\underline{x} - \underline{\tilde{\beta}}\| + |\beta_{n+1}|^2 + \dots + |\beta_n|^2$$

|| für  $\underline{x}$  die Lösung von  $\underline{\tilde{R}}\underline{x} = \underline{\tilde{\beta}}$

$$\begin{bmatrix} * & 0 \\ 0 & \end{bmatrix} \begin{bmatrix} x \\ \end{bmatrix} = \begin{bmatrix} \tilde{\beta} \\ \end{bmatrix}$$

Bsp  $\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$  Normalen Gleichung.

$$\Rightarrow \begin{bmatrix} 3 & 3 \\ 3 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 0 \end{bmatrix} \Rightarrow \begin{matrix} x_1 = 5 \\ x_2 = -3 \end{matrix}$$

QR-Zerlegung

$$\underline{A} = \underline{Q} \underline{R} = \begin{bmatrix} -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & ? \\ -\frac{1}{\sqrt{3}} & 0 & ? \\ -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & ? \end{bmatrix} \begin{bmatrix} -\frac{3}{\sqrt{3}} & -\frac{3}{\sqrt{3}} \\ 0 & -\frac{1}{\sqrt{2}} \\ 0 & 0 \end{bmatrix}$$

numpy.linalg.qr( $\underline{A}$ )

$$\text{LGS} \quad \underline{\tilde{R}}\underline{x} = \underline{Q}^T \underline{b}$$

$$\begin{bmatrix} -\frac{3}{\sqrt{3}} & -\frac{3}{\sqrt{3}} \\ 0 & -\frac{1}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -\frac{6}{\sqrt{3}} \\ \frac{6}{\sqrt{2}} \end{bmatrix} \Rightarrow$$

$$\Rightarrow x_1 = 5, x_2 = -3.$$

Fall 2  $\text{rang}(\underline{A}) < n \Rightarrow \text{SVD!}$

$$\underline{A} = \begin{bmatrix} \underline{U}_1 & \underline{U}_2 \end{bmatrix} \begin{bmatrix} \underline{\Sigma}_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \underline{V}_1^T \\ \underline{V}_2^T \end{bmatrix}$$

$$\underline{\Sigma}_r = \begin{bmatrix} \sigma_1 & \dots & 0 \\ 0 & \dots & \sigma_r \end{bmatrix} \text{ mit } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

$\underline{U}$  orthogonal

$$\|\underline{A}\underline{x} - \underline{b}\|_2 = \left\| \begin{bmatrix} \underline{\Sigma}_r \underline{V}_1^T \underline{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \underline{U}_1^T \underline{b}_1 \\ \underline{U}_2^T \underline{b}_2 \end{bmatrix} \right\|_2 \Rightarrow$$

Minimal für  $\underline{x}$  Lösung von  $\sum_r \underline{v}_r^T \underline{x} = \underline{u}_1^T \underline{b}_1$

$$\underline{x} = \underline{v}_1 \sum_r^{-1} \underline{u}_1^T \underline{b}_1$$

und das Minimum ist  $\| \underline{u}_2^T \underline{b}_2 \|_2$

→ Standard in Codes  $\text{lsq}(\underline{A}, \underline{b})$

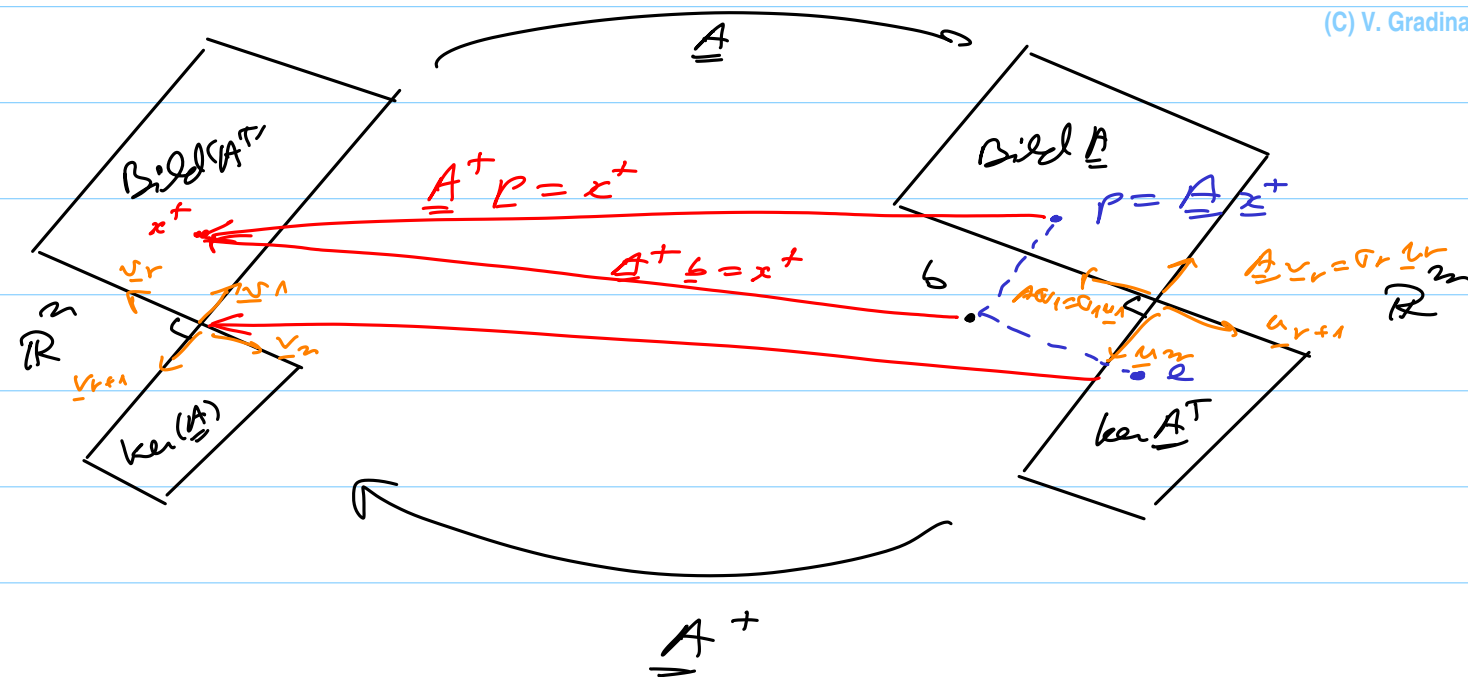
Def Pseudoinverse  $\underline{A}^+ := \underline{V} \underline{\Sigma}^+ \underline{U}^T$

$$\underline{\Sigma}^+ = \begin{bmatrix} \sigma_1^{-1} & & & 0 \\ & \sigma_2^{-1} & & \\ & & \ddots & \\ & & & \sigma_r^{-1} & & 0 \\ 0 & & & & 0 & \ddots & 0 \end{bmatrix} \text{ Moore-Penrose inverse.}$$

$\underline{x} \in \mathbb{R}^n \xrightarrow{\underline{A}} \mathbb{R}^m \underline{Ax}$

$\underline{x}^+ \quad \underline{A}^+ \quad \underline{b}$  die Ausgleichslösung.

$\| \underline{Ax}^+ - \underline{b} \|_2^2 \text{ min!}$



Bem Was ist wenn wir später entdecken,  
1) bestimmte Messungen falsch waren?

$\underline{Q} \underline{R}$ -Methode ist da besser.

scipy. `qr_delete(row)`

2) bestimmte Basisfunktionen gar nicht dabei sein sollen?

`qr_delete(column)`.

### §8.3. Lineare Ausgleichsrechnung mit lin. Nebenbedingungen

Finde  $\underline{x} \in \mathbb{R}^n$  bei gegebenen  $\underline{A} \in \mathbb{R}^{m \times n}$   
 $m \geq n = \text{rang}(\underline{A})$

$$\left\{ \begin{array}{l} \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \text{minimal!} \\ \underline{b} \in \mathbb{R}^m \end{array} \right.$$

$$\left\{ \begin{array}{l} \underline{C}\underline{x} = \underline{d} \end{array} \right. \quad \text{mit } \underline{C} \in \mathbb{R}^{p \times n} \quad \text{mit } p < n$$

$$\underline{d} \in \mathbb{R}^p \quad \text{rang}(\underline{C}) = p$$

$$\left\{ \begin{array}{l} \frac{\partial L}{\partial \underline{x}}(\underline{x}, \underline{m}) = 0 \\ \frac{\partial L}{\partial \underline{m}}(\underline{x}, \underline{m}) = 0 \end{array} \right. \quad \left\{ \begin{array}{l} \underline{A}^T(\underline{A}\underline{x} - \underline{b}) + \underline{C}^T \underline{m} = 0 \\ 0 + \underline{C}\underline{x} - \underline{d} = 0 \end{array} \right.$$

$$\begin{bmatrix} \underline{A}^T \underline{A} & \underline{C}^T \\ \underline{C} & \underline{0} \end{bmatrix} \begin{bmatrix} \underline{x} \\ \underline{m} \end{bmatrix} = \begin{bmatrix} \underline{A}^T \underline{b} \\ \underline{d} \end{bmatrix}$$

Methode 1 Lagrange-Multiplikatoren  $\underline{m} \in \mathbb{R}^p$

$$L(\underline{x}, \underline{m}) := \frac{1}{2} \|\underline{A}\underline{x} - \underline{b}\|_2^2 + \underline{m}^T (\underline{C}\underline{x} - \underline{d})$$

$$\min_{\substack{\underline{x} \in \mathbb{R}^n \\ \underline{m} \in \mathbb{R}^p}} L(\underline{x}, \underline{m})$$

Block-LR-Zerlegung

1) Cholesky-Zerlegung von  $\underline{A}^T \underline{A} = \underline{R}^T \underline{R}$

2) berechne  $\underline{G}$  aus  $\underline{R}^T \underline{G}^T = \underline{C}^T$

3)  $\underline{S}$  aus Cholesky-Zerlegung.

$$\underline{S}^T \underline{S} = -\underline{G} \underline{G}^T$$

$$\begin{bmatrix} A^T A & C^T \\ C & 0 \end{bmatrix} = \begin{bmatrix} R^T & 0 \\ G & S \end{bmatrix} \begin{bmatrix} R & G^T \\ 0 & S^T \end{bmatrix}$$

Cholestky-Zerlegung!

Nachteil:  $\text{cond}(\underline{\underline{A}}^T \underline{\underline{A}}) = \text{cond}(\underline{\underline{A}})^2$  ☹️

Vorteil: falls  $\underline{\underline{A}}, \underline{\underline{C}}$  eine Struktur haben, dann kann man Algorithmen nehmen, die diese Struktur erhalten.

Methode 2 SVD besser für die Kondition, verliert die Struktur

$$\underline{\underline{C}} = \underline{\underline{U}} \begin{bmatrix} \underline{\underline{\Sigma}} & \underline{\underline{0}} \end{bmatrix} \begin{bmatrix} \underline{\underline{V}}_1^H \\ \underline{\underline{V}}_2^H \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}$$

$$\text{Ker}(\underline{\underline{C}}) = \text{Bild}(\underline{\underline{V}}_2).$$

Trick: Definiere  $\underline{\underline{x}}_0 = \underline{\underline{V}}_1 \underline{\underline{\Sigma}}_1^{-1} \underline{\underline{U}}^H \underline{\underline{d}}$

dann suche  $\underline{\underline{x}} = \underline{\underline{x}}_0 + \underbrace{\underline{\underline{V}}_2 \underline{\underline{y}}}_{\substack{p \\ \text{Ker}(\underline{\underline{C}}) = \text{Bild}(\underline{\underline{V}}_2)}} \quad \text{mit } \underline{\underline{y}} \in \mathbb{R}^{n-p}$

$$(\Rightarrow \underline{\underline{C}} \underline{\underline{x}}_0 = \underline{\underline{d}}).$$

$$\|\underline{\underline{A}} \underline{\underline{x}} - \underline{\underline{b}}\|_2^2 = \|\underline{\underline{A}} \underline{\underline{x}}_0 + \underline{\underline{A}} \underline{\underline{V}}_2 \underline{\underline{y}} - \underline{\underline{b}}\|_2^2 = \|\underbrace{\underline{\underline{A}} \underline{\underline{V}}_2}_{\substack{\text{min noch } \underline{\underline{y}}}} \underline{\underline{y}} - \underbrace{(\underline{\underline{b}} - \underline{\underline{A}} \underline{\underline{x}}_0)}_{\text{fest}}\|_2^2$$



Das ist ein standard Ausgleichsproblem!

### §8.4. Totales lineares Ausgleichsproblem

total = Messfehler auch in den Messorte  $\underline{t}_i$

↗  $\underline{A}\underline{x} = \underline{b}$  mit  $\underline{b}$  und  $\underline{A}$  fehlerhaft.

Bem  $\underline{A}\underline{x} = \underline{b}$  hat keine Lösung falls  $\underline{b} \notin \text{Bild}(\underline{A})$ .

Ziel: finde  $\hat{\underline{A}}, \hat{\underline{x}}, \hat{\underline{b}}$  so dass

$$\hat{\underline{A}}\hat{\underline{x}} = \hat{\underline{b}} \quad (\hat{\underline{b}} \in \text{Bild}(\hat{\underline{A}})) \quad \hat{\underline{A}} \in \mathbb{R}^{m \times n}, \hat{\underline{b}} \in \mathbb{R}^m$$

"so nah wie möglich" an  $\underline{A}\underline{x} = \underline{b}$

$$\underline{C} = \begin{bmatrix} \underline{A} & \underline{b} \end{bmatrix} ; \quad \hat{\underline{C}} = \begin{bmatrix} \hat{\underline{A}} & \hat{\underline{b}} \end{bmatrix}$$

$$\|\underline{C} - \hat{\underline{C}}\|_F = \min!$$

Frobenius norm

$$\|\underline{M}\|_F^2 = \sum_{ij} |M_{ij}|^2$$

mit  $\hat{\underline{b}} \in \text{Bild}(\hat{\underline{A}})$

$$\text{rang}(\underline{A}) = n \quad \text{auch} \quad \text{Rang}(\hat{\underline{C}}) = n.$$

Lösung: das ist die Niedrigrangapproximation von  $\underline{C}$

$$\underline{C} = \underline{U} \underline{\Sigma} \underline{V}^H = \sum_{j=1}^{n+1} \sigma_j \underline{u}_j \underline{v}_j^H$$

Definiere 
$$\hat{\underline{C}} = \sum_{j=1}^n \sigma_j \underline{u}_j \underline{v}_j^H$$

↪ optimale Approximation an  $\underline{C}$  in der Menge der Matrizen von Rang  $n$  ist.

Noch dazu wissen wir:  $\underline{v}_j$  orthogonal  $\Rightarrow \hat{\underline{C}} \underline{v}_{n+1} = 0$

Falls  $v_{n+1, n+1} \neq 0$

$$\hat{\underline{A}} \hat{\underline{x}} = \hat{\underline{b}} \Rightarrow \hat{\underline{x}} = - \frac{1}{v_{n+1, n+1}} \underline{v}_{n+1}.$$

Beweis

$$\hat{\underline{A}} \hat{\underline{x}} = \hat{\underline{b}} \Leftrightarrow \hat{\underline{A}} \hat{\underline{x}} - \hat{\underline{b}} = 0 \Leftrightarrow \begin{bmatrix} \hat{\underline{A}} & \hat{\underline{b}} \end{bmatrix} \begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} = 0$$

$$\Leftrightarrow \begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} \in \ker(\hat{\underline{C}})$$

$$\dim \ker(\hat{\underline{C}}) = 1, \quad \underline{v}_{n+1} \in \ker \hat{\underline{C}}$$

$$\begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} = k \underline{v}_{n+1}$$

Definiere  $k = - \frac{1}{v_{n+1, n+1}} \Rightarrow \hat{\underline{x}} = - \frac{1}{v_{n+1, n+1}} \underline{v}_{n+1}.$

## § 8.5 Nichtlineare Ausgleichsrechnung.

Modell  $f(t, \underline{x}) = y$

Parameters, aus Messungen zu bestimmen.

$$f(t_i, \underline{x}) \approx y_i, \quad i = 1, 2, \dots, m$$

$$\underline{F}(\underline{x}) = \begin{bmatrix} f(t_1, \underline{x}) - y_1 \\ f(t_2, \underline{x}) - y_2 \\ \vdots \\ f(t_m, \underline{x}) - y_m \end{bmatrix}; \quad \underline{F}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Finde  $\underline{x}^* \in \mathbb{R}^n$   
 $\|F(\underline{x}^*)\|_2$  Minimal!

$$\Phi(\underline{x}) = \frac{1}{2} \|F(\underline{x})\|_2^2, \quad \Phi: \mathbb{R}^n \rightarrow [0, \infty[$$

Minimum von  $\Phi$ :

Finde  $\underline{x}^* : D\Phi(\underline{x}^*) = 0 \leftarrow$  Nullstellensuche!

mit Newton-Verfahren  $\Rightarrow$  quadratische Konvergenz  
lokal!

$$D\Phi(\underline{x}) = D\left(\frac{1}{2}\underline{F}(\underline{x})^T \underline{F}(\underline{x})\right) = \underline{D}\underline{F}(\underline{x})^T \underline{F}(\underline{x})$$

$$\underline{D}\Phi(\underline{x}^*) = 0 \Leftrightarrow \underline{D}\underline{F}(\underline{x}^*)^T \underline{F}(\underline{x}^*) = 0$$

dafür Newton!

In Newton brauchen wir  $D(D\Phi(\underline{x})) = H\Phi(\underline{x})$   
Hesse Matrix von  $\Phi$

$$H\Phi(\underline{x}) = D\left(\underline{D}\underline{F}(\underline{x})^T \underline{F}(\underline{x})\right) = \underline{D}\underline{F}(\underline{x})^T \underline{D}\underline{F}(\underline{x}) +$$

$$+ \sum_{j=1}^m \underline{F}_j(\underline{x}) \underline{D}^2 \underline{F}_j(\underline{x})$$

Matrix

$$\left(\underline{H}\Phi(\underline{x})\right)_{ik} = \sum_{j=1}^m \left( \frac{\partial \underline{F}_j}{\partial x_k}(\underline{x}) \frac{\partial \underline{F}_j}{\partial x_i}(\underline{x}) + \underline{F}_j(\underline{x}) \frac{\partial^2 \underline{F}_j(\underline{x})}{\partial x_i \partial x_k} \right)$$

Newton-Schritt  $\underline{x}^{(k+1)} := \underline{x}^{(k)} + \underline{\rho}$   
Newton-Korrektur  $\underline{\rho}$  aus LGS

$$\underline{H}\Phi(\underline{x}^{(k)}) \underline{\rho} = - \underline{D}\underline{F}(\underline{x}^{(k)})^T \underline{F}(\underline{x}^{(k)})$$

Alternative: Gauss-Newton-Verfahren.

Idee: linearisiere Lokal in jedem Schritt  
einer Iteration.  $\Rightarrow$  eine Folge von  
lkn. Ausgleichsprobleme

Linearisierung.

$$\arg\min_{\underline{x} \in \mathbb{R}^n} \|\underline{F}(\underline{x})\|_2^2 \simeq \arg\min_{\underline{x} \in \mathbb{R}^n} \|\underline{F}(\underline{x}^0) + \underline{D}\underline{F}(\underline{x}^0)(\underline{x} - \underline{x}^0)\|_2^2$$

$$= \arg\min_{\underline{x} \in \mathbb{R}^n} \|\underline{A}\underline{x} - \underline{b}\|_2^2$$

mit  $\underline{A} = \underline{D}\underline{F}(\underline{x}^0)$  und  $\underline{b} = \underline{F}(\underline{x}^0) - \underline{D}\underline{F}(\underline{x}^0)\underline{x}^0$

$$\underline{x}^{(0)} \xrightarrow{\text{lin LSP}} \underline{x}^{(1)} \xrightarrow{\text{lin LSP}} \underline{x}^{(2)} \rightarrow \dots \rightarrow \underline{x}^{(k)}$$

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} + \underline{\rho}^{(k)}, \quad \underline{\rho}^{(k)} = \underset{\underline{\rho}}{\operatorname{argm}} \| F(\underline{x}^{(k)}) + \underline{D}F(\underline{x}^{(k)})\underline{\rho} \|_2^2$$

Nochteil: lineare Konvergenz.

Professionelle Software: noch raffinierter Alg.  
Levenberg-Marquadt. Verfahren.

$$\min_{\underline{\rho}} \left\{ \| F(\underline{x}^{(k)}) + \underline{D}F(\underline{x}^{(k)})\underline{\rho} \|_2^2 + \lambda^{(k)} \|\underline{\rho}\|_2^2 \right\}$$

$\hookrightarrow \lambda^{(k)} > 0$ , heuristisch gewählt.

Totales Ausgleichsproblem 'software'.  
numpy.odr.



## §9 Eigenwerte

Software:  $\text{eig}(\underline{A})$  kostet  $O(N^3)$  für  $\underline{A} \in \mathbb{R}^{n \times n}$   
 $\text{eigh}(\underline{A})$  kostet  $O(N^2)$  für  $\underline{A}^H = \underline{A}$ .



Bem Die Eigenvektoren (EV) werden  
 typischerweise nicht so gut approximiert)

Bsp  $\begin{cases} \dot{\underline{y}} = \underline{A} \underline{y} \\ \underline{y}(0) = \underline{y}_0 \end{cases} \quad \underline{A} \in \mathbb{R}^{n \times n}$

$N=2, N=3$  ok:  $\underline{y}(t) = e^{\underline{A}t} \underline{y}(0)$

$e^{\underline{A}t} = \sum_{n=0}^{\infty} (\underline{A}t)^n$   

$N \approx 10 \quad \underline{A} = \underline{S} \underline{D} \underline{S}^{-1}$  falls  $\underline{A}$  diagonalisierbar

$\Rightarrow e^{\underline{A}t} = \underline{S} e^{\underline{D}t} \underline{S}^{-1}$    $N_{\text{gross}}$  

Für "kleines"  $N < 20$ :  $e^{\underline{A}t} = \text{expm}(\underline{A}t)$

Algorithmus Padé-Approximation  
 $e^x \approx \frac{p_n(x)}{q_n(x)}$  etc.

Bsp  $\underline{A} = \begin{bmatrix} 5 & 1 & 0 \\ 0 & 5 & 1 \\ 0 & 0 & 5 \end{bmatrix}$  EW: 5 dreifach  
 EV:  $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

$\Rightarrow$  nicht diagonalisierbar

$\begin{bmatrix} \dot{y}_1 \\ \dot{y}_2 \\ \dot{y}_3 \end{bmatrix} = \begin{bmatrix} 5 & 1 & 0 \\ 0 & 5 & 1 \\ 0 & 0 & 5 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$   
 $\dot{y}_1 = 5y_1 + y_2$   
 $\dot{y}_2 = 5y_2 + y_3$   
 $\dot{y}_3 = 5y_3 \Rightarrow y_3(t) = y_3(0) e^{5t}$

$\dot{y}_2 = 5y_2 + y_3(0) e^{5t} \Rightarrow y_2(t) = (y_2(0) + t y_3(0)) e^{5t}$

$\Rightarrow y_1(t) = (y_1(0) + t y_2(0) + \frac{t^2}{2} y_3(0)) e^{5t}$

Falls  $\underline{A} = \begin{bmatrix} & * \\ 0 & \end{bmatrix}$  dann können wir  
 $\underline{\dot{y}} = \underline{A} \underline{y}$

exakt lösen!

Ben

Wenn numerische Berechnungen zu machen sind,  
 dann NIEMALS Jordan's Normalform,  
 da sehr instabil.

Numerik: Schur-Zerlegung.

§ 9.2. Potenzmethoden.  $\underline{A}$   $n \times n$  Matrix

— 1904 Nobel-Preis Physik

Def Rayleigh-Quotienten:

$$\rho_A(\underline{x}) = \frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}}, \quad \underline{x} \neq 0$$

Ben  $\underline{x}$  EV von  $\underline{A}$ :  $\underline{A} \underline{x} = \lambda \underline{x}$

$$\Rightarrow \rho_A(\underline{x}) = \frac{\underline{x}^H \lambda \underline{x}}{\underline{x}^H \underline{x}} = \lambda$$

Ben  $\rho_A(\underline{x}) = \arg \min_{\alpha} \| \underline{A} \underline{x} - \alpha \underline{x} \|_2$

Stationärpunkte von  $\rho_A(\underline{x}) \Rightarrow$  EV.

Taylor für  $\rho_A(\underline{x})$  um EV  $\Rightarrow$

$$\rho_A(\underline{x}) - \rho_A(\text{EV}) = O(\|\underline{x} - \text{EV}\|_2^2) \quad \text{für } \underline{x} \text{ nah an EV.}$$

Direkte Potenzmethode.

Ziel: finde das betragsgrösste  $\in \mathbb{K}$  von  $\underline{A}$   
und ein  $\in V$  dazu.

Annahme  $\underline{A}$  diagonalisierbar

$$\underline{S}^{-1} \underline{A} \underline{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$$

$\in V$  stehen in den Spalten von  $\underline{S}$ ,  $\|\underline{s}_j\|_2 = 1$

Annahme:

$$|\lambda_1| > |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$$

$$\forall \underline{x} = \sum_{j=1}^n a_j \underline{s}_j \quad \text{mit } a_1 \neq 0, \text{ sonst beliebig}$$

$$\underline{A} \underline{x} = \sum_{j=1}^n a_j \underline{A} \underline{s}_j = \sum_{j=1}^n a_j \lambda_j \underline{s}_j$$

$$\underline{A}^2 \underline{x} = \sum_{j=1}^n a_j \lambda_j^2 \underline{s}_j \quad \dots$$

$$\underline{A}^k \underline{x} = \sum_{j=1}^n a_j \lambda_j^k \underline{s}_j = \lambda_1^k \left( a_1 \underline{s}_1 + \sum_{j=2}^n \left( \frac{\lambda_j}{\lambda_1} \right)^k a_j \underline{s}_j \right)$$

$$\left| \frac{\lambda_j}{\lambda_1} \right| < 1 \Rightarrow \left| \frac{\lambda_j}{\lambda_1} \right|^k \ll 1 \rightarrow \downarrow$$

für grosses  $k$  zeigt  $\underline{A}^k \underline{x}$  in der Richtung von  $\underline{s}_1$

$$\text{Somit } \frac{\underline{A}^k \underline{x}}{\|\underline{A}^k \underline{x}\|_2} \rightarrow \pm \underline{s}_1$$

$$\text{Notiere } \underline{x}_k = \underline{A}^k \underline{x} \Rightarrow \rho_{\underline{A}}(\underline{x}_k) = \frac{\underline{x}_k^H \underline{A} \underline{x}_k}{\underline{x}_k^H \underline{x}_k} =$$

$$= \frac{1}{\underline{x}_k^H \underline{x}_k} \left( \underline{x}_k^H \sum_{j=1}^n a_j \lambda_j^{k+1} \underline{s}_j \right) = \lambda_1 + O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$$

Potenzmethode. wähle  $\underline{x}_0$  zufällig,  $\|\underline{x}_0\|=1$

$$\left[ \begin{array}{l} \text{für } k=1, 2, \dots \\ \underline{w} = \underline{A} \underline{x}_{k-1} \\ \underline{x}_k = \frac{\underline{w}}{\|\underline{w}\|} \\ \lambda = \underline{x}_k^H \underline{A} \underline{x}_k \end{array} \right.$$

Theorem Potenzmethode liefert eine Iteration,  
die linear konvergiert gegen  $\lambda_1$   
mit der Rate  $\left| \frac{\lambda_2}{\lambda_1} \right|$

Ben Falls  $\underline{A}$  normal  $\Rightarrow$  EV orthogonal  $\Rightarrow$   
Fehler  $O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right)$

$\Rightarrow$  quadratische Konvergenz!

Ben Das war die erste Idee für Page Rank - Alg.  
von Google!

Ben  $\underline{A} = \underline{A}^H$  OAB von EV.

$$\underline{x} = \sum_{j=1}^n c_j \underline{u}_j = \underline{U} \underline{c}$$

$$\frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}} = \frac{\underline{c}^H \overbrace{\underline{U}^H \underline{A} \underline{U}} \underline{c}}{\underline{c}^H \underline{U}^H \underline{U} \underline{c}} = \frac{\underline{c}^H \underline{\Lambda} \underline{c}}{\underline{c}^H \underline{c}} =$$

$$= \frac{\lambda_1 |c_1|^2 + \lambda_2 |c_2|^2 + \dots + \lambda_n |c_n|^2}{|c_1|^2 + \dots + |c_n|^2}$$

$$p_A(x^k) = \lambda_1 + c \cdot \left( \frac{\lambda_2}{\lambda_1} \right)^{2k} + \dots$$



Ziel: finde das betragskleinste  $\in \mathbb{W}$ :

Annahme:  $\underline{A}$  invertierbar:

$$\underline{A} \underline{x} = \lambda \underline{x} \Rightarrow \underline{x} = \lambda \underline{A}^{-1} \underline{x}$$

$$\underline{A}^{-1} \mid \quad \frac{1}{\lambda} \underline{x} = \underline{A}^{-1} \underline{x}$$

$\Rightarrow$  betragskleinste  $\in \mathbb{W}$ :  $\lambda_n$ :  $\frac{1}{\lambda_n} \underline{x} = \underline{A}^{-1} \underline{x}$

$\frac{1}{\lambda_n}$  betragsgrösste  $\in \mathbb{W}$  von  $\underline{A}^{-1}$  und

diesben  $\in \mathbb{V}$ !

Darum "inverse Potenzmethode" = Potenzmethode für  $\underline{A}^{-1}$ .

Bez Wir berechnen nicht  $\underline{A}^{-1}$  sondern nur eine LU-Zerlegung von  $\underline{A}$  (Strukturerhaltend) nur ein Mal

Dann Löse LGS mit Matrizen  $\underline{L}, \underline{U}$

um  $\underline{A}^{-1} \underline{x}$  zu implementieren (Alg. in Skript)

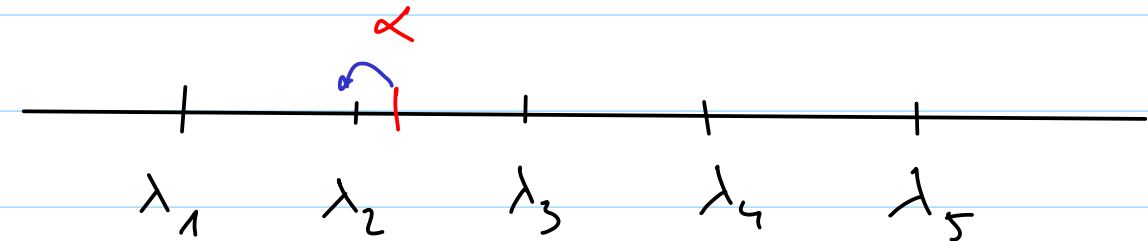
Ziel: Gegeben  $\alpha \in \mathbb{C}$ , finde  $\in \mathbb{W}$  nah an  $\alpha$ :  
 $|\alpha - \lambda| = \min \{ |\alpha - \mu|; \mu \in \mathbb{W} \text{ von } \underline{A} \}$

$$\underline{A} \underline{x} = \lambda \underline{x} \Leftrightarrow \underline{A} \underline{x} - \alpha \underline{I} \underline{x} = (\lambda - \alpha) \underline{x} \Leftrightarrow$$

$$(\underline{A} - \alpha \underline{I}) \underline{x} = (\lambda - \alpha) \underline{x} \Leftrightarrow \frac{1}{\lambda - \alpha} \underline{x} = (\underline{A} - \alpha \underline{I})^{-1} \underline{x}$$

$\Rightarrow$  Iteration mit  $(\underline{A} - \alpha \underline{I})^{-1} \Rightarrow \frac{1}{\lambda - \alpha} \Rightarrow \lambda$

"shifted inverse Iteration"



Bem Die Potenzmethode <sup>ist</sup> schneller wenn  $\alpha \approx \lambda_j$ .

Darum die Idee: wähle  $\alpha$  adaptiv, z.B.

$$\alpha = \rho_A(\underline{x}^{(k-1)}) \quad \text{im } k\text{-ten Schritt}$$

$\Rightarrow$  beschleunigte Konvergenz

Rayleigh-Quotienten-Iteration

(Skript Code & Bsp)

Bem Die Potenzmethoden brauchen einen guten Start-Vektor

Für RQI bekommt man einen guten Start-Vektor mit einige Schritte von shifted inverse Iteration.

$\Rightarrow$  Konvergenzordnung 3 !!!

Für  $\underline{A} = \underline{A}^H$  ;  $\underline{x} = \sum_{j=1}^n c_j \underline{u}_j = \underline{U} \underline{c}$ ,  $u_1, \dots, u_n \in V$

$$\rho_{\underline{A}}(\underline{x}) = \frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}} = \frac{\lambda_1 |c_1|^2 + \dots + \lambda_n |c_n|^2}{\|\underline{x}\|_2^2} \Rightarrow$$

$$\lambda_1 \leq \rho_{\underline{A}}(\underline{x}) \leq \lambda_n \quad \text{falls} \quad \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

$$\rho_{\underline{A}}(\underline{x}) \in [\lambda_1, \lambda_n]$$

$$\lambda_1 = \min_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x}) \quad \text{erreicht für } \underline{c} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\lambda_n = \max_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x}) \quad \underline{c} = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

Bem.  $\underline{x} \perp \underline{u}_1 \Rightarrow c_1 = 0$ ,  $\underline{x} = c_2 \underline{u}_2 + \dots + c_n \underline{u}_n$

$$\rho_{\underline{A}}(\underline{x}) \geq \lambda_2 \quad \text{erreicht für} \quad \underline{c} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\lambda_2 = \min_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

$$\underline{x} \in \mathbb{C}^n$$

$$\underline{x} \perp \underline{u}_1$$

$$\lambda_{n-1} = \max_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

und so weiter

$$\underline{x} \in \mathbb{C}^n$$

$$\underline{x} \perp \underline{u}_n$$

$$\lambda_k = \min_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

$$= \min_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

$$\underline{x} \perp \text{span}\{\underline{u}_1, \dots, \underline{u}_{k-1}\} \quad \underline{x} \in \text{span}\{\underline{u}_k, \underline{u}_{k+1}, \dots, \underline{u}_n\}$$

$$= \max_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

$$= \max_{\underline{x} \in \mathbb{C}^n} \rho_{\underline{A}}(\underline{x})$$

$$\underline{x} \perp \{\underline{u}_{k+1}, \dots, \underline{u}_n\}$$

$$\underline{x} \in \text{span}\{\underline{u}_1, \dots, \underline{u}_k\}$$

Theorem [Courant-Fischer]

$$\lambda_k = \min_{\dim U = k} \max_{x \in U} f_A(x) = \max_{\dim U = n-k+1} \min_{x \in U} f_A(x)$$

Konsequenz

$$\underline{Q}_m \in \mathbb{C}^{n \times m} \text{ mit } \underline{Q}_m^H \underline{Q}_m = \underline{I}$$

erweitere  $\underline{Q}_m$  auf ONB in  $\mathbb{C}^n$

$$\begin{bmatrix} \underline{Q}_m \\ \hat{\underline{Q}}_m \end{bmatrix} = \underline{Q}$$
Theorem [Cauchy]

$$\underline{H} \in \mathbb{C}^{m \times m} \text{ mit EW } \theta_1, \dots, \theta_m$$

$$\underline{A} = \begin{bmatrix} \underline{H} & \underline{B}^H \\ \underline{B} & \underline{R} \end{bmatrix} \quad m < n$$

Dann  $\lambda_k \leq \theta_k \leq \lambda_{k+n-m}$

$$\underline{Q}^H \underline{A} \underline{Q} = \begin{bmatrix} \underline{Q}_m^H \underline{A} \underline{Q}_m & \underline{Q}_m^H \underline{A} \hat{\underline{Q}}_m \\ \hat{\underline{Q}}_m^H \underline{A} \underline{Q}_m & \hat{\underline{Q}}_m^H \underline{A} \hat{\underline{Q}}_m \end{bmatrix}$$

hat dieselbe EW wie  $\underline{A}$

Idee Für  $n$  gross und  $m$  klein, wähle  $\underline{Q}_m$

so dass Bild  $\underline{Q}_m \simeq \text{span}\{ \underbrace{u_1, \dots, u_m}_{\leftarrow \text{EW von } \underline{A}} \}$

$\Rightarrow$  gute Approximation  $\theta_k \simeq \lambda_k$  für  $k=1, \dots, m$

Wie? modifiziertes Gram-Schmidt für  
 $\{ \underline{v}, \underline{A}\underline{v}, \underline{A}^2\underline{v}, \dots, \underline{A}^{m-1}\underline{v} \}$

Krylov-Verfahren (Arnoldi, Lanczos)

### § 9.3 Krylov-Verfahren

für "kleine" Matrizen: "eig" (QR-Alg.) gut

für "grosse" Matrizen ist eig zu langsam,  
 ausserdem QR-Alg. zerstört die Struktur von  $\underline{A}$

Krylov-Verfahren: grosse, dünn besetzte Matrizen.

Def Sei  $0 \neq \underline{z} \in \mathbb{C}^n$ ,  $\underline{A} \in \mathbb{C}^{n \times n}$

$\mathcal{K}_l(\underline{A}, \underline{z}) := \text{span} \{ \underline{z}, \underline{A}\underline{z}, \underline{A}^2\underline{z}, \dots, \underline{A}^{l-1}\underline{z} \} =$

$= \{ p(\underline{A})\underline{z} ; p = \text{Polynom von Grad} \leq l-1 \}$

Krylov-Raum

Suche eine ONB in  $\mathcal{K}_l(\underline{A}, \underline{z})$

$\mathcal{K}_1 = \text{span} \{ \underline{z} \} \subset \text{span} \{ \underline{z}, \underline{A}\underline{z} \} = \mathcal{K}_2 \subset \dots \subset \mathcal{K}_l$

Iterativ: gegeben  $\underline{v}_1, \dots, \underline{v}_l$  ONB in  $\mathcal{K}_l(\underline{A}, \underline{z})$

$\text{span} \{ \underline{v}_1, \dots, \underline{v}_j \} = \mathcal{K}_j(\underline{A}, \underline{z})$  für  $j = 1, 2, \dots, l$

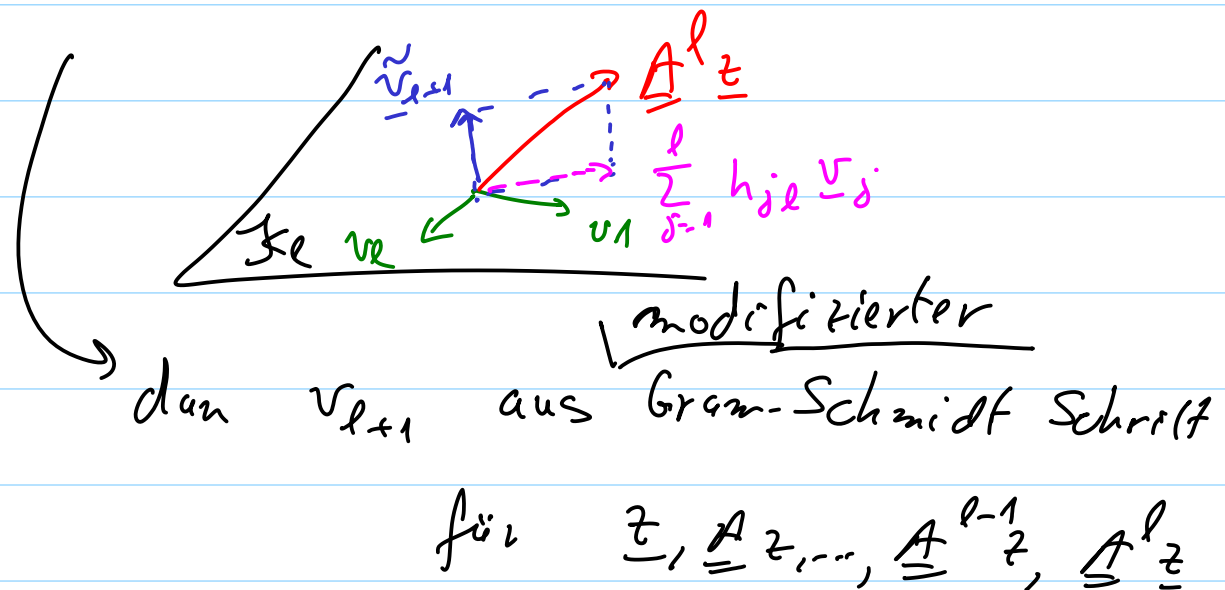
baue  $\underline{v}_1, \dots, \underline{v}_{l+1}$  ONB in  $\mathcal{K}_{l+1}(\underline{A}, \underline{z})$

$\mathcal{K}_{l+1}(\underline{A}, \underline{z}) = \text{span} \{ \underline{z}, \underline{A}\underline{z}, \dots, \underline{A}^{l-1}\underline{z}, \underline{A}^l\underline{z} \}$

entweder  $\underline{A}^l \underline{z} \in \mathcal{K}_l(\underline{A}, \underline{z})$ , d.h.  $\underline{A}^l \underline{z}$  lin. abhängig  
von  $\underline{z}, \dots, \underline{A}^{l-1} \underline{z}$

oder

$$\underline{A}^l \underline{z} \notin \mathcal{K}_l(\underline{A}, \underline{z}) \Rightarrow \underline{A}^l \underline{z} \in \mathcal{K}_{l+1} \setminus \mathcal{K}_l$$



$$\underline{\tilde{v}}_{l+1} = \underline{A}^l \underline{z} - \sum_{j=1}^l h_{j,l} \underline{v}_j$$

$$\underline{v}_l = \frac{\underline{\tilde{v}}_{l+1}}{\|\underline{\tilde{v}}_{l+1}\|}; \text{ dabei sind } h_{j,l} = \underline{v}_j^H \underline{A} \underline{v}_l$$

(ONB in  $\mathcal{K}_l$ )

## Algorithmus [Arnoldi Prozess]

$\underline{z}$  = beliebig.

$$\underline{v}_1 = \underline{z} / \|\underline{z}\|$$

für  $l = 1, 2, \dots, k-1$ :

$$\underline{\tilde{v}} = \underline{A} \underline{v}_l$$

für  $j = 1, 2, \dots, l$ :

$$h_{j,l} = \underline{v}_j^H \underline{\tilde{v}}$$

mod. GS

$$\underline{\tilde{v}} = \underline{\tilde{v}} - h_{j,l} \underline{v}_j$$

$$h_{l+1,l} = \|\underline{\tilde{v}}\|$$

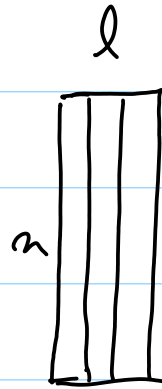
$$\underline{v}_{l+1} = \underline{\tilde{v}} / h_{l+1,l}$$

$l=1$	$l=2$	$l=3$
$h_{11}$	$h_{12}$	$h_{13}$
$h_{21}$	$h_{22}$	$h_{23}$
	$h_{32}$	$h_{33}$
		$h_{43}$

Ben Falls  $h_{l+1,l} = 0 \Rightarrow$  Abbruch der Iteration

$$\hookrightarrow \underline{A} \underline{v}_l \in \mathcal{K}_l(\underline{A}, \tau)$$

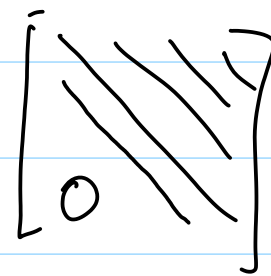
$$\underline{V}_l = [\underline{v}_1 \ \underline{v}_2 \ \dots \ \underline{v}_l] \in \mathbb{C}^{n \times l}$$



$$\tilde{\underline{H}}_l = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ 0 & h_{32} & h_{33} \\ 0 & 0 & h_{43} \end{bmatrix} \in \mathbb{C}^{l+1, l}$$

$$\underline{H}_l \in \mathbb{C}^{l, l}$$

$\underline{H}_l$  obere Hessenbergmatrix  
 $l \times l$



Aus Konstruktion:  $\underline{A} \underline{v}_k = h_{k+1,k} \underline{v}_{k+1} + \sum_{j=1}^k h_{jk} \underline{v}_j$

für  $k=1, 2, \dots, l$

das heisst:

$$\underline{A} \underline{V}_l = \underline{V}_{l+1} \tilde{\underline{H}}_l = \underline{V}_{l+1} \begin{bmatrix} \tilde{h}_{11} & \dots & \tilde{h}_{1,l+1} \\ \vdots & \ddots & \vdots \\ \tilde{h}_{l+1,1} & \dots & \tilde{h}_{l+1,l} \end{bmatrix} + \underline{V}_l \underline{H}_l$$

$$\begin{matrix} n & l & l+1 & l \\ \underline{A} & \underline{V}_l & \underline{V}_{l+1} & \tilde{\underline{H}}_l \end{matrix} =$$

Ben 1)  $\begin{bmatrix} \underline{V}_l^H \\ \underline{V}_l \end{bmatrix} = \begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix}$

$$\underline{V}_l^H \underline{V}_l = \underline{I}_l$$

2)  $\underline{V}_l^H \underline{A} \underline{V}_l = \underbrace{\underline{V}_l^H \underline{V}_{l+1}}_0 \overbrace{\begin{bmatrix} 0 & 0 & \dots & 0 \end{bmatrix}}^{h_{l+1,l}} + \underbrace{\underline{V}_l^H \underline{V}_l}_{\underline{I}_l} \underline{H}_l$

$\Rightarrow \underline{V}_l^H \underline{A} \underline{V}_l = \underline{H}_l$

$\begin{bmatrix} \underline{V}_l^H \end{bmatrix} \begin{bmatrix} \underline{A} \end{bmatrix} \begin{bmatrix} \underline{V}_l \end{bmatrix} = \begin{bmatrix} \underline{H}_l \end{bmatrix}$

3) falls  $h_{l+1,l} = 0 \Rightarrow \underline{K}_{l+1} = \underline{K}_l$  und

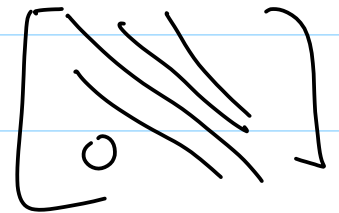
$\underline{A} \underline{V}_l = \underline{V}_l \underline{H}_l$

4) falls  $\underline{A}$  Hermite-symmetrisch:  $\underline{A}^H = \underline{A}$

$\underline{H}_l^H = (\underline{V}_l^H \underline{A} \underline{V}_l)^H = \underline{V}_l^H \underline{A}^H \underline{V}_l = \underline{V}_l^H \underline{A} \underline{V}_l =$

$= \underline{H}_l \Rightarrow \underline{H}_l$  auch Hermite-symmetrisch.

$\underline{H}_l$  Hessenberg



$\Rightarrow \underline{H}_l = \begin{bmatrix} & & 0 \\ & \diagdown & \\ 0 & & \end{bmatrix}$  tridiagonal!

$\beta^H$   $\alpha$   $\beta$

$\Rightarrow$  es reichen die Vektoren  $\alpha, \beta$  um  $\underline{H}_l$  zu speichern

und die innere Schleife im Arnoldi Prozess hat konstante Länge  $l=2$ .



$$\tilde{v}_{l+1} = \underline{A} \underline{v}_l - h_{l,l} \underline{v}_l - h_{l-1,l} \underline{v}_{l-1}$$

Langos-Verfahren  $O(nk)$

Arnoldi-Verfahren  $O(nk^2)$

Allgemeiner Name: Krylov-Raum-Verfahren.

Theorem Falls  $h_{l+1,l} = 0$  und  $h_{j+1,j} \neq 0$  für  $j=1,2,\dots,l-1$ , dann

(1) jeder EW von  $\underline{H}_l$  ist auch EW von  $\underline{A}$

(2) falls  $\underline{A}$  regulär, dann gibt es  $\underline{z} \in \mathbb{C}^l$  so dass

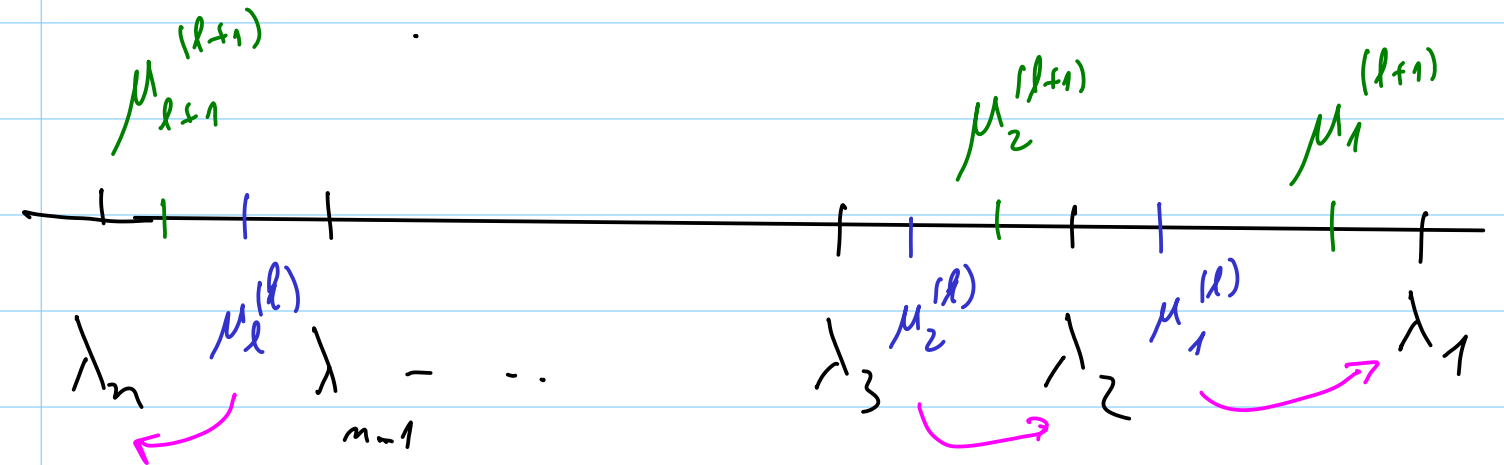
$$\underline{A} \underline{x} = \underline{b} \quad \text{mit} \quad \underline{x} = \underline{V}_l \underline{y}$$

**Theorem 7.4.12.** Seien  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  und  $\mu_1^{(\ell)} \geq \mu_2^{(\ell)} \geq \dots \geq \mu_\ell^{(\ell)}$  die Eigenwerte der Hermite-symmetrischen Matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$ , bzw. von  $\mathbf{H}_\ell = \mathbf{V}_\ell^H \mathbf{A} \mathbf{V}_\ell$  für  $\ell = 1, 2, \dots$ . Dann gelten für  $1 \leq j \leq \ell$  die Ungleichungsketten

$$\lambda_{n-j+1} \leq \mu_{\ell+1-j+1}^{(\ell+1)} \leq \mu_{\ell-j+1}^{(\ell)}$$

und

$$\mu_j^{(\ell)} \leq \mu_j^{(\ell+1)} \leq \lambda_j.$$



Skript → einfache Implementierung

ARPACK ← eigvals

Bsp im Skript.

Langos ~ "ghost EW" da Orthogonalität verloren!

⇒ re-orthogonalisiere

# § 10 Lineare Anfangswertprobleme

1. Fall

$$\begin{cases} \dot{\underline{y}} = \underline{A} \underline{y} \\ \underline{y}(0) = \underline{y}_0 \end{cases} \quad \underline{A} = \underline{S}^{-1} \underline{D} \underline{S}$$

$$\underline{\hat{y}} = \underline{S}^{-1} \underline{y} \Rightarrow \text{entkoppeln}$$

$$\begin{cases} \dot{\hat{y}}_1 = \lambda_1 \hat{y}_1 \\ \dots \\ \dot{\hat{y}}_d = \lambda_d \hat{y}_d \end{cases} \Rightarrow \hat{y}_i(t) = \left( \underline{S}^{-1} \underline{y}_0 \right)_i e^{\lambda_i t} \quad t \in \mathbb{R}$$

$$\underline{y}(t) = \underline{S} \begin{bmatrix} e^{\lambda_1 t} & & \\ & \ddots & \\ & & e^{\lambda_d t} \end{bmatrix} \underline{S}^{-1} \underline{y}_0$$

OK für kleine  $d$ , aber für  $d \geq 5$  instabil.

Für  $d=5, 6, \dots, 20$  besser expm

$\hookrightarrow$

$$\underline{y}(t) \approx \expm(\underline{A}t) \underline{y}_0$$

Padé-Approximation  $d \geq 21$  zulässig

$d$  gross,  $\underline{A}$  dünnbesetzt = Krylov-Verfahren

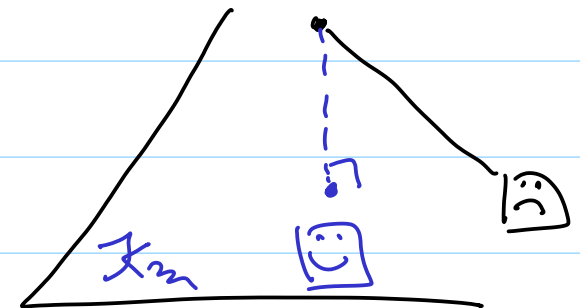
Für  $\underline{A}$ : es gibt  $\underline{V}_m \in \mathbb{C}^{d \times m}$  mit orthonormalen Spalten.

$$\underline{V}_m^H \underline{A} \underline{V}_m = \underline{H}_m \quad m \times m \text{-Matrix (klein)}$$

$\hookrightarrow$  Hessenberg  $m \ll d$ .

$$\begin{cases} \text{Finde } \underline{y}(t) \in \mathbb{R}^d \\ \dot{\underline{y}}(t) - \underline{A} \underline{y}(t) = 0 \\ \underline{y}(0) = \underline{y}_0 \end{cases} \quad \mathcal{K}_m(\underline{A}, \underline{y}_0) = \text{span} \{ \underline{y}_0, \underline{A} \underline{y}_0, \dots, \underline{A}^{m-1} \underline{y}_0 \}$$

$= \text{span} \{ \underline{v}_1, \dots, \underline{v}_m \}$   
ONB!



Darum ersetzen wir das Problem durch

Finde  $\underline{u}_m \in \mathcal{X}_m(\underline{A}, \underline{y}_0)$  so dass

Residuum  $(\dot{\underline{u}}_m(t) - \underline{A} \underline{u}_m(t)) \perp \mathcal{X}_m(\underline{A}, \underline{y}_0)$

$$\underline{u}_m(0) = \underline{y}_0$$

Einsetzen:

$$\langle \underline{w}, \sum_{k=1}^m \dot{c}_k(t) \underline{v}_k - \sum_{k=1}^m c_k(t) \underline{A} \underline{v}_k \rangle = 0$$

für alle  $\underline{w} \in \mathcal{X}_m(\underline{A}, \underline{y}_0)$



Finde  $\underline{u}_m \in \mathcal{X}_m(\underline{A}, \underline{y}_0)$

$$\langle \underline{w}, \dot{\underline{u}}_m(t) - \underline{A} \underline{u}_m(t) \rangle = 0 \text{ für alle } \underline{w} \in \mathcal{X}_m(\underline{A}, \underline{y}_0)$$

$$\underline{u}_m(0) = \underline{y}_0$$

Wähle  $\underline{w} = \underline{v}_1$ :

$$\langle \underline{v}_1, \sum_{k=1}^m \dot{c}_k(t) \underline{v}_k \rangle = \langle \underline{v}_1, \sum_{k=1}^m c_k(t) \underline{A} \underline{v}_k \rangle$$

$$\sum_{k=1}^m \dot{c}_k(t) \underbrace{\langle \underline{v}_1, \underline{v}_k \rangle}_{\delta_{1k}} = \sum_{k=1}^m c_k(t) \underbrace{\langle \underline{v}_1, \underline{A} \underline{v}_k \rangle}_{(H_m)_{1k}}$$

"  
 $\underline{v}_1^H \underline{v}_k = 0$  falls  $k \neq 1$  da ONB  
 $1$  falls  $k = 1$

$$\underline{v}_1^H \underline{A} \underline{v}_k = (H_m)_{1k}$$

$$\underline{u}_m(t) \in \mathcal{X}_m(\underline{A}, \underline{y}_0) = \text{span} \{ \underline{v}_1, \dots, \underline{v}_m \} \Rightarrow$$

$$\underline{u}_m(t) = \sum_{k=1}^m c_k(t) \underline{v}_k ; \quad \underline{c}(t) = \begin{bmatrix} c_1(t) \\ \vdots \\ c_m(t) \end{bmatrix} \in \mathbb{C}^m$$

$$\dot{c}_1(t) = \sum_{k=1}^m (\underline{H}_m)_{1k} c_k(t)$$

Für  $w = \underline{v}_2, \underline{v}_3, \dots, \underline{v}_m$  analog  $\Rightarrow$

$$\begin{cases} \dot{\underline{c}}(t) = \underline{H}_m \underline{c}(t) \\ \underline{y}(0) = \underline{y}_0 \Rightarrow \underline{c}(0) = \|\underline{y}_0\| \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^m \end{cases} \quad \begin{array}{l} m \text{ klein!} \Rightarrow \text{Pade'} \\ \text{geht} \\ \text{schnell} \end{array}$$

$$\underline{c}(t) = \|\underline{y}_0\| e^{\underline{H}_m t} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$\Rightarrow$  Lösung in  $\mathcal{X}_m(\underline{A}, \underline{y}_0)$

$$\underline{u}_m(t) = \sum_{k=1}^m c_k(t) \underline{v}_k = \|\underline{y}_0\| \underline{V}_m e^{\underline{H}_m t} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

2. Fall  $\dot{\underline{y}}(t) = \underline{A} \underline{y}(t) + g(t)$  inhomogene Fall.

Variation der Konstanten

$$\underline{y}(t) = e^{\underline{A}(t-t_0)} \underline{y}_0 + \int_{t_0}^t e^{\underline{A}(t-s)} g(s) ds$$

3. Fall  $\dot{\underline{y}}(t) = \underline{A}(t) \underline{y}(t)$ : Magnus-Integratoren

Unter bestimmte Voraussetzungen:

$$\underline{y}(t) = e^{\underline{\Omega}(t, t_0)} \underline{y}_0 \quad \text{Magnus Entwicklung}$$

$$\text{mit } \underline{\Omega} = \sum_{k=1}^{\infty} \underline{\Omega}_k ; \quad \text{Kommutator } [A, B] = AB - BA$$

$$\Omega_1 = \int_{t_0}^t \mathbf{A}(\tau_1) d\tau_1,$$

$$\Omega_2 = \frac{1}{2} \int_{t_0}^t \int_{t_0}^{\tau_1} [\mathbf{A}(\tau_1), \mathbf{A}(\tau_2)] d\tau_2 d\tau_1,$$

$$\Omega_3 = \frac{1}{12} \int_{t_0}^t \int_{t_0}^{\tau_1} \int_{t_0}^{\tau_2} [[\mathbf{A}(\tau_1), \mathbf{A}(\tau_2)], \mathbf{A}(\tau_3)] + [\mathbf{A}(\tau_1), [\mathbf{A}(\tau_2), \mathbf{A}(\tau_3)]] d\tau_3 d\tau_2 d\tau_1.$$

Idee: Statt  $\dot{\underline{y}} = \underline{A}(t) \underline{y}$  löse  $\dot{\underline{z}} = \underline{\hat{A}}(t) \underline{z}$

wobei  $\underline{\hat{A}}(t) = \sum_{i=1}^p l_i(t) \underline{A}(t_n + c_i \cdot h)$

$l_i(t)$  = Lagrange Polynom in  $t_n + c_i \cdot h$

$c_i$  = Quadraturknoten zwischen  $[0, 1]$

$$t_n + c_i \cdot h \in [t_n, t_{n+1}]$$

$\Rightarrow \underline{\hat{A}} =$  Polynom vom Grad  $s$  in  $t$

$$\underline{\hat{A}}(t_n + c_i \cdot h) = \underline{A}(t_n + c_i \cdot h) \quad \text{für } i=1, 2, \dots, p.$$

Dann Magnus-Entwicklung für

$$\dot{\underline{z}} = \underline{\hat{A}}(t) \underline{z}$$

Vorteil: Integration nur von Polynome in  $t \Rightarrow$  einfach und exakt berechnen.

$\underline{\hat{A}}(t)$  glatt; man kann zeigen,

Rest in der Magnus-Entwicklung ( $\underline{\hat{\Omega}}$ ) ist

$O(t^5)$  nach 4 Terme.

Theorem  $(b_i, c_i)_{i=1,\dots,p}$  Quadraturformel der  
Ordnung  $p \geq 1$   
 $y(t_n+h) - \hat{y}(t_n+h) = \mathcal{O}(h^{p+1})$

**Beispiel 2.8.1** (Magnus-Verfahren 2. Ordnung). Die Mittelpunktsregel auf  $[0, 1]$  ergibt den Schritt von  $t_n$  zu  $t_{n+1}$  als:

$$\mathbf{y}_{n+1} = e^{\mathbf{\Omega}^{[2]}} \mathbf{y}_n \text{ mit } \mathbf{\Omega}^{[2]} = h\mathbf{A}(t_n + \frac{h}{2}),$$

mit dem lokalen Fehler  $\mathcal{O}(h^{2+1})$  und dem globalen Fehler  $\mathcal{O}(h^2)$ .

**Beispiel 2.8.2** (Magnus-Verfahren 4. Ordnung). Die Gauss-Quadratur mit 2 Punkten auf  $[0, 1]$  ergibt den Schritt von  $t_n$  zu  $t_{n+1}$  als:

$$\mathbf{y}_{n+1} = e^{\mathbf{\Omega}^{[4]}} \mathbf{y}_n \text{ mit } \mathbf{\Omega}^{[4]} = \frac{h}{2}(\mathbf{A}_1 + \mathbf{A}_2) - h^2 \frac{\sqrt{3}}{12} [\mathbf{A}_1, \mathbf{A}_2]$$

wobei

$$\begin{aligned} \mathbf{A}_1 &= \mathbf{A} \left( t_n + \left( \frac{1}{2} - \frac{\sqrt{3}}{6} \right) h \right), \\ \mathbf{A}_2 &= \mathbf{A} \left( t_n + \left( \frac{1}{2} + \frac{\sqrt{3}}{6} \right) h \right). \end{aligned}$$

Der lokale Fehler ist  $\mathcal{O}(h^{4+1})$  und der globale Fehler ist  $\mathcal{O}(h^4)$ .

**Beispiel 2.8.3** (Magnus-Verfahren 4. Ordnung - Simpson-Regel). Die Simpson-Regel mit 3 Punkten auf  $[0, 1]$  ergibt den Schritt von  $t_n$  zu  $t_{n+1}$  als:

$$\mathbf{y}_{n+1} = e^{\mathbf{\Omega}^{[4]}} \mathbf{y}_n \text{ mit } \mathbf{\Omega}^{[4]} = \frac{h}{6}(\mathbf{A}_1 + 4\mathbf{A}_2 + \mathbf{A}_3) - h^2 \frac{1}{72} [\mathbf{A}_1 + 4\mathbf{A}_2 + \mathbf{A}_3, \mathbf{A}_3 - \mathbf{A}_1]$$

wobei

$$\begin{aligned} \mathbf{A}_1 &= \mathbf{A}(t_n), \\ \mathbf{A}_2 &= \mathbf{A} \left( t_n + \frac{h}{2} \right), \\ \mathbf{A}_3 &= \mathbf{A}(t_n + h). \end{aligned}$$

Der lokale Fehler ist  $\mathcal{O}(h^{4+1})$  und der globale Fehler ist  $\mathcal{O}(h^4)$ . Eine andere Methode derselben Ordnung ergibt sich mit

$$\mathbf{\Omega}^{[4]} = \frac{h}{6}(\mathbf{A}_1 + 4\mathbf{A}_2 + \mathbf{A}_3) - h^2 \frac{1}{12} [\mathbf{A}_1, \mathbf{A}_3]$$

Ordnung 6  
 ohne Kommutatoren  
 Bsp. } Skript

# §11 Exponentielle Integratoren

$$\begin{cases} \dot{\underline{y}} = \underline{f}(\underline{y}) & \text{mit } \underline{f} \text{ stetig differenzierbar} \\ \underline{y}(0) = \underline{y}_0 & \text{autonom} \end{cases}$$

Idee der Linearisierung:  $\underline{J} := \underline{Df}(\underline{y}_0)$

$$\dot{\underline{y}} = \underbrace{\underline{\partial} \underline{y}}_{\text{linear}} + \underbrace{\underline{f}(\underline{y}) - \underline{J} \underline{y}}_{\underline{g}(\underline{y})} \quad \dot{\underline{y}} = \underline{\partial} \underline{y} + \underline{g}(\underline{y})$$

Variation der Konstanten:

$$\underline{y}(h) = e^{\underline{\partial} h} \underline{y}_0 + \int_0^h e^{\underline{\partial}(h-s)} \underline{g}(\underline{y}(s)) ds$$

Für das Integral ersetze die Unbekannte  $\underline{y}(s)$  durch  $\underline{y}_0 \Rightarrow$

$$\int_0^h e^{\underline{\partial}(h-s)} \underline{g}(\underline{y}(s)) ds \approx \int_0^h e^{\underline{\partial}(h-s)} \underline{g}(\underline{y}_0) ds =$$

$$= h \ell(h \underline{\partial}) \underline{g}(\underline{y}_0)$$

$$\ell(z) \stackrel{\text{def}}{=} \frac{e^z - 1}{z} = \sum_{n=1}^{\infty} \frac{1}{n!} z^{n-1} = \sum_{n=0}^{\infty} \frac{z^n}{(n+1)!}$$

Beweis

$$\int_0^h e^{\underline{\partial}(h-s)} \underline{g}(\underline{y}_0) ds = \int_0^h \sum_{n=0}^{\infty} \frac{1}{n!} \left( \underline{\partial}(h-s) \right)^n \underline{g}(\underline{y}_0) ds =$$

Umtauschen  $\int_0^h (h-s)^n ds$

$$= \sum_{n=0}^{\infty} \frac{1}{n!} \underline{\partial}^n (-1) \frac{(h-s)^{n+1}}{n+1} \Big|_{s=0}^{s=h} \underline{g}(\underline{y}_0) =$$

$$= \sum_{n=0}^{\infty} \frac{h}{(n+1)!} \underline{\partial}^n h^n \underline{g}(\underline{y}_0) = h \ell(h \underline{\partial}) \underline{g}(\underline{y}_0)$$





Verallgemeinerung: exponentielle RK-Verfahren.

$$\underline{\partial} = \underline{D}f(\underline{y})$$

semi-implizite Euler:  $\underline{y}_1 = \underline{y}_0 + \left( \underline{I} - h \underline{\partial} \right)^{-1} h f(\underline{y}_0)$

exponentielle Euler:  $\underline{y}_1 = \underline{y}_0 + \ell(h \underline{\partial}) h f(\underline{y}_0)$

Idee: ersetze  $\frac{1}{1-z}$  durch  $\ell(z) = \frac{e^z - 1}{z}$

$$\Rightarrow \begin{cases} \underline{k}_i = \ell(\alpha h \underline{\partial}) \left( f(\underline{u}_i) + h \underline{\partial} \sum_{j=1}^{i-1} \alpha_{ij} \underline{k}_j \right) \\ \underline{u}_i = \underline{y}_0 + h \sum_{j=1}^{i-1} \alpha_{ij} \underline{k}_j \\ \underline{y}_1 = \underline{y}_0 + h \sum_{i=1}^n b_i \underline{u}_i \end{cases}$$

explizite RK:  $\ell(z) = 1, \quad \underline{\partial} = \underline{0}$   
 RoW  $\ell(z) = \frac{1}{1-z}$

exp. RK:  $\ell(z) = \frac{e^z - 1}{z}$

ODEn.	nicht steif	steif	oszillierend.
Methode	expl. RK; ode45	impl. RK;	exponentielle RK
Stabilität	$h < \frac{1}{\lambda}$	alle $h$	alle $h$
Implementierung	$f(\underline{y})$	$\underline{D}f(\underline{y}_0)$ lösen nicht-lin. Gleichungssystem	$\underline{D}f(\underline{y}_0)$ $\ell(\underline{A}) \underline{b}$ Schnell mit Krylov. <u>exp 4</u>

⊕ Erhaltungseigenschaften wichtig?

⊕ Autonom?

⊕ Ordnung der Dgl.?

# Die Methode des Konjugierten Gradienten

$\underline{A} \in \mathbb{R}^{n \times n}$  symmetrisch, pos. definit.

Ziel  $\underline{A}\underline{x} = \underline{b}$  :  $\underline{x}$  geschickt approximiert.

Definiere "Energie"  $Q(\underline{x}) = \frac{1}{2} \underline{x}^T \underline{A} \underline{x} - \underline{x}^T \underline{b}$

$$\left( \underline{A}\underline{x} = \underline{b} \Leftrightarrow \underline{A}\underline{x} - \underline{b} = \underline{0} \Leftrightarrow \underline{x}^T \underline{A}\underline{x} - \underline{x}^T \underline{b} = 0. \right)$$

Statt Näherung an Lösung von  $\underline{A}\underline{x} = \underline{b}$   
 suche Näherung an Lösung von  $\underset{\underline{x} \in \mathbb{R}^n}{\operatorname{argmin}} Q(\underline{x})$

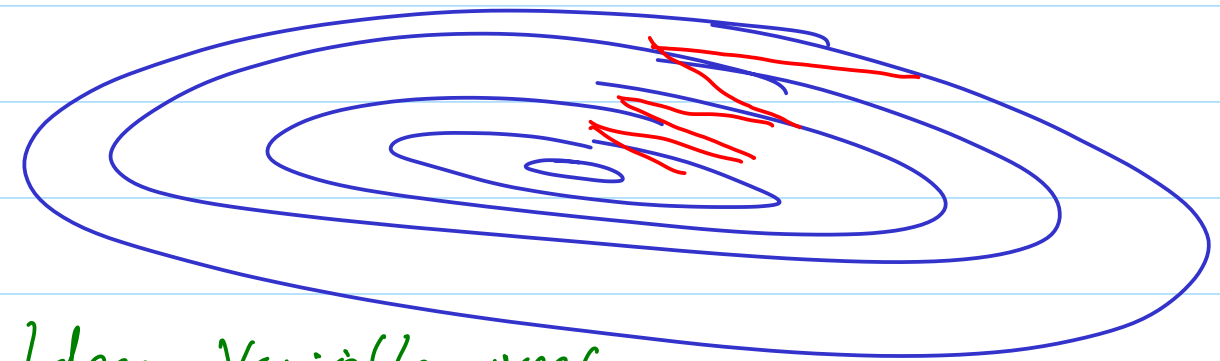
1) Wähle Suchrichtung  $\underline{p}_k \neq \underline{0}$

2) optimiere in dieser Richtung:  $Q(\underline{x}_k + \alpha_k \underline{p}_k) = \min!$

$$\alpha_k = \frac{\underline{p}_k^T \underline{r}_k}{\underline{p}_k^T \underline{A} \underline{p}_k} ; \quad \underline{r}_k = \underline{b} - \underline{A} \underline{x}_k \text{ Residuum.}$$

Wie wählt man die Suchrichtung?

Analysis: Gradienten-Richtung.

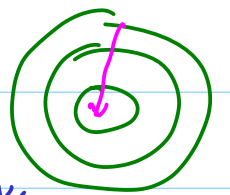


Bessere Idee: Variablenwechsel

$$\underline{x} = \underline{P} \underline{y} \Rightarrow Q(\underline{x}) = \frac{1}{2} \underline{y}^T \underline{D} \underline{y} - (\underline{P}^T \underline{b})^T \underline{y}$$

so dass  $\underline{P}^T \underline{A} \underline{P} = \underline{D} = \text{diagonal.}$

↓  
 das ist nicht eine  $\in \mathbb{W}$ -zerlegung!  
 da  $\underline{P}$  nicht orthogonal!



Stiefel: wie man effizient die Richtungen baut:

$$\underline{x}_0 \text{ gegeben, } \underline{p}_0 = \underline{r}_0 = \underline{b} - \underline{A}\underline{x}_0$$

für  $k=0, 1, 2, \dots, n-1$

$$\alpha_k = \frac{\underline{r}_k^T \underline{p}_k}{\underline{p}_k^T \underline{A} \underline{p}_k}$$

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{p}_k$$

$$\underline{r}_{k+1} = \underline{r}_k - \alpha_k \underline{A} \underline{p}_k$$

$$\beta_k = - \frac{\underline{p}_k^T \underline{A} \underline{r}_{k+1}}{\underline{p}_k^T \underline{A} \underline{p}_k}$$

$$\underline{p}_{k+1} = \underline{r}_{k+1} + \beta_k \underline{p}_k$$

effizienter:

$$\underline{x}_0 \text{ gegeben, } \underline{p}_0 = \underline{r}_0 = \underline{b} - \underline{A}\underline{x}_0$$

$$or = \|\underline{r}_0\|^2$$

für  $k=0, 1, 2, \dots$

$$\underline{w}_k = \underline{A} \underline{p}_k$$

$$\alpha_k = \frac{or}{\underline{p}_k^T \underline{w}_k}$$

$$\underline{x}_{k+1} = \underline{x}_k + \alpha_k \underline{p}_k$$

$$\underline{r}_{k+1} = \underline{r}_k - \alpha_k \underline{w}_k$$

$$r_0 = \|\underline{r}_{k+1}\|^2$$

$$\beta_k = \frac{r_0}{or} \quad \left( \text{also } \frac{\underline{r}_{k+1}^T \underline{r}_k}{\underline{r}_k^T \underline{r}_k} \right)$$

$$\underline{p}_{k+1} = \underline{r}_{k+1} + \beta_k \underline{p}_k$$

CG = conjugate gradient

Theorem 1)  $r_l \in \mathcal{K}_{l+1}(\underline{A}, \underline{r}_0) = \mathcal{K}_{l+1}(\underline{A}, \underline{r}_0)$

2) CG optimal:  $\|\underline{x} - \underline{x}_{l+1}\| = \min_{\tilde{\underline{x}} \in \mathcal{J}_k} \|\underline{x} - \tilde{\underline{x}}\|_{\underline{A}}$

mit  $\mathcal{J}_k = \underline{x}_0 + \mathcal{K}_k(\underline{A}, \underline{r}_0)$

$$\|\underline{x}\|_{\underline{A}}^2 = \underline{x}^T \underline{A} \underline{x}$$

3) Falls  $\underline{A}$  nur  $m \leq n$  verschiedene EW hat,  
dann konvergiert CG in maximal  $m$  Schritten

$$4) \underline{e}_k = \underline{x} - \underline{x}_k \Rightarrow \|\underline{e}_k\|_{\underline{A}} \leq 2 \left( \frac{\sqrt{\text{cond}(\underline{A})} - 1}{\sqrt{\text{cond}(\underline{A})} + 1} \right)^k \|\underline{e}_0\|_{\underline{A}}$$

Begründung der Wahl der Richtung.

Idee: Suche  $\underline{x}_k$  die optimale Approximation in  $\underline{x}_0 + \mathcal{K}_k(\underline{A}, \underline{r}_0)$

$$\underline{r}_k = \underline{b} - \underline{A} \underline{x}_k \perp \mathcal{K}_k(\underline{A}, \underline{r}_0)$$

$$0 = \underline{v}_k^T \underline{r}_k = \underline{v}_k^T (\underline{b} - \underline{A} \underline{x}_k) = \underline{v}_k^T (\underline{b} - \underbrace{\underline{A} \underline{x}_0}_{-\underline{A} \underline{x}_0} - \underline{A} \underline{v}_k \underline{y})$$

$$\underline{x}_k = \underline{x}_0 + \underline{v}_k \underline{y}$$

$$= \underbrace{\underline{v}_k^T (\underline{b} - \underline{A} \underline{x}_0)}_{r_0} - \underbrace{\underline{v}_k^T \underline{A} \underline{v}_k}_{H_k} \underline{y} \Rightarrow$$

$$H_k \underline{y} = \underline{v}_k^T \underline{r}_0 = \|\underline{r}_0\|_{\underline{e}_1}$$

$$H_k = \begin{bmatrix} \parallel \\ 0 \end{bmatrix}$$

Lanczos macht die Iteration kurz!

Bem Bei schlecht konditionierte Matrizen ist CG noch zu langsam!

\* CG in Prinzip  $O(n)$  aber wegen Rundungsfehler  $O(n \log n)$

\* beschleunige das Verfahren indem man die Matrix verändert: vorkonditioniert!

$$\underline{A} \underline{x} = \underline{b} \quad \underline{\hat{A}} \underline{\hat{x}} = \underline{\hat{b}}$$

$$\underline{\hat{A}} = \underline{S} \underline{A} \underline{S}^T \quad \underline{\hat{x}} := \underline{S}^{-T} \underline{x} \quad \underline{\hat{b}} = \underline{S} \underline{b}$$

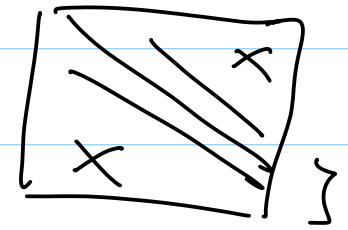
nehme  $\underline{S}$  so dass 1)  $\text{cond}(\underline{\hat{A}}) < \text{cond}(\underline{A})$

2)  $\underline{S} \underline{\hat{x}} = \underline{c}$  soll einfach sein.

v.CG : CG mit  $\underline{S}^{-1} \underline{A}$  statt  $\underline{A}$  !!

typischerweise:  $\text{diag}(\underline{A})$

oder



Bem Noch bessere Methoden:

incomplete LU decomposition,  
 \* algebraisch  $\begin{cases} \text{sparse approx. inverse} \\ \text{alg. multigrid.} \end{cases} \quad O(n)$

\* geometrisch: das Physik. des Problems verwendet.  
 $\begin{cases} \text{domain decomposition method.} \\ \text{geometric multigrid.} \end{cases} \quad O(n)$