

FS19

§1 Quadratur

§1.1. Motivation

Bsp
$$\begin{cases} \dot{y} = f(t, y) & t \in \mathbb{R} \quad y(t) \in \mathbb{R} \\ y(t_0) = y_0 \end{cases}$$

Gesucht: $y: [t_0, T] \rightarrow \mathbb{R}$

Bescheidener: $y(T)$ ist gesucht

$$\int_{t_0}^T dt \Rightarrow y(T) = y(t_0) + \int_{t_0}^T f(t, y(t)) dt$$

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 Planck's Gesetz für die Strahlung eines schwarzen Körpers

(C) V. Gradinaru

T = Temperatur in K, λ = Wellenlänge in μm
 c_1, c_2 = phys. konst.

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

↳ Energieverlust pro Flächeneinheit.

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

↙
 σ = Stefan-Boltzmann-Konstante.

Ziel: berechne σ

$$\sigma T^4 = \int_0^\infty E(\lambda, T) d\lambda \Rightarrow \sigma = \frac{1}{T^4} \int_0^\infty E(\lambda, T) d\lambda$$

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

Gewichte Knoten

Ziel: Wähle Knoten, Gewichte, Strategie für diese Berechnung um den Fehler $|J - Q|$ aber auch die Kosten klein zu halten.

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

z.B.

$$f \approx \text{Polynom} \quad \alpha_0 + \alpha_1 x + \dots + \alpha_n x^n$$

$f \approx$ trigonometrischen Polynom.

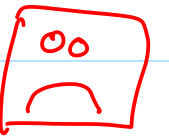
$$\alpha_0 + \alpha_1 e^{i1x} + \alpha_2 e^{i2x} + \dots \\ + \alpha_{-1} e^{-i1x} + \alpha_{-2} e^{-i2x} + \dots$$

Gegeben Knoten x_0, x_1, \dots, x_n
kann man $\alpha_0, \alpha_1, \dots, \alpha_n$ berechnen, sodass
 $p_n(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_n x^n$ erfüllt

$$p_n(x_j) = f(x_j) \quad \text{für } j=0, 1, \dots, n$$

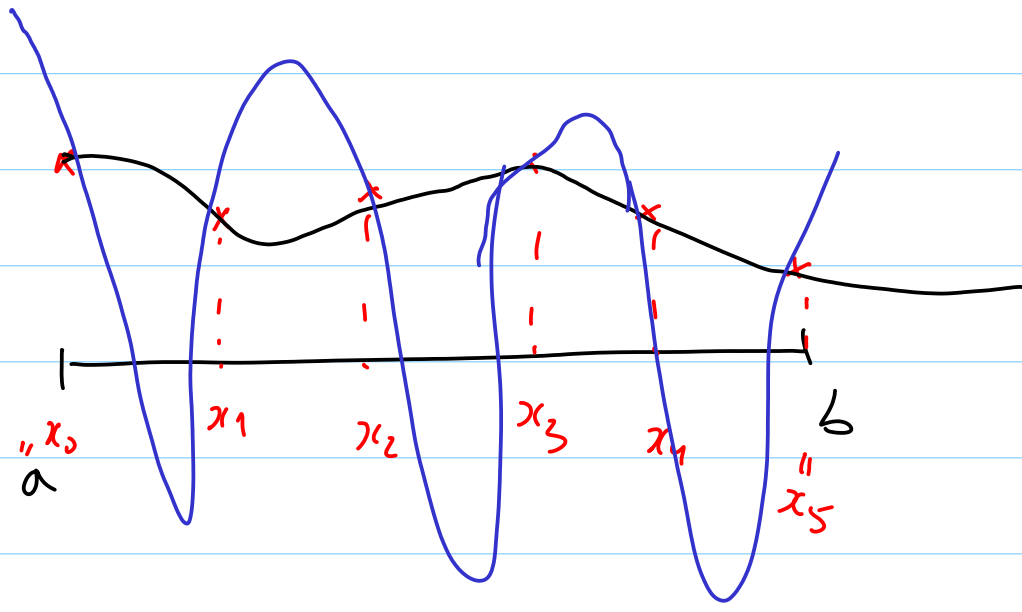
Im Prinzip ist das das LGS:

$$\begin{bmatrix} 1 & x_0 & \dots & x_0^n \\ 1 & x_1 & & x_1^n \\ 1 & x_2 & & x_2^n \\ \vdots & \vdots & & \vdots \\ 1 & x_n & & x_n^n \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix} = \begin{bmatrix} f(x_0) \\ \vdots \\ f(x_n) \end{bmatrix}$$

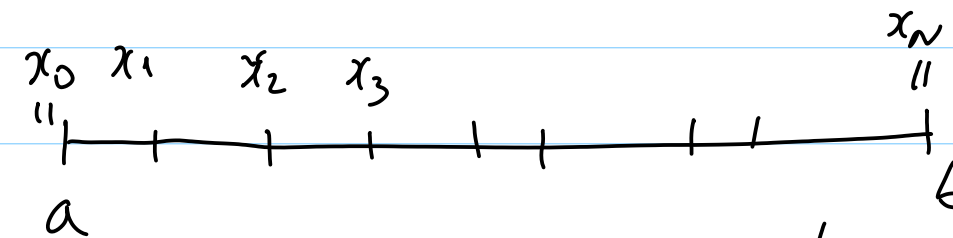


Somit $f(x) \approx p_n(x)$ und dann.

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



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



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

Man kann beweisen:

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

Länge des Intervalls ist wichtig.

Glattheit ist wichtig.

und wende die Quadraturformel auf jedem kleinen Intervall der Länge $h \Rightarrow$

zusammengesetzte Quadraturformel

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

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

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

$$\left. \begin{aligned} \sum_{k=0}^{N-1} \frac{h^{n+1}}{n!} \max_{z \in [x_k, x_{k+1}]} |f^{(n)}(z)| &\leq C \cdot \frac{h^{n+1}}{n!} \sum_{k=0}^{N-1} 1 = C \frac{h^{n+1}}{n!} N \\ &\leq \max_{z \in [a, b]} |f^{(n)}(z)| = C \\ h = \frac{b-a}{N} \Rightarrow N = \frac{b-a}{h} \end{aligned} \right\} \Rightarrow$$

$$\Rightarrow \underline{|J - Q| \leq C \cdot \frac{h^{n+1}}{n!} \cdot \frac{b-a}{h} = C \frac{h^n}{n!} (b-a)}$$

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

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

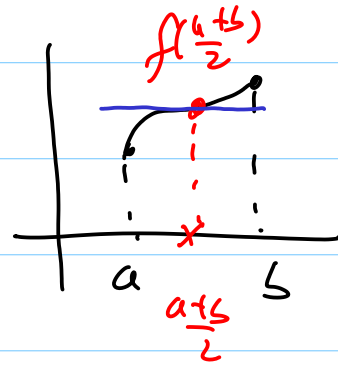
$$\int_a^b p(x) dx = Q(p, a, b) \quad \text{für alle Polynome vom Grad } \leq n$$

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.

$$(MPR) \quad Q^1(f, a, b) = (b-a) f\left(\frac{a+b}{2}\right).$$



Ben 1) Polynome vom Grad 0 und 1 werden exakt integriert

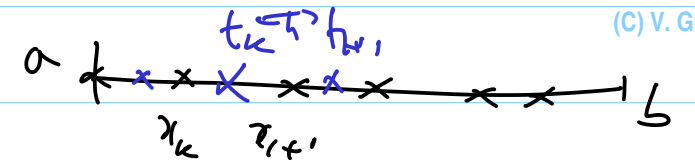
\Rightarrow MPR hat Quadraturordnung 2

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

Ben 2) offene Quadraturformel: Enden $[a, b]$ sind keine Knoten.

MPR:

$$\int_a^b f(x) dx \approx \sum_{k=0}^{n-1} (x_{k+1} - x_k) f\left(\frac{x_k + x_{k+1}}{2}\right)$$



$$t_k = \frac{x_k + x_{k+1}}{2} = x_k + \frac{h}{2} \text{ mit } h = x_{k+1} - x_k = \frac{b-a}{n}$$

$$\int_a^b f(x) dx \approx \frac{b-a}{n} \sum_{k=0}^{n-1} f(t_k)$$

Implementierung.

gegeben $a, b;$

wähle $n;$

$$h = \frac{b-a}{n}; \quad t_k = a + \frac{h}{2}$$

für $k=0, 1, 2, \dots, n-1:$
 $S = S + f(t_k)$

$$t_k = t_k + h$$

$$S = S \cdot h$$

return S

Gegeben: f, a, b

Wähle: N

$$h = (b-a)/N$$

$$t = a + h/2$$

$$S = 0$$

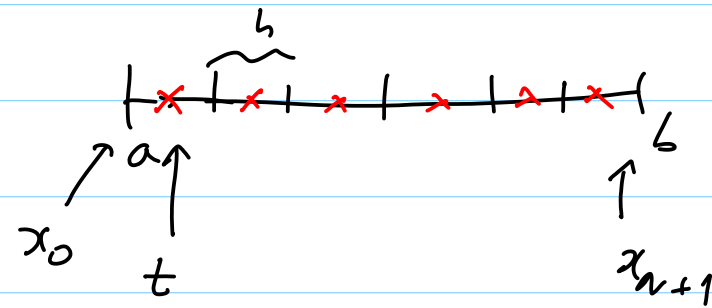
für $k=0, 1, 2, \dots, N-1$:

$$S = S + f(t)$$

$$t = t + h$$

$$S = h \cdot S$$

return S



$$\sum_{k=0}^{N-1} f(t_k)$$

$$x[:] = [x[0], x[1], \dots, x[N]]$$

$$x[:-1] = [x[0], x[1], \dots, x[N-1]]$$

Python: * numpy \rightarrow array $x[-2] = [\dots, x[N-2]]$

* scipy

* matplotlib / ...

$$x[2:-2] =$$

$$= [x[2], \dots, x[N-2]]$$

$$t = \text{linspace}(a + \frac{h}{2}, b, N-1)$$

oder:

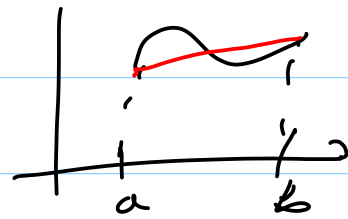
$$x = \text{linspace}(a, b, N+1)$$

$$t = \frac{h}{2} + x[:-1]; \quad ft = f(t); \quad \text{h.mean}(ft)$$

Bsp 2) Trapezregel

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

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



$$= \frac{h}{2} f(x_0) + \boxed{\frac{h}{2} f(x_1)} + \boxed{\frac{h}{2} f(x_2)} + \dots + \boxed{\frac{h}{2} f(x_{n-1})} + \frac{h}{2} f(x_n)$$

$h f(x_1)$ $h f(x_2)$ $h f(x_{n-1})$

Bem 1) Ordnung 2; Fehler $\frac{1}{12} (b-a)^3 f^{(2)}(\tau)$
mit $\tau \in [a, b]$.

2) geschlossene Quadraturformel
Enden von $[a, b]$ sind Knoten.

$$\text{TR: } \int_a^b f(x) dx = \sum_{k=1}^N \int_{t_{k-1}}^{t_k} f(x) dx \approx$$

$$= \sum_{k=1}^N \left(\frac{h}{2} f(x_{k-1}) + \frac{h}{2} f(x_k) \right) =$$

$$= \frac{h}{2} f(x_0) + \boxed{h \sum_{k=1}^{N-1} f(x_k)} + \frac{h}{2} f(x_N)$$

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

$$h \sum_{k=1}^{N-1} f(x[k]) + \frac{h}{2} (f(x[0]) + f(x[N]))$$

Bsp 3) Simpson Regel
 (Polynom von Grad ~~3~~ wird exakt integriert)
 \Rightarrow Ordnung 4

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

Fehler: $\frac{1}{90} \left(\frac{b-a}{2}\right)^5 f^{(4)}(\xi)$ mit $\xi \in (a, b)$

Aufgabe: MPR, TR, Simpson. implementieren.

Anwender:

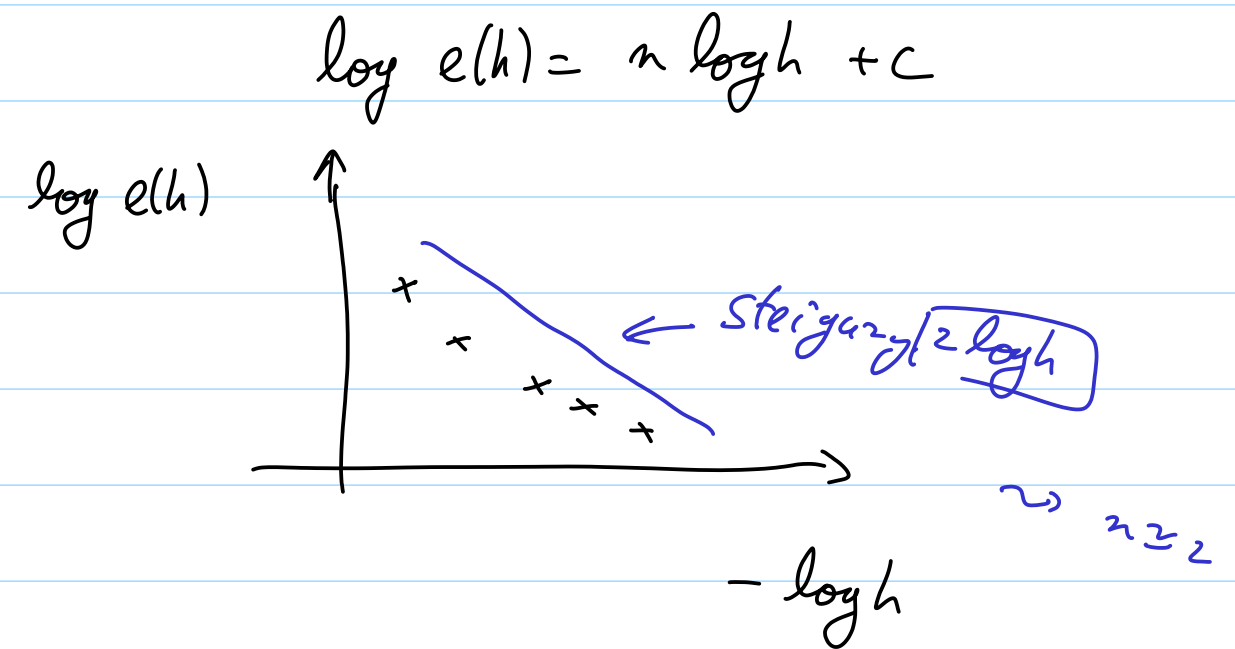
$$\sin(x) \quad \frac{1}{1+(5x)^2}; \quad \sqrt{x}$$

Stefan-Boltzmann Konstante
 (Planck's Gesetz)

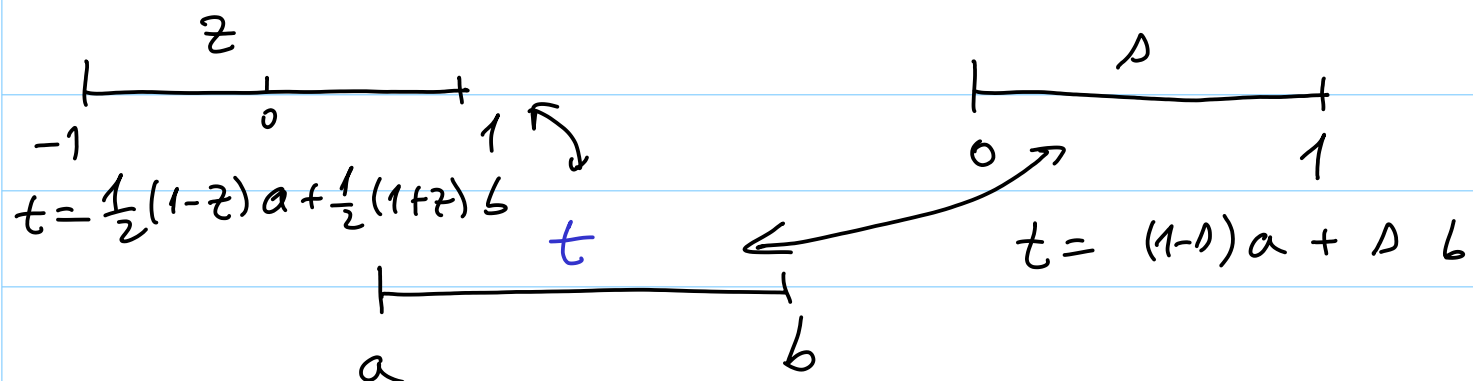
```
from numpy import linspace, array, sum, ...
from scipy import integrate
```

`integrate.quad(f, a, b)`

Glettheit \Rightarrow Fehler $c \cdot h^n = e(h)$



§1.2. Referenzintervalle und Symmetrische Quadraturformel

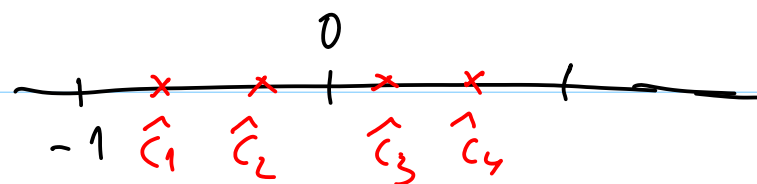


$$\int_a^b f(t) dt = \frac{b-a}{2} \int_{-1}^1 \hat{f}(z) dz \approx \frac{b-a}{2} \sum_{j=1}^n \hat{w}_j \hat{f}(\hat{c}_j)$$

\uparrow
[−1,1]

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

Definition QF auf $[-1,1]$ heißt symmetrisch, falls

$$\hat{c}_k = -\hat{c}_{n+1-k}, \quad \hat{w}_k = \hat{w}_{n+1-k}$$


Theorem Die Quadraturordnung einer symmetrischen QF ist gerade.

Beweis

Annahme: QF exakt für Polynome vom Grad $2m-2$.

Nehme $f(x) = ax^{2m-1}$

$$\int_{-1}^1 f(x) dx = a \int_{-1}^1 x^{2m-1} dx = 0$$

Die QF ist exakt ^{für ax^{2m-1}} falls $Q(f, -1, 1) = 0$

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

$$= -a \sum_{k=1}^n \hat{w}_{n+1-k} \hat{c}_{n+1-k}^{2m-1} =$$

$$= -a \sum_{j=1}^n \hat{w}_j \hat{c}_j^{2m-1} = -Q(f, -1, 1)$$

$$\Rightarrow Q(f, -1, 1) = 0 = \int_{-1}^1 f(x) dx \quad \text{qed.}$$

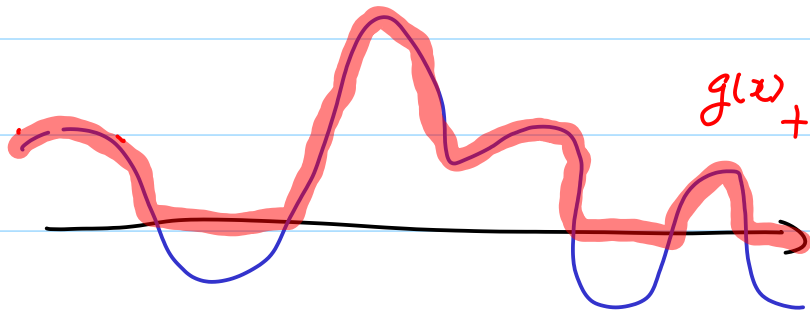
§1.3. Fehler für Quadratur auf $[0,1]$

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

\downarrow Knoten in $[0,1]$
 \downarrow Gewichte

E linear in g . $E(Ag + Bf) = A E(g) + B E(f)$
 $A, B \in \mathbb{R}$ und g, f Funktionen

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



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

für festes z : $(t-z)_+$

für festes t : $(t-z)_+$

als Funktion von der Variable t

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

Theorem

Sei Q eine QF mit Quadraturordnung n
 und sei g n -mal stetig differenzierbar.

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

Beweis Taylor um Punkt 0:

$$g(t) = g(0) + \dots + \frac{t^{n-1}}{(n-1)!} g^{(n-1)}(0) + \underbrace{\int_0^t \frac{(t-z)^{n-1}}{(n-1)!} g^{(n)}(z) dz}_{\text{q(t) Polynom vom Grad n-1}}$$

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

Linearität von E :

$$E(g) = \underbrace{E(q)}_0 + \int_0^1 \underbrace{E(\alpha(z, \cdot))}_{k_n(z)} g^{(n)}(z) dz$$

Bemerkung Wenden wir den Satz auf

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

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

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

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

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

da:

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

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

$$g''(t) = f''(x_0 + ht) h^2$$

...

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

Fehler auf $[x_0, x_0+h]$ ist

$$h^{\textcircled{n+1}} \int_0^1 |k_n(\tau)| f^{(n)}(x_0 + h\tau) d\tau$$

Zusammengesetzt:

$$|E(f)| = \left| \int_a^b f(x) dx - \sum_{k=1}^N \sum_{j=1}^2 b_j \cdot f(x_{k-1} + c_j h) \right| \leq$$

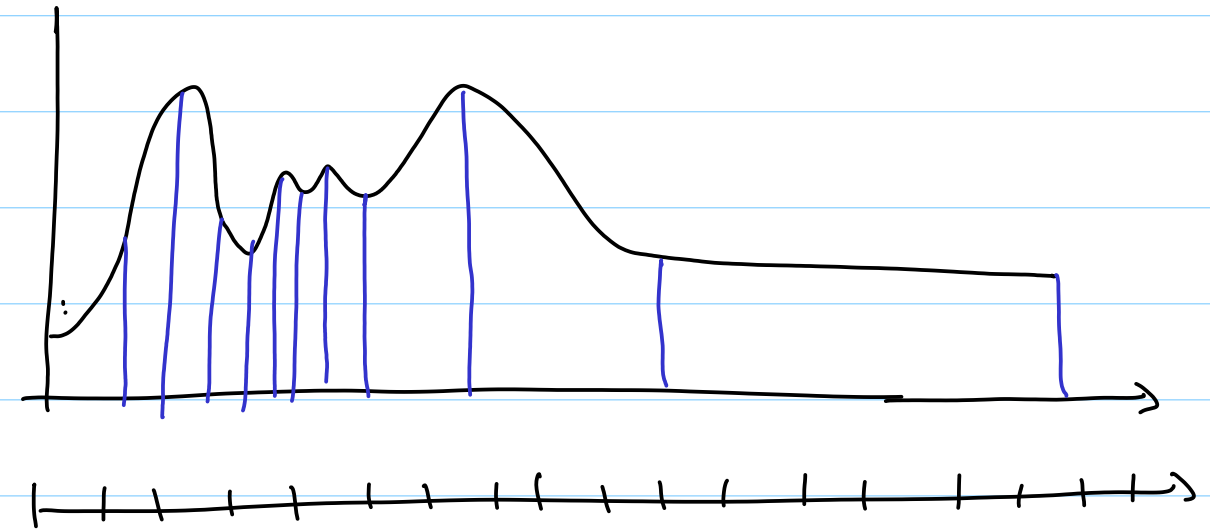
$$\leq c \cdot h^{\textcircled{2}} \max_{x \in [a,b]} |f^{(n)}(x)| \rightarrow 0 \text{ für } h \rightarrow 0$$

\swarrow $x \in [a,b]$ \searrow Glattheit
 \swarrow Ordnung der lokalen QF
 \swarrow kommt aus dem Peano Kern

$$\int_0^1 |k_n(\tau)| d\tau = \text{konstante.}$$

MPR, TR = $\frac{1}{12}$ Simpson $\frac{1}{2880}$

§1.4. Adaptive Quadratur



Optimiere: Anzahl Funktionsauswertungen

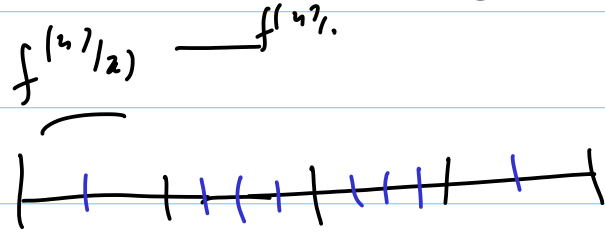
Lokalen Intervall $[x_{k-1}, x_k]$ hat lokalen Fehler ε_k

$$\varepsilon_k = \left| \int_{x_{k-1}}^{x_k} f(x) dx - \sum_{j=1}^2 b_j \cdot f(x_{k-1} + c_j h_k) \right|$$

$$h_k = x_k - x_{k-1} \text{ lokale Intervalllänge}$$

Wissen $\varepsilon_k \sim h^2 \max_{x \in [x_{k-1}, x_k]} |f^{(n)}(x)|$

IDEA: wähle nur dort ein kleines h , wo $|f^{(n)}(x)|$ gross
wie?



Wie schätze ich während der Rechnung
den lokale Fehler ε_k , ohne weitere
Informationen über f ?

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

$$\approx \left| Q^S(f, x_{k-1}, x_k) - Q^T(f, x_{k-1}, x_k) \right| = \tilde{\varepsilon}_k$$

Schätzung des lokalen Fehlers.

(möchte möglichst gleiche lokale Fehler haben)? Verfeinere die Intervalle, wo $\tilde{\varepsilon}_k$ gross ist.

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

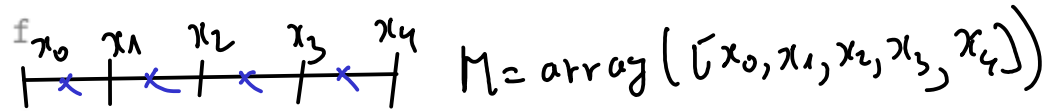
Verwendet wird allerdings das Wert, das
 Q^S gibt

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

```

def adaptquad(f,M,rtol,abstol):
    """
    adaptive quadrature using trapezoid and simpson rules
    Arguments:
    f      handle to function f
    M      initial mesh
    rtol   relative tolerance for termination
    abstol absolute tolerance for termination, necessary in case the exact
    integral value = 0, which renders a relative tolerance meaningless.
    """
    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
    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
    # 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 intervals with large error contributions, recurse.
    return I

```

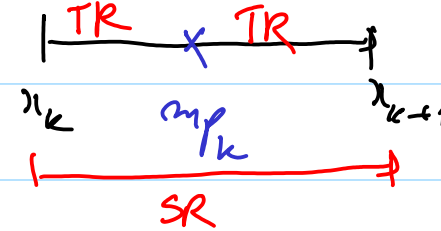


$$h_k = x_{k+1} - x_k \Rightarrow h = [h_0, h_1, h_2, h_3] \text{ array}$$

$$\frac{1}{2}(x_k + x_{k+1}) \Rightarrow [mp_0, mp_1, mp_2, mp_3] \text{ array.}$$

$fx = \text{array}$ der Länge 5

$fm = \text{array}$ der Länge 4



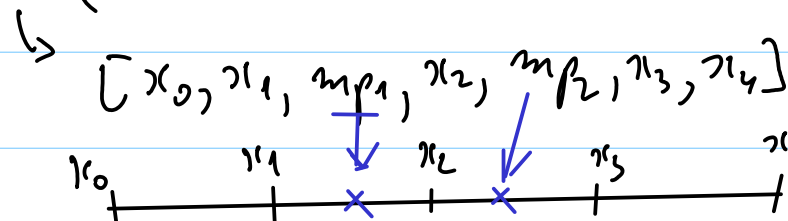
Mittleren Fehler

$$Loc > \frac{9}{10} \cdot \frac{tot}{\# \text{ Intervale}}$$

array. Zahl $\Rightarrow [False, True, True, False]$

$$[0, 1, 1, 0]$$

$$\text{Sort}([x_0, x_1, x_2, x_3, x_4, mp_1, mp_2])$$



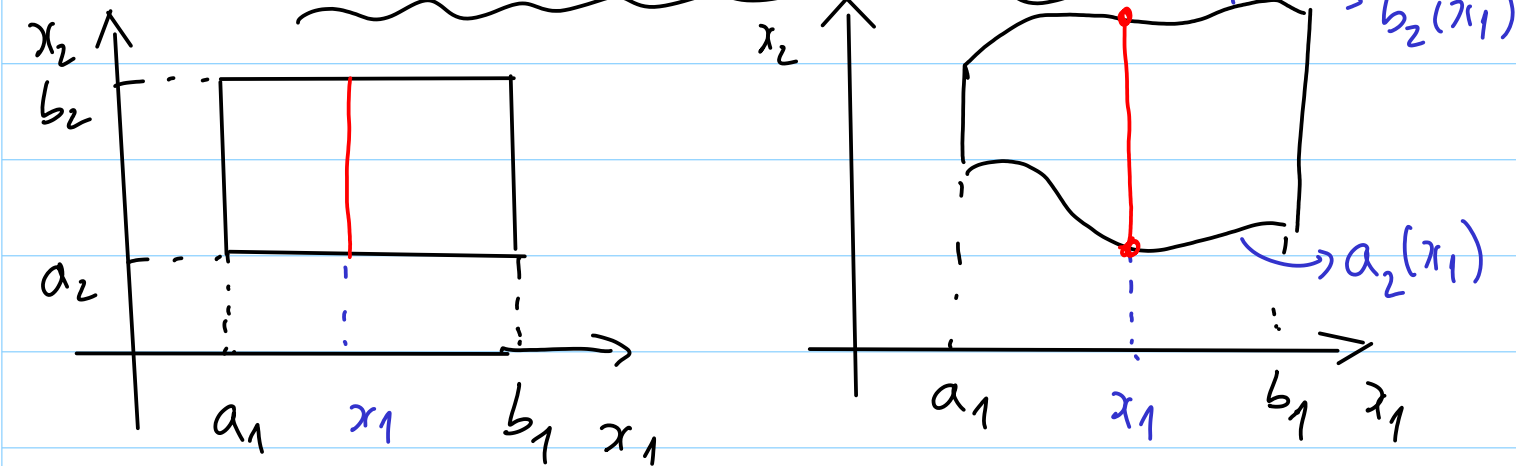
neues Gitter

nonzero \rightarrow welche Indices sind $\neq 0 \Rightarrow 1, 2$

\rightarrow Lese help für nonzero um [0] zu verstehen!

\rightarrow Rekursion!

§1.5. Quadratur in \mathbb{R}^d $d=2,3$



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

$$= \int_{a_1}^{b_1} F(x_1) dx_1 = \sum_{k_1=1}^{N_1} \int_{x_1^{k_1-1}}^{x_1^{k_1}} F(x_1) dx_1 \approx$$

The diagram shows the interval $[a_1, b_1]$ on the x_1 axis, partitioned into subintervals by points $x_1^0, x_1^1, x_1^2, x_1^3, x_1^4$. The points are labeled $x_1^0, x_1^1, x_1^2, x_1^3, x_1^4$, and the interval is labeled a_1 and b_1 .


$$= \frac{b_1 - a_1}{N_1} \sum_{k_1=1}^{N_1} \sum_{j_1=1}^{D_1} \underbrace{F(x_1^{k_1-1} + h_1 c_{j_1}^1)}_{\substack{\uparrow \\ \text{Knoten in } D_{x_1}\text{-Richtung} \\ \text{auf dem Unterintervall } [x_1^{k_1-1}, x_1^{k_1}]}} \cdot w_{j_1}^1$$

$$= \frac{b_1 - a_1}{N_1} \sum_{k_1=1}^{N_1} \sum_{j_1=1}^{D_1} \int_{a_2(x_1^{k_1-1} + h_1 c_{j_1}^1)}^{b_2(x_1^{k_1-1} + h_1 c_{j_1}^1)} f(x_2, x_1^{k_1-1} + h_1 c_{j_1}^1) \underbrace{dx_2}_{\substack{\uparrow \\ \psi F}} \cdot w_{j_1}^1$$

$k_1 = 1, 2, \dots, N_1$

$$= \frac{b_1 - a_1}{N_1} \sum_{k_1=1}^{N_1} \sum_{j_1=1}^{D_1} \frac{b_2 - a_2}{N_2} \sum_{k_2=1}^{N_2} \sum_{j_2=1}^{D_2} f(x_2^{k_2-1} + h_2 c_{j_2}^2, x_1^{k_1-1} + h_1 c_{j_1}^1) \cdot \underbrace{w_{j_2}^2 w_{j_1}^1}_{\substack{\uparrow \\ \text{Knoten in } D_{x_2}\text{-Richtung} \\ \text{auf dem Unterintervall } [x_2^{k_2-1}, x_2^{k_2}]}}$$

$$= \frac{(b_1 - a_1)(b_2 - a_2)}{N_1 N_2} \sum_{k_1=1}^{N_1} \sum_{\substack{j_1=1 \\ j_2=1}}^{N_2} f\left(x_2^{k_2-1} + h_2 c_{j_2}^2, x_1^{k_1-1} + h_1 c_{j_1}^1\right) w_{j_2}^2 w_{j_1}^1$$

in d -Dimensionen $N_1 \cdot N_2 \dots N_d$ 
Auswertungen von f .

$$h_1 h_2 \dots h_d \sum_{k_1} \sum_{k_2} \dots \sum_{k_d} f(\dots) w_{j_d}^d w_{j_{d-1}}^{d-1} \dots w_{j_1}^1$$

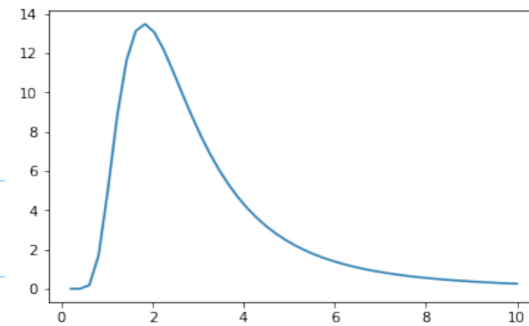
Für grosse d :

- * f glatt \Rightarrow dünne Gitter / sparse grids
- * Monte Carlo ; quasi-Monte Carlo
- * hoch oszillatorische Funktionen.

```
In [1]: from numpy import linspace, exp, inf
        from pylab import plot, show
```

```
In [2]: # Plank's law for blackbody radiation
c1 = 3.7413*10**8 #W \mu m^4/m^2
c2 = 1.4388*10**4 # \mu m K
T = 1000. # T = temperature in K
# lam = wavelength in um
# monochromatic emissive power of a blackbody:
def Eblam(lam, T=1000.):
    return c1/(lam**5*( exp(c2/(lam*T)) -1.))
```

```
In [3]: print('----- plot -----')
y = linspace(0,10)
plot(y, Eblam(y,1600.)/10**4)
show()
```



```
In [4]: print('----- integrate: quad -----')
        from scipy import integrate
        # total energy lost per unit area
        Q, abserr, info = integrate.quad(Eblam,0,inf, full_output=True)
        neval = info['neval']
        # Q = sigma * T^4; sigma = Stefan-Boltzmann constant
        print(Q/T**4, abserr, 'neval=', neval)
```

```
----- integrate: quad -----
5.669294910232968e-08 9.879197389039973e-07 neval= 105
```

```
In [8]: def ga(x):
        return exp(-x**2)

        val, err, info = integrate.quad(ga,-1,1, full_output=True)
        print(val)
        print(err)
        print(info['neval'])
```

```
1.493648265624854
1.6582826951881447e-14
21
```

```
In [9]: def gaabc(x,a,b,c):
        return a*exp(-((x-b)/c)**2)
```

```
val, err, info = integrate.quad(gaabc,-1,1,args=(2,0.5,4), full_output=True)
print(val)
print(err)
print(info['neval'])
```

```
val, err, info = integrate.quad(gaabc,-inf,inf, args=(2,0.5,4), full_output=True)
print(val)
print(err)
print(info['neval'])
```

```
3.8599302795534305
4.2853834698080357e-14
21
14.17963080724413
2.5011577807838655e-08
330
```

```

from scipy.special import jv
f0 = lambda x: jv(0,x)
f1 = lambda x: jv(1,x)
val, err, info = integrate.quad(f0,0,5, full_output=True)
print(val)
print(err)
print(info['neval'])
val, err, info = integrate.quad(f1,0,5, full_output=True)
print(val)
print(err)
print(info['neval'])

```

```

0.7153119177847678
2.47260738289741e-14
21
1.177596771314338
1.8083362065765924e-14
21

```

```

In [11]: from numpy import sqrt
        f = lambda x: 1/sqrt(abs(x))
        integrate.quad(f,-1,1)

```

```
Out[11]: (inf, inf)
```

```
In [12]: integrate.quad(f,-1,1, points=[0]) # can deal with singularity if you tell it
```

```
Out[12]: (3.9999999999999813, 5.684341886080802e-14)
```

```

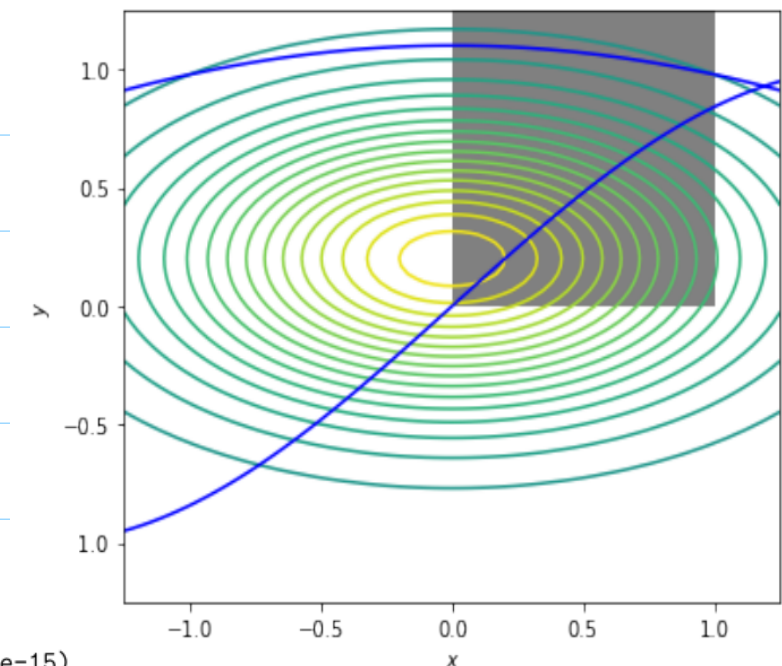
In [13]: from pylab import subplots
        import matplotlib.patches as patches
        from numpy import *
        def f(y,x):
            'y must be the first argument, and x the second.'
            return exp(-x**2 - 3*(y-0.2)**2)

fig, ax = subplots(figsize=(6,5))
x = y = linspace(-1.25,1.25,75)
X, Y = meshgrid(x,y)
c = ax.contour(X,Y,f(Y,X), 16, vmin=-1, vmax=1)
bound = patches.Rectangle((0,0),1,1.25, facecolor='grey')

ax.add_patch(bound)
ax.set_xlabel('$x$')
ax.set_ylabel('$y$')

ay = lambda x: sin(x)
by = lambda x: 0.1+ cos(0.5*x)
plot(x,ay(x),'-b')
plot(x,by(x),'-b')
show()

```



```
In [14]: integrate.dblquad(f,0,1,ay,by)
```

```
Out[14]: (0.24918585141988464, 7.637755883693233e-15)
```


§1.6. Quadratur mit erhöhter Ordnung.

Λ Knoten auf Referenzintervall
 Λ Gewichte
] so bestimmt,

dass Polynome höchstes Grades exakt mit dieser QF integriert werden.

Wie hoch kann so eine Quadraturordnung sein?

2 Λ Unbekannte (Knoten, Gewichte);

2 Λ Gleichungen:

$$\begin{cases}
 \int_0^1 1 dt = Q_\Lambda(1, 0, 1) \\
 \int_0^1 t dt = Q_\Lambda(t, 0, 1) \\
 \dots \left(\frac{1}{3} = b_1 c_1^2 + b_2 c_2^2 + \dots \right) \\
 \int_0^1 t^{2\Lambda-1} dt = Q_\Lambda(t^{2\Lambda-1}, 0, 1)
 \end{cases}$$



Da QF
exakt

Möchte Ordnung $p = \Lambda + m$

Jedes Polynom vom Grad $\leq p-1 = \Lambda+m-1$
soll mit QF exakt integriert werden.

Sei f Polynom vom Grad $\leq p-1 = \Lambda+m-1$

IDEA: teile das Polynom f durch

$$M(x) = (x - c_1)(x - c_2) \dots (x - c_\Lambda)$$

$$\text{Grad}(M) = \Lambda$$

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

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

← da QF exakt

→ ||

$$\sum_{j=1}^{\Lambda} b_j f(c_j) = \sum_{j=1}^{\Lambda} b_j \underbrace{M(c_j)}_{=0} g(c_j) + \sum_{j=1}^{\Lambda} b_j r(c_j) \Rightarrow$$

$$\Rightarrow \int_0^1 M(t)g(t)dt = 0 \quad \text{für alle Polynome } g \text{ mit } \text{Grad}(g) \leq n-1.$$

$f = \text{Polynom vom Grad} \leq n+m-1$

$$\langle M, g \rangle = \int_0^1 M(t)g(t)dt \quad \text{Skalarprodukt im Raum der Polynome}$$

Theorem

Ordnung der QF ist $n+m$ \Leftrightarrow

$$\langle M, g \rangle = 0 \text{ für alle Polynome vom Grad } \leq n-1.$$

$$P_m = \text{span} \{1, t, t^2, \dots, t^{m-1}\} : M \perp P_m$$

Theorem Ordnung einer QF ist höchstens $2n$.

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

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

Nehme $g = M$

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

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

Widerspruch \Rightarrow Wahr: $p \leq 2n$.

Orthogonale Polynome

$w:]a, b[\rightarrow \mathbb{R}$ Gewichtfunktion

stetig, $w(x) > 0$ für alle $x \in]a, b[$

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

Betrachte den linearen Raum:

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

Bem Alle Polynome liegen in V

V : Skalarprodukt

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

Theorem [Gram-Schmidt in LA]

Es existiert eine eindeutige Folge von Polynomen

p_0, p_1, \dots mit

$$p_k(x) = x^k + \text{Polynom von Grad} \leq k-1$$

so dass

$$p_k \perp \text{span} \{p_0, p_1, \dots, p_{k-1}\}.$$

Diese Polynome baut man so: (3-Term Rekursion)

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

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: oft werden sie aber zu $p_k(1) = 1$ "normiert". Darum: $p_2(x) = \frac{3}{2}(x^2 - \frac{1}{3})$

Bem c_1, c_2, \dots, c_n die Nullstellen von (M) von P_n .

$$M = P_n$$

Bsp1) $w(x) \equiv 1$, $a=0, b=1$ $= \int_0^1 f(x)g(x)dx$

Gram-Schmidt \Rightarrow orthogonale Polynome

Legendre Polynome

QF: Gauss-Quadratur

z)

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

QF \Rightarrow Hermite-Quadratur.

Bsp Gauss-Quadratur

$[0, 1]$ Notiere Knoten $x_j, j=1, \dots, n$

$[-1, 1]$ Notiere Knoten $c_j, j=1, \dots, n$.

Schreibe $x_{p_{k-1}} = p_k + q$ mit $\text{Grad } q \leq k-1 \Rightarrow$

$$\Rightarrow \langle x p_k, p_{k-1} \rangle = \langle p_k, p_k \rangle + \langle p_k, 0 \rangle$$

Somit $p_{k+1} \perp p_{k-1} \Leftrightarrow \alpha_{k-1} = -\frac{\langle p_k, p_k \rangle}{\langle p_{k-1}, p_{k-1} \rangle} = -\gamma_{k+1}^2$

Für $j \leq k-2$:

$$p_{k+1} \perp p_j \Leftrightarrow 0 = \langle p_{k+1}, p_j \rangle = \langle \alpha_k p_k, p_j \rangle + \alpha_j \langle p_j, p_j \rangle =$$

$$= \langle p_k, x p_g \rangle + \sum_j \alpha_j \langle p_j^*, p_j \rangle \Rightarrow$$

Grad $j+1 \leq k \Rightarrow 1$

$$\Rightarrow \alpha_j = 0.$$

aber $\langle x P_k, P_{k-1} \rangle = \int_{-1}^1 x P_k(x) P_{k-1}(x) dx = \langle P_k, x P_{k-1} \rangle$

1) $n=1$ auf $[0,1] \Rightarrow p_1(x) = x - \frac{1}{2} \Rightarrow x_1 = \frac{1}{2}, b_1 = 1$
auf $[-1,1] \Rightarrow p_1(x) = x \Rightarrow c_1 = 0, b_1 = 1$

↳ Gauss-Quadratur mit $n=1$ Knoten \equiv MPR.

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

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

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

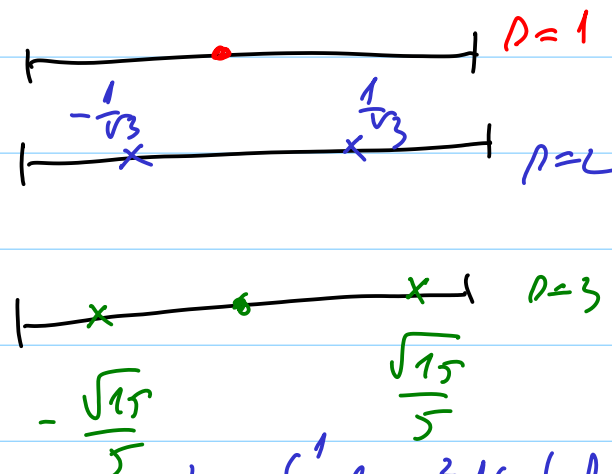
Ordnung $2n = 2 \cdot 2 = 4$ zu $p_3(1) = 1$ "normiert"

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

$x_2 = 0, x_{1,3} = \pm \frac{\sqrt{15}}{5}$

$b_2 = \frac{8}{18}, b_1 = b_3 = \frac{5}{18}$

Ordnung $2n = 2 \cdot 3 = 6$



Ben Gewichte der Gauss-QF sind positiv

Def Lagrange Polynome zu Stützstellen x_0, x_1, \dots, x_n

für $i = 0, 1, 2, \dots, n$

$$l_i(x) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{x - x_j}{x_i - x_j}$$

Bsp x_0, x_1, x_2 :

$$l_0(x) = \frac{x - x_1}{x_0 - x_1} \cdot \frac{x - x_2}{x_0 - x_2}$$

$$l_1(x) = \frac{x - x_0}{x_1 - x_0} \cdot \frac{x - x_2}{x_1 - x_2}$$

$$l_2(x) = \frac{x - x_0}{x_2 - x_0} \cdot \frac{x - x_1}{x_2 - x_1}$$

Ben: Gewichte wurden hier "Hand"-Berechnet via $b_i = \int_{-1}^1 l_i(t)^2 dt$ (l_i = Lagrange Polynom)

Ben 1) $l_i(x_j) = 0$ für alle $i \neq j$

2) $l_i(x_i) = 1$

3) $\text{Grad } l_i = n$

4) $\sum_{i=0}^n l_i(x) = 1$ für alle $x \in \mathbb{R}$.

5) $\sum_{i=0}^n l_i^{(m)}(x) = 0$ für $m \geq 1$.

6) l_0, l_1, \dots, l_n bilden eine Basis im Raum der Polynome vom Grad $\leq n$

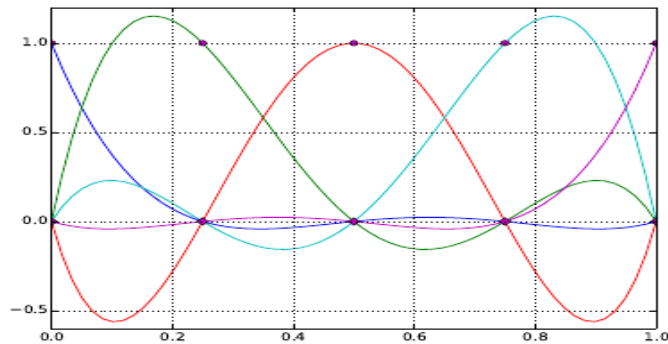


Abb. 8.3.2. Die 5 Lagrange Polynome zu Knoten $0, \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1$.

Seite 289 im Skript

Theorem Die Gewichte der Gauss QF sind positiv.

Beweis Verwende die Lagrange-Polynome zu den Stützstellen. $c_1, c_2, \dots, c_p = \text{Knoten der QF}$

$\text{Grad } l_i = p-1$ für $i=1, 2, \dots, p$ } exakt

$$\int_0^1 f(t) dt \approx \sum_{i=1}^p b_i f(c_i) \quad \Rightarrow$$

$f(t) = l_i(t)^2$ Polynom vom Grad $2(p-1)$

$$0 < \int_0^1 l_i(t)^2 dt = \sum_{j=1}^p b_j l_i(c_j)^2 = b_i \cdot 1$$

\parallel
 0 für $j \neq i$
 1 für $j = i$

Im Prinzip $b_i = \int_0^1 l_i(t)^2 dt$ könnte man verwenden um Gewichte zu rechnen

Es gibt eine bessere Methode.

3-Term Rekurrenz für Polynome:

(Achtung a, b, c sind allgemein, $c \neq$ Knoten t_1, t_2, \dots, t_n)

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

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

$$x \begin{bmatrix} p_0(x) \\ p_1(x) \\ \vdots \\ p_{k-1}(x) \\ \vdots \\ p_{n-2}(x) \\ p_{n-1}(x) \end{bmatrix} = \begin{bmatrix} -\frac{b_1}{a_1} & \frac{1}{a_1} & & & & \\ \frac{c_1}{a_1} & -\frac{b_2}{a_2} & \frac{1}{a_2} & & & \\ & \frac{c_2}{a_2} & -\frac{b_3}{a_3} & \frac{1}{a_3} & & \\ & & \frac{c_k}{a_k} & -\frac{b_k}{a_k} & \frac{1}{a_k} & \\ & & & \ddots & \frac{1}{a_{n-1}} & \\ & & & & \frac{c_n}{a_n} & -\frac{b_n}{a_n} \end{bmatrix} \begin{bmatrix} p_0(x) \\ p_1(x) \\ \vdots \\ p_{k-1}(x) \\ \vdots \\ p_{n-2}(x) \\ p_{n-1}(x) \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1/a_n \end{bmatrix}$$

3-Term Rekursion \Leftrightarrow

$$x \underline{p}(x) = \underline{A} \underline{p}(x) + \frac{1}{a_n} p_n(x) \underline{e}_n$$

Knoten der Gauss-QF sind die Nullstellen von $p_n(x)$

$$\hookrightarrow t_1, t_2, \dots, t_n \Rightarrow p_n(t_j) = 0 \text{ für } j=1, 2, \dots, n$$

$$t_j \underline{p}(t_j) = \underline{A} \underline{p}(t_j) + 0$$

$\Rightarrow t_j$ Eigenwert der Matrix \underline{A}

$\underline{p}(t_j)$ Eigenvektor der Matrix \underline{A}

$$\underline{p}(t_j) = \begin{bmatrix} p_0(t_j) \\ p_1(t_j) \\ \vdots \\ p_{n-1}(t_j) \end{bmatrix}$$

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

\underline{A} symmetrisch.

$\boxed{\text{eigh}(\underline{A})}$

Bem Es gibt ziemlich gute numerische Verfahren um EW-Probleme zu lösen

```

1 from numpy import arange, diag, sqrt
2 from numpy.linalg import eigh
3
4 def gaussquad(n):
5     r"""
6         Compute nodes and weights for Gauss-Legendre quadrature.
7
8         n: Number of node-weight pairs
9     """
10
11     i = arange(n) #i = array([0,1,...,n-1])
12     b = (i+1) / sqrt(4*(i+1)**2 - 1)
13     # now we generate the matrix; it is symmetric
14     J = diag(b, -1) + diag(b, 1)
15     # in order to find the eigenvalues we can use eigh since J is symmetric
16     x, ev = eigh(J)
17     # finally, we apply the formula for the weights
18     w = 2 * ev[0,:]**2
19     return x, w

```

auf $[-1,1]$

```

1 from sympy import *
2
3 x = Symbol("x")
4 n = Symbol("n", positive=True)
5 # First we generate two streams of symbols corresponding to knots and weights
6 xigen = numbered_symbols(prefix="xi", start=1)
7 omegagen = numbered_symbols(prefix="omega", start=1)
8 # Now we specify how many parameters we need. In this case, two knots and two
9 # weights. Therefore we set N = 2
10 N = 2 # Choose <= 3
11 xis = [ next(xigen) for i in range(N) ]
12 wis = [ next(omegagen) for i in range(N) ]
13 # Using sympy one can perform simple symbolic integrations, which deliver the
14 # exact result.
15 # In order to do this we use the function "integrate"
16 # For every n in the interval [0, 2*N-1] we set the condition that the result of
17 # the numerical integration
18 #be equal to the analytical value
19 eqns = [ integrate(x**n, (x, -1, 1)) - (sum([wi*xi**n for xi,wi in
20 zip(xis,wis)])) for n in range(2*N)]
21 pprint(eqns)
22 # Therefore we get a system of equations which we need to solve:
23 #we set eqns = 0 and find the corresponding knots and weights
24 sols = solve(eqns, numerical=True)
25 pprint(sols)

```

Berechnung der Gewichte für Gauss-Quadratur auf $[-1, 1]$

$$\langle P_i, P_k \rangle = \int_{-1}^1 \underbrace{P_i(x) P_k(x)}_{\text{Grad} \leq 2n-2} dx = \sum_{j=1}^n P_i(t_j) P_k(t_j) w_j$$

$$\begin{cases} 0, & i \neq k \\ 1, & i = k \end{cases}$$

Konstruktion: Gram-Schmidt

QF exakt bis Grad $2n-1$.

$$\underline{f}(x) = \begin{bmatrix} P_0(x) \\ P_1(x) \\ \vdots \\ P_{n-1}(x) \end{bmatrix} \quad \text{in } t_1, \dots, t_n \text{ ausgewertet}$$

Notiere

$$\underline{M} = \begin{bmatrix} \underline{f}(t_1) & \underline{f}(t_2) & \dots & \underline{f}(t_n) \end{bmatrix} \Rightarrow$$

$$\underline{I} = \underline{M} \text{diag}(w_1, \dots, w_n) \underline{M}^T$$

$\Rightarrow \underline{M}$ invertierbar

$\underline{M}^{-1} \mid \underline{M}^T$

$$\text{diag}(w_1, \dots, w_n) = (\underline{M}^T \underline{M})^{-1} \Rightarrow$$

$$\text{diag}(w_1, \dots, w_n)^{-1} = \underline{M}^T \underline{M} \Rightarrow$$

$$\begin{aligned} \frac{1}{w_j} &= \underline{f}(t_j)^T \underline{f}(t_j) = \|\underline{f}(t_j)\|^2 = \\ &= \sum_{k=0}^{n-1} P_k(t_j)^2 \end{aligned}$$

oder $w_j = \frac{1}{\|\underline{f}(t_j)\|^2}$

ABER Die Eigenvektoren sind nicht eindeutig.
ausserdem liefert eig normierte EV.

Sei \underline{v}^j ein Eigenvektor \Rightarrow es gibt c Konstante so
dass

$$\underline{v}^j = \tau \underline{f}(t_j) = \tau \begin{bmatrix} p_0(t_j) \\ \vdots \\ p_{n-1}(t_j) \end{bmatrix}$$

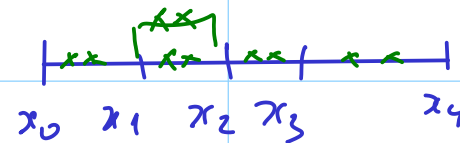
$$1 = \langle p_0, p_0 \rangle = \int_{-1}^1 p_0(t_1) p_0(t_1) dx \Rightarrow p_0(t_1) = \frac{1}{\sqrt{2}}$$

$$(\underline{v}^j)_1 = \tau \cdot p_0(t_1) = \tau \cdot \frac{1}{\sqrt{2}} \Rightarrow \tau = \sqrt{2} (\underline{v}^j)_1$$

$$\underline{f}(t_j) = \frac{1}{\tau} \underline{v}^j = \frac{1}{\sqrt{2} (\underline{v}^j)_1} \underline{v}^j \Rightarrow$$

$$\omega_j^{-1} = \|\underline{f}(t_j)\|^2 = \frac{1}{2} \frac{\|\underline{v}^j\|^2}{(\underline{v}^j)_1^2} \Rightarrow$$

$$\omega_j = \frac{2 (\underline{v}^j)_1^2}{\|\underline{v}^j\|^2}; \text{ eig liefert normiertes } v \Rightarrow \Rightarrow \omega_j = 2 (\underline{v}^j)_1^2$$



Ben 1) Gauss-Knoten sind nicht verschachtelt
 \Rightarrow Teuer bei Adaptivität,...

2) Endpunkte sind keine Knoten \Rightarrow
 Gauss-QF ist offen.

3) Manchmal braucht man ein oder beide
 Endpunkte des Integrationsintervalls als
 Knoten \Rightarrow

1 Knoten fest \Rightarrow max Ordnung $2n-1$ (Radau)

2 Knoten fest \Rightarrow max Ordnung $2n-2$ (Lobatto)

4) Fehler bei Gauss-Quadratur:

$$\int_a^b f(x) dx - \underbrace{\sum_{j=1}^n b_j f(c_j)}_{G_n(f, a, b)} = \frac{f^{(2n)}(\xi)}{(2n)!} \text{ mit } a < \xi < b$$

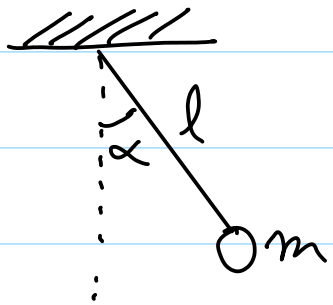
$$\left| \int_a^b f(x) dx - \sum_{k=1}^n G_n(f, x_{k-1}, x_k) \right| \leq c \cdot h^{(2n)} \max_{\xi \in [a, b]} |f^{(2n)}(\xi)|$$

$$\Rightarrow \omega_j = 2 (\underline{v}^j)_1^2$$

§2 Einfache Verfahren für ODEs 1. Ordnung

ODE = ordinary differential equations
gewöhnliche Differentialgleichungen

§2.1. Linearisierung: global und lokal



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

Physik: $\sin \alpha(t) \approx \alpha(t)$

↳ gut nur für $\alpha(t)$ klein.

$$(v) \quad \begin{cases} \ddot{\alpha}(t) = -\frac{g}{l} \sin \alpha(t) & (*) \text{ ODE 2. Ordnung.} \\ \alpha(0) = \alpha_0 & \text{da zweite Ableitung} \\ \dot{\alpha}(0) = \dot{\alpha}_0 & \text{der unbekannten} \\ & \text{Funktion } \alpha \end{cases}$$

$$\dot{\alpha}(t) = \frac{d}{dt} \alpha(t)$$

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

(*) autonom, da t nicht
explizit erscheint

"global": linearisiere die rechte Seite:

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

↳ klein für α klein.

Neues Model:

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

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

Trick: Reduziere die Ordnung der ODE.

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

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

(v) \Leftrightarrow System ODEs 1. Ordnung:

$$\begin{cases} \dot{\alpha} = p \\ \dot{p} = -m g \sin \alpha(t) \\ \alpha(0) = \alpha_0 \\ p(0) = \dot{\alpha}_0 \end{cases}$$

C) V. Gradinaru

$$\underline{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$$

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

$\underline{y}(T) = \text{exakte Lösung zur Endzeit } T$

Zeitgitter:

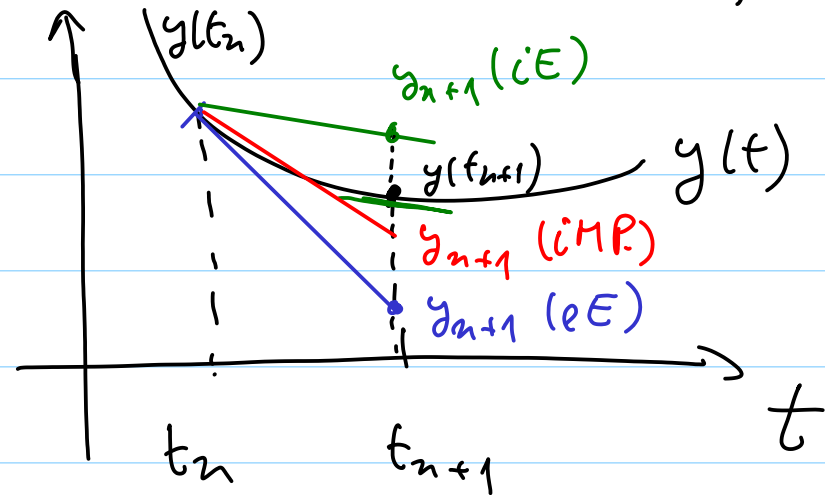
$$0 = t_0 < t_1 < t_2 < \dots < t_n < t_{n+1} < \dots < t_\infty = T$$

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

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

Idee: lokal: linearisiere die Lösung lokal.

Taylor um t_2 für $\underline{y}(t)$



Taylor um t_n : $\underbrace{y}_{y_{n+1}}(t_n+h) = \underbrace{y}_{y_n}(t_n) + h \underbrace{\dot{y}}_{f(t_n, y(t_n))}(t_n) + O(h^2)$

$$(LE) \quad \begin{cases} y_{n+1} = y_n + h f(t_n, y_n) \\ y_0 = \underline{y}(t_0) \text{ gegeben.} \end{cases} \quad \begin{array}{l} \text{expliziter Euler} \\ \text{lokaler Fehler } O(h^2) \end{array}$$

$$n = 0, 1, 2, \dots$$

$$|y(t_{n+1}) - \underline{y}_{n+1}| \leq C \cdot h^2$$

andere Herleitung:

$$\underbrace{f(t_n, \underbrace{y(t_n)}_{\underline{y}_n})}_{\approx \underline{y}_n} = \underbrace{\dot{y}(t_n)}_{\approx \underline{y}_n} = \lim_{h \rightarrow 0} \frac{y(t_n+h) - y(t_n)}{h} \approx \frac{y_{n+1} - y_n}{h}$$

$\Rightarrow (e, \epsilon)$

Taylor um t_{n+1} für $y(t)$:

$$\underbrace{y(t_{n+1}-h)}_{\underline{y}_n} = \underbrace{y(t_{n+1})}_{\underline{y}_{n+1}} - h \underbrace{\dot{y}(t_{n+1})}_{f(t_{n+1}, \underline{y}_{n+1})} + O(h^2) \quad \Rightarrow$$

$$(iE) \begin{cases} \underline{y}_n = \underline{y}_{n+1} - h f(t_{n+1}, \underline{y}_{n+1}) \\ n=0, 1, 2, \dots \quad y_0 = y(t_0) \text{ gegeben} \end{cases} \quad \begin{array}{l} \text{impliziter Euler} \\ \text{lokaler Fehler} \\ O(h^2) \end{array}$$

Bez $f(t, y) = -\frac{g}{l} \sin y$

$$y_n = y_{n+1} - h \left(-\frac{g}{l}\right) \sin y_{n+1} \quad \begin{array}{l} \text{algebraische Gleichung} \\ \text{für } y_{n+1} \end{array}$$

Neue Idee: linearisiere lokal um die Mitte des Intervall:

Taylor $t^* = \frac{1}{2}(t_n + t_{n+1}) = t_n + \frac{1}{2}h$ für $y(t)$:

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

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

$\ominus \Rightarrow$

$$\underbrace{y(t_{n+1}) - y(t_n)}_{\parallel f(t^*, \underline{y}(t^*))} = h \dot{y}(t^*) + O(h^3) \quad \Rightarrow$$

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

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

Bez $y(t^*)$ stört, brauche noch ein Trick!

$$1) \quad \underline{y}(t^*) \approx \frac{1}{2} \left(\underset{\substack{\text{ZZ} \\ \underline{y}_n}}{\underline{y}(t_n)} + \underset{\substack{\text{ZZ} \\ \underline{y}_{n+1}}}{\underline{y}(t_{n+1})} \right) \Rightarrow$$

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

impliziter Mittelpunktsregel / lokaler Fehler
 $O(h^3)$

$$2) \quad \underline{f}(t^*, \underline{y}(t^*)) \approx \frac{1}{2} \left(\underline{f}(t_n, \underline{y}_n) + \underline{f}(t_{n+1}, \underline{y}_{n+1}) \right) \Rightarrow$$

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

impliziter Trapezregel / lokaler Fehler $O(h^3)$

Ben
Das geht so nur wenn $\underline{y}(t)$ glatt genug ist.

Bemerkung Addieren wir die 2 Gleichungen:

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

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

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

Bemerkung Selbe Rechnung um $t_n \Rightarrow$

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

§2.2. Störmer-Verlet Verfahren

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

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

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

expliziter, 2-Schritt-Verfahren

\Rightarrow brauche Startwerte $\underline{y}(t_0) = \underline{y}_0$
 $\underline{y}_1 \approx \underline{y}(t_1)$

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

zweite Möglichkeit: Taylor mit 3 Termen:

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

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

(St-V): local $O(h^3)$ und global $O(h^2)$

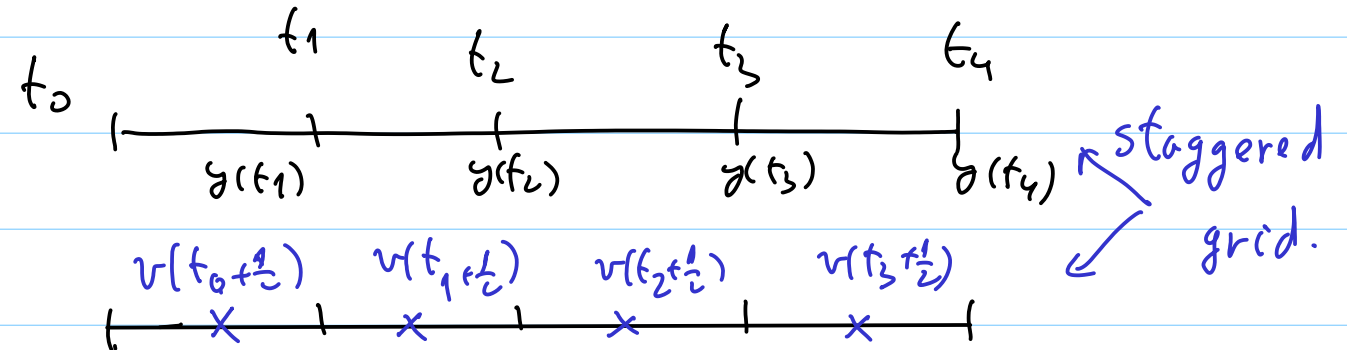
$$|\underline{y}(T) - \underline{y}(t_n)| \leq c \cdot h^2$$

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

(gate Approximation on $\dot{\underline{y}}(t_k + \frac{1}{2}h)$)

(in St-V: \Rightarrow)

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



Besser noch: "velocity-Verlet"-Verfahren:

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

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

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

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

Vorteile: explizit, genauer als (FE), (IE)
erhält die Energie!

§2.3. Vorgehensweise bei Implementierung

Ziel:

Pendelgleichung "lösen": Approximation der Lösung und der Energie mittels vorigen Methoden.

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

1. Schritt: Umschreiben in ODE 1. Ordnung.

$$p = \dot{\alpha}, \quad \underline{y} = \begin{bmatrix} \alpha \\ p \end{bmatrix}, \quad f(\underline{y}) = \begin{bmatrix} p \\ -\frac{g}{l} \sin \alpha \end{bmatrix} = \begin{bmatrix} \dot{y}_1 \\ -\frac{g}{l} \sin y_1 \end{bmatrix}$$

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

Möchte: Lösung & Energie der Lösung zu Zeitpunkten t_2 .

$$0 = t_0 < t_1 < t_2 < \dots < t_N = T$$

Gegeben: $\alpha_0, \dot{\alpha}_0, g, l$ ($\omega^2 = \frac{g}{l}$), $t_0 = 0, T$
Wähle N !

$$\underline{y}(t_0), \underline{y}(t_1), \dots, \underline{y}(t_N)$$

Potentiale Energie: $V(\underline{y}) = -g \cos y_1 + C$

Kinetische Energie: $K(\underline{y}) = \frac{1}{2} \underline{y}_2^2$

Totale Energie $E(\underline{y}) = V(\underline{y}) + K(\underline{y})$

```
def potE(y):
    return -g cos y[1, 0]
```

$$\underline{y} = \begin{bmatrix} x \\ \alpha \end{bmatrix}$$

$\underline{y}(t_0)$

$\underline{y}(t_1)$

```
def kinE(y):
    return 1/2 y[2, 1]^2
```

```
def totE(y):
    return potE(y) + kinE(y)
```

Speicherplatz vorbereiten!

$$\begin{matrix} t_0 & t_1 & \dots & t_N \\ \begin{bmatrix} y(t_0) \\ y(t_1) \end{bmatrix} & \begin{bmatrix} y(t_0) \\ y(t_1) \end{bmatrix} & \dots & \begin{bmatrix} y(t_0) \\ y(t_N) \end{bmatrix} \end{matrix} \Rightarrow \text{array} \begin{matrix} 0 & 1 & \dots & N \\ \hline 0 & 0 & & 0 \\ 0 & 0 & & 0 \end{matrix}$$

(2, N+1) oder

$$y = \text{zeros}(N+1, 2)$$

$$y(t_0) = \text{array}([\alpha_0, \dot{\alpha}_0])$$

$$\begin{matrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{matrix} \begin{matrix} t_0 \\ t_1 \\ \vdots \\ t_N \end{matrix}$$

$$(nA) \quad y[n, 0] = \frac{\dot{\alpha}}{\omega} \sin(\omega t_n) + \alpha_0 \cos(\omega t_n) = g(t_n)$$

$$y[n, 1] = \text{Ableitung davon} = \dot{\alpha} \cos(\omega t_n) - \alpha_0 \omega \sin(\omega t_n) = dg(t_n)$$

$$\left[\begin{array}{l} \text{für } n = 0, 1, 2, \dots, N-1: \\ y[n+1, 0] = g(t_{n+1}) \\ y[n+1, 1] = dg(t_{n+1}) \end{array} \right.$$

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

→ implementiere die rechte Seite f .

def $f(\underline{y})$:

$$\text{return} \begin{pmatrix} y(1) \\ \frac{g}{2} \sin y(2) \end{pmatrix} \quad \left(-\omega^2 \sin \alpha(t) \right)$$

für $n=0, 1, 2, \dots, n-1$

~~$$y[n+1, 0] = y[n, 0] + h f(\underline{y})[0]$$~~

~~$$y[n+1, 1] = y[n+1, 0] + h f(\underline{y})[1]$$~~

$$\underline{y}[n+1] = \underline{y}[n] + h f(\underline{y})$$

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

für $n=0, 1, 2, \dots, n-1$:

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

erste Möglichkeit: allgemein.

Sei die Unbekannte $\underline{z} \Rightarrow (iE) \quad \underline{z} = \underline{a} + h f(\underline{z})$

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

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

$\underline{z} \in \mathbb{R}^d$

Finde die Nullstelle(n) von \underline{F} !

scipy.optimize.fsolve:

fsolve(\underline{F} , \underline{z}_0) ↗ Startwert
(nah der gesuchten
Nullstelle)

Startwert für f_{solve} in einem Schritt von (ϵ) ?

$\underline{z}_0 = \underline{y}_n$ oder Lösung mit lin. Rechte Seite
oder was (ϵ) vorschlägt

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

für $n=0, 1, 2, \dots, N-1$:

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

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

zweite Möglichkeit:

nutze die Form der Gleichung des Pendels.

alg. Problem:

$$\begin{cases} z_1 - a_1 - h z_2 = 0 \end{cases}$$

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

$$\Rightarrow z_1 - a_1 - h \left(a_2 + h \frac{g}{l} \sin z_1 \right) = 0 \quad (\Leftrightarrow)$$

$$\Leftrightarrow G(z) = 0 \quad \text{mit} \quad G(z) = z - a_1 - h a_2 - h^2 \frac{g}{l} \sin z$$

alg. Gleichung in einer Unbekannten $z \in \mathbb{R}$

↖ billiger, einfacher!

$$f_{\text{solve}} \Rightarrow z_1 \Rightarrow z_2$$

§2.4 Fehlerschätzung und Konvergenz

Taylor mit Rest als Integral:

$$f(x) = f(a) + \frac{x-a}{1!} f'(a) + \frac{(x-a)^2}{2!} f''(a) + \dots + \frac{(x-a)^{n-1}}{(n-1)!} f^{(n-1)}(a) + \int_a^x \frac{(x-t)^{n-1}}{(n-1)!} f^{(n)}(t) dt$$

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

$\underline{y}_n \approx \underline{y}(t_n)$ mit $t_n = t_0 + nh$; h = Zeitschrittweite.

↳ vom numerischen Verfahren definiert

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

und Lipschitz:

$$\|f(t, \underline{y}) - f(t, \underline{z})\| \leq L \|\underline{y} - \underline{z}\| \quad \text{für alle } \underline{y}, \underline{z} \in \mathbb{R}^d, t \in [t_0, T]$$

↗ konstante $L \in \mathbb{R}$.

$$\text{Explizite Euler: } \underline{y}_{n+1} = \underline{y}_n + h f(t_n, \underline{y}_n)$$

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

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

Beweis: 3 Schritte.

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

$$\begin{aligned} \underline{y}(t_{n+1}) - \underline{y}_{n+1} &= \underline{y}(t_{n+1}) - (\underline{y}_n + h f(t_n, \underline{y}_n)) \\ &= \underline{y}(t_{n+1}) - \underline{y}(t_n) - h \underbrace{f(t_n, \underline{y}(t_n))}_{\|\dot{\underline{y}}(t_n)\|} \end{aligned}$$

Satz von Taylor mit Rest als Integral:

$$a = t_n, \quad x = a + h = t_{n+1}, \quad f = \underline{y} \quad \Rightarrow$$

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

Möchte aber \int_0^1 , als wie bei der Fehlerschätzung vom Quadrat: Variablenwechsel:

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

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

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

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

2) Fehlerfortpflanzung:

$$\left. \begin{array}{l} (e\bar{E}) \text{ mit Startwert } \underline{z}_n: \quad \underline{z}_{n+1} = \underline{z}_n + h \underline{f}(t_n, \underline{z}_n) \\ \underline{w}_n: \quad \underline{w}_{n+1} = \underline{w}_n + h \underline{f}(t_n, \underline{w}_n) \end{array} \right\}$$

$$\Rightarrow \|\underline{z}_{n+1} - \underline{w}_{n+1}\| \leq \|\underline{z}_n - \underline{w}_n\| + h \|\underline{f}(t_n, \underline{z}_n) - \underline{f}(t_n, \underline{w}_n)\| \leq$$

$$\leq \|\underline{z}_n - \underline{w}_n\| + hL \|\underline{z}_n - \underline{w}_n\| = \|\underline{z}_n - \underline{w}_n\| (1 + hL)$$

3) Fehlerakkumulation:

$$\|\underline{y}(t_n) - \underline{y}_n\| \leq \overset{\text{1. Schritt}}{ch^2} + \overset{\text{2. Schritt}}{ch^2(1+hL)} + \overset{\text{3. Schritt}}{ch^2(1+hL)^2} + \dots + ch^2(1+hL)^{n-1} =$$

$$= ch^2 \frac{(1+hL)^n - 1}{1+hL-1} = ch \frac{(1+hL)^n - 1}{L} \leq$$

$$\leq ch \frac{e^{nhL} - 1}{L} = ch \frac{e^{(t_n - t_0)L} - 1}{L}$$

(imp): lokale Fehler: $t^* = \frac{1}{2}(t_n + t_{n+1})$
 $\bar{y} = \frac{1}{2}(\underline{y}_n + \underline{y}_{n+1})$

$$\underline{y}(t_{n+1}) - \underline{y}(t_n) - h \underline{f}(t^*, \bar{y}) \stackrel{\text{Taylor in } t^*}{=} \overset{\text{ODE}}{=}$$

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

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

§3 Strukturhaltung.

§3.1. Invariante und Hamilton Systeme

Bsp autonome Lotka-Volterra: $d=2$

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

$$\left. \begin{array}{l} \dot{u} = (\alpha - \beta v)u = \left(\frac{\alpha}{v} - \beta\right)uv \\ \dot{v} = (\delta u - \gamma)v = \left(\delta - \frac{\gamma}{u}\right)uv \end{array} \right\} \rightarrow$$

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

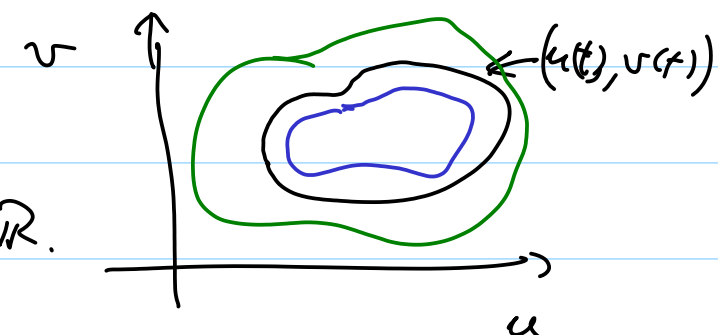
$$\frac{d}{dt} \left(\delta u - \gamma \log u - \alpha \log v + \beta v \right) = 0 \quad \text{für alle } t \Rightarrow \quad I(u(t), v(t))$$

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

$\Rightarrow I(u(t), v(t)) = \text{konstant}$, wenn u, v die Lösungen des L-V-Systems sind.

d.h. die Trajektorien sind Niveaulinien

der Funktion $I: \mathbb{R}^2 \rightarrow \mathbb{R}$.



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

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

partielle Ableitung $\frac{\partial}{\partial x_1} f(\underline{x}) = \text{Ableitung von } f \text{ nach } x_1$
 $\frac{\partial}{\partial x_2}, \dots, \frac{\partial}{\partial x_d}$

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

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

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

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

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

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

\Leftrightarrow

$\underline{\text{grad}} I(\underline{y}) \cdot f(t, \underline{y}) = 0$ für alle $(t, \underline{y}) \in [0, T] \times D$
 wo es eine Lösung gibt.
 (y(t))

Beweis Annahme.

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

\uparrow
 $\underline{y}(t)$ Lösung

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

\uparrow
kettenregel.

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

\Rightarrow : Nehme $\underline{t}_0 \in D$, $(t_0 \in [0, T])$ beliebig.

Verwende die ODE: $\begin{cases} \dot{\underline{y}} = f(t, \underline{y}) \\ \underline{y}(t_0) = \underline{t}_0 \end{cases} \Rightarrow \underline{y}(t) \xrightarrow{I \text{ Invariant}} \underline{y}(t)$

$$\Rightarrow I(\underline{y}(t)) = \text{konstant} \xrightarrow[\text{ODE}]{\frac{d}{dt}} \underline{\text{grad}} I(\underline{y}(t)) \dot{\underline{y}}(t) = 0 = f(t, \underline{y}(t))$$

für alle t $\left\{ \begin{array}{l} \Rightarrow \underline{\text{grad}} I(\underline{y}(t_0)) \dot{\underline{y}}(t_0) = 0 \Rightarrow \\ \text{nehme } t = t_0 \end{array} \right. \underline{\text{grad}} I(\underline{y}(t_0)) \cdot f(t_0, \underline{t}_0) = 0$

\Rightarrow für alle t_0, \underline{t}_0 .

Notation $H: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ Hamilton Funktion.
 $\begin{matrix} \psi & \psi \\ \underline{p} & \underline{q} \\ - & - \end{matrix}$

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

Def Hamiltonische Differentialgleichung.

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

für $j=1, 2, \dots, d$

autonomes Hamilton-System mit der Hamilton-Funktion H .

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

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

$q = \alpha = \text{Winkel}$ $p = \dot{q}$

2) Konservatives Kraftfeld:

d.h.

$$\underline{f}(\underline{x}) = - \text{grad } U(\underline{x}) \quad \underline{x} \in \mathbb{R}^d$$

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

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

$p = m \dot{r} \Rightarrow$ Bewegungsgleichungen
 eines Massenpunktes m
 im konservativen Zentralfeld.

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

vergleiche mit
 Newton's Gleichung,
 $m \ddot{r}(t) = f(r(t)).$

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

Bem $H(p, q)$ Invariante für das Hamilton System.

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

Hamilton System:

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

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

Überprüfe dies für $\underline{I} := H$:

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

$$\text{grad } H(\underline{y}) = \begin{bmatrix} \underline{a} \\ \underline{b} \end{bmatrix} \Rightarrow \begin{bmatrix} \underline{a} & \underline{b} \end{bmatrix} \begin{bmatrix} \underline{0} & -\underline{I} \\ \underline{I} & \underline{0} \end{bmatrix} \begin{bmatrix} \underline{a} \\ \underline{b} \end{bmatrix} =$$

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

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

$\underline{B} \in \mathbb{R}^{d \times d}$, \underline{I} Invariante für $\dot{\underline{y}} = \underline{f}(\underline{y})$ mit \underline{f} differentierbar.

Sei (\underline{y}_n) die numerische Approximation aus (17).

Dann

$$\underline{I}(\underline{y}_n) = \underline{I}(\underline{y}_0) \text{ für alle } n.$$

Beweis

$$\underline{I}(\underline{y}) \in \mathbb{R} \Rightarrow \underline{I}(\underline{y}) = \underline{I}(\underline{y})^T = \left(\frac{1}{2} \underline{y}^T \underline{B} \underline{y} \right)^T =$$

$$= \frac{1}{2} \underline{y}^T \underline{B}^T \underline{y} \Rightarrow$$

$$2\underline{I}(\underline{y}) = \underline{I}(\underline{y}) + \underline{I}(\underline{y})^T = \frac{1}{2} \underline{y}^T (\underline{B} + \underline{B}^T) \underline{y} = \text{konstant.}$$

$\frac{1}{2} \underline{y}^T \underline{A} \underline{y}$ Invariante mit $\underline{A} = \underline{B} + \underline{B}^T$ also symmetrisch.

Invarianz $\Rightarrow \frac{d}{dt} = 0 \Rightarrow \underline{A} \underline{y}(t) \cdot \underline{f}(\underline{y}(t)) = 0$ für $\underline{y}(t)$ Lösung.
 \underline{A} symmetrisch.

$$\begin{aligned} I(\underline{y}_{n+1}) - I(\underline{y}_n) &= \frac{1}{2} (\underline{y}_{n+1} + \underline{y}_n)^T \underline{A} (\underline{y}_{n+1} - \underline{y}_n) = \\ &= \underbrace{\left(\underline{A} \frac{1}{2} (\underline{y}_{n+1} + \underline{y}_n) \right)^T}_{\text{IMP}} \cdot \underbrace{(\underline{y}_{n+1} - \underline{y}_n)}_{\tau} = \\ &= \underline{A} \left(\frac{1}{2} (\underline{y}_{n+1} + \underline{y}_n) \right) \cdot h \underline{f} \left(\frac{1}{2} (\underline{y}_{n+1} + \underline{y}_n) \right) = \\ &= (\underline{A} \underline{\tau}) \cdot h \underline{f}(\underline{\tau}) = 0 \end{aligned}$$

$$\frac{\partial}{\partial \underline{y}_i} (\underline{y}^T \underline{A} \underline{y}) = \dots$$

$$\frac{d}{dt} z I(\underline{y}) = \underline{A} \underline{y}(t) \cdot \underline{\dot{y}}(t) = \underline{A} \underline{y}(t) \cdot \underline{f}(\underline{y}(t))$$

$\xrightarrow[\dot{\underline{y}} = \underline{f}(\underline{y})]{\text{ODE}}$

Da:

$$\begin{aligned} \underline{y}_{n+1}^T \underline{A} \underline{y}_{n+1} - \underline{y}_n^T \underline{A} \underline{y}_n &= \underline{y}_{n+1}^T \underline{A} \underline{y}_{n+1} - \underline{y}_{n+1}^T \underline{A} \underline{y}_n \\ &\quad + \underline{y}_{n+1}^T \underline{A} \underline{y}_n - \underline{y}_n^T \underline{A} \underline{y}_n \\ &= \underline{y}_{n+1}^T \underline{A} \underline{y}_{n+1} - \underline{y}_{n+1}^T \underline{A} \underline{y}_n + \underline{y}_{n+1}^T \underline{A} \underline{y}_n - \underline{y}_n^T \underline{A} \underline{y}_n \end{aligned}$$

$$\begin{aligned} &= \underline{y}_{n+1}^T \underline{A} (\underline{y}_{n+1} - \underline{y}_n) + (\underline{y}_{n+1} - \underline{y}_n)^T \underline{A} \underline{y}_n = (\underline{y}_{n+1}^T + \underline{y}_n^T) \underline{A} (\underline{y}_{n+1} - \underline{y}_n) \\ &\quad \underbrace{\hspace{10em}}_{\text{A symmetrisch.}} \\ &\quad \left((\underline{y}_{n+1} - \underline{y}_n)^T \underline{A} \underline{y}_n \right)^T = \underline{y}_n^T \underline{A} (\underline{y}_{n+1} - \underline{y}_n) \end{aligned}$$

§3.2. Splitting Verfahren

autonome ODE: $\dot{\underline{y}} = \underline{f}(\underline{y})$

Bem Was tun wenn ODE nicht autonom ist?

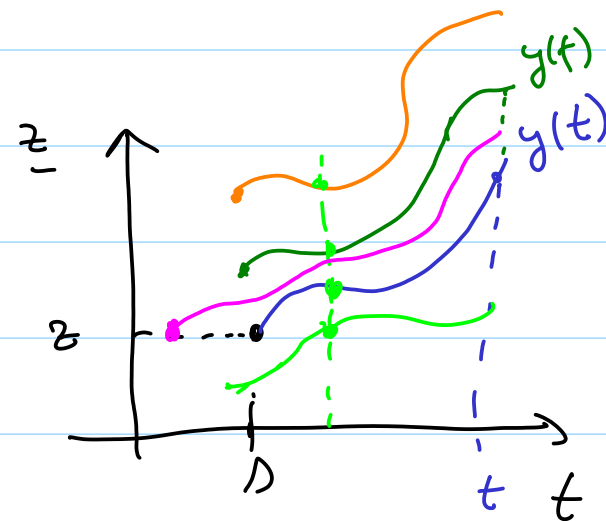
Trick:

$$\underline{z} = \begin{bmatrix} t \\ \underline{y} \end{bmatrix} \Rightarrow \dot{\underline{z}} = \begin{bmatrix} 1 \\ \dot{\underline{y}} \end{bmatrix} = \begin{bmatrix} 1 \\ g(t, \underline{y}) \end{bmatrix} = \begin{bmatrix} 1 \\ g(\underline{z}) \end{bmatrix} = \underline{f}(\underline{z})$$

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

$\Rightarrow \dot{\underline{z}} = \underline{f}(\underline{z})$
mit $\underline{z}(0) = \begin{bmatrix} t_0 \\ \underline{y}(t_0) \end{bmatrix}$

Bem $\begin{cases} \dot{\underline{y}} = g(t, \underline{y}) \\ \underline{y}(1) = \underline{x} \end{cases}$

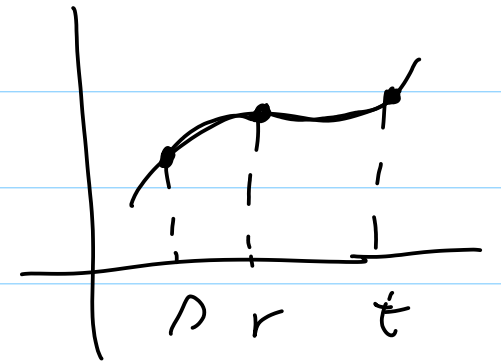


$$\Phi^{1,t} \underline{x} = \underline{y}(t) \text{ mit } \underline{y}(1) = \underline{x}$$

$\Phi^{1,t}: \mathcal{D} \rightarrow \mathcal{D}$ zweiparametrische Familie von Abbildungen. "Fluss der ODE".

1) $\Phi^{t,t} = \text{Identität!}$ $\Phi^{t,t} \underline{x} = \underline{x}$ für alle $\underline{x} \in \mathcal{D}$.

2) $\Phi^{1,t} \underline{x} = \Phi^{r,t} \Phi^{1,r} \underline{x}$



3) Falls ODE autonom ist: $\Phi^{1,t} = \Phi^{0,t-1} = \Phi^{t-1}$
d.h. die Lösung der autonomen ODE

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

ist translationsinvariant.

Beweis:

$$\underline{u}(t) = \underline{y}(t+c) \Rightarrow \frac{d}{dt} \underline{u}(t) = \dot{\underline{y}}(t+c) \cdot 1 = \underline{f}(\underline{y}(t+c)) = \underline{f}(\underline{u}(t))$$

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

Bsp $e^{mh} = \sum_{n=0}^{\infty} \frac{1}{n!} (mh)^n$

↪ auch als Lösung der ODE:

$$\dot{y}(h) = m y(h) ; y(0) = 1$$

Für mehrere solchen linearen, entkoppelten ODE:

$$\begin{cases} \dot{y}_1(t) = m_1 y_1(t) \Rightarrow y_1(h) = e^{m_1 h} \\ \dot{y}_2(t) = m_2 y_2(t) \Rightarrow y_2(h) = e^{m_2 h} \\ \dots \\ \dot{y}_d(t) = m_d y_d(t) \Rightarrow y_d(h) = e^{m_d h} \end{cases}$$

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_d \end{bmatrix} : \quad \dot{\underline{y}} = \text{diag}(m_1, \dots, m_d) \underline{y}$$

$$\underline{y}(h) = \text{diag}(e^{m_1 h}, \dots, e^{m_d h}) \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$= e^{h \text{diag}(m_1, \dots, m_d)} \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$e^{\text{diag}(m_1, \dots, m_d)h} = \sum_{n=0}^{\infty} \frac{1}{n!} \left(\text{diag}(m_1, \dots, m_d)h \right)^n$$

$$\dot{\underline{y}} = \underline{M} \underline{y} \quad \text{mit } \underline{M} \text{ eine beliebige } d \times d \text{-Matrix.}$$

mit Lösung $\underline{y}(t) = e^{\underline{M}t} \underline{y}_0$ wobei:

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

Bsp $\underline{M} = \underline{A} + \underline{B}$

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

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

$$= \left(\underline{I} + \underline{(\underline{A} + \underline{B})h} + \left(\frac{1}{2} \underline{A}^2 + \underline{A}\underline{B} + \frac{1}{2} \underline{B}^2 \right) h^2 \right)$$

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

Man kann beweisen dass $e^{(\underline{A} + \underline{B})h} \approx e^{\underline{A}h} e^{\underline{B}h}$

$$e^{Bh}(e^{Ah}y_0) \cdot y(h) = e^{(A+B)h}y_0$$

$$e^{Bh}y_1 \xrightarrow{e^{Ah}(e^{Bh}y_1)} e^{(A+B)h}y_0$$

Bez Für nicht-lineare autonome ODE erster Ordnung:

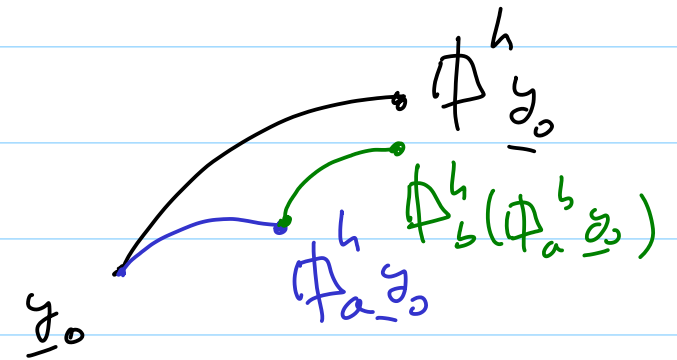
$$\dot{\underline{y}} = \underline{f}(\underline{y}) \quad \text{mit} \quad \underline{f}(\underline{y}) = \underline{f}_a(\underline{y}) + \underline{f}_b(\underline{y})$$

Idee: Wähle $\underline{f}_a, \underline{f}_b$ so dass die ODE:

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

und

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



$$e^{(A+B)h} \approx e^{Bh}(e^{Ah}y_0)$$

$$e^{(A+B)h} \approx e^{Ah}(e^{Bh}y_0)$$

einfach oder exakt lösbar sind.

$$\Psi_1^h := \Phi_b^h \circ \Phi_a^h$$

Lie-Trotter Splitting

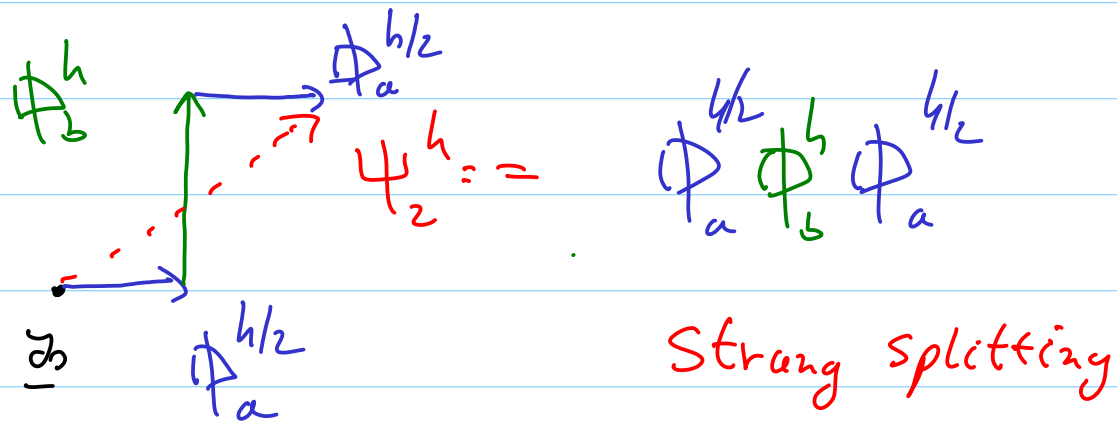
Fehler $O(h)$ zur Endzeit T

Bez Symmetrie schenkt uns eine Ordnung mehr:

$$\Psi_1^h := \Phi_b^h \circ \Phi_a^h$$

Lie-Trotter Splitting

$$\tilde{\Psi}_1^h := \Phi_a^h \circ \Phi_b^h$$



Allgemeines Splittingverfahren:

$$\Psi^h = \prod_{i=1}^n \Phi_b^{b_i h} \Phi_a^{a_i h}$$

mit $\sum_{i=1}^n a_i = 1$, $\sum_{i=1}^n b_i = 1$.

Bsp $n=1$, $a_1 = b_1 = 1$ 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?
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:

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

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

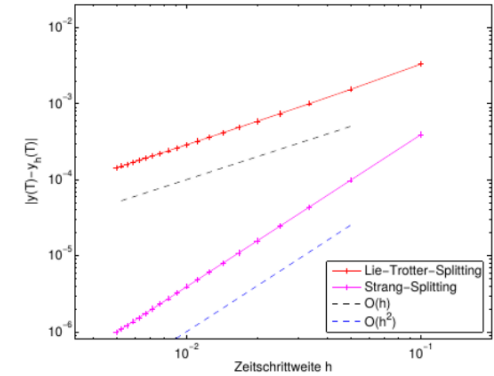
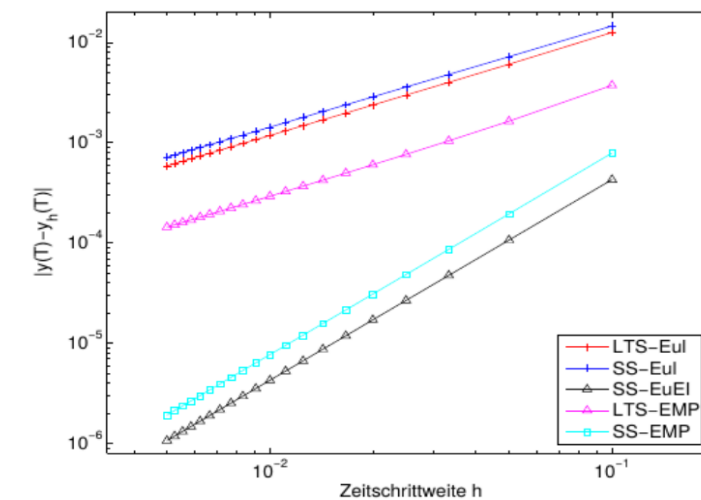


Abb. 2.4.2. Fehler zur Endzeit $T = 1$.



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 Φ_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 $\Psi_{h,a}^h$, $\Psi_{h,f}^h$ und Strang-Splitting

Bsp Splitting Verfahren für Newton's Gleichung.

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

$$\underline{F}(\underline{y}) = \underbrace{\begin{bmatrix} 0 \\ \underline{a}(\underline{r}) \end{bmatrix}}_f + \underbrace{\begin{bmatrix} \underline{v} \\ 0 \end{bmatrix}}_g \quad \text{Startwert} \quad \underline{y}_0 = \begin{bmatrix} \underline{r}_0 \\ \underline{v}_0 \end{bmatrix}$$

(a) $\begin{cases} \dot{\underline{r}} = 0 \\ \dot{\underline{v}} = \underline{a}(\underline{r}) \end{cases}$ Löse von $\underbrace{\text{Startzeit } 0}_{\text{zu Endzeit } h}$: **Exakt!**

$\rightarrow \underline{r}(h) = \underline{r}(0) = \underline{r}_0$

\downarrow

$\dot{\underline{v}} = \underline{a}(\underline{r}_0) \Rightarrow \underline{v}(h) = \underline{a}(\underline{r}_0)h + \underline{v}_0$

$$\Rightarrow \Phi_f^h \begin{bmatrix} \underline{r}_0 \\ \underline{v}_0 \end{bmatrix} = \begin{bmatrix} \underline{r}_0 \\ \underline{v}_0 + h \underline{a}(\underline{r}_0) \end{bmatrix}$$

(b) $\begin{cases} \dot{\underline{r}} = \underline{v} \\ \dot{\underline{v}} = 0 \end{cases} \xrightarrow{\quad} \dot{\underline{r}} = \underline{v}_0 \Rightarrow \underline{r}(h) = \underline{r}_0 + h \underline{v}_0$

$\Rightarrow \underline{v}(h) = \underline{v}(0) = \underline{v}_0$ **exakt!**

$$\Rightarrow \Phi_g^h \begin{bmatrix} \underline{r}_0 \\ \underline{v}_0 \end{bmatrix} = \begin{bmatrix} \underline{r}_0 + h \underline{v}_0 \\ \underline{v}_0 \end{bmatrix}$$

kombinieren wir diese 2 exakten Lösungen:

(1) Lie-Trotter-Splitting:

$$\Psi^h \begin{bmatrix} \underline{r} \\ \underline{v} \end{bmatrix} = \underbrace{\Phi_g^h \circ \Phi_f^h}_{\text{Lie-Trotter-Splitting}} \begin{bmatrix} \underline{r} \\ \underline{v} \end{bmatrix} = \underbrace{\Phi_g^h}_{\text{symplektische Euler-Verfahren}} \begin{bmatrix} \underline{r} \\ \underline{v} + h \underline{a}(\underline{r}) \end{bmatrix} =$$

$$= \begin{bmatrix} \underline{r} + h (\underline{v} + h \underline{a}(\underline{r})) \\ \underline{v} + h \underline{a}(\underline{r}) \end{bmatrix}$$

symplektische Euler-Verfahren.

(2) Strang-Splittung:

$$\begin{aligned}
 \Psi^h \begin{bmatrix} \underline{r} \\ \underline{v} \end{bmatrix} &= \Phi_g^{h/2} \Phi_f^h \Phi_g^{h/2} \begin{bmatrix} \underline{r} \\ \underline{v} \end{bmatrix} = \\
 &= \Phi_g^{h/2} \Phi_f^h \begin{bmatrix} \underline{r} + \frac{h}{2} \underline{v} \\ \underline{v} \end{bmatrix} = \Phi_g^{h/2} \begin{bmatrix} \underline{r} + \frac{h}{2} \underline{v} \\ \underline{v} + h a(\underline{r} + \frac{h}{2} \underline{v}) \end{bmatrix} \\
 &= \begin{bmatrix} \underline{r} + \frac{h}{2} \underline{v} + \frac{h}{2} (\underline{v} + h a(\underline{r} + \frac{h}{2} \underline{v})) \\ \underline{v} + h a(\underline{r} + \frac{h}{2} \underline{v}) \end{bmatrix}
 \end{aligned}$$

Notation:

$$\begin{cases} \underline{r}_{k+\frac{1}{2}} = \underline{r}_k + \frac{1}{2} h \underline{v}_k \\ \underline{v}_{k+1} = \underline{v}_k + h a(\underline{r}_{k+\frac{1}{2}}) \\ \underline{r}_{k+1} = \underline{r}_k + \frac{1}{2} h \underline{v}_{k+1} \end{cases}$$

\equiv ein-Schritt-Formulierung des St-V !!

Bem Trick geht auch für separablen
Hamilton-System

$$H(\underline{p}, \underline{q}) = \underline{T}(\underline{p}) + \underline{V}(\underline{q})$$

Lie-Trotter \Rightarrow
$$\begin{cases} \underline{p}_{n+1} = \underline{p}_n - h \underline{\text{grad}} V(\underline{q}_n) \\ \underline{q}_{n+1} = \underline{q}_n + h \underline{\text{grad}} T(\underline{p}_{n+1}) \end{cases}$$

Symplektische Euler-Verfahren; global $O(h)$.

Strang-Splittung \Rightarrow (St-V) für separablen Hamilton-Systeme.

Bez Was tun für Zeitabhängige rechte Seiten?

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

Autonomisieren!

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

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

Nicht autonom \Rightarrow autonomisieren

Unbekannte $t \Rightarrow \dot{t} = 1$

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

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

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

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

exakt.

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

(a_c) (b_c) aus Splitting $\xrightarrow{(a,b)}$ Methode!

Processing

$$\hat{\Psi}^h = \Pi^h \circ \Psi^h \circ (\Pi^h)^{-1}$$



post-processor

pre-processor

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

Vorteile falls:

+ $\hat{\Psi}^h$ genauer als Ψ^h

+ $\Pi^h, (\Pi^h)^{-1}$ günstig

+ keine/wenige Ausgaben der Lösung vor Endzeit gewünscht!

Bsp Strang-Splitting:

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

$\underbrace{\Phi_a^{h/2} \circ (\Phi_a^{h/2})^{-1}}_{\text{Lie-Trotter.}}$

Leicht gestörte Probleme

$$\dot{\underline{y}} = \underline{f}_a(\underline{y}) + \varepsilon \underline{f}_b(\underline{y}) \quad \text{mit } \underline{\varepsilon} \text{ klein.}$$

optimierte Splitting-Verfahren, klar

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

$$\text{mit } r_1 \geq r_2 + 1 \geq \dots$$

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

§4 Runge-Kutta-Verfahren

§4.1. Grundidee

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

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

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

QF mit Gewichten b_i , Knoten $c_i \in [0, 1] \Rightarrow$

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

$$\underline{y}(t_1) \approx \underline{y}(t_0) + \boxed{h} \sum_{i=1}^n b_i k_i$$

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

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

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

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

Idee: verwende etwas Biliyeres für \nearrow

$$\underline{y}(t_0 + h) \approx \underline{y}_0 + h \underbrace{\underline{f}(t_0, \underline{y}_0)}_{k_1} \quad (e \in)$$

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

explizite Trapezregel

Bsp: QF = Mittelpunktsregel

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

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

$$\begin{cases} \underline{k}_1 := \underline{f}(t_0, \underline{y}_0) \\ \underline{k}_2 = \underline{f}(t_0 + \frac{h}{2}, \underline{y}_0 + \frac{h}{2} \underline{k}_1) \\ \underline{y}_1 = \underline{y}_0 + h \underline{k}_2 \end{cases}$$

explizite Mittelpunktsregel

(iE) \Rightarrow implizite MPR ...

Def Runge-Kutta-Verfahren mit s Stufen:

Gegeben Butcher-Schema

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

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

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

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

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

ein $(s \cdot d) \times (s \cdot d)$ nichtlineares
algebraisches
Gleichungssystem

$$\underline{A} = \begin{bmatrix} 0 & & 0 \\ * & \ddots & \\ & & 0 \end{bmatrix} \quad a_{ij} = 0 \text{ für alle } i \leq j$$

\Rightarrow RK explizit

$$\underline{k}_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^{i-1} a_{ij} \underline{k}_j)$$

$$\underline{A} = \begin{bmatrix} \diagdown & & 0 \\ * & \diagdown & \\ & & \end{bmatrix} \Rightarrow \text{diagonal implizite RK}$$

$$\underline{k}_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^{i-1} a_{ij} \underline{k}_j + h a_{ii} \underline{k}_i)$$

$d \times d$ nicht-lin. alg. Gleichung

$\underline{y} \in \mathbb{R}^d, \underline{k}_i \in \mathbb{R}^d, f: \mathbb{R} \times \mathbb{R}^d \rightarrow \mathbb{R}^d$

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

Def Konsistenzordnung q wenn

lokale Fehler $\|\underline{y}(t_0 + h) - \underline{y}_1\| \leq c \cdot h^{q+1}$

Theorem

RK hat Konsistenzordnung $q \Rightarrow$

QF hat Ordnung q .

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

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

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

$\leq c \cdot h^{q+1} \Rightarrow$

Voraussetzung \nearrow

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

$$\left| \frac{1}{n+1} - \sum_{j=1}^n b_j c_j^n \right| \leq c \cdot h^{q-n} \xrightarrow{h \rightarrow 0} 0 \Rightarrow$$

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 $p_n(t) = t^n$

Konsequenz RK mit n Stufen \Rightarrow max. Konsistenzordnung n

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

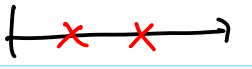
2) RK hat mindestens Konsistenzordnung $q=2$
wenn $\sum_{j=1}^n b_j c_j = \frac{1}{2}$


3) $q=3$ brauchen wir

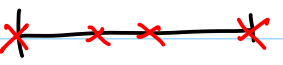
$$\sum_{j=1}^n b_j c_j^2 = \frac{1}{3}, \quad \sum_{j=1}^n b_j \sum_{k=1}^n a_{jk} c_k = \frac{1}{6}$$

usw.

Theorem RK explizit $\Rightarrow q \leq n$

Gauss-Quadratur $\Rightarrow q=2n$ 

Radau-Quadratur 
 \hookrightarrow Radau-Verfahren für ODEs.

Lobatto-Quadratur 
 \hookrightarrow Lobatto-Verfahren für ODEs.

Theorem

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

$$\| \underline{y}(t_i) - \underline{y}_i \| \leq C \cdot h^q \text{ für alle } i=1,2,\dots,n$$

($T = nh$)

§4.2. Kollokation

Def $\tau_1, \tau_2, \dots, \tau_n \in [0, 1]$ verschieden
Kollokationspolynom $u(t)$ von Grad n :

$$\begin{cases} u(t_0) = y_0 & \text{für } i=1, 2, \dots, n \\ \ddot{u}(t_0 + \tau_i h) = f(t_0 + \tau_i h, u(t_0 + \tau_i h)) \end{cases}$$

Bsp 1) $n=1$ Polynom vom Grad 1:

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

mit k so bestimmt dass:

$$\ddot{u}(t_0 + \tau_1 h) = f(t_0 + \tau_1 h, u(t_0 + \tau_1 h))$$

$$\tau_1 = 0 \Rightarrow (LE)$$

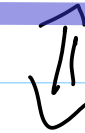
$$\tau_1 = 1 \Rightarrow (RE)$$

$$\tau_1 = \frac{1}{2} \Rightarrow (IMP)$$

2) $n=2$; $\tau_1=0, \tau_2=1 \Rightarrow$ implizite Trapezregel.

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

Theorem Die Kollokation mit k Knoten τ_1, \dots, τ_n



n -Stufiges RKV mit $a_{ij} = \int_0^{\tau_i} l_j(\tau) d\tau$

$$b_i = \int_0^1 l_i(\tau) d\tau$$

wobei

$$l_i(\tau) = \frac{(\tau - \tau_1)(\tau - \tau_2) \dots (\tau - \tau_n)}{(\tau_i - \tau_1)(\tau_i - \tau_2) \dots (\tau_i - \tau_n)} = \prod_{\substack{j=1 \\ j \neq i}}^n \frac{\tau - \tau_j}{\tau_i - \tau_j}$$

Lagrange Polynom.

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

Beweis $k_i = i(t_0 + c_i h)$

$$u(t_0 + \tau h) = \sum_{j=1}^p k_j l_j(\tau) \quad \int_0^{c_i} \Rightarrow$$

hat Grad $p-1$

$$u(t_0 + c_i h) = y_0 + h \sum_{j=1}^p k_j \int_0^{c_i} l_j(\tau) d\tau$$

zwischen Stellen.

a_{ij}

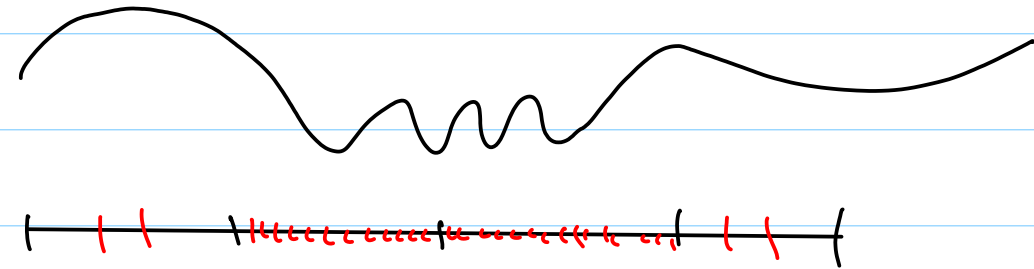
$$\int_0^1 \Rightarrow u(t_0 + h) = y_0 + h \sum_{j=1}^p k_j \int_0^1 l_j(\tau) d\tau$$

y_1

b_j

Konsequenz Kollokationsmethode hat dieselbe Ordnung wie die entsprechende QF .

§4.3. Adaptivität



lokaler Fehler abschätzen:

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

$$u \rightarrow O(h^{p+1})$$

$$\tilde{\Psi}_h \text{ genauer} \rightarrow O(h^{p+2})$$

$$est_k = \left| \tilde{\Psi}^{t, t+h}(y(t_k)) - \Psi^{t, t+h}(y(t_k)) \right| \approx ch^{p+1} = tol$$

$$h^* = h \sqrt[p+1]{\frac{tol}{est_k}}$$

§4.4. Partitionierte RK-Verfahren

System ODE partitioniert:

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

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

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

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

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

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

Bsp 1)

$$\left. \begin{array}{l} (eE) : b_1=1, a_{11}=1 \\ (iE) : b_1=1, \hat{a}_{11}=0 \end{array} \right\} \begin{array}{l} \text{für Newton-Gleichung} \\ \Rightarrow \text{symplektische} \\ \text{Euler-Verfahren.} \end{array}$$

Bsp 2)

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & \frac{1}{2} & \frac{1}{2} \\ \hline 1 & \frac{1}{2} & \frac{1}{2} \end{array}, \quad \begin{array}{c|cc} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ \hline 1 & \frac{1}{2} & \frac{1}{2} \end{array} \Rightarrow \text{PRK für Newton-Gleichung, Störmer-Verlet}$$

Verallgemeinerung vom Störmer-Verlet:

3-stufige Lobatto-Paar

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{5}{24} & \frac{1}{3} & -\frac{1}{24} \\ 1 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \\ \hline 1 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array}, \quad \begin{array}{c|ccc} 0 & \frac{1}{6} & -\frac{1}{6} & 0 \\ \frac{1}{2} & \frac{1}{6} & \frac{1}{3} & 0 \\ 1 & \frac{1}{6} & \frac{5}{6} & 0 \\ \hline 1 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array} \quad O(h^4)$$

Ben Newton! Gleichung.

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

Um schreiben:

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

PRK \Rightarrow RK-Nyström-Verfahren (RK4)

$$\begin{cases} \underline{l}_i = g(t_i, \underline{y}_0 + c_i h \underline{z}_0 + h^2 \sum_{j=1}^n \bar{a}_{ij} \underline{l}_j, \underline{z}_0 + h \sum_{j=1}^n \hat{a}_{ij} \underline{l}_j) \\ \underline{y}_1 = \underline{y}_0 + h \left(\underline{z}_0 + h \sum_{i=1}^n \bar{b}_i \underline{l}_i \right) \text{ mit } \bar{b}_i = \sum_{k=1}^n b_k \hat{a}_{ki} \\ \underline{z}_1 = \underline{z}_0 + h \sum_{i=1}^n \hat{b}_i \underline{l}_i \end{cases}$$

$$\bar{a}_{ij} = \sum_{k=1}^n a_{ik} \hat{a}_{kj}$$

Wenn g nicht von \dot{y} abhängt \Rightarrow braucht man \hat{a}_{kj} nicht.

Ben PRK = 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 klar!

einfacher anzuwenden!

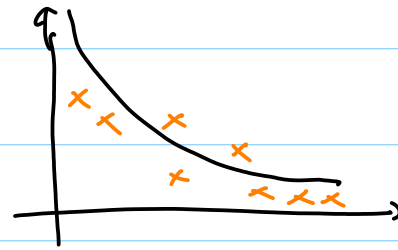
BM42 mit $O(h^4)$

BM63 mit $O(h^6)$

} symplektische Verfahren.

§5 Steife Differentialgleichungen

§5.1. Einführung.



Modelproblem: $\dot{y} = -\lambda y$ mit $\lambda > 0$
 $y(t) = e^{-\lambda t} y_0 \rightarrow 0$ für $t \rightarrow \infty$

Interesse: asymptotisches Verhalten
 der numerischen Lösung sollte
 qualitativ (zumindest) ähnlich der
 exakten Lösung sein.

(QE):

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

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

...

$$y_N = (1 - \lambda h)^N y_0$$

$$|y_N| \rightarrow 0 \text{ nur wenn } |1 - \lambda h| < 1 \Leftrightarrow 0 < h < \frac{2}{\lambda}$$

Def ODE heißt steif, falls explizite Verfahren
 einen Zeitschritt h sehr klein brauchen, kleiner
 als die Genauigkeit verlangt!

$$\text{Bsp } \frac{d}{dt} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \underbrace{\begin{bmatrix} -50 & 49 \\ 49 & -50 \end{bmatrix}}_{\underline{\underline{B}}} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$\dot{\underline{y}} = \underline{\underline{B}} \underline{y}$$

$\underline{\underline{B}}$ symmetrisch \Rightarrow es gibt $\underline{\underline{S}}$ (mit $\underline{\underline{S}} \underline{\underline{S}}^T = \underline{\underline{I}}$)
 dass

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

mit $\underline{\underline{D}} = \text{Diagonalmatrix}$.

$$\dot{\underline{y}} = \underline{\underline{S}} \underline{\underline{D}} \underline{\underline{S}}^T \underline{y} \Rightarrow \underline{\underline{S}}^T \dot{\underline{y}} = \underline{\underline{D}} \underbrace{\underline{\underline{S}}^T \underline{y}}_{\underline{z}} \Rightarrow \dot{\underline{z}} = \underline{\underline{D}} \underline{z}$$

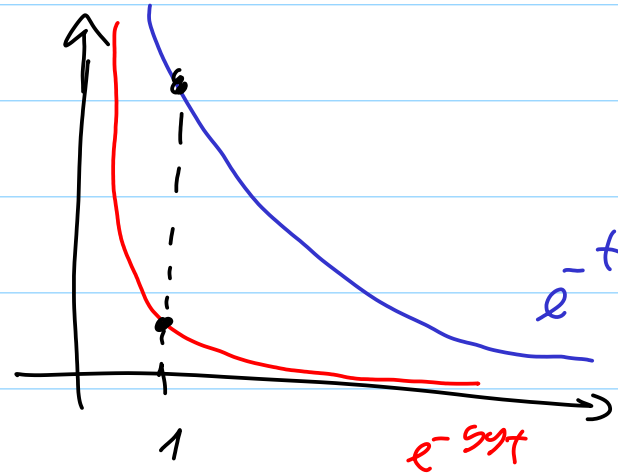
$$\left\{ \begin{array}{l} \dot{z}_1 = \lambda_1 z_1 \\ \vdots \\ \dot{z}_d = \lambda_d z_d \end{array} \right.$$

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

$$\dot{z}_1 = -z_1 \Rightarrow h < 2$$

$$\dot{z}_2 = -99z_2 \Rightarrow h < \frac{2}{99}$$

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



zur Zeit $t=1$: $y_1(1) = e^{-1} + e^{-99}$

sehr klein, irrelevant!

aber ein. expl. Verfahren will ein kleines h

\hookrightarrow bestimmt durch e^{-99t}

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

$$(1 + \lambda h) y_1 = y_0 \quad \left\{ \begin{array}{l} h > 0, \lambda > 0 \end{array} \right. \Rightarrow y_1 = \frac{1}{1 + \lambda h} y_0 \Rightarrow y_N = \left(\frac{1}{1 + \lambda h} \right)^N y_0 \xrightarrow{N \rightarrow \infty} 0$$

Bei implizite Verfahren stellen keine Bedingung an h .

Bsp explizite Trapezregel:

$$(j = -\lambda y)$$

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

$$k_2 = -\lambda y_0 - \lambda^2 h y_0$$

$$y_1 = y_0 + \frac{1}{2} h k_1 + \frac{1}{2} h k_2 = y_0 + (-\lambda h) y_0 + \frac{(\lambda h)^2}{2} y_0$$

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

$$y_N = (S(\lambda h))^N y_0 \rightarrow 0 \text{ nur wenn } |S(\lambda h)| < 1$$

§ 5.2. Stabilität des RK-Verfahrens

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

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

Falls $\lambda = \sigma i \Rightarrow y(t) = e^{\sigma t} y_0 = y_0 (\cos t + i \sin t)$

$$\lambda = -1 + \sigma i \Rightarrow y(t) = e^{-t} y_0 (\cos t + i \sin t) \xrightarrow[t \rightarrow \infty]{} 0$$

Numerisches Verfahren $y_N = S(\lambda h)^N y_0$

Frage: wann $|y_N| \rightarrow 0$ für $N \rightarrow \infty$?

$S(z)$ = Stabilitätsfunktion.

(eE) $S(z) = 1 + z$ (da wir jetzt $\dot{y} = +\lambda y$ notierten, in § 5.1 hatte ich $\dot{y} = -\lambda y$)

(iE) $S(z) = \frac{1}{1 - z}$

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

RK-Verfahren mit ρ Stufen:

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

$$\begin{cases} k_i = f(t_0 + c_i h, y_0 + h \sum_{j=1}^{\rho} a_{ij} k_j) & \text{für } \dot{y} = \lambda y \\ i = 1, 2, \dots, \rho \end{cases}$$

$$\Rightarrow \begin{cases} k_i = \lambda y_0 + \underbrace{\lambda h}_{z} \sum_{j=1}^{\rho} a_{ij} k_j \\ i = 1, 2, \dots, \rho \end{cases} \quad \underline{k} = \begin{bmatrix} k_1 \\ \vdots \\ k_\rho \end{bmatrix}$$

$$\Rightarrow \underline{k} = \lambda y_0 \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} + z \underline{A} \underline{k} \quad \Leftrightarrow$$

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

$\underline{I} - z \underline{A}$ nicht invertierbar \rightarrow

I - z A invertierbar \Rightarrow

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

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

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

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

Stabilitätsfunktion des n-stufiges RKV.

$$y_1 = S(z) y_0 \quad \text{mit } z = \lambda h$$

$|y_N| \rightarrow 0$ für $N \rightarrow \infty$ nur wenn $|S(z)| < 1$

$|y_N| = |y_0|$ falls $|S(z)| = 1$

Theorem Die Stabilitätsfunktion eines n-stufiges RKV ist eine (komplexwertige) rationale Funktion

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

und $Q(z) = 0$ für $z = \frac{1}{\mu}$ mit μ Eigenwert von A.
Falls RKV explizit, dann $Q(z) \equiv 1$.

Beweis explizite RKV: $\underline{A} = \begin{bmatrix} 0 & \dots & 0 \\ * & \ddots & 0 \end{bmatrix}$

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

$$(\underline{I} - z \underline{A})^{-1} = \underline{I} + z \underline{A} + (z \underline{A})^2 + (z \underline{A})^3 + \dots + \underbrace{(z \underline{A})^n}_0 + \underbrace{\dots}_0$$

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

= Polynom vom Grad $n-1$ in z .

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

Cramer Regel.

$$Q(z) = 0 \Leftrightarrow \det(\underline{I} - z \underline{A}) = 0 \Leftrightarrow z = \frac{1}{\mu} \text{ mit } \mu \in W \text{ von } \underline{A}.$$

Konsequenz

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

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

mit wohldefinierten Stabilitätsfunktion $S(z)$.

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

num. Lösung $y_n = S(\lambda h)^n$

RKV hat Konvergenzordnung q :

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

Somit ist $S(z)$ eine Approximation an e^z

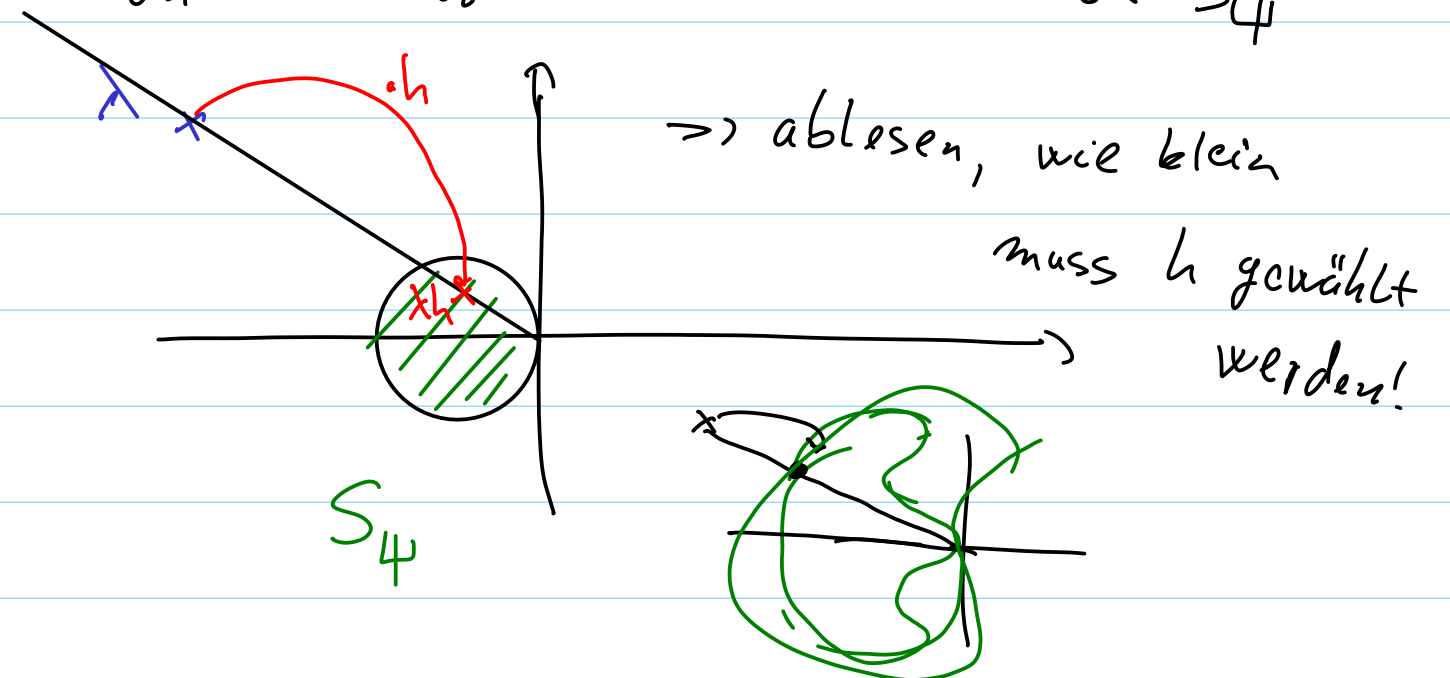
$$|e^z - S(z)| \leq O(|z|^{q+1})$$

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

Definition Stabilitätsgebiet des num. Verfahrens ψ .

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

$$y_n = S(z)^n y_0 \rightarrow 0 \text{ nur wenn } z \in S_\psi$$



Beh RK explizit $\Rightarrow S(z) = \text{Polynom}$ ist \Rightarrow
 Stabilitätsgebiet beschränkt
 \Rightarrow immer eine Schritke an h .

Beh ode45 / dopri5 verwenden explizite RK
 $\Rightarrow S_{\psi}$ beschränkt \Rightarrow brauchen
 kleines h .

Def Ein Verfahren heißt A-stabil falls
 $\{z \in \mathbb{C}; \operatorname{Re} z < 0\} \subset S_{\psi}$

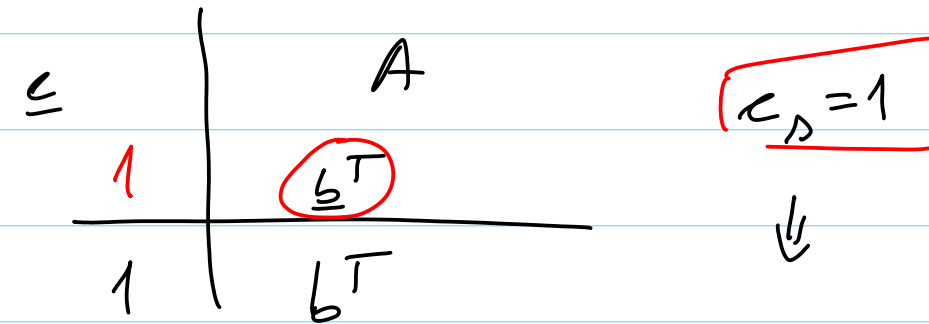
Bsp (iE) ist A-stabil

RK-Gauss Verfahren sind A-stabil.
 (z.B. iMP)

Beh $S(z) \approx e^z$; $z \rightarrow -\infty \Rightarrow e^z \rightarrow 0$
 Frage $S(-\infty) = 0$?

Def Num. Verfahren heißt L-stabil falls
 $\lim_{z \rightarrow -\infty} S(z) = 0$

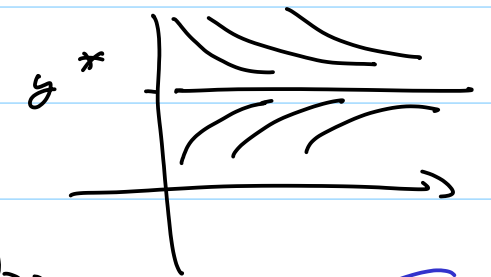
Beh L-stabil falls $\underline{b}^T = a_{\Delta} = \text{letzte Zeile in } \underline{A}$



\Rightarrow Radau-Verfahren (basieren auf Radau-Quadratur)
 L-stabil und höchste Ordnung haben.
 (2D-1)

Radau-Verfahren von Ordnungen 3, 5 \rightarrow Script

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



$\underline{y}^* \in \mathbb{R}^d$ heißt stationär: $\underline{f}(\underline{y}^*) = 0$

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

$\dot{\underline{z}} = \underline{Df}(\underline{y}^*) \underline{z}$ Testproblem!



Ben implizite Rk: in jedem Zeitschritt
ein nicht-lineares algebraisches System
mit d.S. Gleichungen lösen
(n Stufen $\underline{k}_1, \dots, \underline{k}_n \in \mathbb{R}^d$)

TEUER mit fsolve
 $F(\underline{x}) = 0$.

Ben Autonome ODE, in Rk:

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

Idee: linearisiere: \searrow ersetzt durch,

$$\begin{array}{ccccccc} \underline{k}_i & = & \underline{f}(\underline{y}_0) & + & h \underline{D} \underline{f}(\underline{y}_0) & \left(\sum_{j=1}^n a_{ij} \underline{k}_j \right) \\ \uparrow & & \uparrow & & \uparrow & & \underbrace{\sum_{j=1}^n a_{ij} \underline{k}_j}_{\in \mathbb{R}^d} \\ \mathbb{R}^d & & \mathbb{R}^d & & \mathbb{R}^{d \times d} & & \end{array}$$

ein grosses LGS $\mathbb{R}^{dn \times dn}$.

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

Skript: man verliert eine Konvergenzordnung!

(Newton: Ordnung 2)

(linearisierte Rk: Ordnung 1)

$$\text{Radau} \xrightarrow{\text{mit } n=2} + \text{Newton} \Rightarrow O(h^3)$$

$$\text{Radau } (n=2) + \text{linearisierung} \Rightarrow O(h^2)$$

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

$$\text{Bsp } (i \in) \quad \underline{y}_1 = \underline{y}_0 + h \underline{f}(\underline{y}_1)$$

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

ein Schritt Newton für \underline{F} mit Startwert $\underline{z}_0 = \underline{y}_0$

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

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

Wenn wir aber einen besseren Startpunkt für Newton haben, dann vielleicht auf bessere Konvergenz.

Ben Ordnung retten geht bei diagonal-impliziten Verfahren.

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

$$\begin{cases} \underline{k}_i = f\left(\underline{y}_0 + h \sum_{j=1}^i a_{ij} \underline{k}_j\right) \\ i=1,2,\dots,d \end{cases}$$

$$\underline{F}(\underline{k}) = \underline{k} - f\left(\underline{y}_0 + \underline{z} + h a_{ii} \underline{k}\right) \quad \underline{k} \in \mathbb{R}^d$$

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

Newton-Schritt:

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

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.

$$\left(\underline{I} - h a_{ii} \underline{\partial} \right) \underline{k}_i = f\left(\underline{y}_0 + h \sum_{j=0}^{i-1} (a_{ij} + d_{ij}) \underline{k}_j\right) - h \underline{\partial} \sum_{j=1}^{i-1} d_{ij} \underline{k}_j$$

mit $\underline{\underline{\partial}} = \underline{\underline{D}} f \left(\underline{\underline{y}}_0 + h \underbrace{\sum_{j=1}^{i-1} (a_{ij} + d_{ij}) \underline{\underline{k}}_j}_{\substack{\text{wie bei vereinfachten Newton} \\ \text{können wir darauf verzichten} \\ \Rightarrow \text{günstiger.}}} \right)$

\Rightarrow Row 2 > adaptiven impliziten Verfahren.
Row 3
ode23s

§6 Nichtlineare algebraische Gleichungen

(C) V. Grädnaru

§6.1. Einführung

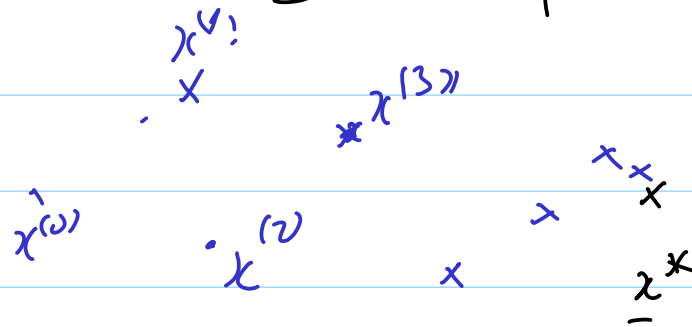
Finde $\underline{x}^* \in \mathbb{R}^d$ so dass $\underline{F}(\underline{x}^*) = 0$.

Bsp

$$d=1 \quad F(x) = x e^x - 1$$

Man baut eine Folge von Approximationen an \underline{x}^*

$$\begin{array}{c} \underline{x}^{(0)}, \underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(k)} \rightarrow \underline{x}^* \\ \uparrow \\ \text{Gegeben} \end{array} \quad \underline{x}^{(k+1)} = \Phi(\underline{x}^{(k)}) \quad \text{Iteration}$$



Def $\underline{x}^{(k+1)} = \Phi(\underline{x}^{(k)})$ heisst linear konvergent nach \underline{x}^* falls es gibt $L < 1$ so dass

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

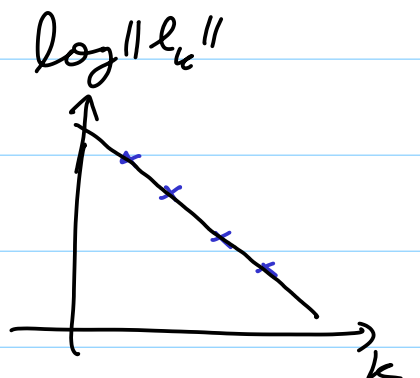
$$\begin{aligned} \underline{\text{Bem}} \quad \|\underline{x}^{(k+1)} - \underline{x}^*\| &\leq L \|\underline{x}^{(k)} - \underline{x}^*\| \leq \\ &\leq L^2 \|\underline{x}^{(k-1)} - \underline{x}^*\| \leq \dots \\ &\leq L^{k+1} \|\underline{x}^{(0)} - \underline{x}^*\| \end{aligned} \quad \left\{ \begin{array}{l} \Rightarrow \text{Konvergenz} \\ \underline{x}^{(k+1)} \rightarrow \underline{x}^* \text{ in } \|\cdot\| \end{array} \right.$$

$$0 < L < 1$$

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

$$\log \|\underline{e}_k\| \leq k \log L + \log \|\underline{e}_0\|$$

Gerade mit Steigung $\log L$.

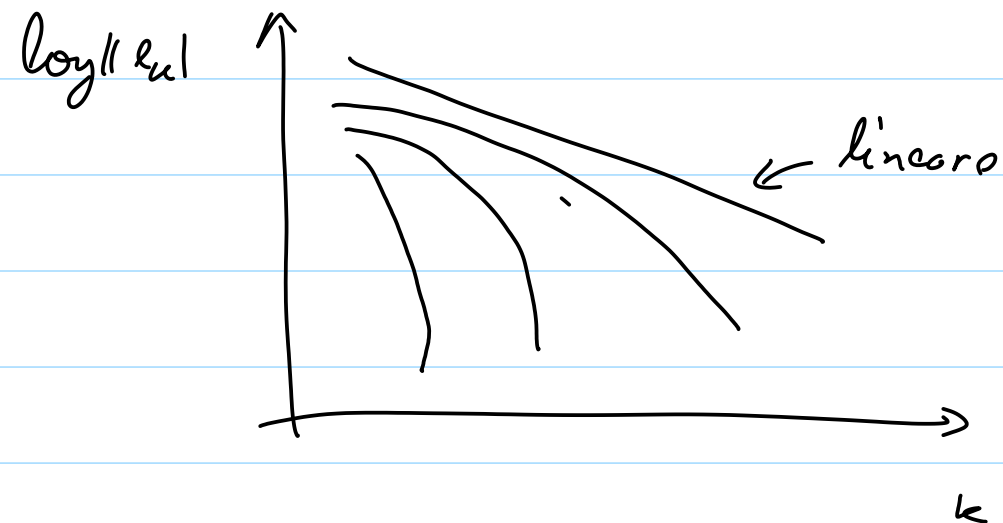


Def Konvergenz der Ordnung p des iteratives
Verfahrens,

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

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

Falls (\underline{x}^k) konvergent \Rightarrow \underline{x}^* $\Rightarrow \underline{\Phi}(\underline{x}^*) = \underline{x}^*$ Lösung

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

$$\Phi_1(x) = e^{-x}$$

Starte mit $x^{(0)}$, $x^{(k+1)} = \Phi_1(x^{(k)})$

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

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

$$\Phi_2(x) = \frac{x+1}{e^x + 1}$$

Starte mit $x^{(0)}$, $x^{(k+1)} = \Phi_2(x^{(k)})$

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

$$\Phi_3(x) = x + 1 - x e^x$$

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

↓
quadratische

↓
keine

k	$ x_1^{(k+1)} - x^* $	$ x_2^{(k+1)} - x^* $	$ x_3^{(k+1)} - x^* $
0	0.067143290409784	0.067143290409784	0.067143290409784
1	0.039387369302849	0.00832287212566	0.108496074240152
2	0.021904078517179	0.00000125374922	0.219330611898582
3	0.012559804468284	0.000000000000003	0.288178118764323
4	0.007078662470882	0.000000000000000	0.723649245792953
5	0.004028858567431	0.000000000000000	0.410183132337935
6	0.002280343429460	0.000000000000000	1.186907542305364
7	0.001294757160282	0.000000000000000	0.146569797006362
8	0.000733837662863	0.000000000000000	0.310516641279937
9	0.000416343852458	0.000000000000000	0.357777386500765
10	0.000236077474313	0.000000000000000	0.974565695952037

Ben Hinreichende Bedingung für lokale lineare Konvergenz

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

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

Wenn $\underline{\phi}(\underline{x}^*) = \underline{x}^*$ für $\underline{x}^* \in U$, dann

konvergiert $\underline{x}^{(k+1)} = \underline{\phi}(\underline{x}^{(k)})$ gegen \underline{x}^* lokal
mindestens linear.

Bez $\phi: U \subset \mathbb{R} \rightarrow \mathbb{R}$ ϕ $(n+1)$ -mal stetig diff^{bar}
 $\phi(x^*) = x^* \in U$

Taylor: $\phi(y) = \phi(x) + \sum_{k=1}^n \frac{1}{k!} \phi^{(k)}(x) (y-x)^k + O(|y-x|^{n+1})$

Theorem Falls $\phi^{(l)}(x^*) = 0$ für $l=1,2,\dots,n \geq 1$
dann konvergiert die Fixpunktiteration

$\underline{x}^{(k+1)} = \underline{\phi}(\underline{x}^{(k)})$ gegen x^* lokal

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

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

$$x^{k+1} - x^* = \phi(x^{(k)}) - \phi(x^*) = \sum_{j=1}^n \frac{1}{j!} \underbrace{\phi^{(j)}(x^*)}_{=0} (x^{(k)} - x^*)^j + O(|x^{(k)} - x^*|^{n+1})$$

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

Bsp $\phi_2(x) = \frac{x+1}{e^{x+1}}$

$$\phi_2'(x) = \frac{e^x + 1 - (x+1)e^x}{(e^{x+1})^2} = \frac{1 - xe^x}{(e^{x+1})^2} \Rightarrow$$

$$\phi_2'(x^*) = \frac{1 - x^* e^{x^*}}{(e^{x^*+1})^2} = 0 \quad \text{für } x^* \text{ die Nullstelle von } xe^x - 1.$$

$\Rightarrow \phi_2$ gibt mindestens Konvergenzordnung 2.

§ 6.3. Abbruchkriterium

Wann brechen wir eine Iteration ab?

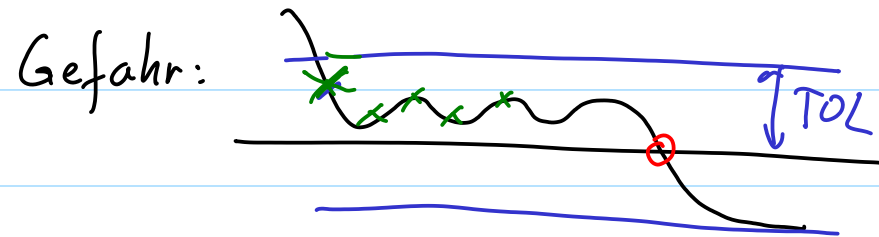
1) noch k_{\max} Schritte: "a priori" Abbruchkriterium.

2) "a posteriori" ideal: $\|\underline{x}^{(k)} - \underline{x}^*\| \leq \text{TOL}$ ☹️

$$\|\underline{x}^{(k+1)} - \underline{x}^{(k)}\| < \text{TOL} \Rightarrow \text{Abbruch}$$

oder verwende Residuum

$$\|F(\underline{x}^{(k)})\| < \text{TOL} \Rightarrow \text{Abbruch}$$



Wenn wir Informationen über die Konvergenzrate der Iteration, dann bessere Kriterien!

Bsp lineare Konvergenz mit der Rate L

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

\uparrow
 $\pm \underline{x}^{(k+1)}$

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

$$\|\underline{x}^{(k+1)} - \underline{x}^*\| \leq \frac{L}{1-L} \|\underline{x}^{(k)} - \underline{x}^{(k+1)}\| \quad \text{da } 0 < L < 1$$

\Rightarrow Abbruchkriterium:

$$\|\underline{x}^{(k)} - \underline{x}^{(k+1)}\| \leq \frac{1-L}{L} \cdot \text{TOL} \Rightarrow \|\underline{x}^{(k+1)} - \underline{x}^*\| \leq \text{TOL}$$

Abbruch!

Falls L unbekannt, verwende ein geschätztes $\tilde{L} > L$

Bezug) Bei der Berechnung Fehler $\varepsilon_k = \|\underline{x}^{(k)} - \underline{x}^*\|$

$$\varepsilon_{k+1} \approx L \varepsilon_k \Rightarrow \log \varepsilon_{k+1} \approx \log L + \log \varepsilon_k$$

$$\Rightarrow \log L \approx \frac{\log \varepsilon_{k+1}}{\log \varepsilon_k}$$

oder.

$$\approx \frac{\|\underline{x}^{(k)} - \underline{x}^{(N)}\|}{\|\underline{x}^{(k-1)} - \underline{x}^{(N)}\|}$$

2) Konvergenzordnung p aus Experiment erraten:

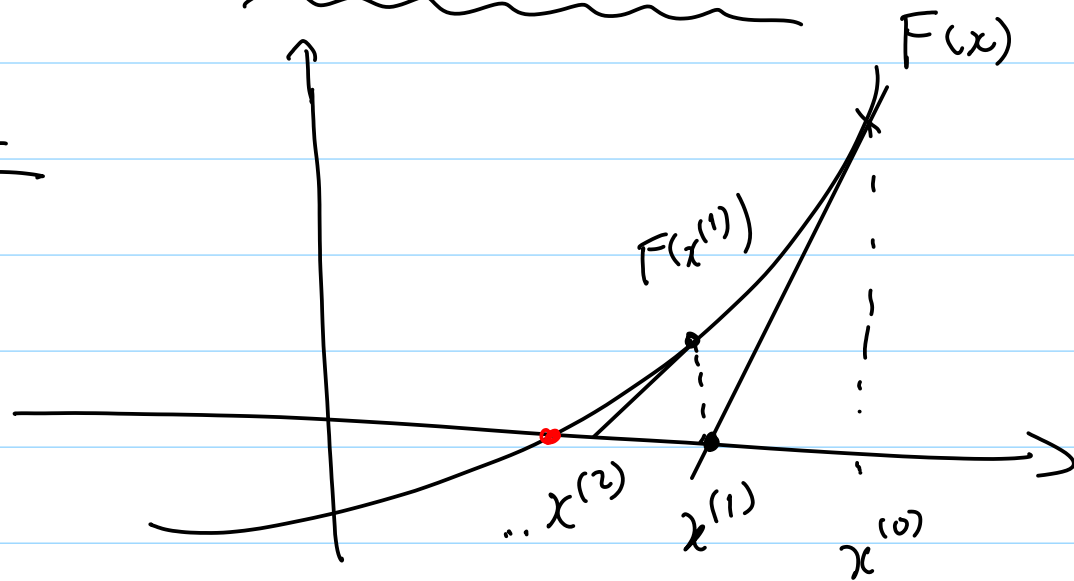
$$\left. \begin{aligned} \varepsilon_{k+1} &\approx C \varepsilon_k^p \Rightarrow \log \varepsilon_{k+1} \approx \log C + p \log \varepsilon_k \\ \log \varepsilon_k &\approx \log C + p \log \varepsilon_{k-1} \end{aligned} \right\} \Rightarrow$$

$$\Rightarrow \log \varepsilon_{k+1} - \log \varepsilon_k \approx p (\log \varepsilon_k - \log \varepsilon_{k-1})$$

$$\Rightarrow p \approx \dots$$

§ 6.4 Newton-Verfahren

Idee



Ersetze die komplizierte Funktion F durch eine einfachere, z.B. durch eine lineare Funktion (\Rightarrow Linearisierung!)

Linearisiere lokal! \Rightarrow Tangente + Neustart

$F \approx \tilde{F}$ mit $\tilde{F}(x)$ einfacher

Newton

$$\tilde{F}(x) = F(x^{(0)}) + \underline{DF}(x^{(0)})(x - x^{(0)})$$

Statt $F(x)=0$ löse $\tilde{F}(x)=0 \Leftrightarrow$

$$F(x^{(0)}) + \underline{DF}(x^{(0)})(x - x^{(0)}) = 0 \quad (\Rightarrow)$$

$\{ \text{LGS in unbekannte } x$

$$\text{LGS: } \underline{DF}(x^{(0)}) \underline{x}^{(1)} = - \underline{F}(x^{(0)}) + \underline{DF}(x^{(0)}) \underline{x}^{(0)}$$

Und starte von diesem $x^{(1)}$ wieder \Rightarrow

Newton-Iteration: gegeben $\underline{x}^{(0)}$

für $k=0,1,2,\dots$:

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \underline{DF}(\underline{x}^{(k)})^{-1} \underline{F}(\underline{x}^{(k)})$$

Bem

Berechne niemals die Inverse einer Matrix.

Sondern löse nur LGS

→ Gauss-Elimination; LU-Zerlegung

→ Cholesky-Zerlegung; QR-Zerlegung

* sympy + lambdify

* Wohl vom Startwert + gedämpftes Newton-Verfahren.

* vereinfachtes Newton-Verfahren.

Theorem

Newton-Verfahren konvergiert mit Ordnung

$p=2$ (Lokal).

Bem In jedem Schritt muss man ein LGS

lösen:

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} - \underline{D} f(\underline{x}^{(k)})^{-1} f(\underline{x}^{(k)})$$

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

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

Beweis $d=1 \rightarrow \underline{x}^{(k+1)} = \underline{x}^{(k)} - f'(\underline{x}^{(k)})^{-1} f(\underline{x}^{(k)})$

Das ist eine Fixpunktiteration

$$\underline{x}^* = \Phi(\underline{x}^*) \text{ mit } \Phi(\underline{x}) = \underline{x} - f'(\underline{x})^{-1} f(\underline{x})$$

Newton Iteration für $f(\underline{x}^*) = 0$

$$\Phi'(\underline{x}) = 1 - \frac{f'(\underline{x}) f'(\underline{x}) - f(\underline{x}) f''(\underline{x})}{(f'(\underline{x}))^2} = 1 - 1 + \frac{f(\underline{x}) f''(\underline{x})}{(f'(\underline{x}))^2} \Rightarrow$$

$$\Rightarrow \Phi'(\underline{x}^*) = \frac{f(\underline{x}^*) f''(\underline{x}^*)}{(f'(\underline{x}^*))^2} = 0 \quad \text{falls } f'(\underline{x}^*) \neq 0. \Rightarrow p=2.$$

$$\underline{D} f(\underline{x}^{(k)}) \underline{\rho}^{(k)} = f(\underline{x}^{(k)}) \quad f(\underline{x}^*) = 0$$

$$\underline{D} f(\underline{x}^{(k)}) \in \mathbb{R}^{d \times d} \rightarrow \text{teuer nach } k \text{ Schritten}$$

$$\underline{\partial}(\underline{x}^{(k)}) = \underline{D} f(\underline{x}^{(k)})$$

$$\underline{\partial} = \underline{\cap} \cup \quad \text{kostet } O(d^3) \text{ Schritte}$$

Idee: Wiederverwendung von \underline{J} für mehrere Newton-Schritte

$\underline{J} = \underline{L} \underline{U}$ ein mal, dann verwende die Faktoren \underline{L} und \underline{U} auch für die nächsten Schritte

(0) $k = 1, 2, 3, \dots$

Vereinfachte Newton:

$$\underline{J} = \underline{D} f(\underline{x}^{(0)}) \quad \rightarrow O(d^3)$$

$$\underline{J} = \underline{L} \underline{U} \quad (\underline{L}, \underline{U} = \text{Lu-factor}(\underline{J}))$$

für $k = 0, 1, 2, \dots$

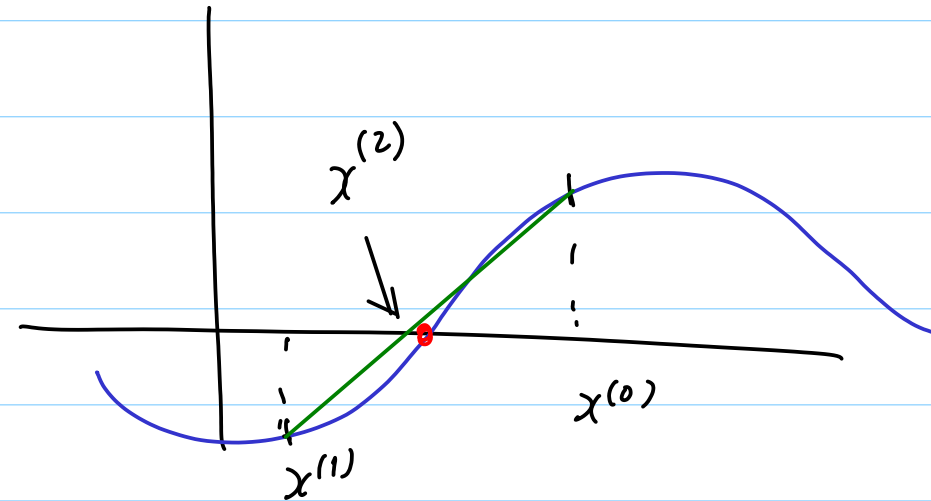
$$\text{löse } \underline{L} \underline{U} \underline{\Delta} = f(\underline{x}^{(k)}) \quad \text{nur } O(d^2)$$

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

Typischerweise bekommt man nur lineare Konvergenz

Bem Was tun wenn $\underline{D} f(\underline{x})$ nicht bekannt?

1D



Idee: verwende statt Tangente eine andere Gerade, z.B. eine Sekante.

Benachrichtigung jetzt 2 Startwerte: $x^{(0)}, x^{(1)}$

$$f'(x^{(k)}) \approx \frac{f(x^{(k)}) - f(x^{(k-1)})}{x^{(k)} - x^{(k-1)}}$$

update: $x^{(k+1)} = x^{(k)} - \Delta$ mit

$$\Delta = \frac{x^{(k)} - x^{(k-1)}}{f(x^{(k)}) - f(x^{(k-1)})} f(x^{(k)})$$

Allgemein: d Dimensionen:

$$\underline{J}_k (\underline{x}^{(k)} - \underline{x}^{(k-1)}) = f(\underline{x}^{(k)}) - f(\underline{x}^{(k-1)})$$

$\swarrow \quad \searrow$
 $d \times d \quad \mathbb{R}^d$

Bez Viele mögliche \underline{J}_k erfüllen das!

Broyden: baue \underline{J}_k iterativ

$$\underline{J}_k = \underline{J}_{k-1} + \frac{1}{\|\underline{x}^{(k)} - \underline{x}^{(k-1)}\|_2^2} \underbrace{f(\underline{x}^{(k)}) (\underline{x}^{(k)} - \underline{x}^{(k-1)})^T}_{\text{Matrix}}$$

Vorteil: man muss f nur ein Mal auswerten!

Broyden-Verfahren:

$\underline{x}^{(0)}$ gegeben

$$\underline{J}_0 = \underline{D}f(\underline{x}^{(0)})$$

für $k = 0, 1, 2, \dots$:

löse $\underline{J}_k \underline{\Delta} = f(\underline{x}^{(k)})$ teuer

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

$$\underline{J}_{k+1} = \underline{J}_k + \frac{1}{\|\underline{\Delta}\|_2^2} f(\underline{x}^{(k+1)}) (-\underline{\Delta})^T$$

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-Morrison-Formel

$$\underline{\underline{\partial}}_{k+1} = \underline{\underline{\partial}}_k + \text{Rang-1-Matrix.}$$

$$\left(\underline{\underline{A}} + \underline{\underline{u}} \underline{\underline{v}}^T \right)^{-1} = \underline{\underline{A}}^{-1} - \frac{\underline{\underline{A}}^{-1} \underline{\underline{u}} \underline{\underline{v}}^T \underline{\underline{A}}^{-1}}{1 + \underline{\underline{v}}^T \underline{\underline{A}}^{-1} \underline{\underline{u}}}$$

(einfach überprüfen!)

Rang-1-update braucht nur $\underline{\underline{A}}^{-1}$

In jedem Schritt:

$$\underline{\underline{\partial}}_{k+1}^{-1} = \underline{\underline{\partial}}_k^{-1} + \frac{\underline{\underline{\partial}}_k^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(k+1)}) \underline{\underline{\rho}}^T \underline{\underline{\partial}}_k^{-1}}{\|\underline{\underline{\rho}}\|^2 - \underline{\underline{\rho}}^T \underline{\underline{\partial}}_k^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(k+1)})}$$

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 $\underline{\underline{L}} \underline{\underline{U}} \underline{\underline{\rho}}^{(0)} = \underline{\underline{f}}(\underline{\underline{x}}^{(0)})$

$$\underline{\underline{x}}^{(1)} = \underline{\underline{x}}^{(0)} - \underline{\underline{\rho}}^{(0)}$$

$$\underline{\underline{\partial}}_1^{-1} = \underline{\underline{\partial}}_0^{-1} + \frac{\underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)}) \underline{\underline{\rho}}^{(0)T} \underline{\underline{\partial}}_0^{-1}}{\|\underline{\underline{\rho}}^{(0)}\|^2 - \underline{\underline{\rho}}^{(0)T} \underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})}$$

$$\underline{\underline{x}}^{(2)} = \underline{\underline{x}}^{(1)} - \underline{\underline{\partial}}_1^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})$$

$$\underline{\underline{\rho}}^{(1)} = \underline{\underline{\partial}}_1^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)}) = \underline{\underline{\rho}}^{(1)}$$

$$= \boxed{\underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})} + \frac{\boxed{\underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})} \underline{\underline{\rho}}^{(0)T} \boxed{\underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})}}{\|\underline{\underline{\rho}}^{(0)}\|^2 - \underline{\underline{\rho}}^{(0)T} \boxed{\underline{\underline{\partial}}_0^{-1} \underline{\underline{f}}(\underline{\underline{x}}^{(1)})}}$$

$\underline{w} = \underline{\partial}_0^{-1} f(\underline{x}^{(1)})$ berechnet man als Lösung
 von $\underline{\partial}_0 \underline{w} = f(\underline{x}^{(1)})$ [ZGS]

$$z = \underline{\partial}^{(0)T} \underline{w} \in \mathbb{R}$$

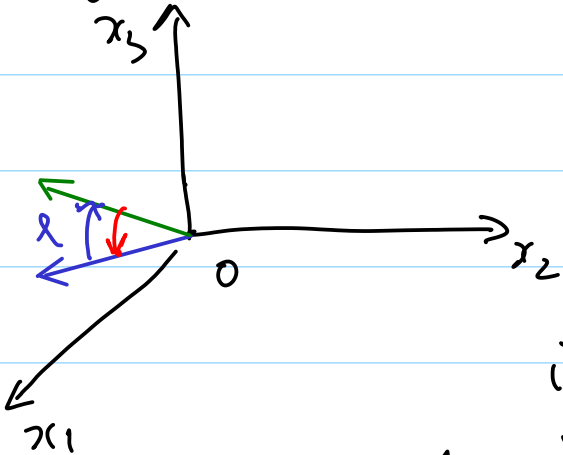
$$\underline{\partial}^{(1)} = \underline{w} + \frac{\underline{w} z}{\|\underline{\partial}^{(0)}\|^2 - z} = \left(1 + \frac{z}{\|\underline{\partial}^{(0)}\|^2 - z}\right) \underline{w}$$

$$\underline{x}^{(2)} = \underline{x}^{(1)} - \underline{\partial}^{(1)}$$

* code

§7. Intermezzo über Lineare Algebra

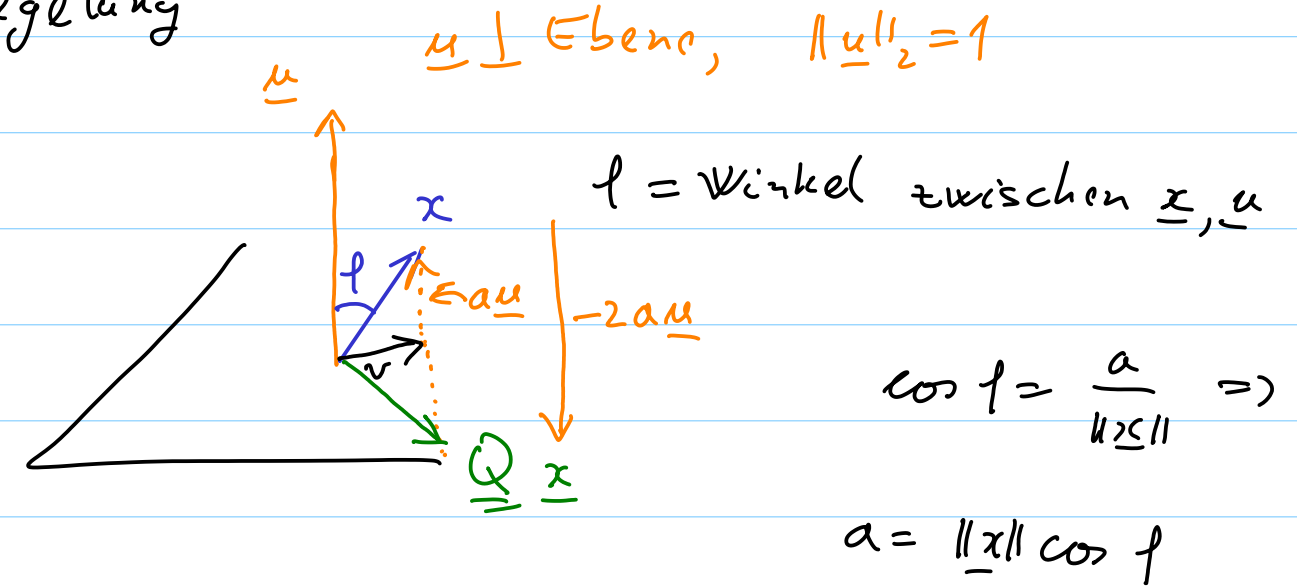
§7.1 Orthogonale Matrizen

Drehung in der x_1, x_3 - Ebene

$$D(\ell) = \begin{bmatrix} \cos \ell & 0 & +\sin \ell \\ 0 & 1 & 0 \\ -\sin \ell & 0 & \cos \ell \end{bmatrix}$$

$$D_{ij}(-\ell) = \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \cos \ell & +\sin \ell \\ & & -\sin \ell & \cos \ell \\ & & & & 1 \end{bmatrix} \begin{matrix} \leftarrow i \\ \leftarrow j \end{matrix}$$

Spiegelung



$$\underline{Q}\underline{x} = \underline{x} + (-2a\underline{u}) = \underline{x} - 2a\underline{u}$$

$$\underline{u}^T \underline{x} = \langle \underline{u}, \underline{x} \rangle = \|\underline{u}\| \|\underline{x}\| \cos \ell = \|\underline{u}\| \cdot a = a$$

$$\begin{aligned} \underline{Q}\underline{x} &= \underline{x} - 2(\underline{u}^T \underline{x}) \underline{u} = \underline{x} - 2\underline{u}(\underline{u}^T \underline{x}) = \underline{x} - 2(\underline{u}\underline{u}^T) \underline{x} \\ &= (\underline{I} - 2\underbrace{\underline{u}\underline{u}^T}) \underline{x} \end{aligned}$$

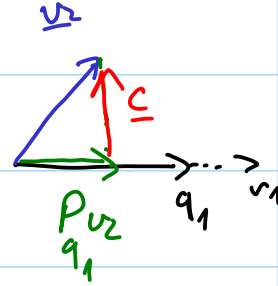
$$\underline{Q} = \underline{I} - 2\underbrace{\underline{u}\underline{u}^T}_{\underline{P}_{\underline{u}}} = \underline{I} - 2\underbrace{\underline{P}_{\underline{u}}}_{\text{Householdermatrix}}$$

§7.2 QR-Zerlegung

(I) via (modifizierten) Gram-Schmidt

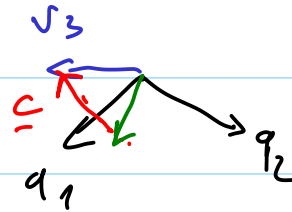
$$\underline{A} = [\underline{v}_1 \ \underline{v}_2 \ \dots \ \underline{v}_n] \rightsquigarrow \underline{q}_1 \dots \underline{q}_n \text{ ONB}$$

$$\underline{q}_1 = \frac{1}{\|\underline{v}_1\|} \underline{v}_1 \Rightarrow \underline{v}_1 = \|\underline{v}_1\| \underline{q}_1 = \langle \underline{q}_1, \underline{v}_1 \rangle \underline{q}_1$$



$$\underline{c}_2 = \underline{v}_2 - P_{\underline{q}_1} \underline{v}_2 = \underline{v}_2 - \langle \underline{q}_1, \underline{v}_2 \rangle \underline{q}_1$$

$$\underline{q}_2 = \frac{1}{\|\underline{c}_2\|} \underline{c}_2$$



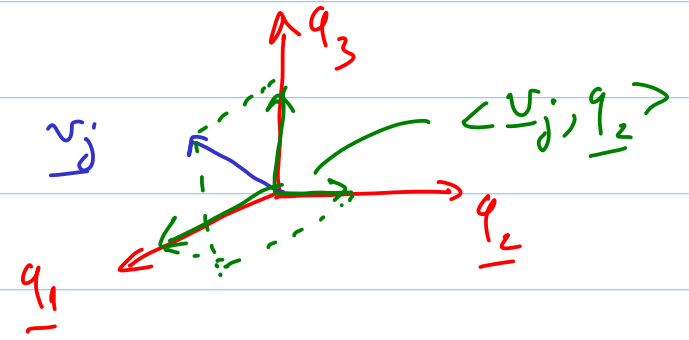
$$\underline{c}_3 = \underline{v}_3 - P_{\text{span}\{\underline{q}_1, \underline{q}_2\}} \underline{v}_3 = \underline{v}_3 - \langle \underline{q}_1, \underline{v}_3 \rangle \underline{q}_1 - \langle \underline{q}_2, \underline{v}_3 \rangle \underline{q}_2$$

$$\underline{q}_3 = \frac{1}{\|\underline{c}_3\|} \underline{c}_3$$

...

$$\underline{c}_k = \underline{v}_k - P_{\text{span}\{\underline{q}_1, \dots, \underline{q}_{k-1}\}} \underline{v}_k = \underline{v}_k - \langle \underline{q}_1, \underline{v}_k \rangle \underline{q}_1 - \dots - \langle \underline{q}_{k-1}, \underline{v}_k \rangle \underline{q}_{k-1}$$

$$\underline{q}_k = \frac{1}{\|\underline{c}_k\|} \underline{c}_k$$



$$\underline{v}_1 = \langle \underline{q}_1, \underline{v}_1 \rangle \underline{q}_1 + 0 \underline{q}_2 + \dots + 0 \underline{q}_n$$

$$\underline{v}_2 = \langle \underline{q}_1, \underline{v}_2 \rangle \underline{q}_1 + \langle \underline{q}_2, \underline{v}_2 \rangle \underline{q}_2 + 0 \underline{q}_3 + \dots + 0 \underline{q}_n$$

$$\underline{v}_3 = \langle \underline{q}_1, \underline{v}_3 \rangle \underline{q}_1 + \langle \underline{q}_2, \underline{v}_3 \rangle \underline{q}_2 + \langle \underline{q}_3, \underline{v}_3 \rangle \underline{q}_3 + 0 \underline{q}_4 + \dots + 0 \underline{q}_n$$

...

$$\underline{v}_k = \langle \underline{q}_1, \underline{v}_k \rangle \underline{q}_1 + \langle \underline{q}_2, \underline{v}_k \rangle \underline{q}_2 + \dots + \langle \underline{q}_k, \underline{v}_k \rangle \underline{q}_k + 0 \underline{q}_{k+1} + \dots + 0 \underline{q}_n$$

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 $\underline{q}_1, \dots, \underline{q}_k$
mit Koeff. $(q_1^H \underline{v}_k), \dots, (q_k^H \underline{v}_k) \in \mathbb{C}$

$$\underline{v}_k = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_k \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} \underline{v}_1 \\ \underline{v}_2 \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 \end{bmatrix} \begin{bmatrix} q_1^H \underline{v}_1 \\ q_2^H \underline{v}_2 \end{bmatrix}$$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \dots & \underline{v}_k \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_k \end{bmatrix} \begin{bmatrix} q_1^H \underline{v}_1 & q_1^H \underline{v}_2 & \dots & q_1^H \underline{v}_k \\ \vdots & \vdots & \ddots & \vdots \\ q_2^H \underline{v}_1 & q_2^H \underline{v}_2 & \dots & q_2^H \underline{v}_k \\ \vdots & \vdots & \ddots & \vdots \\ q_k^H \underline{v}_1 & q_k^H \underline{v}_2 & \dots & q_k^H \underline{v}_k \end{bmatrix}$$

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

↳ Spalten sind orthonormale Vektoren

Nach n Schritten:

$$\begin{bmatrix} \underline{v}_1 & \dots & \underline{v}_n \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \dots & \underline{q}_n \end{bmatrix} \begin{bmatrix} * \\ 0 \end{bmatrix}$$

$i \downarrow$
 $r_{ij} = q_j^H \underline{v}_i$
 $\leftarrow j$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 \end{bmatrix} \underbrace{\begin{bmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{bmatrix}}_{\underline{R}} \underline{R}^{-1}$$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 \end{bmatrix} \begin{bmatrix} \frac{1}{r_{11}} & -\frac{r_{12}}{r_{11}} \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \frac{1}{r_{22}} \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 \end{bmatrix}$$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \underline{v}_3 \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 & \underline{q}_3 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix} \underline{R}^{-1}$$

$$\begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \underline{v}_3 \end{bmatrix} \begin{bmatrix} \frac{1}{r_{11}} & -\frac{r_{12}}{r_{11}} & -\frac{r_{13}}{r_{11}} \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & \frac{1}{r_{22}} & -\frac{r_{23}}{r_{22}} \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 & \underline{q}_3 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \frac{1}{r_{33}} \end{bmatrix} = \begin{bmatrix} \underline{q}_1 & \underline{q}_2 & \underline{q}_3 \end{bmatrix}$$

Bez GS ~ Multiplikation mit oberen Dreiecksmatrizen

$$\underline{A} \underbrace{\underline{R}_1 \underline{R}_2 \dots \underline{R}_n}_{\underline{R}^{-1}} = \underline{Q} \Rightarrow \underline{A} = \underline{Q} \underline{R}$$

$$r_{ij} = \underline{q}_i^H \underline{v}_j$$

$$\underline{R}_1 = \begin{bmatrix} \frac{1}{r_{11}} & -\frac{r_{12}}{r_{11}} & \dots & -\frac{r_{1n}}{r_{11}} \\ 0 & 1 & 0 & 0 \\ \vdots & & \ddots & \\ 0 & 0 & \dots & 0 & 1 \end{bmatrix}$$

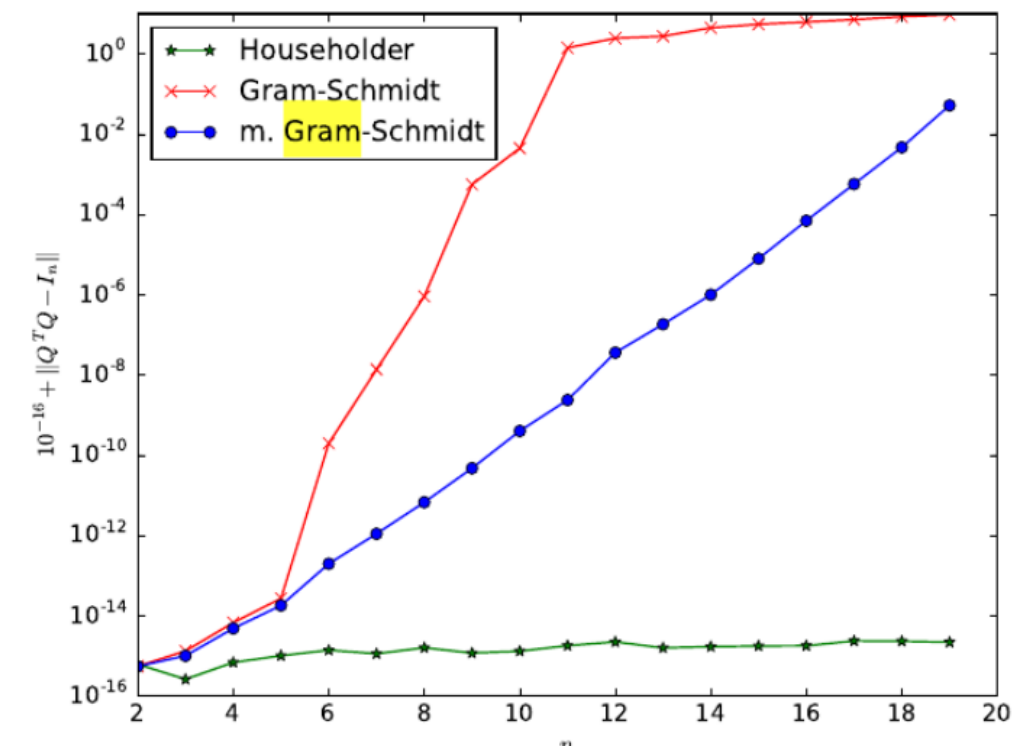
$$\underline{R}_2 = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \frac{1}{r_{22}} & -\frac{r_{23}}{r_{22}} & \dots & -\frac{r_{2n}}{r_{22}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix}$$

$$\underline{R}_3 = \begin{bmatrix} 1 & & & 0 \\ & 1 & & \\ & & \frac{1}{r_{33}} & \dots \\ 0 & & & 1 & \dots & 0 \end{bmatrix}$$

Bem Gram-Schmidt produziert nach n Schritten eine orthogonale Matrix Q in dem man in jedem Schritt mit einer oberen Dreiecksmatrix multipliziert
 \Rightarrow anfällig an Rundungsfehler!

Bem Ein Vorteil vom Gram-Schmidt Algorithmus ist: wenn man nach $k < n$ Schritten aufhört hat man bereits q_1, \dots, q_k orthogonal sodass
 $\text{span}\{q_1, \dots, q_k\} = \text{span}\{v_1, \dots, v_k\}$

$$A = \begin{bmatrix} t_0^{n-1} & \dots & t_0^1 & 1 \\ t_1^{n-1} & \dots & t_1^1 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ t_{19}^{n-1} & \dots & t_{19}^1 & 1 \end{bmatrix}, \text{ mit } t_i = \frac{i}{(\underline{m} - 1)}$$



classisches GS.:

$$P_{q_{k-1}} v_k = \begin{pmatrix} P_{q_{k-1}} & \dots & P_{q_2} & P_{q_1} & v_k \end{pmatrix} \rightarrow \text{alle gleichzeitig}$$

↓ 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 $i=1, 2, \dots, j-1$:

$$r_{ij} = q_i^H \underline{a}_j$$

$$\underline{v}_j = \underline{v}_j - r_{ij} \underline{q}_i$$

$$r_{jj} = \|\underline{v}_j\|$$

$$\underline{q}_j = \frac{\underline{v}_j}{r_{jj}}$$

sofort q_1, \dots, q_{k-1} verwenden.
Projektion aller folgenden
Spalten auf das neu gefundene
 q_i

für $j=1, \dots, n$:

$$\underline{v}_j = \underline{a}_j$$

für $i=1, \dots, n$

$$r_{ji} = \|\underline{v}_i\|; q_i = \frac{1}{r_{ji}} \underline{v}_i$$

für $j=i+1, \dots, n$:

$$r_{ij} = q_i^H \underline{v}_j$$

$$\underline{v}_j = \underline{v}_j - r_{ij} \underline{q}_i$$

Fazit:

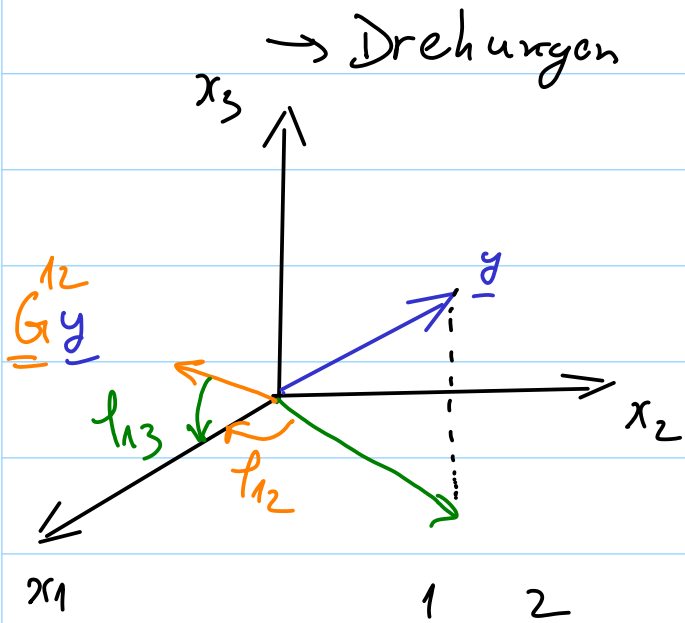
- 1) classisches Gram-Schmidt ist nur für analytische Zwecke geeignet!
- 2) modifiziertes Gram-Schmidt ist nur für kleine Matrizen zu verwenden!
- 3) sonst \underline{II} , \underline{III} die folgen!

\underline{II} $\underline{Q} \underline{R}$ via Rotationen.

$$\begin{matrix} n & m \end{matrix} \begin{bmatrix} \underline{A} \end{bmatrix} = \begin{matrix} m \end{matrix} \begin{bmatrix} \underline{Q} \end{bmatrix} \begin{matrix} n \end{matrix} \begin{bmatrix} \begin{matrix} * \\ * \\ 0 \end{matrix} \end{bmatrix} = \underline{Q} \underline{R}$$

$$\begin{matrix} n & m \end{matrix} \begin{bmatrix} \underline{A} \end{bmatrix} = \begin{matrix} m \end{matrix} \begin{bmatrix} \underline{Q} \end{bmatrix} \begin{matrix} n \end{matrix} \begin{bmatrix} \begin{matrix} * \\ * \\ 0 \end{matrix} \end{bmatrix} = \underline{Q} \underline{R}$$

Idee: Erzeuge die 0 unterhalb der Hauptdiagonale mittels orthogonale Matrizen



$$\underline{A} = \begin{bmatrix} y_1 & \tilde{z}_1 \\ y_2 & \tilde{z}_2 \\ y_3 & \tilde{z}_3 \end{bmatrix}$$

$$\underline{G}^{12}(\theta_{12}) = \begin{bmatrix} \cos \theta_{12} & \sin \theta_{12} & 0 \\ -\sin \theta_{12} & \cos \theta_{12} & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{matrix} 1 \\ 2 \\ \end{matrix}$$

↳ Givens Drehung, die y_2 zu 0 macht

$$\underline{G}^{12}(\theta_{12}) \underline{A} = \begin{bmatrix} * & * \\ 0 & * \\ y_3 & \tilde{z}_3 \end{bmatrix}$$

$$\underline{G}^{13}(\theta_{13}) = \begin{bmatrix} \cos \theta_{13} & 0 & \sin \theta_{13} \\ 0 & 1 & 0 \\ -\sin \theta_{13} & 0 & \cos \theta_{13} \end{bmatrix} \begin{matrix} 1 \\ 3 \\ \end{matrix}$$

↳ Givens Drehung, die y_3 zu 0 macht

$$\underline{G}^{13} \underline{G}^{12} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} * \\ 0 \\ 0 \end{bmatrix}$$

orthogonale Matrix

$$\underline{y} \in \mathbb{R}^n \quad \underline{G}^{12} \dots \underline{G}^{13} \underline{G}^{12} \underline{y} = \begin{bmatrix} * \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Dann kommt \underline{Q}^1 die zweite Spalte von \underline{A} dran:

$$\begin{bmatrix} * & \tilde{z}_1 & * \\ 0 & \tilde{z}_2 & * \\ 0 & \tilde{z}_3 & * \end{bmatrix}$$

\underline{G}^{23} um \tilde{z}_3 zu 0 zu machen

$$\underline{Q}^1 \underline{A} = \begin{bmatrix} * & x & x & \dots & x \\ 0 & x & x & \dots & x \\ 0 & \textcircled{x} & x & \dots & x \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \textcircled{x} & x & \dots & x \end{bmatrix}$$

$$\Rightarrow \underline{Q}^2 = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & \boxed{\text{Drehungen um}} \\ \vdots & \boxed{\begin{bmatrix} x \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ zu erzeugen}} \\ 0 & \vdots & \ddots & \vdots \end{bmatrix}$$

$$\underbrace{\underline{Q}^{n-1} \dots \underline{Q}^2 \underline{Q}^1}_{\text{orthogonale Matrizen}} \underline{A} = \underline{R}$$

$$(\underline{Q}^1)^T \dots (\underline{Q}^{n-1})^T$$

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

$$\Rightarrow \underline{A} = \underbrace{(\underline{Q}^1)^T (\underline{Q}^2)^T \dots (\underline{Q}^{n-1})^T}_{\underline{Q} \text{ orthogonal}} \underline{R}$$

$$\underline{G}^{ij}(l) = \begin{bmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & \cos l & & \sin l & \\ & & 0 & \ddots & 0 & \\ & & -\sin l & & \cos l & \\ 0 & \dots & 0 & & & \ddots & 1 \end{bmatrix} \begin{matrix} \leftarrow i' \\ \\ \\ \leftarrow j' \end{matrix}$$

$$\underline{G}^{ij}(l) \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ x_i \\ x_{i+1} \\ \vdots \\ x_{j-1} \\ x_j \\ x_{j+1} \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} x_1 \\ \vdots \\ x_{i-1} \\ r \\ x_{i+1} \\ \vdots \\ x_{j-1} \\ 0 \\ x_{j+1} \\ \vdots \\ x_n \end{bmatrix}$$

$$r = \sqrt{x_i^2 + x_j^2}$$

$$\cos l = \frac{x_{i'}}{r}$$

$$\sin l = \frac{x_{j'}}{r}$$

Bsp $\begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} \rightsquigarrow \begin{bmatrix} * \\ 0 \\ 0 \end{bmatrix} \quad r = \sqrt{4^2 + (-3)^2} = 5$
 $\cos \varphi = \frac{4}{5}, \sin \varphi = -\frac{3}{5}$

$$G^{12} \begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{4}{5} & -\frac{3}{5} & 0 \\ 3/5 & 4/5 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 \\ -3 \\ 1 \end{bmatrix} = \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix}$$

$$G^{13} \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{5}{\sqrt{26}} & 0 & \frac{1}{\sqrt{26}} \\ 0 & 1 & 0 \\ -\frac{1}{\sqrt{26}} & 0 & \frac{5}{\sqrt{26}} \end{bmatrix} \begin{bmatrix} 5 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} \sqrt{26} \\ 0 \\ 0 \end{bmatrix}$$

$$\underline{G} = \underline{G}^{13} \underline{G}^{12}$$

Ben Diese Methode (Givens-Drehungen) ist günstig falls \underline{A} sehr viele 0-Einträge hat (dünnbesetzte Matrix)

III QR-Zerlegung mit Spiegelungen (standard)

\underline{Q}^1 spiegelt \underline{z} so dass $\underline{Q}^1 \underline{z}$ auf Ox_1 liegt

$$\underline{Q}^1 \underline{z} = \begin{bmatrix} * \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Rightarrow \underline{Q}^1 \underline{A} = \begin{bmatrix} * & x & x & \dots & x \\ 0 & x & - & \dots & \\ \vdots & \vdots & & & \\ 0 & x & - & \dots & \end{bmatrix}$$

↑ Spiegelung in \mathbb{R}^{n-1}
 $x_2 O x_3$

\underline{Q}^2 spiegelt \underline{z} so dass $\underline{Q}^2 \underline{z}$ auf Ox_2 liegt

$$\underline{Q}^2 \underline{Q}^1 \underline{A} = \begin{bmatrix} * & x & x & x & \dots \\ 0 & x & x & x & \dots \\ \vdots & 0 & x & x & \dots \\ 0 & \vdots & \vdots & \vdots & \dots \\ 0 & 0 & x & \dots & \end{bmatrix} \quad \text{usw}$$

$$\underline{Q}^{n-1} \underline{Q}^{n-2} \dots \underline{Q}^2 \underline{Q}^1 \underline{A} = \underline{R} \Rightarrow$$

$$\underline{A} = \underline{Q} \underline{R} \quad \text{mit} \quad \underline{Q} = \underline{Q}^{1T} \underline{Q}^{2T} \dots \underline{Q}^{n-1T}$$

Bsp $\underline{x} = \begin{bmatrix} 2 \\ 2 \\ 1 \end{bmatrix} \Rightarrow \underline{Q}\underline{x} = \begin{bmatrix} 3 \\ 0 \\ 0 \end{bmatrix} = 3 \underline{e}_1$

$$\|\underline{x}\| = \sqrt{2^2 + 2^2 + 1} = 3$$

$$\underline{x} = \|\underline{x}\| \underline{e}_1 + \underline{v} \Rightarrow \underline{v} = \underline{x} - \|\underline{x}\| \underline{e}_1 = \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix}$$

$$\underline{u} = \frac{1}{\|\underline{v}\|} \underline{v} = \frac{1}{\sqrt{6}} \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix}$$

$$\underline{Q} = \underline{I} - 2 \underline{u} \underline{u}^T = \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} - \frac{2}{6} \begin{bmatrix} -1 \\ 2 \\ 1 \end{bmatrix} \begin{bmatrix} -1 & 2 & 1 \end{bmatrix} =$$

$$= \begin{bmatrix} 1 & & \\ & 1 & \\ & & 1 \end{bmatrix} - \frac{2}{6} \begin{bmatrix} 1 & -2 & -1 \\ -2 & 4 & 2 \\ -1 & 2 & 1 \end{bmatrix} =$$

$$= \frac{1}{3} \begin{bmatrix} 2 & 2 & 1 \\ 2 & -1 & -2 \\ 1 & -2 & 2 \end{bmatrix}$$

$$\underline{A} = \begin{bmatrix} 1 & 1 \\ 2 & 0 \\ 2 & 0 \end{bmatrix} \Rightarrow \underline{Q}^1 \underline{A} = \left[\underline{Q}^1 \begin{bmatrix} 1 \\ 2 \\ 2 \end{bmatrix} \quad \underline{Q}^1 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right]$$

$$\underline{Q}^1 = \underline{I} - \frac{2}{12} \begin{bmatrix} -2 \\ 2 \\ 2 \end{bmatrix} \begin{bmatrix} -2 & 2 & 2 \end{bmatrix} \Rightarrow$$

$$\underline{Q}^1 \underline{A} = \begin{bmatrix} 3 & \frac{1}{3} \\ 0 & \frac{2}{3} \\ 0 & \frac{2}{3} \end{bmatrix} \text{ usw.}$$

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

$$\underline{A} = \underline{Q} \begin{bmatrix} \underline{R} \\ 0 \end{bmatrix} \text{ mit } \underline{R} \text{ obere Dreiecksmatrix}$$

und die Elemente auf der Diagonale von \underline{R} sind ≥ 0 .

Bez Falls die Elemente auf der Diagonale von \underline{R} > 0 sind, dann ist \underline{R} die Cholesky-Zerlegung von $\underline{A}^T \underline{A}$

Beweis $\underbrace{\underline{A}^T \underline{A}}_{\text{symm. positiv.}} = \begin{bmatrix} \underline{R}^T & \underline{0} \end{bmatrix} \underbrace{\underline{Q}^T \underline{Q}}_{\underline{I}} \begin{bmatrix} \underline{R} \\ \underline{0} \end{bmatrix} = \underline{R}^T \underline{R}$

Theorem Falls $\text{rang } \underline{A} = r < n$ dann gibt es eine Permutation \underline{P} so dass

$$\underline{A} \underline{P} = \underline{Q} \begin{bmatrix} \underline{R}_{11} & \underline{R}_{12} \\ \underline{0} & \underline{0} \end{bmatrix}$$

r
 $n-r$

mit

\underline{R}_{11} obere Dreiecksmatrix mit Elementen auf der Diagonale > 0 .

\underline{R}_{11} ist eindeutig, \underline{R}_{12} nicht!

Beweis wähle r Spalten lin. unabh. und permutiere sie nach vorne

$$\underline{A} \underline{P} = \begin{bmatrix} \underline{A}_1 & \underline{A}_2 \end{bmatrix} \text{ und wende das vorige Theorem.}$$

r

(Givens-Drehungen) an.

Bez QR-Zerlegung mit orthogonalen Transformationen darf nicht vorzeitig abgebrochen werden.

Bez $\underline{A} = \underline{Q} \underline{R} ; \quad \underline{A} \underline{x} = \underline{b}$

$$\underline{Q}^T \underline{Q} \underline{R} \underline{x} = \underline{b} \Rightarrow \underline{R} \underline{x} = \underline{Q}^T \underline{b}$$

Rückwärts substitution!

§7.3. Singulärwertzerlegung.

$$\underline{A} \in \mathbb{R}^{m \times n}; \mathbb{C}^{n \times n}$$

$$r = \text{rang}(\underline{A})$$

$$\begin{bmatrix} \boxed{A} \end{bmatrix} = \begin{bmatrix} \boxed{U} \end{bmatrix} \begin{bmatrix} \boxed{\Sigma} \end{bmatrix} \begin{bmatrix} \boxed{V^H} \end{bmatrix}$$

Diagram illustrating the SVD decomposition of matrix A (size $m \times n$) into matrices U (size $m \times m$), Σ (size $m \times n$), and V^H (size $n \times n$). The rank r is indicated above U . The matrix Σ is shown with a diagonal of yellow squares. The matrix V^H is shown with a horizontal line separating the first r rows from the remaining $n-r$ rows.

$$\begin{bmatrix} \boxed{A} \end{bmatrix} = \begin{bmatrix} \boxed{U} \end{bmatrix} \begin{bmatrix} \boxed{\Sigma} \end{bmatrix} \begin{bmatrix} \boxed{V^H} \end{bmatrix}$$

Diagram illustrating the SVD decomposition of matrix A (size $m \times n$) into matrices U (size $m \times m$), Σ (size $m \times n$), and V^H (size $n \times n$). The matrix Σ is shown with a diagonal of yellow squares.

Es gibt $\underline{u}_1, \underline{u}_2, \dots, \underline{u}_r, \underline{u}_{r+1}, \dots, \underline{u}_n \in \mathbb{R}^n$ ONB
 $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_r, \underline{v}_{r+1}, \dots, \underline{v}_n \in \mathbb{R}^n$ ONB

$$\underline{u}_i^T \underline{A} \underline{v}_j = \begin{cases} \sigma_i, & \text{falls } i=j \leq r \\ 0, & \text{sonst} \end{cases}$$

$$\underline{U}^T \underline{A} \underline{V} = \underline{\Sigma} = \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & 0 \end{bmatrix}$$

Diagram illustrating the matrix $\underline{\Sigma}$ (size $m \times n$) with a diagonal of yellow squares. The matrix is partitioned into a top $r \times r$ block and a bottom $(m-r) \times (n-r)$ block.

$$\underline{U}_r = [\underline{u}_1 \dots \underline{u}_r], \quad \underline{V}_r = [\underline{v}_1 \dots \underline{v}_r]$$

$\underline{v}_1, \dots, \underline{v}_r$ = rechte Singulärvektoren = EV von $A^T A$
 $\underline{u}_1, \dots, \underline{u}_r$ = linke Singulärvektoren = EV von $A A^T$

$$\underline{A} \underline{v}_1 = \sigma_1 \underline{u}_1, \dots, \underline{A} \underline{v}_r = \sigma_r \underline{u}_r$$

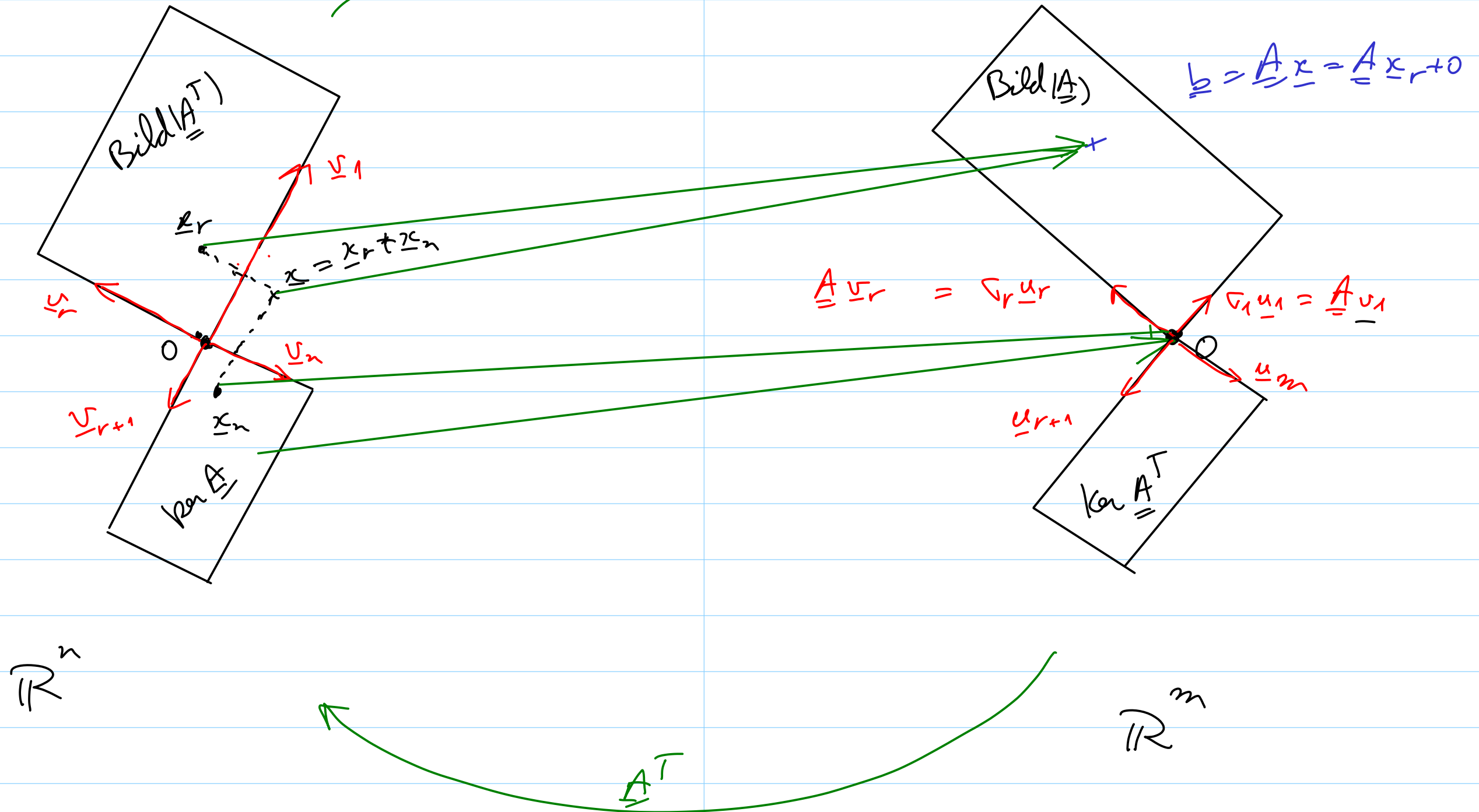
$$\underline{A} \underline{v}_{r+1} = 0, \dots, \underline{A} \underline{v}_n = 0$$

$\text{Ker } \underline{A} = \text{span} \{ \underline{v}_{r+1}, \dots, \underline{v}_n \}; \text{Bild } \underline{A} = \text{span} \{ \underline{u}_1, \dots, \underline{u}_r \}$

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

$$\underline{A} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$\underline{A}$$



Bem $5 \gg 3 \gg 0.1 \gg 10^{-5} \gg 10^{-8} \gg 10^{-12} \gg 5 \cdot 10^{-16} \gg 2 \cdot 10^{-16} \gg 10^{-42} \geq 0 = 0$

$\sigma_1 \quad \sigma_2 \quad \sigma_3 \quad \sigma_4 \quad \sigma_5 \quad \sigma_6 \quad \sigma_7 \quad \sigma_8 \quad \sigma_9 \quad \sigma_{10} = \sigma_{11} = \dots = \sigma_{42}$

Numerischer Rang = ?

Anwendung

1) Speicherplatzreduktion.

SVD: $rm + rn + r = r(m+n+1)$ Plätze

$$\underline{A} = \sum_{k=1}^r \underbrace{u_k \sigma_k v_k^T}_{\text{Rang 1 Matrizen}} \approx \sum_{k=1}^{r_{\min}} u_k \sigma_k v_k^T$$

Rang 1 Matrizen

2) Berechnung der Norm einer Matrix

$$\|\underline{A}\|_2 = \max_{\|\underline{x}\|_2=1} \|\underline{A}\underline{x}\|_2$$

$$\|\underline{A}\underline{x}\|_2^2 = \left\| \underline{U} \sum \underline{V}^T \underline{x} \right\|_2^2 \stackrel{\substack{\underline{U} \text{ orthogonal} \Rightarrow \\ \text{erhält Längen}}}{=} \left\| \sum \underline{V}^T \underline{x} \right\|_2^2$$

$$\|\underline{A}\| = \max_{\|\underline{x}\|_2=1} \|\underline{A}\underline{x}\|_2 = \max_{\|\underline{x}\|_2=1} \left\| \sum \underbrace{\underline{V}^T \underline{x}}_{\underline{z}} \right\|_2 = \max_{\|\underline{z}\|_2=1} \left\| \sum \underline{z} \right\|_2 =$$

$\underline{x} = \underline{V} \underline{z}$

$$= \max_{\|\underline{z}\|_2=1} \left\| \sum \underline{z} \right\|_2 = \max_{\|\underline{z}\|_2=1} \sqrt{\sigma_1^2 |z_1|^2 + \sigma_2^2 |z_2|^2 + \dots + \sigma_r^2 |z_r|^2} \leq$$

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

$$\leq \max_{\|\underline{z}\|_2=1} \sqrt{\sigma_1^2 \|\underline{z}\|_2^2} = \sigma_1 \text{ erreicht für } \underline{z} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

Somit $\|\underline{A}\|_2 = \sigma_1$

3) Kompression.

Daten mit Rausch/ Fehler:

$$\underline{A} = \sum_{j=1}^r \sigma_j \underline{u}_j \underline{v}_j^T \approx \sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^T \text{ mit } k \ll r.$$

\neq Diagonalisierung für $m=n$

Bsp $\underline{A} = \begin{bmatrix} 2 & 100 \\ 0 & 1 \end{bmatrix} \quad \lambda_1=2, \quad \lambda_2=1$

$$\underline{A} = \underline{S} \underline{\Lambda} \underline{S}^{-1} = \underset{\lambda_1}{2} \underset{\underline{S}_1}{\begin{bmatrix} 1 \\ 0 \end{bmatrix}} \underset{\underline{S}_1^T}{\begin{bmatrix} 1 & 100 \end{bmatrix}} + \underset{\lambda_2}{1} \underset{\underline{S}_2}{\begin{bmatrix} 1 \\ -\frac{1}{100} \end{bmatrix}} \underset{\underline{S}_2^T}{\begin{bmatrix} 0 & -100 \end{bmatrix}}$$

Versuch der Approximation mit einer Rang 1-Matrix

Fehler $\underline{A} - \lambda_1 \underline{A}_1 \underline{A}_1^T = \begin{bmatrix} 0 & -100 \\ 0 & 1 \end{bmatrix}$

$$\underline{A} - \lambda_1 \underline{A}_1 \underline{A}_1^T = \begin{bmatrix} 2 & 200 \\ 0 & 0 \end{bmatrix}$$

also so geht es nicht!

Mit SVD:

$$\underline{A} = \underset{\sigma_1}{100.025} \underset{\underline{U}_1}{\begin{bmatrix} 0.9995 \\ 0.001 \end{bmatrix}} \underset{\underline{V}_1^T}{\begin{bmatrix} 0.01999 & 0.9998 \end{bmatrix}} +$$

$$+ \underset{\sigma_2}{0.02} \underset{\underline{U}_2}{\begin{bmatrix} 0.01 \\ 0.9995 \end{bmatrix}} \underset{\underline{V}_2^T}{\begin{bmatrix} -0.9998 & 0.01999 \end{bmatrix}}$$

Fehler $\underline{A} - \sigma_1 \underline{U}_1 \underline{V}_1^T = \begin{bmatrix} 2 \cdot 10^{-4} & 0 \\ -2 \cdot 10^{-2} & 4 \cdot 10^{-4} \end{bmatrix}$ klein

Theorem [Eckart-Young]4) PCA \rightarrow Skript.

Sei $\underline{A} \in \mathbb{C}^{m \times n}$. Für jedes $k \leq \text{Rang}(\underline{A})$
 gibt es eine abgebrochene SVD

$$\underline{A}_k = \sum_{j=1}^k \sigma_j \underline{u}_j \underline{v}_j^H$$

die, die beste Approximation von Rang k an \underline{A} ist,
 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.

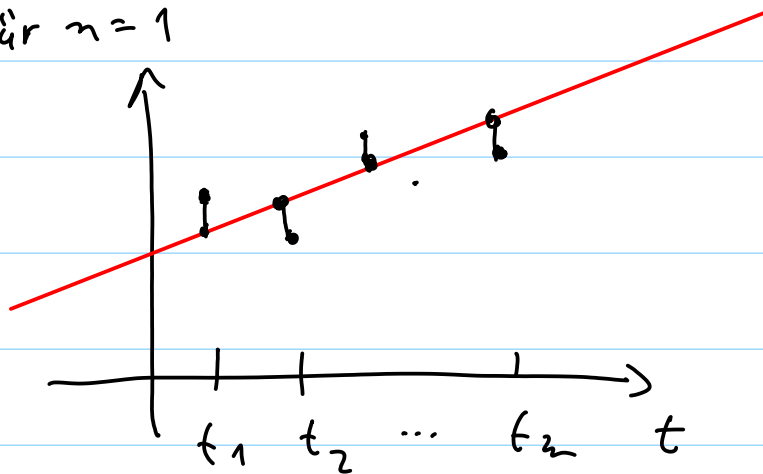
§8 Ausgleichsrechnung

§8.1. Motivation und Normalengleichung

Lineares Model:

$y = \underline{a}^T \underline{t} + c \cdot 1$ mit $\underline{a} \in \mathbb{R}^n$, $c \in \mathbb{R}$ Parameter
 m Messungen: $\underline{t}_1, \underline{t}_2, \dots, \underline{t}_m \in \mathbb{R}^n$ Messpunkte
 $\downarrow \quad \downarrow \quad \dots \quad \downarrow$
 $y_1 \quad y_2 \quad \dots \quad y_m$ Gemessene Werte
 (mit Fehlern)

Für $n=1$



Vorschlag: \underline{a}, c die \underline{q}, p , die das Minimum realisieren:

$$\min_{\substack{\underline{p} \in \mathbb{R}^n \\ \underline{q} \in \mathbb{R}}} \sum_{i=1}^m |y_i - \underline{p}^T \underline{t}_i - 1 \cdot q|^2 = \min_{\underline{x} \in \mathbb{R}^{1+n}} \| \underline{A} \underline{x} - \underline{b} \|_2$$

$$\underline{b} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \quad \underline{x} = \begin{bmatrix} q \\ \underline{p} \end{bmatrix} \in \mathbb{R}^{1+n}, \quad \underline{A} = \begin{bmatrix} 1 & \underline{t}_1^T \\ 1 & \underline{t}_2^T \\ \vdots & \vdots \\ 1 & \underline{t}_m^T \end{bmatrix} \in \mathbb{R}^{m \times (1+n)}$$

$\mathbb{R}^{m \times (1+n)}$

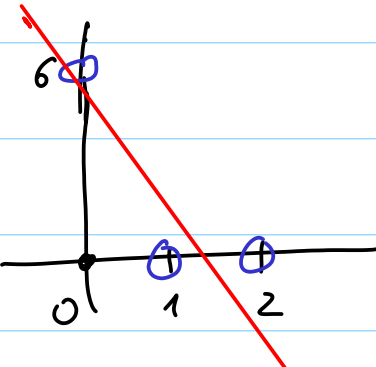
Methode der kleinsten Quadrate
 (least squares Method).

Bsp $n=1$ $t_1=0, t_2=1, t_3=2$

$$y_1=6, y_2=0, y_3=0$$

Modell $y = dt + c_1$

Messwerte: $\begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$, $\underline{A} = \begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix}$



$$\min_{\underline{x} \in \mathbb{R}^2} \|\underline{A}\underline{x} - \underline{b}\|_2$$

alg. Problem:
$$\begin{cases} b_1 = dt_1 + c \\ b_2 = dt_2 + c \\ b_3 = dt_3 + c \end{cases}$$

$\underline{A}\underline{x} = \underline{b}$ hat keine Lösung.

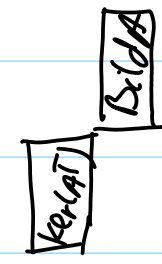
I Algebraische Methode.

$$\underline{A}^T \mid \underline{A}\underline{x} = \underline{b} \Rightarrow \underbrace{\underline{A}^T \underline{A}}_{\text{quadratisch, symmetrisch.}} \underline{x} = \underline{A}^T \underline{b} \quad \text{Normalengleichung}$$

Falls $\text{Rang}(\underline{A}) = n \Rightarrow \underline{A}^T \underline{A}$ invertierbar
kann das LGS eindeutig lösen

$$\text{cond}(\underline{A}) = \frac{\|\underline{A}^{-1}\|}{\|\underline{A}\|} \underset{100}{\text{gross}} \Rightarrow \text{cond}(\underline{A}^T \underline{A}) \underset{10^4}{\text{riesig.}}$$

$$\text{cond}(\underline{A}^T \underline{A}) = \text{cond}(\underline{A})^2 \quad \text{LU} \rightarrow$$

$$\underline{A}: \mathbb{R}^n \rightarrow \mathbb{R}^m; \underline{b} \in \mathbb{R}^m \Rightarrow$$


$$\underline{b} = \underbrace{\underline{p}}_{\text{Bild}(\underline{A})} + \underbrace{\underline{e}}_{\text{Ker}(\underline{A}^T)}, \quad \underline{p} \perp \underline{e}$$

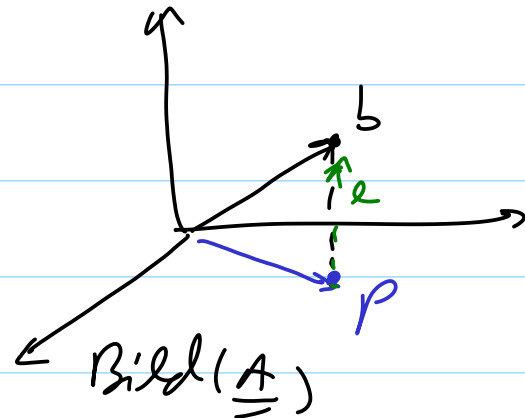
$\underline{p} \in \text{Bild}(\underline{A}) \Rightarrow$ es gibt $\underline{\hat{x}} \in \mathbb{R}^n$ so dass $\underline{A}\underline{\hat{x}} = \underline{p}$

Fehler $\underline{b} - \underline{A}\underline{x} = \underline{b} - \underline{p} = \underline{e} \in \ker(\underline{A}^T) \perp \text{Bild}(\underline{A})$

$$\Rightarrow \|\underline{b} - \underline{A}\underline{x}\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \|\underline{b} - \underline{A}\underline{x}\|_2^2$$

$$\begin{aligned} \underline{A}^T \underline{b} &= \underline{A}^T (\underline{p} + \underline{e}) = \\ &= \underline{A}^T \underline{p} + \underline{A}^T \underline{e} = \underline{A}^T \underline{p} \end{aligned}$$

$$\Rightarrow \underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{p}$$



II $f(\underline{x}) = \|\underline{A}\underline{x} - \underline{b}\|_2^2$; $f: \mathbb{R}^{1+n} \rightarrow [0, \infty[$ Min?

Finde \underline{x} :
 $\nabla f(\underline{x}) = 0$

$$\begin{aligned} f(\underline{x}) &= (\underline{A}\underline{x} - \underline{b})^T (\underline{A}\underline{x} - \underline{b}) = \underline{x}^T \boxed{\underline{A}^T \underline{A}} \underline{x} - \underline{b}^T \underline{A} \underline{x} - \underbrace{\underline{x}^T \underline{A}^T \underline{b}}_{\underline{w}} + \underline{b}^T \underline{b} \\ &= \underline{x}^T \underline{M} \underline{x} - 2 \underline{x}^T \underline{w} + \underline{b}^T \underline{b} = \end{aligned}$$

$$= \sum_{i=1}^n x_i \sum_{j=1}^n m_{ij} x_j - 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} =$$

$$= x_1 \sum_{j=1}^n m_{1j} x_j + \underbrace{\sum_{i=2}^n x_i \sum_{j=1}^n m_{ij} x_j}_{\sum_{j=1}^n x_j \sum_{i=2}^n m_{ij} x_i} - 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} =$$

$$= x_1 m_{11} x_1 + x_1 \sum_{j=2}^n m_{1j} x_j + x_1 \sum_{i=2}^n m_{i1} x_i + \sum_{j=2}^n x_j \sum_{i=2}^n m_{ij} x_i -$$

$$- 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b} =$$

\uparrow
 M symmetrisch.

$$= m_{11} x_1^2 + 2 x_1 \sum_{i=2}^n x_i m_{i1} + \sum_{i,j=2}^n x_i m_{ij} x_j -$$

$$- 2 \sum_{i=1}^n x_i w_i + \underline{b}^T \underline{b}$$

$$\frac{\partial}{\partial x_1} f(\underline{x}) = 2m_1 x_1 + 2 \sum_{i=2}^n x_i m_{i1} - 2w_1 \Rightarrow$$

$$Df(\underline{x}) = 2 \underline{M} \underline{x} - 2 \underline{A}^T \underline{b}$$

kritische Punkte $Df(\underline{x}) = 0 \Leftrightarrow 2 \underline{A}^T \underline{A} \underline{x} - 2 \underline{A}^T \underline{b} = 0$

$$\Leftrightarrow \underline{A}^T \underline{A} \underline{x} = \underline{A}^T \underline{b}$$

Minimum nur wenn $D^2 f(\underline{x}) = 2 \underline{A}^T \underline{A}$

symmetrische pos. def.
nur wenn $\text{Rang}(\underline{A}) = n$.

$$\text{cod}_2(\underline{A}^T \underline{A}) = \text{cod}_2(\underbrace{\underline{V} \underline{\Sigma}^T \underline{U}^T \underline{U} \underline{\Sigma} \underline{V}^T}_{\underline{I}}) =$$

$$= \text{cod}_2(\underline{V} \underline{\Sigma}^2 \underline{V}^T) = \text{cod}_2(\underline{\Sigma}^2) = \frac{\sigma_1^2}{\sigma_r^2} = \text{cod}_2(\underline{A})^2$$

wobei $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r$ die Singulärwerte von \underline{A} sind.

Bsp $f \in L^2 = \left\{ f: I \rightarrow \mathbb{R}; \int_I |f(t)|^2 dt < \infty \right\}$

$f_n \in V_n$, $\dim V_n < \infty$

\hookrightarrow 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

$$\sum_{i=1}^m |f_n(t_i) - y_i|^2 = \text{minimal!}$$

$$\begin{cases} \sum_{j=1}^n b_j(t_i) x_j = y_i \\ i=1, 2, \dots, m \end{cases} \quad A = [b_j(t_i)]_{\substack{i=1, \dots, m \\ j=1, \dots, n}}$$

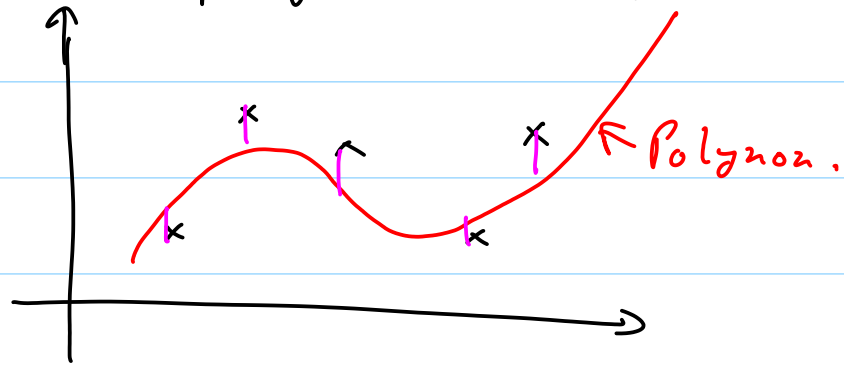
$\hookrightarrow i=1, 2, \dots, m$

Bem Spezielle Wahl:

$V_n = \mathcal{P}_n =$ Polynome vom Grad maximal $n-1$

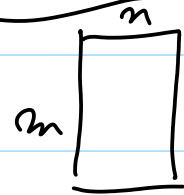
Basis in V_n : Monome $b_j(t) = t^{j-1}$

→ polynomialos fit ; polyfit



§8.2. Lösung mittels orthogonalen Transformationen

Fall 1 $\text{Rang}(\underline{A}) = n$ (voller Rang)



$$\underline{A} = \underline{Q} \underline{R} = \underline{Q} \begin{bmatrix} \tilde{\underline{R}} \\ 0 \end{bmatrix} \begin{matrix} n \\ m-n \end{matrix} \quad \underline{Q} \text{ orthogonal}$$

$$\min_{\underline{x} \in \mathbb{R}^n} \|\underline{A} \underline{x} - \underline{b}\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \|\underline{Q} \underline{R} \underline{x} - \underline{b}\|_2^2 =$$

$\underline{I} = \underline{Q} \underline{Q}^H$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underline{Q} (\underline{R} \underline{x} - \underline{Q}^H \underline{b})\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\underline{R} \underline{x} - \underbrace{\underline{Q}^H \underline{b}}_{\underline{\beta}}\|_2^2 =$$

$$\min_{\underline{x} \in \mathbb{R}^n} \left\| \begin{bmatrix} \tilde{\underline{R}} \\ 0 \end{bmatrix} \underline{x} - \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{n+1} \\ \vdots \\ \beta_m \end{bmatrix} \right\|_2^2 = \min_{\underline{x} \in \mathbb{R}^n} \left\| \begin{bmatrix} \tilde{\underline{R}} \underline{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \tilde{\underline{\beta}} \\ \beta_{n+1} \\ \vdots \\ \beta_m \end{bmatrix} \right\|_2^2 =$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \left\{ \|\tilde{\underline{R}} \underline{x} - \tilde{\underline{\beta}}\|_2^2 + |\beta_{n+1}|^2 + \dots + |\beta_m|^2 \right\}$$

$$\begin{bmatrix} \tilde{\underline{R}} & 0 \\ 0 & 0 \end{bmatrix} \underline{x} = \tilde{\underline{\beta}}$$

$$= \min_{\underline{x} \in \mathbb{R}^n} \|\tilde{\underline{R}} \underline{x} - \tilde{\underline{\beta}}\|_2^2 + |\beta_{n+1}|^2 + \dots + |\beta_m|^2$$

0 " erreicht für \underline{x} = die Lösung von $\tilde{\underline{R}} \underline{x} = \tilde{\underline{\beta}}$

Bsp $\begin{bmatrix} 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 0 \\ 0 \end{bmatrix}$ Normalengleichung $\Rightarrow \begin{bmatrix} 3 & 3 \\ 3 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 6 \\ 0 \end{bmatrix} \Rightarrow \begin{matrix} x_1 = 5 \\ x_2 = -3 \end{matrix}$

Mit QR-Zerlegung:

$$\underline{A} = \underline{Q} \underline{R} = \begin{bmatrix} -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & ? \\ -\frac{1}{\sqrt{3}} & 0 & ? \\ -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & ? \end{bmatrix} \begin{bmatrix} -\frac{3}{\sqrt{3}} & -\frac{3}{\sqrt{3}} \\ 0 & -\frac{2}{\sqrt{2}} \\ 0 & 0 \end{bmatrix}$$

LGS $\begin{bmatrix} -\frac{3}{\sqrt{3}} & -\frac{3}{\sqrt{3}} \\ 0 & -\frac{2}{\sqrt{2}} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} -\frac{6}{\sqrt{3}} \\ 6/\sqrt{2} \end{bmatrix} \Rightarrow \begin{matrix} x_1 = 5 \\ x_2 = -3 \end{matrix}$

Fall 2 $r = \text{Rang}(\underline{A}) < n \Rightarrow \text{SVD!}$

$$\underline{A} = \begin{bmatrix} \overset{r}{\underline{U}_1} & \overset{m-r}{\underline{U}_2} \end{bmatrix} \begin{bmatrix} \underline{\Sigma}_r & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \underline{V}_1^T \\ \underline{V}_2^T \end{bmatrix} \begin{matrix} r \\ n-r \end{matrix}$$

$$\underline{\Sigma}_r = \begin{bmatrix} \sigma_1 & \dots & 0 \\ 0 & \dots & \sigma_r \end{bmatrix} \text{ mit } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0.$$

\underline{U} orthogonal

$$\| \underline{A} \underline{x} - \underline{b} \|_2 = \left\| \begin{bmatrix} \sum_{i=1}^r \underline{V}_i^T \underline{x} \\ 0 \end{bmatrix} - \begin{bmatrix} \underline{U}_1^T \underline{b} \\ \underline{U}_2^T \underline{b} \end{bmatrix} \right\|_2$$

ist minimal für $\underline{x} = \text{Lösung von}$

$$\sum_{i=1}^r \underline{V}_i^T \underline{x} = \underline{U}_1^T \underline{b}$$

$$\text{d.h. } \underline{x} = \underline{V}_1 \underline{\Sigma}_r^{-1} \underline{U}_1^T \underline{b}$$

und das Minimum ist $\| \underline{U}_2^T \underline{b} \|_2$

\rightarrow wird verwendet in standard codes $\text{Lsq}(\underline{A}, \underline{b})$

Def Pseudoinverse einer Matrix
(Moore-Penrose inverse)

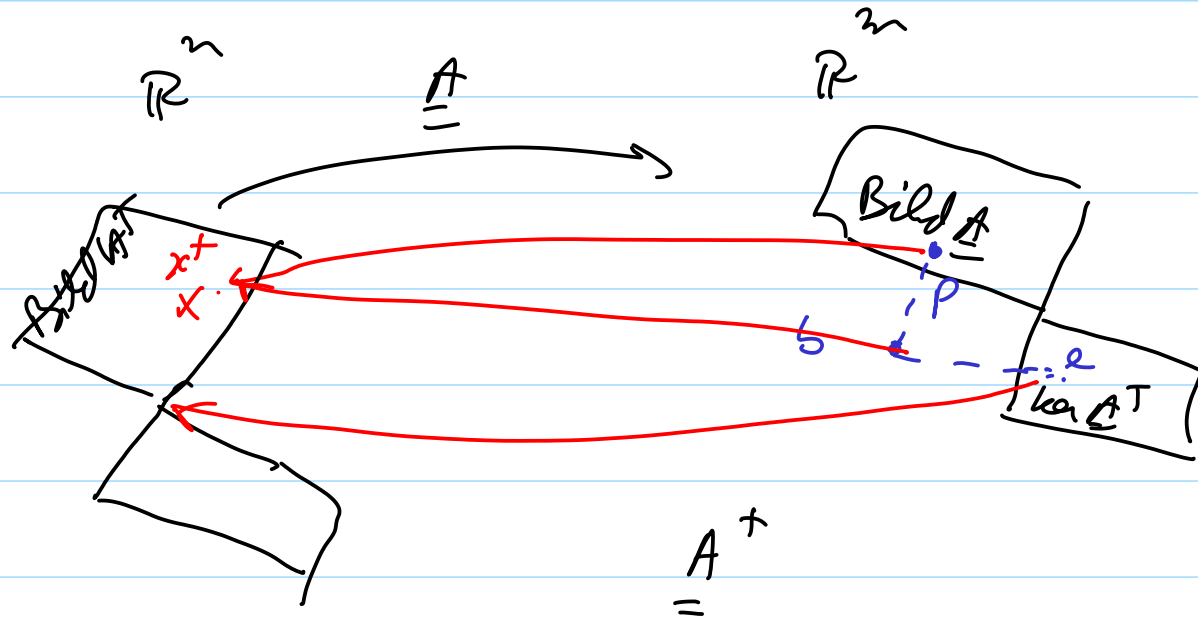
$$\underline{A}^+ := \underline{V} \underline{\Sigma}^+ \underline{U}^T \text{ mit } \underline{\Sigma}^+ = \begin{bmatrix} \sigma_1^{-1} & \dots & \sigma_r^{-1} & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix}$$

$$\mathbb{R}^n \xrightarrow{\underline{A}} \mathbb{R}^m$$

\vdots

$$x^+ \xleftarrow{\underline{A}^+} b$$

"die Ausgleichs Lösung" $\|\underline{A}x^+ - b\|_2^2$



§8.3. Lineare Ausgleichsrechnung mit linearen Nebenbedingungen.

Finde $\underline{x} \in \mathbb{R}^n$ bei gegebenen $\underline{A} \in \mathbb{R}^{m \times n}$, $\underline{b} \in \mathbb{R}^m$
mit $m \geq n = \text{Rang}(\underline{A})$

$$\begin{cases} \|\underline{A}\underline{x} - \underline{b}\|_2^2 = \text{minimal!} \\ \underline{C}\underline{x} = \underline{d} \end{cases} \quad \text{mit } \underline{C} \in \mathbb{R}^{p \times n} \text{ mit } p < n$$

$\underline{d} \in \mathbb{R}^p$ $\text{Rang } \underline{C} = p$

Methoden Lagrange-Multiplikatoren $\underline{m} \in \mathbb{R}^p$

$$L(\underline{x}, \underline{m}) = \frac{1}{2} \|\underline{A}\underline{x} - \underline{b}\|_2^2 + \underline{m}^T (\underline{C}\underline{x} - \underline{d})$$

$$\min_{\underline{x} \in \mathbb{R}^n} \max_{\underline{m} \in \mathbb{R}^p} L(\underline{x}, \underline{m})$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial \underline{x}}(\underline{x}, \underline{m}) = 0 \\ \frac{\partial \mathcal{L}}{\partial \underline{m}}(\underline{x}, \underline{m}) = 0 \end{cases} \quad \begin{cases} \underline{A}^T(\underline{A}\underline{x} - \underline{b}) + \underline{C}^T \underline{m} = 0 \\ 0 + \underline{C}\underline{x} - \underline{d} = 0 \end{cases}$$

$$\begin{bmatrix} \underline{A}^T \underline{A} & \underline{C}^T \\ \underline{C} & \underline{0} \end{bmatrix} \begin{bmatrix} \underline{x} \\ \underline{m} \end{bmatrix} = \begin{bmatrix} \underline{A}^T \underline{b} \\ \underline{d} \end{bmatrix}$$

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 $\underline{S}^T \underline{S} = -\underline{G} \underline{G}^T$

$$\begin{bmatrix} \underline{A}^T \underline{A} & \underline{C}^T \\ \underline{C} & \underline{0} \end{bmatrix} = \begin{bmatrix} \underline{R}^T & \underline{0} \\ \underline{G} & \underline{S} \end{bmatrix} \begin{bmatrix} \underline{R} & \underline{G}^T \\ \underline{0} & \underline{S}^T \end{bmatrix}$$

Block Cholesky-Zerlegung.

Nachteil: $\text{cond}(\underline{A}^T \underline{A}) = \text{cond}(\underline{A})^2$ ☹️

Vorteil: falls $\underline{A}, \underline{C}$ eine Struktur haben, dann können wir eventuell diese in den Algorithmen ausnutzen.

Methode 2 SVD besser für die Kondition, verliert die Struktur.

$$\underline{C} = \underline{U} \begin{bmatrix} \underline{\Sigma} & \underline{0} \end{bmatrix} \begin{bmatrix} \underline{V}_1^H \\ \underline{V}_2^H \end{bmatrix} \begin{matrix} p \\ n-p \end{matrix}$$

$$\ker(\underline{C}) = \text{Bild}(\underline{V}_2)$$

Trick:

$$\text{Definiere } \underline{x}_0 = \underline{V}_1 \underline{\Sigma}_1^{-1} \underline{U}^H \underline{d}$$

dann suche

$$\underline{x} = \underline{x}_0 + \underbrace{\underline{V}_2 \underline{y}}_{\ker(\underline{C}) = \text{Bild}(\underline{V}_2)} \quad \text{mit } \underline{y} \in \mathbb{R}^{n-p}$$

$$\ker(\underline{C}) = \text{Bild}(\underline{V}_2)$$

(und somit automatisch $\underline{C}\underline{x} = \underline{d}$)

$$\| \underline{A} \underline{x} - \underline{b} \|_2^2 = \| \underline{A} \underline{x}_0 + \underline{A} \underline{V}_2 \underline{y} - \underline{b} \|_2^2 =$$

$$= \| \boxed{\underline{A} \underline{V}_2 \underline{y}} - \underbrace{(\underline{b} - \underline{A} \underline{x}_0)}_{\text{fest}} \|_2^2$$

min nach \underline{y} ↑ Standard

(lin. Ausgleichsproblem)

§8.4. Totales lineares Ausgleichsproblem

total = Messfehler auch in den Messorten \underline{t}_i

$\underline{A} \underline{x} = \underline{b}$ mit Fehler in \underline{b} auch in \underline{A}

Ben $\underline{A} \underline{x} = \underline{b}$ hat keine Lösung falls $\underline{b} \notin \text{Bild}(\underline{A})$

Ziel: finde $\hat{\underline{A}}, \hat{\underline{x}}, \hat{\underline{b}}$ sodass $\hat{\underline{A}} \hat{\underline{x}} = \hat{\underline{b}}$ (d.h. $\hat{\underline{b}} \in \text{Bild}(\hat{\underline{A}})$)

$$\hat{\underline{A}}_{m \times n}, \quad \hat{\underline{b}} \in \mathbb{R}^n$$

"so nah wie möglich" an $\underline{A} \underline{x} = \underline{b}$

$$\underline{C} = \begin{bmatrix} \underline{A} & \underline{b} \end{bmatrix}; \quad \hat{\underline{C}} = \begin{bmatrix} \hat{\underline{A}} & \hat{\underline{b}} \end{bmatrix}$$

Finde $\hat{\underline{C}}$ so dass $\| \underline{C} - \hat{\underline{C}} \|_F = \min$ Frobenius norm.

$$\| \underline{M} \|_F^2 = \sum_{i,j} |m_{ij}|^2$$

mit $\hat{\underline{C}} \in \text{Bild}(\hat{\underline{A}})$.

Ben $\text{Rang } \underline{A} = n \Rightarrow \overset{\text{möchte auch}}{\text{Rang } (\hat{\underline{C}})} = n$.

Lösung: Niedrigrangapproximation an \underline{C} :

$$\underline{C} = \underline{U} \underline{\Sigma} \underline{V}^T = \sum_{j=1}^{n+1} \sigma_j \underline{u}_j \underline{v}_j^H$$

$$\text{Definiere } \hat{\underline{C}} = \sum_{j=1}^n \sigma_j \underline{u}_j \underline{v}_j^H$$

↳ optimale Approx. an \underline{C} in der Menge der Matrizen von Rang n ist!

Ausserdem \underline{v}_j sind orthogonal $\Rightarrow \hat{\underline{C}} \underline{v}_{n+1} = 0$

Falls $v_{n+1,n+1} \neq 0 \Rightarrow \hat{\underline{x}} = -\frac{1}{v_{n+1,n+1}} \underline{v}_{n+1}$
 $\hat{\underline{b}} = \hat{\underline{A}} \hat{\underline{x}}$

Beweis

$$\hat{\underline{A}} \hat{\underline{x}} = \hat{\underline{b}} \Leftrightarrow \hat{\underline{A}} \hat{\underline{x}} - \hat{\underline{b}} = 0 \Leftrightarrow \underbrace{\begin{bmatrix} \hat{\underline{A}} & \hat{\underline{b}} \end{bmatrix}}_{\in \mathbb{R}^{(n+1) \times n}} \begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} = 0$$

$$\Rightarrow \begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} \in \ker(\hat{\underline{C}}) \quad \Rightarrow$$

$$\dim \ker(\hat{\underline{C}}) = 1, \underline{v}_{n+1} \in \ker \hat{\underline{C}}$$

$$\begin{bmatrix} \hat{\underline{x}} \\ -1 \end{bmatrix} = k \underline{v}_{n+1}$$

$$\text{Definier } k = -\frac{1}{v_{n+1, n+1}} \Rightarrow \hat{\underline{x}} = -\frac{1}{v_{n+1, n+1}} \underline{v}_{n+1}.$$

§ 8.5. Nichtlineare Ausgleichsrechnung.

$$\text{Modell } f(t, \underline{x}) = y$$

↳ Parameter, aus Messungen zu bestimmen.

$$f(t_i, \underline{x}) \approx y_i, \quad i = 1, 2, \dots, m$$

$$\underline{F}(\underline{x}) = \begin{bmatrix} f(t_1, \underline{x}) - y_1 \\ f(t_2, \underline{x}) - y_2 \\ \vdots \\ f(t_m, \underline{x}) - y_m \end{bmatrix}$$

$$\underline{F}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Finde $\underline{x}^* \in \mathbb{R}^n$ sodass
 $\|\underline{F}(\underline{x}^*)\|_2^2$ Minimal!

$$\Phi(\underline{x}) = \frac{1}{2} \|\underline{F}(\underline{x})\|_2^2, \quad \Phi: \mathbb{R}^n \rightarrow [0, \infty)$$

zu minimieren!

Finde \underline{x}^* : $\nabla \Phi(\underline{x}^*) = 0 \leftarrow$ Nullstellensuche!

→ mit Newton-Verfahren

↳ konvergiert lokal quadratisch.

$$D\phi(\underline{x}) = D\left(\frac{1}{2} \underline{F}(\underline{x})^T \underline{F}(\underline{x})\right) = \underline{D}F(\underline{x})^T \underline{F}(\underline{x}) \quad \left. \vphantom{D\phi(\underline{x})} \right\} \Rightarrow$$

$$D\phi(\underline{x}^*) = 0$$

$$\underline{D}F(\underline{x}^*)^T \underline{F}(\underline{x}^*) = 0$$

~~~~~ dafür müssen wir Newton Verfahren anwenden.

In Newton brauchen wir die Ableitung davon.

$$D(D\phi(\underline{x})) = H\phi(\underline{x}) \text{ Hesse Matrix von } \phi$$

$$H\phi(\underline{x}) = D\left(\underline{D}F(\underline{x})^T \underline{F}(\underline{x})\right) =$$

$$= \underline{D}F(\underline{x})^T \underline{D}F(\underline{x}) + \sum_{j=1}^m \underbrace{F_j(\underline{x})}_{\text{Matrix}} \underline{D}^2 F_j(\underline{x})$$

$$(H\phi(\underline{x}))_{ik} = \sum_{j=1}^m \left( \frac{\partial F_j}{\partial x_k}(\underline{x}) \frac{\partial F_j}{\partial x_i}(\underline{x}) + F_j(\underline{x}) \frac{\partial^2 F_j(\underline{x})}{\partial x_i \partial x_k} \right)$$

Newton-Schritt  $\underline{x}^{(k+1)} = \underline{x}^{(k)} + \underline{\Delta}$  mit  
Newton-Korrektur  $\underline{\Delta}$  aus LGS

$$H\phi(\underline{x}^{(k)}) \underline{\Delta} = -\underline{D}F(\underline{x}^{(k)})^T \underline{F}(\underline{x}^{(k)})$$

Alternative: Gauss-Newton-Verfahren

Idee: linearisiere lokal in jedem Schritt  
einer Iteration  $\Rightarrow$  eine Folge von  
linearen Ausgleichsproblemen

linearisiere

$$\argmin_{\underline{x} \in \mathbb{R}^n} \|\underline{F}(\underline{x})\|_2^2 \approx \argmin_{\underline{x} \in \mathbb{R}^n} \|\underline{F}(\underline{x}^{(0)}) + \underline{D}F(\underline{x}^{(0)}) (\underline{x} - \underline{x}^{(0)})\|_2^2$$

$$= \argmin_{\underline{x} \in \mathbb{R}^n} \|\underline{A} \underline{x} - \underline{b}\|_2^2$$

$$\text{mit } \underline{A} = \underline{D}F(\underline{x}^{(0)}) \text{ und } \underline{b} = \underline{F}(\underline{x}^{(0)}) - \underline{D}F(\underline{x}^{(0)}) \underline{x}^{(0)}$$

$$\underline{x}^{(0)} \xrightarrow{\text{lin LSP}} \underline{x}^{(1)} \xrightarrow{\text{lin LSP}} \underline{x}^{(2)} \rightarrow \dots \rightarrow \underline{x}^{(k)}$$

$$\underline{x}^{(k+1)} = \underline{x}^{(k)} + \underline{\Delta}^{(k)}, \quad \underline{\Delta}^{(k)} = \arg \min_{\underline{\Delta}} \left\| \underline{F}(\underline{x}^{(k)}) + \underline{D}(\underline{x}^{(k)}) \underline{\Delta} \right\|_2^2$$

Nachteil: lineare Konvergenz.

Professionelle Software: Levenberg-Marquadt

## §9 Eigenwerte

### §9.1 Motivation & Grundlagen

\* kein Alg.-die exakt die EW, EV berechnet  
⇒ iterativ

\* Software →  $\text{eig}(\underline{A})$  kostet  $O(N^3)$  für  $N \times N$   
 $\text{eigh}(\underline{A})$  kostet  $O(N^2)$  für  $\underline{A}^H = \underline{A}$

↳ QR-Algorithmus mit shift

\* EV werden typischerweise nicht so gut approximiert!

## §9.2 Potenzmethoden (C) V. Gradinaru $\underline{A}$ $n \times n$ Matrix.

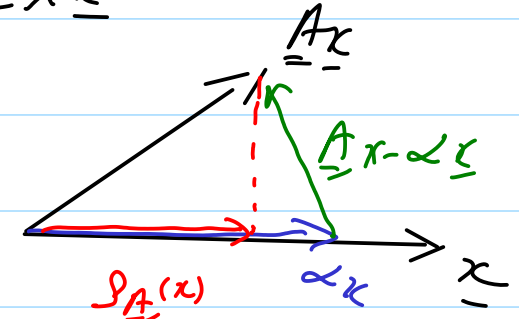
Def Rayleigh-Quotienten:

$$\rho_{\underline{A}}(\underline{x}) = \frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}}, \quad \underline{x} \neq \underline{0}$$

Bem

Falls  $\underline{x}$  EV von  $\underline{A}$ :  $\underline{A} \underline{x} = \lambda \underline{x}$

$$\rho_{\underline{A}}(\underline{x}) = \frac{\underline{x}^H \lambda \underline{x}}{\underline{x}^H \underline{x}} = \lambda$$



Bem

$$\rho_{\underline{A}}(\underline{x}) = \underset{\alpha}{\operatorname{argmin}} \|\underline{A} \underline{x} - \alpha \underline{x}\|_2$$

Ü: Beweis.

$$\text{Bem } D \rho_{\underline{A}}(\underline{x}) = \frac{(\underline{x}^T \underline{x})(\underline{A}^T + \underline{A}) \underline{x} - 2(\underline{x}^T \underline{A} \underline{x}) \underline{x}}{(\underline{x}^T \underline{x})^2} = \frac{(\underline{A}^T + \underline{A}) - 2 \rho_{\underline{A}}(\underline{x}) \underline{I}}{(\underline{x}^T \underline{x})} \underline{x}$$

$$\underline{x} \text{ EV von } \underline{A} \Rightarrow D \rho_{\underline{A}}(\underline{x}) = \frac{(\underline{x}^T \underline{x}) 2 \lambda \underline{x} - 2 \lambda (\underline{x}^T \underline{x}) \underline{x}}{(\underline{x}^T \underline{x})^2} = 0$$

$\Rightarrow EV \Rightarrow$  ein Stationärpunkt von  $f_{\underline{A}}(\cdot)$   
Taylor für  $f_{\underline{A}}(\cdot)$  um ein  $EV$ :

$$f_{\underline{A}}(\underline{x}) - f_{\underline{A}}(EV) = O(\|\underline{x} - EV\|_2^2) \quad \text{für } \underline{x} \text{ nah} \\ \text{an } EV.$$

Direkte Potenzmethode

Ziel: finde das betragsgrösste EW von  $\underline{A}$   
und ein EV dazu

Annahme  $\underline{A}$  diagonalisierbar:  $\underline{S}^{-1} \underline{A} \underline{S} = \text{diag}(\lambda_1, \dots, \lambda_n)$

$$|\lambda_1| \boxed{>} |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n|$$

EV stehen in den Spalten von  $\underline{S}$ ,  $\|\underline{s}_j\|_2 = 1$

$$\mathbb{R}^n \ni \underline{x} = \sum_{j=1}^n a_j \underline{s}_j \quad \text{mit } a_1 \neq 0, \text{ sonst beliebig}$$

$$\underline{A} \underline{x} = \sum_{j=1}^n a_j \underline{A} \underline{s}_j = \sum_{j=1}^n a_j \lambda_j \underline{s}_j$$

$$\underline{A}^2 \underline{x} = \sum_{j=1}^n a_j \lambda_j^2 \underline{s}_j$$

$$\dots$$

$$\underline{A}^k \underline{x} = \sum_{j=1}^n a_j \lambda_j^k \underline{s}_j = \lambda_1^k \left( a_1 \underline{s}_1 + \sum_{j=2}^n \left( \frac{\lambda_j}{\lambda_1} \right)^k a_j \underline{s}_j \right)$$

$$\left| \frac{\lambda_j}{\lambda_1} \right| < 1 \Rightarrow \left| \frac{\lambda_j}{\lambda_1} \right|^k < 1$$

Bem Für grosses  $k$  zeigt  $\underline{A}^k \underline{x}$  in der Richtung  
von  $\underline{s}_1$

$$\frac{\underline{A}^k \underline{x}}{\|\underline{A}^k \underline{x}\|_2} \rightarrow \pm \underline{s}_1$$

Idee: notiere  $\underline{x}_k = \underline{A}^k \underline{x}$  und berechne  $f_{\underline{A}}(\underline{x}_k) =$



$$= \frac{\underline{x}_k^H \underline{A} \underline{x}_k}{\underline{x}_k^H \underline{x}_k} = \frac{1}{\underline{x}_k^H \underline{x}_k} \left( \underline{x}_k^H \sum_{j=1}^n a_j \lambda_j^{k+1} \underline{s}_j \right) =$$

$$= \lambda_1 + O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^k\right)$$

## Potenzmethode

Wähle  $\underline{x}_0$  zufällig,  $\|\underline{x}_0\|=1$

für  $k=1, 2, \dots$

$$\underline{w} = \underline{A} \underline{x}_{k-1}$$

$$\underline{x}_k = \frac{\underline{w}}{\|\underline{w}\|}$$

$$\lambda = \underline{x}_k^H \underline{A} \underline{x}_k$$

Bem. Falls  $\underline{A}$  normal  $\Rightarrow$  EV orthogonal  
 $\Rightarrow$  Fehler  $O\left(\left|\frac{\lambda_2}{\lambda_1}\right|^{2k}\right)$

$\Rightarrow$  quadratische Konvergenz

$\hookrightarrow$  Grundlage für PageRank-Alg. von Google!

Beweis  $\underline{A} = \underline{A}^H$  ONB von EV

$$\underline{x} = \sum_{j=1}^n c_j \underline{u}_j = \underline{U} \underline{c} \quad \text{mit } \underline{U} \text{ unitär.}$$

$$\frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}} = \frac{\underline{c}^H \underline{U}^H \underline{A} \underline{U} \underline{c}}{\underline{c}^H \underline{U}^H \underline{U} \underline{c}} = \frac{\underline{c}^H \underline{A} \underline{c}}{\underline{c}^H \underline{c}} =$$

$$= \frac{\lambda_1 |c_1|^2 + \lambda_2 |c_2|^2 + \dots + \lambda_n |c_n|^2}{|c_1|^2 + \dots + |c_n|^2} = \lambda_1 + \delta \cdot \left|\frac{\lambda_2}{\lambda_1}\right|^{2k} + \dots$$

Theorem Potenzmethode liefert eine Iteration die  
linear konvergiert gegen  $\lambda_1$   
mit der Rate  $\left|\frac{\lambda_2}{\lambda_1}\right|$ .

Ziel Finde das kleinste EW:

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 EW:  $\lambda_n \Rightarrow \frac{1}{\lambda_n} \underline{x} = \underline{A}^{-1} \underline{x}$

$\frac{1}{\lambda_n}$  ist der betragsgrösste EW von  $\underline{A}^{-1}$

$\Rightarrow$  "inverse Potenzmethode" = Potenzmethode für  $\underline{A}^{-1}$

Bem Wir berechnen nicht  $\underline{A}^{-1}$  sondern nur einmal eine LU-Zerlegung von  $\underline{A}$  (strukturhaltend)

Dann löse LGS mit Matrizen  $L, U$  um  $\underline{A}^{-1} \underline{x}$  zu implementieren.

Ziel Gegeben  $\alpha \in \mathbb{C}$ , finde EW nah an  $\alpha$

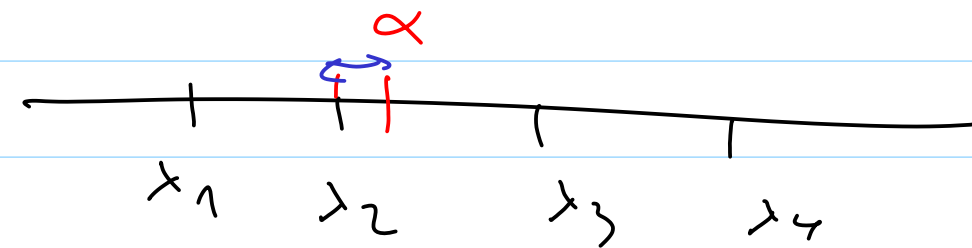
$$|\alpha - \lambda| = \min \{ |\alpha - \mu| \text{ mit } \mu \in \text{EW von } \underline{A} \}$$

$$\underline{A}\underline{x} = \lambda \underline{x} \Leftrightarrow \underline{A}\underline{x} - \alpha \underline{I}\underline{x} = (\lambda - \alpha) \underline{x} \Leftrightarrow$$

$$(\underline{A} - \alpha \underline{I}) \underline{x} = (\lambda - \alpha) \underline{x} \Leftrightarrow \frac{1}{\lambda - \alpha} \underline{x} = (\underline{A} - \alpha \underline{I})^{-1} \underline{x}$$

$\Rightarrow$  Potenzmethode für  $(\underline{A} - \alpha \underline{I})^{-1} \Rightarrow \frac{1}{\lambda - \alpha} \Rightarrow \lambda$

"shifted inverse iteration"



Bem Die Potenzmethode ist schneller wenn  $\alpha \approx \lambda_j$

$\Rightarrow$  Idee: wähle  $\alpha$  adaptiv, z.B.  $\alpha = \rho_{\underline{A}}(x^{(k-1)})$   
im  $k$ -ten Schritt

⇒ beschleunigte Konvergenz

Rayleigh-Quotienten-Iteration.

Ben Wir brauchen immer einen guten Startwert.

z.B. für RQ1 einige Schritte von shifted inverse iteration.

⇒ Konvergenzordnung 3!!!

$$\underline{A} = \underline{A}^H; \quad \underline{x} = \sum_{j=1}^n c_j \underline{u}_j = \underline{U} \underline{c} \quad \text{mit } \underline{u}_1, \dots, \underline{u}_n \in V$$

$\underbrace{\quad}_{\underline{x} \text{ beliebig.}}$

$$f_{\underline{A}}(\underline{x}) = \frac{\underline{x}^H \underline{A} \underline{x}}{\underline{x}^H \underline{x}} = \frac{\lambda_1 |c_1|^2 + \dots + \lambda_n |c_n|^2}{\|\underline{x}\|_2^2}$$

$(\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n)$

$$\lambda_1 \leq f_{\underline{A}}(\underline{x}) \leq \lambda_n$$

$$f_{\underline{A}}(\underline{x}) \in [\lambda_1, \lambda_n]$$

$$\lambda_1 = \min_{\underline{x} \in \mathbb{C}^n} f_{\underline{A}}(\underline{x}) \quad \text{erreicht} \quad \underline{c} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$\lambda_n = \max_{\underline{x} \in \mathbb{C}^n} f_{\underline{A}}(\underline{x}) \quad \text{erreicht} \quad \underline{c} = \begin{bmatrix} 0 \\ \vdots \\ 1 \end{bmatrix}$$

Beh  $\underline{x} \in \text{span}\{\underline{u}_1\}^\perp \Rightarrow c_1 = 0, \underline{x} = c_2 \underline{u}_2 + \dots + c_n \underline{u}_n$

$\downarrow$

$$f_{\underline{A}}(\underline{x}) \geq \lambda_2 \quad \text{erreicht} \quad \text{für} \quad \underline{c} = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\lambda_2 = \min_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \text{span}\{\underline{u}_1\}^\perp}} f_{\underline{A}}(\underline{x}) \quad ; \quad \lambda_{n-1} = \max_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \text{span}\{\underline{u}_n\}^\perp}} f_{\underline{A}}(\underline{x})$$

und so weiter

$$\lambda_k = \min_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \text{span}\{\underline{u}_1, \dots, \underline{u}_{k-1}\}^\perp}} f_{\underline{A}}(\underline{x}) = \min_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \underbrace{\text{span}\{\underline{u}_k, \dots, \underline{u}_n\}}_{\text{span}}}} f_{\underline{A}}(\underline{x}) =$$

$$= \max_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \text{span}\{\underline{u}_1, \dots, \underline{u}_k\}}} f_{\underline{A}}(\underline{x}) = \max_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \in \text{span}\{\underline{u}_{k+1}, \dots, \underline{u}_n\}^\perp}} f_{\underline{A}}(\underline{x})$$

Theorem [Courant-Fisher]

$$\lambda_k = \min_{\dim U = k} \max_{x \in U} \underline{S}_A(x) = \max_{\dim U = n-k+1} \min_{x \in U} \underline{S}_A(x)$$

Konsequenz  $\underline{Q}_m \in \mathbb{C}^{n \times m}$  mit  $\underline{Q}_m^H \underline{Q}_m = I$   
erweitere  $\underline{Q}_m$  auf ONB in  $\mathbb{C}^n$

$$\begin{bmatrix} \underline{Q}_m \\ \hat{\underline{Q}}_m \end{bmatrix} \begin{bmatrix} \hat{\underline{Q}}_m^H \\ \hat{\underline{Q}}_m^H \end{bmatrix} = \hat{\underline{Q}} \in \mathbb{C}^{n \times n}$$

Theorem [Cauchy]

$$\underline{A} = \begin{bmatrix} \underline{H} & \underline{B}^H \\ \underline{B} & \underline{R} \end{bmatrix} \quad \text{mit } \underline{H} \in \mathbb{C}^{m \times m} \quad \text{mit EW } \theta_1, \dots, \theta_m$$

(mit  $m \ll n$ )

$$\underline{Q}^H \underline{A} \underline{Q} = \begin{bmatrix} \underline{Q}_m^H \underline{A} \underline{Q}_m & \underline{Q}_m^H \underline{A} \hat{\underline{Q}}_m \\ \hat{\underline{Q}}_m^H \underline{A} \underline{Q}_m & \hat{\underline{Q}}_m^H \underline{A} \hat{\underline{Q}}_m \end{bmatrix}$$

hat dieselben EW wie  $\underline{A}$

Dann  $\lambda_k \leq \theta_k \leq \lambda_{k+n-m}$

Idee Für  $m \ll n$ , wähle  $\underline{Q}_m$  so dass

$$\text{Bild } \underline{Q}_m \approx \text{span} \{ \underline{u}_1, \dots, \underline{u}_m \}$$

$\searrow \searrow$  EV von  $\underline{A}$

$\Rightarrow$  gute Approximation  $\theta_k \approx \lambda_k$  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 eig zu langsam

ausserdem QR-Alg. zerstört Struktur

Krylov-Verfahren: grosse, dünn besetzte Matrizen

Def Sei  $0 \neq \underline{z} \in \mathbb{C}^n$ ,  $\underline{A} \in \mathbb{C}^{n \times n}$

$$\underline{\mathcal{K}}_l(\underline{A}, \underline{z}) := \text{span}\{ \underline{z}, \underline{A}\underline{z}, \underline{A}^2\underline{z}, \dots, \underline{A}^{l-1}\underline{z} \} =$$

$$= \{ p(\underline{A})\underline{z}; p = \text{Polynom vom Grad} \leq l-1 \}$$

Krylov-Raum

Suche eine ONB in  $\underline{\mathcal{K}}_l(\underline{A}, \underline{z})$

$$\underline{\mathcal{K}}_1 = \text{span}\{ \underline{z} \} \subset \text{span}\{ \underline{z}, \underline{A}\underline{z} \} = \underline{\mathcal{K}}_2 \subset \underline{\mathcal{K}}_3 \subset \dots \subset \underline{\mathcal{K}}_l$$

iterativ: gegeben  $\underline{v}_1, \dots, \underline{v}_l$  ONB in  $\underline{\mathcal{K}}_l(\underline{A}, \underline{z})$   
mit

$$\text{span}\{ \underline{v}_1, \dots, \underline{v}_j \} = \underline{\mathcal{K}}_j(\underline{A}, \underline{z}) \text{ für } j=1, 2, \dots, l$$

baue

$$\underline{v}_1, \dots, \underline{v}_l, \underline{v}_{l+1} \text{ ONB in } \underline{\mathcal{K}}_{l+1}(\underline{A}, \underline{z})$$

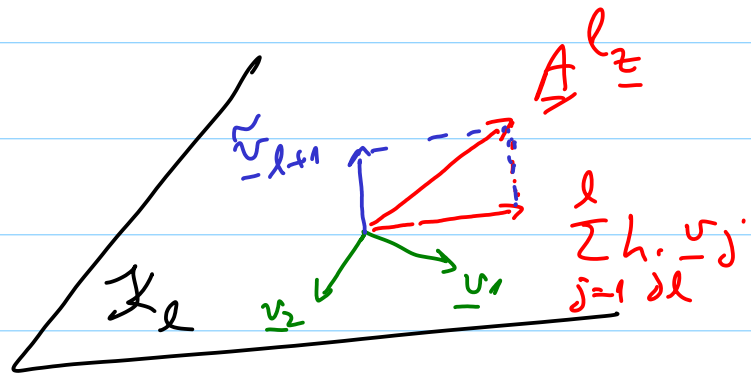
$$\underline{\mathcal{K}}_{l+1}(\underline{A}, \underline{z}) = \text{span}\{ \underline{z}, \underline{A}\underline{z}, \dots, \underline{A}^{l-1}\underline{z}, \underline{A}^l\underline{z} \}$$

entweder  $\underline{A}^l\underline{z} \in \underline{\mathcal{K}}_l(\underline{A}, \underline{z})$ , d.h.

$\underline{A}^l\underline{z}$  lin. abhängig von  $\underline{z}, \underline{A}\underline{z}, \dots, \underline{A}^{l-1}\underline{z}$

oder  $\underline{A}^l\underline{z} \notin \underline{\mathcal{K}}_l(\underline{A}, \underline{z}) \Rightarrow \underline{A}^l\underline{z} \in \underline{\mathcal{K}}_{l+1} \setminus \underline{\mathcal{K}}_l$

dann  $\underline{v}_{l+1}$  aus (modifizierten) Gram-Schmidt



$$\tilde{v}_{l+1} = A^l z_l - \sum_{j=1}^l h_{jl} v_j$$

$$v_{l+1} = \frac{\tilde{v}_{l+1}}{\|\tilde{v}_{l+1}\|}$$

Dabei sind  $h_{jl} = v_j^H A v_l$   
 (da  $v_1, \dots, v_l$  orthonormal zu  $\mathcal{K}_l$ )

## Algorithmus [Arnoldi Prozess]

$z$  = beliebig.

$$v_1 = z / \|z\|$$

für  $l = 1, 2, \dots, k-1$ :

$l=1$

$l=2$

$l=3$

$$\tilde{v} = A v_l$$

für  $j = 1, 2, \dots, l$ :

$$h_{jl} = v_j^H \tilde{v}$$

$h_{11}$

$h_{12}$

$h_{13}$

$h_{22}$

$h_{23}$

$h_{33}$

$h_{21}$

$h_{32}$

$h_{43}$

mod. GS

$$\tilde{v} = \tilde{v} - h_{jl} v_j$$

$$h_{l+1,l} = \|\tilde{v}\|$$

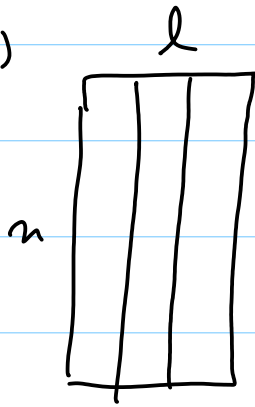
$$v_{l+1} = \tilde{v} / h_{l+1,l}$$

$$\begin{bmatrix} \times & \times & \times & \dots \\ \times & \times & \times & \dots \\ & \times & \times & \\ & & \times & \end{bmatrix}$$

Beim Falls  $h_{l+1,l} = 0 \Rightarrow$  Abbruch der Iteration

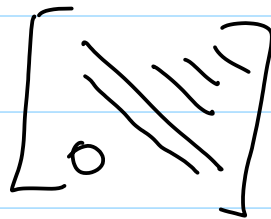
$$\underline{A} \underline{v}_l \in \mathcal{X}_l(\underline{A}, \underline{z})$$

$$\underline{V}_l = \begin{bmatrix} \underline{v}_1 & \underline{v}_2 & \dots & \underline{v}_l \end{bmatrix} \in \mathbb{C}^{n \times l}$$



$$\underline{\tilde{H}}_l = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ 0 & h_{32} & h_{33} \\ \hline 0 & 0 & h_{43} \end{bmatrix} \in \mathbb{C}^{l+1, l}$$

$$\underline{H}_l \in \mathbb{C}^{l, l}$$



$\underline{H}_l$  obere Hessenbergmatrix

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:

$$\underline{A} \underline{V}_l = \underline{V}_{l+1} \underline{\tilde{H}}_l = \begin{bmatrix} \underline{v}_{l+1} \\ 0 \dots 0 \end{bmatrix} \begin{matrix} h_{l+1,l} \\ \vdots \\ h_{1,l} \end{matrix} + \underline{V}_l \underline{H}_l$$

$$\begin{matrix} n & l & l+1 & l \\ \underline{A} & \underline{V}_l & \underline{V}_{l+1} & \underline{\tilde{H}}_{l+1} \end{matrix} = \begin{bmatrix} * & & \\ & * & \\ & 0 & \ddots \end{bmatrix}_{l+1}$$

Bez 1)  $\underline{V}_l^H \underline{V}_l = \underline{I}_l$   $\underline{V}_l^H \underline{V}_l = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}_l$



2)  $\underline{V}_l^H \cdot | \Rightarrow \underline{V}_l^H \underline{A} \underline{V}_l = \underline{V}_l^H \underline{v}_{l+1} \boxed{0 \dots 1 \dots k} + \underline{V}_l^H \underline{V}_l \underline{H}_l$   
 $= \underline{0} + \underline{I}_l \underline{H}_l = \underline{H}_l$

$\Rightarrow \underline{V}_l^H \underline{A} \underline{V}_l = \underline{H}_l$

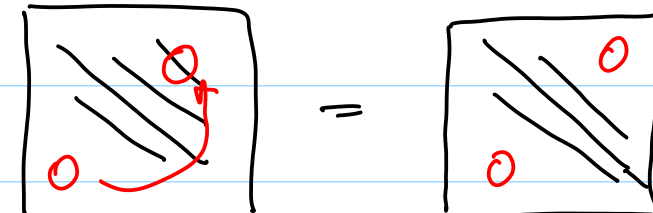
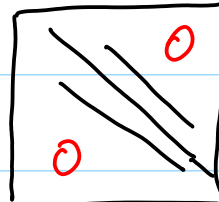
$\boxed{\underline{V}_l^H} \boxed{\underline{A}} \boxed{\underline{V}_l} = \boxed{\underline{H}_l}$

3) falls  $h_{l+1,l} = 0 \Rightarrow \mathcal{K}_{l+1} = \mathcal{K}_l$  und

$\underline{A} \underline{V}_l = \underline{V}_l \underline{H}_l$

4) falls  $\underline{A}^H = \underline{A}$  ( $\underline{A}$  Hermito-symmetrisch):

$\underline{H}_l^H = (\underline{V}_l^H \underline{A} \underline{V}_l)^H = \underline{V}_l^H \underline{A}^H \underline{V}_l = \underline{V}_l^H \underline{A} \underline{V}_l = \underline{H}_l$   
 $\Rightarrow \underline{H}_l$  Hermito-symmetrisch und obere Hessenberg

 =   $\Downarrow$  tridiagonal  
 $\underline{\beta}$   
 $\underline{\alpha}$

$\Rightarrow$  es reicht, die Vektoren  $\underline{\alpha}, \underline{\beta}$  um  $\underline{H}_l$  zu speichern.

und die innere Schleife im Arnoldi-Prozess hat die Länge 2

$\tilde{\underline{v}}_{l+1} = \underline{A} \underline{v}_l - h_{l,l} \underline{v}_l - h_{l-1,l} \underline{v}_{l-1}$

Langos-Verfahren  $O(nk)$

Arnoldi-Verfahren  $O(nk^2)$

Allgemeinen Namen: 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  $\underline{y} \in \mathbb{C}^l$  so dass  $\underline{A} \underline{x} = \underline{b}$  mit  $\underline{x} = \underline{V}_l \underline{y}$

Skript: einfache Implementierung.

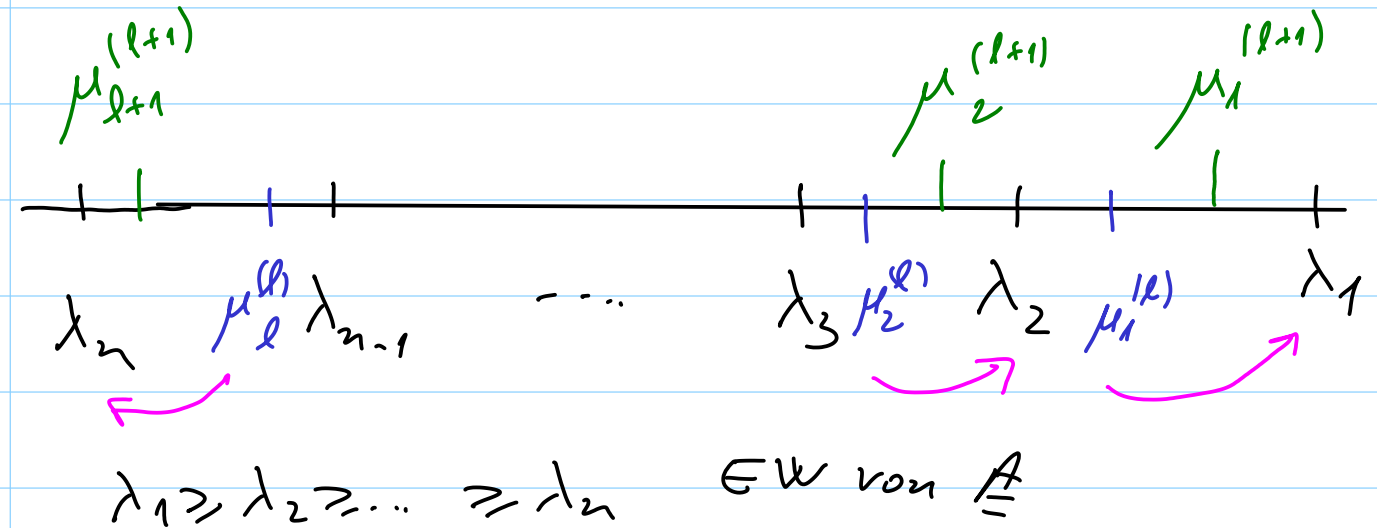
ARPACK  $\rightarrow$  eigvals

**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

$$\mu_j^{(\ell)} \leq \mu_j^{(\ell+1)} \leq \lambda_j.$$



$$\mu_1^{(\ell)} \geq \mu_2^{(\ell)} \geq \dots \geq \mu_\ell^{(\ell)} \text{ EW von } \underline{H}_\ell$$

$$\mu_1^{(\ell+1)} \geq \mu_2^{(\ell+1)} \geq \dots \geq \mu_{\ell+1}^{(\ell+1)} \text{ EW von } \underline{H}_{\ell+1}$$

# §10 Lineare Anfangswertprobleme

1. Fall

$$\begin{cases} \dot{\underline{y}} = \underline{A} \underline{y} \\ \underline{y}(0) = \underline{y}_0 \end{cases} \quad \text{Falls } \underline{A} \text{ diagonalisierbar}$$

$$\underline{A} = \underline{S}^{-1} \underline{D} \underline{S}$$

Variablenwechsel  $\hat{\underline{y}} = \underline{S}^{-1} \underline{y} \Rightarrow$  entkoppeln

$$\begin{cases} \dot{\hat{y}}_1 = \lambda_1 \hat{y}_1 \\ \dots \\ \dot{\hat{y}}_d = \lambda_d \hat{y}_d \end{cases} \Rightarrow \hat{y}_i(t) = (\underline{S}^{-1} \underline{y}_0)_i e^{\lambda_i t} \text{ für } t \in \mathbb{R}$$

$$\underline{y}(t) = \underline{S} \begin{bmatrix} e^{\lambda_1 t} & & \\ & \ddots & \\ & & e^{\lambda_d t} \end{bmatrix} \underline{S}^{-1} \underline{y}_0$$

OK für kleines  $d$  oder für exakt/analytisch diagonalisierbare  $\underline{A}$ , sonst (z.B.  $d \geq 5$ )  $\rightarrow$  instabil.

Für  $d = 5, 6, \dots, 20; 50, 100$ 

$$\underline{y}(t) = \underline{e}^{\underline{A}t} \underline{y}_0 \quad \text{Padé-Approximation}$$

$$\hookrightarrow \expm(\underline{A}t) \underline{y}_0$$

$d$  gross,  $\underline{A}$  dünn besetzt: Krylov-Verfahren.

Krylov: für  $\underline{A}$  gibt es  $\underline{V} \in \mathbb{C}^{d \times m}$  mit orthonormalen Spalten

$$\underline{V}_m^H \underline{A} \underline{V}_m = \underline{H}_m \quad m \times m \quad \text{mit } m \ll d$$

$\hookrightarrow$  obere Hessenberg Matrix.

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

??  $\underline{V}$

Spalten von  $\underline{V}$ 

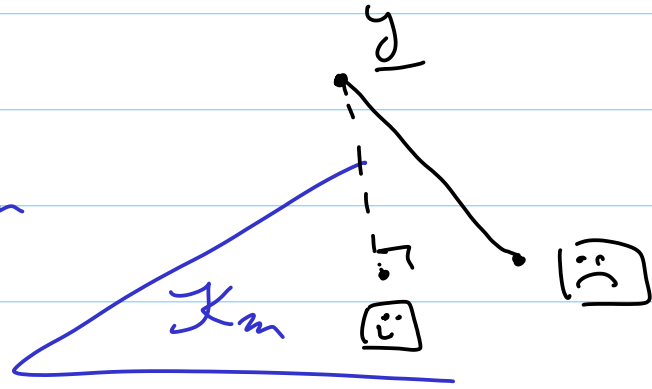
$$\text{span}\{\underline{v}_1, \dots, \underline{v}_m\}$$

ONB in  $\mathbb{R}^m$ 

$$\underline{u}_m(t) \in \mathcal{K}_m(\underline{A}, \underline{y}_0) = \text{span}\{\underline{y}_0, \underline{A}\underline{y}_0, \dots, \underline{A}^{m-1}\underline{y}_0\}$$

$$\dot{\underline{y}}(t) - \underline{A} \underline{y}(t) = 0$$

$$\left( \dot{\underline{u}}_m(t) - \underline{A} \underline{u}_m(t) \right) \perp \mathcal{K}_m$$



$$\langle \underline{w}, \dot{\underline{u}}_m(t) - \underline{A} \underline{u}_m(t) \rangle = 0 \quad \text{für alle } \underline{w} \in \mathcal{K}_m.$$

Ersetze  $\begin{cases} \dot{\underline{y}}(t) = \underline{A} \underline{y}(t) \\ \underline{y}(0) = \underline{y}_0 \end{cases}$  durch  $\begin{cases} \text{finde } \underline{u}_m(t) \in \mathcal{K}_m \text{ so dass} \\ \langle \underline{w}, \dot{\underline{u}}_m(t) - \underline{A} \underline{u}_m(t) \rangle = 0 \\ \text{für alle } \underline{w} \in \mathcal{K}_m \end{cases}$

$$\underline{u}_m(t) \in \mathcal{K}_m = \text{span} \{ \underline{v}_1, \dots, \underline{v}_m \} \Rightarrow$$

$$\underline{u}_m(t) = \sum_{k=1}^m c_k(t) \underline{v}_k$$

$$\underline{c}(t) = \begin{bmatrix} c_1(t) \\ \vdots \\ c_m(t) \end{bmatrix} \in \mathbb{C}^m$$

Einsetzen  $\Rightarrow$

$$\langle \underline{w}, \sum_{k=1}^m \dot{c}_k(t) \underline{v}_k - \sum_{k=1}^m c_k(t) \underline{A} \underline{v}_k \rangle = 0$$

$$\text{für alle } \underline{w} \in \mathcal{K}_m(\underline{A}, \underline{y}_0)$$

$$\text{Wähle } \underline{w} = \underline{v}_1 \Rightarrow$$

$$\langle \underline{v}_1, \sum_{k=1}^m \dot{c}_k(t) \underline{v}_k \rangle = \langle \underline{v}_1, \sum_{k=1}^m c_k(t) \underline{A} \underline{v}_k \rangle$$

$$\sum_{k=1}^m \dot{c}_k(t) \langle \underline{v}_1, \underline{v}_k \rangle = \sum_{k=1}^m c_k(t) \langle \underline{v}_1, \underline{A} \underline{v}_k \rangle$$

$$\underline{v}_1^H \underline{v}_k = 0 \text{ falls } k \neq 1 \\ 1 \text{ für } k=1$$

$$\underline{v}_1^H \underline{A} \underline{v}_k = (H_m)_{1k}$$

$$\Rightarrow \dot{c}_1(t) = \sum_{k=1}^m (H_m)_{1k} c_k(t)$$

Für  $\underline{w} = \underline{v}_2, \dots, \underline{v}_m$  analog  $\Rightarrow$

$$\begin{cases} \dot{\underline{c}}(t) = \underline{H}_m \underline{c}(t) \\ \underline{y}(0) = \underline{y}_0 \Rightarrow \underline{c}(0) = \|\underline{y}_0\| \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^m \end{cases} \quad \begin{array}{l} \text{m klein} \\ \Downarrow \\ \text{kann mit} \\ \text{Pade' lösen!} \end{array}$$

$$\Rightarrow \underline{c}(t) = \exp(\underline{H}_m t) \underline{c}(0)$$

$$\underline{u}_m(t) = \sum_{k=1}^m c_k(t) \underline{v}_k = \|\underline{y}_0\| \underline{v}_m e^{\underline{H}_m t} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

2. Fall  $\dot{\underline{y}}(t) = \underline{A} \underline{y}(t) + \underline{g}(t)$  inhomogener Fall

Variation der Konstanten:

$$\underline{y}(t) = e^{\underline{A}(t-t_0)} \underline{y}_0 + \int_{t_0}^t e^{\underline{A}(t-s)} \underline{g}(s) ds$$

3. Fall

$$\dot{\underline{y}}(t) = \underline{A}(t) \underline{y}(t) \quad \text{Magnus-Integratoren}$$

Unter bestimmten Voraussetzungen

$$\underline{y}(t) = e^{\underline{\Omega}(t, t_0)} \underline{y}_0 \quad \text{Magnus Entwicklung}$$

$$\underline{\Omega} = \sum_{k=1}^{\infty} \underline{\Omega}_k$$

$$[A, B] = AB - BA$$

$$\Omega_1 = \int_{t_0}^t \underline{A}(\tau_1) d\tau_1,$$

$$\Omega_2 = \frac{1}{2} \int_{t_0}^t \int_{t_0}^{\tau_1} [\underline{A}(\tau_1), \underline{A}(\tau_2)] d\tau_2 d\tau_1,$$

$$\Omega_3 = \frac{1}{12} \int_{t_0}^t \int_{t_0}^{\tau_1} \int_{t_0}^{\tau_2} [[\underline{A}(\tau_1), \underline{A}(\tau_2)], \underline{A}(\tau_3)] + [\underline{A}(\tau_1), [\underline{A}(\tau_2), \underline{A}(\tau_3)]] d\tau_3 d\tau_2 d\tau_1.$$

$$[\underline{A}(\tau_1), \underline{A}(\tau_2)] = \underline{A}(\tau_1) \underline{A}(\tau_2) - \underline{A}(\tau_2) \underline{A}(\tau_1)$$

Kommutator

Idee Statt  $\dot{\underline{y}} = \underline{A}(t) \underline{y}$  löse  $\dot{\underline{y}} = \underline{\hat{A}}(t) \underline{y}$

wobei  $\underline{\hat{A}}(t) = \sum_{i=1}^n l_i(t) \underline{A}(t_n + \tau_i h)$

$l_i(t)$  = Lagrange-Polynom in  $t_n + \tau_i h$   
 $\tau_i$  = Quadraturknoten in  $[0, 1]$

$$t_n + \tau_i h \in [t_n, t_n + h]$$

$\underline{\hat{A}}(t)$  = Polynom vom Grad  $n$  in  $t$  auf  $[t_n, t_n + h]$

$$\underline{\hat{A}}(t_n + \tau_i h) = \underline{A}(t_n + \tau_i h) \quad \text{für } i=1, 2, \dots, n$$

Dann Magnus-Entwicklung für  $\dot{\underline{y}} = \underline{\hat{A}}(t) \underline{y}$

Vorteil: Integration nur für Polynome in  $t \Rightarrow$   
 einfach exakt berechnen.

$\underline{\hat{A}}(t)$  glatt; man kann zeigen:

Rest in der Magnus-Entwicklung ( $\hat{\Sigma}$ ) ist  $O(h^5)$  nach 4 Terme.

Theorem  $(b_i, \tau_i)_{i=1, \dots, n}$  Quadraturformel der Ordnung  $p \geq 1$

$$y(t_n + h) - \hat{y}(t_n + h) = O(h^{p+1})$$

\* Siehe Methoden & Beispiel im Skript!

# §11. Exponentielle Integratoren

$$\begin{cases} \dot{\underline{y}} = \underline{f}(\underline{y}) & \text{autonom mit } \underline{f} \text{ stetig differenzierbar} \\ \underline{y}(0) = \underline{y}_0 \end{cases}$$

Idee der Linearisierung:  $\underline{J} = \underline{Df}(\underline{y}_0)$

$$\dot{\underline{y}} = \underbrace{\underline{J}\underline{y}}_{\text{linear}} + \underbrace{\underline{f}(\underline{y}) - \underline{J}\underline{y}}_{\underline{g}(\underline{y})}$$

Variation der Konstanten:

$$\underline{y}(h) = e^{\underline{J}h} \underline{y}_0 + \int_0^h e^{\underline{J}(h-s)} \underline{g}(\underline{y}(s)) ds$$

Ersetze  $\underline{y}(s)$  durch  $\underline{y}_0 \Rightarrow$  "Quadratü"/Approximati.

$$\begin{aligned} \int_0^h e^{\underline{J}(h-s)} \underline{g}(\underline{y}(s)) ds &\approx \int_0^h e^{\underline{J}(h-s)} \underline{g}(\underline{y}_0) ds \\ &= h \underline{f}(\underline{J}\underline{J}) \underline{g}(\underline{y}_0) \end{aligned}$$

$$\underline{f}(z) \stackrel{\text{def}}{=} \frac{e^z - 1}{z} = \sum_{n=1}^{\infty} \frac{1}{n!} z^{n-1} = \sum_{n=0}^{\infty} \frac{z^n}{(n+1)!}$$

für Matrizen:  $\underline{f}(\underline{A}) = \sum_{n=0}^{\infty} \frac{\underline{A}^n}{(n+1)!} = (e^{\underline{A}} - \underline{I}) \underline{A}^{-1}$

Beweis

$$\int_0^h e^{\underline{J}(h-s)} \underline{g}(\underline{y}_0) ds = \int_0^h \sum_{n=0}^{\infty} \frac{1}{n!} (\underline{J}(h-s))^n \underline{g}(\underline{y}_0) ds =$$

$$\int_0^h (h-s)^n ds$$

Umtauschen

$$= \sum_{n=0}^{\infty} \frac{1}{n!} \underline{J}^n (-1) \frac{(h-s)^{n+1}}{n+1} \bigg|_{s=0}^{s=h} \underline{g}(\underline{y}_0) =$$

$$\sum_{n=0}^{\infty} \frac{1}{n!} \frac{\partial^n}{\partial} \frac{h^{n+1}}{n+1} g(\underline{y}_0) = h \underbrace{\sum_{n=0}^{\infty} \frac{1}{(n+1)!} \frac{\partial^n}{\partial} h^n g(\underline{y}_0)}_{= h f(\underline{\partial} h) g(\underline{y}_0)}$$

Bez Definition von  $g(\underline{y}_0) = f(\underline{y}_0) - \underline{\partial} \underline{y}_0$

$$h f(\underline{\partial} h) g(\underline{y}_0) = h f(\underline{\partial} h) f(\underline{y}_0) - h f(\underline{\partial} h) \cdot \underline{\partial} \underline{y}_0$$

$$h f(\underline{\partial} h) \underline{\partial} \underline{y}_0 = h \sum_{n=0}^{\infty} \frac{1}{(n+1)!} \frac{\partial^n}{\partial} h^n \underline{\partial} \underline{y}_0 = e^{\underline{\partial} h} \underline{y}_0 - \underline{y}_0$$

⇒

$$h f(\underline{\partial} h) g(\underline{y}_0) = h f(\underline{\partial} h) f(\underline{y}_0) - e^{\underline{\partial} h} \underline{y}_0 + \underline{y}_0$$

Somit bekommen wir

$$\underline{y}(h) \approx \underbrace{e^{\underline{\partial} h} \underline{y}_0}_{\text{exponentielles Eulerverfahren.}} + h f(\underline{\partial} h) f(\underline{y}_0) - \underbrace{e^{\underline{\partial} h} \underline{y}_0 + \underline{y}_0}$$

$$\Rightarrow \underline{y}(h) \approx \underline{y}_0 + h f(\underline{\partial} h) f(\underline{y}_0)$$

exponentielles Eulerverfahren.

$$\underline{\partial} = \underline{D} f(\underline{y}_0)$$

$$f(h \underline{\partial}) = (e^{h \underline{\partial}} - \underline{I})(h \underline{\partial})^{-1}$$

teuer für grosses  $d$   $O(d^3)$

Bez Für grosses  $d$  kann man das Krylov-Verfahren für  $e^{h \underline{\partial}}$  verwenden!



$$\ell(\underline{A}) \underline{b} = \underline{V}_m \ell(\underline{H}_m) \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

aus Krylov-Verfahren mit  $\mathcal{K}_m(\underline{A}, \underline{b})$

Ben Stabilitätsfunktion  $S(z) = e^z$   
 $\Rightarrow$  das ideale Stabilitätsgebiet!

da exakt für  $\underline{\dot{y}} = \underline{A} \underline{y} + \underline{g}$   
 $\hookrightarrow$  konstante

Verallgemeinerung: exponentielle RK-Verfahren.

$$\underline{\dot{f}} = \underline{D} f(x)$$

semi-implizite Euler:  $\underline{y}_1 = \underline{y}_0 + \boxed{(\underline{I} - h \underline{J})^{-1}} h \underline{f}(\underline{y}_0)$

exponentielle Euler:  $\underline{y}_1 = \underline{y}_0 + \boxed{\ell(h \underline{J})} h \underline{f}(\underline{y}_0)$

Idée: ersetze  $\frac{1}{1-z}$  durch  $\ell(z) = \frac{e^z - 1}{z}$  in Row  $\Rightarrow$

$$\begin{cases} \underline{k}_i = \ell(h \underline{J}) \left( \underline{f}(\underline{u}_i) + h \underline{J} \sum_{j=1}^{i-1} \alpha_{ij} \underline{k}_j \right) \\ \underline{u}_i = \underline{y}_0 + h \sum_{j=1}^{i-1} \alpha_{ij} \underline{k}_j \\ \underline{y}_1 = \underline{y}_0 + h \sum_{i=1}^s b_i \underline{u}_i \end{cases}$$

explizite RK:  $\ell(z) = 1$  und  $\underline{J} = \underline{0}$

Row :  $\ell(z) = \frac{1}{1-z}$

exp. RK :  $\ell(z) = \frac{e^z - 1}{z}$

| ODEs       | nicht steif              | steif                                                                  | ostillierend                                                                           |
|------------|--------------------------|------------------------------------------------------------------------|----------------------------------------------------------------------------------------|
| Methode    | expl. RK<br>ode45        | impl. RK<br>ode23s                                                     | exponentielle RK<br><i>exp4</i>                                                        |
| Stabilität | $h < \frac{1}{\lambda}$  | alle $h$                                                               | alle $h$                                                                               |
| Implement. | <u><math>f(y)</math></u> | <u><math>Df(y)</math></u><br>lösen<br>nicht.-lin.<br>Gleichungssysteme | <u><math>Df(y)</math></u><br>$\ell(\underline{A})\underline{b}$<br>schnell mit Krylov. |

⊕ Erhaltungseigenschaften wichtig?

⊕ Autonom?

⊕ Ordnung der Dgl.?