

# Numerik WS 2021/22

Ida Hönigmann

April 17, 2024

## Contents

<b>1</b>	<b>Grundbegriffe</b>	<b>2</b>
1.1	Gegenstand der Numerischen Mathematik . . . . .	2
1.2	Kondition und Stabilität . . . . .	2
1.3	Verfahrensfehler . . . . .	3
<b>2</b>	<b>Interpolation</b>	<b>4</b>
2.1	Lagrange-Polynominterpolation . . . . .	4
2.2	Cebysev-Knoten . . . . .	6
2.3	Lebesgue-Konstante . . . . .	8
2.4	Auswertung von Interpolationspol. . . . .	9
2.5	Hermite-Polynominterpolation . . . . .	11
2.6	Spline-Interpolation . . . . .	11
2.7	Diskrete und schnelle Fourier-Transformation . . . . .	15
<b>3</b>	<b>Extrapolation</b>	<b>17</b>
3.1	Richardson-Extrapolation . . . . .	17
3.2	Aitken'sches $\Delta^2$ -Verfahren . . . . .	19
<b>4</b>	<b>Numerische Integration</b>	<b>21</b>
4.1	Quadraturformeln . . . . .	21
4.2	Interpolatorische Quadraturformeln . . . . .	23
4.3	Gauss-Quadratur . . . . .	25
<b>5</b>	<b>Iterative Lösung von GLS</b>	<b>27</b>
5.1	Fixpunktprobleme . . . . .	27
5.2	Newton in $\mathbb{R}^d$ . . . . .	30
5.3	Stationäre Iterationsverfahren zur Lösung Linearer GLS . . . . .	33
5.4	Krylov-Verfahren zur Lsg Linearer GLS . . . . .	35
<b>6</b>	<b>Eliminationsverfahren</b>	<b>37</b>
6.1	Dreiecksmatrizen . . . . .	38
6.2	LU-Zerlegung . . . . .	39
6.3	Gauss-Elimination . . . . .	41
6.4	QR-Zerlegung . . . . .	44
6.5	Lineare Ausgleichsprobleme . . . . .	46
<b>7</b>	<b>Eigenwertprobleme</b>	<b>47</b>
7.1	Lineare Algebra + Stabilität . . . . .	47
7.2	Vektoriteration . . . . .	48
7.3	Orthogonale Iteration und QR-Zerlegung . . . . .	54
7.4	Hessenberg-Form einer Matrix . . . . .	56

# 1 Grundbegriffe

## 1.1 Gegenstand der Numerischen Mathematik

Von der Realität bis zur Interpretation einer Simulation ist es ein langer Weg.

- **Mathematisches Modell** versucht mit Hilfe von mathematischen Formeln (idR. Differentialgl.) die Realität zu beschreiben.
- Die wenigsten Lösungen dieser mathematischen Modelle kann man exakt berechnen, d.h. man approximiert die exakte Lösung mittels **numerischer Simulation** am Rechner.
- Diese **numerische Lösung** wird dann interpretiert und man hofft, dass diese Interpretation die Realität beschreibt.

Jede numerische Simulation zerfällt in kleinere **numerische Probleme**, die geeignet zu lösen sind. Die elementarsten numerischen Probleme sind Gegenstand dieser Vorlesung.

**Beispiel 1.** • *Wie approximiert man komplizierte Funktionen mittels einfacher Funktionen (z.B. stückweise Polynome)?*

- *Wie berechnet man Grenzwerte (z.B. Integral, Differential)?*
- *Wie löst man lineare / nichtlineare Gleichungen?*

Jede numerische Simulation ist fehlerbehaftet.

- **Modellfehler:** Das mathematische Modell vereinfacht die Realität.
- **Datenfehler:** Die Eingangsdaten einer Simulation stammen meistens aus physikalischen Messungen und haben daher eine gewisse Mess(un-)genauigkeit.
- **Rundungsfehler:** Auf Rechnern ersetzt die endliche Menge an Gleitkommazahlen das kontinuierliche  $\mathbb{R}$ , d.h. sowohl die Daten als auch die Rechnungen sind rundungsfehlerbehaftet.
- **Verfahrensfehler:** Viele Probleme werden mathematisch in unendlich dimensionalen Räumen oder mit Limiten formuliert. Beides steht im Rechner nicht zur Verfügung und muss diskretisiert werden.

In der Vorlesung liegt unser Hauptaugenmerk auf dem Verfahrensfehler und dem Aufwand zugehöriger Algorithmen.

## 1.2 Kondition und Stabilität

Betrachte ein abstraktes Problem. Werte  $\Phi : X \rightarrow Y$  bei  $x \in X$  aus, wobei  $X, Y$  geeignete normierte Räume sind. Die **Kondition eines Problems** besagt, wie stark Änderungen in  $x$  (z.B. Rundungsfehler) sich auf  $\Phi(x)$  auswirken.

**Definition 1.** Das Problem ist **schlecht konditioniert bzgl. absolutem Fehler**, wenn es eine kleine Störung  $\tilde{x}$  von  $x$  gibt mit  $\|\Phi(x) - \Phi(\tilde{x})\| \gg \|x - \tilde{x}\|$ .

Das Problem ist **schlecht konditioniert bzgl. relativem Fehler**, falls  $x \neq 0 \neq \Phi(x)$  und es ex. eine kleine Störung  $\tilde{x}$  von  $x$  gibt mit  $\frac{\|\Phi(x) - \Phi(\tilde{x})\|}{\|\Phi(x)\|} \gg \frac{\|x - \tilde{x}\|}{\|x\|}$ .

Andernfalls bezeichnet man das Problem als **gut konditioniert (bzgl. abs./rel. Fehler)**.

**Bemerkung 1.** Ist  $\Phi$  stetig differenzierbar, d.h.  $\Phi(x) - \Phi(\tilde{x}) = D\Phi(x)(x - \tilde{x}) + o(\|x - \tilde{x}\|)$  für  $\tilde{x} \rightarrow x$  so beschreibt die Ableitung  $D\Phi(x) \in L(X, Y)$  wie stark sich Änderungen in  $x$  auf den Fehler auswirken.

Deshalb bezeichnet man  $\kappa_{abs}(x) = \|D\Phi(x)\|$ ,  $\kappa_{rel}(x) = \frac{\|D\Phi(x)\| \cdot \|x\|}{\|\Phi(x)\|}$  als **Konditionszahlen (bzgl. abs./rel. Fehler)**, d.h. man ist gut konditioniert für  $\kappa_{abs}, \kappa_{rel}$  vergleichsweise klein.

**Definition 2.** Es sei  $\tilde{\Phi}$  eine algorithmische Umsetzung von  $\Phi$ . Der Algorithmus  $\tilde{\Phi}$  ist **instabil**, wenn es eine kleine Störung  $\tilde{x}$  von  $x$  gibt, sodass

$$\underbrace{\|\Phi(x) - \tilde{\Phi}(\tilde{x})\|}_{\text{tatsächlicher Fehler im Rechner}} \gg \underbrace{\|\Phi(x) - \Phi(\tilde{x})\|}_{\text{unvermeidlicher Fehler}}.$$

Andernfalls ist der Algorithmus stabil.

**Bemerkung 2.** Mir ist bewusst, dass die Symbolik  $\gg$  ("wesentlich größer") ungenauer ist, als Sie es aus anderen Vorlesungen kennen. Aber "schlecht konditioniert" und "instabil" hängt halt an der Genauigkeit der Daten und den Erfordernissen des Nutzers!

In der Vorlesung geht es primär um die Asymptotik, d.h. was könnte im Worst-Case passieren.

**Erinnerung/Warnung:** Die Arithmetik im Rechner erfüllt weder Assoziativität noch Distributivgesetz, d.h. die Reihenfolge (und Formulierung) der Rechenoperatoren spielt eine Rolle für Stabilität.

**Beispiel 2** (schlechte Kondition bei Auslöschung). Als **Auslöschung** bezeichnet man das Phänomen, dass bei Subtraktion zweier annähernd gleicher Zahlen im Rechner die hinteren Ziffern (welche rundungsbehaftet sind) signifikant werden. Der relative Fehler kann sogar beliebig groß werden, d.h.  $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}, (x, y) \mapsto x - y$  hat  $\kappa_{rel}(x, y) = \frac{\sqrt{2} \|(x, y)\|_2}{|x - y|} \gg 0$  für  $x \approx y$ .

**Achtung:** Oft ist ein Problem gut konditioniert, wird aber in Teilprobleme zerlegt (im Algorithmus) sodass der resultierende Algorithmus instabil wird.

**Beispiel 3.** Werte  $\Phi(x) = \frac{1}{x+1} - \frac{1}{x}$  für  $x \gg 0$  aus.

$$\begin{aligned}\Phi'(x) &= -\frac{1}{(x+1)^2} + \frac{1}{x^2} = \frac{(x+1)^2 - x^2}{x^2(x+1)^2} = \frac{2x+1}{x^2(x+1)^2} \\ \Phi(x) &= \frac{x - (x+1)}{x(x+1)} = \frac{-1}{x(x+1)} \\ \kappa_{rel}(x) &= \frac{|\Phi'(x)| \cdot |x|}{|\Phi(x)|} = \frac{(2x+1)}{x^2(x+1)^2} x^2(x+1) = 1 + \frac{x}{x+1} \leq 2\end{aligned}$$

$\implies$  gut konditioniert!

**Beispiel 4.** Es sei  $\|\cdot\|$  eine Norm auf  $\mathbb{K}^n$  und wir verwenden dieselbe Notation für die induzierte Operatornorm  $\|A\| := \sum_{x \in \mathbb{K}^n, x \neq 0} \frac{\|Ax\|}{\|x\|}$  für  $A \in \mathbb{K}^{n \times n}$ .

Ist  $A$  invertierbar, so bezeichnet  $\text{cond}(A) := \|A\| \cdot \|A^{-1}\|$  die **Konditionszahl von  $A$**  (bzgl.  $\|\cdot\|$ ). Betrachtet man das Lösungsproblem  $\Phi : \mathbb{K}^n \rightarrow \mathbb{K}^n, \Phi(b) = A^{-1}b$ , so gilt für die relative Konditionszahl (mit  $x = A^{-1}b$ )

$$\kappa_{rel}(b) = \frac{\|D\Phi(b)\| \cdot \|b\|}{\|\Phi(b)\|} = \frac{\|A^{-1}\| \cdot \|Ax\|}{\|x\|} \leq \text{cond}(A),$$

wobei die letzte Abschätzung **scharf ist**, d.h. es gilt Gleichheit für mindestens ein  $b$  (und ein  $x$ ).

### 1.3 Verfahrensfehler

Im Wesentlichen gibt es zwei Arten von Verfahrensfehlern

- **Abbruchfehler**, wenn ein konvergenter (aber unendlicher) Algorithmus nach endlich vielen Schritten abgebrochen wird.
- **Diskretisierungsfehler**, wenn eine kontinuierliche Größe durch eine diskrete vereinfacht wird, z.B. Differenzenquotienten statt Differenzialquotient.

**Beispiel 5** (Abbruchfehler Heron-Verfahren). Für  $x > 0$  def.  $y_1 := \frac{1}{2}(1+x), y_{n+1} := \frac{1}{2}(y_n + \frac{x}{y_n})$ .

$$\begin{aligned}\implies y_{n+1}^2 - x &= \frac{1}{4}(y_n^2 + 2x + \frac{x^2}{y_n}) - x = \frac{1}{4}(y_n - \frac{x}{y_n})^2 \geq 0 \\ \implies y_{n+1}^2 &\geq x > 0 \text{ und } y_{n+1} - y_n = \frac{1}{2}(y_n + \frac{x}{y_n}) - y_n = \frac{x}{2y_n} - \frac{y_n}{2} = \frac{x - y_n^2}{2y_n} \leq 0 \\ &\implies 0 < \sqrt{x} \leq y_{n+1} \leq y_n \implies y_n \rightarrow y \\ \implies y &= \frac{1}{2}(y + \frac{x}{y}) \implies \frac{1}{2}y^2 = \frac{1}{2}y^2 + \frac{1}{2}x \implies y^2 = x \implies y = \sqrt{x}\end{aligned}$$

$\implies (y_n)$  konvergiert monoton fallend gegen  $\sqrt{x}$ . Sobald  $y_n = \sqrt{x}$ , würde auch  $y_{n+1} = \sqrt{x}$  gelten, d.h. endkonstante Folge.

**Später:** Heron-Verfahren ist tatsächlich **quadratisch konvergent**, d.h. ex.  $C > 0$  mit  $|\sqrt{x} - y_{n+1}| \leq C|\sqrt{x} - y_n|^2$ .  $\implies$  schnelle konvergenz, weil sich korrekte Ziffern pro Schritt verdoppeln!

**Beispiel 6** (Diskretisierungsfehler einseitiger Differenzenquotient). Sei  $f : \mathbb{R} \rightarrow \mathbb{R}$  differenzierbar,  $x \in \mathbb{R}$

Ziel: Approximiere  $\Phi := f'(x)$  durch den einseitigen Diffquot.  $\Phi_h = \frac{f(x+h)-f(x)}{h}$

klar:  $\Phi_h \rightarrow \Phi$  für  $h \rightarrow 0$ , aber die Konvergenz kann beliebig langsam sein. Man interessiert sich in der Numerik auch für Konvergenzraten bzgl. des Diskretisierungsparameters.

Für  $f \in C^2$  (lokal um  $x$ ) gilt nach Mittelwertsatz

$$f'(x) - \frac{f(x+h)-f(x)}{h} = f'(x) - f'(\zeta) = f''(\xi)(x-\zeta)$$

mit Zwischenstellen  $x \leq \xi \leq \zeta \leq x+h$

$$\implies |\Phi - \Phi_h| \leq \|f''\|_{L^\infty(x, x+h)} h = \mathcal{O}(h)$$

d.h. hier Konvergenzrate 1 in  $h$ .

**Definition 3.** Es sei  $\Phi$  eine kontinuierliche Größe mit Diskretisierung  $\Phi_h$  für  $h > 0$ . Dann bezeichnet man eine Abschätzung der Form  $|\Phi - \Phi_h| = \mathcal{O}(h^\alpha)$  als **a-priori Fehlerabschätzung** mit **Konvergenzrate**  $\alpha > 0$  (auch **Konvergenzordnung**).

Natürlich interessiert sich die Numerik für Verfahren, bei denen  $\alpha > 0$  möglich groß ist.

**Beispiel 7** (zentraler Differenzenquotient).  $f : \mathbb{R} \rightarrow \mathbb{R}$  diffbar,  $x \in \mathbb{R}$ ,  $\Phi := f'(x)$  und  $\Phi_h := \frac{1}{2} \left( \frac{f(x+h)-f(x)}{h} + \frac{f(x)-f(x-h)}{h} \right)$

klar:  $\Phi_h \rightarrow \Phi$  für  $h \rightarrow 0$ ,  $|\Phi - \Phi_h| = \mathcal{O}(h)$  sofern  $f \in C^2$  (lokal um  $x$ ).

Für  $f \in C^3$  (lokal um  $x$ ) gilt mit Taylor  $f(x+h) = f(x) + hf'(x) + \frac{h^2}{2}f''(x) + \frac{h^3}{6}f'''(\zeta)$  für geeignete  $x-h \leq \zeta_- \leq x \leq \zeta_+ \leq x+h$

$$\begin{aligned} \implies \frac{f(x+h)-f(x)}{h} &= f'(x) + \frac{h}{2}f''(x) + \frac{h^2}{6}f'''(\zeta_-) \\ \frac{f(x)-f(x-h)}{h} &= f'(x) - \frac{h}{2}f''(x) + \frac{h^2}{6}f'''(\zeta_+) \\ \implies |\Phi - \Phi_h| &= \frac{h^2}{6} \frac{|f'''(\zeta_+) + f'''(\zeta_-)|}{2} = \mathcal{O}(h^2) \end{aligned}$$

d.h. Konvergenzrate  $\alpha = 2$ .

$\implies$  höhere Genauigkeit für gleiches  $h$  bzw. gleiche Genauigkeit für größeres  $h$ .

**Bemerkung 3.** Auslöschung tritt immer auf (insb. bei Diffquot.), aber sie wird abgemildert durch Verfahren höherer Ordnung. Eine andere Möglichkeit für Verfahren höherer Ordnung zur Approximation von  $f'(x)$  ist die Verwendung von Polynomapproximation, d.h.  $f \approx p$  Polynom und berechne  $p'(x) \approx f'(x)$ .

## 2 Interpolation

Bei einem **Interpolationsproblem** sind im einfachsten Fall Paare  $(x_j, y_j)$  gegeben und eine "einfache" Funktion  $p$  mit  $p(x_j) = y_j \forall j$  gesucht, z.B. Polynome, Splines (= stückweise Polynome), rationale Funktionen (= Quotienten von Polynomen). Verwandt, aber mathematisch schwieriger sind **Approximationsprobleme**. Dabei ist eine Funktion  $f$  und eine Norm  $\|\cdot\|$  gegeben, und es wird eine einfache Funktion  $p$  gesucht, die  $\|f-p\|$  in dieser Klasse einfacher Fkt. minimiert. Oft ist dabei die Funktion  $f$  nur implizit gegeben, d.h. unbekannt.

### 2.1 Lagrange-Polynominterpolation

**Problemstellung:** Gegeben sind  $n+1$  reelle **Stützstellen**  $a \leq x_0 < \dots < x_n \leq b$  und **Funktionswerte**  $y_0, \dots, y_n \in \mathbb{K}$ . Die **Lagrange-Interpolationsaufgabe** sucht ein Polynom  $p \in \mathbb{P}_n = \{p(x) = \sum_{j=0}^n a_j x^j \mid a_0, \dots, a_n \in \mathbb{K}\}$  vom Grad  $n$  mit  $p(x_j) = y_j \forall j = 0, \dots, n$

**Lemma 1.** 1.  $\mathbb{P}_n$  ist  $\mathbb{K}$ -Vektorraum mit  $\dim \mathbb{P}_n = n+1$ .

2. Die **Monome**  $p_j(x) = x^j, j = 0, \dots, n$  sind eine Basis von  $\mathbb{P}_n$ .

3. Die **Lagrange-Polynome**  $L_j(x) = \prod_{k=0, k \neq j}^n \frac{x-x_k}{x_j-x_k} \in \mathbb{P}_n$  erfüllen  $L_j(x_k) = \delta_{jk}$  für alle  $j, k = 0, \dots, n$  und bilden eine Basis von  $\mathbb{P}_n$ .

4. Die **Newton-Polynome**  $q_j(x) = \prod_{k=0}^{j-1} (x-x_k) \in \mathbb{P}_j$  für  $j = 0, \dots, n$  bilden eine Basis von  $\mathbb{P}_n$ .

*Proof.* klar:  $\mathbb{P}_n$  ist  $\mathbb{K}$ -Vektorraum,  $\dim \mathbb{P}_n \leq n+1$ ,  
 zz:  $\{L_0, \dots, L_n\} \subseteq \mathbb{P}_n$  lin. unab.  
 Sei  $\mu_0, \dots, \mu_n \in \mathbb{K}$  mit  $\sum_{j=0}^n \mu_j L_j(x) = 0 \forall x$   
 Für  $x = x_k$  folgt

$$0 = \sum_{j=0}^n \mu_j \underbrace{L_j(x_k)}_{=\delta_{jk}} = \mu_k$$

$\implies$  lin. unab. laut Def.  $\implies \dim \mathbb{P}_n \geq n+1 \implies$  Monome + Lagrange Pol. bilden Basis von  $\mathbb{P}_n$ .  
 zz:  $\{q_0, \dots, q_n\} \subseteq \mathbb{P}_n$  lin. unab.

Seien  $\mu_0, \dots, \mu_n \in \mathbb{K}$  mit  $\sum_{j=0}^n \mu_j \underbrace{q_j(x)}_{=\prod_{k=0}^{j-1} (x-x_k)} = 0$ .

Für  $x = x_0$  folgt  $\mu_0 q_0(x) = 0 \implies \mu_0 = 0$ .

Für  $x = x_1$  folgt  $\mu_1 \underbrace{q_1(x)}_{\neq 0} = 0$ , also  $\mu_1 = 0$ . Induktives Vorgehen zeigt  $\mu_j = 0 \forall j$ . □

**Satz 1** (Eindeutigkeit + Existenz). *Betrachte Lagrange-Interpolation zu Stützstellen  $a \leq x_0 < \dots < x_n \leq b$  und Funktionswerten  $y_0, \dots, y_n \in \mathbb{K}$ . Dann existiert ein eindeutiges  $p \in \mathbb{P}_n$  mit  $p(x_j) = y_j \forall j$ . Dieses wird gegeben durch  $p = \sum_{j=0}^n y_j L_j$ . Ist  $\{q_0, \dots, q_n\} \subseteq \mathbb{P}_n$  eine Basis von  $\mathbb{P}_n$  und  $p = \sum_{j=0}^n \lambda_j q_j$ , so löst  $\lambda = (\lambda_0, \dots, \lambda_n)$  das lineare Gleichungssystem*

$$\underbrace{\begin{pmatrix} q_0(x_0) & \dots & q_n(x_0) \\ \vdots & & \vdots \\ q_0(x_n) & \dots & q_n(x_n) \end{pmatrix}}_{=:A} \lambda = \begin{pmatrix} y_0 \\ \vdots \\ y_n \end{pmatrix}$$

Die Matrix  $A$  ist regulär, d.h.  $\lambda$  ist die eindeutige Lösung.

*Proof.* Da  $L_j(x_k) = \delta_{jk} \forall j, k$  ist offensichtlich, dass  $p = \sum_{j=0}^n \mu_j L_j$  genau dann das Interpolationsproblem löst, wenn  $\mu_j = y_j \forall j$ .  $\implies$  Eindeutigkeit + Existenz

Def. Lösungsoperator  $\mathcal{P} : \mathbb{K}^{n+1} \rightarrow \mathbb{P}_n$  durch  $(\mathcal{P}y)(x_j) = y_j \forall j = 0, \dots, n \forall y \in \mathbb{K}^{n+1}$

$\implies$  wohldef, bijektiv

Def. Auswertungsoperator  $\mathcal{A} : \mathbb{P}_n \rightarrow \mathbb{K}^{n+1}, p \mapsto (p(x_0), \dots, p(x_n))$

$\implies$  wohldef, linear

$\mathcal{P} \circ \mathcal{A} = \text{Identität}, \mathcal{A} \circ \mathcal{P} = \text{Identität},$

$\implies \mathcal{A} = \mathcal{P}^{-1}, \mathcal{P} = \mathcal{A}^{-1}$

$\implies A$  ist die darstellende Matrix  $\mathcal{A}$ .  $\implies A$  ist regulär, da  $\mathcal{A}$  bijektiv, linear. □

**Bemerkung 4.** Die Konditionszahl  $\text{cond}(A)$  der sogenannten **Vandermonde-Matrix**  $A$  hängt stark von der Wahl der Basis ab. Für die Lagrange-Polynome wäre  $A$  die Identität. Für die Monome ist  $\text{cond}(A)$  in der Regel indiskutabel schlecht (hängt an der Wahl der  $x_j$ ). Die Basiswahl beeinflusst auch die Besetzungsstruktur der Matrix.

**Beispiel 8.** Die Newton-Basis führt auf eine untere Dreiecksmatrix

$$\begin{pmatrix} 1 & 0 & \dots & 0 \\ 1 & q_1(x_1) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & q_1(x_n) & \dots & q_n(x_n) \end{pmatrix}$$

d.h. das lineare GLS kann in  $\mathcal{O}(n^2)$  statt  $\mathcal{O}(n^3)$  gelöst werden.

**Lemma 2** (Horner-Schema). Sei  $p(x) = \sum_{j=0}^n \lambda_j \underbrace{q_j(x)}_{\prod_{k=0}^{j-1} (x-x_k)}$ . Für einen Auswertungspunkt  $x \in \mathbb{R}$  betrachte

- $y = \lambda_n$
- for  $k = n-1 : -1 : 0$
- $y = (x - x_k)y + \lambda_k$

• end

$\implies$  Der Algorithmus berechnet in  $3n$  Operationen den Funktionswert  $y = p(x)$ .

**Satz 2** (Interpolationsfehlerdarstellung). Sei  $f \in \mathcal{C}^{n+1}[a, b]$ ,  $0 \leq m \leq n$ ,  $p \in \mathbb{P}_n$  mit  $p(x_j) = f(x_j) \forall j = 0, \dots, n$ , wobei  $a \leq x_0 < \dots < x_n \leq b$

$$\implies f^{(m)}(x) - p^{(m)}(x) = \frac{f^{(n+1)}(\xi)}{(n+1-m)!} \prod_{l=0}^{n-m} (x - \zeta_l),$$

wobei  $\xi = \xi(m, x)$  und  $\zeta_l = \zeta_l(m, x, x_0, \dots, x_n)$  in  $[a, b]$

Für  $m = 0$  gilt  $\zeta_l = x_l \forall l$ .

*Proof.*  $e := f - p \in \mathcal{C}^{n+1}[a, b]$

$\implies e$  hat mindestens  $n+1$  Nullstellen (bei  $x_l$ )  $\implies e'$  hat mindestens  $n$  Nullstellen  $\implies e^{(m)}$  hat mindestens  $n+1-m$  Nullstellen  $a < \zeta_0 < \dots < \zeta_{n-m} < b$

o.B.d.A.  $x \notin \{\zeta_0, \dots, \zeta_{n-m}\}$

Def.  $F(y) := e^{(m)}(x)w(y) - e^{(m)}(y)w(x)$  mit  $w(x) := \prod_{l=0}^{n-m} (y - \zeta_l)$

$\implies F$  hat  $n+2-m$  Nullstellen  $\implies F^{(n+1-m)}$  hat mind. 1 Nullstelle  $\xi$

$$0 = F^{(n+1-m)}(\xi) = \underbrace{e^{(m)}(x)}_{=f^{(m)}(x)-p^{(m)}(x)} \underbrace{w^{(n+1-m)}(\xi)}_{=(n+1-m)!} - \underbrace{e^{(n+1)}(\xi)}_{=f^{(n+1)}(\xi)} \underbrace{w(x)}_{=\prod_{l=0}^{n-m} (x-\zeta_l)}$$

□

**Korollar 1** (Interpolationsfehler-Abschätzung). Seien  $f \in \mathcal{C}^{n+1}[a, b]$  reell- oder komplexwertig,  $a \leq x_0 < \dots < x_n \leq b$ ,  $p \in \mathbb{P}_n$  mit  $p(x_j) = f(x_j) \forall j = 0, \dots, n$ ,  $0 \leq m \leq n$

$$\implies \|f^{(m)} - p^{(m)}\|_{L^\infty(a,b)} \leq C_{\mathbb{K}} \frac{\|f^{(n+1)}\|_{L^\infty(a,b)}}{(n+1-m)!} (b-a)^{n+1-m}$$

mit  $C_{\mathbb{K}} = 1$  für reellwertiges  $f$ ,  $C_{\mathbb{K}} = 2$  für komplexwertiges  $f$ .

*Proof.* klar für  $\mathbb{K} = \mathbb{R}$ .

Für  $\mathbb{K} = \mathbb{C}$ , betrachte  $\operatorname{Re}(f), \operatorname{Im}(f) \in \mathcal{C}^{n+1}[a, b]$ .

□

**Bemerkung 5.** Aus der Fehlerabschätzung und der Konvergenz der Exponentialreihe  $\exp(y) = \sum_{k=0}^{\infty} \frac{y^k}{k!}$  folgt, dass der Interpolationsfehler "schnell" konvergiert, sofern sich die Ableitungen  $\|f^{(k)}\|_{L^\infty(a,b)}$  gut verhalten (z.B.  $\|f^{(k)}\|_{L^\infty(a,b)} \leq M < \infty$ ).

**Bemerkung 6.** Für  $m = 1$ ,  $a = x$  und  $b = x + h$  folgt

$$|f'(x) - p'(x)| \leq \|f' - p'\|_{L^\infty(x, x+h)} \leq C_{\mathbb{K}} \frac{\|f^{(n+1)}\|_{L^\infty(x, x+h)}}{n!} h^n$$

d.h. besser als die Differenzquotienten aus Kapitel 1.

## 2.2 Cebyshev-Knoten

**Definition 4.** Sei  $f : \mathbb{R} \rightarrow \mathbb{R}$  eine Funktion und  $x \in \mathbb{R}$ . Man nennt  $x$  eine  **$n$ -fache Nullstelle von  $f$** , gdw.  $f(x) = 0$  und  $f$  ist lokal um  $x$   $(n-1)$ -mal diffbar mit  $f^{(k)}(x) = 0 \forall k = 1, \dots, n-1$ . Wir schreiben  $n(f, x) \in \mathbb{N}_0$  für die Vielfachheit.

**Lemma 3.** Sei  $p \in \mathbb{P}_n$  mit Nullstellen  $x_1 < \dots < x_k$  und  $N := \sum_{j=1}^k n(f, x_j) > n$ .

$\implies p = 0$ , d.h. ein nicht-triviales Polynom vom Grad  $n$  hat  $\leq n$  Nullstellen, wobei diese mit Vielfachheit gezählt werden.

*Proof.* Induktion nach  $n$ .

Ind.anf.:  $n = 0$ , d.h.  $p$  ist konstant mit mind. einer Nullstelle  $\implies p = 0$  ✓

Ind.hyp: Die Aussage gelte für alle Polynome  $q \in \mathbb{P}_{n-1}$ .

$p \in \mathbb{P}_n$  hat Nullstellen  $x_1 < \dots < x_k$  und  $N = \sum_{j=1}^k n(p, x_j) > n$

$\implies p' \in \mathbb{P}_{n-1}$  hat Nullstellen  $\zeta_1 < \dots < \zeta_{k-1}$  mit  $x_j < \zeta_j < x_{j+1}$  (nach MWS) und bei allen  $x_j$  mit  $n(p, x_j) > 1$ .

Für die Nullstellen von  $p' \in \mathbb{P}_{n-1}$  gilt also

$$\sum_{j=1}^{k-1} \underbrace{n(p', \zeta_j)}_{\geq 1} + \sum_{j=1}^k \max\{n(p, x_j) - 1, 0\} \geq -1 + \underbrace{\sum_{j=1}^k (\max\{n(p, x_j) - 1, 0\} + 1)}_{\substack{\geq n(p, x_j) \\ = N > n}} > n - 1$$

$$\implies p' = 0 \implies p \text{ konstant} \implies p = 0. \quad \square$$

**Bemerkung 7.** Aus der linearen Algebra wissen wir, dass sich jedes Polynom  $p \in \mathbb{P}_n$  mit Nullstelle  $x_0$  in der Form  $p(x) = q(x)(x - x_0)$  schreiben lässt mit  $q \in \mathbb{P}_{n-1}$ , sog. **Polynomdivision**.

Ferner gilt der **Fundamentalsatz der Algebra**: Für jedes  $p \in \mathbb{P}_n$  existieren  $x_1, \dots, x_n \in \mathbb{C}$  und  $\lambda \in \mathbb{C}$  mit  $p(x) = \lambda \prod_{j=1}^n (x - x_j)$ . Offensichtlich ist diese Aussage viel stärker als "mein Lemma".

**Ziel:** Für  $m = 0$  gilt für alle  $x \in [a, b]$

$$|f(x) - p(x)| \leq C_{\mathbb{K}} \frac{\|f^{(n+1)}\|_{L^\infty(a,b)}}{(n+1)!} \prod_{l=0}^n |x - x_l|,$$

wenn  $f$  glatt und  $p \in \mathbb{P}_n$  Lagrange-Interpolationspolynom zu  $x_j$ .

Nun wollen wir die  $x_j$  so wählen, dass  $\max_{x \in [a,b]} \prod_{l=0}^n |x - x_l|$  minimal wird.

**Definition 5.** Für  $n \in \mathbb{N}_0$  definiere die **Cebysev-Polynome (der ersten Art)** durch  $T_n(t) := \cos(n \arccos t)$  auf  $[-1, 1]$ .

**Lemma 4.** 1.  $T_n(\cos(\Phi)) = \cos(n\Phi) \forall 0 \leq \Phi \leq \pi \forall n \in \mathbb{N}_0$

2. Auf  $[-1, 1]$  gilt  $T_0(t) = 1, T_1(t) = t, T_{n+1}(t) = 2tT_n(t) - T_{n-1}(t) \forall n \in \mathbb{N}$

3.  $T_n \in \mathbb{P}_n[-1, 1]$  mit Leitkoeffizient  $2^{n-1}$  für  $n \geq 1$

4.  $\|T_n\|_{L^\infty(-1,1)} = 1$

5.  $T_n$  hat in  $[-1, 1]$  genau  $n+1$  lokale Extrema  $T_n(s_j^{(n)}) = (-1)^j$  mit  $s_j^{(n)} = \cos\left(\frac{j\pi}{n}\right)$  für  $j = 0, \dots, n$

6.  $T_n$  hat in  $[-1, 1]$  genau  $n$  einfache Nullstellen  $T_n(t_j^{(n)}) = 0, t_j^{(n)} = \cos\left(\frac{(2j-1)\pi}{2n}\right)$  für  $j = 1, \dots, n$

*Beweis nur die sog. Drei-Term-Rekursion in (2).* Whl: Additionstheorem des Cosinus:

$$\begin{aligned} \cos(x) + \cos(y) &= 2 \cos\left(\frac{x+y}{2}\right) \cos\left(\frac{x-y}{2}\right) \\ t &= \cos(\Theta), x := (n+1)\Theta, y := (n-1)\Theta \end{aligned}$$

$$\begin{aligned} \implies \frac{x+y}{2} &= n\Theta, \frac{x-y}{2} = \Theta \\ \implies T_{n+1}(t) + T_{n-1}(t) &= \cos(\underbrace{(n+1)\Theta}_x) + \cos(\underbrace{(n-1)\Theta}_y) = 2 \cos\left(\underbrace{\frac{x+y}{2}}_{n\Theta}\right) \cos\left(\underbrace{\frac{x-y}{2}}_{\Theta}\right) \\ &= \underbrace{2 \cos\left(\frac{x+y}{2}\right)}_{T_n(t)} \underbrace{\cos\left(\frac{x-y}{2}\right)}_{=t} \end{aligned}$$

□

**Satz 3** (Optimalität der Chebyshev-Knoten). Betrachte die affine Transformation  $\Psi : [-1, 1] \rightarrow [a, b], \Psi(t) = \frac{1}{2}\{(a+b) + t(b-a)\}$ .

Seien  $t_1^{(n+1)}, \dots, t_{n+1}^{(n+1)}$  die Nullstellen von  $T_{n+1}$ .

$$\implies \min_{x_0, \dots, x_n \in [a,b]} \max_{x \in [a,b]} \prod_{j=0}^n |x - x_j| = \max_{x \in [a,b]} \prod_{j=0}^n |x - \Psi(t_{j+1}^{(n+1)})| = \left(\frac{b-a}{2}\right)^{n+1} \frac{1}{2^n}.$$

Die  $\Psi(t_{j+1}^{(n+1)})$  für  $j = 0, \dots, n$  heißen **Chebyshev-Knoten in  $[a, b]$** .

*Proof.* 1. zz:  $\max_{t \in [-1,1]} \prod_{j=0}^n |t - t_{j+1}^{(n+1)}| = \frac{1}{2^n}$

$$\text{Lemma (iii) + (vi)} \implies T_{n+1}(t) = 2^n \prod_{j=0}^n (t - t_{j+1}^{(n+1)})$$

$$\text{Lemma (iv)} \implies 1 = \|T_{n+1}\|_{L^\infty(-1,1)} = \max_{t \in [-1,1]} 2^n \prod_{j=0}^n |t - t_{j+1}^{(n+1)}|$$

2. zz.  $\frac{1}{2^n} \leq \inf_{t_0, \dots, t_n \in [-1, 1]} \max_{t \in [-1, 1]} \prod_{j=0}^n |t - t_j|$  (dann folgt die Behauptung für  $[a, b] = [-1, 1]$ )

Annahme: Ex.  $t_0, \dots, t_n \in [-1, 1]$  mit  $w(t) := \prod_{j=0}^n (t - t_j)$  erfüllt

$$\|w\|_{L^\infty(-1,1)} = \max_{t \in [-1,1]} \prod_{j=0}^n |t - t_j| \leq \frac{1}{2^n}.$$

Definiere  $p := \underbrace{\frac{1}{2^n} T_{n+1}}_{\in \mathbb{P}_{n+1}} - \underbrace{w}_{\in \mathbb{P}_{n+1}} \in \mathbb{P}_n$ . Ferner  $\frac{1}{2^n} T_{n+1}(s_j^{(n+1)}) = \frac{(-1)^j}{2^n}$  und  $|w(s_j^{n+1})| < \frac{1}{2^n}$ .

$\implies p$  hat  $n+1$  Vorzeichenwechsel  $\implies n+1$  Nullstellen  $\implies p = 0 \implies w = \frac{1}{2^n} T_{n+1} \nmid$

3. klar:  $\Psi$  ist Bijektion von  $[-1, 1]$  auf  $[a, b]$

$$\Psi(t) - \Psi(t_{j+1}^{(n+1)}) = \frac{1}{2} \{(t - t_{j+1}^{(n+1)})(b - a)\}$$

$$\begin{aligned} \implies \max_{x \in [a,b]} \prod_{j=0}^n |x - \Psi(t_{j+1}^{(n+1)})| &= \max_{t \in [-1,1]} \prod_{j=0}^n |\Psi(t) - \Psi(t_{j+1}^{(n+1)})| = \\ &= \left(\frac{b-a}{2}\right)^{n+1} \underbrace{\max_{t \in [-1,1]} \prod_{j=0}^n |t - t_{j+1}^{(n+1)}|}_{=\frac{1}{2^n}} \end{aligned}$$

□

## 2.3 Lebesgue-Konstante

**Satz 4.** Seien  $a \leq x_0 < \dots < x_n \leq b$  Stützstellen mit zugehörigen Lagrange-Polynomen  $L_0, \dots, L_n \in \mathbb{P}_n$ .

Def.  $I_n : \mathcal{C}[a, b] \rightarrow \mathbb{P}_n$ ,  $I_n f := \sum_{j=0}^n f(x_j) L_j$ .

$\implies I_n$  ist eine lineare Projektion auf  $\mathbb{P}_n$  mit Operatornorm

$$\|I_n\| := \sup_{f \in \mathcal{C}[a,b], f \neq 0} \frac{\|I_n f\|_{L^\infty(a,b)}}{\|f\|_{L^\infty(a,b)}} = \max_{x \in [a,b]} \sum_{j=0}^n |L_j(x)| =: \Lambda(x_0, \dots, x_n).$$

Die Zahl  $\Lambda(x_0, \dots, x_n)$  heißt **Lebesgue-Konstante**.

*Proof.*  $I_n$  wohldef., linear ✓

Für  $p \in \mathbb{P}_n$  gilt  $I_n p = p$ , da Polynominterpolation eine eindeutige Lsg. hat.

Für  $f \in \mathcal{C}[a, b]$  mit  $f \neq 0$  gilt

$$\frac{\|I_n f\|_{L^\infty(a,b)}}{\|f\|_{L^\infty(a,b)}} = \max_{x \in [a,b]} \left| \sum_{j=0}^n \frac{f(x_j)}{\|f\|_{L^\infty(a,b)}} L_j(x) \right| \leq \Lambda(x_0, \dots, x_n).$$

Um Gleichheit zu zeigen, wähle  $x \in [a, b]$  mit  $\sum_{j=0}^n |L_j(x)| = \Lambda(x_0, \dots, x_n)$ . Wähle  $f \in \mathcal{C}[a, b]$  als Polygonzug mit  $\|f\|_{L^\infty(a,b)} \leq 1$  und  $f(x_j) = \text{sign}(L_j(x))$ .  $\implies$  Gleichheit bei obiger Abschätzung. □

**Bemerkung 8.** Derselbe Beweis zeigt für  $\tilde{I}_n : \mathbb{K}^{n+1} \rightarrow \mathbb{P}_n$ ,  $\tilde{I}_n(y_0, \dots, y_n) := \sum_{j=0}^n y_j L_j$ , dass  $\|\tilde{I}_n(y_0, \dots, y_n)\|_{L^\infty(a,b)} \leq \Lambda \max_{j=0, \dots, n} |y_j|$  mit Gleichheit für spezielle  $y_j = \text{sign}(L_j(x))$ , wenn  $x \in [a, b]$  mit  $\sum_{j=0}^n |L_j(x)| = \max_{\tilde{x} \in [a,b]} \sum_{j=0}^n |L_j(\tilde{x})|$ .  
 $\implies \tilde{I}$  ist der (lineare) Lösungsoperator der Pol.int.

$$\implies \|\tilde{I}_n(y_0, \dots, y_n) - I_n(\tilde{y}_0, \dots, \tilde{y}_n)\|_{L^\infty(a,b)} \leq \Lambda \max_{j=0, \dots, n} |y_j - \tilde{y}_j|$$

d.h.  $\Lambda$  ist die abstrakte Konditionszahl der Polynominterpolation.

Abschließend einige Bemerkungen zum Bestapproximationsproblem.

**Lemma 5.**  $X$  normierter Raum,  $Y \leq X$  endlich-dim. Teilraum,  $x \in X$

$\implies$  Ex.  $y \in Y$  mit  $\|x - y\|_X = \min_{\tilde{y} \in Y} \|x - \tilde{y}\|_X$



*Proof.* Wähle Folge  $(y_n)_{n \in \mathbb{N}} \subseteq Y$  mit

$$\begin{aligned} \lim_{n \rightarrow \infty} \|x - y_n\|_X &= \inf_{\tilde{y} \in Y} \|x - \tilde{y}\|_X \\ \implies \|y_n\|_X &\leq \underbrace{\|x - y_n\|_X}_{\text{glm. beschränkt wegen Konvergenz}} + \|x\|_X \implies \sup_{n \in \mathbb{N}} \|y_n\|_X \leq M < \infty \end{aligned}$$

Da  $Y$  endl.-dim., gilt der Satz von Bolzano-Weierstraß, d.h.  $(y_n)_{n \in \mathbb{N}}$  hat eine konvergente Teilfolge. O.B.d.A. ex.  $y \in Y$  mit  $\|y - y_n\|_X \rightarrow 0$  für  $n \rightarrow \infty$ .

$$\|x - y\|_X = \lim_{n \rightarrow \infty} \|x - y_n\|_X = \inf_{\tilde{y} \in Y} \|x - \tilde{y}\|_X. \quad \square$$

**Bemerkung 9.** Mit Satz über Lebesgue-Konstante und dem Lemma gilt für alle  $q \in \mathbb{P}_n$

$$\begin{aligned} \underbrace{\|f - I_n f\|_{L^\infty(a,b)}}_{\geq \min_{q \in \mathbb{P}_n} \|f - q\|_{L^\infty(a,b)}} &\leq \|f - q\|_{L^\infty(a,b)} + \underbrace{\|I_n(f - q)\|_{L^\infty(a,b)}}_{\leq \Lambda \|f - q\|_{L^\infty(a,b)}} \\ \implies \min_{q \in \mathbb{P}_n} \|f - q\|_{L^\infty(a,b)} &\leq \|f - I_n f\|_{L^\infty(a,b)} \leq (1 + \Lambda) \min_{q \in \mathbb{P}_n} \|f - q\|_{L^\infty(a,b)} \end{aligned}$$

**Bemerkung 10.** Nach Satz von Weierstraß gilt  $\lim_{n \rightarrow \infty} \min_{p \in \mathbb{P}_n} \|f - p\|_{L^\infty(a,b)} = 0 \forall f \in \mathcal{C}[a,b]$ .

Nach Satz von Faber gilt allerdings, dass es für jede Folge von Stützstellen  $(x_0^{(n)}, \dots, x_n^{(n)})_{n \in \mathbb{N}}$  eine Funktion  $f \in \mathcal{C}[a,b]$  mit der Eigenschaft, dass  $\|f - I_n^{(n)} f\|_{L^\infty(a,b)}$  divergiert!

Insbesondere muss also  $\Lambda_n^{(n)} \rightarrow \infty$  gelten!

**Bemerkung 11.** Für äquidistante Stützstellen divergiert  $\Lambda_n$  exponentiell schnell. Für Chebyshev-Knoten gilt allerdings  $\Lambda_n = \mathcal{O}(\log n)$ .

**Bemerkung 12.** Der **Remez-Algorithmus** berechnet (in unendlich vielen Schritten) ein Polynom  $q \in \mathbb{P}_n$  mit  $\|f - q\|_{L^\infty(a,b)} = \min_{p \in \mathbb{P}_n} \|f - p\|_{L^\infty(a,b)} \forall f \in \mathcal{C}[a,b]$ .

Startwert ist dafür der Chebyshev-Interpoland.

Der **Alternantensatz von Chebyshev** zeigt, dass das Bestapprox.polynom  $q \in \mathbb{P}_n$  bzgl.  $\|\cdot\|_{L^\infty(a,b)}$  in der Tat eindeutig ist.

## 2.4 Auswertung von Interpolationspol.

**Satz 5** (Neville-Verfahren). Seien  $a \leq x_0 < \dots < x_n \leq b$  Stützstellen mit Funktionswerten  $y_j \in \mathbb{K}$  und  $p \in \mathbb{P}_n$  mit  $p(x_j) = y_j \forall j = 0, \dots, n, x \in [a,b]$  Auswertungspunkt.

Für  $j, m \in \mathbb{N}_0$  mit  $j + m \leq n$ , definiere  $p_{j,m} \in \mathbb{P}_m$  als eind. Int.polynom mit  $p(x_k) = y_k \forall k = j, \dots, j + m$

$$\begin{aligned} p(x) &= p_{0,n}(x) \\ p_{j,0}(x) &= y_j \\ p_{j,m}(x) &= \underbrace{\frac{(x - x_j)p_{j+1,m-1}(x) - (x - x_{j+m})p_{j,m-1}}{x_{j+m} - x_j}}_{=: q(x), q \in \mathbb{P}_m} \end{aligned}$$

*Proof.*  $q(x_j) = y_j, q(x_{j+m}) = y_{j+m}, q(x_k) = y_k, k = j + 1, \dots, n - m + 1 \implies q = p_{j,m}$   $\square$

Dieser Satz führt auf das induktive **Neville-Schema**

$$\begin{array}{ccccccc} y_0 & = & p_{0,0}(x) & \rightarrow & p_{0,1}(x) & \rightarrow & \dots & p_{0,n}(x) = p(x) \\ & & & \nearrow & & & & \\ y_1 & = & p_{1,0}(x) & \rightarrow & p_{1,1}(x) & \nearrow & & \\ y_2 & = & p_{2,0}(x) & \nearrow & & & & \\ & \vdots & & & & & & \\ y_{n-1} & = & p_{n-1,0}(x) & \rightarrow & p_{n-1,1}(x) & & & \\ y_n & = & p_{n,0}(x) & \nearrow & & & & \end{array}$$

**Bemerkung 13.** • Das Neville-Verfahren ist ein sog. **Einschritt-Verfahren**, d.h. eine "neue Spalte" nur mit Hilfe der vorausgegangenen Spalte berechnet.

• Wenn man "von oben nach unten rechnet", ist kein zusätzlicher Speicher nötig. In diesem Fall sollte man die "Diagonale" speichern.

- Man kann im Neville-Verfahren dann leicht einen neuen Punkt  $(x_{n+1}, y_{n+1})$  hinzunehmen und erhält  $p_{0,n+1}(x)$ , indem man nur die neue Diagonale rechnet.

**Algorithmus 1** (Neville). *Input:* Stützstellen  $a \leq x_0 < \dots < x_n \leq b$ , Funktionswerte  $y_0, \dots, y_n \in \mathbb{K}$ , Auswertungspunkt  $x \in \mathbb{R}$

- for  $m = 1 : n$
- for  $j = 0 : n - m$
- $y_j = \frac{(x - x_j)y_{j+1} - (x - x_{j+m})y_j}{x_{j+m} - x_j}$
- end
- end

*Output:*  $y_0 = p(x)$ , wobei  $p \in \mathbb{P}_n$  mit  $p(x_j) = y_j \forall j$

klar: Speicherbedarf  $n + 1$  (überschreiben von  $y$ -Vektor), Arithmetischer Aufwand  $\frac{7}{2}n(n + 1)$ .

**Definition 6.** Sei  $p = \sum_{j=0}^n \lambda_j x^j \in \mathbb{P}_n$ . Dann bezeichnet man  $\lambda_n$  als **führenden Koeffizienten von  $p$  bzgl.  $\mathbb{P}_n$** .

Falls  $j = 0$  oder  $(\lambda_j \neq 0 \text{ und } \lambda_k = 0 \forall k > j)$ , so bezeichnet man  $\lambda_j$  als **Leitkoeffizient von  $p$** .

**Satz 6** (Newtons Dividierte Differenzen). Seien  $a \leq x_0 < \dots < x_n \leq b$  Stützstellen,  $y_j \in \mathbb{K}, p \in \mathbb{P}_n$  mit  $p(x_j) = y_j \forall j = 0, \dots, n$ . Für  $j, m \in \mathbb{N}_0$  mit  $j + m \leq n$  definiere

$$y_{j,0} := y_j \qquad y_{j,m} := \frac{y_{j+1,m-1} - y_{j,m-1}}{x_{j+m} - x_j}$$

$\Rightarrow$

1.  $y_{j,m}$  ist der führende Koeff. von  $p_{j,m} \in \mathbb{P}_m$  aus dem Neville-Verfahren.

2. Mit  $\lambda_j := y_{0,j}$  gilt  $p(x) = \sum_{j=0}^n \lambda_j \underbrace{\prod_{k=0}^{j-1} (x - x_k)}_{=q_j \in \mathbb{P}_j}$  d.h. die dividierten Differenzen geben die Koeffizienten

des Int.pol. bzgl. Newton-Basis.

*Proof.*  $q_k := p_{0,k} - p_{0,k-1} \in \mathbb{P}_n$  mit führendem Koeff.  $y_{0,k}$  und Nullstellen  $x_0, \dots, x_{n-1}$

$\Rightarrow q_k = y_{0,k} \prod_{j=0}^{k-1} (x - x_j)$  nach Pol.div.

$$\Rightarrow p = p_{0,k} = p_{0,0} + \sum_{k=1}^n \underbrace{(p_{0,k} - p_{0,k-1})}_{=q_k} = y_{0,0} + \sum_{k=1}^n y_{0,k} \prod_{j=0}^{k-1} (x - x_j) = \sum_{k=0}^n \underbrace{y_{0,k}}_{=\lambda_k} \prod_{j=0}^{k-1} (x - x_j)$$

□

Schema der dividierten Differenzen

$$\begin{array}{ccccccc} y_0 & = & y_{0,0} & \searrow & & & \\ y_1 & = & y_{1,0} & \rightarrow & y_{0,1} & & \\ & & \vdots & & & & \\ & & \vdots & \searrow & & & \\ & & \vdots & \rightarrow & y_{1,1} & \rightarrow & y_{0,2} \\ & & \vdots & & & & \ddots \\ y_{n-1} & = & y_{n-1,0} & \searrow & & & \\ y_n & = & y_{n,0} & \rightarrow & y_{n-1,1} & \rightarrow & \dots \quad y_{0,n} \end{array}$$

$\Rightarrow$  arithmetischer Aufwand  $3 \frac{n(n+1)}{2}$ , um alle  $y_{0,j}$  zu berechnen.

**Bemerkung 14.** • Die dividierten Differenzen sind ein Einschrittverfahren.

- Wenn man den  $y$ -Vektor überschreibt, braucht man keinen zusätzlichen Speicher.
- Das Verfahren löst das Vandermonde-System für die Newton-Basis, ohne die Matrix aufzustellen.

- Die Auswertung von  $p(x)$  erfolgt mit Horner-Schema und Aufwand  $3n$  pro  $x \in \mathbb{K}$ .

**Algorithmus 2.** Input: Stützstellen  $x_0 < \dots < x_n$ , Funktionswerte  $y_0, \dots, y_n \in \mathbb{K}$

- for  $m = 1 : n$
- for  $j = n - m : -1 : 0$
- $y_{j,m} := \frac{y_{j+m} - y_{j,m-1}}{x_{j+m} - x_j}$
- end
- end

Output: Koeffizienten des Interpol. pl.  $p \in \mathbb{P}_n$   $y_0, \dots, y_n$  bzgl. Newton-Basis.

**Bemerkung 15.** Will man das Interpolationspolynom  $p(x)$  an  $N$  Stellen auswerten, so gilt für den Gesamtaufwand: Aufwand(Neville) =  $\frac{7}{2}Nn(n+1)$ , Aufwand(Div. Diff. + Horner) =  $\underbrace{\frac{3}{2}n(n+1)}_{\text{div. Diff.}} + \underbrace{3Nn}_{\text{Horner}}$ .

Es gilt immer: Aufwand(Div. Diff. + Horner)  $\leq$  Aufwand(Neville). Wenn man sich den Fortpflanzungsfehler anschaut dann sieht man aber, dass Neville weniger anfällig ist für Auslöschung.

In der Praxis verwendet man deshalb Neville für kleine  $N$  und Div. Diff. + Horner für große  $N$ .

## 2.5 Hermite-Polynominterpolation

**Satz 7** (Wohlgestelltheit). Gegeben seien Stützstellen  $a \leq x_0 < \dots < x_n \leq b$ , Funktionswerte  $y_j^{(k)} \in \mathbb{K}$  für  $j = 0, \dots, n$  und  $k = 0, \dots, n_j \in \mathbb{N}_0$  (Lagrange  $n_j = 0 \forall j$ ), Def  $N := \left(\sum_{j=0}^n (n_j + 1)\right) - 1$   
 $\implies$  Ex. eind.  $p \in \mathbb{P}_N$  mit  $p^{(k)}(x_j) = y_j^{(k)} \forall j = 0, \dots, n, \forall k = 0, \dots, n_j$ , wobei  $p^{(0)} = p$ .

*Proof.* Betrachte den Auswertungsoperator  $\mathcal{A} : \mathbb{P}_N \rightarrow \mathbb{K}^{N+1}$ ,  $\mathcal{A}p := (p(x_0), \dots, p^{(n_0)}(x_0), p(x_1), \dots, p^{(n_1)}(x_1), \dots, p^{(n_n)}(x_n))$

klar:  $\mathcal{A}$  ist linear und  $\dim \mathbb{P}_N = N + 1$

$\implies \mathcal{A}$  ist  $\underbrace{\text{bijektiv}}_{\text{=Behauptung}}$ , gdw.  $\mathcal{A}$   $\underbrace{\text{injektiv}}_{\text{zu zeigen!}}$  (oder  $\mathcal{A}$  ist surj.)  $\underbrace{\hspace{1cm}}_{\text{"schwierig"}}$

Sei  $p \in \mathbb{P}_N$  mit  $\mathcal{A}p = 0$ , d.h.  $x_j$  eine  $(n_j + 1)$ -fache Nullstelle von  $p \forall j$ .  $\implies p \in \mathbb{P}_N$  hat  $\sum_{j=0}^n (n_j + 1) = N + 1$  viele Nst. (bzgl. Vielfachheit)  $\implies p = 0$ .  $\square$

**Bemerkung 16.** • Der vorausgegangene Beweis ist das "normale Beweisprinzip" für lineare Interpolationsaufgaben. Klar: Man kann die Interpolationsaufgabe insb. lineares Gleichungssystem (äquivalent) formulieren.

- Neville-Verfahren und dividierte Differenzen lassen sich auch für das Hermite-Interpolationsproblem formulieren.
- Analog zu Lagrange (dieselbe Basis) kann man Fehlerdarstellung und Fehlerabschätzung beweisen, z.B.

$$|f(x) - p(x)| \leq C_{\mathbb{K}} \frac{\|f^{(N+1)}\|_{L^\infty(a,b)}}{(N+1)!} \prod_{j=0}^n |x - x_j|^{n_j+1}$$

$$C_{\mathbb{C}} = \sqrt{2} \text{ (vorher 2)}$$

## 2.6 Spline-Interpolation

Die Polynominterpolation erfordert hohe Glätte an  $f$ , um Fehlerabschätzung zu kriegen. Alternativ kann man deshalb stückweise Polynome betrachten (sog. Splines), um Verfahren und Fehlerkontrolle zu haben, falls  $f$  nicht so glatt ist.

**Beispiel 9** (affiner Interpolationspline). Zu Stützstellen  $a = x_0 < x_1 < \dots < x_n = b$  und  $f \in C[a, b]$  ist  $s \in C[a, b]$  mit

- $s|_{[x_{j-1}, x_j]} \in \mathbb{P}_1 \forall j = 1, \dots, n$
- $s(x_j) = f(x_j) \forall j = 0, \dots, n$

$\implies$  Offensichtlich eindeutig  $s(x) = f(x_{j-1}) \frac{x - x_j}{x_{j-1} - x_j} + f(x_j) \frac{x - x_{j-1}}{x_j - x_{j-1}} \forall j \forall x \in [x_{j-1}, x_j]$

**Lemma 6.** Zu  $f \in \mathcal{C}[a, b] \cap \mathcal{C}^2[x_{j-1}, x_j] \forall j = 1, \dots, n$  sei  $s \in \mathcal{C}[a, b]$  der affine Interpolationsspline.

Def  $h : [a, b] \rightarrow \mathbb{R}_{>0}, h|_{[x_{j-1}, x_j]} := x_j - x_{j-1}$  **lokale Netzweite**

$$\implies \|f - s\|_{L^\infty(a, b)} \leq \frac{C_{\mathbb{K}}}{8} \|h^2 f''\|_{L^\infty(a, b)}$$

*Proof.* Sei  $x \in [x_{j-1}, x_j]$

$$\begin{aligned} \implies |f(x) - s(x)| &\leq C_{\mathbb{K}} \frac{\|f''\|_{L^\infty(x_{j-1}, x_j)}}{2} \underbrace{|(x - x_{j-1})(x - x_j)|}_{\text{maximal für } x = \frac{x_{j-1} + x_j}{2}} \\ \implies \|f - s\|_{L^\infty(x_{j-1}, x_j)} &\leq C_{\mathbb{K}} \frac{\|f''\|_{L^\infty(x_{j-1}, x_j)}}{2} \cdot \frac{(x_j - x_{j-1})^2}{4} = \frac{C_{\mathbb{K}}}{8} \|h^2 f''\|_{L^\infty(x_{j-1}, x_j)} \end{aligned}$$

□

**Lemma 7.** Zu  $f \in \mathcal{C}[a, b]$  und  $s \in \mathcal{C}[a, b]$  affiner Int.spline

$$\implies (i) \|f - s\|_{L^2(a, b)} \leq \|hf'\|_{L^2(a, b)}, \text{ sofern } f \in \mathcal{C}^1[x_{j-1}, x_j] \text{ für alle } j$$

$$(ii) \|f - s\|_{L^2(a, b)} \leq \|h^2 f''\|_{L^2(a, b)}, \text{ sofern } f \in \mathcal{C}^2[x_{j-1}, x_j] \text{ für alle } j$$

*Proof.* zz: (i) elementweise für  $I_j = [x_{j-1}, x_j]$

$$\begin{aligned} F &:= f - s \in \mathcal{C}^1[x_{j-1}, x_j], h_j := x_j - x_{j-1} \\ F(x_{j-1}) = 0 &\implies \int_{I_j} |F(x)|^2 = \int_{I_j} \left| \int_{x_{j-1}}^{x_j} F' dx \right|^2 \leq h_j^2 \|F'\|_{L^2(I_j)}^2 \\ &\implies \|F\|_{L^2(I_j)} \leq h_j \|F'\|_{L^2(I_j)} \\ F(x_{j-1}) = 0 = F(x_j) &\implies \int_{I_j} F' dx = 0 \end{aligned}$$

$s'|_{I_j}$  konstant

$$\begin{aligned} \implies s'|_{I_j} &= \frac{1}{h_j} \int_{I_j} f' dx \implies \langle f', s' \rangle_{L^2(I_j)} = \|s'\|_{L^2(I_j)}^2 \\ \implies \|F'\|_{L^2(I_j)}^2 &= \|f'\|_{L^2(I_j)}^2 - 2 \operatorname{Re} \langle f', s' \rangle_{L^2(I_j)} + \|s'\|_{L^2(I_j)}^2 = \|f'\|_{L^2(I_j)}^2 - \|s'\|_{L^2(I_j)}^2 \\ \implies \|f - s\|_{L^2(I_j)} &= \|F\|_{L^2(I_j)} \leq h_j \|F'\|_{L^2(I_j)} \leq h_j \|f'\|_{L^2(I_j)} = \|hf'\|_{L^2(I_j)} \end{aligned}$$

zz: (ii) elementweise

$$\begin{aligned} \operatorname{Re} F(x_{j-1}) = 0 = \operatorname{Re} F(x_j) &\implies \operatorname{Re} F'(\zeta) = 0 \text{ für } x_{j-1} < \zeta < x_j \\ |\operatorname{Re} F'(x)| &= \left| \int_J \operatorname{Re} F'' dt \right| \leq h_j^{\frac{1}{2}} \|\operatorname{Re} F''\|_{L^2(I_j)} \end{aligned}$$

analog für  $\operatorname{Im} F'$

$$\begin{aligned} \implies \int_{I_j} |F'(x)|^2 &\leq h_j^2 \|F''\|_{L^2(I_j)}^2 \implies \|F'\|_{L^2(I_j)} \leq h_j \|F''\|_{L^2(I_j)} \\ \implies \|f - s\|_{L^2(I_j)} &\leq h_j \|F'\|_{L^2(I_j)} \leq h_j^2 \underbrace{\|F''\|_{L^2(I_j)}}_{= f'' \text{ auf } I_j} = \|h^2 f''\|_{L^2(I_j)} \end{aligned}$$

$$(3) \|f - s\|_{L^2(a, b)}^2 = \sum_{j=1}^n \|f - s\|_{L^2(I_j)}^2 \leq \sum_{j=1}^n \|h^2 f''\|_{L^2(I_j)}^2 = \|h^2 f''\|_{L^2(a, b)}^2$$

□

**Definition 7.** Es sei  $\Delta = (x_0, \dots, x_n)$  eine **Zerlegung von**  $[a, b]$ , d.h.  $a = x_0 < \dots < x_n = b$ . Zu gegebenen  $p, q \in \mathbb{N}_0$  heißt  $s : [a, b] \rightarrow \mathbb{K}$  **Spline vom Grad  $p$  mit Glattheit  $q$** , gdw.  $s \in \mathcal{C}^q[a, b]$  mit  $s|_{[x_{j-1}, x_j]} \in \mathbb{P}_p \forall j = 1, \dots, n$ .

Schreibweise  $s \in \mathbb{S}_q^p(\Delta)$  bzw.  $s \in \mathbb{S}^p(\Delta)$ , falls  $q = p - 1$ .

**Bemerkung 17.** Die wichtigsten Beispiele sind **affine Splines**  $\mathbb{S}^1(\Delta)$ , **quadratische Splines**  $\mathbb{S}^2(\Delta)$ , **kubische Splines**  $\mathbb{S}^3(\Delta)$ .

**Bemerkung 18.** Für Wohlgestellttheit von Interpolationsaufgaben muss man  $\dim \mathbb{S}_q^p(\Delta)$  bestimmen.

**Beispiel 10.**  $\dim \mathbb{S}_1^p(\Delta) = \underbrace{n}_{\text{Anz. Intervalle}} \underbrace{(p+1)}_{=\dim \mathbb{P}_p} \hat{=} \text{global unstetige stückweise Polynome}$

$\dim \mathbb{S}_0^p(\Delta) = n(p+1) - (n-1) \hat{=} \text{global stetige stw. Polynome}$

$\implies$  Bestimmung von Basen (und Dimension) von  $\mathbb{S}_q^p(\Delta)$  ist komplizierter als bei Polynomräumen, lässt sich aber über sog. B-Splines bewerkstelligen.

**Bemerkung 19.** Bei Splines hat man mehrere Möglichkeiten, um den Fehler  $\|f - s\|$  zu verringern

- **h-Methode**, d.h. die Zerlegung wird verfeinert
- **p-Methode**, d.h. man erhöht den Polynomgrad
- **r-Methode**, d.h. man erhält  $\#\Delta = n$ , aber man verschiebt die Stützstellen.

Zusätzlich kann man alles lokal kombinieren, z.B.

- **hp-Methode**: Verfeinere  $\Delta$ , wo  $f$  unglatt ist, und erhöhe  $p$ , wo  $f$  glatt ist.

**Satz 8.** Sei  $\Delta = (x_0, \dots, x_n)$  Zerlegung von  $[a, b]$

$\implies \dim \mathbb{S}^p(\Delta) = n + p$  und  $\mathcal{B} := \{x^j, \max\{x - x_k, 0\}^p \mid j = 0, \dots, p, k = 1, \dots, n-1\}$  ist eine Basis.

*Proof.* klar:  $\mathcal{B} \subseteq \mathbb{S}^p(\Delta)$

(1) zz.  $\mathcal{B}$  ist lin. unabh.

Seien  $\lambda_j, \mu_k \in \mathbb{K}$  mit  $0 = \sum_{j=0}^p \lambda_j x^j + \sum_{k=1}^{n-1} \mu_k q_k$

klar: Auf  $[a, b]$  gilt  $q_k|_{[x_0, x_1]} = 0 \forall k = 1, \dots, n-1 \implies \lambda_j = 0 \forall j = 0, \dots, p$ , da Monome lin. unabh. auf  $[x_0, x_1]$ .

klar: Auf  $[x_1, x_2]$  gilt  $q_k|_{[x_1, x_2]} = 0 \forall k = 2, \dots, n-1 \implies \mu_1 q_1 = 0$  auf  $[x_1, x_2]$ ,  $q_1(x_2) \neq 0 \implies \mu_1 = 0$

Induktives Vorgehen zeigt  $\mu_k = 0 \forall k = 1, \dots, n-1$

(2) zz:  $\mathbb{S}^p(\Delta) \subseteq \langle \mathcal{B} \rangle$

Sei  $s \in \mathbb{S}^p(\Delta)$ . Sei  $z_1 \in \mathbb{P}_p$  mit  $s|_{[x_0, x_1]} = z_1|_{[x_0, x_1]}$ .

Für  $k = 2, \dots, n$  definiere sukzessive

$$z_k := z_{k-1} + \frac{s(x_k) - z_{k-1}(x_k)}{q_{k-1}(x_n)} q_{k-1} \in \text{span}(\mathcal{B})$$

zz.  $z_n = s$

Betrachte Residuum  $r := s - z_n$

klar:  $r(x_k) = 0 \forall k = 0, \dots, n$

Auf  $[x_0, x_1]$  gilt  $q_k|_{[x_0, x_1]} = 0 \forall k = 1, \dots, n-1 \implies z_n|_{[x_0, x_1]} = z_1|_{[x_0, x_1]} = s|_{[x_0, x_1]} \implies r|_{[x_0, x_1]} = 0$ .

Auf  $[x_1, x_2]$  gilt  $r \in \mathbb{P}_p$  mit  $r(x_1) = 0 = r(x_2)$ .  $r \in \mathcal{C}^{p-1}[a, b]$ , insb.  $(p-1)$ -mal stetig diffbar in  $x_1$ , d.h.  $x_1$  ist  $p$ -fache Nst. von  $r$ .  $\implies p+1$  viele Nst. auf  $[x_1, x_2] \implies r = 0$  auf  $[x_1, x_2]$ .

Dasselbe Vorgehen auf allen Intervallen zeigt  $r = 0$  auf  $[a, b]$ .  $\square$

**Bemerkung 20.** Sei  $\Delta = (x_0, \dots, x_1)$  Zerlegung von  $[a, b]$  und  $s \in \mathbb{S}^p(\Delta)$  mit  $s(x_j) = y_j \forall j = 0, \dots, n$ . Dann sind nur  $n+1$  Bedingungen fixiert, aber  $\dim \mathbb{S}^p(\Delta) = n+p$ , d.h. es fehlen noch  $p-1$  Bedingungen für Wohlgestelltheit.

Diese werden in der Regel als Randbedingungen formuliert, d.h. an Bedingungen an Ableitungen von  $s$  in  $a$  und  $b$ . Um diese symmetrisch zu stellen, müssen wir annehmen, dass  $p-1 = 2r$  gerade ist.

- (H) **Hermite-Randbedingungen** (oder **vollständige Randbedingungen**): Es werden  $s^{(j)}(a), s^{(j)}(b)$  für  $j = 1, \dots, r$  zusätzlich vorgegeben.
- (N) **Natürliche Randbedingungen**:  $r \leq n$  und  $s^{(j)}(a) = 0 = s^{(j)}(b)$  für  $j = r+1, \dots, 2r$
- (P) **Periodische Randbedingungen**:  $s^{(j)}(a) = s^{(j)}(b)$  für  $j = 1, \dots, 2r$

**Satz 9.** 1. Zu  $p-1 = 2r$  und Randbedingungen (H), (P), (N) existiert (jeweils) ein eindeutiger Interpolationsspline  $s \in \mathbb{S}^p(\Delta)$ , der diese Randbedingung und  $s(x_j) = y_j \forall j = 0, \dots, n$  wobei  $\Delta = (x_0, \dots, x_n)$  und  $y_j \in \mathbb{K}$  gegeben.

2. Erfüllt  $g \in \mathcal{C}^{(r+1)}[a, b]$  dieselben Interpolations- und Randbedingungen wie  $s$ , so gilt

$$\|g^{(r+1)} - s^{(r+1)}\|_{L^2(a,b)}^2 = \|g^{(r+1)}\|_{L^2(a,b)}^2 - \|s^{(r+1)}\|_{L^2(a,b)}^2$$

d.h. der Spline erfüllt die Minimaleigenschaft

$$\|s^{(r+1)}\|_{L^2(a,b)} \leq \|g^{(r+1)}\|_{L^2(a,b)}$$

3. Falls  $s$  ( $N$ ) erfüllt, so muss  $g$  nur die Interpolationsbedingungen erfüllen!

**Bemerkung 21.** Für  $p = 2$  minimieren kubische Interpolationssplines also die Krümmungsenergie  $\|s''\|_{L^2(a,b)}^2$ . Daher kommt auch der Name Spline (dt: Biegestab)

*Proof.* (1) zz. (i) unter der Voraussetzung, dass (ii) gilt.

Seien  $s, \tilde{s} \in \mathbb{S}^p(\Delta)$  mit denselben Interpolations- und Randbedingungen. Aus (ii) folgt dann

$$\|\tilde{s}^{(r+1)} - s^{(r+1)}\|_{L^2(a,b)} = 0 \implies \rho := \tilde{s} - s \in \mathbb{P}_r$$

(H)  $\rho(a) = 0$  und  $\rho^{(j)}(a) = 0$  für  $j = 1, \dots, r \implies a$  ist  $(r+1)$ -fache Nullstelle  $\implies \rho = 0$

(N)  $\rho(x_j) = 0$  für  $j = 0, \dots, n$  und  $r \leq n \implies$  mehr als  $r+1$  Nullstellen  $\implies \rho = 0$

(P)  $\rho^{(j)}(a) = \rho^{(j)}(b)$  für  $j = 1, \dots, 2r$ .  $\rho \in \mathbb{P}_r \implies \rho^{(r-1)} \in \mathbb{P}_1$  und  $\rho^{(r-1)}(a) = \rho^{(r-1)}(b) \implies \rho^{(r-1)} \in \mathbb{P}_0 \implies \rho \in \mathbb{P}_{r-1}$ . Induktiv  $\rho \in \mathbb{P}_0$  mit  $\rho(a) = 0 \implies \rho = 0$ .

Jetzt folgt mit Standardargumentation, dass der lineare Auswertungsoperator  $\mathcal{A} : \mathbb{S}^p(\Delta) \rightarrow \mathbb{K}^{n+p}$  bijektiv ist und aus Dimensionsgründen bijektiv. Insbesondere ist der Lösungsoperator  $\mathcal{L} = \mathcal{A}^{-1}$  bijektiv.

(ii) Betrachte  $|x - y|^2 = |x|^2 - 2\operatorname{Re}x\bar{y} + |y|^2 = |x|^2 - |y|^2 - 2\operatorname{Re}(x - y)\bar{y}$

$$\implies \|g^{(r+1)} - s^{(r+1)}\|_{L^2(a,b)}^2 = \|g^{(r+1)}\|_{L^2(a,b)}^2 - \|s^{(r+1)}\|_{L^2(a,b)}^2 - 2\operatorname{Re} \underbrace{\int_a^b (g^{(r+1)} - s^{(r+1)}) \overline{s^{(r+1)}} dx}_{\text{zz: } = 0}$$

$$\text{Wh: } \int_a^b F'G + FG'dx = [FG]_a^b$$

$$\begin{aligned} \int_a^b (g^{(r+1)} - s^{(r+1)}) \overline{s^{(r+1)}} dx &= \left[ (g^{(r)} - s^{(r)}) \overline{s^{(r+2)}} dx \right]_a^b - \int_a^b (g^{(r)} - s^{(r)}) \overline{s^{(r+2)}} dx = \\ &= \left[ (g^{(r)} - s^{(r)}) \overline{s^{(r+2)}} dx \right]_a^b - \left[ (g^{(r-1)} - s^{(r-1)}) \overline{s^{(r+2)}} dx \right]_a^b + \int_a^b (g^{(r-1)} - s^{(r-1)}) \overline{s^{(r+3)}} dx = \\ &= \left( \sum_{j=0}^{r-1} (-1)^j \left[ (g^{(r-j)} - s^{(r-j)}) \overline{s^{(r+1+j)}} \right]_a^b \right) + (-1)^j \int_a^b (g' - s') \underbrace{\overline{s^{(2r+1)}}}_{= \overline{s^{(p)}} \leftarrow \text{ist konstant auf } [x_{j-1}, x_j]} dx \end{aligned}$$

$$\begin{aligned} \int_a^b (g' - s') \overline{s^{(p)}} dx &= \sum_{j=1}^n \underbrace{\int_{x_{j-1}}^{x_j} (g' - s') \overline{s^{(p)}} dx}_{= \overline{s^{(p)}} [x_{j-1}, x_j] \int_{x_{j-1}}^{x_j} (g' - s') dx} \\ &= \underbrace{\int_{x_{j-1}}^{x_j} (g' - s') dx}_{= [g-s]_{x_{j-1}}^{x_j} = 0} \end{aligned}$$

(H)  $(g^{(r-j)} - s^{(r-j)})(a/b) = 0 \implies \sum(\cdot) = 0$

(N)  $s^{(r+1+j)}(a/b) = 0 \implies \sum(\cdot) = 0$

(P)  $[g^{(r-j)}]_a^b = 0, [s^{(r-j)}]_a^b = 0 \implies \sum(\cdot) = 0$  □

**Bemerkung 22.** Man kann mit Hilfe der sog.  $B$ -Splines Basen von allgemeineren Spline-Räumen konstruieren:

Sei  $(t_j)_{j \in \mathbb{Z}}$  monoton steigend mit  $\lim_{j \rightarrow \pm\infty} t_j = \pm\infty$ .

Def  $B_{j,0} := \Phi_{[t_j, t_{j+1})}$  für  $p = 0$ ,  $B_{j,p} := w_{j,p} B_{j,p-1} + (1 - w_{j,p}) B_{j+1,p-1}$  für  $p \geq 1$  mit  $w_{j,p}(x) = \frac{x - t_j}{t_{j+p} - t_j}$  für  $t_j < t_{j+p}$  und  $w_{j,p}(x) = 0$  für  $t_j = t_{j+p}$

- Man kann zeigen, dass für  $t_j = x_j$  mit  $\Delta = (x_0, \dots, x_n)$  Zerlegung von  $[a, b]$  gilt  $\mathbb{S}^p(\Delta) = \operatorname{span}\{B_{j,p}|_{[a,b]} : j = p, \dots, n-1\}$  und diese  $B_{j,p}$  bilden auch eine Basis (da  $n+p$  viele).
- Ferner haben die  $B_{j,p}$  "gute" Eigenschaften, z.B.  $B_{j,p} > 0$ ,  $\operatorname{supp}(B_{j,p}) = [t_j, t_{j+p+1}]$ ,  $\sum_{j \in \mathbb{Z}} B_{j,p} = 1$ .
- Die Vielfachheit eines Knotens  $t_j = \dots = t_{j+k}$  reduziert die Differenzierbarkeit der Splines bei  $t_k$  um  $k$ , d.h.  $\mathcal{C}^{p-(k+1)}$ -Differenzierbarkeit.

Übung: Sei  $s \in \mathbb{S}^p(\Delta)$  interpolierend  $s(x_j) = y_j$  mit natürlichen Randbedingungen ( $\rightsquigarrow$  "naives" Gleichungssystem mit Basis aus Satz wobei  $(n+3) \times (n+3)$  Gleichungssysteme mit vollbesetzter Vandermonde-Matrix). Auf  $[x_{j-1}, x_j]$  machen wir den Ansatz  $s|_{[x_{j-1}, x_j]} = a_0^{(j)} + a_1^{(j)}(x - x_j) + a_2^{(j)}(x - x_j)^2 + a_3^{(j)}(x - x_j)^3$  mit unbekannten Koeffizienten  $a_0^{(j)}, a_1^{(j)}, a_2^{(j)}, a_3^{(j)} \forall j = 1, \dots, n$ .

$$\implies a_0^{(j)} = y_j, a_1^{(j)} = \frac{y_j - y_{j-1}}{h_j} + \frac{h_j}{3}(2a_2^{(j)} + a_2^{(j-1)}), a_3^{(j)} = \frac{a_2^{(j)} - a_2^{(j-1)}}{3h_j}, h_j := x_j - x_{j-1} \text{ und } a_2^{(j)} = 0$$

TODO Matrix 07 27:32

**Bemerkung 23.** Für interpolierende Splines kann man Fehlerabschätzungen der Form  $\|f - s\|_{L^\infty(a,b)} \leq Ch^{p+1} \|f^{(p+1)}\|_{L^\infty(a,b)}$  wobei  $h := \max_j (x_j - x_{j-1})$ . Aber  $C > 0$  hängt an der Glattheit der Splines. Das ist technisch zu beweisen und hängt (natürlich) auch an den Randbdg.

## 2.7 Diskrete und schnelle Fourier-Transformation

Das Ziel dieses Abschnitts ist die Präsentation und Analyse des FFT-Algorithmus ("fast Fourier transform"), der in der Numerik einer der wichtigsten und vielfältigsten Algorithmen ist. Unsere "Motivation" ist die einfachste Anwendung.

**Definition 8.** Wir bezeichnen  $\mathbb{T}_n := \{\sum_{j=0}^{n-1} \lambda_j \exp(ijx) | \lambda_0, \dots, \lambda_{n-1} \in \mathbb{C}\}$  als Raum der **trigonometrischen Polynome**.

**Satz 10.** 1.  $\dim \mathbb{T}_{n-1} = n$

2. Zu Stützstellen  $0 \leq x_0 < \dots < x_{n-1} \leq 2\pi$  und Funktionswerten  $y_j \in \mathbb{C}$  ex. eind.  $p \in \mathbb{T}_{n-1}$  mit  $p(x_j) = y_j \forall j = 0, \dots, n-1$

*Proof.* klar:  $\dim \mathbb{T}_{n-1} \leq n$

Sei  $p(x) = \sum_{j=0}^{n-1} \lambda_j \exp(ijx) \in \mathbb{T}_n$  mit  $p(x_j) = 0 \forall j = 0, \dots, n-1$

zz:  $\lambda_j = 0 \forall j = 0, \dots, n-1$  (dann lin. unabh. und Eind. des Int.pol.)

Def  $z_k := \exp(ix_k) \in \mathbb{C}, \exp(ijx_k) = z_k^j$

$\implies 0 = p(x_k) = \sum_{j=0}^{n-1} \lambda_j z_k^j =: \tilde{p}(z_k), \tilde{p} \in \mathbb{P}_{n-1} \implies \tilde{p} = 0$ , da  $\tilde{p}$   $n$  Nullstellen hat  $\implies \lambda_j = 0 \forall j = 0, \dots, n-1$ .  $\square$

Im restlichen Abschnitt betrachten wir **äquidistante Stützstellen** (oder **uniforme Stützstellen**)  $x_k = \frac{2\pi k}{n}$  für  $k = 0, \dots, n-1$ . Dies führt auf zusätzliche Struktur der Vandermonde-Matrix, die durch FFT genutzt werden kann.

**Satz 11.** Seien  $x_k = \frac{2\pi k}{n}, y_k \in \mathbb{C}$  gegeben für  $k = 0, \dots, n-1$ . Sei  $p \in \mathbb{T}_n$  das eind. trig. Int.pol. mit  $p(x_j) = y_j \forall k$ .

Sei  $p(x) = \sum_{j=0}^{n-1} \lambda_j \exp(ijx)$  mit Koeff.  $\lambda_j \in \mathbb{C}$ .

Def  $w_n := \exp(-\frac{2\pi i}{n})$   $n$ -te **Einheitswurzel** und  $V_n \in \mathbb{C}^{n \times n}, V_n := (w_n^{jk})_{j,k=0}^{n-1}$  **Fourier-Matrix** (oder: **DFT-Matrix**)

$\implies$

1.  $\frac{1}{n} V_n y = \lambda$  mit  $y = (y_0, \dots, y_{n-1})$ , d.h.  $\lambda_j = \frac{1}{n} \sum_{k=0}^{n-1} w_n^{jk} y_k$

2.  $\frac{1}{\sqrt{n}} V_n$  ist symmetrisch und orthogonal, d.h.  $\left(\frac{1}{\sqrt{n}} V_n\right)^{-1} = \frac{1}{\sqrt{n}} \overline{V_n}$ .

3. Insb. ist  $W = \overline{V_n}$  die Vandermonde-Matrix der trigonometrischen Interpolation

*Proof.* (iii)  $\lambda$  ist Lösung des Vandermonde-Systems  $W\lambda = y$  mit

$$W = \begin{pmatrix} p_0(x_0) & \cdots & p_{n-1}(x_0) \\ \vdots & & \vdots \\ p_0(x_{n-1}) & \cdots & p_{n-1}(x_{n-1}) \end{pmatrix} \text{ mit } p_j = \exp(ijx)$$

$W_{jk} = p_k(x_j) = \exp(ikx_j) = \exp\left(\frac{2\pi i}{n} jk\right) = w_n^{-jk} = (\overline{w_n})^{jk} \implies W = \overline{V_n}$  symmetrisch.

Beweis (i), (ii): Sei  $W^{(k)} = (W_{jk})_{j=0}^{n-1}$   $k$ -te Spalte von  $W$ .

$$W^{(k)} \cdot W^{(k)} = \sum_{j=0}^{n-1} W_{jk} \overline{W_{jk}} = \sum_{j=0}^{n-1} \underbrace{|W_{jk}|^2}_{=|w_n|^{2jk}=1} = n$$

Für  $k+l$  gilt

$$W^{(k)} \cdot W^{(l)} = \sum_{j=0}^{n-1} W_{jk} \overline{W_{jl}} = \sum_{j=0}^{n-1} \underbrace{w_n^{j(l-k)}}_{=(w_n^{(l-k)})^{j=0}} = \frac{1 - (w_n^{(l-k)})^n}{1 - w_n^{(l-k)}} \\ (w_n^{(l-k)})^n = (w_n^n)^{l-k} = 1$$

$\implies \frac{1}{\sqrt{n}} W$  ist symmetrisch und orthogonal und  $W^{-1} = \frac{1}{\sqrt{n}} \left(\frac{1}{\sqrt{n}} W\right)^{-1} = \frac{1}{n} \overline{W} = \frac{1}{n} V_n$ .  $\square$

**Definition 9.** Die Abbildung  $\mathcal{F} : \mathbb{C}^n \rightarrow \mathbb{C}^n, \mathcal{F}_n(x) = V_n x$  bezeichnet man als **diskrete Fourier-Transformation (DFT) der Länge  $n$** . Die Inverse  $\mathcal{F}_n^{-1} : \mathbb{C}^n \rightarrow \mathbb{C}^n, \mathcal{F}_n^{-1}(x) = \frac{1}{n} \bar{V}_n x$  heißt **diskrete Fourier-Rücktransformation**.

**Bemerkung 24.** Die Skalierung von  $\mathcal{F}_n$  ist in der Literatur uneinheitlich. Manchmal  $\frac{1}{n} V_n, \frac{1}{\sqrt{n}} V_n$ , hier  $\mathcal{F}_n(x) = V_n x$ .

**Satz 12.** Für  $p \in \mathbb{N}$  und  $n = 2^p$  und  $m = \frac{n}{2}$  sei  $w_n := \exp(-\frac{2\pi i}{n})$ . Betrachte Permutation  $\sigma_n : \mathbb{C}^n \rightarrow \mathbb{C}^n, \sigma_n(x) = (\underbrace{x_1, x_3, \dots, x_{n-1}}_{\text{ungerade Indizes}}, \underbrace{x_2, x_4, \dots, x_n}_{\text{gerade Indizes}})$ .

Für  $x \in \mathbb{C}^n$  definiere  $a, b \in \mathbb{C}^m$  durch  $a_j := x_j + x_{j+m}, b_j = (x_j - x_{j+m})w_n^{j-1}$  für  $j = 1, \dots, m = \frac{n}{2}$ .

$\Rightarrow \sigma_n(\mathcal{F}_n(x)) = (\mathcal{F}_m(a), \mathcal{F}_m(b))$ , d.h. Auswertung von  $\mathcal{F}_n(x)$  wird auf Fourier-Trafo halber Länge  $\mathcal{F}_m(a), \mathcal{F}_m(b)$  zurückgeführt (+ Vertauschen).

*Proof.*

$$\begin{aligned}\mathcal{F}_n(y_0, \dots, y_{n-1}) &= \left( \sum_{l=0}^{n-1} w_n^{jl} y_l \mid j = 0, \dots, n-1 \right) \\ \mathcal{F}_n(x_1, \dots, x_n) &= \left( \sum_{l=0}^{n-1} w_n^{(j-1)l} x_{l+1} \mid j = 1, \dots, n \right)\end{aligned}$$

(1) zz:  $(\mathcal{F}_n(x))_{2j-1} = (\mathcal{F}_m(a))_j$

klar:  $w_n^2 = w_m, w_m^m = 1$

$$\begin{aligned}\Rightarrow (\mathcal{F}_n(x))_{2j-1} &= \sum_{l=0}^{n-1} w_n^{(2j-2)l} x_{l+1} = \sum_{l=0}^{m-1} \underbrace{(w_n^{2(j-1)l})}_{=w_m^{(j-1)l}} x_{l+1} + \underbrace{w_n^{2(j-1)(l+m)}}_{=w_m^{(j-1)l} \underbrace{w_m^{2(j-1)m}}_{=1}} x_{l+m+1} = \\ &\sum_{l=0}^{m-1} w_m^{(j-1)l} \underbrace{(x_{l+1} + x_{l+m+1})}_{=a_{l+1}} = (\mathcal{F}_m(a))_j\end{aligned}$$

(2) zz:  $(\mathcal{F}_n(x))_{2j} = (\mathcal{F}_m(b))_j$

klar:  $w_n^{(2j-1)l} = w_n^{2(j-1)l} w_n^l = w_m^{(j-1)l} w_n^l, w_n^{(2j-1)(l+m)} = w_n^{2(j-1)l} w_n^l \underbrace{w_n^{2(j-1)m}}_{=(w_n^{n/2})^{2j+1} = -1} = -w_m^{(j-1)l} w_n^l$

$$\Rightarrow (\mathcal{F}_n(x))_{2j} = \sum_{l=0}^{n-1} w_n^{(2j-1)l} x_{l+1} = \sum_{l=0}^{m-1} \underbrace{(w_n^{(2j-1)l} x_{l+1} + w_n^{(2j-1)(l+m)} x_{l+m+1})}_{=w_m^{(j-1)l}} \underbrace{(x_{l+1} - x_{l+m+1}) w_n^l}_{=b_{l+1}} = (\mathcal{F}_m(b))_j$$

□

**Bemerkung 25.** In Matrixform kann man die FFT als Faktorisierung schreiben:

$$V_n = P_n \begin{pmatrix} V_{n/2} & 0 \\ 0 & V_{n/2} \end{pmatrix} \begin{pmatrix} I_{n/2} & I_{n/2} \\ D_{n/2} & -D_{n/2} \end{pmatrix}$$

mit  $P_n = (e_1, e_3, \dots, e_{n-1}, e_2, e_4, \dots, e_n) \in \mathbb{R}^{n \times n}$ ,  $I_{n/2}$  Identität,  $D_{n/2} = \text{diag}(w_n^0, \dots, w_n^{n-1}) \in \mathbb{C}^{n \times n}$

**Korollar 2.** Als **Fast-Fourier-Transformation (FFT)** bezeichnet man die rekursive Berechnung von  $\mathcal{F}_n(x)$  mit  $n = 2^p$  mittels der Rekursion aus dem Satz. Die gesamte Rekursion benötigt weniger als  $\frac{3}{2} n \log_2 n$  arithmetische Operationen plus die Berechnung von  $w_n^l$  für  $l = 0, \dots, n-1$ .

*Proof.* (1) Wegen  $w_{n/2}^l = w_n^{2l}$  reicht es, alle  $w_n^l$  für  $l = 0, \dots, n-1$  und das maximale  $n$  zu berechnen.

(2) Beweis des arithm. Aufwands durch Induktion nach  $p$ .

- $A(p)$  = Anzahl Additionen/Subtraktionen
- $M(p)$  = Anzahl Multiplikationen



Beh:  $A(p) = p^{2^p}, M(p) \leq \frac{1}{2}p^{2^p}$  (dann fertig, da Aufwand( $\mathcal{F}_n(x)$ ) =  $A(p) + M(p) \leq \frac{2}{3}p^{2^p} = \frac{3}{2}n \log_2 n$ )

Ind.anf.  $p = 1, n = 2, \dots, m-1 \implies A(p) = 2, M(p) = 0 \checkmark$

Ind.schritt: Die Aussage gelte für  $p$ , dann  $A(p+1) = 2 \underbrace{A(p)}_{=p^{2^p}} + 2 \cdot 2^p = (p+1)2^{p+1} \checkmark$ .  $M(p+1) = 2 \underbrace{M(p)}_{\leq \frac{1}{2}p^{2^p}} + 2^p \leq$

$$p^{2^p} + 2^p = \frac{1}{2}(p+1)2^{p+1} \checkmark \quad \square$$

**Bemerkung 26.** Die FFT ist eine schnelle Matrix-Vektor-Multiplikation für vollbesetzte Matrizen, die eine gewisse Struktur haben, ohne die Matrix jeweils voll aufzubauen, d.h. Speicherbedarf  $\mathcal{O}(n)$ .

UE:  $A \in \mathbb{K}^{n \times n}$  **zirkulant**, d.h.

$$A = \begin{pmatrix} a_0 & a_{n-1} & \cdot & a_1 \\ a_1 & a_0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ a_{n-1} & a_{n-2} & \cdots & a_0 \end{pmatrix}$$

$\implies V_n A V_n^{-1} = D = \text{diag}(p(1), p(w_n), \dots, p(w_n^{n-1}))$  mit  $p(x) = \sum_{j=0}^n a_j x^j$

$\implies$  klar:  $A$  ist regulär, gdw.  $p(w_n^j) \neq 0 \forall j = 0, \dots, n-1$

Sei  $A$  zusätzlich regulär,  $b \in \mathbb{K}^k$ , Ziel: Löse  $Ax = b$

$\implies A = V_n^{-1} D V_n, A^{-1} = V_n^{-1} D^{-1} V_n \implies x = V_n^{-1} D^{-1} V_n b = \frac{1}{n} \overline{V_n} D^{-1} V_n b, \overline{V_n} \overline{y}$

Frage: Wie berechnet man die Diagonale von  $D$  effizient?

naiv:  $\mathcal{O}(n^2)$ , clever:  $V_n A = D V_n$ , insb.  $V_n a = D(1, \dots, 1)^T = \text{diag}(D) \implies$  FFT liefert  $\text{diag}(D)$ .

**Bemerkung 27.** • Die FFT verdankt ihren Namen dem Zusammenhang mit der Fourier-Transformation:

Für  $f \in L^2[0, 2\pi]$  definiere  $\hat{f}(k) := \frac{1}{2\pi} \int_0^{2\pi} f(t) e^{-ikt} dt$  und  $S_n(t) := \sum_{k=-n}^n \hat{f}(k) e^{ikt}$

$$\implies \lim_n \|f - S_n\|_{L^2(0, 2\pi)} = 0 \text{ und } \frac{1}{2\pi} \|f\|_{L^2(0, 2\pi)}^2 = \sum_{k=-\infty}^{\infty} |\hat{f}(k)|^2$$

- Die trig. Interpolation ist auch für  $p(x) := \sum_{k=-n}^n c_k \exp(ikt)$  mit  $p(x_j) = y_j \forall j = -n, \dots, n$  eind. lösbar, denn  $p(x) = \sum_{l=0}^{2n} c_{l-n} \underbrace{\exp(i(l-n)x)}_{=\exp(ilk) \exp(-inx) \neq 0}$

### 3 Extrapolation

Bei einem Interpolationsproblem versucht man einen Funktionswert  $\Phi(x)$  zu approximieren mit  $x \in [x_0, x_n]$  und  $x_0 < \dots < x_n$  Stützstellen. Bei der Extrapolation ist der einzige Unterschied, dass  $x \notin [x_0, x_n]$ .

#### 3.1 Richardson-Extrapolation

Problemstellung:  $\Phi : [0, 1] \rightarrow \mathbb{K}$  stetig und  $\Phi(h)$  berechenbar für  $h > 0$ , Ziel: Approximiere  $\Phi(0)$

**Algorithmus 3** (Richardson Extrapolation). Input: Die stetige Funktion  $\Phi$  habe eine asymptotische Entwicklung

$$\Phi(h) = \Phi(0) + \sum_{j=1}^n a_j h^{\alpha_j} + \mathcal{O}(h^{\alpha(n+1)})$$

für  $h > 0$  und  $n \in \mathbb{N}$ , wobei  $\alpha > 0$  bekannt, aber  $\Phi(0)$  und  $a_1, \dots, a_n$  unbekannt.

Prozedur: Berechne  $\Phi(h_j)$  für  $h_0 > \dots > h_n$  und werte das eindeutige  $p_n \in \mathbb{P}_n$  mit  $p_n(h_j^\alpha) = \Phi(h_j) \forall j = 0, \dots, n$  mittels Neville-Verfahren bei  $h = 0$  aus, d.h.  $\Phi(0) \approx p_n(0)$ .

**Bemerkung 28.** Meist versucht man kein fixes  $n$ , sondern iteriert, bis  $p_{n+1}(0) \approx p_n(0)$ . Dabei nutzt die Implementierung, dass man das Neville-Verfahren "leicht" um den Knoten  $(h_{n+1}^\alpha, \Phi(h_{n+1}))$  erweitern kann.

**Bemerkung 29.** Aus der asymp. Entwicklung folgt, dass  $|\Phi(0) - \Phi(h)| = \mathcal{O}(h^\alpha)$ , sofern  $a_1 \neq 0$ , und  $\alpha > 0$  ist nicht verbesserbar, d.h. maximale Konvergenzordnung.

**Beispiel 11** (einseitiger Diff.quotient).  $\Phi(h) = \frac{f(x+h) - f(x)}{h}$  für  $h > 0$  und  $\Phi(0) = f'(x)$ .

$$\text{Taylor} \implies f(x+h) = f(x) + \sum_{j=1}^{n+1} \frac{f^{(j)}(x)}{j!} h^j + \mathcal{O}(h^{n+2})$$

$$\begin{aligned} \implies \underbrace{\frac{f(x+h) - f(x)}{h}}_{=\Phi(h)} &= \underbrace{f'(x)}_{=\Phi(0)} + \underbrace{\sum_{j=2}^{n+1} \frac{f^{(j)}(x)}{j!} h^{j-1}}_{=\sum_{k=1}^n \underbrace{\frac{f^{(k+1)}(x)}{(k+1)!}}_{=a_k} h^k} + \mathcal{O}(h^{n+1}) \end{aligned}$$

und  $\alpha = 1$ .

**Beispiel 12** (zentraler Diff.quotient).  $\Phi(h) = \frac{f(x+h)-f(x-h)}{2h}$  für  $h > 0$  und  $\Phi(0) = f'(x)$ .

$$\text{Taylor: } f(x \pm h) = \sum_{j=0}^{2n+2} \frac{f^{(j)}(x)}{j!} (\pm h)^j + \mathcal{O}(h^{2n+3})$$

$$\begin{aligned} \implies \Phi(h) &= \frac{1}{2n} \left( \sum_{j=0}^{2n+2} \frac{f^{(j)}(x)}{j!} \underbrace{(h^j - (-h)^j)}_{=0 \text{ für } j \text{ gerade, } =2h^j \text{ für } j \text{ ungerade}} + \mathcal{O}(h^{2n+3}) \right) = \\ &= \frac{1}{2n} \left( \sum_{l=1}^{n+1} \frac{f^{(2l-1)}(x)}{(2l-1)!} h^{2(l-1)} + \mathcal{O}(h^{2n+3}) \right) = \underbrace{f'(x)}_{=\Phi(0)} + \underbrace{\sum_{l=2}^{n+1} \frac{f^{(2l-1)}(x)}{(2l-1)!} h^{2(l-1)}}_{=\sum_{k=1}^n \frac{f^{(2(k+1)-1)}(x)}{(2(k+1)-1)!} h^{2k} + \mathcal{O}(h^{2(n+1)})} + \mathcal{O}(h^{2n+2}) \end{aligned}$$

$\implies$  asymp. Entwicklung mit  $\alpha = 2$ .

**Beispiel 13** (Romberg-Verfahren). Sei  $f \in \mathcal{C}[a, b]$ ,  $\Phi(0) = \int_a^b f dx$ .  $\Phi(h) = \frac{h}{2} \left( f(a) + 2 \sum_{j=1}^{n-1} f(x_j) + f(b) \right)$  mit  $h = \frac{b-a}{n}$ ,  $x_j := a + jh$ , sog. **summierte Trapezregel**

$\implies \Phi(h) = \Phi(0) + \sum_{j=1}^n a_j h^{2j} + \mathcal{O}(h^{2(n+1)})$  gilt, sog. **Euler-Maclaurin'sche Summenformel**.

**Romberg-Verfahren** = Anwendung von Richardson-Extrapol. auf summierte Trapezregel.

**Satz 13.** Mit  $C > 0, \alpha > 0$  erfülle  $\Phi : [0, 1] \rightarrow \mathbb{K}$  die Entwicklung  $\Phi(h) = \Phi(0) + \sum_{j=1}^n a_j h^{\alpha j} + a_{n+1}(h)$  mit  $|a_{n+1}(h)| \leq C h^{\alpha(n+1)} \forall h > 0$ .

Es sei  $0 < q < 1$  und  $h_k = q^k$  für  $k = 0, \dots, n$ . Sei  $p \in \mathbb{P}_n$  mit  $p(h_j^\alpha) = \Phi(h_j) \forall j = 0, \dots, n$ .

$$\implies |\Phi(0) - p(0)| \leq M q^{\alpha \frac{n(n+1)}{2}} \text{ mit } M = C \underbrace{\frac{1}{1-q^\alpha} \exp\left(\frac{2q^\alpha}{(1-q^\alpha)^2}\right)}_{\text{unabhängig vom Restglied}}$$

**Bemerkung 30.** • Unter der Voraussetzung des Satzes gilt  $|\Phi(0) - \Phi(h_n)| = \mathcal{O}(h_n^\alpha) = \mathcal{O}(q^{\alpha n})$  bei naiver Realisierung;  $|\Phi(0) - p(0)| = \mathcal{O}(q^{\alpha n^2})$  durch Extrapolation.  $\implies$  wesentlich kleineres  $n$  nötig, um dieselbe Genauigkeit zu erhalten. Man sagt "Extrapolation mindert Auslöschung".

• Achtung: Die Richardson-Extrapolation ist nur sinnvoll, wenn man  $\alpha > 0$  kennt!

*Proof.* (1) zz:  $|\Phi(0) - p(0)| \leq C \sum_{l=0}^n q^{\alpha l(n+1)} |L_l(0)|$  mit  $L_l(x) = \prod_{k=0, k \neq l}^n \frac{x - x_l}{x_k - x_l}$ ,  $x_k = h_k^\alpha = q^{\alpha k}$  klar:

$$\begin{aligned} p(x) &= \sum_{l=0}^n \underbrace{\Phi(h_l)}_{\text{asymp. Entwicklung } L_l(x)} = \Phi(0) \underbrace{\sum_{l=0}^n L_l(x)}_{=1} + \sum_{j=1}^n a_j \underbrace{\sum_{l=0}^n \underbrace{(h_l^{\alpha j})}_{=x_l^j}}_{=x^j} + \sum_{l=0}^n a_{n+1}(h_l) L_l(x) = \\ &= \Phi(0) + \sum_{j=1}^n a_j x^j + \sum_{l=0}^n a_{n+1}(h) L_l(x) \end{aligned}$$

$$\implies |p(0) - \Phi(0)| \leq \sum_{l=0}^n \underbrace{|a_{n+1}(h_l)|}_{\leq C h_l^{\alpha(n+1)} = C q^{\alpha l(n+1)}} |L_l(0)|$$

$$(2) \text{ zz: } |L_l(0)| = q^{-l\alpha(n+1)} q^{+\frac{\alpha n(n+1)}{2}} \prod_{k=0, k \neq l}^n \frac{1}{|1 - q^{(k-l)\alpha}|}$$

$$|L_l(0)| = \prod_{k=0, k \neq l}^n \frac{|x_k|}{|x_l - x_k|} = \frac{\prod_{k \neq l} q^{\alpha k}}{\prod_{k \neq l} |q^{\alpha l} - q^{\alpha k}|}$$

$$\prod_{k \neq l} q^{\alpha k} = q^{-lk} \prod_{k=0}^n q^{\alpha k} = q^{-lk} q^{\alpha \sum_{k=0}^n k} = q^{-lk} q^{\alpha \frac{n(n+1)}{2}}$$

$$\prod_{k \neq l} |q^{\alpha l} - q^{\alpha k}| = \prod_{k \neq l} (q^{\alpha l} |1 - q^{\alpha(k-l)}|) = \underbrace{\left( \prod_{k \neq l} q^{\alpha l} \right)}_{= q^{\alpha l n}} \left( \prod_{k \neq l} |1 - q^{\alpha(k-l)}| \right)$$

(3) zz:  $\prod_{k=0, k \neq l}^n \frac{1}{|1 - q^{(k-l)\alpha}|} \leq q^{\alpha \frac{l(l+1)}{2}} \prod_{k=1}^n \frac{1}{(1 - q^{k\alpha})^2}$   
 Betrachte  $\{k - l | k = 0, \dots, n \text{ mit } k \neq l\} = \{-l, \dots, -1\} \cup \{1, \dots, n - l\}$

$$\begin{aligned} \Rightarrow \prod_{k=0, k \neq l}^n |1 - q^{(k-l)\alpha}| &= \underbrace{\left( \prod_{k=1}^{n-l} |1 - q^{k\alpha}| \right)}_{\geq \prod_{k=1}^{n-l} (1 - q^{k\alpha})} \left( \prod_{k=1}^l |1 - q^{-k\alpha}| \right) \geq \\ &= \prod_{k=1}^n (1 - q^{k\alpha}) \prod_{k=1}^l q^{-k\alpha} \prod_{k=1}^l (1 - q^{k\alpha}) = \prod_{k=1}^n (1 - q^{k\alpha})^2 q^{-\alpha \frac{l(l+1)}{2}} \end{aligned}$$

(4) zz:  $\prod_{k=1}^n \frac{1}{(1 - q^{k\alpha})^2} \leq \exp\left(\frac{2q^\alpha}{(1 - q^\alpha)^2}\right)$

$$\log \left( \prod_{k=1}^n \frac{1}{(1 - q^{k\alpha})^2} \right) = 2 \sum_{k=1}^n \log \underbrace{\left( \frac{1}{1 - q^{k\alpha}} \right)}_{\substack{= 1 + \frac{q^{k\alpha}}{1 - q^{k\alpha}} \\ \leq \frac{q^{k\alpha}}{1 - q^{k\alpha}}}} \leq 2 \sum_{k=1}^n \frac{q^{k\alpha}}{1 - q^{k\alpha}} \leq 2 \frac{1}{1 - q^\alpha} \sum_{k=1}^n \underbrace{q^{k\alpha}}_{\leq \frac{q^\alpha}{1 - q^\alpha}} \leq 2 \frac{q^\alpha}{(1 - q^\alpha)^2}$$

(5)

$$\begin{aligned} |p(0) - \Phi(0)| &\stackrel{(1)}{\leq} C \sum_{l=0}^n q^{l\alpha(n+1)} \underbrace{|L_l(0)|}_{\stackrel{(2)}{\leq} q^{-l\alpha(n+1)} q^{\alpha \frac{n(n+1)}{2}} \prod_{k \neq l} \frac{1}{|1 - q^{(k-l)\alpha}|}} \leq \\ &C q^{\alpha \frac{n(n+1)}{2}} \sum_{l=0}^n \underbrace{\prod_{k \neq l} \frac{1}{|1 - q^{(k-l)\alpha}|}}_{\stackrel{(3)}{\leq} q^{\alpha \frac{l(l+1)}{2}} \prod_{k=1}^n \frac{1}{(1 - q^{k\alpha})^2}} \leq \left[ C \exp\left(\frac{2q^\alpha}{(1 - q^\alpha)^2}\right) \underbrace{\sum_{l=0}^n q^{\alpha \frac{l(l+1)}{2}}}_{\leq \frac{1}{1 - q^\alpha}} \right] q^{\alpha \frac{n(n+1)}{2}} \\ &\stackrel{(4)}{\leq} \exp\left(\frac{2q^\alpha}{(1 - q^\alpha)^2}\right) q^{\alpha \frac{n(n+1)}{2}} \end{aligned}$$

□

### 3.2 Aitken'sches $\Delta^2$ -Verfahren

Beim  $\Delta^2$ -Verfahren handelt es sich um ein Verfahren zur **Konvergenzbeschleunigung**, d.h. sei  $(x_n)_{n \in \mathbb{N}} \subset \mathbb{K}$  eine bekannte, konvergente Folge mit unbekanntem Limes  $x = \lim_{n \rightarrow \infty} x_n$

Ziel: Konstruiere eine Folge  $(y_n)_{n \in \mathbb{N}}$  mit  $\lim_{n \rightarrow \infty} \frac{x - y_n}{x - x_n} = 0$ , d.h.  $(y_n)$  konvergiert schneller gegen  $x$ .

**Bemerkung 31.** Sei  $(x_n)_{n \in \mathbb{N}}$  eine geometrisch konvergente Folge mit Limes  $x$ , d.h. ex.  $q \in \mathbb{K}$  mit  $|q| < 1$  und  $x - x_{n+1} = q(x - x_n) \forall n \in \mathbb{N}$

klar:  $\lim_{n \rightarrow \infty} x_n = x$

Betrachte Differenzenoperator  $\Delta y_n := y_{n+1} - y_n$

zz:  $x = x_n - \frac{(\Delta x_n)^2}{\Delta^2 x_n} = x_n - \frac{(x_{n+1} - x_n)^2}{\Delta(x_{n+1} - x_n)} = x_n - \frac{(x_{n+1} - x_n)^2}{x_{n+2} - 2x_{n+1} + x_n}$ , d.h. Limes  $x$  kann aus 3 Folgengliedern  $x_n, x_{n+1}, x_{n+2}$  exakt berechnet werden.

Proof.

$$x_{n+2} - 2x_{n+1} + x_n = (x_{n+2} - x) - 2(x_{n+1} - x) + (x_n - x) = \underbrace{[q^2 - 2q + 1]}_{\neq 0} \underbrace{(x_n - x)}_{\neq 0 \text{ o.B.d.A.}}$$

$$\left(x_n - \frac{(\Delta x_n)^2}{\Delta^2 x_n}\right) - x = (x_n - x) - \frac{(q-1)^2(x_n - x)^2}{(q-1)^2(x_n - x)} = 0$$

□

**Satz 14.** Sei  $(x_n)_{n \in \mathbb{N}} \subset \mathbb{K}$  mit  $x_n \neq x \in \mathbb{K}$  und  $x_{n+1} - x = (q + \delta_n)(x_n - x) \forall n \in \mathbb{N}$  mit  $q \in \mathbb{K}, |q| < 1, (\delta_n)_{n \in \mathbb{N}} \subset \mathbb{K}$  mit  $\lim_n \delta_n = 0$ .

$\implies$

1.  $\lim_n x_n = x$

2. Ex.  $n_0 \in \mathbb{N}$ , sodass  $y_n := x_n - \frac{(\Delta x_n)^2}{\Delta^2 x_n} \in \mathbb{K}$  wohldef. für  $n \geq n_0$

3.  $\lim_n \frac{x - y_n}{x - x_n} = 0$

Proof. (i) Sei  $0 < |q| < \kappa < 1$ . Wegen  $\lim_n \delta_n = 0$ , ex.  $\tilde{n}_0 \in \mathbb{N}$  mit  $|q + \delta_n| \leq \kappa \forall n \geq \tilde{n}_0$

$$\implies |x_{n+1} - x| \leq \underbrace{|q + \delta_n|}_{\leq \kappa} |x_n - x| \implies \lim_n |x_{n+1} - x| = 0 \implies x = \lim_n x_n$$

(ii)

$$x_{n+2} - 2x_{n+1} + x_n = (x_n - x) \underbrace{[(q + \delta_{n+1})(q + \delta_n) - 2(q + \delta_n) + 1]}_{= \underbrace{(q-1)^2}_{\neq 0} + \underbrace{(\delta_n \delta_{n+1} + q(\delta_n + \delta_{n+1} - 2\delta_n))}_{=: \epsilon_n \rightarrow 0}}$$

(iii)

$$y_n - x = (x_n - x) - \frac{(x_{n+1} - x)^2}{(x_n - x)[(q-1)^2 + \epsilon_n]} =$$

$$(x_n - x) - \frac{(q + \delta_n - 1)^2(x_n - x)^2}{(x_n - x)[(q-1)^2 + \epsilon_n]} = (x_n - x) \left(1 - \frac{(q-1 + \delta_n)^2}{(q-1)^2 + \epsilon_n}\right)$$

$$\implies \frac{y_n - x}{x_n - x} = 1 - \frac{(q-1 + \delta_n)^2}{(q-1)^2 + \epsilon_n} \rightarrow 0$$

□

**Bemerkung 32.** Die Voraussetzungen von Aitken sind defakto für jedes numerische Verfahren erfüllt, d.h. bevor man nichts macht und naiv  $(x_n)$  betrachtet, macht man immer Aitken.

**Beispiel 14** (einseitiger Diff.quot.).  $f \in \mathcal{C}^2(\mathbb{R}), z \in \mathbb{R}, h_n > 0$

Taylor  $\implies f(z + h) = f(z) + h_n f'(z) + \frac{h_n^2}{2} f''(\zeta_n)$  mit  $z < \zeta_n < z + h_n$

$$\implies \underbrace{f'(z)}_{=: x} - \underbrace{\frac{f(z + h) - f(z)}{h_n}}_{=: x_n} = h_n \left( -\frac{f''(\zeta_n)}{2} \right) = h_n \left( -\frac{f''(z)}{2} + \underbrace{\frac{f''(z) - f''(\zeta_n)}{2}}_{=: \epsilon_n} \right)$$

$\epsilon_n \rightarrow 0$  für  $h_n \rightarrow 0$

In der Praxis  $h_n := 2^{-n} h_0$ .

Für  $f''(x) \neq 0$

$$\implies x - x_{n+1} = h_{n+1} \left( -\frac{f''(z)}{2} + \epsilon_{n+1} \right) = \underbrace{\frac{1}{2}}_{=: q} \underbrace{\frac{-\frac{f''(z)}{2} + \epsilon_{n+1}}{-\frac{f''(z)}{2} + \epsilon_n}}_{=: 1 + \underbrace{\frac{\epsilon_{n+1} - \epsilon_n}{-\frac{f''(z)}{2} + \epsilon_n}}_{=: \frac{\delta_n}{q}}} (x - x_n)$$

## 4 Numerische Integration

Im ganzen Kapitel seien  $a, b \in \mathbb{R}$  mit  $a < b$ . Ferner sei  $w \in L^1(a, b)$  eine Gewichtsfunction mit  $w(x) > 0$  fast überall.

Ziel: Approximiere  $Qf := \int_a^b f w dx$  für  $f \in \mathcal{C}[a, b]$

**Bemerkung 33.** Soll eine Funktion  $g \in L^1(a, b)$  numerisch integriert werden, so zerlegt man  $g = fw$  mit  $f$  dem glatten Anteil von  $g$  und  $w$  dem singulären Anteil, z.B.  $g(x) = \frac{\sin x}{\sqrt{1-x^2}}$ , dann  $f(x) = \sin(x), w(x) = \frac{1}{\sqrt{1-x^2}}$  analog  $g(x) = x \log x$ , dann  $f(x) = x, w(x) = \log x$ .

### 4.1 Quadraturformeln

**Definition 10.** Gegeben seinen **Stützstellen** (oder: **Quadraturknoten**)  $a \leq x_0 < \dots < x_n \leq b$  und **Gewichte**  $w_0, \dots, w_n \in \mathbb{K}$ . Dann bezeichnet man  $Q_n f = \sum_{j=0}^n w_j f(x_j)$  als **Quadraturformel (der Länge  $n$ )**.  $Q_n$  hat **Exaktheitsgrad**  $m \in \mathbb{N}_0$ , gdw.  $Q_n p = Q p$  für alle  $p \in \mathbb{P}_m$ , d.h. Polynome vom Grad  $m$  werden exakt integriert.

**Lemma 8.** 1.  $Q, Q_n$  sind linear und stetig auf  $\mathcal{C}[a, b]$  mit Operatornorm  $\|Q\| := \sup_{f \in \mathcal{C}[a, b], f \neq 0} \frac{\|Qf\|}{\|f\|_{L^\infty(a, b)}} = \|w\|_{L^1(a, b)}, \|Q_n\| = \sum_{j=0}^n |w_j|$

2. Der Exaktheitsgrad von  $Q_n \leq 2n + 1$

3. Ist  $Q_n$  exakt auf  $\mathbb{P}_{2n+1}$ , so gibt es kein  $p \in \mathbb{P}_{2n+2} \setminus \mathbb{P}_{2n+1}$  mit  $Q_n p = Q p$

4. Ist  $Q_n$  exakt auf  $\mathbb{P}_0$ , so gilt  $\sum_{j=0}^n w_j = \|w\|_{L^1(a, b)}$

5. Für  $w(x) = 1$  und  $Q_n$  exakt auf  $\mathbb{P}_1$ , so gilt  $\sum_{j=0}^n w_j = b - a, \sum_{j=1}^n w_j x_j = \frac{b^2 - a^2}{2}$

**Bemerkung 34.** Oft verwendet man banale Identitäten wie (iv), (v) um zu testen, ob eine Quadratur korrekt implementiert ist.

*Proof.* (i) klar:  $|Qf| \leq \|f\|_{L^\infty(a, b)} \|w\|_{L^1(a, b)}$

$\implies \|Q\| \leq \|w\|_{L^1(a, b)}$

klar:  $|Q_n f| \leq \sum_{j=0}^n |w_j| |f(x_j)| \leq \|f\|_{L^\infty(a, b)} \sum_{j=0}^n |w_j|$

$\implies \|Q_n\| \leq \sum_{j=0}^n |w_j|$

Wähle einen Polygonzug  $f \in \mathcal{C}[a, b]$  mit  $\|f\|_{L^\infty(a, b)} \leq 1$  mit  $w_j f(x_j) = |w_j|$ , d.h.  $f(x_j) = \text{sign} w_j$  also  $|f(x_j)| \leq 1$

$\implies Q_n f = \sum_{j=0}^n w_j f(x_j) \implies \|Q_n\| = \sum_{j=0}^n |w_j|$

(ii) Wähle  $p(x) = \prod_{j=0}^n (x - x_j)^2$ .  $p \in \mathbb{P}_{2n+2}$  und  $p > 0$  für

$\implies Qf = \int_a^b \underbrace{pw}_{>0} dx > 0 = Q_n p$

$\implies p$  wird nicht exakt integriert  $\implies \text{Exaktheit}(Q_n) \leq 2n + 1$

(iii)  $R := Q - Q_n$  linear. Sei  $\{p_0, \dots, p_{2n+1}\} \subseteq \mathbb{P}_{2n+1}$  Basis. Falls  $p_{2n+2} \in \mathbb{P}_{2n+2} \setminus \mathbb{P}_{2n+1}$ , so ist  $\{p_0, \dots, p_{2n+2}\} \subseteq \mathbb{P}_{2n+2}$  Basis.

Es gilt  $R = 0$  auf  $\mathbb{P}_{2n+2}$  gdw.  $R(p_j) = 0 \forall j = 0, \dots, 2n + 2$

(iv)  $\|w\|_{L^1(a, b)} = Q1 = Q_n 1 = \sum_{j=0}^n w_j$

(v) Für  $w = 1$  gilt  $\|w\|_{L^1(a, b)} = b - a, \underbrace{Qx}_{= \int_a^b x dx = \frac{b^2}{2} - \frac{a^2}{2}} = Q_n x = \sum_{j=0}^n w_j x_j$ , da  $Q_n$  exakt auf  $\mathbb{P}_1$  und  $x$  ein

Monom  $\in \mathbb{P}_1$ . □

Achtung: Meistens verwendet Quadratur den Laufindex  $j = 0, \dots, n$  (d.h.  $n + 1$  Stützstellen, Exaktheit  $\leq 2n + 1$ ). Manchmal wird aber  $j = 1, \dots, n$  betrachtet, d.h.  $n$  Stützstellen, Exaktheit  $\leq 2n - 1$ .

**Bemerkung 35.** In der Literatur (z.B. Abramowitz oder Secrest-Strand) sind Quadraturformeln auf Standardintervallen tabelliert, z.B.  $[0, 1], [-1, 1]$ . Um Quadraturformeln auf  $[a, b]$  zu erhalten, verwendet man in der

Regel eine affine Transformation:

$$\begin{aligned}
\Phi : [-1, 1] &\rightarrow [a, b], \Phi(t) = \frac{1}{2}\{a + b + t(b - a)\} \\
\Rightarrow \int_a^b f w dx &= \int_{-1}^1 f(\Phi(x)) \underbrace{w(\Phi(x))}_{=: \tilde{w}(x)} \underbrace{|\det D\Phi(x)|}_{\frac{b-a}{2}} dx = \\
&\frac{b-a}{2} \tilde{Q}(f \circ \Phi) \approx \frac{b-a}{2} \tilde{Q}_n(f \circ \Phi) = \\
&\sum_{j=0}^n \underbrace{\frac{b-a}{2} \tilde{w}_j}_{=: w_j} f(\underbrace{\Phi(\tilde{x}_j)}_{=: x_j}) = \sum_{j=0}^n w_j f(x_j) =: Q_n f
\end{aligned}$$

analog für  $\Phi : [0, 1] \rightarrow [a, b], \Phi(t) = a + t(b - a)$

klar: Falls  $f \in \mathbb{P}_m$ , dann  $f \circ \Phi \in \mathbb{P}_m \Rightarrow \text{Exaktheit}(\tilde{Q}_n) = \text{Exaktheit}(Q_n)$

**Satz 15** (Fehlerabschätzung + Konvergenz). Sei  $Q_n f = \sum_{j=0}^n w_j^{(n)} f(x_j^{(n)})$  eine Quadraturformel der Länge  $n$  und Exaktheit  $m$ .

$$\Rightarrow \|Qf - Q_n f\| \leq (\|w\|_{L^1(a,b)} + \sum_{j=0}^n |w_j^{(n)}|) \min_{p \in \mathbb{P}_m} \|f - p\|_{L^\infty(a,b)}$$

Ferner sind äquivalent:

1.  $Qf = \lim_{n \rightarrow \infty} Q_n f \forall f \in \mathcal{C}[a, b]$
2.  $Qp = \lim_{n \rightarrow \infty} Q_n p \forall p \in \mathbb{P} := \bigcup_{n \in \mathbb{N}} \mathbb{P}_n$  und  $\sup_{n \in \mathbb{N}} \sum_{j=0}^n |w_j^{(n)}| < \infty$ .

*Proof.* Sei  $p \in \mathbb{P}_m$  mit  $\|f - p\|_{L^\infty(a,b)} = \min_{\tilde{p} \in \mathbb{P}_m} \|f - \tilde{p}\|_{L^\infty(a,b)}$

$$\|Qf - Q_n f\| < \|Qf - Qp\| + \|Q_n p - Q_n f\| \leq \|Q\| \|f - p\|_{L^\infty(a,b)} + \|Q_n\| \|f - p\|_{L^\infty(a,b)}$$

zz: (ii  $\Rightarrow$  i)

Sei  $\epsilon > 0$ . Nach Weierstrass ex.  $m \in \mathbb{N}$  und  $p \in \mathbb{P}_m$  mit  $\|f - p\|_{L^\infty(a,b)} \leq \epsilon$

$$\begin{aligned}
\Rightarrow \|Qf - Q_n f\| &\leq \underbrace{\|Qf - Qp\|}_{\leq \epsilon} + \underbrace{\|Qp - Q_n p\|}_{\rightarrow 0 \text{ für } n \rightarrow \infty} + \underbrace{\|Q_n p - Q_n f\|}_{\leq \|Q_n\| \|f - p\|_{L^\infty(a,b)} \leq \epsilon} \leq (\|w\|_{L^1(a,b)} + M)\epsilon + \|Qp - Q_n p\| \\
&\Rightarrow \limsup_n \|Qf - Q_n f\| \leq (\|w\|_{L^1(a,b)} + M)\epsilon \forall \epsilon > 0 \\
&\Rightarrow 0 \leq \liminf_n \|Qf - Q_n f\| \leq \limsup_n \|Qf - Q_n f\| = 0
\end{aligned}$$

(i  $\Rightarrow$  ii) mittels **Satz von Banach-Steinhaus**:

$X, Y$  Banach-Räume,  $T_n \in L(X, Y)$

Dann sind äquivalent:

- $\sup_{n \in \mathbb{N}} \|T_n\|_{L(X,Y)} < \infty$  glm. Beschränktheit
- $\forall x \in X : \sup_{n \in \mathbb{N}} \|T_n x\|_Y < \infty$  pktw. Beschränktheit

jetzt  $X \in \mathcal{C}[a, b], Y = \mathbb{K}, T_n = Q_n$

d.h.  $(Q_n f)_{n \in \mathbb{N}}$  konvergent nach (i), also punktweise Beschränktheit  $\Rightarrow \infty > \sup_n \|Q_n\| = \sup_n \sum_{j=0}^n |w_j^{(n)}|$ .  $\square$

**Bemerkung 36.** 1. Die Implikation (i  $\Rightarrow$  ii) ist nur mathematisch interessant. Der Satz von Banach-Steinhaus ist einer der Fundamentalsätze der Funktionalanalysis.

2. Die Implikation (ii  $\Rightarrow$  i) ist praktisch relevant. Die Eigenschaft  $\lim_n Q_n p = Qp \forall p \in \mathbb{P}$  gilt für alle interpolatorischen Quadraturformeln.

3. Die zentrale Bedingung  $\sup_n \sum_{j=0}^n |w_j^{(n)}| < \infty$  ist kritisch, aber für alle Gauss-Quadraturen erfüllt.  $\Rightarrow$  Konvergenz gilt immer für Gauss-Quadratur.

## 4.2 Interpolatorische Quadraturformeln

**Definition 11.** Zu gegebenen Stützstellen  $a \leq x_0 < \dots < x_n \leq b$  bezeichnet man  $Q_n f := \sum_{j=0}^n \underbrace{\left( \int_a^b L_j w dx \right)}_{=: w_j = Q(L_j)} f(x_j)$

als **interpolatorische Quadraturformel** (oder **Interpolationsquadratur**), wobei  $L_j x := \prod_{k=0}^n \frac{x-x_k}{x_j-x_k}$  die Lagrange-Polynome sind.

**Satz 16.** Für  $Q_n f = \sum_{j=0}^n w_j f(x_j)$  sind äquivalent:

1.  $Q_n$  ist interpolatorisch
2. Für  $f \in C[a, b]$  mit Lagrange-Interpolationspolynom  $p \in \mathbb{P}_n$  (d.h.  $p(x_j) = f(x_j) \forall j = 0, \dots, n$ ) gilt  $Q_n f = Qp$
3.  $\text{Exaktheitsgrad}(Q_n) \geq n$ .

*Proof.* (i)  $\iff$  (ii).

Betrachte  $p = \sum_{j=0}^n f(x_j) L_j \implies Qp = \sum_{j=0}^n f(x_j) \int_a^b L_j w dx$

(i  $\implies$  iii) klar  $Q_n(L_j) = Q(L_j)$ , da  $L_j(x_k) = \delta_{jk}$

$\implies R = Q - Q_n$  ist Null auf Basis  $\{L_0, \dots, L_n\} \implies R = 0$  auf  $\mathbb{P}_n \implies \text{Exaktheit}(Q_n) \geq n$ .

(iii  $\implies$  i)  $\checkmark$

□

**Bemerkung 37.** Die Gewichte einer Interpolationsquadratur kann man durch Lösen eines lin. GLS berechnen. Sei  $\{p_0, \dots, p_n\} \subseteq \mathbb{P}_n$  Basis, sei  $Q_n f := \sum_{j=0}^n w_j f(x_j)$  eine Int. quadratur

$$\implies \underbrace{\begin{pmatrix} p_0(x_0) & \dots & p_0(x_n) \\ \vdots & & \vdots \\ p_n(x_0) & \dots & p_n(x_n) \end{pmatrix}}_{\text{Transponierte der Vandermonde-Matrix aus der Interpolation, also regulär}} \begin{pmatrix} w_0 \\ \vdots \\ w_n \end{pmatrix} = \begin{pmatrix} \int_a^b p_0 w dx \\ \vdots \\ \int_a^b p_n w dx \end{pmatrix}$$

Transponierte der Vandermonde-Matrix aus der Interpolation, also regulär

**Beispiel 15.** • **abgeschlossene Newton-Cotes-Formeln**  $x_j := a + j \frac{b-a}{n}$  für  $j = 0, \dots, n$

- **affine Newton-Cotes-Formeln**  $x_j := a + (j+1) \frac{b-a}{n+2}$  für  $j = 0, \dots, n$
- **Moolavrin-Formeln**  $x_j := a + (j + \frac{1}{2}) \frac{b-a}{n+1}$
- **Clebsch-Curtis-Formeln** Verwende die Extrema oder die Nullstellen des Chebyshev-Polynoms.

**Bemerkung 38.** Ein Satz von Kusunoki besagt, dass für äquidistante Stützstellen immer gilt  $\sum_{j=0}^n |w_j^{(n)}| \rightarrow \infty$ , d.h. es gibt stetige Funktionen  $f \in C[a, b]$  mit  $Q_n f \not\rightarrow Qf$  für  $n \rightarrow \infty$ .

Bei Clebsch-Curtis kann man  $w_j^{(n)} \geq 0$  zeigen, also  $\sum_{j=0}^n |w_j^{(n)}| \stackrel{!}{=} \sum_{j=0}^n w_j^{(n)} = b-a$

$\implies$  Konvergenz  $Q_n f \rightarrow Qf \forall f \in C[a, b]$

**Bemerkung 39.** Einige der abg. Newton-Cotes-Formeln haben auch eigene Namen.

- $n = 1$  Trapezregel,
- $n = 2$  Simpson-Regel,
- $n = 3$  Newton'sche  $\frac{3}{8}$ -Regel,
- $n = 4$  Milne-Regel

**Bemerkung 40.** Man verwendet die NC-Formeln in der Praxis nur für  $n \leq 8$ , da für  $n \geq 9$  negative Gewichte auftreten, was zur Auslöschung führt.

**Bemerkung 41.** Sei  $\tilde{Q}_n$  eine Quadratur auf  $[0, 1]$  zu  $w = 1$ . Sei  $Q_n^j$  die induzierte Quadratur auf  $[a_j, b_j]$  mit  $a = a_0 < b_0 = a_1 < b_1 = \dots < b_N = b$ . Mit der Zerlegung

$$Qf = \int_a^b f dx = \sum_{j=0}^N \underbrace{\int_{a_j}^{b_j} f dx}_{\approx Q_n^j f}$$

erhalte eine sogenannte **summierte Quadraturformel**  $Q_{nN} f := \sum_{j=0}^N Q_n^j f$ .

Falls  $\text{Exaktheitsgrad}(\tilde{Q}_n) \geq 0$  und  $\max_{j=0, \dots, N} (b_j^N - a_j^N) \rightarrow 0$ , so folgt  $Q_{nN} f \rightarrow Qf$  für  $N \rightarrow \infty, \forall f \in C[a, b]$  und sie kriegen sogar a-priori Fehlerabschätzungen!

**Satz 17.** Sie Stützstellen seien symmetrisch, d.h.  $x_j = a + b - x_{n-j} \forall j = 0, \dots, n$ . Das Gewicht sei symmetrisch, d.h.  $w(x) = w(a + b - x) \forall x \in [a, b]$   
 $\implies$

1. Die Gewichte der zugehörigen Interpolationsquad.  $Q_n$  sind symmetrisch, d.h.  $w_j = w_{n-j} \forall j = 0, \dots, n$
2. Falls  $n$  gerade, so gilt Exaktheit( $Q_n$ )  $\geq n + 1$

*Proof.* (i) Betrachte  $\tilde{Q}_n f := \sum_{j=0}^n w_{n-j} f(x_j)$

zz: Exaktheit( $\tilde{Q}_n$ )  $\geq n$  (dann  $\tilde{Q}_n$  interpolatorisch, d.h.  $\{w_j\}$  eindeutig durch  $\{x_j\}$ )

Betrachte  $p_k(x) := (x - \frac{a+b}{2})^k$ ,  $\tilde{p}_k(x) := p_k(a + b - x) = (\frac{a+b}{2} - x)^k = (-1)^k p_k(x)$

$$\tilde{Q}_n p_k = \sum_{j=0}^n w_{n-j} p_k(x_j) = \sum_{l=0}^n w_l \overbrace{p_k(x_{n-l})}^{=\tilde{p}_k(x_l)} = Q_n \tilde{p}_k = Q \tilde{p}_k \forall k = 0, \dots, n$$

$$Q \tilde{p}_k = \int_a^b p_k(a + b - x) \underbrace{w(x)}_{=w(a+b-x)} dx = \int_a^b p_k(y) w(y) dy = Q p_k$$

$\implies \tilde{Q}_n p_k = Q p_k \forall k = 0, \dots, n \implies \tilde{Q}_n$  interpolatorisch.

(ii) zz:  $Q_n p_{n+1} = Q p_{n+1} = 0$ ,  $x_{\frac{1}{2}} = \frac{a+b}{2}$

$$Q_n p_{n+1} = \sum_{j=0}^n w_j p_{n+1}(x_j) = \sum_{j=1}^{\frac{n}{2}} w_{\frac{n}{2}-j} (x_{\frac{n}{2}-j} - \frac{a+b}{2})^{n+1} + \sum_{j=1}^{\frac{n}{2}} \underbrace{w_{\frac{n}{2}+j}}_{=w_{\frac{n}{2}-j}} \left( x_{\frac{n}{2}+j} - \frac{a+b}{2} \right)^{n+1} = 0$$

$\overbrace{n+1}^{\text{ungerade}}$

$$Q p_{n+1} = \int_a^b p_{n+1} w dx = - \int_a^b \underbrace{\tilde{p}_{n+1}(x)}_{p_{n+1}(a+b-x)} \underbrace{w(x)}_{w(a+b-x)} dx = - \int_a^b p_{n+1}(y) w(y) dy = -Q p_{n+1}$$

$\implies Q p_{n+1} = 0$

□

**Korollar 3** (Konkrete Fehlerabschätzungen). Sei  $w(x) = 1$ ,  $C_{\mathbb{K}} = 1$  für  $f$  reellwertig,  $C_{\mathbb{K}} = \sqrt{2}$  für  $f$  komplexwertig.

1. Die Trapezregel  $Q_1 f = \frac{b-a}{2} (f(a) + f(b))$  erfüllt  $|Q f - Q_1 f| \leq C_{\mathbb{K}} \frac{(b-a)^3}{12} \|f''\|_{L^\infty(a,b)} \forall f \in \mathcal{C}^2[a, b]$
2. Die Simpson-Regel  $Q_2 f = \frac{b-a}{6} (f(a) + 2f(\frac{a+b}{2}) + f(b))$  erfüllt  $|Q f - Q_2 f| \leq C_{\mathbb{K}} \frac{(b-a)^5}{2880} \|f^{(4)}\|_{L^\infty(a,b)} \forall f \in \mathcal{C}^4[a, b]$

*Proof.*  $Q_1, Q_2$  sind abg. Newton-Cotens-Formeln, also interpolatorisch.

1. Sei  $p \in \mathbb{P}_1$  mit  $p(a) = f(a), p(b) = f(b)$

$$\implies |Q f - Q_1 f| = |Q f - Q p| \leq \int_a^b \underbrace{|f(x) - p(x)|}_{\leq C_{\mathbb{K}} \frac{\|f''\|_{L^\infty(a,b)}}{2!} |x-a||x-b|} dx \leq C_{\mathbb{K}} \frac{\|f''\|_{L^\infty(a,b)}}{2!} \underbrace{\int_a^b (x-a)(b-x) dx}_{=\frac{(b-a)^3}{6}}$$

2. Sei  $p \in \mathbb{P}_3$  mit  $p(a) = f(a), p(b) = f(b), p(\frac{a+b}{2}) = f(\frac{a+b}{2})$  und  $p'(\frac{a+b}{2}) = f'(\frac{a+b}{2})$

$$\implies |Q f - Q_2 f| = |Q f - \underbrace{Q_2 p}_{=Q p}| = |Q f - Q p| \leq \int_a^b |f(x) - p(x)| dx \leq$$

$$C_{\mathbb{K}} \frac{\|f^{(4)}\|_{L^\infty(a,b)}}{4!} \overbrace{\int_a^b (x-a)(b-x) \left(\frac{a+b}{2} - x\right)^2 dx}^{=\frac{(b-a)^5}{120}}$$

□

**Bemerkung 42.** Auch mit Fehlerabschätzung aus Konvergenzsatz + Abschätzung des Bestapprox.fehlers durch Interpolationsfehler bekommt man konkrete Fehlerabschätzungen, allerdings schlechtere Konstanten  $\rightsquigarrow$  Trapezregel  $\frac{1}{4}$  statt  $\frac{1}{12}$ .



### 4.3 Gauss-Quadratur

Ziel: Konstruiere (eindeutige) Quadraturformel  $Q_n$  mit Exaktheit( $Q_n$ ) =  $2n + 1$

klar:  $Q_n$  muss interpolatorisch sein.

**Bemerkung 43.** Die Analysis in diesem Abschnitt geht auch für ein unbeschränktes Intervall, z.B.  $(0, \infty), (-\infty, \infty)$  sofern  $\int_a^b |x|^n w(x) < \infty \forall n \in \mathbb{N}_0$ .

Betrachte Innenproduktraum  $H := \{f : (a, b) \rightarrow \mathbb{R} \text{ integrierbar} : \|f\| < \infty\}$  mit  $\|f\| = \langle f, f \rangle^{\frac{1}{2}}, \langle f, g \rangle = \int_a^b fg w dx$

klar:  $\langle \cdot, \cdot \rangle$  Skalarprodukt ( $\leadsto L^2(a, b; w dx) = H$ )

**Lemma 9** (Gram-Schmidt-Orthogonalisierung). Sei  $(x^n)_{n \in \mathbb{N}}$  die Folge der Monome in  $H$ . Definiere induktiv  $p_0 := x^0 = 1, p_n := x^n - \sum_{k=0}^{n-1} \frac{\langle x^n, p_k \rangle}{\|p_k\|^2} p_k$

$\implies (p_n)_{n \in \mathbb{N}}$  sind orthogonal bzgl.  $\langle \cdot, \cdot \rangle$  und insb.  $\{p_0, \dots, p_n\} \subseteq \mathbb{P}_n$  Basis und alle  $p_n$  haben Leitkoeffizient

1. Die Polynome  $p_n$  heißen **Orthogonalpolynome**.

**Bemerkung 44.** Für  $(a, b) = (-1, 1), w = 1$  erhält man die **Legendre-Polynome**. Für  $(a, b) = (-1, 1), w(x) = \frac{1}{\sqrt{1-x^2}}$  erhält man **Cebysev-Polynome**. Für  $(a, b) = (0, \infty)$  und  $w(x) = e^{-x}$  erhält man **Lagrange-Polynome**.

**Lemma 10.** Es seien  $x_0, \dots, x_n \in \mathbb{C}$  die gemäß Vielfachheit gezählten Nullstellen des Orth.pol.  $p_{n+1} \in \mathbb{P}_{n+1}$

1. alle Nullstellen sind einfach und liegen in  $(a, b)$

2. Mit den Lagrange-Polynomen  $L_j(x) = \prod_{k=0}^n \frac{x-x_k}{x_j-x_k}$  gilt  $x_j = \frac{\langle x L_j, L_j \rangle}{\|L_j\|^2}$

*Proof.* 1. Seien  $x_0, \dots, x_k \in (a, b)$  alle Nullstellen von  $p_{n+1}$ , die ungeraden Vielfachheit haben und in  $(a, b)$  liegen, bzw.  $k = -1$ , falls keine solchen existieren.

$$q(x) := \prod_{j=0}^k (x - x_j), q \in \mathbb{P}_{k+1}$$

$\implies r := qp_{n+1} \neq 0$  hat nur Nst. gerader Vielfachheit in  $(a, b) \implies r \geq 0$  in  $(a, b)$  oder  $r \leq 0$  in  $(a, b)$

Annahme:  $k < n \implies 0 = \langle q, p_{n+1} \rangle = \int_a^b r w dx \implies r w = 0$  f.ü. Widerspruch zu  $r \neq 0 \neq w$  fast überall. Also  $k = n$ .

2. Polynomdivision  $p_{n+1} = (x - x_j)q$  mit  $q \in \mathbb{P}_n$

$$\implies \underbrace{\langle p_{n+1}, q \rangle}_{=0} = \langle xq, q \rangle - x_j \langle q, q \rangle \implies x_j = \frac{\langle xq, q \rangle}{\langle q, q \rangle} \text{ und } q = \prod_{k=0, k \neq j}^n (x - x_k) = c L_j$$

□

**Satz 18** (Existenz + Eindeutigkeit der Gauss-Quadratur). 1. Ex. eind Quadraturformel  $Q_n f = \sum_{j=0}^n w_j f(x_j)$  mit Exaktheitsgrad( $Q_n$ ) =  $2n + 1$

2. Die Knoten sind die Nullstellen von Orth.pol.  $p_{n+1} \in \mathbb{P}_{n+1}$

3. Die Gewichte erfüllen  $w_j = \int_a^b w L_j dx = \int_a^b w L_j^2 dx > 0$

4.  $|Qf - Q_n f| \leq C_{\mathbb{K}} \frac{\|f^{(2n+2)}\|_{L^\infty(a,b)}}{(2n+2)!} \int_a^b w(x) \prod_{j=0}^n (x - w_j)^2 dx \forall f \in \mathcal{C}^{2n+2}(a, b)$

*Proof.* 1. Existenz: Wähle Nullstellen  $x_0, \dots, x_n$  von  $p_{n+1}$ , definiere  $w_j = \int_a^b w L_j dx$

$\implies$  Exaktheitsgrad( $Q_n$ )  $\geq n$

Sei  $q \in \mathbb{P}_{2n+1}$ . zz:  $Q_n q = Qq$

Polynomdivision  $\implies q = p_{n+1} \alpha + \beta$  mit  $\alpha, \beta \in \mathbb{P}_n$

$$Q_n q = Q_n \beta = Q \beta = \langle \beta, 1 \rangle + \underbrace{\langle p_{n+1}, \alpha \rangle}_{=0} = \underbrace{\langle \beta + \alpha p_{n+1}, 1 \rangle}_{=q} = Qq$$

2. Eindeutigkeit: Sei  $\tilde{Q}_n f = \sum_{j=0}^n \tilde{w}_j f(\tilde{x}_j)$  eine weitere Quadraturformel mit Exaktheit( $\tilde{Q}_n$ ) =  $2n + 1$

zz:  $\tilde{x}_j \in \{x_0, \dots, x_n\} \forall j = 0, \dots, n$

(dann folgt  $\{\tilde{x}_0, \dots, \tilde{x}_n\} = \{x_0, \dots, x_n\}$  und damit  $\tilde{Q}_n = Q_n$ )

Sei  $j \in \{0, \dots, n\}$ .

$$\begin{aligned}
q(x) &:= \left( \prod_{k=0}^n (x - x_k) \right) \left( \prod_{k=0, k \neq j}^n (x - \tilde{x}_k) \right) \in \mathbb{P}_{2n+1} \\
&\implies 0 = Q_n q = Qq = \underbrace{\tilde{Q}_n q}_{\neq 0} = \tilde{w}_j q(\tilde{x}_j) \\
q(\tilde{x}_j) &= \left( \prod_{k=0}^n (\tilde{x}_j - x_k) \right) \underbrace{\left( \prod_{k=0, k \neq j}^n (\tilde{x}_j - \tilde{x}_k) \right)}_{\neq 0} \\
&\implies \tilde{x}_j \in \{x_0, \dots, x_n\}.
\end{aligned}$$

3.  $w_j = \int_a^b w L_j^2 dx > 0$

$$w_j = \int_a^b L_j w dx = Q(L_j) = Q_n(L_j) = \sum_{k=0}^n w_k \underbrace{L_j(x_k)^2}_{=\delta_{jk}} = Q_n(L_j^2) = Q(L_j^2) = \int_a^b L_j^2 w dx$$

4. zz. Fehlerabschätzung

Wähle  $q \in \mathbb{P}_{2n+1}$  mit  $\underbrace{q(x_j) = f(x_j)}_{\text{um interpolatorisch}}, \quad \underbrace{q'(x_j) = f'(x_j)}_{\text{zusätzliche Freiheit für Verbesserung}} \quad \forall j = 0, \dots, n$

$$\begin{aligned}
&\implies |f(x) - q(x)| \leq C_{\mathbb{K}} \frac{\|f^{(2n+2)}\|_{L^\infty(a,b)}}{(2n+2)!} \prod_{j=0}^n (x - x_j)^2 \\
&\implies |Qf - \underbrace{Q_n f}_{=Q_n q = Qq}| = |Q(f - q)| \leq \int_a^b |f(x) - q(x)| w(x) dx
\end{aligned}$$

□

**Lemma 11** (3-Term-Rekursion). Die Orth.pol.  $(p_n)_{n \in \mathbb{N}_0}$  erfüllen

$$\begin{aligned}
p_0(x) &= 1, p_1(x) = x - \beta_0, p_{n+1}(x) = (x - \beta_n)p_n(x) - \gamma_n^2 p_{n-1}(x) \quad \forall n \geq 1 \\
\text{mit reellen Koeff. } \beta_n &= \frac{\langle x p_n, p_n \rangle}{\|p_n\|^2}, \gamma_n = \frac{\|p_n\|}{\|p_{n-1}\|}
\end{aligned}$$

*Proof.* durch Induktion nach  $n$ .

Ind.anf.  $n = 0, 1$ :

$$\text{Erinnerung } p_0(x) = x^0 = 1, p_n(x) = x^k - \sum_{j=0}^{k-1} \frac{\langle x^k, p_j \rangle}{\|p_j\|^2} p_j \stackrel{k=1}{=} x - \frac{\langle x^1, p_0 \rangle}{\|p_0\|^2} p_0 = x - \frac{\langle x, 1 \rangle}{\|1\|^2} 1$$

$$\text{Def } q_{n+1}(x) := (x - \beta_n)p_n(x) - \gamma_n^2 p_{n-1}(x) \in \mathbb{P}_{n+1}, \text{Leitkoeff}(q_{n+1}) = 1 = \text{Leitkoeff}(p_{n+1})$$

$$\implies p_{n+1} - q_{n+1} \in \mathbb{P}_n$$

$$\text{zz: } \langle q_{n+1}, q \rangle = 0 \quad \forall q \in \mathbb{P}_n \quad (\text{dann } \langle p_{n+1} - q_{n+1}, \underbrace{p_{n+1} - q_{n+1}}_{\in \mathbb{P}_n} \rangle = 0 \text{ also } p_{n+1} - q_{n+1} = 0)$$

$$\text{zz: } \langle q_{n+1}, p_j \rangle = 0 \quad \forall j = 0, \dots, n$$

$$\text{Sei } j \in \{0, \dots, n-2\}:$$

$$\begin{aligned}
\langle q_{n+1}, p_j \rangle &= \langle p_n, \underbrace{(x - \beta_n)p_j}_{\in \mathbb{P}_{n-1}} \rangle - \gamma_n^2 \langle p_{n-1}, \underbrace{p_j}_{\in \mathbb{P}_{n-2}} \rangle = 0 \\
&\quad \underbrace{\hspace{10em}}_{=0} \quad \underbrace{\hspace{10em}}_{=0}
\end{aligned}$$

Sei  $j = n-1$ :

$$\begin{aligned}
\langle q_{n+1}, p_{n-1} \rangle &= \langle p_1, x p_{n-1} \rangle - \beta_n \underbrace{\langle p_n, p_{n-1} \rangle}_{=0} - \gamma_n^2 \underbrace{\langle p_{n-1}, p_{n-1} \rangle}_{=\langle p_n, p_n \rangle \text{ Def. } \gamma_n} = \langle p_n, \underbrace{x p_{n-1} - p_n}_{\in \mathbb{P}_{n-1}} \rangle = 0
\end{aligned}$$

Sei  $j = n$ :

$$\begin{aligned}
\langle q_{n+1}, p_n \rangle &= \langle x p_n, p_n \rangle - \underbrace{\beta_n \langle p_n, p_n \rangle}_{=\langle x p_n, p_n \rangle \text{ Def. } \beta_n} - \gamma_n^2 \underbrace{\langle p_{n-1}, p_n \rangle}_{=0} = 0
\end{aligned}$$

□

Übung: Mit den Konstanten  $\gamma_n, \beta_n$  der 3-Term-Rekursion betrachte

$$A = \begin{pmatrix} \beta_0 & -\gamma_1 & 0 & \dots & 0 \\ -\gamma_1 & \beta_1 & \gamma_2 & \dots & 0 \\ \vdots & & & & \vdots \\ 0 & \dots & \dots & -\gamma_n & \beta_n \end{pmatrix} \in \mathbb{R}_{\text{sym}}^{(n+1) \times (n+1)} \quad v^{(k)} = \begin{pmatrix} \tau_0 p_0(x_k) \\ \vdots \\ \tau_n p_n(x_k) \end{pmatrix} \in \mathbb{R}^{n+1}$$

mit  $x_0, \dots, x_n$  Nullstellen von  $p_{n+1}$ ,  $\tau_j = (-1)^j \left( \prod_{k=1}^j \gamma_k \right)^{-1}$

$$\implies Av^{(k)} = x_k v^{(k)}, \text{ d.h. die } x_n \text{ sind genau die } w_k := \int_a^b L_j w dx = \frac{\|w\|_{L^1(a,b)}}{\|v^{(k)}\|_2^2}$$

DH: Um eine Gauss-Quadratur zu berechnen, muss man alle EW und alle EV der Matrix  $A$  bestimmen (d.h. das EW-Problem vollständig lösen).

## 5 Iterative Lösung von GLS

Ziel:

- Wenn man nichtlineare GLS lösen will, so muss regelmäßig eine Folge von linearen GLS lösen (z.B. Newton).
- Man kann lineare GLS lösen, indem man iterativ Matrix-Vektor-Produkte ausrechnet, insb. muss man die Matrix nicht speichern (z.B. FFT, dividierte Diff.)

### 5.1 Fixpunktprobleme

**Definition 12.** Ein **Iterationsverfahren** ist ein Tripel  $(X, \Phi, x^*)$  mit  $X$  metrischer Raum,  $\Phi : X \rightarrow X$ ,  $\Phi(x^*) = x^*$ , d.h.  $x^*$  ist ein **Fixpunkt** von  $\Phi$ . Zu einem **Startwert**  $x_0 \in X$ , sei  $x_{k+1} := \Phi(x_k) \forall k \in \mathbb{N}_0$  die erzeugte **Iteriertenfolge**  $(x_n)_{n \in \mathbb{N}_0}$ .

**Bemerkung 45.** 1. Existiert  $x := \lim_{n \rightarrow \infty} x_n$  und ist  $\Phi$  stetig bei  $x$ , so ist  $x = \Phi(x)$ , dann  $x = \lim_n x_{n+1} = \lim_n \Phi(x_n) = \Phi(x)$ .

2. Ist  $X$  normiert und die Lösung von  $F(x^*) = 0$  gesucht mit  $F : X \rightarrow X$ , so formuliert man dies i.d.R. als Fixpunktproblem, z.B.  $x^* = \Phi(x^*) := x^* \pm F(x^*)$ .

**Satz 19** (Banachscher Fixpunktsatz).  $X$  vollständig metrischer Raum,  $0 < q < 1$  und  $\Phi : X \rightarrow X$  mit  $d(\Phi(x), \Phi(y)) \leq qd(x, y)$

$\implies$

1. Ex. eind.  $x^* \in X$  mit  $\Phi(x^*) = x^*$
2. Für alle  $x_0 \in X$  und  $x_{k+1} := \Phi(x_k) \forall k \in \mathbb{N}_0$  gilt  $\lim_n x_k = x^*$
3. Für alle  $k \in \mathbb{N}_0$  gilt:

- $d(x_k, x^*) \leq qd(x_{k-1}, x^*)$
- $d(x_k, x^*) \leq \frac{q}{1-q} d(x_k, x_{k-1}) \leq \frac{q^k}{1-q} d(x_1, x_0)$
- $d(x_k, x_{k-1}) \leq (1+q)d(x_{k-1}, x^*)$

**Proof.** 1. Eindeutigkeit Fixpunkt: Seien  $x^*, y^* \in X$  mit  $\Phi(x^*) = x^*, \Phi(y^*) = y^*$

$$\implies d(x^*, y^*) = d(\Phi(x^*), \Phi(y^*)) \leq qd(x^*, y^*) \implies d(x^*, y^*) = 0 \implies x^* = y^*$$

2. gezeigt: Falls  $(x_k)_{k \in \mathbb{N}}$  konvergiert, ist  $x^* = \lim_n x_k$  ein Fixpunkt.

3. zz:  $(x_k)_{k \in \mathbb{N}}$  für alle Startwerte  $x_0 \in X$  eine Cauchy-Folge ist.

Für  $m \leq n$  gilt

$$d(x_m, x_n) \leq \sum_{k=m}^{n-1} \underbrace{d(x_k, x_{k+1})}_{=d(\Phi(x_{k-1}), \Phi(x_k)) \leq qd(x_{n-1}, x_k) \leq q^{k-1}d(x_0, x_1)} \leq \left( \sum_{k=m}^{n-1} q^k \right) d(x_0, x_1) \leq q^m \frac{1}{1-q} d(x_0, x_1) \rightarrow 0, m \rightarrow \infty.$$

4. Abschätzungen:

$$\begin{aligned} d(x_k, x^*) &= d(\Phi(x_{k-1}), \Phi(x^*)) \leq q \underbrace{d(x_{k-1}, x^*)}_{\leq d(x_{k-1}, x_k) + d(x_k, x^*)} \\ \implies d(x_k, x^*)(1-q) &\leq q \underbrace{d(x_{k-1}, x_k)}_{\leq q^{k-1}d(x_0, x_1)} \leq q^k d(x_0, x_1) \end{aligned}$$

$$\text{und } d(x_k, x_{k-1}) \leq \underbrace{d(x_k, x^*)}_{\leq qd(x_{k-1}, x^*)} + d(x_{k-1}, x^*) \leq (1+q)d(x_{k-1}, x^*)$$

□

**Definition 13.** Ein Iterationsverfahren  $(X, \Phi, x^*)$  heißt

- **global konvergent**, gdw.  $\forall x_0 \in X : x^* = \lim_{n \rightarrow \infty} x_n$  mit  $(x_n)_{n \in \mathbb{N}_0}$  der Iteriertenfolge  $x_{n+1} := \Phi(x_n) \forall n$
- **lokal konvergent**, gdw.  $\exists \epsilon > 0 \forall x_0 \in \underbrace{U_\epsilon(x^*)}_{:= \{y \in X \mid d(x, y) < \epsilon\}} : x^* = \lim_n x_n$
- **linear konvergent** (oder: mit Konvergenzordnung  $p = 1$ ), gdw.  $\exists q \in (0, 1) \exists \epsilon > 0 \forall x_0 \in U_\epsilon(x^*) \forall n \in \mathbb{N}_0 : d(x^*, x_{n+1}) \leq qd(x^*, x_n)$
- **von Konvergenzordnung**  $p > 1$ , gdw.  $\exists C > 0 \forall \epsilon > 0 \forall x_0 \in U_\epsilon(x^*) \forall n \in \mathbb{N}_0 : d(x^*, x_{n+1}) \leq Cd(x^*, x_n)^p$

Die Menge  $U_\epsilon(x^*)$  nennt man auch **Konvergenzbereich**.

**Beispiel 16.** Ist  $\Phi : X \rightarrow X$  eine (strikte) Kontraktion auf einem vollständig metrischen Raum mit Fixpunkt  $x^* \in X$ , so ist  $(X, \Phi, x^*)$  global linear konv.

**Lemma 12.** Sei  $(X, \Phi, x^*)$  ein Iterationsverfahren mit Konvergenzordnung  $p \geq 1$ . Dann ist  $(X, \Phi, x^*)$  lokal konvergent und in jeder Konvergenzordnung  $1 \leq \tilde{p} \leq p$ .

*Proof.* 1.  $p = 1 \implies$  lokale konvergenz

Wähle  $0 < q < 1$  und  $\epsilon > 0$  gemäß Definition. Sei  $x_0 \in U_\epsilon(x^*)$ . Dann  $d(x^*, x_n) \leq q^n \underbrace{d(x^*, x_0)}_{\in \mathbb{R}}$

2. Konvergenzordnung  $p > 1 \implies$  lineare konvergenz mit  $q = \frac{1}{2}$ . Seien  $C > 0, \epsilon > 0$  gemäß Def. gewählt.

Wähle  $\delta := \min\{\epsilon, (\frac{1}{2C})^{1/(p-1)}\}$ . Sei  $x_0 \in U_\delta(x^*)$

Beh.  $d(x_n, x^*) \leq \underbrace{2^{-n}}_{\leq 1} \underbrace{d(x_0, x^*)}_{< \delta} < \delta \forall n \in \mathbb{N}_0$  (und  $d(x_n, x^*) \leq C \underbrace{d(x_{n-1}, x^*)^{p-1}}_{\leq \delta^{p-1} \leq 1/(2C)} d(x_{n-1}, x^*) \leq \frac{1}{2} d(x_{n-1}, x^*) \forall n \in \mathbb{N}$ )

□

Beweis der Beh. durch Induktion, klar  $n = 0$

$$d(x_{n+1}, x^*) \leq Cd(x_n, x^*)^p \stackrel{\text{IV}}{\leq} C2^{-np} \underbrace{d(x_0, x^*)^p}_{\leq \delta^{p-1}d(x_0, x^*) \leq \frac{1}{2C}d(x_0, x^*)} \leq \frac{2^{-np}}{2} d(x_0, x^*) = \underbrace{2^{-(np+1)}}_{\leq 2^{-(n+1)}} d(x_0, x^*)$$

3. Konvergenzordnung  $p > 1 \implies$  Konvergenzord.  $1 < \tilde{p} < p$

$$d(x^*, x_n) \leq C \underbrace{d(x^*, x_n)^{p-\tilde{p}}}_{< \delta^{p-\tilde{p}}} d(x^*, x_n)^{\tilde{p}} \underbrace{\hspace{1cm}}_{C(p, \tilde{p}, \delta)}$$

□

**Satz 20.** Sei  $(\mathbb{R}, \Phi, x^*)$  ein Iterationsverfahren und  $\Phi$  lokal  $m$ -mal stetig differenzierbar um Fixpunkt  $x^*$ .

$\implies$

1. Falls  $m = 1$  und  $|\Phi'(x^*)| < 1$ , so ist  $(\mathbb{R}, \Phi, x^*)$  linear konvergent
2. Falls  $\Phi^{(k)}(x^*) = 0 \forall k = 0, \dots, m-1$ , so hat  $(\mathbb{R}, \Phi, x^*)$  Konvergenzordnung  $m$
3. Gilt (i) oder (ii) und  $\Phi^{(m)}(x^*) \neq 0$ , so hat  $(\mathbb{R}, \Phi, x^*)$  nicht Ordnung  $m+1$ .

4. Gilt  $|\Phi'(x^*)| > 1$ , so ist die Iteriertenfolge i.a. nicht konvergent, denn

$$\exists C > 1 \exists \epsilon > 0 \forall x \in U_\epsilon(x^*) : |x * - \Phi(x)| \geq C|x * - x|$$

*Proof.* (i) + (ii): Taylor  $\implies$

$$\begin{aligned} \Phi(x) &= \sum_{k=0}^m \frac{\Phi^{(k)}(x^*)}{k!} (x - x^*)^k + o(|x - x^*|^m) = x^* + \frac{\Phi^{(m)}(x^*)}{m!} (x - x^*)^m + o(|x - x^*|^m) \\ &\implies \lim_{x \rightarrow x^*} \frac{\Phi(x) - x^*}{(x - x^*)^m} = \frac{\Phi^{(m)}(x^*)}{m!} \\ &\implies \forall \epsilon > 0 \exists \delta > 0 \forall x \in U_\delta(x^*) : \left| \frac{\Phi(x) - x^*}{(x - x^*)^m} - \frac{\Phi^{(m)}(x^*)}{m!} \right| \leq \epsilon \\ &\implies \left| \frac{\Phi(x) - x^*}{(x - x^*)^m} \right| \leq \left( \frac{|\Phi^{(m)}(x^*)|}{m!} + \epsilon \right) =: C(m, \epsilon) \\ &\implies \exists \delta > 0 \forall x \in U_\delta(x^*) : \underbrace{|\Phi(x) - x^*|}_{x_{n+1}} \leq C(m, \epsilon) \underbrace{|x - x^*|}_{x_n}^m \end{aligned}$$

$\implies$  (ii) für  $m > 1$ . Für  $m = 1$  wähle  $\epsilon > 0$  mit  $C(m, \epsilon) < 1$ .

(iii) Analog

$$\begin{aligned} \left| \frac{\Phi(x) - x^*}{(x - x^*)^m} \right| &\geq \left| \frac{\Phi^{(m)}(x^*)}{m!} \right| - \epsilon =: \tilde{C}(m, \epsilon) \\ &\implies |x - x^*| \tilde{C}(m, \epsilon) \leq |\Phi(x) - x^*| \end{aligned}$$

Wählt man  $\epsilon > 0$  mit  $\tilde{C}(m, \epsilon) > 0$ , so kann das Verfahren nicht Ordnung  $m + 1$  haben.

(iv) Analog  $|\Phi(x) - x^*| \geq \tilde{C}(1, C)|x - x^*|$  mit  $\tilde{C}(1, \epsilon) > 1$  für  $\epsilon > 0$  klein genug.  $\square$

**Beispiel 17.** Die nichtlineare Gleichung  $x^2 + \exp(x) = 2$  hat eine eindeutige Lsg.  $x^* > 0$ . Es gibt mehrere naive Fixpunktformulierungen:

$$\begin{aligned} \Phi_1(x) &= x \pm (x^2 + \exp(x) - 2), & \Phi'_1(x^*) &\approx 1 \pm 2,79 \\ \Phi_2(x) &= \sqrt{2 - \exp(x)}, & \Phi'_2(x^*) &\approx 1,59 \\ \Phi_3(x) &= \log(2 - x^2), & \Phi'_3(x^*) &\approx 0,63 \end{aligned}$$

**Beispiel 18** (Newton-Verfahren). Gegeben  $f : \mathbb{R} \rightarrow \mathbb{R}$  diffbar mit  $f(x^*) = 0 \neq f'(x^*)$ .

Vorgeben: Gegeben  $x_n$ , berechne  $x_{n+1}$  als Nullstelle der Tangente an  $x_n$ , d.h.  $0 = f(x_n) + f'(x_n)(x_{n+1} - x_n) \implies x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$

$\implies \Phi(x) = x - f'(x)^{-1}f(x)$  Iterationsvorschrift

Falls  $f$  2x stetig diffbar gilt  $\Phi'(x) = 1 - \frac{f'(x)^2 - f(x)f''(x)}{f'(x)^2} \implies \Phi'(x^*) = 0$ , d.h. das Newton-Verfahren ist lokal von Ordnung  $p = 2$  (quadratisch konvergent).

**Beispiel 19** (Heron-Verfahren). Gegeben  $z > 0$ , definiere  $x_1 := \frac{1}{2}(1 + z)$ ,  $x_{n+1} := \frac{1}{2}(x_n + \frac{z}{x_n}) \forall n \in \mathbb{N}$ . gezeigt:  $\lim_n x_n =: x = \sqrt{z}$  (monoton fallend)

Betrachte  $x = \sqrt{z} \iff f(x) := x^2 - z = 0 \implies \Phi(x) = x - \frac{f(x)}{f'(x)} = x - \frac{x^2 - z}{2x} = \frac{1}{2}(x + \frac{z}{x})$ , d.h. spezielles Newton-Verf. (sogar lokal konv.).

**Bemerkung 46.** Im allgemeinen ist das Newton-Verfahren nicht global konvergent, sondern nur lokal konvergent, z.B.  $f(x) = \arctan(x)$  hat eindeutige Nst.  $x = 0$ . Klar:  $f$  ist beliebig nett (glatt, strikt monoton mit  $f'(x) = \frac{1}{1+x^2}$ ). Aber man kann zeigen, dass Newton divergent für jeden Startwert  $x_0$  mit  $|x_0| > y$  und  $y$  löst  $2y = (1 + y^2) \arctan(y)$ ,  $y \approx 1,37$

**Bemerkung 47.** Jedes iterative Verfahren zeigt i.a. 3 Phasen

- **vorasymptotische Phase:** Üblicherweise wird der Startwert  $x_0 \in X$  zufällig gewählt, d.h. es ist unklar, ob die Iteriertenfolge konvergiert.
- **asymptotische Phase:** Die Iterierte  $x_{k+1}, \dots, x_l$  konvergieren mit Konvergenzordnung gemäß Theorie.
- **nachasymptotische Phase:** Aufgrund der Rechnerarithmetik und Auslöschung zeigen  $x_{l+1}, \dots$  keine Konvergenz mehr.

Klarerweise will man die Iteration also nach  $x_l$  abbrechen, aber wie erkennt man  $x_l$ ?

Eine mögliche Realisierung für Nullstellensuche ist folgende:

Input:  $x_0 \in X, \tau_{abs}, \tau_{rel} > 0$  Toleranzen und  $K, L \in \mathbb{N}$  maximale Iterationszahlen:

1. Berechne  $x_1, \dots, x_k$  bis

- entweder  $k = K$  (dann Fehlerabbruch)
- oder  $|f(x_n)| \leq \max\{\tau_{abs}, \tau_{rel}|f(x_0)|\}$

2. Berechne  $x_{k+1}, \dots, x_l$  bis

- entweder  $l = L$  (dann Fehlerabbruch)
- oder  $|f(x_l)| \leq \tau_{abs}$  und  $|x_l - x_{l-1}| \leq \max\{\tau_{abs}, \tau_{rel}|x_l|\}$

z.B.  $\tau_{rel} = 10^{-6}, \tau_{abs} = 10^{-12}, K = 20, L = 10$  (zumindest für Newton)

## 5.2 Newton in $\mathbb{R}^d$

**Lemma 13.** Sei  $A \in \mathbb{R}^{d \times d}$  regulär,  $B \in \mathbb{R}^{d \times d}$ ,  $\|\cdot\|$  eine Norm auf  $\mathbb{R}^d$  und  $\|M\| := \sup_{x \in \mathbb{R}^d \setminus \{0\}} \frac{\|Mx\|}{\|x\|}$  die induzierte Operatornorm. Dann gilt  $\|B - A\| < \frac{1}{\|A^{-1}\|} \implies B$  ist regulär und  $\|B^{-1}\| < \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|B - A\|}$ , d.h.  $\{B \in \mathbb{R}^{d \times d} | B \text{ regulär}\} \subset \mathbb{R}^{d \times d}$  offen

*Proof.* 1. Sei  $M \in \mathbb{R}^{d \times d}$  mit  $\|M\| < 1 \implies \sum_{j=0}^{\infty} M^j$  (absolut) konvergent (da  $\|AB\| \leq \|A\|\|B\|$ )

$$\implies \left( \sum_{j=0}^{\infty} M^j \right) (I - M) = \sum_{j=0}^{\infty} M^j - \sum_{j=1}^{\infty} M^j = I$$

$$\implies (I - M) \text{ invertierbar, } (I - M)^{-1} = \sum_{j=0}^{\infty} M^j \text{ und } \|(I - M)^{-1}\| \leq \sum_{j=0}^{\infty} \|M\|^j = \frac{1}{1 - \|M\|}$$

2.  $M := A^{-1}(A - B), \|M\| < 1, I - M = A^{-1}(A - (A - B)) = A^{-1}B$  regulär  $\implies B$  regulär

$$\|B^{-1}\| < \left\| \overbrace{B^{-1}A}^{(I-M)^{-1}} \right\| \|A^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|M\|} \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\|\|A - B\|}$$

□

**Satz 21** ((gedämpftes) Newton-Verfahren).  $\|\cdot\|$  Norm auf  $\mathbb{R}^d$ ,  $\Omega \subseteq \mathbb{R}^d$  offen,  $F \in \mathcal{C}^2(\Omega; \mathbb{R}^d), x_* \in \Omega$  mit  $F(x_*) = 0$  und  $DF(x_*) \in \mathbb{R}^{d \times d}$  regulär,  $0 < \lambda \leq 1$  und  $\lambda \leq \lambda_n \leq 1 \forall n \in \mathbb{N}$   
 $\implies \exists \epsilon > 0$  mit folgenden Eigenschaften:

1. Für alle  $x_0 \in U_\epsilon(x_*)$  ist die induktiv definierte Folge  $x_{n+1} := x_n - \lambda_n DF(x_n)^{-1} F(x_n) \forall n \in \mathbb{N}_0$  wohldefiniert mit  $x_n \in U_\epsilon(x_*) \forall n \in \mathbb{N}_0$ , d.h. das (gedämpfte) Newton-Verfahren ist wohldefiniert.
2. Ex.  $0 < q < 1$  mit  $\|x_* - x_n\| \leq q \|x_* - x_{n-1}\| \forall n$ , d.h. das gedämpfte Newton-Verfahren konvergiert linear.
3. Falls  $\lambda_n = 1 \forall n$ , so ex.  $C > 0$  mit  $\|x_* - x_n\| \leq C \|x_* - x_{n-1}\|^2 \forall n$ , d.h. das klass. NV konvergiert quadratisch.

*Proof.* O.B.d.A.  $\|\cdot\| = \|\cdot\|_2$

1.  $DF : \Omega \rightarrow \mathbb{R}^{d \times d}$  stetig und  $\{B \in \mathbb{R}^{d \times d} \text{ regulär}\} \subseteq \mathbb{R}^{d \times d}$  ist offen.

$\implies \exists \epsilon > 0$  mit  $DF(x)$  regulär  $\forall x \in U_\epsilon(x_*)$  und  $\overline{U_\epsilon(x_*)} \subseteq \Omega$

$$\implies M := \sup_{x \in U_\epsilon(x_*)} \|DF(x)^{-1}\|_2 < \infty$$

$$\tilde{M} := \sup_{x \in U_\epsilon(x_*)} \left( \sum_{j,k,l=1}^d \left| \frac{\partial^2 F_j}{\partial x_k \partial x_l}(x) \right|^2 \right)^{\frac{1}{2}} < \infty$$

2. zz:  $\|F(y) - F(x) - DF(x)(y - x)\|_2 \leq \frac{\tilde{M}}{2} \|y - x\|_2^2 \forall x, y \in U_\epsilon(x^*)$

Betrachte  $f(t) := F(x + t(y - x))$

$$\begin{aligned}
f'(t) &= DF(x + t(y - x))(y - x) \\
f'_j(t) &= \sum_{k=1}^d \frac{\delta F_j}{\delta x_k}(x + t(y - x))(y_k - x_k) \\
f''_j(t) &= \sum_{k=1}^d (y_k - x_k) D \frac{\delta F_j}{\delta x_k}(x + t(y - x))(y - x) = \sum_{k,l=1}^d (y_k - x_k) \frac{\delta^2 F_j}{\delta x_k \delta x_l}(x + t(y - x))(y_l - x_l) \\
\Rightarrow \int_0^1 (1-t) f''(t) dt &= \int_0^1 f'' dt - \underbrace{\int_0^1 t f''(t) dt}_{=[t f'_1(t)]_{t=0}^1 - \int_0^1 f' dt} = f'(1) - f'(0) - f'(1) + f(1) - f(0) \\
&\Rightarrow F(y) - F(x) - DF(x)(y - x) = \int_0^1 (1-t) f''(t) dt \\
\Rightarrow \|F(y) - F(x) - DF(x)(y - x)\|_2 &\leq \int_0^1 (1-t) \|f''(t)\|_2 dt \leq \tilde{M} \|y - x\|_2^2 \frac{1}{2} \\
\|f''(t)\|_2^2 &= \sum_{j=1}^d |f''_j(t)|^2 \leq \|y - x\|_2^4 \sum_{j,k,l} \underbrace{\left| \frac{\delta^2 F_j}{\delta x_k \delta x_l}(x + \dots) \right|^2}_{\leq \tilde{M}^2}
\end{aligned}$$

Bemerkung: Beweis von (1), (2) geht für jede Norm, aber mit anderer Konstante  $\tilde{M}$ .

3. zz:  $\|x * -y\|_2 \leq \{(1 - \lambda) + \frac{M\tilde{M}}{2} \|x * -x\|_2\} \|x * -x\|_2 \forall x \in U_\epsilon(x^*), \forall \tilde{\lambda} \in [\lambda, 1], y := x - \tilde{\lambda} DF(x)^{-1} F(x)$

$$\begin{aligned}
x * -y &= x * -x - \tilde{\lambda} DF(x)^{-1} \underbrace{\{F(x^*) - F(x)\}}_{=0} = \\
&= (1 - \tilde{\lambda})(x * -x) - \tilde{\lambda} DF(x)^{-1} \{(F(x^*) - F(x)) - DF(x)(x * -x)\} \\
\Rightarrow \|x * -y\|_2 &\leq (1 - \tilde{\lambda}) \|x * -x\|_2 + \tilde{\lambda} \|DF(x)^{-1}\|_2 \frac{\tilde{M}}{2} \|x * -x\|_2^2 \leq (1 - \lambda) \|x * -x\|_2 + \frac{M\tilde{M}}{2} \|x * -x\|_2^2
\end{aligned}$$

4. Wähle  $\epsilon$  ggf noch kleiner als in den vorausgegangenen Schritten, damit  $q := (1 - \lambda) + \frac{M\tilde{M}}{2} \epsilon < 1$

$\Rightarrow$  gedämpftes Newton ist wohldef, linear konvergent und  $x_n \in U_\epsilon(x^*)$ , wenn  $x_0 \in U_\epsilon(x^*)$  und (ungedämpftes) Newton ist wohldef., mit  $x_n \in U_\epsilon(x^*)$  für  $x_0 \in U_\epsilon(x^*)$  und quadratische Konvergenz mit  $C = \frac{M\tilde{M}}{2}$ .

□

**Korollar 4.** Für das Newton-Verfahren gilt mit  $x_0 \in U_\epsilon(x^*)$  und  $\epsilon$  klein genug:

$$C^{-1} \|x * -x_n\| \leq \overbrace{\|F(x_n)\|}^{\text{Residuum}} \leq C \|x * -x_n\|$$

und insb.  $\|F(x_n)\| \leq C^4 \|F(x_{n-1})\|^2 \forall n \in \mathbb{N}$  wobei  $C > 0$  unabhängig.

*Proof.* Nach Newton-Beweis gilt

$$\begin{aligned}
\|y - x - DF(x)^{-1} \{F(y) - F(x)\}\| &\leq \frac{M\tilde{M}}{2} \|y - x\|^2 \\
\Rightarrow \|x - x * \| &\leq \underbrace{\frac{M\tilde{M}}{2} \|x - x * \|^2}_{\leq \frac{1}{2} \|x - x * \|, \text{ falls } \frac{M\tilde{M}}{2} \epsilon \leq \frac{1}{2}} + \|DF(x^*)^{-1} \{F(x) - \underbrace{F(x^*)}_{=0}\}\| \\
&\Rightarrow \frac{1}{2} \|x - x * \| \leq M \|F(x)\|
\end{aligned}$$

Ferner

$$\begin{aligned} \|F(x)\| &\leq \underbrace{\|DF(x^*)\|}_{\leq M} \underbrace{\|DF(x^*)^{-1}\{F(x) - F(x^*)\}\|}_{=0} \leq M \frac{3}{2} \|x - x^*\| \\ &\leq \|x - x^*\| + \underbrace{\frac{M\tilde{M}}{2} \|x - x^*\|^2}_{\leq \frac{1}{2} \|x - x^*\|} \end{aligned}$$

(bisher alles unabhängig von Newton)

Mit Newton:

$$\begin{aligned} \|F(x_k)\| &\leq C \underbrace{\|x^* - x_k\|}_{\leq \tilde{C} \|x^* - x_{n-1}\|^2} \leq C^3 \tilde{C} \|F(x_{n-1})\|^2 \\ &\leq \tilde{C} \underbrace{\|x^* - x_{n-1}\|^2}_{\leq C^2 \|F(x_{n-1})\|^2} \end{aligned}$$

□

**Satz 22** (Zweck der Dämpfung).  $\Omega \subseteq \mathbb{R}^d$  offen,  $F \in \mathcal{C}^2(\Omega, \mathbb{R}^d)$  mit  $DF(x)$  regulär  $\forall x \in \Omega$ ,  $K \subset \Omega$  kompakt  
 $\implies$  Ex.  $\lambda_{max}, \gamma > 0$  mit

$$\|F(\underbrace{x - \lambda DF(x)^{-1} F(x)}_{\text{gedämpfter Newton-Schritt}})\|_2^2 \leq \underbrace{(1 + \gamma \lambda^2 - 2\lambda)}_{\leq q < 1 \text{ für } \lambda \text{ klein}} \|F(x)\|_2^2 \forall x \in K \forall 0 < \lambda \leq \lambda_{max}$$

*Proof.* Betrachte  $g(t) := \|F(x - t\lambda p)\|_2^2$  mit  $p := DF(x)^{-1} F(x)$ ,  $\lambda$  freier Parameter

$$\begin{aligned} \implies g(1) &= g(0) + g'(0) + \int_0^1 (1-t)g''(t)dt \\ f(y) &:= \|F(y)\|_2^2 = F(y) \cdot F(y) \\ \implies Df(y) &= 2F(y)^T DF(y) \\ g(t) &= f(x - t\lambda p) \\ g'(t) &= Df(x - t\lambda p)(-\lambda p) = -2F(x - t\lambda p)^T DF(x - t\lambda p)(\lambda p) \\ g'(0) &= -2\lambda F(x)^T \underbrace{DF(x)p}_{=F(x)} = -2\lambda \|F(x)\|_2^2 \\ g''(t) &= +\lambda p \cdot D^2 f(x - t\lambda p)\lambda p \\ \implies \|g''(t)\| &\leq \sup_t \|D^2 f(x - t\lambda p)\|_2 \|\lambda p\|_2^2 \\ &\leq C < \infty \\ \implies \underbrace{\|F(x - \lambda DF(x)^{-1} F(x))\|_2^2}_{=g(p)} &\leq (1 - 2\lambda + \frac{1}{2}CM^2\lambda^2) \|F(x)\|_2^2 \end{aligned}$$

klar:  $D = \sup_{x \in K} \|F(x)\|_2 < \infty$ ,  $M = \sup_{x \in K} \|DF(x)^{-1}\|_2 < \infty$ ,  $dist(K, \delta\Omega) := \inf\{\|x - y\|_2 | x \in K, y \in \delta\Omega\} > 0$ ,  $\lambda_{max} := \min\{1, \frac{1}{2DM} dist(K, \delta\Omega)\} > 0$ ,  $\tilde{K} := \{x - \lambda p | x \in K, p \in \mathbb{R}^d, \|p\|_2 \leq DM, \lambda \in [0, \lambda_{max}]\} \subseteq \mathbb{R}^d$  kompakt.

Sogar  $\tilde{K} \subseteq \Omega$

$$\implies \sup_{\tilde{x} \in \tilde{K}} \|D^2 f(\tilde{x})\|_2 =: C < \infty$$

□

**Algorithmus 4** (gedämpftes Newton-Verfahren). *Input:*  $F : \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $x_0 \in \mathbb{R}^d$ ,  $0 < q < 1$ ,  $0 < \lambda_{max} < 1$

*Initialisierung:*  $\lambda_0 := 1$

*Schleife:* For  $k = 0, 1, 2, \dots$  iteriere

- Berechne  $A := DF(x_k)$ ,  $b := -F(x_k)$
- Löse  $Ap_k = b$
- Iteriere die Berechnung  $x_{k+1} := x_k + \lambda_k p_k$  bis
  - entweder  $\|F(x_{k+1})\|_2 < \|F(x_k)\|_2 \rightarrow ok$
  - oder  $\lambda_k := q\lambda_k < \lambda_{min} \rightarrow \text{Fehlerabbruch}$



- danach  $\lambda_{k+1} := \min\{1, \frac{\lambda_k}{q}\}$
- Fertig, falls  $x_{k+1}$  hinreichend dicht bei  $x^*$  mit  $F(x^*) = 0$

**Bemerkung 48.** In der Praxis ist die Berechnung von  $DF(x_n)$  teuer und man verwendet billigere Approximationen, sog. **Quasi-Newton-Verfahren** z.B. Sekantenverfahren in 1D (erste Ableitung durch Differenzenquotient) oder Broyden-Verfahren in  $\mathbb{R}^d$ .

### 5.3 Stationäre Iterationsverfahren zur Lösung Linearer GLS

Gegeben  $A \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n$

Gesucht  $x^* \in \mathbb{K}^n$  mit  $Ax^* = b$

Vorgehen: Definiere  $M \in \mathbb{K}^{n \times n}$ ,  $c \in \mathbb{K}^n$ , wähle  $x_0 \in \mathbb{K}^n$  und betrachte  $x_{k+1} := \Phi(x_k)$  mit  $\Phi(x) = Mx + c$ .

**Bemerkung 49.** Die iterative Lsg. eines linearen GLS ist dann sinnvoll, wenn  $A$  schwach besetzt ist (d.h. nur  $\mathcal{O}(n)$  nicht-null Einträge), aber ohne Struktur für Eliminationsverfahren (d.h. Gauss bräuchte  $\mathcal{O}(n^2)$  Speicher) oder wenn  $\text{cond}(A)$  groß ist (da sich iterative Verfahren "selbst stabilisieren"). Üblicherweise brauchen iterative Verfahren nur eine (schnelle) Matrix-Vektor-Multiplikation (z.B. FFT).

**Lemma 14.** Für  $M \in \mathbb{K}^{n \times n}$  mit **Spektrum**  $\sigma(M) = \{\lambda \in \mathbb{C} \text{ EW von } M\}$  und **Spektralradius**  $\rho(M) = \max_{\lambda \in \sigma(M)} |\lambda|$  gilt  $\rho(M) = \inf\{\|M\| : \|\cdot\| \text{ Norm auf } \mathbb{C}^n \text{ und } \|M\| := \sup_{x \in \mathbb{C}^n \setminus \{0\}} \frac{\|Mx\|}{\|x\|}\}$

*Proof.*  $\leq$ : Sei  $\|\cdot\|$  Norm auf  $\mathbb{C}^n$ , Sei  $\lambda \in \sigma(M)$ ,  $x \in \mathbb{C}^n \setminus \{0\}$  mit  $Mx = \lambda x$   
 $\implies |\lambda| \|x\| = \|\lambda x\| = \|Mx\| \leq \|M\| \|x\| \implies |\lambda| \leq \|M\| \implies \rho(M) \leq \|M\|$ .  
 $\geq$ : Ex.  $T \in \mathbb{C}^{n \times n}$  regulär mit

$$R := T^{-1}MT = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & \dots & r_{2n} \\ & & \ddots & \vdots \\ 0 & & & r_{nn} \end{pmatrix}$$

und  $\sigma(M) = \sigma(R) = \{r_{11}, \dots, r_{nn}\}$

zz:  $\forall \epsilon > 0 \exists \|\cdot\|_\epsilon$  Norm auf  $\mathbb{C}^n$  :  $\|M\|_\epsilon \leq \underbrace{\max_j |r_{jj}|}_{=\rho(M)} + \epsilon$

Sei  $\epsilon > 0$ , definiere  $\|x\|_\epsilon := \|D_\epsilon^{-1}T^{-1}x\|_\infty$ ,  $D_\epsilon := \text{diag}(1, \epsilon, \dots, \epsilon^{n-1})$

$$\|M\|_\epsilon = \sup_{x \in \mathbb{C}^n \setminus \{0\}} \frac{\|D_\epsilon^{-1}T^{-1}Mx\|_\infty}{\underbrace{\|D_\epsilon^{-1}T^{-1}x\|_\infty}_{=y}} = \sup_{y \in \mathbb{C}^n \setminus \{0\}} \frac{\|D_\epsilon^{-1}T^{-1}MTD_\epsilon y\|_\infty}{\|y\|_\infty} = \sup_{y \in \mathbb{C}^n \setminus \{0\}} \frac{\|D_\epsilon^{-1}RD_\epsilon y\|_\infty}{\|y\|_\infty} = \|D_\epsilon^{-1}RD_\epsilon\|_\infty$$

$$D_\epsilon^{-1}RD_\epsilon = \begin{pmatrix} r_{11} & \epsilon r_{12} & \dots & \epsilon^{n-1} r_{1n} \\ & r_{22} & \dots & \epsilon^{n-2} r_{2n} \\ & & \ddots & \vdots \\ 0 & & & r_{nn} \end{pmatrix}$$

UE: Operatornorm zu  $\|\cdot\|_\infty$  ist die Zeilensummennorm

$$\implies \|M\|_\epsilon = \|D_\epsilon^{-1}RD_\epsilon\|_\infty = \max_{j=1, \dots, n} \sum_{k=1}^n |(D_\epsilon^{-1}RD_\epsilon)_{jk}| = \max_{j=1, \dots, n} \left( |r_{jj}| + \epsilon \underbrace{\sum_{k=j+1}^n \epsilon^{k-(j-1)} |r_{jk}|}_{\leq C} \right) \leq \rho(M) + \epsilon C$$

□

**Bemerkung 50.** Der Beweis zeigt, dass das Infimum über Normen auf  $\mathbb{R}^n$  gebildet werden kann, wenn  $M$  trigonalisierbar über  $\mathbb{R}$ .

**Satz 23** (globale Konvergenz). Für  $M \in \mathbb{K}^{n \times n}$  sind äquivalent:

1.  $\rho(M) < 1$
2. Für alle  $c \in \mathbb{K}^n$  und alle  $x_0 \in \mathbb{K}^n$  konvergiert die Iteriertenfolge  $x_{n+1} := \Phi(x_n)$ ,  $\Phi(x) := Mx + c$

In diesem Fall ist  $x^* := \lim_n x_n$  sogar eindeutig und unabhängig von  $x_0$ , und  $(\mathbb{K}^n, \Phi, x^*)$  konvergiert global linear für eine geeignete Norm auf  $\mathbb{K}^n$ .

(Beweis zeigt, dass Ex. von einem  $c \in \mathbb{K}^n, x_0 \in \mathbb{K}^n$  bereits ausreicht).

*Proof.* (i)  $\implies$  (ii) Nach Lemma existiert  $\|\cdot\|$  auf  $\mathbb{C}^n$  mit  $\|M\| < 1$ . Sei  $c \in \mathbb{K}^n$ .

$$\implies \|\Phi(x) - \Phi(y)\| = \|M(x - y)\| \leq \underbrace{\|M\|}_{=q < 1} \|x - y\|.$$

$\implies$  Behauptung folgt aus Banachschem Fixpunktsatz

gezeigt (ii) für  $\mathbb{K} = \mathbb{C}$

$\mathbb{K} = \mathbb{R}$ : auch Ok, da  $\Phi : \mathbb{R}^n \rightarrow \mathbb{R}^n$ , d.h.  $x^* \in \mathbb{R}^n$

Beweis (ii)  $\implies$  (i) für  $\mathbb{K} = \mathbb{C}$

Sei  $\|\cdot\|$  Norm auf  $\mathbb{C}^n$ . Wähle  $x_0 \in \mathbb{C}^n$  mit  $-x^* + x_0 \in \mathbb{C}^n \setminus \{0\}$  Eigenvektor zum Eigenwert  $\lambda \in \sigma(M)$  mit  $|\lambda| = \rho(M)$

$$\begin{aligned} \implies \underbrace{\|x_k - x^*\|}_{\rightarrow 0} &= \|M(x_{k-1} - x^*)\| = \underbrace{\|M^k(x_0 - x^*)\|}_{=|\lambda|^k} \underbrace{\|x_0 - x^*\|}_{\neq 0} \implies \rho(M) < 1. \\ &= \underbrace{|\lambda|^k}_{=\rho(M)^k} \underbrace{\|x_0 - x^*\|}_{\neq 0} \end{aligned}$$

(ii)  $\implies$  (i) für  $\mathbb{K} = \mathbb{R}$  zz: (ii) gilt auch für  $\mathbb{C}$

Sei  $c \in \mathbb{C}^n, x_0 \in \mathbb{C}^n$ , def  $x_{k+1} = \Phi(x_k)$  und  $a_k := \operatorname{Re} x_k, b_k := \operatorname{Im} x_k$

$$\implies a_{k+1} = \Phi(a_k), \text{ denn } \operatorname{Re} \Phi(x) = M(\operatorname{Re} x) + \operatorname{Re} c, b_{k+1} = \Phi(b_k)$$

(ii) für  $\mathbb{R} \implies a_k \rightarrow a^*, b_k \rightarrow b^* \implies x_n \rightarrow x^* := a^* + ib^*$  □

**Beispiel 20** (Richardson-Iteration).  $A \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n, \Phi(x) := \underbrace{(I - \lambda A)}_{=M} x + \underbrace{\lambda b}_{=c}$  mit  $\lambda \in \mathbb{K}$  geeignet.

klar:  $\Phi(x) = x \iff Ax = b$

benötigt:  $\rho(I - \lambda A) < 1$

**Beispiel 21** (Jacobi-Iteration).  $D \in \mathbb{K}^{n \times n}$  Diagonale von  $A$ , d.h.  $D_{jk} = \begin{cases} A_{jj} & \text{für } j = k \\ 0 & \text{sonst} \end{cases}$

vorausgesetzt:  $A_{jj} \neq 0 \forall j$ , damit  $D$  invertierbar

$$\Phi(x) := \underbrace{-D^{-1}(A - D)x}_{=M} + \underbrace{D^{-1}b}_{=c}$$

klar:  $\Phi(x) = x \iff Dx = -(A - D)x + b \iff Ax = b$

**Beispiel 22** (Gauss-Seidel-Iteration).  $L, U \in \mathbb{K}^{n \times n}$  mit  $L_{jk} = \begin{cases} A_{jk} & \text{für } j \geq k \\ 0 & \text{sonst} \end{cases}, U_{jk} = \begin{cases} A_{jk} & \text{für } j < k \\ 0 & \text{sonst} \end{cases}$

vorausgesetzt:  $A_{jj} \neq 0 \forall j$ , damit  $L$  invertierbar

$$\Phi(x) := \underbrace{-L^{-1}Ux}_{=M} + \underbrace{L^{-1}b}_{=c}$$

$\Phi(x) = x \iff Lx = -Ux + b \iff Ax = b$

**Satz 24** (Konvergenz bei strikter Diagonaldominanz). Sei  $A \in \mathbb{K}^{n \times n}$  **strikt diagonaldominant**, d.h.  $\sum_{k=1, k \neq j}^n |A_{jk}| < |A_{jj}| \forall j = 1, \dots, n$

$\implies A$  ist regulär und Jacobi- und Gauss-Seidel-Verfahren sind wohldefiniert und konvergent mit  $\|M^{(GS)}\|_\infty \leq \|M^{(J)}\|_\infty < 1$

*Proof.* klar:  $A$  ist regulär (UE) und  $A_{jj} \neq 0 \forall j$ , d.h. JV und GSV sind wohldef.

$$1. M^{(J)} = -D^{-1}(A - D); \|M^{(J)}\|_\infty = \max_{j=1, \dots, n} \sum_{k=1, k \neq j}^n \underbrace{|M_{jk}^{(J)}|}_{= \frac{|A_{jk}|}{|A_{jj}|}} < 1$$

2. Von nun an  $|\cdot|$  und  $\leq$  komponentenweise

$$\begin{aligned} D^{-1}L &= I + D^{-1}(L - D) \text{ und } \rho(D^{-1}(L - D)) = 0 \\ \implies (D^{-1}L)^{-1} &= (I - (-D^{-1}(L - D)))^{-1} = \sum_{k=0}^{\infty} (-D^{-1}(L - D))^k \\ \implies |(D^{-1}L)^{-1}| &\leq \sum_{k=0}^{\infty} |D^{-1}(L - D)|^k = (I - |D^{-1}(L - D)|)^{-1} \end{aligned}$$

$$3. M^{(J)} = -D^{-1}(A - D) = -(D^{-1}(L - D) + D^{-1}U)$$

$$\begin{aligned} &\Rightarrow |M^{(J)}| = |D^{-1}(L - D)| + |D^{-1}U| \\ &\Rightarrow |D^{-1}U| = (|M^{(J)}| - I) + (I - |D^{-1}(L - D)|) \end{aligned}$$

$$4. \text{ zz: } \|M^{(GS)}\|_\infty \leq \|M^{(J)}\|_\infty$$

$$\begin{aligned} |M^{GS}| &= |L^{-1}U| = |(D^{-1}L)^{-1}D^{-1}U| \leq |(D^{-1}L)^{-1}||D^{-1}U| \leq \\ &(I - |D^{-1}(L - D)|)^{-1}[ (|M^{(J)}| - I) + (I - |D^{-1}(L - D)|) ] = (I - |D^{-1}(L - D)|)^{-1}(|M^{(J)} - I|) + I \end{aligned}$$

Def  $e = (1, \dots, 1)$

$$\begin{aligned} |M^{GS}|e &\leq \underbrace{(I - |D^{-1}(L - D)|)^{-1}}_{\geq I} \underbrace{(|M^{(J)}|e - e)}_{\substack{\leq \|M^{(J)}\|_\infty e - e = (\|M^{(J)}\|_\infty - 1)e \\ \in \mathbb{R}_{<0}}} + e \leq (\|M^{(J)}\|_\infty - 1)e + e = \|M^{(J)}\|_\infty e \\ &\Rightarrow \|M^{GS}\|_\infty = \| |M^{GS}|e \|_\infty \leq \| \|M^{(J)}\|_\infty e \|_\infty = \|M^{(J)}\|_\infty \end{aligned}$$

□

**Korollar 5.** Sei  $A \in \mathbb{K}^{n \times n}$  mit  $A^T$  strikt diagonaldominant  $\Rightarrow A$  regulär, und das Jacobi-Verfahren ist wohldef und konvergent.

*Proof.*  $A^T$  strikt diagonaldominant  $\Rightarrow A_{jj} \neq 0 \forall j$ , d.h. Jacobi wohldef. und  $A$  regulär, da  $\det(A^T) = \det(A)$

UE: Die  $l_1$ -Norm  $\|x\|_1 = \sum_j |x_j|$  induziert als Operatornorm die Spaltensummennorm  $\|M\|_1 = \max_{k=1, \dots, n} \sum_{j=1}^n |M_{jk}|$   
 $\|M\|_1 = \|M^T\|_\infty$  und analog zum letzten Beweis:  $A^T$  strikt diagonaldominant  $\Rightarrow \|M^J\|_1 < 1 \Rightarrow \rho(M^{(J)}) < 1$  □

## 5.4 Krylov-Verfahren zur Lsg Linearer GLS

Gauss-Seidl und Jacobi benötigen Zugriff auf Matrixkoeffizienten.

Ziel: Löse  $Ax = b$  mit  $A \in \mathbb{K}^{n \times n}$  regulär nur durch Verwendung der Matrix-Vektor-Multiplikation, ohne auf Einträge explizit zuzugreifen.

**Lemma 15** (Krylov-Räume).  $A \in \mathbb{K}^{n \times n}$  regulär,  $n \in \mathbb{N} \setminus \{0\}$ ,  $x^* \in \mathbb{K}^n$  mit  $Ax^* = b$ . Zu  $l \in \mathbb{N}$  definiere  $\mathcal{K}_l := \mathcal{K}_l(A, b) := \text{span} \underbrace{\{b, Ab, \dots, A^{l-1}b\}}_{=\{A^j b | j=0, \dots, l-1\}}$  **Krylov-Räume**.

Dann sind äquivalent:

1.  $\dim \mathcal{K}_{l+1} \leq l$
2.  $\mathcal{K}_l = \mathcal{K}_{l+1}$
3.  $A(\mathcal{K}_l) \subseteq \mathcal{K}_l$
4.  $x^* \in \mathcal{K}_l$

*Proof.* (i  $\Rightarrow$  ii): Wähle  $m < l$  minimal mit  $\{b, Ab, \dots, A^m b\}$  linear abhängig.  $\Rightarrow \exists \lambda_j \in \mathbb{K}$  mit  $A^m b = \sum_{j=0}^{m-1} \lambda_j A^j b$

zz.  $\mathcal{K}_{l+1} \subseteq \mathcal{K}_l$ , d.h. zz:  $A^l b \in \mathcal{K}_l$   
 $A^l b = A^{l-m}(A^m b) = \sum_{j=0}^{m-1} \lambda_j \underbrace{A^{j+l-m} b}_{\in \mathcal{K}_l} \in \mathcal{K}_l$

(ii  $\Rightarrow$  iii)  $A(\underbrace{A^j b}_{\text{Basiselm von } \mathcal{K}_l}) \in \mathcal{K}_{l+1} = \mathcal{K}_l \forall j = 0, \dots, l-1 \Rightarrow A(\mathcal{K}_l) \subseteq \mathcal{K}_l$

(iii  $\Rightarrow$  iv)  $A : \mathcal{K}_l \rightarrow \mathcal{K}_l$  linear, wohldef., injektiv  $\Rightarrow$  bijektiv,  $b \in \mathcal{K}_l \Rightarrow \underbrace{A^{-1}b}_{=x^*} \in \mathcal{K}_l$

(iv  $\Rightarrow$  i):  $x^* \in \mathcal{K}_l = \text{span}\{b, \dots, A^{l-1}b\} \Rightarrow b = Ax^* \in \text{span}\{Ab, \dots, A^l b\} \Rightarrow \underbrace{\{b, Ab, \dots, A^l b\}}_{\text{erzeugt } \mathcal{K}_{l+1}} \text{ lin. abh.}$

$\Rightarrow \dim \mathcal{K}_{l+1} \leq l$  □

**Beispiel 23** (Krylov-Verfahren). • **CG-Verfahren:** Für  $A \in \mathbb{K}^{n \times n}$  SPD (selbstadjungiert  $A^H = A$ , positiv definit) berechnet CG die Iterierten  $x_l \in \mathcal{K}_l = \mathcal{K}_l(A, b)$  mit  $\|x * -x_l\|_A = \min_{y_l \in \mathcal{K}_l} \|x * -y_l\|_A$  mit  $\|y\|_A := (y^H A y)^{\frac{1}{2}}$

• **CGNR-Verfahren:** Für  $A \in \mathbb{K}^{n \times n}$  regulär berechne  $x_l \in \mathcal{K}_l := \mathcal{K}_l(\underbrace{A^H A}_{SPD}, A^H b)$  mit  $\|x * -x_l\|_{A^H A} =$

$$\min_{y_l \in \mathcal{K}_l} \|x * -y_l\|_{A^H A}$$

$$\text{betrachte } \|x * -y_l\|_{A^H A}^2 = (x * -y_l)^H A^H A (x * -y_l) = \|\underbrace{b}_{=Ax*} - Ay_l\|_2^2$$

• **GMRES-Verfahren:** Für  $A \in \mathbb{K}^{n \times n}$  regulär berechne  $x_l \in \mathcal{K}_l = \mathcal{K}_l(A, b)$  mit  $\|b - Ax_l\|_2 = \min_{y_l \in \mathcal{K}_l} \|b - Ay_l\|_2 (= \min_{z_l \in A(\mathcal{K}_l)} \|b - z_l\|_2)$

$\Rightarrow$  Alle Verfahren sind wohldefiniert (sogar mit eindeutigen  $x_l$ !) und brechen nach endlich vielen Schritten mit  $x_* = x_{l_*}$  ab und  $l_* \leq n$  (zumindest theoretisch!)

Im Folgenden betrachten wir nur noch CG (bzw. CGNR)  $\rightsquigarrow$  VO "Iterative Lösung großer Gleichungssysteme".

**Lemma 16** (Orthogonalprojektion).  $X$  Hilbert-Raum,  $Y \leq X$  Teilraum mit  $\dim Y = n$ ,  $\{y_1, \dots, y_n\} \subseteq Y$  Orthonormalbasis.

Definiere  $\mathcal{P} : X \rightarrow Y$ ,  $\mathcal{P}x := \sum_{j=1}^n \langle y_j, x \rangle y_j$

Dann gilt:

1.  $\mathcal{P}$  ist linear mit  $\mathcal{P}y = y \forall y \in Y$  und  $\langle x - \mathcal{P}x, y \rangle = 0 \forall x \in X \forall y \in Y$ , sog. **Orthogonalprojektion auf  $Y$** .
2.  $\|x - \mathcal{P}x\|_X = \min_{y \in Y} \|x - y\|$  und  $\|x - y\|_X^2 = \|x - \mathcal{P}x\|_X^2 + \|\mathcal{P}x - y\|_X^2 \forall x \in X \forall y \in Y$ , d.h.  $y = \mathcal{P}x$  ist der eindeutige Minimierer.

*Proof.* Für  $y \in Y$  existieren  $\lambda_j \in \mathbb{K}$  mit  $y = \sum_{j=1}^n \lambda_j y_j$

$$\begin{aligned} \Rightarrow \langle y_k, y \rangle &= \sum_{j=1}^n \lambda_j \underbrace{\langle y_k, y_j \rangle}_{=\delta_{jk}} = \lambda_k \\ \Rightarrow y &= \sum_{j=1}^n \langle y_j, y \rangle y_j = \mathcal{P}y \\ &= \sum_j \langle y_j, x \rangle y_j \\ \langle y_k, x - \mathcal{P}x \rangle &= \langle y_k, x \rangle - \underbrace{\langle y_k, \mathcal{P}x \rangle}_{=\sum_{j=1}^n \langle y_j, x \rangle \langle y_k, y_j \rangle = \langle y_k, x \rangle} = 0 \forall k = 1, \dots, n \\ \Rightarrow \langle x - \mathcal{P}x, y \rangle &= 0 \forall y \in Y \\ a := x - \mathcal{P}x, b := \underbrace{\mathcal{P}x - y}_{\in Y} &\Rightarrow \langle a, b \rangle_X = 0 \\ \Rightarrow \underbrace{\|a + b\|_X^2}_{=\|x - y\|_X^2} &= \|a\|_X^2 + \|b\|_X^2 = \|x - \mathcal{P}x\|_X^2 + \|\mathcal{P}x - y\|_X^2 \end{aligned}$$

□

**Satz 25.** Sei  $A \in \mathbb{K}^{n \times n}$  SPD und  $b, x_* \in \mathbb{K}^n$ ,  $Ax_* = b$ . Zu  $l \in \mathbb{N}_0$  sei  $\mathcal{K}_l = \text{span}\{A^j b \mid j = 0, \dots, l-1\}$  und  $x_l \in \mathcal{K}_l$  mit  $\|x * -x_l\|_A = \min_{y_l \in \mathcal{K}_l} \|x * -y_l\|_A$ . Sei  $l_* \leq n$  minimal mit  $x_* \in \mathcal{K}_{l_*}$ .

Definiere  $r_l := b - Ax_l$ , das sog. Residuum  $\Rightarrow \{r_0, \dots, r_{l-1}\} \subseteq \mathcal{K}_l$  Basis  $\forall l < l_*$  und Gram-Schmidt liefert  $\{d_0, \dots, d_{l-1}\} \subseteq \mathcal{K}_l$  orthogonale Basis bzgl.  $\langle \cdot, \cdot \rangle_A$ .

$\Rightarrow d_0 = b, d_{l+1} = r_{l+1} + \beta_l d_l$  mit  $\beta_l = \frac{\|r_{l+1}\|_2^2}{\|r_l\|_2^2}$ ,  $x_0 = 0, r_0 = b$  und  $x_{l+1} = x_l + \alpha_l d_l, r_{l+1} = r_l - \alpha_l A d_l$  mit  $\alpha_l = \frac{\|r_l\|_2^2}{\|d_l\|_A^2} \forall l < l_*$ .

**Algorithmus 5.**  $A \in \mathbb{K}^{n \times n}$  SPD,  $b \in \mathbb{K}^n$

Initialisierung:  $r_0 := b, d_0 := b, x_0 := 0$

Für alle  $l = 0, 1, 2, \dots$  iteriere:

1. Abbruch, falls  $r_l = 0$  (d.h.  $x_* = x_l$ )

2. Def  $\alpha_l := \frac{\|r_l\|_2^2}{d_l^H Ad_l}$  und  $x_{l+1} := x_l + \alpha_l d_l, r_{l+1} := r_l - \alpha_l Ad_l$

3. Def.  $\beta_l := \frac{\|r_{l+1}\|_2^2}{\|r_l\|_2^2}$  und  $d_{l+1} := r_{l+1} + \beta_l d_l$

Output:  $x^* = x_l$  und  $l = l^*$  minimal mit  $x^* \in \mathcal{K}_{l^*}$ .

**Bemerkung 51.** CG-Algorithmus braucht pro Schritt eine Matrix-Vektor-Multiplikation zur Berechnung  $Ad_l$ , danach nur Vektor- und Skalar-Operationen.

$\Rightarrow$  Aufwand =  $\mathcal{O}(n) + \mathcal{O}(\text{Matrix-Vektor-Multiplikation})$  pro Schritt.

Proof. 1.

$$\langle r_k, y \rangle_2 = \langle b - Ax_k, y \rangle_2 = \langle A(x^* - x_k), y \rangle_2 = \underbrace{\langle x^* - x_k, y \rangle_A}_{\text{Minimum, d.h. } x_k = \mathcal{P}_{\mathcal{K}_l} x^*} = 0 \forall y \in \mathcal{K}_l$$

$\Rightarrow$  Induktiv:  $\{r_0, \dots, r_{l-1}\} \subseteq \mathcal{K}_l$  Basis von  $\mathcal{K}_l$  und orthogonal bzgl.  $\langle \cdot, \cdot \rangle_2$  ✓

Gram-Schmidt:  $d_0 := r_0 := b, d_{k+1} = r_{k+1} - \sum_{j=0}^k \frac{\langle r_{k+1}, d_j \rangle_A}{\|d_j\|_A^2} d_j$

$\Rightarrow \{d_0, \dots, d_j\}$  orthogonal bzgl.  $\langle \cdot, \cdot \rangle_A$  und  $\text{span}\{d_0, \dots, d_j\} = \text{span}\{r_0, \dots, r_j\}$  und  $\langle r_{k+1}, d_j \rangle_A = \langle r_{k+1}, \underbrace{Ad_j}_{\in \mathcal{K}_{j+1}} \rangle_2$

$r_k \in \mathcal{K}_{k+1}$  mit  $r_k \perp \mathcal{K}_k$  in  $\langle \cdot, \cdot \rangle_2, d_k \in \mathcal{K}_{k+1}$  mit  $d_k \perp \mathcal{K}_k$  in  $\langle \cdot, \cdot \rangle_A$

$$\Rightarrow d_{k+1} = r_{k+1} - \underbrace{\frac{\langle r_{k+1}, d_k \rangle_A}{\|d_k\|_A^2}}_{=: \tilde{\alpha}_k} d_k$$

2.

$$x_{k+1} = \sum_{j=0}^k \frac{\langle x^*, d_j \rangle_A}{\|d_j\|_A^2} d_j = x_k + \underbrace{\frac{\langle x^*, d_k \rangle_A}{\|d_k\|_A^2}}_{=: \tilde{\alpha}_k} d_k$$

$$r_{k+1} = b - Ax_{k+1} = b - (x_k + \tilde{\alpha}_k d_k) = r_k - \tilde{\alpha}_k d_k$$

3. zz  $\tilde{\alpha}_k = \frac{\langle x^*, d_k \rangle_A}{\|d_k\|_A^2} \stackrel{!}{=} \frac{\|r_k\|_2^2}{\|d_k\|_A^2} = \alpha_k$

$$\|r_k\|_2^2 = \langle r_k, r_k \rangle_2 \stackrel{(1)}{=} \langle r_k, d_k \rangle_2 = \langle b - Ax_k, d_k \rangle_2 = \langle A(x^* - x_k), d_k \rangle_2 = \overline{\langle x^*, x_k, d_k \rangle_A} = \overline{\langle x^*, d_k \rangle_A} = \langle x^*, d_k \rangle_A$$

4. zz:  $\tilde{\beta}_k = -\frac{\langle r_{k+1}, d_k \rangle_A}{\|d_k\|_A^2} = \frac{\|r_{k+1}\|_2^2}{\|r_k\|_2^2} = \beta_k$

$$\|r_{k+1}\|_2^2 = -\alpha_k \langle r_{k+1}, Ad_k \rangle_2 = -\alpha_k \langle r_{k+1}, d_k \rangle_A = \overline{\tilde{\beta}_k} \|r_k\|_2^2$$

□

**Bemerkung 52.** Man kann zeigen, dass im CG-Verfahren stets gilt  $\|x^* - x_{k+1}\|_A \leq q \|x^* - x_k\|_A$  mit  $q = \left(1 - \frac{1}{\text{cond}_2(A)}\right)^{1/2}$ ,  $\text{cond}_2(A) = \|A\|_2 \|A^{-1}\|_2$

$\Rightarrow$  CG passt auch in den Rahmen von Banachschen Fixpunktsatz, d.h. man a priori und a posteriori Fehlerabschätzungen.

Beweisidee.  $\|x^* - x_{k+1}\|_A^2 = \min_{y_{k+1} \in \mathcal{K}_{k+1}} \|x^* - y_{k+1}\|_A^2 \leq \min_{t \in \mathbb{R}} \|x^* - (x_k + tr_k)\|_A^2$

$\Rightarrow \tilde{x}_{k+1} = x_k + t_{\min} r_k \Rightarrow \|x^* - \tilde{x}_{k+1}\|_A \leq q \|x^* - x_k\|_A \rightsquigarrow t$  ausrechnen! □

## 6 Eliminationsverfahren

**Eliminationsverfahren** oder **direkte Löser** sind Algorithmen zur Lösung von  $Ax = b$ , die in endlich vielen Schritten die exakte Lösung  $x$  berechnen (zumindest bei exakter Arithmetik), z.B. Gauss-Elimination.

**Bemerkung 53.** Aufgrund der Rechnerarithmetik und unvermeidlichen Fehlern in der rechten Seite gilt  $\frac{\|x - \tilde{x}\|}{\|x\|} \leq \text{cond}(A) \frac{\|b - \tilde{b}\|}{\|b\|}$  mit  $\tilde{b} \approx b$  Approx. der RHS,  $Ax = b, A\tilde{x} = \tilde{b}$ , wobei  $\tilde{x}$  ist berechnete Lsg.

**Bemerkung 54.** In der Numerik werden nie Inverse berechnet, sondern immer Gleichungen gelöst, z.B.  $y = BA^{-1}z$

$\Rightarrow$  Zunächst lösen  $Ax = z$  (d.h.  $x = A^{-1}z$ ), dann berechnen  $y = Bx$ .

## 6.1 Dreiecksmatrizen

**Definition 14.**  $L \in \mathbb{K}^{n \times n}$  heißt **untere Dreiecksmatrix**, gdw.  $L_{jk} = 0$  für  $k > j$ .

$U \in \mathbb{K}^{n \times n}$  heißt **obere Dreiecksmatrix**, gdw.  $U_{jk} = 0$  für  $j > k$ .

**Bemerkung 55.** Falls  $A$  eine Dreiecksmatrix, so gilt  $\sigma(A) = \{A_{11}, \dots, A_{nn}\}$  und  $\det(A) = \prod_{j=1}^n A_{jj}$ . Insb.  $A$  regulär  $\iff \forall j = 1, \dots, n : A_{jj} \neq 0$ .

**Algorithmus 6** (Lösen einer oberen Dreiecksmatrix). *Input:*  $U \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n$

- for  $j = n : -1 : 1$
- $x_j = \left( b_j - \sum_{k=j+1}^n U_{jk} x_k \right) / U_{jj}$
- end

*Output:*  $x \in \mathbb{K}^n$

**Lemma 17.** Algorithmus ist wohldef. und berechnet in  $n^2$  arithmetischen Operationen die Lösung  $x \in \mathbb{K}^n$  von  $Ux = b$ .

*Proof.*

$$Ux = b \iff \sum_{k=1}^n \underbrace{U_{jk}}_{=0 \text{ für } j > k} x_k = b_j \forall j \iff \sum_{k=j}^n U_{jk} x_k = b_j \forall j$$

$$= \underbrace{U_{jj}}_{\neq 0} x_j + \sum_{k=j+1}^n U_{jk} x_k$$

Aufwand pro  $j$ :  $n - j$  Multiplikationen + Subtraktionen, 1 Division

$$\implies \text{Aufwand} = \sum_{j=1}^n \{2(n - j) + 1\} = n + 2 \underbrace{\sum_{k=0}^{n-1} k}_{= \frac{(n-1)n}{2}} = n^2$$

□

**Algorithmus 7** (Lösen einer unteren Dreiecksmatrix). *Input:*  $L \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n$

- for  $j = 1 : n$
- $x_j = \left( b_j - \sum_{k=1}^{j-1} L_{jk} x_k \right) / L_{jj}$
- end

*Output:*  $x \in \mathbb{K}^n$

**Lemma 18.** Algorithmus ist wohldef. und berechnet in  $n^2$  arithm. Op. die Lösung  $x \in \mathbb{K}^n$  von  $Lx = b$ .

*Proof.*

$$Lx = b \iff \sum_{k=1}^n \underbrace{L_{jk}}_{=0 \text{ für } k > j} x_k = b_j \forall j \iff \sum_{k=1}^j L_{jk} x_k = b_j \forall j$$

$$= \underbrace{L_{jj}}_{\neq 0} x_j + \sum_{k=1}^{j-1} L_{jk} x_k$$

□

**Lemma 19.** Sei  $\mathcal{U} := \{U \in \mathbb{K}^{n \times n} \text{ obere } \Delta\text{-Matrix}\}$

$\implies$

1.  $A, B \in \mathcal{U} \implies AB \in \mathcal{U}$
2.  $A \in \mathcal{U}$  regulär  $\implies B := A^{-1} \in \mathcal{U}$  und  $B_{jj} = A_{jj}^{-1} \forall j$

*Proof.* 1. Seien  $A, B \in \mathcal{U}, C := AB$

$$C_{jl} = \sum_{k=1}^n \underbrace{A_{jk}}_{=0 \text{ für } j > k} \underbrace{B_{kl}}_{=0 \text{ für } k > l} = \sum_{k=j}^l A_{jk} B_{kl}$$

$$\implies C_{jl} = 0 \text{ für } j > l \implies C \in \mathcal{U}.$$

2. Sei  $A \in \mathcal{U}$  regulär. Falls  $B := A^{-1} \in \mathcal{U}$  existiert, so folgt  $1 = (AB)_{jj} = A_{jj}B_{jj}$ , d.h.  $B_{jj} = A_{jj}^{-1}$ .

Sei  $b^{(l)} \in \mathbb{K}^n$   $l$ -te Spalte von  $B$ , d.h.  $Ab^{(l)} = e_l \rightsquigarrow$  Man kann  $b^{(l)}$  berechnen durch

- for  $j = n : -1 : 1$
- $b_j^{(l)} = \left( \delta_{jl} - \sum_{k=j+1}^n A_{jk}b_k^{(l)} \right) / A_{jj}$
- end

Beh:  $b_j^{(l)} = 0 \forall j > l$

Ind.anf.  $j = n > l$ :  $b_j^{(l)} = \delta_{jl}/A_{jj} = 0$

Ind.schritt: Aussage gelte bis  $j$ , zz:  $j-1 > l$

$$b_{j-1}^{(l)} = \left( \underbrace{\delta_{j-1,l}}_{=0} - \sum_{k=j}^n A_{jk} \underbrace{b_k^{(l)}}_{=0 \text{ nach Ind.}} \right) / A_{jj} = 0$$

$$\implies B_{jl} = b_j^{(l)} = 0 \forall j > l \implies B \in \mathcal{U}.$$

□

**Korollar 6.** Sei  $\mathcal{L} := \{L \in \mathbb{K}^{n \times n} \text{ untere } \triangle\text{-Matrix}\}$

$\implies$

$$1. A, B \in \mathcal{L} \implies AB \in \mathcal{L}$$

$$2. A \in \mathcal{L} \text{ regulär} \implies B := A^{-1} \in \mathcal{L} \text{ und } B_{jj} = A_{jj}^{-1} \forall j$$

*Proof.*  $\mathcal{L} = \{U^T | U \in \mathcal{U}\}$ .

□

## 6.2 LU-Zerlegung

Im ganzen Abschnitt sei  $A \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n$ .

**Definition 15.** Eine Faktorisierung  $A = LU$  mit  $L \in \mathcal{L}, U \in \mathcal{U}$  heißt **LU-Zerlegung von A**.

**Bemerkung 56.** Falls LU-Zerlegung existiert, so sind  $L, U$  regulär. Es braucht  $n$  zusätzliche Bedingungen, damit LU-Zerlegung eindeutig sein kann (denn  $A$  hat  $n^2$  Einträge und  $L, U$  jeweils  $\frac{n(n+1)}{2}$ )

Die Lösung von  $Ax = b$  erhält man durch

- Löse  $Ly = b$
- Löse  $Ux = y$

$$\implies Ax = L \underbrace{Ux}_{=y} = b \text{ in } 2n^2 \text{ arithm. Op.}$$

**Satz 26** (Existenz + Eindeutigkeit). Es sind äquivalent

1. Alle Untermatrizen  $A_l = (A_{jk})_{j,k=1}^l$  sind regulär
2. Ex. LU-Zerlegung  $A = LU$ .

In diesem Fall existiert eine eindeutige LU-Zerlegung mit  $L_{jj} = 1 \forall j$ , sog. **normalisierte LU-Zerlegung**.

*Proof.*

$$A = \begin{pmatrix} A_l & * \\ * & * \end{pmatrix} = \begin{pmatrix} L_l & 0 \\ * & * \end{pmatrix} \begin{pmatrix} U_l & * \\ 0 & * \end{pmatrix} = LU$$

(ii  $\implies$  i) Es gilt  $A = LU$  und  $L, U$  regulär  $\implies L_l, U_l$  regulär  $\implies A_l$  regulär.

(i  $\implies$  ii) zz: Eindeutige Existenz der norm. LU-Zerl. durch Induktion nach  $n$ .

Ind.anf.  $n = 1$  ✓

Ind.schritt: Aussage gelte für  $n-1$ , d.h.  $A_{n-1}$

Wir machen einen Ansatz

$$A = \begin{pmatrix} A_{n-1} & \beta \\ \alpha^T & a_{nn} \end{pmatrix} \stackrel{!}{=} \begin{pmatrix} L_{n-1} & 0 \\ \gamma^T & 1 \end{pmatrix} \begin{pmatrix} U_{n-1} & \delta \\ 0 & u_{nn} \end{pmatrix}$$

zz:  $\gamma, \delta \in \mathbb{K}^{n-1}, u_n n \in \mathbb{K}$  existiert eindeutig.

$L_{n-1} \delta = \beta$  hat eind. Lsg.

$\gamma^T U_{n-1} = \alpha^T \iff U_{n-1}^T \gamma = \alpha$  hat eind. Lsg.

$\gamma^T \delta + u_{nn} = a_{nn} \iff u_{nn} = a_{nn} - \gamma^T \delta$  hat eind. Lsg. □

**Beispiel 24.** Die Matrix  $\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$  hat keine LU-Zerlegung.

**Beispiel 25.** Positiv definite Matrizen (d.h.  $\langle Ax, x \rangle > 0 \forall x \in \mathbb{K}^n \setminus \{0\}$ ) und strikt diagonaldominante Matrizen haben eine LU-Zerlegung.

Erinnerung: Beim Lösen von  $Ax = b$  gilt  $\frac{\|x - \tilde{x}\|}{\|\tilde{x}\|} \leq \text{cond}(A) \frac{\|b - \tilde{b}\|}{\|\tilde{b}\|}$ , und diese Abschätzung ist scharf. Wenn man mittels LU-Zerlegung löst, gilt also  $\frac{\|x - \tilde{x}\|}{\|\tilde{x}\|} \leq \text{cond}(U) \frac{\|y - \tilde{y}\|}{\|\tilde{y}\|} \leq \text{cond}(U) \text{cond}(L) \frac{\|b - \tilde{b}\|}{\|\tilde{b}\|}$ .

**Bemerkung 57.** Das GLS  $Ax = b$  mittels LU-Zerlegung ist also instabil, falls  $\text{cond}(U) \text{cond}(L) \gg \text{cond}(A)$ . Dies ist der sog. **Standardfehler der Numerik**: Man zerlegt ein Problem  $\Phi = \Phi_1 \circ \Phi_2$  in Teilprobleme, sodass eines der Teilprobleme  $\Phi_j$  wesentlich schlechter konditioniert ist als  $\Phi$ .

**Beispiel 26.**  $A = \begin{pmatrix} \epsilon & 1 \\ 1 & 0 \end{pmatrix}$  mit  $\epsilon$  klein.  $\implies A^{-1} = \begin{pmatrix} 0 & 1 \\ 1 & -\epsilon \end{pmatrix}$  mit  $\|A\|_\infty = 1 + \epsilon = \|A^{-1}\|_\infty, \text{cond}_\infty(A) = (1 + \epsilon)^2$  und

$$A = \underbrace{\begin{pmatrix} 1 & 0 \\ \frac{1}{\epsilon} & 1 \end{pmatrix}}_{=L} \underbrace{\begin{pmatrix} \epsilon & 1 \\ 0 & -\frac{1}{\epsilon} \end{pmatrix}}_{=U} \text{ mit } L^{-1} = \begin{pmatrix} 1 & 0 \\ -\frac{1}{\epsilon} & 1 \end{pmatrix}, U^{-1} = \begin{pmatrix} \frac{1}{\epsilon} & 1 \\ 0 & -\epsilon \end{pmatrix}$$

$$\implies \text{cond}_\infty(L) = \left(\frac{1}{\epsilon} + 1\right)^2, \text{cond}_\infty(U) = \frac{1}{\epsilon} \left(1 + \frac{1}{\epsilon}\right)$$

**Bemerkung 58.** Falls  $A \in \mathbb{K}^{n \times n}$  SPD ist, so existiert eine spezielle LU-Zerlegung mit  $A = LL^H$  mit eindeutigen  $L \in \mathcal{L}$  mit  $L_{jj} > 0 \forall j$ , sog. **Cholesky-Zerlegung**. Es gilt  $\text{cond}_2(L) = \text{cond}_2(L^H) = \sqrt{\text{cond}_2(A)}$ , d.h. Cholesky-Zerlegung ist stabiles Verfahren zur Lösung von  $Ax = b$ .

**Algorithmus 8** (Crout). Input:  $A \in \mathbb{K}^{n \times n}$  mit LU-Zerlegung

- for  $i = 1 : n$
- for  $k = i : n$
- $U_{ik} = A_{ik} - \sum_{j=1}^{i-1} L_{ij} L_{jk}$
- end
- for  $k = i + 1 : n$
- $L_{ki} = \left(A_{ki} - \sum_{j=1}^{i-1} L_{kj} U_{ji}\right) / U_{ii}$
- end
- end

**Satz 27.** Der Crout-Algorithmus ist wohldef. und berechnet in  $\frac{2}{3}n^3 + \mathcal{O}(n^2)$  arithm. Op. die nicht-trivialen Einträge der normalisierten LU-Zerlegung. Man kann  $A_{ik}$  durch  $U_{ik}$  und  $A_{ki}$  durch  $L_{ki}$  überschreiben, d.h. es wird kein zusätzlicher Speicher benötigt.

Idee des Algorithmus = Parkettierung von  $A$ .

$$\text{Proof. Für } i \leq k \implies A_{ik} = \sum_{j=1}^n \underbrace{L_{ij}}_{=0 \text{ für } j > i} U_{jk} = \sum_{j=1}^i L_{ij} U_{jk} = \underbrace{L_{ii}}_{=1} U_{ik} + \sum_{j=1}^{i-1} L_{ij} U_{jk}$$

$$\text{Für } i > k \implies A_{ki} = \sum_{j=1}^n L_{kj} \underbrace{U_{ji}}_{=0 \text{ für } j > i} = \sum_{j=1}^i L_{kj} U_{ji} = L_{ki} U_{ii} + \sum_{j=1}^{i-1} L_{kj} U_{ji}$$

Aufwand für fixes  $i$ :  $(n - i + 1)(i - 1)2 + (n - i)[(i - 1)2 + 1]$



Gesamtaufwand:

$$\begin{aligned} \sum_{i=1}^n \{(n-(i-1))(i-1)2 + (n-(i-1))(i-1)2 - (i-1)2 + n-i\} &= 4 \sum_{j=0}^{n-1} (n-j)j - 2 \sum_{j=0}^{n-1} j + \sum_{j=0}^{n-1} j = \\ &= 4n \sum_{j=0}^{n-1} j - 4 \sum_{j=0}^{n-1} j^2 - \sum_{j=0}^{n-1} j = \underbrace{2n^3}_{=\frac{6}{3}n^3} - 4 \underbrace{\frac{2n^3}{6}}_{=\frac{4}{3}n^3} + \mathcal{O}(n^2) = \frac{2}{3}n^3 + \mathcal{O}(n^2) \\ &= \underbrace{\frac{(n-1)n}{2}}_{=\frac{(n-1)n(2n-1)}{6}} - \underbrace{\frac{n(n-1)}{2}}_{=\mathcal{O}(n^2)} = \mathcal{O}(n^2) \end{aligned}$$

□

**Bemerkung 59.** Mit Hilfe des Crout-Algorithmus kann man sehen, dass gewisse Struktur von  $A$  bei LU-Zerlegung erhalten bleibt, z.B. Bandstruktur oder Skyline-Struktur.

### 6.3 Gauss-Elimination

**Algorithmus 9.** Input:  $A^{(1)} = A \in \mathbb{K}^{n \times n}$  regulär,  $b^{(1)} := b \in \mathbb{K}^n$

1. Schritt: Erhalte  $A^{(2)} \in \mathbb{K}^{n \times n}$ , indem erste Zeile von  $A^{(1)}$  unverändert bleibt und in allen folgenden Zeilen der Eintrag  $A_{i1}^{(1)}$  eliminiert wird, d.h. def  $L_{i1} := \frac{A_{i1}^{(1)}}{A_{11}^{(1)}}$ ,  $A_{ij}^{(2)} = A_{ij}^{(1)} - L_{i1}A_{1j}^{(1)}$ ,  $b_i^{(2)} = b_i^{(1)} - L_{i1}b_1^{(1)}$

$$A^{(2)} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ 0 & a_{22}^{(2)} & \dots & a_{2n} \\ \vdots & \vdots & & \vdots \\ 0 & a_{n2}^{(2)} & \dots & a_{nn}^{(n)} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & & & \\ -L_{21} & 1 & & \\ \vdots & & \ddots & \\ -L_{n1} & & & 1 \end{pmatrix}}_{=:L^{(1)}} A^{(1)}$$

$k$ -ter Schritt:  $L_{ik} := \frac{A_{ik}^{(k)}}{A_{kk}^{(k)}}$ ,  $A_{ij}^{(k+1)} = A_{ij}^{(k)} - L_{ik}A_{kj}^{(k)}$ ,  $b_i^{(k+1)} = b_i^{(k)} - L_{ik}b_k^{(k)}$

$$\Rightarrow A^{(k+1)} = \begin{pmatrix} a_{11} & & & & * \\ & a_{22}^{(2)} & & & \\ & & \ddots & & \\ & & & a_{k+1,k+1}^{(k+1)} & \dots & a_{k+1,n}^{(k+1)} \\ & & & \vdots & & \vdots \\ 0 & & & a_{n,k+1}^{(k+1)} & \dots & a_{nn}^{(k+1)} \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -L_{k+1,k} & 1 & \\ & & \vdots & & \ddots \\ 0 & & -L_{n,k} & 0 & & 1 \end{pmatrix}}_{=:L^{(k)}} A^{(k)}$$

$$b^{(k+1)} = L^{(k)}b^{(k)}$$

Output nach  $(n-1)$  Schritten:  $U = A^{(n-1)}$  obere Dreiecksmatrix und Vektor  $y := y^{(n-1)}$ ,  $L$  normalisierte untere Dreiecksmatrix.

**Satz 28.** Das Gauss-Verfahren ist genau dann durchführbar, wenn  $A$  eine LU-Zerlegung hat. In diesem Fall ist  $A = LU$  die eindeutige normalisierte LU-Zerlegung und  $y = L^{-1}b$ . Man erhält also die Lösung von  $Ax = b$  durch Lösen von  $Ux = y$ .

*Proof.* 1.  $A$  besitze eine LU-Zerlegung.

zz:  $A_{kk}^{(k)} \neq 0 \forall k = 1, \dots, n-1$  (dann Gauss durchführbar)

Induktion nach  $k$ , klar:  $A_{11} = A_{11}^{(1)} \neq 0$  (da LU-Zerlegung existiert)  $A^{(k)} = \underbrace{L^{(k-1)} \dots L^{(1)}}_{\text{reguläre Dreiecksmatrix}} A^{(1)}$

$\Rightarrow$  Untermatrix  $\underbrace{A_k^{(k)}}_{\text{regulär}} = \underbrace{L_k^{(k-1)} \dots L_k^{(n)}}_{\text{regulär}} \underbrace{A_k^{(n)}}_{\text{regulär, da LU existiert}}$

$$A_k^{(k)} = \begin{pmatrix} A_{11} & \dots & \dots & A_{1n} \\ & A_{22}^{(2)} & \dots & A_{2n}^{(2)} \\ & & \ddots & \\ 0 & & & A_{kk}^{(k)} \end{pmatrix} \Rightarrow A_{kk}^{(k)} \neq 0$$

2. Das Gauss-Verfahren sei durchführbar, d.h.  $U = A^{(n)} = \underbrace{L^{(n-1)} \dots L^{(1)}}_{=: L^{-1}} A$  ist obere Dreiecksmatrix

klar:  $A = LU$ ,  $L$  untere Dreiecksmatrix

zz:  $L_{ik}$  aus Algorithmus bilden  $L$ .

Def.  $l_k = (0, \dots, 0, L_{k+1,k}, \dots, L_{nk}) \in \mathbb{K}^n, e_k \in \mathbb{K}^n$  Einheitsvektor

$$\begin{aligned} \Rightarrow L^{(k)} &= \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & -L_{k+1,k} & 1 & \\ & & \vdots & & \ddots \\ 0 & & -L_{k,k} & 0 & & 1 \end{pmatrix} = Id - l_k e_k^T \\ &\Rightarrow L^{(k)}(Id + l_k e_k^T) = Id - \underbrace{l_k e_k^T l_k}_{=0} e_k^T \\ &\Rightarrow (L^{(k)})^{-1} = Id + l_k e_k^T \end{aligned}$$

Ind. beh.  $(L^{(1)})^{-1} \dots (L^{(k)})^{-1} = Id + \sum_{j=1}^k l_j e_j^T$

Ind.anf.  $k = 1 \checkmark$

Ind.schritt:

$$\begin{aligned} (L^{(1)})^{-1} \dots (L^{(k+1)})^{-1} &= (Id + \sum_{j=1}^k l_j e_j^T)(Id + l_{k+1} e_{k+1}^T) = Id + \sum_{j=1}^{k+1} l_j e_j^T + \sum_{j=1}^k l_j \underbrace{e_j^T l_{k+1}}_{=0} e_{k+1}^T \\ \Rightarrow L &= (L^{(n-1)} \dots L^{(1)})^{-1} = (L^{(1)})^{-1} \dots (L^{(n-1)})^{-1} = Id + \sum_{j=1}^{n-1} l_j e_j^T = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ L_{ik} & & 1 \end{pmatrix} \end{aligned}$$

□

**Algorithmus 10** (Implementierung). *Input:*  $A \in \mathbb{K}^{n \times n}$  mit  $LU$ -Zerlegung,  $b \in \mathbb{K}^n$

- for  $k = 1 : n - 1$
- $i = k + 1 : n$
- $L_{ik} = \frac{A_{ik}^{(k)}}{A_{kk}^{(k)}}$
- $b_i^{k+1} = b_i - L_{ik} b_k^{(k)}$
- for  $j = k + 1 : n$
- $A_{ij}^{(k+1)} = A_{ij}^{(k)} - L_{ik} A_{kj}^{(k)}$
- end
- end
- end

*Output:* Nicht-triviale Einträge der normalisierten  $LU$ -Zerlegung sowie modifizierte rechte Seite  $b^{(n)} = L^{-1}b$ .

**Lemma 20.** Gauss-Algorithmus benötigt  $\frac{2}{3}n^3 + \mathcal{O}(n^2)$  arith. Operationen. Eine Implementierung darf die oben Indizes weglassen und  $A$  und  $b$  überschreiben (dann kein zusätzlicher Speicher nötig), d.h. nicht-triviale  $L_{ik}$  auf  $A_{ik}$  speichern.

**Algorithmus 11** (Gauss-Elimination mit Zeilenvertauschung). • Im  $k$ -ten Schritt bestimme Index  $p = p(k) \in \{k, \dots, n\}$  mit  $|A_{pk}^{(k)}| = \max_{i=k, \dots, n} |A_{ik}^{(k)}|$

- Vertausche Zeilen  $p$  und  $k$  in  $(A^{(k)}, b^{(k)})$  und erhalte  $(\tilde{A}^{(k)}, \tilde{b}^{(k)})$
- Führe Eliminationsschritt für  $(\tilde{A}^{(k)}, \tilde{b}^{(k)})$  und erhalte  $(A^{(k+1)}, b^{(k+1)})$

**Bemerkung 60.** Das Verfahren wird üblicherweise ohne Vertauschung im Speicher über einen Permutationsvektor realisiert, d.h. anfangs  $\pi = (1, \dots, n)$  und Vertauschen vertauscht nur  $\pi(k)$  und  $\pi(p)$ .

**Algorithmus 12.** Input:  $A \in \mathbb{K}^{n \times n}$  regulär,  $b \in \mathbb{K}^n$

Initialisierung:  $\pi := (1, \dots, n) \in \mathbb{N}^n$

- for  $k = 1 : n - 1$
- Finde  $p \in \{k, \dots, n\}$  mit  $|A_{pk}^{(k)}| = \max_{i=k, \dots, n} |A_{ik}^{(k)}|$
- Vertausche  $\pi(k)$  und  $\pi(p)$
- for  $i = k + 1 : n$
- $L_{\pi(i), k} := A_{\pi(i), k}^{(k)} / A_{\pi(k), k}^{(k)}$
- $b_{\pi(i)}^{(k+1)} := b_{\pi(i)}^{(k)} - L_{\pi(i), k} b_{\pi(k)}^{(k)}$
- for  $j = k + 1 : n$
- $A_{\pi(i), j}^{(k+1)} := A_{\pi(i), j}^{(k)} - L_{\pi(i), k} A_{\pi(k), j}^{(k)}$
- end
- end
- end

Output: Matrizen  $L, U \in \mathbb{K}^{n \times n}$  mit nicht-trivialen Einträgen  $U_{ij} := A_{\pi(i), j}^{(n-1)}$  für  $i \leq j$ ,  $L_{ij} := A_{\pi(i), j}^{(n-1)}$  für  $i > j$  modifiziere rechte Seite  $y_i := b_{\pi(i)} \forall i = 1, \dots, n$  und Permutationsvektor  $\pi$ .

**Bemerkung 61.** klar: Aufwand ist  $\frac{2}{3}n^3 + \mathcal{O}(n^2)$ , da Matrixnormsuche  $\mathcal{O}(n - k)$  in jedem Durchlauf der  $k$ -Schleife.

**Bemerkung 62.** Falls man  $A$  überschreibt, muss man zusätzlich das Lösen von  $Ux = y$  mit Permutationsvektor realisieren.

**Satz 29.** Für jedes reguläre  $A \in \mathbb{K}^{n \times n}$  und  $b \in \mathbb{K}^n$  ist die Gauss-Elimination mit Zeilenvertauschung durchführbar und berechnet die normalisierte LU-Zerlegung  $PA = LU$ , wobei  $P \in \{0, 1\}^{n \times n}$  die Matrix zum abschleifen Permutationsvektor ist, d.h.  $Pe_i = e_{\pi(i)} \forall i$ . Ferner gilt  $|L_{ij}| \leq 1 \forall i, j$  und  $b^{(n)} = L^{-1}Pb$ .

**Bemerkung 63.** Aus Kenntnis von  $PA = LU$  folgt

- Löse  $Ly = Pb$
- Löse  $Ux = y$

$\Rightarrow \underbrace{LU}_{PA} x = Pb \Rightarrow Ax = b$ . Das erste Lösen wird durch das Verfahren erledigt!

*Proof.* Im  $k$ -ten Schritt ist  $p^{(k)} = (e_1, e_2, \dots, e_{k-1}, e_p, e_{k+1}, \dots, e_{p-1}, e_k, e_{p+1}, \dots, e_n)$  die Matrix der Zeilenvertauschung.

klar:  $p^{(k)}$  regulär,  $p^{(k)} = (p^{(k)})^{-1}$   
 $\Rightarrow A^{(k+1)} = L^{(k)} p^{(k)} A^{(k)}, b^{(k+1)} = L^{(k)} p^{(k)} b^{(k)}$

1. zz: Wohldefiniertheit

Gauss mit Zeilenvertauschung ist nicht wohldef, falls ex.  $k \in \{1, \dots, n-1\}$  mit  $\max_{i=k, \dots, n} |A_{ik}^{(k)}| = 0$   
 $\Rightarrow$  ersten  $k$  Spalten von  $A^{(k)}$  linear abhängig.

Andererseits gilt  $\text{rang}(A^{(n)}) = \text{rang}(A^{(n-1)}) = \dots = \text{rang}(A)$

Widerspruch aus Existenz von minimalen  $k$  liefert wohldef.

2. zz:  $|L_{ij}| \leq 1 \forall i, j$  denn  $L_{ik} = \frac{A_{\pi(i), k}^{(k)}}{A_{\pi(k), k}^{(k)}}$  und  $|A_{\pi(k), k}^{(k)}| = \max_{i=k, \dots, n} |A_{ik}^{(k)}|$

3. zz:  $PA = LU$

$$\begin{aligned}
A^{(1)} &= A \\
A^{(2)} &= L^{(1)} P^{(1)} A^{(1)} \\
A^{(3)} &= L^{(2)} P^{(2)} A^{(2)} = L^{(2)} P^{(2)} L^{(1)} P^{(1)} A = L^{(2)} (P^{(2)} L^{(1)} \underbrace{P^{(2)}}_{=Id} P^{(1)}) A \\
A^{(4)} &= L^{(3)} P^{(3)} A^{(2)} = L^{(3)} P^{(3)} L^{(2)} (P^{(2)} L^{(1)} P^{(2)}) (P^{(2)} P^{(1)}) A = \\
&= L^{(3)} (P^{(3)} L^{(2)} P^{(2)}) (P^{(3)} P^{(2)} L^{(1)} P^{(2)} P^{(3)}) (P^{(3)} P^{(2)} P^{(1)}) A \\
\Rightarrow A^{(n)} &= \hat{L}^{(n-1)} \dots \hat{L}^{(1)} P A \text{ und } b^{(n)} = \hat{L}^{(n-1)} \dots L^{(1)} P b \text{ mit } P = P^{(n-1)} \dots P^{(1)} \\
\hat{L}^{(k)} &= P^{(n-1)} \dots P^{(k+1)} L^{(k)} P^{(k+1)} \dots P^{(n-1)} \text{ und } U := A^{(n)} \text{ obere Dreiecksmatrix.} \\
\text{Wh: } L^{(k)} &= Id - l_k e_k^T, l_k = \underbrace{(0, \dots, 0, *, \dots, *)}_{k \text{ viele}} \\
\text{klar: } P^{(j)} (Id - l_k e_k^T) P^{(j)} &= Id - \underbrace{P^{(j)} l_k}_{=\tilde{l}_k = (0, \dots, 0, *, \dots, *)} \underbrace{e_k^T P^{(j)}}_{=e_k^T \text{ f\"ur } j > k} \\
\Rightarrow \hat{L}^{(k)} &= Id - \hat{l}_k e_k^T \text{ mit } \hat{l}_k = \underbrace{(0, \dots, 0, *, \dots, *)}_{k \text{ viele}} \\
\Rightarrow L &= (\hat{L}^{(n+1)} \dots \hat{L}^{(n)})^{-1} \text{ normalisierte untere } \triangle\text{-Matrix} \Rightarrow U = L^{-1} P A \Rightarrow LU = P A.
\end{aligned}$$

□

## 6.4 QR-Zerlegung

**Definition 16.** Zu  $A \in \mathbb{K}^{m \times n}$  mit  $m \geq n$  heit  $A = QR$  mit  $Q \in \mathbb{K}^{m \times m}$  unitr und  $R \in \mathbb{K}^{m \times n}$  **verallgemeinerte obere Dreiecksmatrix** (d.h.  $R_{jk} = 0$  fr  $j > k$ ) **QR-Zerlegung von A**.

**Bemerkung 64.** In der Literatur wird manchmal auch (quivalent!)  $A = \tilde{Q} \tilde{R}$  mit  $\tilde{R} \in \mathbb{K}^{n \times n}$  obere Dreiecksmatrix und  $\tilde{Q} \in \mathbb{K}^{m \times n}$  mit orthogonalen Spalten als QR-Zerlegung definiert.

**Bemerkung 65.** Falls  $m = n$  und  $A$  regulr, so gilt fr die QR-Zerlegung  $A = QR$ , dass  $\text{cond}_2(A) = \text{cond}_2(R)$ , da  $\|Q\|_2 = 1 = \|Q^{-1}\|_2$ , und  $Ax = b$  ist quivalent zu  $Rx = Q^H b$ .  $\Rightarrow$  QR-Zerlegung gibt stabile Lsungsstrategie fr jede regulre Matrix.

**Lemma 21** (Householder-Transformation). 1. Zu  $w \in \mathbb{K}^n$  definiere  $W := ww^H \in \mathbb{K}^{n \times n}$

$$\Rightarrow W = W^H \text{ und } Wx = (w^H x)w \forall x \in \mathbb{K}^n$$

2. Ist  $\|w\|_2 = 1$ , so ist die **Householder-Transformation**  $H := Id - 2ww^H \in \mathbb{K}^{n \times n}$  unitr und  $H^H = H^{-1} = H$ .

*Proof.* (ii)  $H^2 = Id - 4ww^H + 4w \underbrace{w^H w}_{= \|w\|_2^2 = 1} w^H = Id$ . □

**Bemerkung 66.** Geometrisch sind die Householder-Transformationen Spiegelungen an der Ebene  $E = \{x \in \mathbb{K}^n : x^H w = 0\}$  mit Normalvektor  $w$ . Man kann fr  $x \in \mathbb{K}^n \setminus \text{span}\{e_1\}$  den Vektor  $w$  so whlen, dass  $Hx \in \text{span}\{e_1\}$  liegt.

**Lemma 22.** Sei  $x \in \mathbb{K}^n \setminus \text{span}\{e_1\}$  und  $\lambda \in \mathbb{K}$  mit  $|\lambda| = 1$  und  $x_1 = \lambda|x|_2$ .

Definiere  $w := \frac{x + \sigma e_1}{\|x + \sigma e_1\|_2}$  mit  $\sigma := \lambda|x|_2$

$$\Rightarrow \|w\|_2 = 1 \text{ und } Hx := (Id - 2ww^H)x = -\sigma e_1$$

*Proof.* klar:  $w$  ist wohldef. und  $\|w\|_2 = 1$ , da  $x \notin \text{span}\{e_1\}$

$$\begin{aligned}
\|x + \sigma e_1\|_2^2 &= \|x\|^2 + 2\text{Re} \underbrace{(\sigma e_1)^H x}_{=\bar{\sigma} x_1 = \bar{\lambda} x_1 \underset{=|x_1|}{\in \mathbb{R}}} + \underbrace{|\sigma|^2}_{=\|x\|_2^2} = 2(x + \sigma e_1)^H x \\
Hx &= x - \underbrace{2(w^H x)w}_{=\frac{(x + \sigma e_1)^H x}{\|x + \sigma e_1\|_2^2} (x + \sigma e_1)}_{=1} = -\sigma e_1
\end{aligned}$$

□

**Algorithmus 13** (Householder-Verfahren zur Berechnung der QR-Zerlegung). *Input:*  $A \in \mathbb{K}^{m \times n}$  mit  $m \geq n$   
1. Schritt: Falls erste Spalte  $A_{(1)}^{(0)}$  mit  $A^{(0)} := A$  in  $\text{span}\{e_1\}$  liegt, def.  $H := Id$ . Ansonsten wähle Householder-Trafo mit  $HA_{(1)}^{(0)} \in \text{span}\{e_1\}$ .

Def  $A^{(1)} := Q^{(1)}A$  mit  $Q^{(1)} := H$

2. Schritt: Betrachte  $B_2 \in \mathbb{K}^{(m-1) \times (n-1)}$  mit  $B_2 = \begin{pmatrix} A_{22}^{(1)} & \dots & A_{2n}^{(1)} \\ \vdots & & \vdots \\ A_{m2}^{(1)} & \dots & A_{mn}^{(1)} \end{pmatrix}$

Falls erste Spalte  $B_{2,(1)} \in \text{span}\{e_1\} \subseteq \mathbb{K}^{m-1}$ , wähle  $H := Id$ , andernfalls wähle Householder-Trafo  $H$  mit  $HB_{2,(1)} \in \text{span}\{e_1\}$ .

Def  $Q^{(2)} = \begin{pmatrix} 1 & 0 \\ 0 & H \end{pmatrix}$  unitär,  $A^{(2)} := Q^{(2)}A^{(1)}$

$k$ -ter Schritt: Betrachte  $B_k \in \mathbb{K}^{(m-(k-1)) \times (n-(k-1))}$  mit  $B_k = \begin{pmatrix} A_{kk}^{(k-1)} & \dots & A_{kn}^{(k-1)} \\ \vdots & & \vdots \\ A_{mk}^{(k-1)} & \dots & A_{mn}^{(k-1)} \end{pmatrix}$

Falls erste Spalte  $B_{k,(1)} \in \text{span}\{e_1\} \subseteq \mathbb{K}^{m-(k-1)}$ , wähle  $H := Id \in \mathbb{K}^{(k-1) \times (k-1)}$ , andernfalls wähle Householder-Trafo mit  $HB_{k,(1)} \in \text{span}\{e_1\}$ .

Def  $Q^{(k)} = \begin{pmatrix} Id_{k-1} & 0 \\ 0 & H \end{pmatrix}$  unitär,  $A^{(k)} := Q^{(k)}A^{(k-1)}$

Output: Nach  $n$  Schritten ist  $R := A^{(n)}$  eine verallgemeinerte obere Dreiecksmatrix und  $Q := Q^{(n)} \dots Q^{(1)}$  unitär.

**Satz 30.** Das Householder-Verfahren berechnet für  $A \in \mathbb{K}^{m \times n}$  mit  $m \geq n$  eine QR-Zerlegung. Falls  $m = n$ , so werden nur  $n - 1$  Schritte gebraucht.

Proof.  $R := A^{(n)} = Q^{(n)}A^{(n-1)} = Q^{(n)} \dots Q^{(1)}A$

$$\begin{aligned} Q^{(j)} &= \begin{pmatrix} Id & 0 \\ 0 & H \end{pmatrix} \implies Q^{(j)}Q^{(j)} = Id \\ \implies A &= \underbrace{(Q^{(n)} \dots Q^{(1)})^{-1}}_{=Q^{(1)}} R = QR \\ &= \underbrace{(Q^{(1)})^{-1}}_{=Q^{(1)}} \dots \underbrace{(Q^{(n)})^{-1}}_{=Q^{(n)}} \end{aligned}$$

□

**Bemerkung 67.** Bei der Implementierung müssen Matrix-Matrix-Multiplikationen berechnet werden.

$$A^{(n+1)} = \begin{pmatrix} Id_k & 0 \\ 0 & H_{m-k} \end{pmatrix} \begin{pmatrix} U & X \\ 0 & B_{k+1} \end{pmatrix} = \begin{pmatrix} U & X \\ 0 & H_{m-k}B_{k+1} \end{pmatrix}$$

mit  $U \in \mathbb{K}^{k \times k}$  obere  $\Delta$ -Matrix,  $X \in \mathbb{K}^{k \times (n-k)}$  i.a. voll besetzt  $B_{k+1} \in \mathbb{K}^{(m-k) \times (n-k)}$  i.a. voll besetzt.

Man darf  $H_{m-k}B_{k+1}$  nicht als Matrix-Matrix-Produkt realisieren. Stattdessen nutzt man die Struktur  $H_{m-k}B_{k+1} = (Id - 2ww^H)B_{k+1} = B_{k+1} - ww^T$  mit  $v := 2B_{k+1}^T \bar{w}$ , d.h. "quadratischer" Aufwand statt "kubischer" Aufwand pro Schritt. Die Realisierung erfordert zusätzlichen Speicher. In der Regel speichert man die Diagonalelemente in einem Zusatzvektor  $(A_{11}^{(1)}, A_{22}^{(2)}, \dots, A_{nn}^{(n)})$ . Dann kann man den Householder-Vektor im unteren Dreieck speichern.

Unter diesen Voraussetzungen gilt für  $m = n$  der Gesamtaufwand  $\frac{4}{3}n^3 + \mathcal{O}(n^2)$ .

$Ax = b \iff PAx = Pb$  mit  $P$  regulär und  $\underbrace{\text{cond}_2(PA)}_{\approx 1} \leq \text{cond}_2(A)$  und  $P$  "billig" und Matrix-Vektor-Mult.

$Py$  und  $PA$  SPD bzgl. geeignetem Skalarprodukt.

**Satz 31** (Eindeutigkeit von QR-Zerlegung). Sei  $m = n$  und  $A \in \mathbb{K}^{n \times n}$  regulär und  $\sigma \in \mathbb{K}^n$  mit  $|\sigma_j| = 1$   
 $\implies$  Ex. eind. Zerlegung  $A = QR$  mit  $Q$  unitär und  $R$  verallg. obere  $\Delta$ -Matrix auf  $R_{jj} = \sigma_j |R_{jj}| \forall j$

**Satz 32.** 1. Existenz: Wähle  $A = \tilde{Q}\tilde{R}$  Zerlegung.  $A = \underbrace{(\tilde{Q}D^{-1})}_{\text{unitär}} \underbrace{(D\tilde{R})}_{\text{obere verall. } \Delta\text{-Matrix}}$  mit  $D \in \mathbb{K}^{n \times n}$  diagonal

$$\text{mit } D_{jj} = \underbrace{\frac{|\tilde{R}_{jj}|}{\tilde{R}_{jj}}}_{\text{wohldef., da } \tilde{R} \text{ regulär}} \sigma_j, |D_{jj}| = 1$$

$$(D\tilde{R})_{jj} = D_{jj} = \tilde{R}_{jj} \stackrel{!}{=} \sigma_j |\tilde{R}_{jj}|$$

$\Rightarrow Q := \tilde{Q}D^{-1}, R := D\tilde{R}$  zeigt Existenz

2. Eindeutigkeit: Seien  $QR = A = \tilde{Q}\tilde{R}$  zwei  $QR$ -Zerlegungen mit  $R_{jj} = \sigma_j |R_{jj}|, \tilde{R}_{jj} = \sigma_j |\tilde{R}_{jj}|$

$$\Rightarrow D := \underbrace{Q^{-1}\tilde{Q}}_{\text{unitär}} = \underbrace{\tilde{R}R^{-1}}_{\text{obere } \Delta\text{-Matrix}}$$

zz:  $D$  diagonal per Induktion

1. Spalte  $|D_{11}| = 1$ , insb.  $D_{1k} = 0 \forall k > j$

$j$ -te Spalte: analog

$$\Rightarrow D_{jj} = \frac{R_{jj}}{\tilde{R}_{jj}} = \frac{\sigma_j |R_{jj}|}{\sigma_j |\tilde{R}_{jj}|} \text{ und } |D_{jj}| = 1 \Rightarrow R_{jj} = \tilde{R}_{jj}.$$

## 6.5 Lineare Ausgleichsprobleme

Gegeben  $A \in \mathbb{K}^{m \times n}, b \in \mathbb{K}^m$

Das **lineare Ausgleichsproblem (LAP)** sucht  $x \in \mathbb{K}^n$  mit  $\|Ax - b\|_2 = \min_{y \in \mathbb{K}^n} \|Ay - b\|_2$ .

**Beispiel 27.** Gegeben  $(a_j, b_j)$  für  $j = 1, \dots, m$  finde  $p(t) = \sum_{k=0}^n x_k t^k \in \mathbb{P}_n$  mit  $\sum_{j=1}^m |p(a_j) - b_j|^2 = \min_{q \in \mathbb{P}_n} \sum_{j=1}^m |q(a_j) - b_j|^2$  (z.B.  $n = 1$  entspricht Ausgleichsgerade).

$\Rightarrow A = (a_j^k)_{j=1, \dots, m, k=0, \dots, n} \in \mathbb{K}^{m \times (n+1)}$  "Vandermonde-Matrix"

$Ay = (q(a_j))_{j=1, \dots, m} \in \mathbb{K}^m$  mit  $q(t) := \sum_{k=0}^n y_k t^k$

$$\Rightarrow \min_{q \in \mathbb{P}_n} \sum_{j=1}^m |q(a_j) - b_j|^2 = \min_{y \in \mathbb{K}^{n+1}} \|Ay - b\|_2^2$$

**Satz 33.** 1. Für beliebige  $m, n \in \mathbb{N}, A \in \mathbb{K}^{m \times n}, b \in \mathbb{K}^m$  hat LAP eine Lösung  $x \in \mathbb{K}^n$ .

2.  $x \in \mathbb{K}^n$  löst LAP  $\iff x$  löst die **Gauss'sche Normalgleichung**  $A^H A x = A^H b$

3. Für  $m \geq n = \text{rang}(A)$  hat LAP eine eindeutige Lsg.

*Proof.* 1.  $\text{Bild}(A)^\perp = \text{Kern}(A^H)$

$$\begin{aligned} \text{Bild}(A)^\perp &= \{y \in \mathbb{K}^m : \forall v \in \text{Bild}(A) : v^H y = 0\} = \\ &= \{y \in \mathbb{K}^m : \forall x \in \mathbb{K}^n : \underbrace{(Ax)^H y}_{=x^H A^H y} = 0\} = \{y \in \mathbb{K}^m : A^H y = 0\} = \text{Kern}(A^H) \end{aligned}$$

2.  $\mathbb{K}^m = \text{Bild}(A) + \text{Bild}(A)^\perp$ , d.h.  $b = v + w$  mit eind.  $v \in \text{Bild}(A), w \in \text{Bild}(A)^\perp$ , insb.  $v = Ax$  mit  $x \in \mathbb{K}^n$ .

$$\Rightarrow A^H b = A^H Ax + \underbrace{A^H w}_{=0} \Rightarrow \text{GNG hat mind. eine Lsg.}$$

3. zz:  $x$  löst LAP  $\Rightarrow x$  löst GNG

$$\|Ax - b\|_2^2 = \|\underbrace{Ax - v}_{\in \text{Bild}(A)}\|_2^2 + \|\underbrace{w}_{\in \text{Bild}(A)^\perp}\|_2^2. \text{ Falls } x \text{ LAP löst, muss } Ax = v \Rightarrow A^H Ax = A^H b \text{ wie oben.}$$

4. zz:  $x$  löst GNG  $\Rightarrow x$  löst LAP.

Sei  $y \in \mathbb{K}^n$

$$\|Ay - b\|_2^2 = \|\underbrace{Ay - Ax}_{\in \text{Bild}(A)}\|_2^2 + \|\underbrace{Ax - b}_{\text{löst GNG, } \in \text{Kern}(A^H) = \text{Bild}(A)^\perp}\|_2^2 \geq \|Ax - b\|_2^2$$

5.  $m \geq n = \text{rang}(A) \Rightarrow A^H A \in \mathbb{K}^{n \times n}$  SPD

$$< A^H Ax, x >_2 = < Ax, Ax >_2 = \|Ax\|_2^2 > 0 \forall x \neq 0.$$

□

**Bemerkung 68.** Sei  $A = QR \in \mathbb{K}^{m \times n}$  mit  $m \geq n = \text{rang}(A)$ . Partitioniere  $R = \begin{pmatrix} \tilde{R} \\ 0 \end{pmatrix}$  mit  $\tilde{R} \in \mathbb{K}^{n \times n}$  obere

$\Delta$ -Matrix,  $b = \begin{pmatrix} v \\ v \end{pmatrix}$  mit  $v \in \mathbb{K}^n$

$\Rightarrow \|Ax - b\|_2^2 = \|QRx - b\|_2^2 = \|Rx - Q^H b\|_2^2 = \|\tilde{R}x - v\|_2^2 + \|r\|_2^2$ , d.h. erhalte Lsg. von LAP durch Lösen von  $\tilde{R}x = v$ .

**Bemerkung 69.** Im selben Fall könnte man Cholesky verwenden von  $A^H A$ , aber  $\text{cond}_2(A^H A) = \text{cond}_2(A)^2$  (für  $m = n$ ), aber  $\text{cond}_2(A) = \text{cond}_2(\tilde{R})$

D.h. Cholesky wäre ggf. keine stabile Strategie zur Lsg des LAP.

## 7 Eigenwertprobleme

### 7.1 Lineare Algebra + Stabilität

**Satz 34** (Jordan-Form). Zu  $A \in \mathbb{K}^{n \times n}$  ex.  $X \in \mathbb{C}^{n \times n}$  regulär mit  $J := X^{-1}AX = \begin{pmatrix} J_1 & & \\ & \ddots & \\ & & J_p \end{pmatrix}$  blockdiag-

onal und  $J_i = \begin{pmatrix} \lambda_i & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i \end{pmatrix}$ , wobei  $\lambda_i$  die EW von  $A$ . Die  $J_i$  heißen **Jordan-Blöcke**.

**Bemerkung 70.** Die **algebraische Vielfachheit** eines EW  $\lambda_i$  ist die Vielfachheit der Nullstelle im char. Polynom (= die Summe der Dimensionen der Jordan-Blöcke zu  $\lambda_i$ ). Die **geometrische Vielfachheit** ist die Dimension des Eigenraums zu  $\lambda_i$  (= Anzahl der Jordan-Blöcke zu  $\lambda_i$ )

**Korollar 7.** Mit der Jordan-Form  $J = X^{-1}AX$  mit Jordan-Blöcken  $J_i$  definiere  $\tilde{A} := X\tilde{J}X^{-1}$  mit  $\tilde{J} = \begin{pmatrix} \tilde{J}_1 & & \\ & \ddots & \\ & & \tilde{J}_p \end{pmatrix}$  und  $\tilde{J}_i := \begin{pmatrix} \lambda_i + \epsilon_{i_1} & 1 & & \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ & & & \lambda_i + \epsilon_{i_m} \end{pmatrix}$  mit  $\epsilon_{i_j} > 0$

$\implies$  Falls alle  $\lambda_i + \epsilon_{i_j}$  unterschiedlich, so ist  $\tilde{A}$  diagonalisierbar und  $\|A - \tilde{A}\| \leq C \max_{i,j} |\epsilon_{i_j}|$ , wobei  $C$  nur von  $\|\cdot\|$  abhängt und von  $X$ .

**Bemerkung 71.** Korollar zeigt, dass die diagonalisierbare Matrizen dicht sind im Raum aller Matrizen. Insbesondere kann man also die Jordan-Form numerisch nicht berechnen.

**Satz 35** (Schur-Zerlegung). Sei  $A \in \mathbb{K}^{n \times n}$ . Für  $\mathbb{K} = \mathbb{R}$  gelte ferner  $\sigma(A) := \{\lambda \in \mathbb{C} \mid \exists v \in \mathbb{C}^n \setminus \{0\}, Av = \lambda v\} \subseteq \mathbb{R}$ .

$\implies$  Ex.  $Q \in \mathbb{K}^{n \times n}$  unitär/orthogonal mit  $R := Q^H A Q$  rechte obere  $\Delta$ -Matrix, insb.  $\sigma(R) = \sigma(A)$ .

**Korollar 8** (Spektralzerlegung). Sei  $A \in \mathbb{K}^{n \times n}$  mit  $A = A^H$ .  $\implies \sigma(A) \subseteq \mathbb{R}$  und ex.  $Q \in \mathbb{K}^{n \times n}$  unitär/orthogonal mit  $D := Q^H A Q$  diagonal.

*Proof.* Für  $Av = \lambda v$  folgt  $\lambda \|v\|_2^2 = \lambda v^H v = v^H Av = (Av)^H v = \bar{\lambda} \|v\|_2^2$

$\implies R = Q^H A Q, R^H = Q^H A^H Q = Q^H A Q = R$  diagonal. □

**Satz 36** (Bauer-Fike). Sei  $\|\cdot\|$  induzierte Operatornorm mit  $\|D\| := \sup_{\tilde{x} \in \tilde{K}^n \setminus \{0\}} \frac{\|D\tilde{x}\|}{\|\tilde{x}\|} = \max_j |D_{jj}|$  für alle  $D \in \mathbb{K}^{n \times n}$  Diagonalmatrix, z.B.  $\|\cdot\|_p$ . Sei  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $T \in \mathbb{K}^{n \times n}$  regulär mit  $D := T^{-1}AT$  diagonal,  $\tilde{A} \in \mathbb{K}^{n \times n}$

$\implies \forall \tilde{\lambda} \in \sigma(\tilde{A}) : \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \underbrace{\text{cond}(T)}_{= \|T\| \|T^{-1}\|} \|A - \tilde{A}\|$

**Korollar 9.** Sei  $A = A^H \in \mathbb{K}^{n \times n}, \tilde{A} \in \mathbb{K}^{n \times n}$

$\implies \forall \tilde{\lambda} \in \sigma(\tilde{A}) : \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \|A - \tilde{A}\|_2$

*Bauer-Fike.* Sei  $\tilde{\lambda} \in \sigma(\tilde{A})$ , o.B.d.A.  $\tilde{\lambda} \notin \sigma(A)$ .

Sei  $v \in \mathbb{K}^n \setminus \{0\}$  mit  $\tilde{A}v = \tilde{\lambda}v$

$$\implies 0 = (\tilde{A} - \tilde{\lambda})v = \underbrace{(A - \tilde{\lambda})}_{\text{regulär}} v + (\tilde{A} - A)v$$

$$\implies v = -(A - \tilde{\lambda})^{-1}(\tilde{A} - A)v$$

$$\implies 1 = \frac{\|(A - \tilde{\lambda})^{-1}(\tilde{A} - A)v\|}{\|v\|} \leq \|(A - \tilde{\lambda})^{-1}\| \|\tilde{A} - A\|$$

$$(A - \tilde{\lambda})^{-1} = (TDT^{-1} - \tilde{\lambda})^{-1} = (T \underbrace{(D - \tilde{\lambda})}_{\text{regulär}} T^{-1})^{-1} = T(D - \tilde{\lambda})^{-1}T^{-1}$$

$$\implies \|(A - \tilde{\lambda})^{-1}\| \leq \underbrace{\|T\| \|T^{-1}\|}_{= \text{cond}(T)} \underbrace{\|(D - \tilde{\lambda})^{-1}\|}_{\text{diagonal}} = \max_j \frac{1}{|D_{jj} - \tilde{\lambda}|} = \frac{1}{\min_j |D_{jj} - \tilde{\lambda}|} = \frac{1}{\min_{\lambda \in \sigma(D) = \sigma(A)} |\lambda - \tilde{\lambda}|}$$

□

## 7.2 Vektoriteration

Ziel: Entwickle iterative Verfahren, basierend auf Matrix-Vektor-Multiplikationen, die ein EW-EV-Paar berechnen.

**Lemma 23.** Für  $X, Y \leq \mathbb{K}^n$  Unterräume definiere  $d(X, Y) := \begin{cases} 1 & \text{für } \dim X \neq \dim Y \\ \|\mathbb{P}_X - \mathbb{P}_Y\|_2 & \text{für } \dim X = \dim Y \end{cases}$  wobei

$\mathbb{P}_Z : \mathbb{K}^n \rightarrow Z$  Orth.proj. auf  $Z \leq \mathbb{K}^n$  bzgl.  $\|\cdot\|_2$

$\implies d(\cdot, \cdot)$  ist eine Metrik auf den Unterräumen von  $\mathbb{K}^n$  mit  $d(X, Y) = d(X^\perp, Y^\perp) \forall X, Y \leq \mathbb{K}^n$ .

*Proof.* • Definitheit:  $\forall X, Y \leq \mathbb{K}^n : d(X, Y) = 0 \implies X = Y$

Seien  $X, Y \leq \mathbb{K}^n$  mit  $\mathbb{P}_X = \mathbb{P}_Y$ . Für  $x \in X$  gilt  $x = \mathbb{P}_X x = \mathbb{P}_Y x \in Y$ , d.h.  $X \subseteq Y$  analog  $Y \subseteq X$ .

• Symmetrie:  $\forall X, Y \leq \mathbb{K}^n : d(X, Y) = d(Y, X)$

• Dreiecksungleichung:  $\forall X, Y, Z \leq \mathbb{K}^n : d(X, Y) \leq d(X, Z) + d(Z, Y)$

1. Fall  $\dim X \neq \dim Y : d(X, Y) = 1 \leq d(X, Z) + d(Z, Y)$

2. Fall  $\dim X = \dim Y \neq \dim Z : d(X, Y) \leq \|\mathbb{P}_X\|_2 + \|\mathbb{P}_Y\|_2, \|z\|_2^2 = \|\mathbb{P}_X z\|_2^2 + \|(1 - \mathbb{P}_X)z\|_2^2, \mathbb{P}_X z = \sum_{j=1}^m (v_j^H z) v_j, m = \dim X$ .

3. Fall  $\dim X = \dim Y = \dim Z \implies d(X, Y) \leq \|\mathbb{P}_X - \mathbb{P}_Z\|_2 + \|\mathbb{P}_Z - \mathbb{P}_Y\|_2 = d(X, Z) + d(Z, Y)$ .

4. Fall  $\dim X^\perp = n - \dim X$ , d.h.  $\dim X = \dim Y \iff \dim X^\perp = \dim Y^\perp$

O.B.d.A  $\dim X = \dim Y$

$$d(X, Y) = \|\underbrace{(1 - \mathbb{P}_X)}_{=\mathbb{P}_{X^\perp}} - \underbrace{(1 - \mathbb{P}_Y)}_{=\mathbb{P}_{Y^\perp}}\|_2 = d(X^\perp, Y^\perp).$$

□

**Lemma 24** (Rayleigh-Quotient).  $A \in \mathbb{K}^{n \times n}, x \in \mathbb{K}^n \setminus \{0\}$

•  $x$  EV zu EW  $\lambda \in \mathbb{K} \implies \lambda = \frac{x^H A x}{\|x\|_2^2}$   
sog. *Rayleigh-Quotient*

• Sei  $\lambda \in \mathbb{K}$  mit  $\|Ax - \lambda x\|_2 = \min_{\mu \in \mathbb{K}} \|Ax - \mu x\|_2$   
 $\implies \lambda = \frac{x^H A x}{\|x\|_2^2}$

*Proof.*  $\min_{\mu \in \mathbb{K}} \|Ax - \mu x\|_2 = \min_{y \in \text{span}\{x\}} \|Ax - y\|_2$

$\implies$  Minimum wird eindeutig in  $y = \underbrace{\mathbb{P}_X(Ax)}_{=\frac{x^H A x}{\|x\|_2^2} x}$  angenommen mit  $X = \text{span}\{x\} \implies \lambda = \frac{x^H A x}{\|x\|_2^2}$

□

**Lemma 25** (Residuum als Fehlerkontrolle).  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $x \in \mathbb{K}^n$  mit  $\|x\|_2 = 1, \tilde{\lambda} \in \mathbb{K}, r = r(\tilde{\lambda}, x) := Ax - \tilde{\lambda}x$  **Residuum**. Dann gilt

1.  $\min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \text{cond}_2(T) \|r\|_2$ , sofern  $T \in \mathbb{K}^{n \times n}$  regulär mit  $D := T^{-1} A T$  diagonal.

2.  $A = A^H \implies \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \|r\|_2$

3.  $\tilde{\lambda} := x^H A x, A = A^H, \lambda \in \sigma(A)$  mit  $|\lambda - \tilde{\lambda}| = \min_{\lambda' \in \sigma(A)} |\lambda' - \tilde{\lambda}|$   
 $\implies |\lambda - \tilde{\lambda}| \leq C \|r\|_2^2, C := \frac{2}{\min_{\lambda' \in \sigma(A) \setminus \{\tilde{\lambda}\}} |\lambda' - \tilde{\lambda}|}$

*Proof.* 1.  $\tilde{A} := A - r x^H, \implies \tilde{A} x = Ax - r \underbrace{x^H x}_{=\|x\|_2^2=1} = \tilde{\lambda} x$

$$\text{Bauer-Fike} \implies \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \text{cond}_2(T) \underbrace{\|A - \tilde{A}\|_2}_{=\sup_{y \in \mathbb{K} \setminus \{0\}} \frac{\|r x^H y\|_2}{\|y\|_2} = \|r\|_2}$$

2. Wähle  $T$  orthogonal/unitär ✓

3. Beweis in 2 Schritten:



(a) Sei  $\tilde{\lambda} \in (\alpha, \beta) \subseteq \mathbb{R}$  und  $\sigma(A) \cap (\alpha, \beta) = \emptyset$

zz:  $0 < (\beta - \tilde{\lambda})(\tilde{\lambda} - \alpha) \leq \|r\|_2^2$

Sei  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  ONB aus EV zu  $A$ ,  $Av_j = \lambda_j v_j$ .

Sei  $x = \sum_j \mu_j v_j$  mit geeigneten  $\mu_j \in \mathbb{K}$

$$\begin{aligned} \implies \langle (A - \alpha)x, (A - \beta)x \rangle_2 &= \sum_{j,k} \bar{\mu}_j (\lambda_j - \alpha) \mu_k (\lambda_k - \beta) \underbrace{\langle v_j, v_k \rangle_2}_{\delta_{jk}} = \sum_j |\mu_j|^2 \underbrace{(\lambda_j - \alpha)(\lambda_j - \beta)}_{\geq 0} \geq 0 \\ &< (A - \alpha)x, (A - \beta)x \rangle_2 = \langle \underbrace{(A - \tilde{\lambda})x}_{=r} + (\tilde{\lambda} - \alpha)x, \underbrace{(A - \tilde{\lambda})x}_r + (\tilde{\lambda} - \beta)x \rangle_2 = \\ &\|r\|_2^2 + \tilde{\lambda} - \alpha \underbrace{\langle x, r \rangle_2}_{\substack{x^H(Ax - (x^H Ax)x) = x^H Ax - \underbrace{x^H x}_{=1} x^H Ax = 0}} + (\tilde{\lambda} - \beta) \langle r, x \rangle_2 + (\tilde{\lambda} - \alpha)(\tilde{\lambda} - \beta) \|x\|_2^2 \end{aligned}$$

(b) Sei  $\tilde{\lambda} \in (a, b) \subseteq \mathbb{R}$  mit  $\sigma(A) \cap (a, b) = \{\lambda\}$  zz: Behauptung.

1. Fall  $\lambda = \tilde{\lambda}$  ✓

2. Fall  $a < \tilde{\lambda} < \lambda$

Wähle  $\alpha = a, \beta = \lambda$

$$\implies (\lambda - \tilde{\lambda})(\tilde{\lambda} - a) \leq \|r\|_2^2 \implies |\lambda - \tilde{\lambda}| = \lambda - \tilde{\lambda} \leq \frac{1}{\lambda - a} \|r\|_2^2$$

Durch Wahl von  $a, \tilde{\lambda} - a \geq \frac{1}{2} \min_{\lambda' \in \sigma(A), \lambda' \neq \lambda} |\lambda' - \lambda|$

(c)  $\lambda < \tilde{\lambda} < b$

Wähle  $\alpha = \lambda, \beta = b$

$$\implies (b - \tilde{\lambda})(\tilde{\lambda} - \lambda) \leq \|r\|_2^2 \implies |\tilde{\lambda} - \lambda| = \tilde{\lambda} - \lambda \leq \frac{1}{b - \tilde{\lambda}} \|r\|_2^2 \rightsquigarrow \text{analog zu zuvor.}$$

□

**Algorithmus 14** (Power-Iteration). *Input:*  $A \in \mathbb{K}^{n \times n}, x_0 \in \mathbb{K}^n \setminus \{0\}$

- Für  $l = 0, 1, 2, \dots$  (solange wie  $\|Ax_l - \mu_l x_l\|_2$  "zu groß")
- $y_{l+1} := Ax_l$
- $x_{l+1} := \frac{y_{l+1}}{\|y_{l+1}\|_2}$  % approximativer EV
- $\mu_{l+1} := x_{l+1}^H Ax_{l+1}$  % approximativer EW

**Satz 37** (Konvergenz der Power-It.). Sei  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  Basis aus EV,  $Av_j = \lambda_j v_j$  mit  $|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots$

Sei  $x_0 = \sum_j \alpha_j v_j$  mit  $\alpha_1 \neq 0$

$\implies$  Power-Iteration ist wohldefiniert und ex.  $l_0 \in \mathbb{N}$  und  $C > 0$  mit

$$|\mu_l - \lambda_1| \leq C \left| \frac{\lambda_2}{\lambda_1} \right|^l, d(\text{span}\{x_l\}, \text{span}\{v_1\}) \leq C \left| \frac{\lambda_2}{\lambda_1} \right|^l \quad \forall l \geq l_0$$

Falls  $A = A^H$ , so gilt  $l_0 = 0$  und  $|\mu_2 - \lambda_1| \leq C \left| \frac{\lambda_2}{\lambda_1} \right|^{2l}$ .

*Proof.* 1.  $x_0 = \sum_j \alpha_j v_j \implies A^l x_0 = \sum_j \alpha_j \lambda_j^l v_j$  insb. Power-It. ist wohldef, da  $A^l x_0 \neq 0$ .

$$\implies x_l = c_l A^l x_0 \text{ mit } c_l \neq 0$$

$$\begin{aligned} \implies x_l &= c_l \sum_j \alpha_j \lambda_j^l v_j = c_l \alpha_1 \lambda_1^l \left( v_1 + \underbrace{\sum_{j \geq 2} \frac{\alpha_j}{\alpha_1} \left( \frac{\lambda_j}{\lambda_1} \right)^l v_j}_{=: \epsilon_l \in \mathbb{K}^n} \right) \\ \implies \|\epsilon_l\|_2 &\leq \sum_{j \geq 2} \left| \frac{\alpha_j}{\alpha_1} \right| \underbrace{\left| \frac{\lambda_j}{\lambda_1} \right|^l}_{\leq \left| \frac{\lambda_2}{\lambda_1} \right|^l} \|v_j\|_2 =: C \left| \frac{\lambda_2}{\lambda_1} \right|^l \end{aligned}$$

2.

$$\begin{aligned}
d(\underbrace{\text{span}\{x_l\}}_{=:X}, \underbrace{\text{span}\{v_1\}}_{=:Y}) &= \|\mathbb{P}_X - \mathbb{P}_Y\|_2 = \max_{z \in \mathbb{K}^n \setminus \{0\}} \frac{\| \frac{(x_l^H z)}{\|x_l\|_2^2} x_l - \frac{v_1^H z}{\|v_1\|_2^2} v_1 \|_2}{\|z\|_2} \\
\frac{x_l^H z}{\|x_l\|_2^2} x_l &= \frac{(v_1 + \epsilon_l)^H z}{\|v_1 + \epsilon_l\|_2^2} (v_1 + \epsilon_l) \\
\frac{x_l^H z}{\|x_l\|_2^2} x_l - \frac{v_1^H z}{\|v_1\|_2^2} v_1 &= \frac{(v_1 + \epsilon_l)^H z}{\|v_1 + \epsilon_l\|_2^2} (v_1 + \epsilon_l) - \frac{(v_1 + \epsilon_l)^H z}{\|v_1\|_2^2} (v_1 + \epsilon_l) + \frac{(v_1 + \epsilon_l)^H z}{\|v_1\|_2^2} (v_1 + \epsilon_l) - \frac{v_1^H z}{\|v_1\|_2^2} v_1 \leq \\
&\quad \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2 (v_1 + \epsilon_l)^H \|z\|_2}{\|v_1\|_2^2 \|v_1 + \epsilon_l\|_2^2} (v_1 + \epsilon_l) + \underbrace{\frac{\|(v_1 + \epsilon_l)^H z (v_1 + \epsilon_l) - v_1^H z v_1\|_2}{\|v_1\|_2^2}}_{= \frac{\|\epsilon_l^H z (v_1 + \epsilon_l) + v_1^H z \epsilon_l\|_2}{\|v_1\|_2^2}} \\
d(X, Y) &\leq \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2}{\|v_1\|_2^2} + \underbrace{\frac{\|\epsilon_l\|_2 \|v_1 + \epsilon_l\|_2 + \|v_1\|_2 \|\epsilon_l\|_2}{\|v_1\|_2}}_{=\mathcal{O}(\|\epsilon_l\|_2)} = \underbrace{\frac{|-2\text{Re}v_1^H \epsilon_l + \|\epsilon_l\|_2^2|}{\|v_1\|_2^2}}_{\mathcal{O}(\|\epsilon_l\|_2)} + \mathcal{O}(\|\epsilon_l\|_2)
\end{aligned}$$

$$3. \mu_l = x_l^H A x_l = \frac{(v_1 + \epsilon_l)^H A (v_1 + \epsilon_l)}{\|v_1 + \epsilon_l\|_2^2} = \frac{v_1^H A v_1 + v_1^H A \epsilon_l + \epsilon_l^H A v_1 + \epsilon_l^H A \epsilon_l}{\|v_1 + \epsilon_l\|_2^2} \text{ und } v_1^H A v_1 = \lambda_1 \|v_1\|_2^2$$

$$\begin{aligned}
\|v_1 + \epsilon_l\|_2 &\geq \|v_1\|_2 - \underbrace{\|\epsilon_l\|_2}_{\rightarrow 0} \leq \frac{1}{4} \|v_1\|_2 \forall l \geq l_0 \\
\mu_l - \lambda_1 &= \lambda_1 \underbrace{\left( \frac{\|v_1\|_2^2}{\|v_1 + \epsilon_l\|_2^2} - 1 \right)}_{= \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2}{\|v_1 + \epsilon_l\|_2^2} \leq \tilde{C} \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2}{\|v_1\|_2}} = \mathcal{O}(\|\epsilon_l\|_2^2) \\
&= \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2}{\|v_1 + \epsilon_l\|_2^2} \leq \tilde{C} \frac{\|v_1\|_2^2 - \|v_1 + \epsilon_l\|_2^2}{\|v_1\|_2} = \mathcal{O}(\|\epsilon_l\|_2^2)
\end{aligned}$$

4. Falls  $A = A^H$ , wähle  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  ONB aus EV.

$$\implies \|v_1 + \epsilon_l\|_2^2 = \|v_1\|_2^2 + \|\epsilon_l\|_2^2, A v_1 \in \text{span}\{v_1\}, A \epsilon_l \in \text{span}\{v_2, \dots, v_n\}$$

$$\text{Orthogonalität} \implies \mu_l - \lambda_1 = \mathcal{O}(\|\epsilon_l\|_2^2)$$

□

**Algorithmus 15** (Inverse Iteration).  $A \in \mathbb{K}^{n \times n}, x_0 \in \mathbb{K}^n \setminus \{0\}, \lambda \in \mathbb{K} \setminus \sigma(A)$

- Für  $l = 0, 1, 2, \dots$  (solange wie  $\|A x_l - \mu_l x_l\|_2$  "zu groß")
- Löse  $(A - \lambda) y_{l+1} = x_l$
- $x_{l+1} := \frac{y_{l+1}}{\|y_{l+1}\|_2}$
- $\mu_{l+1} := x_{l+1}^H A x_{l+1}$

also Power-Iteration für  $(A - \lambda)^{-1}$

**Korollar 10.**  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  EV-Basis,  $A v_j = \lambda_j v_j$  mit  $\frac{1}{|\lambda_1 - \lambda|} \gtrsim \frac{1}{|\lambda_2 - \lambda|} \geq \dots$

Sei  $x_0 = \sum_j \alpha_j v_j$  mit  $\alpha_j \neq 0$

$\implies$  Inverse Iteration wohldef. und ex  $C > 0, l_0 \in \mathbb{N}$  mit

$$|\mu_l - \lambda_1| \leq C \underbrace{\left| \frac{\lambda_2 - \lambda}{\lambda_1 - \lambda} \right|^l}_{=: q < 1}, d(\text{span}\{x_l\}, \text{span}\{v_1\}) \leq C q^l \forall l \geq l_0$$

Falls  $A = A^H$ , so gilt  $l_0 = 0$  und  $|\mu_2 - \lambda_1| \leq C q^{2l}$ .

*Proof.*  $A v_j = \lambda_j v_j, (A - \lambda) v_j = (\lambda_j - \lambda) v_j, \frac{1}{\lambda_j - \lambda} v_j = (A - \lambda)^{-1} v_j$

□

**Bemerkung 72.** • Aussagen über Konvergenz  $\mu_l$  für Power-It. und inverse It. gelten auch, falls  $\lambda_1$  ein mehrfacher EW ist. Aber inv. It scheitert, falls  $\lambda_1 \neq \lambda_2$ , aber  $|\lambda_1| = |\lambda_2|$ . Analoges gilt für inv. It. Dort aber  $\lambda$  wählbar mit  $\left| \frac{1}{\lambda_1 - \lambda} \right| \gtrsim \left| \frac{1}{\lambda_2 - \lambda} \right|$ , sofern  $\lambda_1 \neq \lambda_2$ .

- In der Praxis wählt man  $x_0 \neq 0$  zufällig und dann gilt (mit Wahrscheinlichkeit 1)  $\alpha_1 \neq 0$  (oder erfüllt durch Rundungsfehler im Verfahren).

**Lemma 26.**  $\implies d(X, Y) := \|\mathbb{P}_X - \mathbb{P}_Y\|_2 \stackrel{!}{=} \|(1 - \mathbb{P}_X)\mathbb{P}_Y\|_2 = \sup_{x \in X \setminus \{0\}} \inf_{y \in Y} \frac{\|x - y\|_2}{\|x\|_2}$

*Proof.* 1. zz:  $D(X, Y) := \|(1 - \mathbb{P}_Y)\mathbb{P}_X\|_2 = \sup_{x \in X \setminus \{0\}} \inf_{y \in Y} (\dots) \leq d(X, Y)$

$$\text{klar: } D(X, Y) = \|(\mathbb{P}_X - \mathbb{P}_Y)\mathbb{P}_X\|_2 \leq \underbrace{\|\mathbb{P}_X - \mathbb{P}_Y\|_2}_{=d(X, Y)} \underbrace{\|\mathbb{P}_X\|_2}_{=1}$$

2. TODO

3. TODO

4. TODO

$$\begin{aligned} \|(\mathbb{P}_X - \mathbb{P}_Y)z\|_2^2 &= \underbrace{\|\mathbb{P}_X(\mathbb{P}_X - \mathbb{P}_Y)z\|_2^2}_{=\mathbb{P}_X(1-\mathbb{P}_Y)^2} + \underbrace{\|(1 - \mathbb{P}_X)(\mathbb{P}_X - \mathbb{P}_Y)z\|_2^2}_{=-(1-\mathbb{P}_X)\mathbb{P}_Y^2} \leq \\ &\|\mathbb{P}_X(1 - \mathbb{P}_Y)\|_2^2 \|(1 - \mathbb{P}_Y)z\|_2^2 + \|(1 - \mathbb{P}_X)\mathbb{P}_Y\|_2^2 \|\mathbb{P}_Y z\|_2^2 \leq \\ &\max\left\{ \underbrace{\|\mathbb{P}_X(1 - \mathbb{P}_Y)\|_2^2}_{=D(X, Y)^2}, \underbrace{\|(1 - \mathbb{P}_X)\mathbb{P}_Y\|_2^2}_{=D(Y, X)^2} \underbrace{(\|(1 - \mathbb{P}_Y)z\|_2^2 + \|\mathbb{P}_Y z\|_2^2)}_{=\|z\|_2^2} \right\} \end{aligned}$$

$\|(1 - \mathbb{P}_Y)x\|_2 = \min_{y \in Y} \|x - y\|_2 \forall x \in \mathbb{K}^n$  und min wird eindeutig für  $y = \mathbb{P}_Y x$  angenommen!

5. zz:  $D(X, Y) = D(Y, X)$  (dann folgt  $d(X, Y) = D(X, Y)$ )

$$D(X, Y)^2 = \sup_{x \in X \setminus \{0\}} \inf_{y \in Y} \frac{\|x - y\|_2^2}{\|x\|_2^2} = \sup_{x \in X, \|x\|_2 \leq 1} \inf_{y \in Y, \|y\|_2 \leq 1} \underbrace{\|x - y\|_2^2}_{=\|x\|_2^2 + \|y\|_2^2 - 2\operatorname{Re}(x^H y)}$$

1. Fall:  $\dim X = 1 = \dim Y$ , wähle  $x \in X, \|x\|_2 = 1, y \in Y, \|y\|_2 = 1$

$$\implies D(X, Y)^2 = \sup_{s \in \mathbb{K}} \inf_{t \in \mathbb{K}} (|s|^2 + |t|^2 - 2\operatorname{Re}(s t x^H y)) \implies s, t \text{ vertauschbar!} \implies D(X, Y) = D(Y, X)$$

2. Fall:  $\dim X = k = \dim Y$

Seien  $\hat{X}, \hat{Y} \in \mathbb{K}^{n \times k}$  ONB (als Matrix geschrieben) zu  $X$  bzw.  $Y$ .

Singulärwertzerlegung  $\implies \hat{X}^H \hat{Y} = U \Sigma V^H$  mit  $U, V \in \mathbb{K}^{k \times k}$  unitär/orthogonal,  $\Sigma \in \mathbb{R}_{\geq 0}^{k \times k}$  diagonal

$$\begin{aligned} \implies D(X, Y)^2 &= \sup_{\alpha \in \mathbb{K}^k, \|U^H \alpha\|_2 \leq 1} \inf_{\beta \in \mathbb{K}^k, \|\beta\|_2 \leq 1} \underbrace{\|\hat{X} \alpha - \hat{Y} \beta\|_2^2}_{=\|\hat{X} \alpha\|_2^2 + \|\hat{Y} \beta\|_2^2 - 2\operatorname{Re}(\hat{X} \alpha)^H (\hat{Y} \beta)} = \\ &\sup_{\tilde{\alpha} \in \mathbb{K}^k, \|\tilde{\alpha}\|_2 \leq 1} \inf_{\tilde{\beta} \in \mathbb{K}^k, \|\tilde{\beta}\|_2 \leq 1} (\|\tilde{\alpha}\|_2^2 + \|\tilde{\beta}\|_2^2 - 2\operatorname{Re}(\tilde{\alpha}^H \tilde{\Sigma} \tilde{\beta})) \end{aligned}$$

$$\implies \tilde{\alpha}, \tilde{\beta} \text{ vertauschbar} \implies D(Y, X) = D(X, Y)$$

□

**Algorithmus 16** (Rayleigh-Iteration). *Input:*  $A = A^H \in \mathbb{K}^{n \times n}, x_0 \in \mathbb{K}^n$  mit  $\|x_0\|_2 = 1, \mu_0 := x_0^H A x_0$

- Für  $l = 0, 1, 2, \dots$  (solange wie  $\|Ax_l - \mu_l x_l\|_2$  "zu groß")
- Löse  $(A - \mu_l)y_{l+1} = x_l$  (Inverse It. mit  $\lambda = \mu_l$ )
- Def.  $x_{l+1} := \frac{y_{l+1}}{\|y_{l+1}\|_2}$
- $\mu_{l+1} := x_{l+1}^H A x_{l+1}$

**Satz 38.** Sei  $A = A^H \in \mathbb{K}^{n \times n}, \lambda \in \sigma(A)$  einfacher EW mit EV  $v \in \mathbb{K}^n \setminus \{0\}$

$\implies$  Ex.  $C > 0$  und  $\epsilon_0 > 0$ , sodass für alle  $0 < \epsilon \leq \epsilon_0$  und alle  $x_0 \in \mathbb{K}^n$  mit  $\|x_0\|_2 = 1$  gilt  
 $d(\operatorname{span}\{x_0\}, \operatorname{span}\{v\}) \leq \epsilon \implies (|\mu_0 - \lambda| \leq C\epsilon^2, |\mu_1 - \lambda| \leq C\epsilon^6, d(\operatorname{span}\{x_1\}, \operatorname{span}\{v\}) \leq C\epsilon^3 \leq (C\epsilon_0^2)\epsilon \leq \epsilon)$   
d.h. Rayleigh-Iteration ist (so etwas wie) lokal kubisch konvergent.

*Proof.* O.B.d.A  $d(\operatorname{span}\{x_0\}, \operatorname{span}\{v\}) > 0$

1. zz: Ex.  $v_1 \in \text{span}\{v\}$  mit  $\|v\|_2 = 1, \|x_0 - v_1\|_2 \leq 2\epsilon$

klar:  $\inf_{w \in \text{span}\{v\}} \|x_0 - w\|_2 \leq \sup_{x \in \text{span}\{x_0\}, \|x\|_2 \neq 0} \inf_{w \in \text{span}\{v\}} \frac{\|x - w\|_2^2}{\|x\|_2^2} = d(\text{span}\{x_0\}, \text{span}\{v\}) \leq \epsilon$ .

Wähle  $w \in \text{span}\{v\}$  mit  $\|x_0 - w\|_2 \leq \epsilon$

Def.  $v_1 := \frac{w}{\|w\|_2}$ , da  $\|x_0\|_2 = 1$ , also  $w \neq 0$  für  $\epsilon < 1$

$$\begin{aligned} \implies \|x_0 - v_1\|_2 &\leq \underbrace{\|x_0 - w\|_2}_{\leq \epsilon} + \underbrace{\left\|w - \frac{w}{\|w\|_2}\right\|_2}_{\leq \|w\|_2 - 1} \leq \|w - x_0\|_2 \leq \epsilon \\ &= \frac{1}{\|w\|_2} \underbrace{\| \|w\|_2 w - w \|_2}_{= \|(\|w\|_2 - 1)w\|_2 = (\|w\|_2 - 1)\|w\|_2} \leq \|w\|_2 - 1 \leq \|w - x_0\|_2 \leq \epsilon \end{aligned}$$

2. klar:  $v_k$  EV von  $A$

Ergänze zu ONB  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  aus EV zu  $A$ ,  $Av_j = \lambda_j v_j$  mit  $\lambda_1 = \lambda$ .

Wähle  $\alpha \in \mathbb{K}^n$  mit  $x_0 = (1 + \alpha_1)v_1 + \sum_{j \geq 2} \alpha_j v_j$

$$\implies \|\alpha\|_2^2 = \|\sum_{j \geq 2} \alpha_j v_j\|_2^2 = \|x_0 - v_1\|_2^2 \leq 4\epsilon^2$$

3. zz:  $|\lambda_1 - \mu_0| \leq 8\rho(A)\epsilon^2$

$$\begin{aligned} \mu_0 &= x_0^H A x_0 = |1 + \alpha_1|^2 \lambda_1 + \sum_{j \geq 2} |\alpha_j|^2 \lambda_j \\ \implies |\lambda_1 - \mu_0| &\leq \underbrace{|1 - |1 + \alpha_1|^2|}_{\leq 4\epsilon^2} \underbrace{|\lambda_1|}_{\leq \rho(A)} + \sum_{j \geq 2} \underbrace{|\alpha_j|^2}_{\leq 4\epsilon^2} \underbrace{|\lambda_j|}_{\leq \rho(A)} \\ 1 &= \|x_0\|_2^2 = |1 + \alpha_1|^2 + \sum_{j \geq 2} |\alpha_j|^2 \end{aligned}$$

4. zz: Für  $y_1 = (A - \mu_0)^{-1} x_0 =: \sum_{j \geq 2} \beta_j v_j$  gilt  $(|\beta_1| \geq \frac{1}{16\rho(A)}\epsilon^{-2}, \sum_{j \geq 2} |\beta_j|^2 \leq \frac{4}{\Delta_1^2}\epsilon^2, \Delta_1 := \min_{j \neq 1} |\lambda_1 - \lambda_j| > 0)$

$$\sum_{j \geq 1} \beta_j (\lambda_j - \mu_0) v_j = (A - \mu_0) y_1 = x_0 = (1 + \alpha_1) v_1 + \sum_{j \geq 2} \alpha_j v_j$$

$$\implies \text{lin. unabh. zeigt } (1 + \alpha_1) = \beta_1 (\lambda_1 - \mu_0), \alpha_j = \beta_j (\lambda_j - \mu_0) \forall j \geq 2$$

$$\implies |\beta_1| = \frac{|1 + \alpha_1|}{|\lambda_1 - \mu_0|} \geq \frac{1 - |\alpha_1|}{|\lambda_1 - \mu_0|} \geq \frac{1}{2} \frac{1}{|\lambda_1 - \mu_0|} \geq \frac{1}{16\rho(A)} \epsilon^2$$

$$\text{und } |\beta_j| = \frac{|\alpha_j|}{|\lambda_j - \mu_0|} \leq \frac{2}{\Delta_1} |\alpha_j|$$

$$\sum_{j \geq 2} |\beta_j|^2 \leq \frac{4}{\Delta_1^2} \sum_{j \geq 2} |\alpha_j|^2 \leq \frac{16}{\Delta_1^2} \epsilon^2 \leq \frac{16}{\Delta_1^2} \epsilon^2 \leq 4\epsilon^2$$

5.

$$\begin{aligned} d(\text{span}\{x_1\}, \text{span}\{v\}) &= \sup_{x \in \text{span}\{x_1\}, x \neq 0} \inf_{w \in \text{span}\{v\}} \frac{\|x - w\|_2}{\|x\|_2} = \sup_{s \in \mathbb{K} \setminus \{0\}} \inf_{t \in \mathbb{K}} \frac{\|s y_1 - t v_1\|_2}{\|s y_1\|_2} = \\ &= \inf_{t \in \mathbb{K}} \frac{\|y_1 - t v_1\|_2}{\|y_1\|_2} \leq \frac{\left(\sum_{j \geq 2} |\beta_j|^2\right)^{\frac{1}{2}}}{|\beta_1|} \leq \frac{4 \cdot 16\rho(A)}{\Delta_1} \epsilon^3 \end{aligned}$$

6.

$$\begin{aligned} \mu_1 &= \frac{y_1^H A y_1}{\|y_1\|_2^2} = \frac{(\beta_1 v_1 + e)^H A (\beta_1 v_1 + e)}{\|\beta_1 v_1 + e\|_2^2} = \frac{|\beta_1|^2 \lambda_1 + e^H A e}{\|\beta_1 v_1 + e\|_2^2} \\ |\mu_1 - \lambda_1| &= \left| \frac{|\beta_1|^2 - \|\beta_1 v_1 + e\|_2^2}{\|\beta_1 v_1 + e\|_2^2} \right| |\lambda_1| + \left| \frac{e^H A e}{\|\beta_1 v_1 + e\|_2^2} \right| = 2 \frac{\|e\|_2^2}{|\beta_1|^2 + \|e\|_2^2} \rho(A) \leq 2\rho(A) \frac{\sum_{j \geq 2} |\beta_j|^2}{|\beta_1|^2} \leq 2\rho(A)^3 \frac{16}{\Delta_1^2} \epsilon^6 \leq 256 \end{aligned}$$

□

**Satz 39** (Singulärwertzerlegung). Zu  $A \in \mathbb{K}^{m \times n}$  ex. unitäre/orthogonale Matrizen  $U \in \mathbb{K}^{m \times m}, V \in \mathbb{K}^{n \times n}$  und eine verallgemeinerte Diagonalmatrix  $\Sigma \in \mathbb{R}^{m \times n}$  mit  $\Sigma_{jk} = \sigma_j \delta_{jk}$  mit  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}} \geq 0$  und  $A = U \Sigma V^H$ , sog. **Singulärwertzerlegung**.

Es gelten:

1. Die Matrix  $\Sigma$  ist eindeutig und  $\sigma_j^2 \in \sigma(A^H A)$
2. Mit  $\text{rang}(A) = r$  gilt  $\sigma_1 \geq \dots \sigma_r > 0 = \sigma_{r+1} = \dots$
3.  $\|A\|_2 = \sigma_1$

*Proof.* (ii) trivial, da  $\text{rang}(A) = \text{rang}(\Sigma) = \max\{j | \sigma_j \neq 0\}$

$$(iii) \|A\|_2 = \sqrt{\rho(A^H A)} = \sqrt{\sigma_1^2} = \sigma_1$$

$$(i) A = U \Sigma V^H \implies A^H A = (V \underbrace{\Sigma^H}_{=\Sigma^T} \underbrace{U^H}_{=Id}) U \Sigma V^H = V \Sigma^T \Sigma V^H$$

$$\implies \sigma(\underbrace{\Sigma^T \Sigma}_{=D \in \mathbb{R}^{n \times n} \text{ Diagonalmatrix}}) = \sigma(A^H A) \subseteq \mathbb{R}_{\geq 0} \implies \sigma(\Sigma^T \Sigma) = \{\underbrace{D_{jj}}_{=\sigma_j^2} | j = 1, \dots, n\}$$

O.B.d.A. sortieren  $\sigma_1 \geq \sigma_2 \geq \dots$ , also eindeutig. insb.  $\Sigma$  eindeutig.

Existenz:  $A^H A \in \mathbb{K}^{n \times n}$  selbstadjungiert, positiv semidefinit  $\implies$  Ex.  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  ONB mit  $A^H A v_j = \mu_j v_j, \mu_j \geq 0, \mu_1 \geq \dots \geq \mu_n$ .

Sei  $r \in \mathbb{N}$  mit  $\mu_1 \geq \dots \geq \mu_r > 0 = \mu_{r+1} = \dots = \mu_n$ .

Def.  $\sigma_j := \sqrt{\mu_j}, S := \text{diag}(\sigma_1, \dots, \sigma_r) \in \mathbb{R}^{r \times r}$  regulär.

Def.  $V := (v_1, \dots, v_n) = (V_1, V_2)$  unitär/orth. mit  $V_1 \in \mathbb{K}^{n \times r}, V_2 \in \mathbb{K}^{n \times (n-r)}$

klar:  $A^H A V_1 = V_1 S^2$

Def.  $U_1 := A V_1 S^{-1} \in \mathbb{K}^{m \times r}$

$$\implies U_1^H U_1 = S^{-1} V_1^H \underbrace{A^H A V_1}_{=V_1 S^2} S^{-1} = S^{-1} \underbrace{V_1^H V_1}_{=Id} S^2 S^{-1}, \text{ d.h. Spalten von } U_1 \text{ sind orthonormal.}$$

Ergänze  $U := (U_1 U_2) \in \mathbb{K}^{m \times m}$  unitär/orthogonal

$$\begin{aligned} \implies U^H A V &= \begin{pmatrix} U_1^H \\ U_2^H \end{pmatrix} A (V_1 V_2) = \begin{pmatrix} U_1^H A V_1 & U_1^H A V_2 \\ U_2^H A V_1 & U_2^H A V_2 \end{pmatrix} = \Sigma \\ U_1^H A V_1 &= (S^{-1} V_1^H \underbrace{A^H A V_1}_{=V_1 S^2}) S^{-1} = S^{-1} \underbrace{V_1^H V_1}_{=Id} S^2 S^{-1} = S \\ U_2^H A V_1 &= \underbrace{U_2^H U_1}_{=0} S = 0 \end{aligned}$$

$V_2$  sind EV zum EW  $\mu = 0 \implies A^H A V_2 = 0$

$$\implies \|A v_2 x\|_2^2 = x^H V_2^H A^H A V_2 x = 0 \forall x \implies A V_2 = 0 \implies \Sigma = \begin{pmatrix} S & 0 \\ 0 & 0 \end{pmatrix} \quad \square$$

**Korollar 11** (Pseudo-Inverse, Moore-Penrose-Inverse). 1. Zu  $A \in \mathbb{K}^{m \times n}$  existiert eindeutiges  $A^+ \in \mathbb{K}^{n \times m}$  mit

- $(A^+ A)^H = A^+ A$
- $(A A^+)^H = A A^+$
- $A A^+ A = A$
- $A^+ A A^+ = A^+$

2. Ist  $A \in \mathbb{K}^{n \times n}$  regulär, so gilt  $A^+ = A^{-1}$

3. Falls  $A_{jk} = \sigma_j \delta_{jk} \in \mathbb{R}$ , so gilt  $A_{kj}^+ = \sigma_j^+ \delta_{jk}$  mit  $\sigma_j^+ = \begin{cases} \sigma_j^{-1} & \text{für } \sigma_j \neq 0 \\ 0 & \text{sonst} \end{cases}$

4. Mit  $SWZ A = U \Sigma V^H$  gilt  $A^+ = V \Sigma^+ U^H$

*Proof.* Eindeutigkeit: Seien  $B, C$  Pseudo-Inverse

$$\begin{aligned} \implies B &= B A B = B \underbrace{(A C A)}_{=A} B = B A \underbrace{(C A C)}_{=C} \underbrace{(A C A)}_{=A} B = (B A)^H (C A)^H C (A C)^H (A B)^H = \\ &\underbrace{A^H B^H A^H}_{=(A B A)^H = A^H} C^H C C^H \underbrace{A^H B^H A^H}_{=(A B A)^H = A^H} = \underbrace{A^H C^H}_{=(C A)^H = C A} C \underbrace{C^H A^H}_{=(A C)^H = A C} = \underbrace{C A C}_{=C} A C = C A C = C \end{aligned}$$

$\square$

**Bemerkung 73.** Betrachte Lösungsmenge  $\mathcal{L} := \{x \in \mathbb{K}^n | \underbrace{A^H A x = A^H b}_{\text{Gauß'sche Normalgleichung}}\} \neq \emptyset$  zu  $A \in \mathbb{K}^{m \times n}, b \in \mathbb{K}^m$

zum Linearen Ausgleichsproblem.

$\implies A^+ b \in \mathbb{L}$  und eindeutig mit  $\|A^+ b\|_2 = \min_{x \in \mathbb{L}} \|x\|_2$  sog. **Minimum-Norm-Lösung von LAP.**  
(folgt durch Ausrechnen!)

### 7.3 Orthogonale Iteration und QR-Zerlegung

**Bemerkung 74.** Ab jetzt schreiben wir eine Basis  $\{x_1, \dots, x_k\} \subseteq \hat{X} \leq \mathbb{K}^n$  als Matrix  $X := (x_1, \dots, x_k) \in \mathbb{K}^{n \times k}$

**Algorithmus 17** (orthogonale Iteration).  $A \in \mathbb{K}^{n \times n}$ ,  $X_0 \in \mathbb{K}^{n \times k}$  Basis von  $\hat{X}_0 \leq \mathbb{K}^n$

- Für  $l = 0, 1, 2, \dots$
- Berechne (reduzierte) QR-Zerlegung mit  $Q_l \in \mathbb{K}^{n \times k}$  mit orthonormalen Spalten  $R_l \in \mathbb{K}^{k \times k}$  rechte obere Dreiecksmatrix
- Def.  $X_{l+1} := AQ_l$

**Satz 40.**  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  Basis mit  $Av_j = \lambda_j v_j$  mit  $|\lambda_1| \geq \dots \geq |\lambda_k| \geq |\lambda_{k+1}| \leq \dots \leq |\lambda_n|$ . Sei  $\hat{X}_0 \leq \mathbb{K}^n$  Unterraum mit Basis  $X_0 \in \mathbb{K}^{n \times k}$  und  $\hat{X}_0 \cap \text{span}\{v_{k+1}, \dots, v_n\} = \{0\}$

$\implies$  Ex.  $C > 0$  mit  $d(A^l \hat{X}_0, \text{span}\{v_1, \dots, v_k\}) \leq C \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l, \max_{\tilde{\lambda} \in \sigma(Q_l^H A Q_l)} \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq C \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l$ , d.h. simultane Approximation der ersten  $k$  EW und zugehöriger EV.

*Proof.* 1. zz: Für  $\hat{X}_l := A^l \hat{X}_0$  gilt  $\dim \hat{X}_l = k$

Sei  $\{x_1, \dots, x_l\} \subseteq \hat{X}_0$  Basis. zz:  $\{A^l x_1, \dots, A^l x_n\}$  lin. unabh.

Seien  $\alpha_j \in \mathbb{K}$  mit  $0 = \sum_{j=1}^k \alpha_j A^l x_j = A^l \left( \sum_{j=1}^k \alpha_j x_j \right) \implies \sum_{j=1}^k \alpha_j x_j \in \text{kern}(A^l) \supseteq \text{kern}(A)$

$x \in \text{kern} A^l, A$  diagonalisierbar  $A^l = T^{-1} D^l T, y := Tx \implies D^l y = 0 \implies D_{jj}^l y_j = 0 \forall j \implies D_{jj} y_j = 0 \forall j \implies Ax = T^{-1} D \underbrace{Tx}_{=0} = 0 \implies \sum_j \alpha_j x_j \in \text{kern}(A) \subseteq \text{span}\{v_{k+1}, \dots, v_n\}$

$\sum_{j=1}^k \alpha_j x_j \in \hat{X}_0 \cap \text{span}\{v_{k+1}, \dots, v_n\} \implies \sum_{j=1}^k \alpha_j x_j = 0 \implies \alpha_j = 0 \forall j$

2. Für  $x = \sum_{j=1}^n \alpha_j v_j \in \mathbb{K}^n$  def.  $\|x\| := \sum_{j=1}^n |\alpha_j|, |||x||| := \sum_{j=1}^k |\alpha_j|$

zz:

- $\|\cdot\|$  ist eine Norm auf  $\mathbb{K}^n$
- $|||\cdot|||$  ist eine äquivalente Norm auf  $\hat{X}_0$

zz:  $|||\cdot|||$  definiert auf  $\hat{X}_0$ . Es gelte  $|||x||| = 0$  und  $x \in \hat{X}_0$  zz:  $x = 0$

$x = \sum_{j=k+1}^n \alpha_j v_j \in \text{span}\{v_{k+1}, \dots, v_n\} \cap \hat{X}_0 = \{0\}$

3. zz:  $d(\underbrace{A^l \hat{X}_0}_{=\hat{X}_l}, \underbrace{\text{span}\{v_1, \dots, v_n\}}_{=:V}) \leq C \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l$

Notation:  $x = \sum_{j=1}^n \alpha_j(x) v_j \in \hat{X}_0, v = \sum_{j=1}^k \beta_j(v) v_j \in V$

$$\begin{aligned} d(\hat{X}_l, V) &= \sup_{x \in \hat{X}_l, x \neq 0} \inf_{v \in V} \frac{\|x - v\|_2}{\|x\|_2} = \sup_{x \in \hat{X}_0, x \neq 0} \inf_{u \in V} \frac{\|A^l x - v\|_2}{\|A^l x\|_2} \stackrel{\substack{= \frac{\|A^l x - v\|}{\|A^l x\|} \text{ Normäquivalenz}}}{=} \\ &= \sup_{x \in \hat{X}_0, x \neq 0} \inf_{v \in V} \frac{\sum_{j=1}^k |\lambda_j^l - \alpha_j(x) - \beta_j(v)| + \sum_{j=k+1}^n |\lambda_j^l \alpha_j(x)|}{\sum_{j=1}^k |\lambda_j^l \alpha_j(x)|} \leq \\ &= \sup_{x \in \hat{X}_0, x \neq 0} \frac{\sum_{j=k+1}^n |\lambda_j^l \alpha_j(x)|}{\sum_{j=1}^k |\lambda_j^l \alpha_j(x)|} \leq \sup_{x \in \hat{X}_0, x \neq 0} \frac{|\lambda_{k+1}|^l \sum_{j=1}^n |\alpha_j(x)|}{|\lambda_k|^l \sum_{j=1}^k |\alpha_j(x)|} = \sup_{x \in \hat{X}_0, x \neq 0} \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l \underbrace{\frac{\|x\|}{|||x|||}}_{=1} \end{aligned}$$

4.  $\hat{X}, \hat{Y} \leq \mathbb{K}^n, \dim X = \dim Y = k$

$Q \in \mathbb{K}^{n \times k}$  ONB von  $\hat{X}, U := \mathbb{P}_{\hat{Y}} Q \in \mathbb{K}^{n \times k}$  spaltenweise

$$d(\hat{X}, \hat{Y}) = \sup_{x \in \hat{X}, x \neq 0} \inf_{y \in \hat{Y}} \frac{\|x - y\|_2}{\|x\|_2} = \sup_{\alpha \in \mathbb{K}^k, \alpha \neq 0} \frac{\|Q\alpha - \mathbb{P}_{\hat{Y}}(Q\alpha)\|_2}{\|Q\alpha\|_2} = \sup_{\alpha \in \mathbb{K}^k, \alpha \neq 0} \frac{\|Q\alpha - \overbrace{\mathbb{P}_{\hat{Y}} Q}^{=U} \alpha\|_2}{\|Q\alpha\|_2} = \|Q - U\|_2$$

5. Ergänze  $Q_l \in \mathbb{K}^{n \times k}$  zu orth./unitäre Matrix  $\tilde{Q}_l = (Q_l, Q'_l)$  mit  $Q'_l \in \mathbb{K}^{n \times (n-k)}$

$$\implies \tilde{Q}_l^H A \tilde{Q}_l = \begin{pmatrix} Q_l^H A Q_l & Q_l^H A Q'_l \\ Q'^H_l A Q_l & Q'^H_l A Q'_l \end{pmatrix} =: \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

$$\text{zz: } \|Q'^H_l A Q_l\|_2 \leq 2d(\hat{X}_l, V) \|A\|_2$$

$$\text{Def. } U_1 := \mathbb{P}_V Q_l \in \mathbb{K}^{n \times k}, U_2 := \mathbb{P}_{V^\perp} Q'_l \in \mathbb{K}^{n \times (n-k)}$$

$$A_{21} = 0, \max_{\tilde{\lambda} \in \sigma(\tilde{A})} \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \underbrace{\text{cond}_2(T) \|A - \tilde{A}\|_2}_{= \|A_{21}\|_2}$$

$$\|Q'_l A Q_l\|_2 \leq \|(Q'_l - U_2)^H\|_2 \|A\|_2 \|Q_l\|_2 + \|U_2^H\|_2 \|A\|_2 \|(Q_l - U_1)\|_2 + U_2^H A U_1 \leq \|Q'_l - U_2\|_2 \|A\|_2$$

6. gezeigt:  $\tilde{Q}_l^H A Q_l = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$  mit  $\|A_{21}\|_2 \leq 2\|A\|_2 d(\hat{X}_l, V) \leq \tilde{C} \left| \frac{\tilde{\lambda}_{k+1}}{\lambda_k} \right|^l$

$$\text{Definiere } \tilde{A} := \begin{pmatrix} A_{11} & A_{12} \\ 0 & A_{22} \end{pmatrix}$$

$$\text{zz: } \sigma(\tilde{A}) = \sigma(A_{11}) \cup \sigma(A_{22})$$

$$\text{"} \subseteq \text{" } \begin{pmatrix} A_{11}x + A_{12}y \\ A_{22}y \end{pmatrix} = \tilde{A} \begin{pmatrix} x \\ y \end{pmatrix} = \lambda \begin{pmatrix} x \\ y \end{pmatrix} \text{ mit } \begin{pmatrix} x \\ y \end{pmatrix} \neq 0, \lambda \in \mathbb{K}$$

$$1. \text{ Fall: } y = 0 \implies x \neq 0 \implies A_{11}x = \lambda x \implies \lambda \in \sigma(A_{11})$$

$$2. \text{ Fall: } y \neq 0 \implies A_{22}y = \lambda y \implies \lambda \in \sigma(A_{22})$$

$$\text{"} \supseteq \text{" } 1. \text{ Fall: } \lambda \in \sigma(A_{11}) \implies A_{11}x = \lambda x \text{ mit } x \neq 0. \text{ Wähle } y = 0 \implies \tilde{A} \begin{pmatrix} x \\ 0 \end{pmatrix} = \lambda \begin{pmatrix} x \\ 0 \end{pmatrix} \implies \lambda \in \sigma(\tilde{A})$$

$$2. \text{ Fall: } \lambda \in \sigma(A_{22}) \setminus \sigma(A_{11}) \implies A_{22}y = \lambda y \text{ mit } y \neq 0$$

$$\text{zz: Ex. } x \text{ mit } A_{11}x + A_{12}y = \lambda x \iff A_{12}y = \underbrace{-(A_{11} - \lambda)x}_{\text{regulär}} \iff x = -(A_{11} - \lambda)^{-1} A_{12}y \implies \lambda \in \sigma(\tilde{A}).$$

7.  $\|A - \tilde{A}\|_2 = \|A_{21}\|_2 \leq \tilde{C} \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l$  und  $\tilde{\lambda} \in \underbrace{\sigma(Q_l^H A Q_l)}_{= A_{11}} \subseteq \sigma(\tilde{A})$

$$\text{Bauer-Fike} \implies \max_{\tilde{\lambda} \in \sigma(Q_l^H A Q_l)} \min_{\lambda \in \sigma(A)} |\lambda - \tilde{\lambda}| \leq \text{cond}_2(T) \|A - \tilde{A}\|_2.$$

□

**Korollar 12.**  $A \in \mathbb{K}^{n \times n}$  diagonalisierbar,  $\{v_1, \dots, v_n\} \subseteq \mathbb{K}^n$  Basis mit  $Av_j = \lambda_j v_j$  und  $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_k| \geq |\lambda_{k+1}|$ ,  $\underbrace{\text{span}\{e_1, \dots, e_k\}}_{\text{Eigenvektoren}} \cap \text{span}\{v_{k+1}, \dots, v_n\} = \{0\} \forall k$ .

$$Q_l \in \mathbb{K}^{n \times n} \text{ aus arth. Iteration zu } X_0 = Id \in \mathbb{K}^{n \times n}$$

$$A_l := Q_l^H A Q_l \text{ (klar } \sigma(A_l) = \sigma(A))$$

$\implies$  Ex.  $C > 0$ , sodass für alle  $k$  gilt:  $C^{-1} \sum_{j=k+1}^n |(A_l)_{jk}| \leq \|A_l(k+1:n, 1:k)\|_2 \leq C \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l$ , d.h.  $A_l$  konvergiert gegen rechte obere  $\Delta$ -Matrix.

*Proof.* Voraussetzung erfüllt Voraussetzungen für orthogonale Iteration für jedes  $k$ , d.h. orthogonale Iteration mit  $k = n$  macht simultan orthogonale Iteration für alle  $k = 1, \dots, n$ . Für  $k$  fix, partitioniere  $Q_l = (Q_{l,k}, Q'_{l,k})$  mit  $Q_{l,k} \in \mathbb{K}^{n \times k}$ ,  $Q'_{l,k} \in \mathbb{K}^{n \times (n-k)}$

$$A_l = \begin{pmatrix} Q_{l,k}^H A Q_{l,k} & Q_{l,k}^H A Q'_{l,k} \\ Q'^H_{l,k} A Q_{l,k} & Q'^H_{l,k} A Q'_{l,k} \end{pmatrix} = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

$$\implies \|\underbrace{A_{21=A_l(k+1:n, 1:k)}}_{\|A_{21}\|_2}\|_2 \leq C \left| \frac{\lambda_{k+1}}{\lambda_k} \right|^l$$

$$\sum_{j=k+1}^n |(A_l)_{jk}| \leq \underbrace{\sum_{i=1}^n \sum_{j=k+1}^n |(A_l)_{ji}|}_{=: \|A_{21}\| \text{ Norm auf } \mathbb{K}^{(n-k) \times k}} = \|A_{21}\|_2$$

□

Herleitung des QR-Verfahrens

$$\begin{aligned}
 A Q_l &= X_{l+1} = Q_{l+1} R_{l+1} \\
 \Rightarrow A_l &= Q_l^H A Q_l = \underbrace{Q_l^H Q_{l+1}}_{=: \tilde{Q}_{l+1}} R_{l+1} \text{ ist QR-Zerlegung von } A_l \\
 \Rightarrow A_{l+1} &= Q_{l+1}^H A \underbrace{Q_{l+1}}_{=: Q_l \tilde{Q}_{l+1}} = \tilde{Q}_{l+1}^H \underbrace{Q_l^H A Q_l}_{=: A_l} \tilde{Q}_{l+1} = R_{l+1} \tilde{Q}_{l+1} \\
 &\quad \underbrace{\hspace{10em}}_{=: R_{l+1}}
 \end{aligned}$$

**Algorithmus 18** (QR-Zerlegung).  $A_0 := A \in \mathbb{K}^{n \times n}$

- Für  $l = 0, 1, 2, \dots$  (bis  $A_l$  dicht an oberer  $\triangle$ -Matrix)
- Berechne QR-Zerlegung  $A_l = Q_{l+1} R_{l+1}$
- Definiere  $A_{l+1} := R_{l+1} Q_{l+1}$

**Bemerkung 75.** 1. Unter den Voraussetzungen des Korollars konvergiert QR-Verfahren und  $\sigma(A_{l+1}) = \sigma(A_l) = \sigma(A)$

2. Naives Vorgehen braucht  $\mathcal{O}(n^3)$  Operationen pro Schritt (für QR-Zerlegung und Matrix-Matrix-Mult.). Tatsächlich reichen eine Initialisierung mit Aufwand  $\mathcal{O}(n^3)$  und danach sind alle Schritte  $\mathcal{O}(n^2)$  (bzw.  $\mathcal{O}(n)$ , falls  $A = A^H$ )

## 7.4 Hessenberg-Form einer Matrix

**Definition 17.** Matrix  $B \in \mathbb{K}^{n \times n}$  heißt **obere Hessenberg-Matrix**, gdw.  $B_{jk} = 0 \forall j > k + 1$

Ziel: Zu  $A \in \mathbb{K}^{n \times n}$  konstruiere  $Q \in \mathbb{K}^{n \times n}$  unitär/orthogonal, sodass  $B := Q^H A Q$  obere Hessenberg-Matrix und insb.  $\sigma(B) = \sigma(A)$

Danach: Effiziente Realisierung von QR-Verfahren für  $B$ .

Idee: Verwende Householder-Spiegelung  $H = Id - 2ww^H = H^H \in \mathbb{K}^{n-k}$

$$\Rightarrow \begin{pmatrix} Id & 0 \\ 0 & H \end{pmatrix} \begin{pmatrix} U & V \\ W & X \end{pmatrix} \begin{pmatrix} Id & 0 \\ 0 & H \end{pmatrix} = \begin{pmatrix} U & V \\ HW & HX \end{pmatrix} \begin{pmatrix} Id & 0 \\ 0 & H \end{pmatrix} = \begin{pmatrix} U & VH \\ HW & HXH \end{pmatrix}$$

Beachte:

- $U$ -Block bleibt unverändert
- $w$  wählbar, sodass eine Spalte von  $HW$  in  $\text{span}\{e_1\}$
- Nullspalte in  $W$  bleibt Nullspalte  $HW$

**Algorithmus 19.** Input:  $A_0 := A \in \mathbb{K}^{n \times n}$

1. Schritt:  $A_0 \rightsquigarrow A_1 = Q_1 A_0 Q_1$

2. Schritt:  $A_1 \rightsquigarrow A_2 = Q_2 A_1 Q_2$

Nach  $n - 2$  Schritten erhalte  $B = A_{n-2} = Q_{n-2} A_{n-1} Q_{n-2}$  obere Hessenberg-Matrix.

**Satz 41.** Algorithmus berechnet  $B = Q^H A Q$  mit  $B$  obere Hessenberg-Matrix und  $Q \in \mathbb{K}^{n \times n}$  unitär/orthogonal.

*Proof.*  $\Rightarrow B = A_{n-2} = \underbrace{Q_{n-2} Q_{n-3} \dots Q_2 Q_1}_{=(Q_1 \dots Q_{n-2})^H = Q^H} A \underbrace{Q_1 Q_2 \dots Q_{n-2}}_{=Q}$  □

**Bemerkung 76.** Asymptotischer Aufwand (bei cleverer Realisierung der Householder-Matrizen) ist  $\mathcal{O}(n^3) = \frac{10}{3}n^3 + \mathcal{O}(n^2)$



**Definition 18.** Zu  $t, c, s \in \mathbb{R}$  mit  $c^2 + s^2 = 1$  und  $t = 0$  für  $\mathbb{K} = \mathbb{R}$  definiere die **Givens-Rotation**

$$G_{kj}^{tcs} := \begin{pmatrix} 1 & & & & & & & \\ & \ddots & & & & & & \\ & & 1 & & & & & \\ & & & c & - & - & - & (+e^{it}s) \\ & & & | & 1 & & & | \\ & & & | & & \ddots & & | \\ & & & | & & & 1 & | \\ & & (-e^{it}s) & - & - & - & - & c \\ & & & & & & & 1 \\ & & & & & & & \ddots \\ & & & & & & & 1 \end{pmatrix}$$

**Lemma 27.** 1.  $G_{kj}^{tcs}$  ist unitär/orthogonal.

2.  $(G_{kj}^{tcs} A)$  hat nur  $j$ -te und  $k$ -te Zeile anders als  $A$ .
3.  $(AG_{kj}^{tcs})$  hat nur  $j$ -te und  $k$ -te Spalte anders als  $A$ .  
(nämlich jeweils linear kombiniert).

**Lemma 28.**  $A \in \mathbb{K}^{n \times n}, A_{kj} \neq 0, B := G_{kj}^{tcs} A$  mit  $c := \frac{|A_{jj}|}{\sqrt{|A_{jj}|^2 + |A_{kj}|^2}}, s := \frac{|A_{kj}|}{\sqrt{|A_{jj}|^2 + |A_{kj}|^2}}$  und  $\lambda := e^{it} = \frac{\text{sign}(A_{jj})}{\text{sign}(A_{kj})}$   
 $\implies B_{kj} = 0$  und nur Zeile  $j$  und  $k$  anders als bei  $A$ .

*Proof.*

$$\begin{aligned} B_{kj} &= \sum_l (G_{kj}^{tcs})_{kl} A_{lj} = (G_{kj}^{tcs})_{kk} A_{kj} + (G_{kj}^{tcs})_{kj} A_{jj} = c A_{kj} + \underbrace{e^{it}s}_{=\lambda} A_{jj} = \\ &= c A_{kj} + \frac{|A_{jj}|}{\text{sign}(A_{kj})} s = \frac{|A_{jj}| |A_{kj}|}{(\dots)} - \frac{|A_{jj}| |A_{kj}|}{(\dots) \text{sign}(A_{kj})} = 0 \end{aligned}$$

□

**Algorithmus 20** (QR-Zerlegung einer oberen Hessenberg-Matrix). *Input:*  $A = A_0 \in \mathbb{K}^{n \times n}$  obere Hessenberg-Matrix

1. Schritt:  $A_1 := G_{21} A_0$  mit  $G_{21} \in \{\text{Id}, \text{geeinete Givens-Rotation}\}$
2. Schritt:  $A_2 := G_{32} A_1$

Nach  $n - 1$  Schritten erhalte  $R := G_{n,n-1} A_{n-2}$  obere  $\triangle$ -Matrix und  $Q := G_{21}^H \dots G_{n,n-1}^H$  unitär/orthogonal

**Satz 42.** Sei  $A \in \mathbb{K}^{n \times n}$  obere Hessenberg-Matrix

1. Alg. berechnet in  $\mathcal{O}(n^2)$  Operationen eine QR-Zerlegung  $A = QR$ .
2.  $B := RQ$  ist eine obere Hessenberg-Matrix und kann in  $\mathcal{O}(n^2)$  Operationen berechnet werden.
3.  $A = A^H \implies B = B^H$  und (i) + (ii) können in  $\mathcal{O}(n)$  Operationen durchgeführt werden, da  $R$  obere Bandbreite 2 hat

*Proof.* 1.  $R = G_{n,n-1} A_{n-2} = G_{n,n-1} G_{n-1,n-2} A_{n-3} = \underbrace{G_{n,n-1} \dots G_{21}}_{=Q^H} A \implies QR = A$

Habe  $\mathcal{O}(n)$  Schritte, jeder Schritt macht Linearkombinationen von 2 Zeilen, d.h.  $\mathcal{O}(n)$  pro Schritt  $\implies \mathcal{O}(n^2)$  insgesamt.

2.  $B = RQ = \underbrace{RG_{21}^H}_{=B_1} \underbrace{G_{32}^H \dots G_{n,n-1}^H}_{=B_2}$ , d.h.  $n - 1$  Schritte zur Berechnung und  $\mathcal{O}(n)$  pro Schritt, also  $\mathcal{O}(n^2)$  insgesamt.

klar:  $B = RQ$  ist Hessenberg.

3.  $A = A^H$  Hessenberg  $\implies$  tridiagonal, d.h. pro Schritt maximal 3 Einträge pro Zeile linear kombinieren, d.h.  $\mathcal{O}(1)$  pro Schritt, also  $\mathcal{O}(n)$  zur Berechnung  $A = QR$   
 zz:  $R$  hat obere Bandbreite 2  
 $\implies \mathcal{O}(n)$  für Berechnung  $RQ$ , da  $R$  maximal 3 Elte. pro Zeile.

□

**Bemerkung 77.**  $B = RQ = G_{n,n-1} \dots G_{2,1} A G_{2,1}^H \dots G_{n,n-1}^H$ , d.h. berechne  $R$  aus  $A$ , ohne  $QR$ -Zerlegung explizit zu berechnen!