

# Mitschrift Einf. in die Num. Mathematik

Vorlesung SS10 Prof. Dr. H. Yserentant

Technische Universität Berlin

## Vorwort

Dies ist eine Mitschrift, *kein* Skript! Für eventuelle Fehler übernimmt niemand die Verantwortung, sollten allerdings welche entdeckt werden freue ich mich über Hinweise.

WICHTIG: Auch dient die Mitschrift nicht als Ersatz für die Vorlesung, wichtige Bemerkungen und Erklärungen des Dozenten tauchen nicht auf. Ich empfehle euch also trotzdem regelmäßig zur Vorlesung zu gehen.

Tipp: Nicht immer gleich ausdrucken wenns online ist, im aktuellen Teil sind immer massenhaft Tippfehler die ich erst dann korrigiere wenn sie mir auffallen(d.h. wenn ich die Woche darauf damit arbeite).

## Inhaltsverzeichnis

### **1. Rechnerarithmetik und Rundungsfehler**

### **2. Skalare Gleichungen**

2.1 Das lokalisieren von Nullstellen mit Vorzeichenwechsel

2.2 Das Newton-Verfahren

2.3 Allgemeine Iterationsverfahren

2.4 Die Regula falsi

### **3. Lineare Gleichungssysteme**

3.1 Normen

3.2 Die Kondition von Matrizen

3.3 Die LR-Zerlegung einer Matrix

3.4 Eine Fehleranalyse für den Gaußschen Algorithmus

3.5 Symmetrisch positiv-definite Matrizen

3.6 Eliminationsverfahren für schwach besetzte Matrizen

3.7 Iterationsverfahren

### **4. Nichtlineare Gleichungssysteme**

4.1 Der Banachsche Fixpunktsatz

4.2 Das Newton Verfahren

### **5. Lineare Ausgleichsprobleme**

5.1 Lineare Ausgleichsprobleme in unitären Räumen

5.2 Lineare Ausgleichsprobleme im  $\mathbb{R}^n$

5.3 Singulärwertzerlegung und Pseudoinverse

5.4 Das Gauß-Newton-Verfahren

### **6. Interpolation**

6.1 Interpolation durch Polynom

6.2 Beste Approximation und optimale Stützstellen

6.3 Interpolation durch Splines

### **7. Numerische Integration**

7.1 Zusammengesetzte Quadraturformeln

7.2 Interpolatorische Quadraturformeln

7.3 Orthogonalpolynome

7.4 Gaußsche Quadraturformeln

### **8. Gewöhnliche Differentialgleichungen**

# 1. Rechnerarithmetik und Rundungsfehler

## Normalisierte Gleitkommazahl

$$\pm \underbrace{0, \overset{=0}{\downarrow} 12345}_{\text{Mantisse}} \cdot 10^{\overbrace{-14}^{\text{Exponent}}}$$

$\pm 0,12345E - 14$  und die Null!

## Unsere Beispiele:

Dezimalsystem, 5 Dezimalstellen

$$\pm 0,12345E + 122$$

bei großer Exponent. (schließt "underflow" und "overflow" aus)

## IEEE-Standard

Dualsystem

32 bit (real) : Mantisse 23 bit  
Exponent 8 bit  
Vorzeichen 1 bit

64 bit (double) : Mantisse 52 bit  
Exponent 11 bit  
Vorzeichen 1 bit

## Rundungsoperator

$G$  = Menge der Gleitkommazahlen

$\text{rd} : \mathbb{R} \rightarrow G$  Rundungsoperator

$$|x - \text{rd}(x)| \leq |x - g|, g \in G$$

## **Beispiel 1.1:**

$$\text{rd}(7,12341234\dots) = \text{rd}(0,712341234\dots \cdot 10^1) = 0,71234E + 1$$

$$\text{rd}(7,12345678\dots) = \text{rd}(0,712345678\dots \cdot 10^1) = 0,71235E + 1$$

$$\text{rd}(0,99999999\dots) = 0,10000E + 1$$

## relativer Fehler:

$(x \neq 0)$

$$\left| \frac{\text{rd}(x) - x}{x} \right| \leq \left| \frac{5 \cdot 10^{-6}}{\text{Mantisse}} \right| \leq 5 \cdot 10^{-5}$$

(bei 5 Dezimalstellen)

$$\text{rd}(x) = x(1 + \varepsilon) \text{ , } |\varepsilon| \leq 5 \cdot 10^{-5} \text{ (gilt auch für } \text{rd}(0) = 0)$$

### ideale Gleitkommaoperationen

$x, y$  Gleitkommazahlen

$$x \oplus y = \text{rd}(x + y)$$

$$x \ominus y = \text{rd}(x - y)$$

$$x \odot y = \text{rd}(x \cdot y)$$

$$x \oslash y = \text{rd}(x/y), y \neq 0$$

### reale Maschinenoperationen

$x, y$  Gleitkommazahlen

$$x \oplus y = (x + y)(1 + \varepsilon_1)$$

$$x \ominus y = (x - y)(1 + \varepsilon_2)$$

$$x \odot y = (x \cdot y)(1 + \varepsilon_3)$$

$$x \oslash y = (x/y)(1 + \varepsilon_4)$$

$$|\varepsilon_i| = |\varepsilon_i(x, y)| \leq \varepsilon^*$$

$\varepsilon^*$  Maschinengenauigkeit

### Bemerkung:

Diese Beziehungen sind Grundlage jeder Rundungsfehleranalyse.

### Vorsicht:

Das Assoziativgesetz der Addition und Multiplikation und das Distributivgesetz gelten nicht exakt.

### **Beispiel 1.2**

Berechnung des Skalarprodukts  $s = \sum_{i=1}^n x_i y_i$  nach dem Algorithmus.

$$s_1 = x_1 y_1$$

$$s_{k+1} = s_k + x_{k+1} y_{k+1}, k = 1, \dots, n - 1$$

$$s = s_n$$

### Behauptung:

Für das berechnete Skalarprodukt  $\tilde{s}$  gilt

$$\tilde{s} = s + f, |f| \leq \underbrace{\left( (1 + (2 + \varepsilon^*)\varepsilon^*)^n - 1 \right)}_{=2n\varepsilon^* + O(\varepsilon^{*2})} \sum_{i=1}^n |x_i y_i|$$

Beweis: Sei  $\tilde{s}_k = s_k + f_k$  leite Rekursion für den Fehler  $f_k$  her!

Es gilt:

$$\tilde{s}_1 = x_1 y_1 (1 + \varepsilon_1) = x_1 y_1 + f_1 \text{ mit } f_1 = \varepsilon_1 x_1 y_1, |\varepsilon_1| \leq \varepsilon^*$$

Weiter ist  $\tilde{s}_k = s_k + f_k$

$$\tilde{s}_{k+1} := (\tilde{s}_k + x_{k+1} y_{k+1} \varepsilon_1)(1 + \varepsilon_2) = s_{k+1} + f_{k+1}$$

$$f_{k+1} = \varepsilon_2 s_{k+1} + (1 + \varepsilon_2) f_k + x_{k+1} y_{k+1} \varepsilon_2 (1 + \varepsilon_2), |\varepsilon_1|, |\varepsilon_2| \leq \varepsilon^*$$

Daraus folgt:

$$|f_1| \leq \varepsilon^* \sum_{i=1}^n |x_i y_i| \leq (2 + \varepsilon^*) \varepsilon^* \sum_{i=1}^n |x_i y_i|$$

$$|f_{k+1}| \leq (1 + (2 + \varepsilon^*)\varepsilon^*) |f_k| + \varepsilon^* (2 + \varepsilon^*) \sum_{i=1}^{k+1} |x_i y_i|$$

Mit Induktion folgt daraus

$$|f_k| \leq \left( (1 + (2 + \varepsilon^*)\varepsilon^*)^k - 1 \right) \sum_{i=1}^k |x_i y_i|$$

□

Interpolation:

Der relative Fehler kann sehr groß werden, falls

$$\left| \sum_{i=1}^n x_i y_i \right| \ll \sum_{i=1}^n |x_i y_i|$$

Es gilt aber

$$\tilde{s} = \sum_{i=1}^n (1 + \varepsilon_i) x_i y_i \text{ mit } |\varepsilon_i| \leq (1 + (2 + \varepsilon^*)\varepsilon^*)^n - 1$$

Das berechnete Ergebnis ist gleich dem exakten Ergebnis bei leicht gestörten Eingangsdaten.

Der Algorithmus ist rückwärtsstabil!

günstigere Strategie binärer Baum

$$\begin{array}{ccccccc}
 & x_1 y_1 & & x_2 y_2 & & x_3 y_3 & & x_4 y_4 \\
 & \searrow & & \swarrow & & \searrow & & \swarrow \\
 x_1 y_1 & + & & x_2 y_2 & & x_3 y_3 & + & x_4 y_4 \\
 & & & \searrow & & \swarrow & & \\
 & & & (x_1 y_1 + x_2 y_2) & + & (x_3 y_3 + x_4 y_4) & & 
 \end{array}$$

## 2. Skalare Gleichungen

gegeben:  $f : [a, b] \rightarrow \mathbb{R}$

gesucht:  $\bar{x} \in [a, b]$  mit  $f(\bar{x}) = 0$

### Definition 2.1

Der Punkt  $a \leq \bar{x} \leq b$  ist eine  $m$ -fache Nullstelle der stetigen Funktion  $f : [a, b] \rightarrow \mathbb{R}$ , wenn sich  $f$  in der Form  $f(x) = (x - \bar{x})^m g(x)$ ,  $g(\bar{x}) \neq 0$  mit einer stetigen Funktion  $g : [a, b] \rightarrow \mathbb{R}$  darstellen lässt.

### Beispiel 2.2

$$f(x) = 2x^3 - 4x^2 + 2x = 2x(x - 1)^2$$

$\bar{x} = 0$  ist einfache Nullstelle,

$\bar{x} = 1$  ist doppelte Nullstelle

Bemerkung: Vorzeichenwechsel finden nur in Nullstellen ungerader Vielfachheit statt.

### Satz 2.3

Die Funktion  $f$  sei in einem Intervall um den Punkt  $x = \bar{x}$   $m$ -mal stetig differenzierbar. Dann ist  $\bar{x}$  eine  $m$ -fache Nullstelle von  $f$ , wenn  $f(\bar{x}) = f'(\bar{x}) = \dots = f^{(m-1)}(\bar{x}) = 0$  und  $f^{(m)}(\bar{x}) \neq 0$  ist.

Beweis: Nach dem Taylorschen Satz ist

$$f(x) = \sum_{k=0}^{m-1} \frac{1}{k!} f^{(k)}(\bar{x})(x - \bar{x})^k + (x - \bar{x})^m g(x)$$

$$\text{mit } g(x) = \frac{1}{(m-1)!} \int_0^1 (1 - \theta)^{m-1} f^{(m)}(\bar{x} + \theta(x - \bar{x})) d\theta$$

$$\text{und } g(\bar{x}) = \frac{1}{m!} f^{(m)}(\bar{x})$$

□

### Satz 2.4

Die Funktion  $f$  sei auf einer Umgebung ihrer  $m$ -fachen Nullstelle mindestens  $(m + 1)$ -mal stetig differenzierbar.

Dann ist  $\bar{x}$  eine  $(m - 1)$ -fache Nullstelle von  $f'(x)$  und eine einfache Nullstelle von  $F(x) = \frac{f(x)}{f'(x)}$ .

Beweis: Die Funktion  $f$  hat die Darstellung  $f(x) = (x - \bar{x})^m g(x)$  mit der stetig differenzierbaren Funktion

$$g(x) = \frac{1}{(m-1)!} \int_0^1 (1-\theta)^{m-1} f^{(m)}(\bar{x} + \theta(x - \bar{x})) d\theta$$

Daher ist  $f'(x) = m(x - \bar{x})^{m-1} g(x) + (x - \bar{x})^m g'(x) = (x - \bar{x})^{m-1} h(x)$   
 $h(x) = m g(x) + (x - \bar{x}) g'(x)$ . Wegen  $h(\bar{x}) = m g(\bar{x}) \neq 0$  ist  $\bar{x}$  daher eine  $(m - 1)$ -fache Nullstelle von  $f'(x)$ .

$$\text{Weiter gilt } F(x) = \frac{f(x)}{f'(x)} = \frac{(x - \bar{x})^m g(x)}{(x - \bar{x})^{m-1} h(x)} = (x - \bar{x}) \frac{g(x)}{h(x)}$$

Wegen  $g(\bar{x})/h(\bar{x}) \neq 0$  ist  $\bar{x}$  daher einfache Nullstelle von  $F$ .

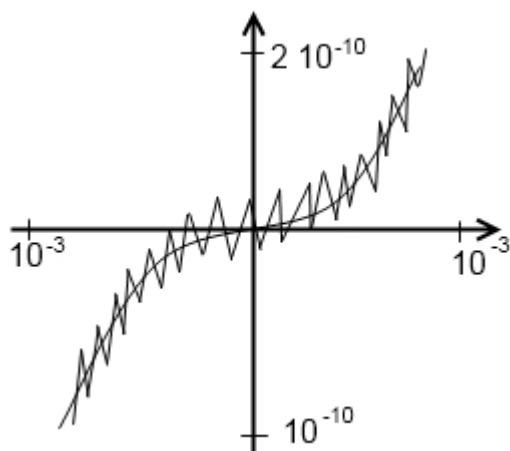
□

Soweit die Theorie! Zur Praxis:

### Beispiel 2.5:

$$f(x) = e^x - \left(1 + x + \frac{1}{2}x^2\right) \approx x^3 \left(\frac{1}{6} + \frac{x}{24} + \frac{x^2}{120}\right)$$

$\bar{x} = 0$  ist eine 3-fache Nullstelle.





## 2.1 Das Lokalisieren von Nullstellen mit Vorzeichenwechsel

Mache vom Zwischenwertsatz Gebrauch!

### Satz 2.6 (Bisektion)

Sei  $f : [a, b] \rightarrow \mathbb{R}$  stetig und  $f(a_0)f(b_0) < 0$ . Die Folge  $(x_n)_{n=0}^{\infty}$ ,

deren Glieder durch  $x_n = (a_n + b_n)/2$  mit

$a_{n+1} = a_n, b_{n+1} = x_n$ , falls  $f(a_n)f(x_n) \leq 0$ ,

$a_{n+1} = x_n, b_{n+1} = b_n$ , sonst

gegeben sind, erfüllt:

(i)  $\lim_{n \rightarrow \infty} x_n = \bar{x}, f(\bar{x}) = 0$

(ii)  $|x_n - \bar{x}| \leq \frac{1}{2}(b_n - a_n) = \frac{1}{2^{n+1}}(b_0 - a_0)$ .

Jedes Intervall  $[a_n, b_n]$  enthält die Nullstelle  $\bar{x}$ .

Vorteile:

-absolut sicher, falls man Rundungsfehler vernachlässigt

-Konvergenzgeschwindigkeit unabhängig von  $f$

Nachteil:

-sehr langsame Konvergenz verglichen mit anderen Verfahren

## 2.2 Das Newton-Verfahren

Idee: Gegeben sei eine Näherung  $x_0$  für  $\bar{x}$ . Linearisiere  $f$  in  $x_0$  :

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0)$$

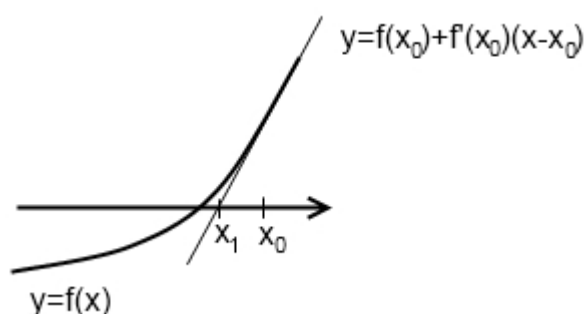
Berechne eine verbesserte(?) Näherung  $x_1$  als Lösung der

linearisierten Gleichung:

$$f(x_0) + f'(x_0)(x - x_0) = 0$$

Das ergibt  $x_1 = x_0 - \frac{f(x_0)}{f'(x_0)}$

Berechne ausgehend von  $x_1$  die Näherung  $x_2, \dots$



## Newton-Verfahren

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

iteratives Verfahren  $x_{n+1} = f(x_n)$ .

Setze  $x_n$  in die Iterationsvorschrift ein, um  $x_{n+1}$  zu erhalten.

### Beispiel 2.8:

$$f(x) = x^5 - 3 = 0, \bar{x} = \sqrt[5]{3}$$

Die Anzahl der sicheren Stellen scheint in jedem Schritt um einen festen Faktor zuzunehmen.

### Satz 2.9:

Sei  $\bar{x}$  eine einfache Nullstelle der zweimal stetig differenzierbaren Funktion  $f$ . Es sei  $\delta > 0$  und für alle  $x$  mit  $|x - \bar{x}| \leq \delta$  sei  $f(x)$  definiert und  $f'(x) \neq 0$ .

$$\text{Sei } c_1 := \min_{|x-\bar{x}| \leq \delta} |f'(x)| > 0, c_2 := \max_{|x-\bar{x}| \leq \delta} |f''(x)| > 0$$

Ist dann  $|x - \bar{x}| \leq \min\left(\delta, \frac{c_1}{c_2}\right)$ , so ist die Folge der Iterierten

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}, n = 0, 1, 2, \dots \text{ wohldefiniert.}$$

Die  $x_n$  konvergieren gegen  $\bar{x}$ . Es gilt:

$$|x_{n+1} - \bar{x}| \leq \frac{1}{2} \frac{c_2}{c_1} |x_n - \bar{x}|^2.$$

Beweis: Sei  $|x_n - \bar{x}| \leq \min\left(\delta, \frac{c_1}{c_2}\right)$ .

Dann ist  $x_{n+1}$  wohldefiniert und erfüllt die Gleichung

$$0 = f(x_n) + f'(x_n)(x_{n+1} - x_n).$$

Nach dem Taylorschen Satz ist

$$0 = f(\bar{x}) = f(x_n) + f'(x_n)(\bar{x} - c_n) + \frac{1}{2} f''(\eta)(\bar{x} - x_n)^2 \text{ oder}$$

$$x_{n+1} - \bar{x} = \frac{1}{2} \frac{f''(\eta)}{f'(x_n)} (x_n - \bar{x})^2$$

Daraus folgt

$$|x_{n+1} - \bar{x}| \leq \frac{1}{2} \frac{c_2}{c_1} |x_n - \bar{x}|^2$$

Wegen  $|x_n - \bar{x}| \leq \frac{c_1}{c_2}$  ist speziell

$$|x_{n+1} - \bar{x}| \leq \frac{1}{2} \frac{c_2}{c_1} |x_n - \bar{x}| |x_n - \bar{x}| \leq \frac{1}{2} |x_n - \bar{x}| \leq \min\left(\delta, \frac{c_1}{c_2}\right)$$

Daher ist die Folge der Iterierten wohldefiniert und

$$\lim_{n \rightarrow \infty} x_n = \bar{x}$$

gesucht: Nullstelle  $\bar{x}$  eine Funktion  $f : [\bar{x} - \delta, \bar{x} + \delta] \rightarrow \mathbb{R}$

Annahmen:  $f'(x) \neq 0$  auf  $[\bar{x} - \delta, \bar{x} + \delta]$  (wesentlich!)

$f''(x)$  existiert und ist beschränkt auf  $[\bar{x} - \delta, \bar{x} + \delta]$ .

(kann eingeschränkt werden)

Newton Verfahren:

Linearisiere  $f$  um Näherung  $\bar{x} - \delta < x_n < \bar{x} + \delta$  :

$$f(x) \approx f(x_n) + f'(x_n)(x - x_n)$$

Löse  $f(x_n) + f'(x_n)(x - x_n) = 0$  nach  $x$  auf.

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

Folgerungen aus Satz 2.9

-Konvergenz, falls  $|x_n - \bar{x}| < \delta, \delta$  klein genug.

-asymptotisch sehr schnelle Konvergenz

$$|x_{n+1} - \bar{x}| \leq K|x_n - \bar{x}|^2 = \underbrace{K}_{\rightarrow 0}|x_n - \bar{x}||x_n - \bar{x}|$$

$$\limsup_{n \rightarrow \infty} \frac{|x_{n+1} - \bar{x}|}{|x_n - \bar{x}|^2} \leq K \text{ " (lokal) quadratische Konvergenz"}$$

### Satz 2.10

Ist  $f : [a, b] \rightarrow \mathbb{R}$  zweimal stetig differenzierbar,  $f(a) < 0, f(b) > 0$

und  $f'(x) > 0$  und  $f''(x) \geq 0$  für alle  $a \leq x \leq b$ , so konvergiert

das Newton Verfahren für Startwerte  $\bar{x} \leq x_0 \leq b$  monoton gegen die

einzig Nullstelle  $\bar{x}$  von  $f$  in  $[a, b]$ .

Beweis: Man beweise, dass für die Schrittfunktion  $f(x) := x - \frac{f(x)}{f'(x)}$

für  $\bar{x} \leq x \leq b$  die Abschränkung  $\bar{x} \leq \varphi(x) \leq x$  gilt.

Die Folge der Iterierten konvergiert dann monoton gegen eine

(die!) Nullstelle von  $f$ .

□

Frage: Verhalten bei mehrfachen Nullstellen?

### Beispiel 2.11

Die Funktion  $f(x) = (\sin x)^3$  hat an der Stelle  $\bar{x} = \pi = 3,14159265358\dots$

eine dreifache Nullstelle. Die Iterationsvorschrift des Newton-Verfahrens

lautet hier:

$$x_{n+1} = x_n - \frac{(\sin x_n)^3}{3(\sin x_n)^2 \cos x_n} = x_n - \frac{1}{3} \frac{\sin(x_n)}{\cos(x_n)}$$

Beobachtung: Das Newtonverfahren konvergiert für mehrfache Nullstellen nicht mehr quadratisch.

Modifikation: Wende das Newtonverfahren auf die transformierte Funktion  $F(x) = \frac{f(x)}{f'(x)}$  aus Satz 2.4 an:

$$x_{n+1} = x_n - \frac{F(x_n)}{F'(x_n)} = x_n - m(x_n) \frac{f(x_n)}{f'(x_n)}$$

$$\frac{1}{m(x)} = 1 - \frac{f(x)f''(x)}{(f'(x))^2}$$

### Beispiel 2.12:

Für die Funktion  $f(x) = (\sin x)^3$  erhält man  $x_{n+1} = x_n - \frac{\tan x_n}{1 + \tan x_n}$

### Satz 2.13

Der Faktor  $m(x_n)$  strebt gegen die Vielfachheit der Nullstelle  $\bar{x}$ .

Beweis: Sei  $j$  die Vielfachheit der Nullstelle. Dann folgt aus der Darstellung

$$f(x) = (x - \bar{x})^j g(x), g(\bar{x}) \neq 0, \text{ wegen}$$

$$f'(x) = j(x - \bar{x})^{j-1} g(x) + (x - \bar{x})^j g'(x) = (x - \bar{x})^{j-1} [j g(x) + (x - \bar{x}) g'(x)]$$

$$f''(x) = (x - \bar{x})^{j-2} [j(j-1)g(x) + j(x - \bar{x})g'(x) + (x - \bar{x})^2 g''(x)]$$

$$\lim_{x \rightarrow \bar{x}} m(x) = \lim_{x \rightarrow \bar{x}} \left( 1 - \frac{f(x)f''(x)}{(f'(x))^2} \right)^{-1} = \left( 1 - \frac{g(\bar{x})j(j-1)g(\bar{x})}{(j g(\bar{x}))^2} \right)^{-1} = j$$

□

## 2.3 Allgemeine Iterationsverfahren

Verfahren der Form  $x_{n+1} = \varphi(x_n)$ ,  $n = 0, 1, 2, \dots$   $x_0$  gegeben zur Brechnung eines Fixpunkts  $\bar{x} = \varphi(\bar{x})$ .

allgemeine Konvergenzaussage

Banachscher Fixpunktsatz

Lässt sich unter sehr viel allgemeineren Voraussetzungen beweisen!

### Satz 2.14 (Banachscher Fixpunktsatz)

Die Funktion  $\varphi : [a, b] \rightarrow \mathbb{R}$  bilde dass Intervall  $[a, b]$  in sich ab.

Sie sei Lipschitz-stetig auf  $[a, b]$ , für alle  $x, y \in [a, b]$  sei also

$$|\varphi(x) - \varphi(y)| \leq L|x - y| \text{ mit einer festen Konstanten } L. \text{ Es sei } L < 1.$$

Dann besitzt  $\varphi$  genau einen Fixpunkt  $\bar{x}$  in  $[a, b]$ . Für alle  $x \in [a, b]$  strebt die durch  $x_{n+1} = \varphi(x_n)$ ,  $n = 0, 1, 2, \dots$  definierte Folge gegen  $\bar{x}$ . Es gelten die Abschätzungen  $|x_n - x| \leq \frac{L}{1-L}|x_n - x_{n-1}|$ ,  
 $|x_n - \bar{x}| \leq \frac{L^n}{1-L}|x_1 - x_0|$ .

Beweis:

Schritt 1: Für alle  $k$  ist  $|x_{k+1} - x_k| \leq L^k|x_1 - x_0|$ , denn:

$k = 0$  klar,  $k \rightarrow k + 1$ :

$$|x_{k+2} - x_{k+1}| = |\varphi(x_{k+1}) - \varphi(x_k)| \leq L|x_{k+1} - x_k| \leq L L^k|x_1 - x_0|$$

Schritt 2: Für alle  $n$  und  $k$  ist  $|x_{n+k} - x_n| \leq \frac{L^n}{1-L}|x_1 - x_0|$ , denn:

$$\begin{aligned} |x_{n+k} - x_n| &= \left| \sum_{j=0}^{k-1} (x_{n+j+1} - x_{n+j}) \right| \leq \sum_{j=0}^{k-1} |x_{n+j+1} - x_{n+j}| \\ &\leq \sum_{j=0}^{k-1} L^{n+j}(x_1 - x_0) = L^n \left( \sum_{j=0}^{k-1} L^j \right) (x_1 - x_0) \\ &\leq L^n \left( \sum_{j=0}^{\infty} L^j \right) (x_1 - x_0) = L^n \frac{1}{1-L} (x_1 - x_0) \end{aligned}$$

Schritt 3: Die Folge der  $x_n$  strebt gegen einen Fixpunkt  $\bar{x}$  von  $\varphi$ , denn:

Nach Schritt 2 ist sie eine Cauchyfolge in  $\mathbb{R}$ . Sie konvergiert als solche gegen ein  $\bar{x} \in \mathbb{R}$ . Da die  $x_n$  in  $[a, b]$  liegen und  $[a, b]$  abgeschlossen ist liegt der Grenzwert  $\bar{x}$  in  $[a, b]$ .

Die  $\varphi$  als Lipschitz-stetige Funktion stetig ist, ist

$$\bar{x} = \lim_{n \rightarrow \infty} x_n = \lim_{n \rightarrow \infty} x_{n+1} = \lim_{n \rightarrow \infty} \varphi(x_n) = \varphi(\lim_{n \rightarrow \infty} x_n) = \varphi(\bar{x}),$$

$\bar{x}$  ist also Fixpunkt von  $\varphi$ .

Schritt 4: Für alle  $n$  ist  $|x_n - \bar{x}| \leq \frac{L^n}{1-L}|x_1 - x_0|$ ,

$$|x_n - \bar{x}| \leq \frac{L}{1-L}|x_n - x_{n-1}|$$

denn: Die erste Abschätzung erhält man, wenn man in der Abschätzung aus Schritt 2  $k$  gegen Unendlich streben lässt.

Die zweite Abschätzung erhält man aus den ersten, wenn man mit  $x_{n-1}$  statt  $x_0$  startet.

Schritt 5: Es gibt nur einen Fixpunkt, denn:

$$\text{Ist } \bar{x} = \varphi(\bar{x}), \bar{\bar{x}} = \varphi(\bar{\bar{x}}), \text{ so ist } |\bar{x} - \bar{\bar{x}}| = |\varphi(\bar{x}) - \varphi(\bar{\bar{x}})| \leq L|\bar{x} - \bar{\bar{x}}|$$

was wegen  $L < 1$  nur für  $|\bar{x} - \bar{\bar{x}}| = 0$  möglich ist.

□

Bemerkung: Ist  $\varphi$  stetig differenzierbar, und  $L = \max_{a \leq x \leq b} |\varphi'(x)|$

so gilt nach dem Mittelwertsatz für alle  $x, y \in [a, b]$

$$|\varphi(x) - \varphi(y)| = |\varphi'(\eta)(x - y)| \leq L|x - y|.$$

### Satz 2.15

Die Funktion  $\varphi$  sei auf einer Umgebung ihres Fixpunkts  $\bar{x}$   $q$ -mal stetig differenzierbar. Es sei  $q \geq 2$  und  $\varphi^{(k)}(\bar{x}) = 0, k = 1, \dots, q - 1$ . Dann ist die durch  $x_{n+1} = \varphi(x_n), n = 0, 1, 2, \dots$  gegebene Folge für alle genügend nahe an  $\bar{x}$  gelegenen Startwerte wohldefiniert und konvergent gegen  $\bar{x}$ . Überdies ist  $|x_{n+1} - \bar{x}| \leq c|x_n - \bar{x}|^q$  mit einer festen Konstante  $c$ .

Beweis: Nach dem Taylorschen Satz ist

$$\begin{aligned} \varphi(x) &= \sum_{k=1}^{n-1} \frac{1}{k!} \varphi^{(k)}(\bar{x})(x - \bar{x})^k + \frac{1}{q!} \varphi^{(q)}(\eta)(x - \bar{x})^q \\ &= \bar{x} + \frac{1}{q!} \varphi^{(q)}(\eta)(x - \bar{x})^q \end{aligned}$$

und damit für alle  $x$  aus einer festen  $\delta$ -Umgebung um  $\bar{x}$

$$|\varphi(x) - \bar{x}| \leq c|x - \bar{x}|^q, c = \max_{|\eta - \bar{x}| \leq \delta} \left| \frac{1}{q!} \varphi^{(q)}(\eta) \right|$$

Daraus folgt für alle  $x$  mit

$$|x - \bar{x}| \leq \min\left(\delta, \left(\frac{1}{2c}\right)^{\frac{1}{q-1}}\right) =: \delta^*$$

die Abschätzung

$$|\varphi(x) - \bar{x}| \leq \underbrace{c|x - \bar{x}|^{q-1}}_{\leq c\frac{1}{2c}} |x - \bar{x}| \leq \frac{1}{2}|x - \bar{x}|$$

Damit bildet  $\varphi$  die  $\delta^*$ -Umgebung von  $\bar{x}$  in sich ab. Für  $|x_0 - \bar{x}| \leq \delta^*$

sind die  $x_n$  wohldefiniert und es ist  $|x_{n+1} - \bar{x}| \leq \left(\frac{1}{2}\right)^n |x_1 - \bar{x}|$ .

Die  $x_n$  konvergieren also gegen  $\bar{x}$  und es ist

$$|x_{n+1} - \bar{x}| = |\varphi(x_n) - \bar{x}| \leq c|x_n - \bar{x}|^q$$

□

### Satz 2.16(modifiziertes Newtonverfahren)

Die Funktion sei auf einer Umgebung ihrer  $m$ -fachen Nullstelle hinreichend oft stetig differenzierbar. Dann konvergiert die durch

$$x_{n+1} = x_n - m \frac{f(x_n)}{f'(x_n)}, n = 0, 1, 2, \dots$$

gegebene Folge für alle hinreichend nahe bei  $\bar{x}$  gelegenen  $x_0$

quadratisch gegen  $\bar{x}$ .

Beweis: Untersuche  $\varphi(x) = x - m \frac{f(x)}{f'(x)}$

Wegen  $f(x) = (x - \bar{x})^m g(x), g(\bar{x}) \neq 0$ , gilt

$$\begin{aligned}\varphi(x) &= x - m \frac{(x-\bar{x})^m g(x)}{m(x-\bar{x})^{m-1}g(x) + (x-\bar{x})^m g'(x)} \\ &= x - (x-\bar{x})h(x), h(x) = \frac{mg(x)}{mg(x) + (x-\bar{x})g'(x)}.\end{aligned}$$

Es ist  $f(\bar{x}) = \bar{x}$  und wegen  $\varphi'(x) = 1 - h(x) - (x-\bar{x})h'(x)$

$$\varphi'(\bar{x}) = 1 - h(\bar{x}) = 0$$

□

## 2.4 Die Regula falsi

Gesucht:

schnell konvergierendes, ableitungsfreies Verfahren

Idee: Ersetze im Newtonverfahren

$$x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

die Ableitung  $f'(x_n)$  durch den Differenzenquotienten

$$\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$$

Dies ergibt die Regula falsi

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n)$$

### Beispiel 2.17

$$f(x) = x^5 - 3, x_0 = 1, x_1 = 2$$

$$x_{10} = 1,245730940$$

### Satz 2.18

Die Funktion  $f$  sei auf einer  $\delta$ -Umgebung ihrer Nullstelle  $\bar{x}$  zweimal stetig differenzierbar.

Für  $|x - \bar{x}| < \delta$  sei  $|f'(x)| \geq c_1 > 0$ ,  $|f''(x)| \leq c_2 \neq 0$ .

Die Umgebung  $U$  von  $\bar{x}$  sei definiert durch  $U = \left\{ x \mid |x - \bar{x}| \leq \min\left(\delta, \frac{c_1}{c_2}\right) \right\}$

Liegen dann  $x_0$  und  $x_1$  in  $U$ , so brich die durch

$$x_{n+1} = x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n), n = 1, 2, 3, \dots$$

entweder mit der exakten Lösung  $\bar{x}$  ab, oder ist für alle  $n$  definiert und konvergiert gegen  $\bar{x}$ . In beiden Fällen liegen alle  $x_n$  in  $U$ .

Es gibt eine Nullfolge  $(\Delta_n)$  mit  $|x_n - \bar{x}| \leq \frac{2c_1}{c_2} \Delta_n$  für alle berechenbaren

$$x_n \text{ und } \lim_{n \rightarrow \infty} \frac{\Delta_{n+1}}{\Delta_n} = 1, q = \frac{1+\sqrt{5}}{2} = 1,6180\dots$$

Beweis: Nach dem Hauptsatz der Differential- und Integralrechnung

ist wegen  $f(\bar{x}) = 0$

$$f(x) = (x - \bar{x})g(x), g(x) = \int_0^1 f'(\bar{x} + \theta(x - \bar{x}))d\theta$$

Ist nun  $x_{n-1}, x_n \in U$  und  $x_{n-1} \neq x_n$ , so ist  $x_{n+1}$  wegen

$$\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}} = f'(\eta) \neq 0 \text{ wohldefiniert und mit Zwischenstellen } \alpha_n \text{ und } \beta_n.$$

$$\begin{aligned} x_{n+1} - \bar{x} &= x_n - \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n) - \bar{x} \\ &= \frac{(x_n - \bar{x})(f(x_n) - f(x_{n-1})) - (x_n - x_{n-1})f(x_n)}{f(x_n) - f(x_{n-1})} \\ &= \frac{(x_{n-1} - \bar{x})f(x_n) - (x_n - \bar{x})f(x_{n-1})}{f(x_n) - f(x_{n-1})} \\ &= \frac{(x_{n-1} - \bar{x})(x_n - \bar{x})g(x_n) - (x_n - \bar{x})(x_{n-1} - \bar{x})g(x_{n-1})}{f(x_n) - f(x_{n-1})} \\ &= \frac{g(x_n) - g(x_{n-1})}{f(x_n) - f(x_{n-1})} (x_{n-1} - \bar{x})(x_n - \bar{x}) \\ &= \frac{g'(\alpha_n)(x_n - x_{n-1})}{f'(\beta_n)(x_n - x_{n-1})} (x_{n-1} - \bar{x})(x_n - \bar{x}) \\ &= \frac{g'(\alpha_n)}{f'(\beta_n)} (x_{n-1} - \bar{x})(x_n - \bar{x}) \end{aligned}$$

Wegen  $g'(x) = \int_0^1 f'(\bar{x} + \theta(x - \bar{x}))\theta d\theta$  ist

$$|g'(\alpha_n)| \leq \int_0^1 c_2 \theta d\theta = \frac{1}{2} c_2, \text{ so dass wir die Abschätzung}$$

$$|x_{n+1} - \bar{x}| \leq \frac{c_2}{2c_1} |x_{n-1} - \bar{x}| |x_n - \bar{x}| \quad (*)$$

erhalten. Wegen  $|x_{n-1} - \bar{x}| \leq \frac{c_1}{c_2}$  und  $x_n \in U$  folgt aus (\*)

$$|x_{n+1} - \bar{x}| \leq \frac{1}{2} |x_n - \bar{x}| \text{ und damit } x_{n+1} \in U. \text{ Ist } x_{n+1} \neq x_n,$$

so ist die Ausgangssituation wiederhergestellt und man kann

die Iteration fortsetzen. Ist hingegen  $x_{n+1} = x_n$ , so ist

$$0 = \frac{x_n - x_{n-1}}{f(x_n) - f(x_{n-1})} f(x_n) = \frac{1}{f'(\eta_n)} f(x_n), \text{ also } x_n \text{ Nullstelle,}$$

d.h.  $x_n = \bar{x}$ . Das Verfahren bricht also mit der exakten Nullstelle ab.

Um die Konvergenz nachzuweisen, setzen wir  $\Delta_0 = \frac{c_2}{2c_1} |x_0 - \bar{x}|$ ,

$$\Delta_1 = \frac{c_2}{2c_1} |x_1 - \bar{x}| \text{ und } \Delta_{n+1} = \Delta_{n-1} \Delta_n, n = 1, 2, 3, \dots$$

Durch Induktion folgt dann aus (\*)  $|x_n - \bar{x}| \leq \frac{2c_1}{c_2} \Delta_n$ .

Wegen  $x_0, x_1 \in U$  ist  $\Delta_0 \leq \frac{1}{2}, \Delta_1 \leq \frac{1}{2}$ . Durch Induktion folgt

erst  $\Delta_n \leq \frac{1}{2}$  und dann  $\Delta_n \leq \left(\frac{1}{2}\right)^n$ . Damit ist die Konvergenz

der  $x_n$  gegen  $\bar{x}$  bewiesen.

Wir setzen nun  $\delta_n = \ln \Delta_n$ . Wegen  $\Delta_{n+1} = \Delta_{n-1} \Delta_n$  ist dann

$$\delta_{n+1} = \delta_n + \delta_{n-1}, n = 1, 2, 3, \dots$$

Die allgemeine Lösung dieser Differenzgleichung ist



$$\delta_n = a \left( \frac{1+\sqrt{5}}{2} \right)^n + b \left( \frac{1-\sqrt{5}}{2} \right)^n$$

mit reellen Konstanten  $a$  und  $b$ . Diese Lösung findet zudem mit dem Ansatz  $\delta_n = \lambda^n$ . Wegen

$$\left| \frac{1-\sqrt{5}}{2} \right| < 1 \text{ folgt daraus, } \lim_{n \rightarrow \infty} \left( \delta_{n+1} - \frac{1+\sqrt{5}}{2} \delta_n \right) = 0$$

Setzt man  $q = \frac{1+\sqrt{5}}{2}$ , ist dies gleichbedeutend mit

$$\lim_{n \rightarrow \infty} \frac{\Delta_{n+1}}{\Delta_n^q} = 1$$

□

### Aufwand Newton Verfahren

2 Funktionsauswertungen pro Schritt

### Aufwand Regula falsi

1 Funktionsauswertung pro Schritt

Doppelschritt:

$$\lim_{n \rightarrow \infty} \frac{\Delta_{n+2}}{\Delta_n^2} = \lim_{n \rightarrow \infty} \frac{\Delta_{n+2}}{\Delta_{n+1}^q} \left( \frac{\Delta_{n+1}}{\Delta_n} \right)^2 = 1$$

effektive Konvergenzordnung

$$\left( \frac{1+\sqrt{5}}{2} \right)^2 = 2,6180\dots$$

Die Regula falsi ist asymptotisch schneller als das Newtonverfahren!

## 3. Lineare Gleichungssysteme

### 3.1 Normen

Normen dienen dazu, die Länge von Vektoren und damit indirekt den Abstand zwischen Vektoren zu messen.

#### **Definition 3.1:**

Eine reellwertige Funktion  $\|\cdot\|$  auf dem  $\mathbb{R}^n$  heißt Norm auf dem  $\mathbb{R}^n$ , falls für alle  $x, y \in \mathbb{R}^n$  und alle  $\alpha \in \mathbb{R}$  gilt:

- (i)  $\|x\| \geq 0$ ,  $\|x\| = 0 \iff 0$
- (ii)  $\|\alpha x\| = |\alpha| \|x\|$  (Homogenität)
- (iii)  $\|x + y\| \leq \|x\| + \|y\|$  (Dreiecksungleichung)

Die Größe  $\|x - y\|$  ist dann der Abstand von  $x$  und  $y$  bezüglich dieser Norm.

**Beispiel 3.2:**

Sei  $x \in (x_1, \dots, x_n) \in \mathbb{R}^n$ . Dann ist die euklidische Norm oder auch 2-Norm von  $x$  definiert durch:

$$\|x\| = \left(\sum_{i=1}^n x_i^2\right)^{\frac{1}{2}}$$

Sie induziert den euklidischen Abstand  $\|x - y\| = \left(\sum_{i=1}^n (x_i - y_i)^2\right)^{\frac{1}{2}}$

Mit dem Skalarprodukt  $\langle x, y \rangle = \sum_{i=1}^n x_i y_i$  gilt  $\|x\| = \langle x, x \rangle^{\frac{1}{2}}$ .

Für alle Vektoren  $x, y$  ist  $\langle x, y \rangle \leq \|x\| \|y\|$   
(Schwarzsche Ungleichung)

**Beispiel 3.3:**

Die Maximumnorm auf dem  $\mathbb{R}^n$  ist definiert durch  $\|x\|_{\infty} = \max_{i=1, \dots, n} |x_i|$

**Beispiel 3.4:**

Die 1-Norm auf dem  $\mathbb{R}^n$  ist gegeben durch  $\|x\|_1 = \sum_{i=1}^n |x_i|$ .

Bemerkung: Es gibt viele andere (und wesentlich kompliziertere) Normen auf dem  $\mathbb{R}^n$ .

So ist z.B. für jede symmetrische positiv definite  $(n \times n)$ -Matrix  $Q$  durch

$$\langle x, y \rangle = x^T Q y$$

ein Skalarprodukt auf dem  $\mathbb{R}^n$  gegeben, das die Norm  $\|x\| = \langle x, x \rangle^{\frac{1}{2}}$  induziert.

**Satz 3.5:**

Sind  $\|\cdot\|_1$  und  $\|\cdot\|_2$  zwei beliebige (!) Normen auf dem  $\mathbb{R}^n$  so gibt es Konstanten  $c_1, c_2 > 0$  mit  $c_1 \|x\|_1 \leq \|x\|_2 \leq c_2 \|x\|_1$  für alle  $x \in \mathbb{R}^n$ .

Beweis: Es genügt die behauptung für die Maximumnorm und eine weitere beliebige Norm  $\|\cdot\|$  auf dem  $\mathbb{R}^n$  zu zeigen.

Wir zeigen zunächst die Ungleichung  $\|x\| \leq c \|x\|_{\infty}$ . seien  $e_1, \dots, e_n$  die Einheitsvektoren. Nach der Dreiecksungleichung ist dann

$$\text{für alle } x \in \mathbb{R}^n \quad \|x\| = \left\| \sum_{i=1}^n x_i e_i \right\| \leq \sum_{i=1}^n \|x_i e_i\| = \sum_{i=1}^n |x_i| \|e_i\|,$$

also  $\|x\| \leq c\|x\|_\infty$  mit  $c = \sum_{i=1}^n \|e_i\|$ .

Zum Beweis der Umkehrung betrachten wir zunächst die Funktion

$f : \mathbb{R}^n \rightarrow \mathbb{R} : x \mapsto \|x\|$ . Für alle  $x, y \in \mathbb{R}^n$  ist

$\|x\| = \|(x - y) + y\| \leq \|x - y\| + \|y\|$ , und damit  $\|x\| - \|y\| \leq \|x - y\|$ .

Entsprechend ist  $\|y\| - \|x\| \leq \|x - y\|$ , also  $|\|x\| - \|y\|| \leq \|x - y\|$

Damit ist für alle  $x, y \in \mathbb{R}^n$   $|f(x) - f(y)| = |\|x\| - \|y\|| \leq \|x - y\| \leq c\|x - y\|_\infty$ ,

so dass  $f$  stetig ist.

Die Menge  $K = \{x \in \mathbb{R}^n \mid \|x\|_\infty = 1\}$  ist kompakt. Daher gibt es ein

$x^* \in K$  mit  $f(x^*) \leq f(x)$  für alle  $x \in K$ . Dies bedeutet  $f(x^*) \leq \left\| \frac{1}{\|x^*\|_\infty} x^* \right\|$   
 $= \frac{1}{\|x^*\|_\infty} \|x^*\|$  oder mit  $c' = f(x^*) = \|x^*\| > 0$ ,  $c_1 \|x\|_\infty \leq \|x\|$  für alle  $x \neq 0$

□

Bemerkung: Praktisch ist diese Aussage oft von nur sehr geringem Wert, da der Satz keine Information über die Größe der Konstanten liefert und die Konstanten in der oberen und unteren Abschätzung leicht um Größenordnungen differieren können.

### Beispiel 3.6:

Für alle  $x \in \mathbb{R}^n$  ist  $\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1$  mit der Maximumnorm  $\|x\|_\infty$  aus Beispiel 3.3, der 1-Norm aus Beispiel 3.4 und der euklidischen Norm  $\|x\|_2$ .

Umgekehrt gilt:

$$\|x\|_2 \leq \sqrt{n} \|x\|_\infty$$

$$\|x\|_1 \leq \sqrt{n} \|x\|_2$$

$$\|x\|_1 \leq n \|x\|_\infty$$

Alle Konstanten sind optimal.

### Definition 3.7:

Einer gegebenen Vektornorm  $\|\cdot\|$  ist die Matrixnorm  $\|A\| := \max_{x \neq 0} \frac{\|Ax\|}{\|x\|}$

auf dem Raum der reellen  $(n \times x)$ -Matrizen zugeordnet:

Für alle  $x \in \mathbb{R}^n$  ist damit  $\|Ax\| \leq \|A\| \|x\|$ , und es gibt  $x^* \neq 0$  mit

$$\|Ax^*\| = \|A\| \|x^*\|.$$

Eigenschaften:

$$(i) \|A\| \geq 0, \|A\| = 0 \iff A = 0$$

$$(ii) \|\alpha A\| = |\alpha| \|A\|$$

$$(iii) \|A + B\| \leq \|A\| + \|B\|$$

(deshalb Norm)

$$(iv) \|AB\| \leq \|A\| \|B\|$$

invertierbare Matrizen:

$$\|x\| = \|A^{-1}Ax\| \leq \|A^{-1}\| \|Ax\|, \text{ daher:}$$

$$\frac{1}{\|A^{-1}\|} \|x\| \leq \|Ax\| \leq \|A\| \|x\|$$

Existiert  $A^{-1}$ , so ist

$$\frac{1}{\|A^{-1}\|} = \min_{x \neq 0} \frac{\|Ax\|}{\|x\|}$$

### Satz 3.8

Der Maximumnorm aus Beispiel 3.3 entspricht die Matrixnorm

$$\|A\|_{\infty} = \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \text{ der Matrix } A = (a_{ij}), \text{ die Zeilensummennorm.}$$

Die der 1-Norm aus Beispiel 3.4 zugeordnete Norm ist die

$$\text{Spaltensummennorm } \|A\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^n |a_{ij}|$$

Beweis: (für die Maximumnorm)

$$\text{Für alle } x \in \mathbb{R}^n \text{ ist } \|Ax\|_{\infty} = \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}x_j|$$

$$\leq \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \|x\|_{\infty}$$

Damit ist  $\|A\|_{\infty} \leq \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}|$ . Um zu zeigen, dass die Abschätzung

$$\text{scharf ist, nehmen wir an, dass } \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| = \sum_{j=1}^n |a_{kj}|$$

ist und definieren wir  $x^*$  Komponentenweise durch

$$x_j^* = \begin{cases} 1, & \text{falls } a_{kj} \geq 0 \\ -1, & \text{falls } a_{kj} < 0 \end{cases}$$

$$\text{Dann ist } \|Ax^*\|_{\infty} \geq \left| \sum_{j=1}^n a_{kj} x_j^* \right| = \sum_{j=1}^n |a_{kj}| = \left( \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}| \right) \|x_j^*\|_{\infty}$$

$$\text{also } \|A\|_{\infty} \geq \max_{i=1, \dots, n} \sum_{j=1}^n |a_{ij}|.$$

□

**Beispiel 3.9:**

Die Matrix  $A = \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix}$  hat die Zeilensummennorm:

$$\|A\|_\infty = \max\{|1| + |2|, |3| + |4|\} = 7$$

Ihre Inverse  $A^{-1} = \begin{pmatrix} -2 & 1 \\ \frac{3}{2} & -\frac{1}{2} \end{pmatrix}$

hat die Zeilensummennorm

$$\|A^{-1}\|_\infty = \max\{|-2| + |1|, |\frac{3}{2}| + |-\frac{1}{2}|\} = 3$$

**Satz 3.10**

Der euklidischen Norm als Vektornorm ist die Spektralnorm als Matrixnorm zugeordnet. Die Spektralnorm der Matrix  $A$  ist die Wurzel aus dem größten Eigenwert der Matrix  $A^\top A$ .

Beweis: Wegen  $(A^\top A)^\top = A^\top (A^\top)^\top = A^\top A$  ist  $A^\top A$  eine symmetrische Matrix.

Es gibt daher nach dem Spektralsatz der Linearen Algebra eine Orthonormalbasis  $x_1, \dots, x_n$  des  $\mathbb{R}^n$  mit  $A^\top A x_i := \lambda_i x_i$ .

Es ist  $\lambda_i \langle x_i, x_i \rangle = \langle x_i, \lambda_i x_i \rangle = \langle x_i, A^\top A x_i \rangle = \langle A x_i, A x_i \rangle = \|A x_i\|_2^2 \geq 0$

Für den Vektor  $x = \sum_{i=1}^n \alpha_i x_i$  gilt

$$\|A x\|_2^2 = \langle A x, A x \rangle = \langle x, A^\top A x \rangle = \left\langle \sum_{i=1}^n \alpha_i x_i, \sum_{j=1}^n \alpha_j \lambda_j x_j \right\rangle$$

$$= \sum_{i=1}^n \lambda_i \alpha_i^2 \leq \left( \max_{i=1, \dots, n} \lambda_i \right) \sum_{i=1}^n \alpha_i^2 = \left( \max_{i=1, \dots, n} \lambda_i \right) \|x\|_2^2$$

Damit ist  $\|A\|_2 \leq \left( \max_{i=1, \dots, n} \lambda_i \right)^{\frac{1}{2}}$

Da für  $x = x_i$

$\|A x\|_2^2 = \lambda_i = \lambda_i \|x_i\|_2^2$  ist, ist  $\|A\|_2 \geq \left( \max_{i=1, \dots, n} \lambda_i \right)^{\frac{1}{2}}$ , also

$$\|A\|_2 = \left( \max_{i=1, \dots, n} \lambda_i \right)^{\frac{1}{2}}$$

□

Bemerkung: Die Spektralnorm von  $A^{-1}$  ist  $\|A^{-1}\|_2 = \frac{1}{\left( \min_{i=1, \dots, n} \lambda_i \right)^{1/2}}$ .

**Beispiel 3.11:**

Für die Matrix aus Beispiel 3.9 ist

$$A^T A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} 10 & 14 \\ 14 & 20 \end{pmatrix}$$

Die Matrix  $A^T A$  hat das charakteristische Polynom

$$\begin{vmatrix} 10 - \lambda & 14 \\ 14 & 20 - \lambda \end{vmatrix} = (10 - \lambda)(20 - \lambda) - 14^2 = \lambda^2 - 30\lambda + 4$$

Die Eigenwerte der Matrix  $A^T A$  sind daher

$$\lambda_{1,2} = 15 \pm \sqrt{15^2 - 4} = 15 \pm \sqrt{221}$$

Damit:

$$\|A\|_2 = \sqrt{15 + \sqrt{221}}$$

**3.2 Die Kondition von Matrizen**

Ziel: Ein Maß für die Empfindlichkeit eines linearen Gleichungssystems gegenüber Störungen der Koeffizientenmatrix und der rechten Seite.

Bezeichnungen: Sei  $\hat{x}$  eine Näherungslösung des linearen Gleichungssystems  $Ax = b$  mit der exakten Lösung  $x$ .

Der Fehler  $e$  und das Residuum  $r$  sind dann

$$e = x - \hat{x}, r = b - A\hat{x}$$

**Satz 3.12**

Es ist  $\frac{1}{\|A\|\|A^{-1}\|} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \|A\|\|A^{-1}\| \frac{\|r\|}{\|b\|}$

Diese Abschätzungen sind scharf in dem Sinn, dass sie für bestimmte  $b$  und  $\hat{x}$  angenommen werden.

Beweis: Wegen  $\|r\| = \|Ae\| \leq \|A\|\|e\|$  ist  $\frac{\|r\|}{\|A\|} \leq \|e\|$ .

Wegen  $\|x\| = \|A^{-1}b\| \leq \|A^{-1}\|\|b\|$  ist  $\frac{1}{\|A^{-1}\|\|b\|} \leq \frac{1}{\|x\|}$ .

Mit diesen beiden Abschätzungen folgt die weitere Abschätzung

$$\frac{1}{\|A\|\|A^{-1}\|} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|}.$$

Die Abschätzung wird angenommen, wenn

$$\|r\| = \|Ae\| = \|A\|\|e\|, \|x\| = \|A^{-1}b\| = \|A^{-1}\|\|b\|$$

Nach Definition von  $\|A\|$  und  $\|A^{-1}\|$  gibt es derartige Vektoren

$b^*$  und  $e^*$ . Für dieses  $b^*$  und  $\hat{x}^* = A^{-1}b^* - e^*$  wird die Abschätzung

angenommen; die Norm des fehlers  $e^*$  kan dabei beliebig gewählt werden.

Wegen  $\|b\| = \|Ax\| \leq \|A\|\|x\|$  ist  $\frac{1}{\|x\|} \leq \|A\| \frac{1}{\|b\|}$

Wegen  $\|e\| = \|A^{-1}r\| \leq \|A^{-1}\|\|r\|$  folgt  $\frac{\|e\|}{\|x\|} \leq \|A\|\|A^{-1}\| \frac{\|r\|}{\|b\|}$

Die Abschätzung wird angenommen, falls  $\|b\| = \|Ax\| = \|A\|\|x\|$ ,

$\|e\| = \|A^{-1}r\| = \|A^{-1}\|\|r\|$ .

Nach Definition der Matrixnorm gibt es derartige Vektoren  $x^*$  und  $r^*$ , wobei  $\|r^*\|$  als beliebig klein gewählt werden kann.

Für  $b^* = Ax^*$  und  $\hat{x} = x^* - A^{-1}r^*$  geht die Ungleichung wieder in eine Gleichung über.

□

### Definition 3.13

Die Größe  $\kappa(A) = \|A\|\|A^{-1}\|$  heißt Konditionszahl oder Kondition der invertierbaren  $(n \times n)$ -Matrix  $A$  bezgl. der gegebenen auf dem  $\mathbb{R}^n$

Kurzfassung Satz 3.12:

$$\frac{1}{\kappa(A)} \frac{\|r\|}{\|b\|} \leq \frac{\|e\|}{\|x\|} \leq \kappa(A) \frac{\|r\|}{\|b\|}$$

Bemerkung:

Für alle Matrizen  $A$  gilt wegen  $1 = \|I\| = \|AA^{-1}\| \leq \|A\|\|A^{-1}\|$

stets  $\kappa(A) \geq 1$ . Ist  $\kappa(A) \approx 1$ , spricht man von einer gut konditionierten Matrix, und ist  $\kappa(A) \gg 1$  ist  $A$  schlecht konditioniert.

### Satz 3.14:

Löst man statt des exakten Systems  $Ax = b$  das gestörte System  $\hat{A}\hat{x} = b$ , so gilt  $\frac{\|x - \hat{x}\|}{\|\hat{x}\|} \leq \kappa(A) \frac{\|A - \hat{A}\|}{\|A\|}$ .

Beweis: Es ist  $x = A^{-1}b = A^{-1}\hat{A}\hat{x} = A^{-1}(A + \hat{A} - A)\hat{x} = \hat{x} + A^{-1}(\hat{A} - A)\hat{x}$

und damit  $\|x - \hat{x}\| \leq \|A^{-1}\| \|A - \hat{A}\| \|\hat{x}\|$ , das heißt

$$\frac{\|x - \hat{x}\|}{\|\hat{x}\|} \leq \|A^{-1}\| \|A - \hat{A}\| \frac{\|\hat{x}\|}{\|\hat{x}\|} = \kappa(A) \frac{\|A - \hat{A}\|}{\|A\|}$$

Sei nun  $\|A^{-1}y\| = \|A^{-1}\|\|y\|$  und  $\hat{A} = A + \delta I, \hat{x} = \frac{1}{\delta}y$ .

Dann ist  $\|x - \hat{x}\| = \|A^{-1}(\hat{A} - A)\hat{x}\| = \|A^{-1}y\| = \|A^{-1}\|\|y\|$

$= \|A^{-1}\| \|A - \hat{A}\| \|\hat{x}\|$  so daß in diesem Fall Gleichheit gilt.

□

**Satz 3.15:**

Für alle invertierbaren Matrizen  $A$  ist  $\min\left\{\frac{\|A-B\|}{\|A\|} \mid B \text{ singular} \right\} \geq \frac{1}{\kappa(A)}$

Für die Maximumnorm, die 1-Norm und die euklidische Norm gilt Gleichheit.

Beweis: Ist  $B$  singular, geht es eine  $x \neq 0$  mit  $Bx = 0$ . Für dieses  $x$  ist

$$\begin{aligned} \|x\| &= \|A^{-1}Ax\| \leq \|A^{-1}\| \|Ax\| = \|A^{-1}\| \|(A-B)x\| \\ &\leq \|A^{-1}\| \|A-B\| \|x\| \text{ Daher ist } 1 \leq \|A^{-1}\| \|A-B\| \text{ und damit} \\ \frac{1}{\kappa(A)} &= \frac{1}{\|A\| \|A^{-1}\|} \leq \frac{\|A-B\|}{\|A\|} \end{aligned}$$

Um zu zeigen, dass für die Maximumnorm auch Gleichheit einreten kann, sei  $v$  ein Vektor mit  $\|v\| = 1$ ,  $\|A^{-1}v\| = \|A^{-1}\| \|v\|$ .

Weiter sei  $y = A^{-1}v$ ,  $\|y\| = |y_m|$ , sowie mit dem  $m$ -ten Einheitsvektor  $e_m$   $z = \frac{1}{y_m}e_m$ ,  $B = A - vz^T$

Dann ist  $B$  wegen  $y \neq 0$  und  $By = AA^{-1}v - \frac{1}{y_m}y_mv = 0$  singular.

Für alle Vektoren  $x$  gilt

$$\begin{aligned} \|(A-B)x\| &= \left\| \frac{1}{y_m}x_mv \right\| = \frac{1}{|y_m|} |x_m| \|r\| = \frac{1}{\|y_m\|} |x_m| = \frac{1}{\|A^{-1}v\|} |x_m| \\ &= \frac{1}{\|A^{-1}\| \|v\|} |x_m| = \frac{1}{\|A^{-1}\|} |x_m| \text{ und damit} \\ \|A-B\| &\leq \frac{1}{\|A^{-1}\|}, \text{ d.h. } \frac{\|A-B\|}{\|A\|} \leq \frac{1}{\kappa(A)}. \end{aligned}$$

□

**3.3 Die LR-Zerlegung einer Matrix**

Problem: Wie löst man lineare Gleichungssysteme der Form  $Ax = b$  mit einer nichtsingulären  $(n \times n)$ -Matrix  $A$  und bekannter rechter Seite  $b$ ?

naheliegende Methode: Gaußscher Algorithmus

Durch schrittweise Elimination der Unbekannten  $x_1, \dots, x_n$ .

**Beispiel 3.16**

$$\left( \begin{array}{ccc|c} 2 & 1 & 1 & 7 \\ 4 & 1 & 0 & 6 \\ -2 & 2 & 1 & 5 \end{array} \right)$$

$$\left( \begin{array}{ccc|c} & A & & b \end{array} \right)$$

Vielfaches der 1. Zeile von Zeile 2 und 3 abziehen, so dass in der 1. Spalte von Zeile 2 und 3 die Null steht.



$$\left( \begin{array}{ccc|c} 2 & 1 & 1 & 7 \\ 0 & -1 & -2 & -8 \\ 0 & 3 & 2 & 12 \end{array} \right)$$

geeignetes Vielfaches der 2. Zeile von 3. Zeile abziehen:

$$\left( \begin{array}{ccc|c} 2 & 1 & 1 & 7 \\ 0 & -1 & -2 & -8 \\ 0 & 0 & -4 & -12 \end{array} \right)$$

Ende Eliminationsphase.

Dreieckssystem auflösen:

$$-4x_3 = -12$$

$$-x_2 = -8 + 2x_3$$

$$2x_1 = 7 - x_2 - x_3$$

$$x_3 = 3$$

$$x_2 = 2$$

$$x_1 = 1$$

Bemerkung:

(1) Der Gaußsche Algorithmus ist ein starres Rechenschema.

(2) Dieses Rechenschema kann zusammenbrechen, ohne dass die Matrix singular ist.

### Beispiel 3.17

Die erste Gleichung des eindeutig lösbaren Systems

$$\left( \begin{array}{cc|c} 0 & 1 & 1 \\ 1 & 1 & 2 \end{array} \right)$$

kann nicht genutzt werden, um die erste Variable zu eliminieren!

Abhilfe: Gleichungen in systematischer Weise vertauschen:

Spaltenpivotsuche.

### Beispiel 3.18

Gegeben wieder das System

$$\left( \begin{array}{ccc|c} 2 & 1 & 1 & 7 \\ 4 & 1 & 0 & 6 \\ -2 & 2 & 1 & 5 \end{array} \right)$$

Zunächst wird die Gleichung mit dem betragsmäßig größten Quotienten vor der ersten Variablen mit der ersten Gleichung vertauscht.

$$\left( \begin{array}{ccc|c} 4 & 1 & 0 & 6 \\ 2 & 1 & 1 & 7 \\ -2 & 2 & 1 & 5 \end{array} \right)$$

Jetzt wird ein geeignetes Vielfaches der 1. Gleichung von der 2. bzw 3. Gleichung abgezogen und aus diesen Gleichungen so die 1. Variable eliminiert.

$$\left( \begin{array}{ccc|c} 4 & 1 & 0 & 6 \\ 0 & \frac{1}{2} & 1 & 4 \\ 0 & \frac{3}{2} & 1 & 8 \end{array} \right)$$

Nun wird die 3. mit der 2. Gleichung vertauscht, denn mit den verbliebenen Gleichungen hat die 3. Gleichung den größten Koeffizienten vor der 2. Variablen.

$$\left( \begin{array}{ccc|c} 4 & 1 & 0 & 6 \\ 0 & \frac{3}{2} & 1 & 8 \\ 0 & \frac{1}{2} & 1 & 4 \end{array} \right)$$

Eliminierung der 2. Variablen aus der 3. Gleichung ergibt

$$\left( \begin{array}{ccc|c} 4 & 1 & 0 & 6 \\ 0 & \frac{3}{2} & 1 & 8 \\ 0 & 0 & \frac{4}{5} & \frac{12}{5} \end{array} \right)$$

Damit ist die Eliminationsphase beendet. Auflösen des Dreieckssystems ergibt  $x_3 = 3, x_2 = 2, x_1 = 1$ .

□

## Gaußscher Algorithmus ohne Pivotsuche

Der Gaußsche Algorithmus formt das Ausgangssystem  $Ax = b$  oder  $A^{(0)}x = b^{(0)}$  schrittweise um in Systeme  $A^{(k)}x = b^{(k)}, k = 1, \dots, n-1$  mit unveränderter Lösung  $x$

Für die Einträge  $a_{ij}^{(k)}$  von  $A^{(k)}, k \geq 1$ , gilt  $a_{ij}^{(k)} = a_{ij}^{(k-1)}$  für  $i = 1, \dots, k$

$a_{ij}^{(k)} = 0$  für  $i = k+1, \dots, n, j = 1, \dots, k$ .

Die restlichen Einträge ergeben sich gemäß:

Für  $k = 1, \dots, n-1$  berechne

Für  $i = k+1, \dots, n$  berechne

Für  $j = k+1, \dots, n$  berechne

$$a_{ij}^{(k)} = a_{ij}^{(k-1)} - \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}} a_{kj}^{(k-1)}$$

$$b_i^{(k)} = b_i^{(k-1)} - \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}} b_k^{(k-1)}$$

Endergebnis ist das Dreieckssystem

$$\left( \begin{array}{cccc|c} a_{11}^{(n-1)} & a_{12}^{(n-1)} & \cdots & a_{1n}^{(n-1)} & b_1^{(n-1)} \\ 0 & a_{22}^{(n-1)} & & \vdots & b_2^{(n-1)} \\ \vdots & \ddots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & a_{nn}^{(n-1)} & b_n^{(n-1)} \end{array} \right)$$

Dieses Dreieckssystem lässt sich einfach auflösen:

$$x_n = \frac{1}{a_{nn}^{(n-1)}} b_n^{(n-1)}$$

Für  $i = n - 1, \dots, 1$  berechne

$$x_i = \frac{1}{a_{ii}^{(n-1)}} \left\{ b_i^{(n-1)} - \sum_{j=i+1}^n a_{ij}^{(n-1)} x_j \right\}$$

Bemerkungen:

1) Im Rechner können die Elemente  $a_{ij}^{(k-1)}$  der Matrix  $A^{(k-1)}$  mit den Elementen  $a_{ij}^{(k)}$  der Matrix  $A^{(k)}$  überschrieben werden.

2) Man kann die Eliminationsfaktoren  $l_i = \frac{a_{ik}^{(k-1)}}{a_{kk}^{(k-1)}}$  auf den Platz der mit

Null überschriebenen Matrixelemente  $a_{ik}^{(k-1)}$  abspeichern und die rechte Seite dann nachträglich bearbeiten.

Dies ist insbesondere vorteilhaft, wenn das Gleichungssystem für mehrere rechte Seiten gelöst werden soll.

3) Die Eliminationsfaktoren  $l_{ik}$  sollten außerhalb der innersten Schleife berechnet werden.

## Matrixdarstellung

Führt man die Eliminationsfaktoren

$$l_{ik} = \left. \begin{array}{l} a_{ik}^{(k-1)} \\ a_{kk}^{(k-1)} \end{array} \right\} \begin{array}{l} k = 1, 2, \dots, n-1 \\ i = k+1, \dots, n \end{array}$$

und die Eliminationsmatrizen

$$G_k|_{ij} = \begin{cases} l_{ik}, & j = k, i = k+1, \dots, n \\ 0, & \text{sonst} \end{cases}$$

ein, so ist

$$A^{(k)} = (I - G_k)A^{(k-1)}, b^{(k)} = (I - G_k)b^{(k-1)}$$

Beweis: Da man die  $k$ -te Spalte von  $G_k$  von Null verschieden ist, gilt

$$G_k A^{(k-1)}|_{ij} = \sum_{m=1}^n G_k|_{im} A^{(k-1)}|_{mj} = G_k|_{ik} A^{(k-1)}|_{kj}$$

Für  $i = 1, \dots, k$  ist  $G_k|_{ik} = 0$  und daher  $G_k A^{(k-1)}|_{ij}, j = 1, \dots, n$

Für  $i = k+1, \dots, n$  ist  $G_k A^{(k-1)}|_{ij} = l_{ik} a_{kj}^{(k-1)}, j = 1, \dots, n$

□

**Satz 3.19**

Ist der Gaußsche Algorithmus für die  $(n \times n)$ -Matrix  $A$  durchführbar und ergeben sich dabei die Eliminationsmatrizen  $G_1, \dots, G_{n-1}$

so gilt: (i)  $(I - G_{n-1})(I - G_{n-2}) \cdots (I - G_1)A = R$

(ii)  $G_k G_l = 0, k \leq l$

(iii)  $(I - G_k)^{-1} = I + G_k$

(iv)  $((I - G_{n-1}) \cdots (I - G_1))^{-1} = I + \sum_{k=1}^{n-1} G_k = I + L$

Dabei ist  $I$  die  $(n \times n)$ -Einheitsmatrix,

$$R = \begin{pmatrix} a_{11}^{(n-1)} & \cdots & a_{1n}^{(n-1)} \\ \vdots & \ddots & \vdots \\ 0 & \cdots & a_{nn}^{(n-1)} \end{pmatrix}$$

eine obere Dreiecksmatrix und

$$L = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ l_{21} & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ l_{n1} & \cdots & l_{n,n-1} & 1 \end{pmatrix}$$

Beweis:

(i) Nach Konstruktion ist  $A^{(k)} = (I - G_k)A^{(k-1)}$  (siehe oben)

(ii) Für alle Vektoren  $x$  mit  $x|_k = 0$  ist  $G_k x = 0$

(iii)  $(I + G_k)(I - G_k) = I - G_k + G_k - G_k G_k = I$

$(I - G_k)(I + G_k) = I + G_k - G_k - G_k G_k = I$

(iv)  $((I - G_{n-1}) \cdots (I - G_1))^{-1}$

$= (I - G_1)^{-1} \cdots (I - G_{n-1})^{-1}$

$= (I + G_1) \cdots (I + G_{n-1})$  wegen (ii)

$= I + \sum_{k=1}^{n-1} G_k$  wegen (ii)

□

**Folgerung:**

Der Gaußsche Algorithmus liefert bei exakter Rechnung eine  $LR$ -Zerlegung

$$A = LR$$

Für Matrix  $A$  in eine untere  $\Delta$ -Matrix mit Einsen in der Diagonale und eine obere Dreiecksmatrix  $R$ , die  $LR$ -Zerlegung von  $A$ .

Das Gleichungssystem  $Ax = b$  lässt sich dann wie folgt lösen:

- 1) Löse  $Ly = b$
- 2) Löse  $Rx = y$

### Spaltenpivotsuche

Führt man eine Spaltenpivotsuche durch, so liefert der Gaußsche Algorithmus eine Zerlegung  $LR = PA$  einer Matrix  $PA$ , deren Zeilen gegenüber  $A$  vertauscht sind.

Begründung:

Der Gaußsche Algorithmus mit Spaltenpivotsuche liefert Eliminationsmatrizen  $G_1, \dots, G_{n-1}$  Permutationsmatrizen,  $P_1, \dots, P_{n-1}$  und eine obere Dreiecksmatrix  $R$  mit

$$(I - G_{n-1})P_{n-1} \cdots (I - G_1)P_1 A = R$$

Die Permutationsmatrix  $P_i$  vertauscht die  $i$ -te Komponente eines Vektors mit einer Komponente  $k \geq i$ . Daher ist  $P_i^{-1} = P_i$  und wegen

$$(I - G_i)^{-1} = I + G_i, A = P_i(I + G_1) \cdots P_{n-1}(I + G_{n-1})R$$

Sei nun  $Q_k = P_{n-1} \cdots P_k, P = P_{n-1} \cdots P_1$

Wegen  $P_k P_k = I$  gilt dann

$$\begin{aligned} Q_k P_k (I + G_k) &= Q_{k+1} P_k P_k (I + G_k) = Q_{k+1} (I + G_k) \\ &= (I + Q_{k+1} G_k Q_{k+1}^{-1}) Q_{k+1} \end{aligned}$$

Da die Spalten  $k+1, \dots, n$  von  $G_k$  ausschließlich mit Nullen besetzt sind, gilt

$$G_k Q_{k+1}^{-1} = G_k P_{k+1} \cdots P_{n-1} = G_k \text{ und damit}$$

$$Q_k P_k (I + G_k) = (I + Q_{k+1} G_k) Q_{k+1}$$

Die Matrix  $I + Q_{k+1} G_k$  ist eine Matrix mit gleicher Besetzungsstruktur wie  $I + G_k$ . Daher ist

$$\begin{aligned} L &:= P P_1 (I + G_1) \cdots P_{n-1} (I + G_{n-1}) \\ &= (I + Q_2 G_1) \cdots (I + Q_{n-1} G_{n-2}) (I + G_{n-1}) \\ &= I + Q_2 G_1 + \cdots + Q_{n-1} G_{n-2} + G_{n-1} \end{aligned}$$

eine untere Dreiecksmatrix mit Einsen in der Diagonale und

$$PA = LR$$

### Beispiel 3.20

ohne Spaltensuche

$$\begin{pmatrix} 2 & 1 & 1 \\ 4 & 1 & 0 \\ -2 & 2 & 1 \end{pmatrix} \rightarrow \left( \begin{array}{c|cc} 2 & 1 & 1 \\ \hline 4 & 1 & 0 \\ -2 & 2 & 1 \end{array} \right) \rightarrow \left( \begin{array}{c|cc} 2 & 1 & 1 \\ \hline 2 & -1 & -2 \\ -1 & -3 & -4 \end{array} \right)$$

$$\Rightarrow \begin{pmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ -1 & -3 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 \\ 0 & -1 & -2 \\ 0 & 0 & -4 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 1 \\ 4 & 1 & 0 \\ -2 & 2 & 1 \end{pmatrix}$$

mit Spaltensuche

$$\begin{pmatrix} 2 & 1 & 1 \\ 4 & 1 & 0 \\ -2 & 2 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 4 & 1 & 0 \\ 2 & 1 & 1 \\ -2 & 2 & 1 \end{pmatrix} \rightarrow \left( \begin{array}{ccc|ccc} 4 & 1 & 0 & & & \\ \hline \frac{1}{2} & & & \frac{1}{2} & & \\ -\frac{1}{2} & & & \frac{5}{2} & & 1 \end{array} \right)$$

$$\Rightarrow \left( \begin{array}{ccc|ccc} 4 & 1 & 0 & & & \\ \hline -\frac{1}{2} & & & \frac{5}{2} & & \\ \frac{1}{2} & & & \frac{1}{2} & & 1 \end{array} \right) \leftarrow \text{erste Spalte mitvertauschen!}$$

$$\rightarrow \left( \begin{array}{ccc|ccc} 4 & 1 & 0 & & & \\ \hline -\frac{1}{2} & & & \frac{5}{2} & & \\ \frac{1}{2} & & & \frac{1}{5} & & \frac{4}{5} \end{array} \right) \Rightarrow \begin{pmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ \frac{1}{2} & \frac{1}{5} & 0 \end{pmatrix} \begin{pmatrix} 4 & 1 & 0 \\ 0 & \frac{5}{2} & 1 \\ 0 & 0 & \frac{4}{5} \end{pmatrix} = \begin{pmatrix} 4 & 1 & 0 \\ -2 & 2 & 1 \\ 2 & 1 & 1 \end{pmatrix}$$

Aufwand für die Berechnung der LR-Zerlegung

Anzahl der Multiplikationen und Divisionen zur Elimination der  $k$ -ten Variablen:

$(n - k)$  Berechnung der  $l_{ik}, i = k + 1, \dots, n$

$(n - k)^2$  Berechnung der  $a_{ij}^{(k)}, i, j = k + 1, \dots, n$

Gesamtzahl:  $(k = n - k)$

$$\sum_{k=1}^{n-1} ((n - k) + (n - k)^2) = \sum_{k=1}^{n-1} (k + k^2) = \frac{1}{3}n^3 - \frac{1}{3}n$$

Aufwand für die Lösung der Dreieckssysteme

$$\sum_{k=2}^n (k - 1) \quad |Ly = b$$

$$+ \sum_{k=1}^{n-1} (n - k) + n \quad |Rx = y$$

Gesamt:  $n^2$

Aufwand für die Lösung von  $Ax = b$

$$\frac{1}{3}n^3 - \frac{1}{3}n + n^2$$

Aufwand für die Berechnung von  $A^{-1}$ :

$$\left(\frac{1}{3}n^3 - \frac{1}{3}n\right) + n \cdot n^3 = \frac{4}{3}n^3 - \frac{1}{3}n$$

Aufwand für die Lösung von  $k$  Systemen mit gleicher Matrix  $A$

$(\frac{1}{3}n^3 - \frac{1}{3}n) + kn^2$  mit  $LR$  Zerlegung

$(\frac{1}{3}n^3 - \frac{1}{3}n) + n \cdot n^2 + kn^2$  mit  $A^{-1}$

Merke: Man berechnet nie die Inverse einer Matrix, es sei denn man benötigt sie explizit.

### 3.4 Eine Fehleranalyse für den Gaußschen Algorithmus

Frage: Wie stark können Rundungsfehler das Ergebnis verfälschen?

Idee: Rückwärtsanalyse

Zeige, dass die berechnete Lösung  $\bar{x}$  des Systems  $Ax = b$  die

exakte Lösung eines leicht(?) gestörten Systems  $\hat{A}\hat{x} = \hat{b}$  ist.

Schätze die Größe der Störung  $\hat{A} - A, \hat{b} - b$  ab.

Ihr Einfluss hängt von der Kondition von  $A$  ab.

Bezeichnung: Ist  $A$  eine reelle Matrix mit den Matrixelementen  $a_{ij}$ , so

bezeichnet  $|A|$  die Matrix gleicher Dimension mit den Matrixelementen  $|a_{ij}|$ .

Bezeichnung:  $A \leq B \iff a_{ij} \leq b_{ij}, i, j = 1, \dots, n$

#### Satz 3.21

Ist der Gaußsche Algorithmus(ohne Pivotsuche) für die Matrix  $A$  durchführbar,

so liefert er eine  $LR$  Zerlegung  $\hat{L}\hat{R} = A + H$  einer gestörten Matrix  $A + H$ .

Die Störungsmatrix  $H$  genügt der komponentenweisen Abschätzung

$|H| \leq C(n)\varepsilon^*(|A| + |\hat{L}||\hat{R}|)$ , wobei die Konstante  $C(n)$  gegeben ist durch

$C(n) = \frac{(1+(2+\varepsilon^*)\varepsilon^*)^{n-1}}{\varepsilon^*} \approx 2n$  und  $\varepsilon^*$  die Maschinengenauigkeit bezeichnet.

#### Satz 3.22

Die berechnete Lösung der beiden Dreieckssysteme  $\hat{L}y = b, \hat{R}x = y$

sind die exakten Lösungen zweier gestörter Dreieckssysteme

$(\hat{L} + F)\hat{y} = b, (\hat{R} + G)\hat{x} = \hat{y}$  mit Matrizen  $F$  und  $G$  mit

$|F| \leq C(n) \frac{\varepsilon^*}{1-(2+\varepsilon^*)\varepsilon^*} |\hat{L}|, |G| \leq C(n) \frac{\varepsilon^*}{1-(2+\varepsilon^*)\varepsilon^*} |\hat{R}|$

Folgerung:

Die Auflösung von Dreieckssystemen ist ein gutartiger Prozeß.

Ob die  $LR$ -Zerlegung selbst gutartig ist, hängt davon ab, ob man die

Größenordnung der Matrixelemente von  $\hat{L}$  und  $\hat{R}$  kontrollieren kann.

Spaltenpivotisierung

$$\widehat{L}\widehat{R} = PA + H$$

$$|H| \leq C(n)\varepsilon^*(|PA| + |\widehat{L}||\widehat{R}|)$$

Begründung:

Der Gaußsche Algorithmus ohne Pivotsuche angewandt auf  $PA$  liefert identische Matrizen  $\widehat{L}, \widehat{R}$  wie der Gaußsche Algorithmus mit Pivotsuche angewandt auf  $A$ .

**Satz 3.23**

Die Spaltenpivotsuche gilt für die Matrixelemente der Dreiecksmatrizen  $L$  und  $R$  bei exakter Rechnung. Die Abschätzung

$$|l_{ik}| \leq 1, |r_{kl}| \leq 2^{k-1}a, \text{ wobei } a = \max_{k,l} |a_{kl}| \text{ sein soll.}$$

Beweis: (einfache) Induktion über  $k$ .

Bemerkung: Diese Schranken sind scharf, wie das Beispiel der Matrix

$$\begin{pmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ -1 & \cdots & -1 & 1 \end{pmatrix}$$

zeigt.

Meistens ist diese Abschätzung für die Elemente von  $R$  aber viel zu pessimistisch. Der Gaußalgorithmus verhält sich fast immer gutartig.

**3.5 Symmetrisch positiv-definite Matrizen****Definition 3.24**

Eine  $(m \times n)$ -Matrix  $A$  heißt symmetrisch, falls  $A^T = A$  ist.

Sie heißt positiv definit, falls  $\langle x, Ax \rangle > 0$  für alle  $x \in \mathbb{R}^n, x \neq 0$  gilt,

oder wenn es äquivalent dazu eine Konstante  $\delta > 0$  gibt mit  $\langle x, Ax \rangle \geq \delta \|x\|_2^2$  für alle  $x \in \mathbb{R}^n$ .



**Satz 3.25**

Jede symmetrische und positiv definite Matrix  $A$  besitzt eine  $LR$ -Zerlegung von der Form  $A = LR = LDL^T$  mit einer unteren Dreiecksmatrix  $L$  mit Einsen in der Diagonale und einer Diagonalmatrix  $D$  mit positiven Diagonalelementen.

Beweis: Wegen  $\langle e_1, Ae_1 \rangle = a_{11}$  ist  $a_{11} > 0$ . Sei  $G_1$  die Gaußsche

Eliminationsmatrix mit

$$\left( \begin{array}{c|ccc} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n1}^{(1)} & \cdots & a_{nn}^{(1)} \end{array} \right)$$

Dann ist wegen  $a_{ij} = a_{j1}$

$$(I - G_1)A(I - G_1)^T = \left( \begin{array}{c|ccc} a_{11} & 0 & \cdots & 0 \\ 0 & a_{22}^{(1)} & \cdots & a_{2n}^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & a_{n1}^{(1)} & \cdots & a_{nn}^{(1)} \end{array} \right) = \begin{pmatrix} a_{11} & 0 \\ 0 & B \end{pmatrix}$$

Die  $(n-1) \times (n-1)$ -Restmatrix  $B$  ist wegen  $((I - G_1)A(I - G_1)^T)^T$

$$\begin{aligned} &= ((I - G_1)^T)^T A^T (I - G_1)^T = (I - G_1)A^T(I - G_1)^T \\ &= (I - G_1)A(I - G_1)^T \text{ symmetrisch.} \end{aligned}$$

Ist  $x = (x_2, \dots, x_n)^T \in \mathbb{R}^{n-1}$  und  $\tilde{x} \in \mathbb{R}^n$  gegeben durch

$$\begin{aligned} \tilde{x} &= (0, x_2, \dots, x_n)^T, \text{ so gilt } \langle x, Bx \rangle = \langle \tilde{x}, (I - G_1)A(I - G_1)^T \tilde{x} \rangle \\ &= \langle (I - G_1)^T \tilde{x}, A(I - G_1)^T \tilde{x} \rangle \end{aligned}$$

Da  $(I - G_1)^T \tilde{x}$  genau dann der Nullvektor ist, wenn  $x = 0$  ist, ist für  $x \neq 0$   $\langle x, Bx \rangle > 0$ . Damit Restmatrix  $B$  wieder symmetrisch und positiv definit.

Durch Induktion folgt aus diesen Überlegungen, dass es Eliminationsmatrizen

$G_1, \dots, G_{n-1}$  gibt derart, dass

$$D := (I - G_{n-1}) \cdots (I - G_1)A(I - G_1)^T \cdots (I - G_{n-1})^T$$

eine Diagonalmatrix mit positiven Diagonalelementen ist. Wegen

$$L := ((I - G_{n-1}) \cdots (I - G_1))^{-1} = I + \sum_{k=1}^{n-1} G_k \text{ und}$$

$$\begin{aligned} &((I - G_1)^T \cdots (I - G_{n-1})^T)^{-1} = (((I - G_{n-1}) \cdots (I - G_1))^T)^{-1} \\ &= (((I - G_{n-1}) \cdots (I - G_1))^{-1})^T = L^T \end{aligned}$$

folgt daraus die Existenz der Zerlegung.

□

Folgerung:

Eine symmetrische Matrix ist genau dann positiv definit, wenn sich der Gaußsche Algorithmus ohne Spaltenpivotsuche für sie als durchführbar erweist und auf eine obere Dreiecksmatrix  $R(=DL^T)$  mit positiven Diagonalelementen liefert.

Wegen  $\langle x, LDL^T x \rangle = \langle L^T x, DL^T x \rangle > 0$

**Satz 3.26**

Ist die Matrix  $A$  symmetrisch und positiv definit und liefert der Gaußsche Algorithmus bei exakter Rechnung die  $LR$ -Zerlegung  $A = LR$ , so gilt für die Elemente der Matrix  $|L||R|$  die Abschätzung

$$|L||R|_{ij} \leq \sqrt{a_{ii}a_{jj}}$$

Beweis: Sei  $R = DL^T$ ,  $C = LD^{\frac{1}{2}}$ . Dann ist  $A = LR = LDL^T = CC^T$  und  $|L||R| = |C||C|^T$ . Für die Elemente der Matrix  $C = (c_{ij})$  gilt

$$\sum_{j=1}^n c_{ij}^2 = CC^T|_{ii} = a_{ii}$$

Daher ist

$$\begin{aligned} |L||R|_{ij} &= |C||C|^T|_{ij} = \sum_{k=1}^n |c_{ik}||c_{jk}| \\ &\leq \left(\sum_{k=1}^n |c_{ik}|^2\right)^{1/2} \left(\sum_{k=1}^n |c_{jk}|^2\right)^{1/2} \quad (\text{Schwarzsche Ungleichung}) \\ &= \sqrt{a_{ii}a_{jj}} \quad \square \end{aligned}$$

Folgerung:

Der Gaußsche Algorithmus ist für symmetrische-positiv definite Matrizen nicht nur immer durchführbar, sondern auch numerisch extrem stabil.

Bemerkung:

Die Zerlegung  $A = LBL^T$  einer symmetrisch positiv definiten Matrix  $A$  lässt sich durch Koeffizientenvergleich auf direktem Weg berechnen.

Gegenüber der allgemeinen Form des Gaußschen Algorithmus halbieren sich dabei Rechenaufwand und Speicherplatzbedarf.

Anzahl Punktoperationen:

$$\frac{1}{6}n^3 + \mathcal{O}(n^2)$$

Cholesky Zerlegung:

$$A = CC^T \quad (C = LD^{1/2} \text{ untere Dreiecksmatrix})$$

Lässt sich in gleicher Weise durch Koeffizientenvergleich bestimmen.

Rechenaufwand:

$$\frac{1}{6}n^3 + \mathcal{O}(n^2) \text{ Punktoperationen, } n \text{ Quadratwurzeln.}$$

### 3.6 Eliminationsverfahren für schwach besetzte Matrizen

Schwach besetzte Matrix:

sehr wenige Nichtnullelemente im Vergleich zur Gesamtzahl der Matrixelemente.

typisch  $\mathcal{O}(n)$  Nichtnullelemente bei  $n \times n$ -Matrix.

$$n = 1\,000, 10\,000, 100\,000, 1\,000\,000$$

#### Direkte Verfahren

Nur die Nichtnullelemente von  $A$ ,  $L$  und  $R$  werden abgespeichert.

Von besonderer Bedeutung sind in diesem Zusammenhang die benutzten Datenstrukturen. Für wesentliche Problemklassen wie z.B. symmetrische positiv definite Matrizen kann man auf eine Pivotsuche verzichten.

Dann läuft ein direktes Verfahren in vier Phasen ab.

Phase 1: Analysephase

In dieser Phase wird eine Permutationsmatrix  $P$  bestimmt, so dass die  $LR$ -bzw.  $LDL^T$ -Zerlegung der Matrix  $PAP^T$  möglichst wenig Nichtnullelemente enthält.

Dabei wird allein von der Besetzungsstruktur, d.h. der Position der Nichtnullelemente von  $A$  ausgegangen.

Algorithmen zur Bestimmung solcher Permutationsmatrizen beruhen auf Graphentheoretischen Überlegungen.

Phase 2: symbolische Zerlegung

In dieser Phase wird die Position der Nichtnullelemente von  $L$  und  $R$  bzw. nur von  $L$  bestimmt und eine entsprechende Datenstruktur aufgebaut.

Phase 3: numerische Zerlegung

In dieser Phase werden  $L$  und  $R$  bzw.  $L$  und  $D$  berechnet.

Bemerkung: Phase 2 und Phase 3 sind am aufwändigsten. In der Regel wächst der Speicherplatz und Rechenaufwand überproportional zur Anzahl der Unbekannten.

Phase 4: Lösungsphase

Auflösung der Dreieckssysteme.

### Definition 3.27 (Bandmatrix)

Eine  $(n \times n)$ -Matrix  $A = (a_{ik})$  heißt Bandmatrix der Bandbreite  $m$ ,  $m < n$ , falls  $a_{ik} = 0$  für  $|i - k| > m$  gilt.

Im Fall  $m = 1$  spricht man von einer Tridiagonalmatrix.

### Satz 3.28

Besitzt eine Bandmatrix  $A$  der Bandbreite  $m$  eine  $LR$ -Zerlegung und ist  $A$  invertierbar, so sind auch  $L$  und  $R$  Bandmatrizen der Bandbreite  $m$ .

Beweis: (und Konstruktion von  $L$  und  $R$ )

durch Koeffizientenvergleich.

□

Folgerung:

Berechnet man eine solche  $LR$ -Zerlegung durch Koeffizientenvergleich, braucht man  $nm^2 - \frac{2}{3}m^3 + mn$  Punktoperationen.

Für  $m \ll n$  ist dies eine sehr große Ersparnis gegenüber den  $\frac{1}{3}n^3 - \frac{1}{3}n$  Punktoperationen bei voller Abspeicherung der Matrix.

bei Spaltenpivotsuche: ähnliche Einsparungen

### 3.7 Iterationsverfahren

Problem: Bei direkten Verfahren (auch wenn sie auf schwach besetzte Matrizen zugeschnitten sind) wächst die Rechenzeit und Speicherplatzaufwand in der Regel überproportional zur Anzahl der Unbekannten.

Ausweg: Iterative Verfahren

gegeben:  $A$  nichtsinguläre ( $n \times n$ ) Matrix

$B$  nichtsinguläre ( $n \times n$ ) Matrix

Iterationsverfahren:

$$x^{(i+1)} = x^{(i)} + B^{-1}(b - Ax^{(i)})$$

zur Lösung des Gleichungssystems

$$Ax = b$$

mit der exakten Lösung  $\bar{x}$ .

wichtige Bemerkung:

Im Gegensatz etwa zum Gaußschen Eliminationsverfahren sind iterative Verfahren keine universellen Verfahren. Zur Konstruktion einer guten Näherungsinversen  $B^{-1}$  für  $A$  muss die Problemstruktur ausgenutzt werden.

Beispiel:

Mehrgitterverfahren für diskretisierte partielle Differentialgleichungen.

#### Beispiel 3.29

Die Matrix  $A$  habe die Darstellung

$$A = L + D + R \text{ mit}$$

$$L|_{ij} = 0 \text{ für } j \geq i$$

$$L|_{ij} = 0 \text{ für } j \leq i$$

$$D|_{ij} = 0 \text{ für } i \neq j$$

$D$  sei nichtsingulär. Dann ist  $x^{(i+1)} = x^{(i)} + D^{-1}(b - Ax^{(i)})$  das

Gesamtschritt oder Jacobi-Verfahren.

In Koordinatenschreibweise:

$$x_k^{(i+1)} = x_k^{(i)} + \frac{1}{a_{kk}} \left( b_k - \sum_{j=1}^n a_{kj} x_j^{(i)} \right)$$

$$= \frac{1}{a_{kk}} \left( b_k - \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj} x_j^{(i)} \right)$$

**Beispiel 3.30**

(Bezeichnungen wie im Beispiel 3.29)

Das verfahren

$x^{(i+1)} = x^{(i)} + (L + D)^{-1}(b - Ax^{(i)})$  ist das Einzelschritt- oder

Gauß-Seidel-Verfahren.

Wegen  $Dx^{(i+1)} = b - Lx^{(i+1)} - Rx^{(i)}$  lautet es in Komponentenschreibweise

$$x_k^{(i+1)} = \frac{1}{a_{kk}}(b_k - \sum_{j < k} a_{kj}x_j^{(i+1)} - \sum_{j > k} a_{kj}x_j^{(i)}) \quad (\text{leere Summe}=0)$$

einfacher zu realisieren als das Jacobi-Verfahren!

**Satz 3.31**

Für die Fehler  $x^{(i)} - \bar{x}$  gilt die Rekursion:

$$x^{(i+1)} - \bar{x} = (I - B^{-1}A)(x^{(i)} - \bar{x})$$

Die Matrix  $I - B^{-1}A$  ist die Fehlerfortpflanzungsmatrix.

Beweis: klar

□

**Satz 3.32**

Gilt in der durch die Vektornorm  $\|\cdot\|$  induzierte Matrixnorm

$$\|I - B^{-1}A\| \leq q < 1 \text{ so ist } \|x^{(i+1)} - \bar{x}\| \leq q\|x^{(i)} - \bar{x}\| \leq q^{(i)}\|x^{(0)} - \bar{x}\|,$$

speziell also  $\|x^{(i)} - \bar{x}\| \rightarrow 0$ .

Beweis: klar

□

**Beispiel 3.33**

Gibt es eine Konstante  $\theta < 1$  und

$$\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \leq \theta |a_{ii}|, \quad i = 1, \dots, n$$

so gilt bezüglich der Matrixnorm als Vektornorm für die Fehlerfortpflanzungsmatrix des Jacobiverfahrens  $\|I - D^{-1}A\| \leq \theta$ .

Beweis: Für  $i = 1, \dots, n$  ist

$$|x_i - \frac{1}{a_{ii}} \sum_{j=1}^n a_{ij} x_j|$$

$$= |\frac{1}{a_{ii}} \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j|$$

$$\leq (\frac{1}{|a_{ii}|} \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|) \|x\|$$

$$\leq \theta \|x\|$$

$$\Rightarrow \|(I - D^{-1}A)x\| \leq \theta \|x\|$$

□

### Beispiel 3.34

Unter den Voraussetzungen aus Beispiel 3.33

gilt für die Fehlerfortpflanzungsmatrix des Gauß-Seidel Verfahrens bezüglich der Maximumnorm  $\|I - (L + D)^{-1}A\| \leq \theta$ .

Beweis: Wir zeigen durch Induktion über  $k$

$$|x_k^{(i+1)} - \bar{x}_k| \leq \theta \|x^{(i)} - \bar{x}\|$$

Es ist

$$Dx^{(i+1)} = b - Lx^{(i+1)} - Rx^{(i)}$$

$$D\bar{x} = b - L\bar{x} - R\bar{x}, \text{ also}$$

$$x^{(i+1)} - \bar{x} = D^{-1}L(\bar{x} - x^{(i+1)}) + D^{-1}R(\bar{x} - x^{(i)})$$

oder in Komponentenschreibweise

$$x_k^{(i+1)} - \bar{x}_k = \frac{1}{a_{kk}} \sum_{j < k} a_{kj} (\bar{x}_j - x_j^{(i+1)}) + \frac{1}{a_{kk}} \sum_{j > k} a_{kj} (\bar{x}_j - x_j^{(i)})$$

Daraus folgt

$$|x_k^{(i+1)} - \bar{x}_k| \leq \left[ \frac{1}{|a_{kk}|} \sum_{j < k} |a_{kj}| \theta + \frac{1}{|a_{kk}|} \sum_{j > k} |a_{kj}| \right] \|\bar{x} - x^{(i)}\|$$

$$\leq \left[ \frac{1}{|a_{kk}|} \sum_{j \neq k} |a_{kj}| \right] \|\bar{x} - x^{(i)}\|$$

und endlich

$$\|x^{(i+1)} - \bar{x}\| \leq \theta \|x^{(i)} - \bar{x}\|$$

Die  $x^{(i)} - \bar{x}$  ein beliebiger Vektor sein kann folgt daraus

$$\|I - (L + D)^{-1}A\| \leq \theta.$$

Verfahren:  $x^{(i+1)} = x^{(i)} + B^{-1}(b - Ax^{(i)})$

Fehlerfortpflanzung:

$$x^{(i+1)} - \bar{x} = (I - B^{-1}A)(x^{(i)} - \bar{x})$$

### Satz 3.35

Sind  $\|\cdot\|$  und  $\|\|\cdot\|\|$  Normen auf dem  $\mathbb{R}^n$  für die mit paarweise verschiedenen Konstanten  $K_1, K_2$ ,  $K_1\|x\| \leq K_2\|\|x\|\|$  ( $x \in \mathbb{R}^n$ ) gilt, und gilt bezüglich der durch die Norm  $\|\|\cdot\|\|$  induzierte Matrixnorm  $\|\|I - B^{-1}A\|\| \leq q < 1$ , so ist  $\|x^{(i)} - \bar{x}\| \leq \frac{K_2}{K_1} q^i \|x^{(i)} - \bar{x}\|$ .

Beweis:

$$\|x^{(i)} - \bar{x}\| \leq \frac{1}{K_1} \|\|x^{(i)} - \bar{x}\|\| \leq \frac{1}{K_1} q^i \|\|x^{(i)} - \bar{x}\|\| \leq \frac{1}{K_1} q^i K_2 \|x^{(i)} - \bar{x}\|$$

□

Bemerkung: Nach Satz 3.5 existieren solche Konstanten  $K_1$  und  $K_2$  immer; praktisch ist die Größenordnung von  $K_2/K_1$  entscheidend.

Ziel: Charakterisierung der Konvergenz

Spektralradius  $\rho(A)$  einer Matrix  $A$ :

Betrag des größten Eigenwertes  $\lambda \in \mathbb{C}$  von  $A$ .

### Satz 3.36

Zu jeder  $(n \times n)$ -Matrix  $A$  und jedem  $\varepsilon > 0$  gibt es eine Vektornorm auf dem  $\mathbb{R}^n$ , so dass für die induzierte Matrixnorm  $\|\|A\|\| \leq \rho(A) + \varepsilon$  gilt.

Beweis: Sei  $A = T^{-1}JT$  und

$$J = \begin{pmatrix} J_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & J_n \end{pmatrix}, J_i = \begin{pmatrix} \lambda_i & 1 & & 0 \\ & \ddots & \ddots & \\ & & \ddots & 1 \\ 0 & & & \lambda_i \end{pmatrix}$$

(Jordansche Normalform)

Sei weiter  $D = \text{diag}(1, \varepsilon, \dots, \varepsilon^{n-1})$

Dann ist

$$D^{-1}TAT^{-1}D = \begin{pmatrix} J'_1 & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & J'_n \end{pmatrix}, J'_i = \begin{pmatrix} \lambda_i & \varepsilon & & 0 \\ & \ddots & \ddots & \\ & & \ddots & \varepsilon \\ 0 & & & \lambda_i \end{pmatrix}$$



Wählt man jetzt als Vektornorm  $\|x\| := \|D^{-1}Tx\|_\infty$ , so ist wegen

$$D^{-1}T(Ax) = (D^{-1}TAT^{-1}D)(D^{-1}Tx)$$

die zugehörige Matrixnorm die Zeilensummennorm von

$$D^{-1}TAT^{-1}D, \text{ also } \|A\| \leq \max_i |\lambda_i| + \varepsilon = \rho(A) + \varepsilon$$

□

### Satz 3.37

Ist  $\rho(I - B^{-1}A) < \delta < 1$ , so gilt in jeder Norm auf dem  $\mathbb{R}^n$

$$\lim_{i \rightarrow \infty} \delta^{-i} \|x^{(i)} - \bar{x}\| = 0$$

Ist umgekehrt  $\rho(I - B^{-1}A) \geq 1$  so konvergiert  $x^{(i)}$  nicht für alle Startvektoren  $x^{(i)}$  gegen  $\bar{x}$ .

Beweis: Sei  $\rho(I - B^{-1}A) < q < \delta < 1$ . Dann gibt es nach Satz 3.36 eine

Norm  $\|\cdot\|$  auf dem  $\mathbb{R}^n$ , so daß bezüglich der induzierten Matrixnorm

$$\|I - B^{-1}A\| \leq q \text{ ist. Nach Satz 3.35 folgt } \|x^{(i)} - \bar{x}\| \leq q^i \|x^{(i)} - \bar{x}\|.$$

Da sämtliche Normen auf dem  $\mathbb{R}^n$  äquivalent sind, gibt es eine

$$\text{Konstante } K \text{ mit } \|x\| \leq K \|x\| \text{ für alle } x, \text{ also } \|x^{(i)} - \bar{x}\| \leq K q^i \|x^{(0)} - \bar{x}\|$$

$$\text{Daraus folgt wegen } q < \delta \lim_{i \rightarrow \infty} \delta^{-i} \|x^{(i)} - \bar{x}\| = 0$$

Ist  $\rho(I - B^{-1}A) \geq 1$  so gibt es  $z \in \mathbb{C}, z \neq 0$ , und ein  $\lambda \in \mathbb{C}$  mit

$$(I - B^{-1}A)z = \lambda z, |\lambda| \geq 1$$

Damit kann  $(I - B^{-1}A)^i z = \lambda^i z$  nicht gegen 0 streben.

Ist  $z = u + iv, n, v \in \mathbb{R}^n$ , so gilt für ein  $e \in \{u, v\}$

$$(I - B^{-1}A)e \not\rightarrow 0 \text{ für } i \rightarrow \infty.$$

## 4. Nichtlineare Gleichungssysteme

$$\sin x_1 - x_2 = 0$$

$$x_1 - \cos x_2 = 0$$

### 4.1 Der Banachsche Fixpunktsatz

#### Satz 4.1 (BFP)

Die Funktion  $\varphi : G \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$  bilde die abgeschlossene Menge  $G$  in sich ab. Sei sie Lipschitz-stetig bezüglich der Norm  $\|\cdot\|$ , d.h. für alle  $x, y \in G$  sei  $\|\varphi(x) - \varphi(y)\| \leq L\|x - y\|$  mit einer Konstanten  $L < 1$ .

Dann besitzt  $\varphi$  genau einen Fixpunkt  $\bar{x} \in G$ . Für alle  $x_0 \in G$  strebt da durch  $x_{k+1} = \varphi(x_k)$ ,  $k = 0, 1, 2, \dots$  gegebene Folge gegen  $\bar{x}$ .

Es gelten die a posteriori Abschätzung

$$\|x_k - \bar{x}\| \leq \frac{L}{1-L} \|x_k - x_{k-1}\|$$

und die a priori Abschätzung

$$\|x_k - \bar{x}\| \leq \frac{L^k}{1-L} \|x_1 - x_0\|.$$

Beweis: genau wie Satz 2.14

□

Bemerkung: Ist  $G$  konvex, ist die Schrittfunktion  $\varphi$  auf einer offenen Obermenge von  $G$  stetig differenzierbar und ist  $L := \sup_{x \in G} \|\varphi'(x)\| < 1$ ,

so gilt nach dem Hauptsatz der Differential- und Integralrechnung

$$\|\varphi(x) - \varphi(y)\| = \left\| \int_0^1 \varphi'(x + t(y-x))(y-x) dt \right\|$$

$$\leq \int_0^1 \|\varphi'(x + t(y-x))(y-x)\| dt$$

$$\leq \int_0^1 \|\varphi'(x + t(y-x))\| \|y-x\| dt$$

$$\leq \int_0^1 L \|y-x\| dt$$

$$= L \|y-x\|$$

□

**Beispiel 4.2**

Iterationsverfahren der Form  $x_{k+1} = \varphi(x_k)$ ,  $\varphi(x) = x + B^{-1}(b - Ax)$   
zur Lösung des linearen Gleichungssystems  $Ax = b$ .

In diesem Falls ist  $L = \|I - B^{-1}A\|$

weiteres Beispiel:

Fixpunktiteration der Form

$$x_{k+1} = x_k - f'(x_0)^{-1}f(x_k)$$

zur Lösung des Systems  $f(x) = y$

In diesem Fall ist  $\varphi(x) = x - f'(x_0)^{-1}(f(x) - y)$  und  
 $\varphi'(x) = I - f'(x_0)^{-1}f'(x)$

anderes Beispiel:

$$\varphi(x) = a + \varepsilon\psi(x)$$

**4.2 Das Newton Verfahren**

gegeben:  $f : G \subseteq \mathbb{R}^n \rightarrow \mathbb{R}^n$

gesucht:  $\bar{x} \in G$  mit  $f(\bar{x}) = 0$

Idee: lokale Linearisierung

$$0 = f(\bar{x}) = f(x_k) + f'(x_k)(\bar{x} - x_k) + \dots$$

$$0 = f(x_k) + f'(x_k)(\bar{x} - x_k)$$

$$x_{k+1} = x_k - f'(x_k)^{-1}f(x_k)$$

rechnerisch:

$$f'(x_k)(x_{k+1} - x_k) = -f(x_k)$$

(Löse lineares Gleichungssystem für  $\delta_2 = x_{k+1} - x_k$ )

**affine Invarianz**

$A$  invertierbare Matrix,  $g(x) = Af(x)$

$$x - g'(x)^{-1}g(x)$$

$$= x - (Af'(x))^{-1}(Af(x))$$

$$= x - (f'(x)^{-1}A^{-1})(Af(x))$$

$$= x - f'(x)^{-1}(A^{-1}A)f(x)$$

$$= x - f'(x)^{-1}f(x)$$

**Beispiel 4.3**

$$\sin x_1 - x_2 = 0$$

$$x_1 - \cos x_2 = 0$$

$$\text{Kurzform: } f(x) = 0 \text{ mit } f(x) = \begin{pmatrix} f_1(x_1, x_2) \\ f_2(x_1, x_2) \end{pmatrix} = \begin{pmatrix} \sin x_1 - x_2 \\ x_1 - \cos x_2 \end{pmatrix}$$

Es ist

$$f'(x) = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{pmatrix} = \begin{pmatrix} \cos x_1 & -1 \\ 1 & \sin x_2 \end{pmatrix}$$

$$f'(x)^{-1} = \frac{1}{1 + \cos x_1 \sin x_2} \begin{pmatrix} \sin x_2 & 1 \\ -1 & \cos x_1 \end{pmatrix}$$

Iterationsvorschrift:

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} - \frac{1}{1 + \cos x_1 \sin x_2} \begin{pmatrix} \sin x_2 & 1 \\ -1 & \cos x_1 \end{pmatrix} \begin{pmatrix} \sin x_1 - x_2 \\ x_1 - \cos x_2 \end{pmatrix}$$

Beobachtung: sehr schnelle Konvergenz

□

**Satz 4.4**

Die Funktion  $f$  sei auf einer Kugel  $\mathcal{K}$  mit ihrer Nullstelle  $\bar{x}$  als Mittelpunkt definiert und dort stetig differenzierbar.

$\forall x \in \mathcal{K}$  sei  $f'(x)$  nichtsingulär.

Es geben eine Konstante  $\omega > 0$  mit  $\|f'(x)^{-1}[f'(y) - f'(x)]\| \leq \omega \|y - x\|$  für alle  $x, y \in \mathcal{K}$ . Für ein gegebenes  $x_0 \in \mathcal{K}$  ist  $q := \frac{\omega}{2} \|x_0 - \bar{x}\| < 1$ .

Dann ist die Folge der Newton Iterierten  $x_{k+1} = x_k - f'(x_k)^{-1}f(x_k)$ ,  $k = 0, 1, 2, \dots$ , wohldefiniert. Sie konvergiert gegen  $\bar{x}$ .

Alle  $x_k$  liegen in  $\mathcal{K}$  und es gilt die Abschätzung

$$\|x_{k+1} - \bar{x}\| \leq \frac{\omega}{2} \|x_k - \bar{x}\|^2$$

Beweis: Für  $x \in \mathcal{K}$  sei  $\varphi(x) := x - f'(x)^{-1}f(x)$ .

Hält man  $x \in \mathcal{K}$  fest und wendet den Hauptsatz der Differential- und Integralrechnung auf die Komponenten der Hilfsfunktion

$\psi(t) = f(x + t(\bar{x} - x))$ ,  $0 \leq t \leq 1$  an, so erhält man

$$f(\bar{x}) = f(x) + \int_0^1 f'(x + t(\bar{x} - x))(\bar{x} - x) dt$$

Wegen  $f(\bar{x}) = 0$  folgt daraus

$$f(x) = \int_0^1 f'(x + t(\bar{x} - x))(x - \bar{x}) dt$$

$$\begin{aligned}
\text{Daher ist } \bar{x} - \varphi(x) &= \bar{x} - x + f'(x)^{-1}f(x) \\
&= f'(x)^{-1}[f(x) - f'(x)(x - \bar{x})] \\
&= \int_0^1 f'(x)^{-1}[f'(x + t(\bar{x} - x)) - f'(x)](x - \bar{x}) dt
\end{aligned}$$

Daraus folgt

$$\begin{aligned}
\|\varphi(x) - \bar{x}\| &\leq \int_0^1 \|f'(x)^{-1}[f'(x + t(\bar{x} - x)) - f'(x)](x - \bar{x})\| dt \\
&\leq \int_0^1 \|f'(x)^{-1}[f'(x + t(\bar{x} - x)) - f'(x)]\| \|x - \bar{x}\| dt \\
&\leq \int_0^1 \omega \|x + t(\bar{x} - x) - x\| \|x - \bar{x}\| dt \\
&= \int_0^1 \omega t \|\bar{x} - x\| \|x - \bar{x}\| dt \\
&= \frac{\omega}{2} \|x - \bar{x}\|^2
\end{aligned}$$

Für alle  $x \in \mathcal{K}$  ist also

$$\|\varphi(x) - \bar{x}\| \leq \frac{\omega}{2} \|x - \bar{x}\|^2 \quad (1)$$

Ist daher  $x_k \in \mathcal{K}$  und gilt

$$\frac{\omega}{2} \|x_k - \bar{x}\| \leq q \quad (2)$$

so erfüllt  $x_{k+1} = \varphi(x_k)$  wegen (1)

$$\|x_{k+1} - \bar{x}\| \leq \frac{\omega}{2} \|x_k - \bar{x}\|^2 \leq q \|x_k - \bar{x}\|$$

Damit ist die Folge der  $x_k$  wohldefiniert. Alle  $x_k$  liegen in  $\mathcal{K}$  und erfüllen (2).

Aus (1) folgt weiter  $\|x_{k+1} - \bar{x}\| \leq \frac{\omega}{2} \|x_k - \bar{x}\|^2$

Das bedeutet  $\frac{\omega}{2} \|x_k - \bar{x}\| \leq q^{2^k}$  also Konvergenz der  $x_k$  gegen  $\bar{x}$ .

□

Folgerung: lokal quadratische Konvergenz

Bemerkung: Die Konstante  $\omega$  ist so etwas wie eine Lipschitzkonstante von  $f'$ .

Sie ist affin invariant: Sei also  $g(x) = Af(x)$ ,  $A$  invertierbar, nichtsingulär.

Dann gilt wegen  $g'(x) = Af'(x)$

$$g'(x)^{-1}[g(y) - g(x)] = f'(x)^{-1}[f(y) - f(x)]$$

globales Konvergenzverhalten:

sehr kompliziert

### Beispiel 4.5

Nullstellen des konvergenten Polynoms  $z^3 - 1$ :

äquivalentes reelles System für  $x, y$  ( $z = x + iy$ )

$$x^3 - 3xy^2 - 1 = 0$$

$$-3xy^2 + y^3 = 0$$

#### Vergrößerung des Einzugsbereichs

Benutze eine Dämpfungsstrategie, um zu garantieren, dass eine vorgegebene, durch ein Skalarprodukt induzierte Norm des Residuums  $f(x_k)$  von Schritt zu Schritt hinreichend stark abnimmt. Solche Strategien beruhen auf der

Beobachtung, dass die Funktion

$$\psi(\theta) = \|f(x - \theta f'(x)^{-1} f(x))\|^2$$

bei vorgegebener durch ein Skalarprodukt induzierter Norm die Ableitung

$$\psi'(0) = -2\|f(x)\|^2$$

hat. Eine ganz einfache (nicht die beste) Version eines solchen Verfahrens sieht wie folgt aus:

Start: (Schritt 0)

Gebe  $x_0$  vor. Berechne  $f(x_0), \|f(x_0)\|$

#### Schritt k+1:

$x_k, f(x_k), \|f(x_k)\|$  sind bekannt

1) Berechne  $f'(x_k)$

2) Berechne die Korrekturrichtung  $d$  als Lösung des GLS  $f'(x_k)d = -f(x_k)$

Setze  $\theta = 1$

3) Berechne  $x_k + \theta d, f(x_k + \theta d), \|f(x_k + \theta d)\|$

4)  $\|f(x_k + \theta d)\| \leq (1 - \frac{\theta}{2})\|f(x_k)\|$ ?

ja: Setze  $x_{k+1} = x_k + \theta d$

nein: Halbiere  $\theta$  und springe auf 3)

#### Konvergenz:

bei Start in kompakter Niveaumenge  $\{x \mid \|f(x)\| \leq \|f(x_0)\|\} =: N(x_0)$ , falls  $f'(x)$  auf  $N(x_0)$  nichtsingulär ist.

affin invariante Version:

$x_k, f'(x_k)^{-1}f(x_k), \|f'(x_k)^{-1}f(x_k)\|_2$  gegeben.

1) setze  $\theta = 1$

2) Berechne mit  $d := -f'(x_k)^{-1}f(x_k), x_k + \theta d, f(x_k + \theta d)$  sowie

$\bar{d} = f'(x_k)^{-1}f(x_k + \theta d), \|\bar{d}\|_2$

3)  $\|\bar{d}\|_2 \leq \|d\|_2$ ?

ja: Setze  $x_{k+1} = x_k + \theta d, f(x_{k+1}) = f(x_k + \theta d)$

Berechne  $f'(x_{k+1})$  und  $f'(x_{k+1})^{-1}f(x_{k+1}), \|f'(x_{k+1})^{-1}f(x_{k+1})\|$

nein: halbiere  $\theta$  und gehen auf 2)

Abwandlungen des Newtonverfahrens:

1)  $f'(x)$  wird über mehrere Schritte konstant gehalten

2)  $f'(x)$  wird nicht exakt, sondern durch numerische Differentiation berechnet (Differenzenquotienten)

3)  $f'(x)^{-1}$  wird in singulären Punkten durch  $f'(x)^+$  (Pseudoinverse) ersetzt

$x_{k+1} = x_k - f'(x_k)^+ f(x_k)$

4) Rangmodifikationsverfahren

$x_{k+1} = x_k - H_k f(x_k)$

$\text{rang}(H_k - H_{k+1}) = 1$  oder  $2$

5) näherungsweise Lösung des linearen Gleichungssystems

$f'(x_k)(x_{k+1} - x_k) = -f(x_k)$

## 5. Lineare Ausgleichsprobleme

### 5.1 Lineare Ausgleichsprobleme in unitären Räumen

#### Definition 5.1

Sei  $V$  ein reeller Vektorraum. Eine Abbildung  $\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  heißt Skalarprodukt auf  $V$ , wenn für alle  $x, y, z \in V$  und alle  $\alpha, \beta \in \mathbb{R}$  die folgenden Bedingungen erfüllt sind.

- (i)  $\langle x, y \rangle = \langle y, x \rangle$
- (ii)  $\langle \alpha x + \beta y, z \rangle = \alpha \langle x, z \rangle + \beta \langle y, z \rangle$
- (iii)  $\langle x, x \rangle \geq 0, \langle x, x \rangle = 0 \iff x = 0$

Zusammen mit einem solchen Skalarprodukt wird  $V$  zu einem unitären Raum.

#### Satz 5.2

Ist  $\langle \cdot, \cdot \rangle$  ein Skalarprodukt auf dem reellen Vektorraum  $V$ , so gilt für alle  $x, y \in V$  die Schwarze Ungleichung.

$$\langle x, y \rangle \leq \sqrt{\langle x, x \rangle \langle y, y \rangle}$$

Durch  $\|x\| := \sqrt{\langle x, x \rangle}$  wird eine Norm auf  $V$  induziert.

Es ist  $\|x - y\|^2 + \|x + y\|^2 = 2\|x\|^2 + 2\|y\|^2$  und

$$\langle x, y \rangle = \frac{1}{4}(\|x + y\|^2 - \|x - y\|^2)$$

#### Beispiel 5.3

Der  $\mathbb{R}^n$  und dem Standardskalarprodukt

$$\langle x, y \rangle = \sum_{i=1}^n x_i y_i$$

und der durch es induzierten euklidischen Norm

$$\|x\|_2 = \left( \sum_{i=1}^n x_i^2 \right)^{\frac{1}{2}}$$

#### Beispiel 5.4

Der Raum  $C[a, b]$  aller stetigen Funktionen  $f : [a, b] \rightarrow \mathbb{R}$  mit dem

Skalarprodukt  $\langle f, g \rangle = \int_a^b f(x)g(x)dx$  und der  $L_2$ -Norm

$$\|f\| = \left( \int_a^b f(x)^2 dx \right)^{\frac{1}{2}}$$



**Beispiel 5.5**

Der Raum  $C_0^1[a, b]$  aller stetig differenzierbaren Funktionen  $f : [a, b] \rightarrow \mathbb{R}$  und  $f(a) = 0, f(b) = 0$  und dem Skalarprodukt.

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

**Definition 5.6**

Sei  $V$  ein reeller Vektorraum mit der Norm  $\|\cdot\|$  und  $M$  ein linearer Teilraum von  $V$ . Dann heißt  $y^* \in M$  ein Element bester Approximation an  $x \in V$  bezüglich der gegebenen Norm, wenn für  $y \in M$   $\|x - y^*\| \leq \|x - y\|$  gilt.

**Satz 5.7**

Die Norm werde durch ein Skalarprodukt  $\langle \cdot, \cdot \rangle$  induziert. Dann ist  $y^* \in M$  genau dann ein Element bester Approximation an  $x \in V$ , wenn für alle  $v \in M$   $\langle x - y^*, v \rangle = 0$  ist. Es gibt höchstens ein Element bester Approximation.

Beweis: Sei  $y^* \in M$  ein Element bester Approximation an  $x \in V$ . Sei  $v \in M$ .

Dann nimmt die skalare Funktion  $F(\alpha) := \|x - (y^* + \alpha v)\|^2 = \|(x - y^*) - \alpha v\|^2 = \|x - y^*\|^2 - 2\langle x - y^*, v \rangle \alpha + \|v\|^2 \alpha^2$  wegen  $y^* + \alpha v$  in  $M$  ein globales Minimum an. Wegen  $F'(\alpha) = -2\langle x - y^*, v \rangle + 2\|v\|\alpha$  und  $F'(0) = -2\langle x - y^*, v \rangle$  folgt daraus  $\langle x - y^*, v \rangle = 0$ .

Sei nun umgekehrt  $\langle x - y^*, v \rangle = 0$  für alle  $v \in M$ . Dann ist für alle

$$v \in M : \|x - (y^* + v)\|^2 = \|x - y^*\|^2 - 2\langle x - y^*, v \rangle + \|v\|^2 = \|x - y^*\|^2 + \|v\|^2 \text{ also } \|x - (y^* + v)\| > \|x - y^*\| \text{ für } v \neq 0.$$

Damit ist  $y^*$  ein Element bester Approximation und

$$\|x - y^*\| < \|x - y\| \text{ für alle } y \neq y^* \text{ aus } M. \square$$

**Satz 5.8**

Ist  $M$  endlichdimensional und bilden die Vektoren  $v_1, \dots, v_n$  eine Basis

von  $M$ , so ist  $y^* = \sum_{i=1}^n a_i v_i$  genau dann ein Element bester Approximation

an  $x \in V$ , wenn die Koeffizienten  $a_i$  das Normalgleichungssystem

$$\sum_{i=1}^n \langle v_i, v_j \rangle a_i = \langle x, v_j \rangle, j = 1, \dots, n \text{ erfüllen.}$$

Der Koeffizientenmatrix dieses Systems ist symmetrisch und positiv definit.

Das Normalgleichungssystem ist damit eindeutig lösbar.

Beweis: Die charakterisierende Bedingung aus Satz 5.7 ist genau dann erfüllt, wenn für  $j = 1, \dots, n$

$$\langle x - \sum_{i=1}^n \alpha_i v_i, v_j \rangle = 0$$

gilt, das heißt wenn  $\alpha_i$  das Gleichungssystem

$$\sum_{i=1}^n \langle v_i, v_j \rangle \alpha_i = \langle x, v_j \rangle, \quad j = 1, \dots, n \text{ erfüllen.}$$

Die Koeffizientenmatrix  $A$  mit  $A|_{ij} = \langle v_j, v_i \rangle$  dieses Systems ist wegen der Symmetrie des Skalarprodukts symmetrisch. Für alle Vektoren

$$\alpha = (\alpha_1, \dots, \alpha_n)^T \in \mathbb{R}^n \text{ gilt } \alpha^T A \alpha = \sum_{i,j=1}^n \langle v_i, v_j \rangle \alpha_i \alpha_j$$

$$= \langle \sum_{i=1}^n \alpha_i v_i, \sum_{j=1}^n \alpha_j v_j \rangle = \left\| \sum_{i=1}^n \alpha_i v_i \right\|^2 \geq 0 \text{ und genau dann}$$

$$\alpha^T A \alpha = 0, \text{ wenn } \sum_{i=1}^n \alpha_i v_i = 0 \text{ ist. Da die}$$

$v_i$  als Basis linear unabhängig sind, ist dies genau dann der Fall,

wenn  $\alpha_1, \dots, \alpha_n = 0$  oder  $\alpha = 0$  ist.

□

### Beispiel 5.9

Sei  $V = C[0, 1]$ , versehen mit dem Skalarprodukt  $\langle f, g \rangle = \int_0^1 f(x)g(x)dx$ .

Sei  $M$  der Raum aller Polynome  $n$ -ten Grades. Zu einer gegebenen Funktion  $f : [0, 1] \rightarrow \mathbb{R}$  ist das Polynom  $P = P^*$  vom Grad  $n$  gesucht

für das der Abstand  $\|f - P\| = \left( \int_0^1 |f(x) - P(x)|^2 dx \right)^{1/2}$  minimal wird.

Eine Basis des Raums aller Polynome  $n$ -ten Grades bilden die Monome  $\varphi_i(x) = x^i, i = 0, \dots, n$ . Stellt man das optimale Polynom  $P^*$  in der Form

$$P^*(x) = \sum_{i=0}^n \alpha_i x^i = \sum_{i=0}^n \alpha_i \varphi_i(x) \text{ dar, müssen die } \alpha_i \text{ das System}$$

$$\sum_{i=0}^n \langle \varphi_i, \varphi_j \rangle \alpha_i = \langle f, \varphi_j \rangle, \quad j = 0, \dots, n \text{ erfüllen.}$$

Wegen  $\langle \varphi_i, \varphi_j \rangle = \int_0^1 x^i x^j dx = \frac{1}{i+j+1}$  ist die Koeffizientenmatrix dieses

Systems die Hilbertmatrix  $H$  der Dimension  $n + 1$  mit den Koeffizienten

$$H|_{ij} = \frac{1}{i+j+1}$$

Die Hilbertmatrizen sind berüchtigt für ihre schlechte Kondition.

Das beschriebene Gleichungssystem ist daher ungeeignet, um das gegebene Ausgleichsproblem zu lösen.

Ausweg: man muss zu einer besser konditionierten Basis übergehen.

Ideal ist eine Orthonormalbasis mit Basiselementen  $\varphi_i, i = 0, \dots, n$

mit  $\langle \varphi_i, \varphi_j \rangle = \delta_{ij}$ . Dann ist  $P^*(x) = \sum_{i=0}^n \langle f, \varphi_i \rangle \varphi_i$ .

## Fourieranalyse

Gesucht ist die beste Approximation einer  $2\pi$ -periodischen Funktion

$f: \mathbb{R} \rightarrow \mathbb{C}$  durch ein trigonometrisches Polynom

$$\sum_{k=-n}^n a_k e^{ikx}$$

$n$ -ten Grades im Sinne der durch das Skalarprodukt

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(x) \overline{g(x)} dx$$

induzierten Norm. Die Monome  $f_k(x) = e^{ikx}$  bilden eine Basis dieses Raumes.

Wegen  $\langle \varphi_k, \varphi_l \rangle = 2\pi \delta_{kl}$  ist  $\sum_{k=-n}^n \hat{f}(k) e^{ikx}, \hat{f}(k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx$

die optimale Lösung. Die Koeffizienten  $\hat{f}(k), k \in \mathbb{Z}$ , sind die

Fourierkoeffizienten der Funktion  $f$ .

## Beispiel 5.10

Sei  $V$  der Vektorraum  $C[a, b]$  mit dem Skalarprodukt

$$\langle f, g \rangle = \int_a^b f(x) g(x) dx. \text{ Mit } h = \frac{b-a}{n}, n \in \mathbb{N} \text{ fest, sei}$$

$$x_k = a + kh = a + \frac{b-a}{n} k, k = 0, \dots, n$$

Der Raum  $M$  besteht aus allen Funktionen aus  $C[a, b]$ , die auf den

Teilintervallen  $[x_k, x_{k+1}], k = 0, \dots, n-1$  linear sind.

Eine Funktion  $M$  ist durch ihre Werte in den Punkten  $x_0, \dots, x_n$

eindeutig festgelegt. Eine Basis von  $M$  bilden die Funktionen  $\varphi_0, \dots, \varphi_n$

mit  $\varphi_i(x_k) = \delta_{ik}$ . Diese Funktionen haben lokale Träger.

Damit ist  $\langle \varphi_i, \varphi_j \rangle = 0$ , falls  $|i - j| > 1$  ist. Die Koeffizientenmatrix

ist damit eine Tridiagonalmatrix. Sie ist

$$A = h \begin{pmatrix} \frac{1}{3} & \frac{1}{6} & 0 & \cdots & 0 \\ \frac{1}{6} & \frac{2}{3} & \ddots & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & \ddots & \frac{2}{3} & \frac{1}{6} \\ 0 & \cdots & 0 & \frac{1}{6} & \frac{1}{3} \end{pmatrix}$$

Sie ist sehr gut Konditioniert. Ihre Zeilensummennorm ist

$$\|A\| = \frac{h}{6} + \frac{2}{3}h + \frac{h}{6} = h$$

Es ist

$$\frac{h}{3}x_1 = Ax|_1 - \frac{h}{6}x_2$$

$$\frac{2}{3}hx_i = Ax|_i - \frac{h}{6}x_{i-1} - \frac{h}{6}x_{i+1}, i = 1, \dots, n-1$$

$$\frac{h}{3}x_n = Ax|_n - \frac{h}{6}x_{n-1}$$

Daraus folgt

$$|x_1| \leq \frac{3}{h}\|Ax\|_\infty + \frac{1}{2}\|x\|_\infty$$

$$|x_i| \leq \frac{3}{2h}\|Ax\|_\infty + \frac{1}{2}\|x\|_\infty$$

$$|x_n| \leq \frac{3}{h}\|Ax\|_\infty + \frac{1}{2}\|x\|_\infty$$

und damit

$$\|x\|_\infty \leq \frac{3}{h}\|Ax\| + \frac{1}{2}\|x\| \text{ oder}$$

$$\|x\| \leq \frac{6}{h}\|Ax\|, \text{ d.h. } \|A^{-1}\| \leq \frac{6}{h}$$

$$\kappa_\infty(A) = \|A\|\|A^{-1}\| \leq 6$$

## 5.2 Lineare Ausgleichsprobleme im $\mathbb{R}^n$

### Die QR Zerlegung einer Matrix

Beispiel: Gesucht ist die beste Approximation einer Funktion  $f$  im Sinne des Abstandsmaßes

$$\left(\sum_{i=1}^n (f(x_i) - \varphi(x_i))^2\right)^{\frac{1}{2}}$$

durch eine Funktion  $\varphi$  in einem endlich dimensionalen Teilraum.

Abstrakter: gesucht ist die beste Approximation des Vektors

$$(f(x_1), \dots, f(x_n))^T \in \mathbb{R}^n \text{ durch einen Vektor } (\varphi(x_1), \dots, \varphi(x_n))^T$$

aus einem gegebenen Teilraum des  $\mathbb{R}^n$ .

allgemein: Gegeben sei eine  $(m \times n)$ -Matrix  $A$  (mit  $m$  Zeilen und  $n \leq m$  Spalten) eines Rangs  $\leq n$  und ein Vektor  $b \in \mathbb{R}^m$ . Gesucht ist die beste Approximation von  $b$  durch eine Linearkombination der Spalten von  $A$ .  
Summe des euklidischen Skalarprodukts, also ein Vektor  $x^* \in \mathbb{R}^n$  mit

$$\|Ax^* - b\|_2 \leq \|Ax - b\|_2, x \in \mathbb{R}^n.$$

#### Satz 5.11

Hat die Matrix  $A$  den Rang  $n$ , sind also die Spalten von  $A$  linear unabhängig, so ist  $x^*$  eindeutig bestimmt und die Lösung des Gleichungssystems

$$A^\top Ax^* = A^\top b.$$

Beweis: Seien  $a_1, \dots, a_n \in \mathbb{R}^m$  die Spalten von  $A$  und  $M$  der von diesen Spalten aufgespannte Teilraum des  $\mathbb{R}^m$ , d.h. das Bild des  $\mathbb{R}^n$  unter  $A$ .

Gesucht sind dann die eindeutig bestimmten Koeffizienten  $x_i$  in der

Linearkombination  $\sum_{i=1}^n x_i a_i \in M$ , für die für alle  $v \in M$

$$\left\| \sum_{i=1}^n x_i a_i - v \right\|_2 \leq \|v - b\|_2 \text{ ist.}$$

Nach Satz 5.8 erfüllen die  $x_i$  das Normalgleichungssystem

$$\sum_{i=1}^n \langle a_i, a_j \rangle x_i = \langle b, a_j \rangle, j = 1, \dots, n$$

das in Matrixform mit  $x = \langle x_1, \dots, x_n \rangle^\top$

$$A^* Ax = A^* b$$

lautet

□

Bemerkung: Numerisch ist es nicht immer günstig diese Charakterisierung zu nutzen.

### Beispiel 5.12

Für die  $(6 \times 5)$  Matrix

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ \varepsilon & 0 & \cdots & \cdots & 0 \\ 0 & \varepsilon & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & \cdots & \cdots & 0 & \varepsilon \end{pmatrix}$$

vom Rang 5 ist

$$A^\top A = \begin{pmatrix} 1 + \varepsilon^2 & 1 & \cdots & \cdots & 1 \\ 1 & 1 + \varepsilon^2 & \ddots & & \vdots \\ \vdots & \ddots & 1 + \varepsilon^2 & \ddots & \vdots \\ \vdots & & \ddots & 1 + \varepsilon^2 & 1 \\ 1 & \cdots & \cdots & 1 & 1 + \varepsilon^2 \end{pmatrix}$$

Für  $\varepsilon = \frac{1}{2} \cdot 10^{-5}$  ist  $\varepsilon^2 = \frac{1}{4} \cdot 10^{-10}$  und deshalb in 10-stelliger Arithmetik:

$$\text{gl}(A^T A) = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 \end{pmatrix}$$

eine Matrix vom Rang 1.

□

#### Weiteres Problem:

Der Rang von  $A$  kann kleiner als  $n$  sein. Die optimale Approximation von  $b$  durch die Spalten von  $A$  ist dann nach wie vor eindeutig, nicht aber die Koeffizientenmatrix.

#### **Satz 5.13**

Unter allen Vektoren  $x \in \mathbb{R}^n$  mit  $\|Ax - b\|_2 = \min_{y \in \mathbb{R}^n} \|Ay - b\|_2$

gibt es genau einen Vektor  $x^*$  kürzester euklidischer Länge.

Beweis: Sei  $M \subseteq \mathbb{R}^m$  der Raum, der von den Spalten von  $A$  aufgespannt wird.

Dann gibt es nach Satz 5.8 genau ein  $y^*$  mit  $\|y^* - b\|_2 \leq \|y - b\|_2$  für alle  $y \in M$ . Es gilt genau dann  $\|Ax^* - b\|_2 \leq \|Ax - b\|_2$  für alle  $x \in \mathbb{R}^n$ ,

wenn  $Ax^* = y^*$  ist. Ist  $x_0^*$  ein derartiges  $x^*$ , so ist die Menge aller dieser  $x^*$  von der Form  $x^* = x_0^* - v, v \in \text{Kern}(A)$ . Der Kern von  $A$  ist ein endlichdimensionaler Teilraum des  $\mathbb{R}^n$ .

Die Länge von  $x^*$  wird also genau dann minimal, wenn  $v$  die eindeutig bestimmte beste Approximation von  $x_0^*$  durch ein Element von  $\text{Kern}(A)$  ist.

Aufgabe: Man löse das Ausgleichsproblem  $\|Ax - b\|_2 \rightarrow \min$  in numerisch stabiler Weise!

#### **Definition 5.14**

Eine  $QR$ -Zerlegung der  $(m \times n)$ -Matrix  $A$  ist eine Zerlegung von  $A$  in das Produkt  $A = QR$  einer orthogonalen  $(m \times m)$ -Matrix  $Q$  und eine obere Dreiecksmatrix der Dimension  $m \times n$ .

Beobachtung 1: Die ersten  $k$  Spalten,  $k \leq n$  von  $A$  lassen sich als Linearkombination der ersten  $k$  Spalten von  $Q$  darstellen.

Beobachtung 2:  $\text{Rang}(A) = \text{Rang}(R)$

Beobachtung 3: Hat  $A$  vollen Rang, so sind die ersten  $n$  Spalten von  $Q$  bis auf das Vorzeichen eindeutig bestimmt.

Berechnung der ersten  $n$  Spalten von  $Q$ :

durch Orhtogonalisierung nach Gram-Schmidt?

leider instabil!

Lösung des Ausgleichsproblems mittels einer  $QR$ -Zerlegung von  $A$ :

$$\begin{aligned}\|Ax - b\|_2 &= \|QRx - b\|_2 \\ &= \|Q(Rx - Q^\top b)\|_2 \quad (\text{Wegen } QQ^\top = I) \\ &= \|Rx - Q^\top b\|_2 \quad (\text{wegen } \|Qv\|_2 = \|v\|_2)\end{aligned}$$

Sei nun

$$R = \begin{pmatrix} R_1 & \\ & 0 \end{pmatrix} \begin{matrix} n \\ m-n \end{matrix}, Q^\top b = \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \begin{matrix} n \\ m-n \end{matrix}$$

Dann ist genau  $\|Ax^* - b\|_2 \leq \|Ax - b\|_2$  für alle  $x \in \mathbb{R}^n$ , wenn#

$\|R_1 x^* - c_1\|_2 \leq \|R_1 x - c_1\|_2$  für alle  $x \in \mathbb{R}^n$  und genau dann

$\|Ax^* - b\|_2 \geq \|c_2\|_2$ , wenn  $R_1 x^* = c_1$  ist.

Bemerkung: Ist  $A$  eine  $(n \times n)$ -Matrix  $A = QR$  eine  $QR$  Zerlegung von  $A$ , so kann man das lineare Gleichungssystem  $Ax = b$  wie folgt lösen:

- 1) Löse  $Qy = b$ , d.h. setze  $y = Q^\top b$
- 2) Löse  $Rx = y$

### Definition 5.15

Sei  $v \in \mathbb{R}^m$  und  $\|v\|_2 = 1$ . Dann definiert  $Hx = x - 2\langle v, x \rangle v$ ,  $x \in \mathbb{R}^m$ , die zu  $v$  gehörige Householdertransformation.

Sie hat die Matrixdarstellung  $H = I - 2vv^\top$

geometrische Interpretation:

$H$  ist die Spiegelung an der Hyperebene

$$\mathcal{H} = \{x \in \mathbb{R}^m \mid \langle v, x \rangle = 0\}$$

Für alle  $x \in \mathbb{R}^m$  ist nämlich

$$x = \underbrace{x - \langle v, x \rangle v}_{\in \mathcal{H}} + \underbrace{\langle v, x \rangle v}_{\in \mathcal{H}^\perp}$$

$$Hx = \underbrace{x - \langle v, x \rangle v}_{\in \mathcal{H}} - \underbrace{\langle v, x \rangle v}_{\in \mathcal{H}^\perp}$$

**Satz 5.16**

Householdertransformationen sind orthogonal und selbstadjungiert:  
Ist  $\|v\|_2 = 1$ , so gelten für die Matrix  $H = I - 2vv^\top$  die Beziehungen  
 $H^\top H = I, H^\top = H$

Beweis:

$$H|_{ij} = \delta_{ij} - 2v_i v_j = H|_{ij}$$

$$\begin{aligned} H^\top H x &= H H x = (x - 2\langle v, x \rangle v) - 2\langle x - 2\langle v, x \rangle v, v \rangle v \\ &= x - 2\langle v, x \rangle v - 2\langle v, x \rangle v + 4\langle v, x \rangle \langle v, v \rangle v = x \end{aligned}$$

Wegen  $\langle v, v \rangle = 1$ .

□

**Satz 5.17**

Sei  $a \in \mathbb{R}^m, a \neq 0$ , gegeben. Sei  $w = a + \sigma \|a\|_2 e_1$  mit  $\sigma = +1$ , falls  $a|_1 = a_1 \geq 0$  ist, und  $\sigma = -1$ , falls  $a|_1 < 0$  ist. Dann ist  
 $\|w\|_2^2 = 2\|a\|_2(\|a\|_2 + |a_1|)$ . Setzt man nun  $H = I - 2vv^\top, v = \frac{1}{\|w\|_2} w$ ,  
so gilt:  $Ha = -\sigma \|a\|_2 e_1$

Beweis:  $\|w\|_2^2 = \langle a + \sigma \|a\|_2 e_1, a + \sigma \|a\|_2 e_1 \rangle$

$$\begin{aligned} &= \langle a, a \rangle + 2\sigma \|a\| \langle a, e_1 \rangle + \sigma^2 \|a\|^2 \langle e_1, e_1 \rangle \\ &= \|a\|^2 + 2|a_1| \|a\| + \|a\|^2 \\ &= 2\|a\|(\|a\| + |a_1|) \\ &\geq 2\|a\|^2 \\ &> 0 \end{aligned}$$

$$\begin{aligned} Ha &= a - 2\langle a, v \rangle v \\ &= a - 2 \frac{\langle a, v \rangle}{\|w\|_2^2} w \\ &= a - 2 \frac{\langle a, a \rangle + \sigma \|a\| \langle a, e_1 \rangle}{2\|a\|(\|a\| + |a_1|)} w \\ &= a - 2 \frac{\|a\|^2 + |a_1| \|a\|}{2\|a\|(\|a\| + |a_1|)} w \\ &= a - w \\ &= a - (a + \sigma \|a\| e_1) \\ &= -\sigma \|a\| e_1 \end{aligned}$$

□



Bemerkung: Bei der Berechnung dieser Transformation kann keine Auslöschung durch Addition von Zahlen unterschiedlichen Vorzeichens auftreten!

### Berechnung der QR-Zerlegung einer Matrix

Falls die erste Spalte von  $A$  ungleich 0 ist, wähle eine Householder-Transformation  $H$ , die sie in ein Vielfaches des 1. Einheitsvektors transformiert.

(Ist die erste Spalte = 0, setze man  $H = I - 2vv^T, v = 0$ )

Schritt  $k + 1, k < n$

Man wähle eine Householdermatrix  $H_{k+1} = I - 2vv^T$  mit einem Vektor  $v = (0, \dots, 0, v_{k+1}, \dots, v_m)^T$ , die die erste Spalte der Restmatrix in einen Vektor  $x = (x_1, \dots, x_k, x_{k+1}, 0, \dots, 0)^T$  transformiert.

(Gehe nach Satz 5.17 vor, falls die erste Spalte der Restmatrix  $\neq 0$  ist, sonst setze  $v = 0$ )

Ergebnis: Orthogonale Matrizen  $H_1, \dots, H_n$  der Form  $H_k = I - 2vv^T$  mit  $\|v\|_2 = 1$  oder  $\|v\|_2 = 0$  und eine obere Dreiecksmatrix  $R$  und  $H_n \cdots H_1 A = R$  oder  $A = (H_n \cdots H_1)^T R = H_1 \cdots H_n R = QR$

Bemerkung:  $Q$  wird nicht explizit berechnet. Man speichert Stattdessen die Vektoren  $v$  ab, die die Householdertransformationen bestimmen.

Speicheraufwand:

$mn + n$  Speicherplätze für die Faktoren von  $Q$  und für  $R$ .

Aufwand  $k$ -ter Schritt:

Berechnung von  $H = I - 2vv^T$  und  $Ha = -\text{sgn}(a_1)\|a\|_2 e_1$ :

$\|a\|_2$  :  $m - k + 1$  Multiplikationen, 1 Quadratwurzel

$w = a + \text{sgn}(a_1)\|a\|_2 e_1$  : -

$\|w\|_2 = (2\|a\|_2(\|a\| + |a_1|))^{1/2}$  : 1 Quadratwurzel

$r = \frac{1}{\|w\|_2} w$  :  $m - k + 1$  Divisionen

$\Sigma = 2(m - k + 1)$  Punktoperationen, 2 Quadratwurzeln

Multiplikation der restlichen  $n - k$  Spalten mit  $H$

$Ha = a - 2\langle a, v \rangle v$  jeweils  $2(m - k + 1)$  Punktoperationen

Gesamtaufwand:

$$\sum_{k=1}^n [2(m-k+1) + 2(m-k+1)(n-k)]$$

$$= (m - \frac{n-1}{3})(n^2 + n)$$

numerische Stabilität:

ganz hervorragend! Bei quadratischen Matrizen viel besser als  $LR$ -Zerlegung, ja sogar besser als bei der Zerlegung symmetrischer positiv definiten Matrizen.

Lösung linearer Gleichungssysteme

$$Ax = b \iff x = R^{-1}Q^T b$$

Aufwand:

$$QR\text{-Zerlegung: } \frac{2}{3}n^3 + n^2 + \mathcal{O}(N)$$

$$Q^T b =: y \quad : n^2$$

$$R^{-1}y \quad : \frac{1}{2}n^2$$

asymptotisch doppelt so hoch wie bei der  $LR$ -Zerlegung

aber:

ungemein stabil!

Householder-Orthogonalisierung mit Pivotsuche:

$$H_n \cdots H_1 A P_1 \cdots P_{n-1} = R$$

$$A = H_1 \cdots H_n R P_{n-1} \cdots P_1$$

Wähle die Transpositionsmatrix  $P_{k+1}$  (Vertauschung zweier Komponenten) so, dass die erste Spalte der Restmatrix die längste Spalte der Restmatrix ist.

## 5.3 Singulärwertzerlegung und Pseudoinverse

Problem: Wie findet man den Vektor kürzester euklidischer Länge, der  $\|Ax - b\|_2$  minimiert?

### Satz 5.18

Zu jeder  $(m \times n)$ -Matrix  $A, n \leq m$  gibt es eine orthogonale  $(m \times m)$ -Matrix  $U$  und eine orthogonale  $(n \times n)$ -Matrix  $V$  sowie eine Diagonalmatrix  $D$  mit  $A = UDV^T$ . Diese Zerlegung heißt eine Singulärwertzerlegung von  $A$ .

$$D = \text{diag}(\sigma_1, \dots, \sigma_n)$$

Beweis: Seien  $x \in \mathbb{R}^n$  und  $y \in \mathbb{R}^m$  Vektoren mit  $\|x\|_2 = 1, \|y\|_2 = 1$  und  $Ax = ry$ , wobei  $\sigma = \max_{v \neq 0} \frac{\|Av\|_2}{\|v\|_2}$  die Norm von  $A$  sein soll.

Sei  $x, v_2, \dots, v_n$  eine orthonormalbasis des  $\mathbb{R}^n$  und  $V_1$  die  $(n \times n)$ -Orthogonalmatrix, deren Spalten diese Vektoren sind.

Sei  $y, u_2, \dots, u_m$  eine Orthonormalbasis des  $\mathbb{R}^m$  und  $U_1$  die  $(m \times m)$ -Orthogonalmatrix, deren Spalten diese Vektoren sind. Dann ist wegen  $Ax = \sigma y$ :

$$U_1^\top AV_1 = \begin{pmatrix} \sigma & | & w^\top \\ 0 & | & B \end{pmatrix} \begin{matrix} 1 \\ m-1 \end{matrix}$$

Bezeichne  $\begin{pmatrix} \sigma \\ w \end{pmatrix}$  die Aufspaltung des Vektors  $w$  in 2 Teile

$$\text{Wegen } \left\| \begin{pmatrix} \sigma \\ w \end{pmatrix} \right\|_2^2 = \sigma^2 + \|w\|_2^2$$

$$\|V_1^\top AV_1 \begin{pmatrix} \sigma \\ w \end{pmatrix}\|_2^2 = \left\| \begin{pmatrix} \sigma^2 + w^\top w \\ * \end{pmatrix} \right\|_2^2 \geq \sigma^2 + \|w\|_2^2 \text{ ist}$$

$$\|A\|_2^2 = \|U_1^\top AV_1\|_2^2 \geq \sigma^2 + \|w\|_2^2$$

Wegen  $\|A\|_2 = \sigma$  folgt daraus  $\|w\| = 0$ , also

$$U_1^\top AV_1 = \begin{pmatrix} \sigma & | & 0 \\ 0 & | & B \end{pmatrix} \begin{matrix} 1 \\ m-1 \end{matrix}$$

Durch Induktion folgt daraus die Existenz einer orthogonalen  $(m \times m)$ -Matrix  $U$ , einer orthogonalen  $(n \times n)$ -Matrix  $V$  und einer Diagonalmatrix  $D$  mit  $U^\top AV = D$ .

Bemerkung: Der Spektralsatz besagt, dass sich jede symmetrische Matrix  $A \in \mathbb{R}^{n \times n}$  in der Form  $A = UDU^\top$  einer Diagonalmatrix  $D$  und einer orthogonalen Matrix  $U$  schreiben lässt.

Daher lässt sich Satz 5.18 als Verwandter des Spektralsatzes deuten.

### Definiton 5.19

Ist  $A = UDV^\top$  eine Singulärwertzerlegung der  $(m \times n)$  Matrix  $A$ , so ist die  $(n \times m)$  Matrix  $A^+ = VD^+U^\top$  eine Pseudoinverse von  $A$ .

Also ist  $D = \begin{pmatrix} D_1 \\ 0 \end{pmatrix}$  mit  $D_1 = \text{diag}(\sigma_1, \dots, \sigma_n)$  so ist  $D^+$  gegeben durch

$$D^+ = (D_1^+, 0) \text{ mit } D_1^+ = \text{diag}(\mu_1, \dots, \mu_n)$$

$$\text{und } \mu_i = \begin{cases} \frac{1}{\sigma_i}, & \sigma_i \neq 0 \\ 0, & \text{sonst} \end{cases}$$

Jetzt können wir die definitive Lösung unsere Ausgleichsproblems angeben.

### Satz 5.20

Ist  $A^+$  eine Pseudoinverse der  $(m \times n)$  Matrix  $A$  und  $b \in \mathbb{R}^m$ , so ist  $x^* = A^+b$  der Vektor  $x$  kürzester euklidischer Länge mit

$$\|Ax - b\|_2 = \inf_{y \in \mathbb{R}^n} \|Ay - b\|_2$$

Beweis: Wegen  $\|Ax - b\|_2 = \|UDV^\top x - b\|_2 = \|U(DV^\top x - U^\top b)\|_2$   
 $= \|DV^\top x - U^\top b\|_2$

ist  $x^*$  genau dann die gesuchte Lösung kürzester Länge, wenn  $y^* = V^\top x^*$  den Vektor kürzester Länge mit  $\|Dy - U^\top b\|_2 = \min_{z \in \mathbb{R}^n} \|Dz - U^\top b\|_2$  ist.

Wegen  $\|Dy - U^\top b\|_2 = \min!$  Gilt also  $y^* = D^+U^\top b$  und somit  
 $x^* = Vy^* = VD^+U^\top b = A^+b$

□

### Satz 5.21

Jede  $(m \times n)$  Matrix  $A$ ,  $n \leq m$  besitzt genau eine Pseudoinverse.

Beweis: Nach Satz 5.20 ist für jedes  $b \in \mathbb{R}^n$   $x^* = A^+b$  ein Vektor  $x$  kürzester Länge mit  $\|Ax - b\|_2 = \inf_{y \in \mathbb{R}^n} \|Ay - b\|_2$  und nach Satz 5.13 ist dieser Vektor eindeutig bestimmt.

□

Bemerkung: Für  $n > m$  ist die Pseudoinverse von  $A^\top$  eindeutig und damit auch von  $A$ .

### Alternative Darstellung der Singulärwertzerlegung

$$A = UDV^\top = \sum_{i=1}^n \sigma_i u_i v_i^\top$$

falls  $U = [u_1, \dots, u_n, *, \dots, *]$

$$V = [v_1, \dots, v_n]$$

$$D = \text{diag}(\sigma_1, \dots, \sigma_n)$$

### Annahme:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1}, \dots, \sigma_n = 0$$

Dann gilt

$$A = \sum_{i=1}^r \sigma_i u_i v_i^\top$$

$$A^+ = \sum_{i=1}^r \frac{1}{\sigma_i} v_i u_i^\top$$

Frobeniusnorm:  $\|A\|_F = \left(\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2\right)^{1/2}$

### Satz 5.22

Die bestmögliche Approximation der  $(m \times n)$  Matrix  $A$  mit der SWZ

$$A = \sum_{i=1}^{\min(n,m)} \sigma_i u_i v_i^\top \text{ mit } \sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min(n,m)} \geq 0 \text{ durch eine Matrix}$$

$B$  vom Rang  $k, k \leq \min(n, m)$ , im Sinne der Frobeniusnorm ist.

$$B^* = \sum_{i=1}^k \sigma_i u_i v_i^\top.$$

Beweis: Es ist  $\|A - B\|_F = \|UDV^\top - B\|_F = \|U(D - U^\top BV)V^\top\|_F$ .

Da für orthogonale Matrizen  $U$  und  $V$

$$\|UA\|_F = \|A\|_F, \|AV\|_F = \|A\|_F$$

gilt, folgt  $\|A - B\|_F = \|D - U^\top BV\|_F$

Die Beste Approximation von  $D$  durch eine Matrix vom Rang  $k$  ist

offensichtlich  $U^\top BV = \sum_{i=1}^k \sigma_i e_i e_i^\top$  woraus

$$B = \sum_{i=1}^k \sigma_i U e_i e_i^\top V^\top = \sum_{i=1}^k \sigma_i u_i v_i^\top \text{ folgt.}$$

□

### Berechnung einer SWZ

#### Schritt 1

Reduktion der Matrix  $A$  auf Bidiagonalform mit Hilfe von Householder Matrizen. Finiter Prozess

#### Schritt 2

Berechnung der SWZ der Bidiagonalmatrix mit dem QR-Algorithmus  
iterativer Prozess

Householderschritt 1

$$\begin{pmatrix} x & x & x & x \\ x & x & x & x \\ x & x & x & x \\ x & x & x & x \end{pmatrix} \xrightarrow{L} \begin{pmatrix} x & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \end{pmatrix} \xrightarrow{R} \begin{pmatrix} x & x & 0 & 0 \\ 0 & x & x & x \\ 0 & x & x & x \\ 0 & x & x & x \end{pmatrix} \xrightarrow{L} \begin{pmatrix} x & x & 0 & 0 \\ 0 & x & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{pmatrix} \xrightarrow{R} \begin{pmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{pmatrix} \\
\xrightarrow{L} \begin{pmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & 0 & x \end{pmatrix} \xrightarrow{L} \begin{pmatrix} x & x & 0 & 0 \\ 0 & x & x & 0 \\ 0 & 0 & x & x \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

$L, R$  : Multiplikation mit Householdermatrizen von links bzw. rechts.

Anwendung: "information retrieval"

Finde aus einer Sammlung von  $n$  Dokumenten die Dokumente heraus, die Informationen zu einem gegebenen Thema enthalten. Dazu muss die Information durch Schlagworte oder Suchterme kodifiziert werden.

Vektorraummodell

Stelle eine Kollektion von  $m$  möglichen Suchtermen zusammen.

Dokumente wie Anfragen werden als Vektoren im  $\mathbb{R}^m$  aufgefasst.

Die Komponente  $a_i$  eines solchen Vektors gibt an, wie oft der Suchterm, im gegebenen Dokument vorkommt, bzw. ob dieser Term gesucht wird.

Datenbasis

Eine  $(m \times n)$  Matrix  $A$ ; die Spalten von  $A$  entsprechen den einzelnen Dokumenten in der Datenbasis.

Bsp.: Internet Suchmaschine

$m = 300.000$  (Worte in engl. Sprache)

$n \approx 3 \cdot 10^9$  (Seiten im Netz)

ideale Suchmaschine

Berechne zur Anfrage  $q \in \mathbb{R}^m$  den Zeilenvektor  $q^\top A \in \mathbb{R}^n$ .

Liste die von Null verschiedenen Komponenten des Vektors aus.

Faktisch unmöglich!

Stattdessen: Ersetze  $A$  durch Rang  $k$  Approximation  $A_k = \sum_{i=1}^k \sigma_i u_i v_i^\top$  und

berechne  $q^\top A_k = \sum_{i=1}^k \sigma_i (q^\top u_i) v_i^\top$

## 5.4 Das Gauß-Newton-Verfahren

### nichtlineares Ausgleichsproblem

Sei  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $n \leq m$  gegeben. Finde  $x^* \in \mathbb{R}^n$  mit  
 $\|f(x^*)\|_2 \leq \|f(x)\|_2, x \in \mathbb{R}^n$

#### Spezialfall

$f(x) = Ax - b$ , lin. Ausgleichsproblem

#### Linearisierung um Näherungslösung $x_k$ :

$$f(x) \approx f(x_k) + f'(x_k)(x - x_k)$$

#### Gauß-Newton-Verfahren

Finde neue Näherung  $x_{k+1}$  als Lösung von

$$\|f(x_k) + f'(x_k)(x - x_k)\|_2 = \min!$$

#### formal:

$$x_{k+1} = x_k - f'(x_k)^+ f(x_k)$$

## 6. Interpolation

”Darstellung von Funktionen durch ihre Werte in vorgegebenen Punkten.”

#### gegeben:

eine von reellen Parametern  $\alpha_0, \dots, \alpha_n$  abhängige Funktionenklasse

$$\varphi(x; \alpha_0, \dots, \alpha_n)$$

und  $n + 1$  Stützstellen  $x_0, \dots, x_n$  im Definitionsbereich von  $\varphi$ .

#### gesucht:

(die?) Werte  $\alpha_0, \dots, \alpha_n$  mit  $\varphi(x_i; \alpha_0, \dots, \alpha_n) = f_i$  für  $i = 0, \dots, n$  für vorgegebene Werte  $f_i (= f(x_i))$

#### hier:

linearer Ansatzfunktionenraum

$$\varphi(x; \alpha_0, \dots, \alpha_n) = \sum_{i=1}^n \alpha_i \varphi_i(x)$$

$\varphi_0, \dots, \varphi_n$  gegeben.  $\rightarrow$  lineares Interpolationsproblem

Interpolation durch Polynome und stückweise Polynome (Splines)

Existenz?

Eindeutigkeit?

Fehler?

### Satz 6.1

Folgende Aussagen sind äquivalent

i) Zu beliebig vorgegebenen Werten  $f_i$  gibt es eindeutig bestimmte Werte  $\alpha_i$  mit

$$\sum_{i=0}^n \alpha_i \varphi_i(x_j) = f_j, \quad j = 0, \dots, n$$

ii) Zu beliebig vorgegebenen Werten  $f_i$  gibt es Werte mit

$$\sum_{i=0}^n \alpha_i \varphi_i(x_j) = f_j, \quad j = 0, \dots, n$$

iii) Aus

$$\sum_{i=0}^n \alpha_i \varphi_i(x_j) = 0, \quad j = 0, \dots, n$$

folgt  $\alpha_0, \dots, \alpha_n = 0$ .

Beweis: Für die Werte  $\alpha_0, \dots, \alpha_n$  gilt genau dann

$$\sum_{i=0}^n \alpha_i \varphi_i(x_j) = f_j$$

für  $j = 0, \dots, n$ , wenn die  $\alpha_i$  das lineare Gleichungssystem

$$\begin{pmatrix} \varphi_0(x_0) & \cdots & \varphi_n(x_0) \\ \vdots & \ddots & \vdots \\ \varphi_0(x_n) & \cdots & \varphi_n(x_n) \end{pmatrix} \begin{pmatrix} \alpha_0 \\ \vdots \\ \alpha_n \end{pmatrix} = \begin{pmatrix} f_0 \\ \vdots \\ f_n \end{pmatrix}$$



erfüllen. Damit lässt sich die Behauptung auf eine fundamentale Strukturaussage über lineare Gleichungssysteme zurückführen.  $\square$

## 6.1 Interpolation durch Polynome

### Satz 6.2

Die Stützstellen  $x_0, \dots, x_n \in \mathbb{R}$  seien paarweise verschieden. Dann gibt es ein eindeutig bestimmtes Polynom

$$P(x) = \sum_{i=0}^n \alpha_i x^i$$

vom Grad  $\leq n$ , das in den Stützstellen  $x_j$  vorgegebene Werte  $f_j$  annimmt. Dieses Polynom besitzt die Darstellung

$$P(x) = \sum_{j=0}^n f_j l_j(x)$$

mit den Polynomen

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

den *Lagrangeschen Grundpolynomen* zu den Stützstellen  $x_0, \dots, x_n$ .

Beweis: Es ist  $l_j(x_j) = 1$  und  $l_j(x_k) = 0$  für  $k \neq j$ . Damit ist

$$\sum_{j=0}^n f_j l_j(x_k) = \sum_{j=0}^n f_j \delta_{jk} = f_k, \quad k = 0, \dots, n$$

also eine Lösung des Interpolationsproblems gefunden. Die Eindeutigkeit folgt aus Satz 6.1.  $\square$

### Beispiel 6.3

Das Interpolationspolynom durch die Wertepaare  $(-5, f(-5)), (-1, f(-1)), (5, f(5)), (1, f(1))$  und der Funktion

$f(x) = \frac{1}{1+x^2}$  lautet

$$P(x) = \frac{1}{26} l_0(x) + \frac{1}{2} l_1(x) + \frac{1}{2} l_2(x) + \frac{1}{26} l_3(x)$$

mit

$$l_0(x) = \frac{(x+1)(x-1)(x-5)}{(-5+1)(-5-1)(-5-5)} = \frac{1}{240}(x+1)(x-1)(x-5).$$

anderer Beweis für die Eindeutigkeit:

Gilt

$$\sum_{i=0}^n \alpha_i x_j^i = \sum_{i=0}^n \tilde{\alpha}_i x_j^i = f_j, \quad j = 0, \dots, n$$

so hat das Polynom

$$\sum_{i=0}^n (\alpha_i - \tilde{\alpha}_i) x^i$$

$n+1$  Nullstellen  $x_1, \dots, x_n$ . Nach dem Fundamentalsatz der Algebra ist daher  $\alpha_i - \tilde{\alpha}_i = 0$ ,  $i = 1, \dots, n$ .

Newtonsche Darstellung

$$P(x) = a_0 + \sum_{i=1}^n a_i \prod_{j=0}^{i-1} (x - x_j)$$

vom Rechenaufwand her wesentlich günstiger!

Voraussetzung: Für

$$P_m(x) = a_0 + \sum_{i=1}^m a_i \prod_{j=0}^{i-1} (x - x_j)$$

gelte

$$P_m(x_k) = f_k, \quad k = 0, 1, \dots, m.$$

Beobachtung:

Dann gilt für jedes Polynom  $P(x) = P_m(x) + a \prod_{j=0}^m (x - x_j)$  für  $k = 0, \dots, m$ :  $P(x_k) = f_k$ .

Bestimme daher  $a$  so, dass  $P(x_{m+1}) = f_{m+1}$  ist: also

$$a = a_{m+1} = \frac{f_{m+1} - P_m(x_{m+1})}{\prod_{j=0}^m (x_{m+1} - x_j)}.$$

Mit  $a_0 = f_0$  folgt die Existenz und die Eindeutigkeit der Darstellung.

Schreibweise:  $a_i = f[x_0, \dots, x_i]$

Die dividierte Differenz  $f[x_0, \dots, x_m]$  ist der Leitkoeffizient des Polynoms

$$P_m(x) = f[x_0] + \sum_{i=1}^m f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j).$$

### Satz 6.4

Das Interpolationsproblem zu den paarweise verschiedenen Stützstellen  $x_0, \dots, x_n$  und den Werten  $f_0, \dots, f_n$  besitzt die Newtonsche Darstellung

$$f[x_0] + \sum_{i=1}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j).$$

Die dividierten Differenzen  $f[x_0, \dots, x_i]$  lassen sich ausgehend von

$$f[x_i] = f_i, \quad i = 0, \dots, n$$

rekursiv über die Formel

$$f[x_k, \dots, x_i] = \frac{f[x_{k+1}, \dots, x_i] - f[x_k, \dots, x_{i-1}]}{x_i - x_k},$$

$0 \leq k < i \leq n$ , berechnen.

Beweis: Sei  $Q_1$  das Interpolationspolynom zu den Punkten  $x_0, \dots, x_{m-1}$  und  $Q_2$  das Interpolationspolynom zu den Punkten  $x_1, \dots, x_m$ . Dann ist

$$P(x) = \frac{1}{x_m - x_0} \{ (x - x_0)Q_2(x) + (x_m - x)Q_1(x) \}$$

das Interpolationspolynom zu den Punkten  $x_0, \dots, x_m$ . Sein Leitkoeffizient ist einerseits  $f[x_0, \dots, x_m]$  und andererseits

$$\frac{f[x_1, \dots, x_m] - f[x_0, \dots, x_{m-1}]}{x_m - x_0}$$

woraus die angegebene Rekursionsformel folgt.  $\square$

Bemerkung:

Die dividierte Differenz  $f[x_0, \dots, x_m]$  ist von der Anordnung der Stützstellen unabhängig.

Beweis: Sie ist der Leitkoeffizient des Polynoms

$$\sum_{i=0}^m f_i \prod_{\substack{j=0 \\ j \neq i}}^m \frac{x - x_j}{x_i - x_j}$$

und damit gegeben durch

$$\sum_{i=0}^m f_i \prod_{\substack{j=0 \\ j \neq i}}^m \frac{1}{x_i - x_j}.$$

□

### Dreiecksschema zur Berechnung der dividierten Differenzen

$$\begin{array}{l|l} x_0 & f[x_0] \\ x_1 & f[x_1] \quad f[x_0, x_1] = \frac{f[x_1] - f[x_0]}{x_1 - x_0} \\ x_2 & f[x_2] \quad f[x_1, x_2] = \frac{f[x_2] - f[x_1]}{x_2 - x_1} \quad f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0} \\ x_3 & f[x_3] \quad f[x_2, x_3] = \frac{f[x_3] - f[x_2]}{x_3 - x_2} \quad f[x_1, x_2, x_3] = \frac{f[x_2, x_3] - f[x_1, x_2]}{x_3 - x_1} \quad f[x_0, x_1, x_2, x_3] \end{array}$$

neues Wertepaar  $\Rightarrow$  neue Zeile hinzufügen.

### Beispiel 6.5 ....

Auswertung mit dem Horner-Schema:

$$P(x) = f[x_0] + \sum_{i=1}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j)$$

$$S_n = f[x_0, \dots, x_n]$$

$$S_i = f[x_0, \dots, x_i] + (x - x_i)S_{i+1}, \quad i = n - 1, \dots, 0$$

$$P(x) = S_0$$

### Satz 6.6

Die Funktion  $f : [a, b] \rightarrow \mathbb{R}$  sei  $(n + 1)$ -mal stetig differenzierbar. Das Polynom  $P$   $n$ -ten Grades interpoliere  $f$  in den  $n + 1$  Punkten  $a = x_0 < x_1 < \dots < x_n = b$

Dann gilt in  $x \in [a, b]$  für den Fehler  $f(x) - P(x) = \frac{f^{(n+1)}(\eta)}{(n+1)!} \prod_{j=1}^n (x - x_j)$

mit einer Zwischenstelle  $\eta \in (a, b)$ , die von der Wahl der Stützstellen der Funktion  $f$  und dem betrachteten Punkt  $x$  abhängt.

Beweis: Sei  $x$  fest und ohne Einschränkung  $x \neq x_j$  für  $j = 0, \dots, n$ .

Betrachte die Hilfsfunktion  $g(t) = (f(t) - P(t)) - \left(\prod_{j=0}^n \frac{t-x_j}{x-x_j}\right)(f(x) - P(x))$ .

Die Funktion  $g$  ist wie  $f$  selbst  $(n + 1)$ -mal differenzierbar und verschwindet in den  $n + 2$  Punkten  $x$  und  $x_0, \dots, x_n$ .

Nach dem Satz von Rolle verschwindet  $g'(t)$  dann in mindestens  $n + 1$  von einander verschiedenen Punkten, die zwischen je zwei Nullstellen von  $g(t)$  liegen.

Entsprechend verschwindet  $g''$  in mindestens  $n$  von einander verschiedenen Punkten aus  $(a, b)$  und  $g^{(n+1)}$  schließlich noch in einem Punkt  $\eta \in (a, b)$ .

Wegen  $g^{(n+1)}(t) = f^{(n+1)}(t) - \frac{(n+1)!}{\prod_{j=0}^n (x-x_j)}(f(x) - P(x))$  folgt daraus die

Behauptung.

□

Folgerung:

$$|f(x) - P(x)| \leq \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} \left| \prod_{j=0}^n (x - x_j) \right|, \text{ falls } \|f^{(n+1)}\|_\infty = \sup_{a < i < b} |f^{(n+1)}(t)|$$

existiert. Falls  $|x - x_j| \leq Kh$  ist, folgt  $|f(x) - P(x)| = \mathcal{O}(h^{n+1})$  für  $h \rightarrow 0+$ .

### Beispiel 6.7

Die Funktion  $f(x) = \ln(x)$  sei in  $x_j = 1 + jh, j = 0, 1, 2, \dots$  tabelliert.

stückweise lineare Interpolation

$$\frac{x_{j+1}-x}{x_{j+1}-x_j} f(x_j) + \frac{x-x_j}{x_{j+1}-x_j} f(x_{j+1}) \text{ für } x_j \leq x \leq x_{j+1}$$

$$\begin{aligned} & \left| \frac{f''(\eta)}{2!} (x - x_j)(x - x_{j+1}) \right| \leq \frac{1}{2} \left| \left( \frac{x_j + x_{j+1}}{2} - x_j \right) \left( \frac{x_j + x_{j+1}}{2} - x_{j+1} \right) \right| |f''(\eta)| \\ & = \frac{1}{8} (x_{j+1} - x_j)^2 |f''(\eta)| = \frac{1}{8} (x_{j+1} - x_j)^2 \left| -\frac{1}{\eta^2} \right| \leq \frac{h^2}{8} \end{aligned}$$

### stückweise quadratische Inteprolation

Interpolationspolynom zweiten Grades mit den Stützstellen  $x_{j-1}, x_j, x_{j+1}$

für  $x_j - \frac{h}{2} \leq x \leq x_j + \frac{h}{2}$ .

$$\begin{aligned} & \left| \frac{f^{(3)}(\eta)}{3!} (x - x_{j-1})(x - x_j)(x - x_{j+1}) \right| \leq \frac{h^3}{16} |f^{(3)}(\eta)| \\ & = \frac{h^3}{16} \left| \frac{2}{\eta^3} \right| \leq \frac{h^3}{8} = h \cdot \frac{h^2}{8} \end{aligned}$$

### Satz 6.8

Die Funktion  $f$  sei auf  $[a, b]$   $m$ -mal differenzierbar. Die Stützstellen  $x_0, \dots, x_n \in [a, b]$  seien von einander verschieden. Es sei  $f[x_2] = f(x_k)$  für  $k = 0, 1, \dots, m$ . Dann ist  $f[x_0, \dots, x_m] = \frac{f^{(m)}(\eta)}{m!}, \eta \in (a, b)$ .

Beweis: Einerseits gilt nach Satz 6.3 für das Interpolationspolynom  $P$  zu den Stützstellen  $x_0, \dots, x_m$  und jeden von diesen Stützstellen verschiedenen Punkt  $x$ :

$$\begin{aligned} f(x) &= f[x_0] + \sum_{i=1}^n f[x_0, \dots, x_i] \prod_{j=0}^{i-1} (x - x_j) + f[x_0, \dots, x_n, x] \prod_{j=0}^n (x - x_j) \\ &= P(x) + f[x_0, \dots, x_n, x] \prod_{j=0}^n (x - x_j) \end{aligned}$$

Andererseits ist nach Satz 6.6

$$f(x) = P(x) + \frac{f^{(n+1)}(\eta)}{(n+1)!} \prod_{j=0}^n (x - x_j)$$

Daraus folgt

$$f[x_0, \dots, x_n, x] = \frac{f^{(n+1)}(\eta)}{(n+1)!}$$

also die Behauptung für  $m = n + 1$  und  $x_{k+1} = x$ .

### Satz 6.9

Die Funktion  $f : [a, b] \rightarrow \mathbb{R}$  sei  $(n + 1)$ -mal differenzierbar und  $P$  sei das Interpolationspolynom von  $f$  in  $a = x_0 < x_1 < \dots < x_n = b$ . Dann gilt für  $m \leq n$  und beliebiges  $x \in [a, b]$

$$f^{(m)}(x) - P^{(m)}(x) = \frac{f^{(n+1)}(\eta)}{(n-m+1)!} \prod_{j=0}^{n-m} (x - \theta_j)$$

mit Punkten  $\eta, \theta_j \in (a, b)$ .

Beweis: Da  $f(x) - P(x)$  in den  $n + 1$  Punkten  $x_0, \dots, x_n$  verschwindet, gibt es nach dem Satz von Rolle voneinander verschiedene Punkte  $\theta_0, \dots, \theta_{n-m}$  mit  $f^{(m)}(\theta_i) - P^{(m)}(\theta_i) = 0$ .

Mit der Hilfsfunktion

$$g(t) = [f^{(m)}(t) - P^{(m)}(t)] - \left( \prod_{j=0}^{n-m} \frac{t - \theta_j}{x - \theta_j} \right) [f^{(m)}(x) - P^{(m)}(x)]$$

kann man argumentieren wie im Beweis von Satz 6.6 und findet durch mehrmalige Anwendung des Satzes von Rolle die Fehlerdarstellung.

□

Folgerung:

$$|f^{(m)}(x) - P^{(m)}(x)| \leq \frac{\|f^{(n+1)}\|}{(n-m+1)!} \left| \prod_{j=0}^{n-m} (x - \theta_j) \right|$$

$$|f^{(m)}(x) - P^{(m)}(x)| = 0(h^{n-m+1}) \text{ falls } |x - x_j| = \mathcal{O}(h)$$

Verlust einer Fehlerordnung pro Differentiationsordnung!

Bemerkung:

Wertet man die Ableitungen von Interpolationspolynomen an relativ zur Lage der Interpolationspunkte festen Stellen aus, erhält man Differenzenformeln zur Approximation von Ableitungen.

### Beispiel 6.10

Differenziert man das Interpolationspolynom

$$P(t) = \frac{(t-x)(t-x-h)}{(x-h-x)(x-h-x-h)}y(x-h) + \frac{(t-x)(t-x-h)}{(x-x+h)(x-x-h)}y(x) + \frac{(t-x+h)(t-x)}{(x+h-x)(x+h-x)}y(x+h)$$

von  $y(t)$  zu den Stützstellen  $x-h, x$  und  $x+h$ , so erhält man die Näherung

$$P'(t) = \dots$$

für  $y'(t)$ , speziell also die Näherung  $P'(x) = \frac{y(x+h) - y(x-h)}{2h}$  für  $y'(x)$ .

Der Fehler ist von der Größenordnung  $\mathcal{O}(h^2)$ .

□

Bemerkung: Neben den Funktionswerten kann man auch die Ableitungen der Funktion mit interpolieren: dann spricht man von HERMITE-Interpolation.

### Beispiel 6.11

Gesucht ist ein Polynom  $P$  vom Grad  $\leq 2n + 1$  mit

$$P(x_k) = f(x_k), k = 0, \dots, n$$

$$P'(x_k) = f'(x_k), k = 0, \dots, n$$

**Beispiel 6.12**

Gesucht ist ein Polynom  $P$   $n$ -ten Grades mit

$$P^{(k)}(x_k) = f^{(k)}(x_0), k = 0, \dots, n$$



## 6.2 Beste Approximation und optimale Stützstellen

### Globales Verhalten

Bei Vergrößerung der Stützstellenzahl neigen Interpolationspolynome zu Überschwingeffekten.

### Beispiel 6.13

$$f(x) = |x|, \quad -2 \leq x \leq 2$$

$$x_j = -2 + j \frac{4}{14}, \quad j = 0, \dots, 14$$

### Beispiel 6.14

$$f(x) = \frac{1}{1+x^2}, \quad -5 \leq x \leq 5$$

$$x_j = -5 + j \frac{10}{14}, \quad j = 0, \dots, 14$$

### Definition 6.15

Für eine stetige Funktion  $f : [a, b] \rightarrow \mathbb{R}$  ist

$$E_n(f) = \inf \left\{ \max_{a \leq x \leq b} |f(x) - P(x)| \mid P \text{ Polynom } n\text{-ten Grades} \right\}$$

Ein Polynom  $P^*$  vom Grade  $n$  heißt Polynom bester Approximation an  $f$ , wenn  $\max_{a \leq x \leq b} |f(x) - P^*(x)| = E_n(f)$  gilt.  $E_n$  heißt Fehler der besten Approximation.

### Satz 6.16

Zu jeder stetigen Funktion  $f$  existiert genau ein eindeutig bestimmtes Polynom  $n$ -ten Grades bester Approximation an  $f$ .

Bemerkung: Dieses Polynom hängt nichtlinear von  $f$  ab.

### Satz 6.17

Interpoliert das Polynom  $P_n$  die Funktion  $f$  in den Stützstellen

$$a \leq x_0 < x_1 < \dots < x_n \leq b \text{ so gilt } \max_{a \leq x \leq b} |f(x) - P(x)| \leq (1 + \lambda_n) E_n(f)$$

Dabei ist  $\lambda_n = \max_{a \leq x \leq b} \sum_{j=0}^n |l_j(x)|$ ,  $l_j(x) = \prod_{\substack{i=0 \\ i \neq j}}^n \frac{x-x_i}{x_j-x_i}$  die Lebesgue-Konstante

des Interpolationsoperators.

Beweis: Für alle Polynome  $P$   $n$ -ten Grades ist

$$\begin{aligned}
f(x) - P_n(x) &= f(x) - P(x) + P(x) - P_n(x) \\
&= f(x) - P(x) + \sum_{i=0}^n (P(x_i) - P_n(x_i))l_i(x) \\
&= f(x) - P(x) + \sum_{i=0}^n (P(x_i) - f(x_i))l_i(x)
\end{aligned}$$

Damit

$$\begin{aligned}
|f(x) - P_n(x)| &\leq |f(x) - P(x)| + \sum_{i=0}^n |P(x_i) - f(x_i)| |l_i(x)| \\
&\leq |1 + \sum_{i=0}^n |l_i(x)|| - \max_{a \leq x \leq b} |f(x) - P(x)|
\end{aligned}$$

Bilde nun das Infimum über alle  $P$ . Dann folgt die Behauptung.

□

$$\|f - P_n\| \leq (1 - \lambda_n)E_n(f)$$

Frage1: Wie verhält sich  $E_n(f)$ ?

Frage2: Wie verhält sich  $\lambda_n$ ?

Für welche Stützstellenwahl wird  $\lambda_n$  minimal?

### Satz 6.18 (Weierstraß)

Für alle stetigen Funktionen  $f$  gilt:

$$\lim_{n \rightarrow \infty} E_n(f) = 0$$

### Satz 6.19 (Jackson)

Ist  $f$   $r$ -mal stetig differenzierbar, so gilt für  $n \geq r + 1$ :

$$E_n(f) \leq \frac{\pi}{2} \left(\frac{b-1}{2}\right)^r \|f^{(r)}\|_{\infty} \left(\frac{1}{n+1}\right)^r.$$

### Satz 6.20

Es gibt eine von  $n$  unabhängige Konstante  $c > 0$  mit  $\lambda_n \geq \frac{2}{\pi} \ln(n) - c$  für jede beliebige Stützstellenwahl.

### Satz 6.21

Bei konstantem Stützstellenabstand wächst  $\lambda_n$  exponentiell.

### Satz 6.22

Wählt man als Stützstellen die auf das Intervall Transformierten Tschebyscheff-Knoten.

$$x_j = \frac{b-a}{2} \cos\left(\frac{2j+1}{2} \frac{\pi}{n+1}\right) + \frac{a+b}{2}$$

$j = 0, \dots, n$ , so gilt  $\lambda_n \leq \frac{2}{\pi} \ln(n) + 1$

Bemerkung: Interpoliert man eine einmal stetig differenzierbare Funktion  $f : [a, b] \rightarrow \mathbb{R}$  an die Tschebyscheff-Knoten von Satz 6.22, so Konvergieren die Interpolationspolynome  $P_n$  für  $n \rightarrow \infty$  in der Supremumsnorm gegen  $f$ .

Es gilt nach Satz 6.17  $\max_{a \leq x \leq b} |P_n(x) - f(x)| \leq (1 - \lambda_n) E_n(f)$  und nach

Satz 6.19  $E_n(f) \leq C \cdot \frac{1}{n+1}$ . Zusammen mit Satz 6.20 also

$$\|P_n - f\|_\infty \leq \left(1 + \frac{2}{\pi} \ln(n) + 1\right) - C \frac{1}{n+1} \rightarrow 0$$

für  $n \rightarrow \infty$ .

## Berechnung des Interpolationspolynoms

$$g(t) = f\left(\frac{b-a}{2} \cos(t) + \frac{a+b}{2}\right), t \in \mathbb{R}$$

ist eine  $2a$ -periodische Funktion. Berechne die Koeffizienten  $a_r$  des trigonometrischen Interpolationspolynoms

$$(S_n g)(t) = \sum_{k=0}^n a_k \cos(kt)$$

von  $g$  in den Punkten

$$t_j = \frac{2j+1}{2} \frac{\pi}{n}, j \in \mathbb{Z}$$

mit Hilfe der schnellen Fouriertransformation(FFT). Dazu benötigt man asymptotisch  $n \log(n)$  Operationen. Dann ist

$$(*) P(x) = \sum_{k=0}^n a_k T_k\left(\frac{2}{b+a}x - \frac{b+a}{b-a}\right)$$

das gesuchte Interpolationspolynom, wobei

$$T_k = \cos(k \arccos(x)), -1 \leq x \leq 1$$

das Tschebyscheffpolynom  $k$ -ten Grades ist. Die Tschebyscheff Polynome

genügen der Rekursion  $T_0(x) = 1, T_1(x) = x, T_k(x) = 2x T_{k-1}(x) - T_{k-2}(x)$

die sich mit Hilfe der trig. Identität  $\cos\alpha + \cos\beta = 2\cos\left(\frac{\alpha+\beta}{2}\right)\cos\left(\frac{\alpha-\beta}{2}\right)$

leicht beweisen lässt.

### Auswertung von (\*)

Setze  $x' = \frac{2}{b-a}x - \frac{b+a}{b-a}, T_0 = 1, T_1 = x', s_1 = a_1 x' + a_0$

Berechne für  $k = 2, 3, \dots, n$

$$T_k = 2x' T_{k-1} - T_{k-2}$$

$$S_k = S_{k-1} + a_k T_k.$$

Dann ist  $P(x) = S_n$

Wegen  $|x'| < 1$  ist diese Rekursion im Unterschied zum Horner-Schema stabil und benötigt ebenfalls nur  $\mathcal{O}(n)$  Operationen.

## 6.3 Interpolation durch Splines

### Definition 6.24

Im Intervall  $[a, b]$  seien  $N + 1$  Punkte  $a = x_0 < x_1 < \dots < x_N = b$  gegeben. Eine auf  $[a, b]$  definierte Funktion  $S$  heißt genau dann  $m$ -mal stetig differenzierbare Splinefunktion  $k$ -ten Grades zu der gegebenen Unterteilung von  $[a, b]$ , wenn die Einschränkungen von  $S$  auf jedes Intervall  $[x_i, x_{i+1}]$ ,  $i = 0, \dots, N - 1$ . Ein Polynom  $k$ -ten Grades ist und wenn  $S$   $m$ -mal stetig differenzierbar ist.

Bemerkung: Sinnvoll ist nur der Fall  $m < k$ . Für  $m \geq k$  wäre  $S$  wieder ein Polynom und durch die stückweise Definition gewonnene Flexibilität ginge verloren.

hier: Nur 2 mal stetig differenzierbar kubische Splines. Für viele Anwendungszwecke, bspw. in der Grafischen Datenverarbeitung, optimal.

### Satz 6.25

Erfüllt die zweimal stetig differenzierbare Splinefunktion  $S$  die Interpolationsbedingungen  $S(x_j) = f_j$ ,  $j = 0, 1, \dots, N$ , so hat sie auf dem Intervall  $x_j \leq x \leq x_{j+1}$  die Darstellung.

$$S(x) = f_j \frac{x_{j+1}-x}{x_{j+1}-x_j} + f_{j+1} \frac{x-x_j}{x_{j+1}-x_j} - \frac{1}{6} S''(x_j) \frac{(x_{j+1}-x)(x-x_j)}{x_{j+1}-x_j} ((x_{j+1}-x) + (x_{j+1}-x_j)) - \frac{1}{6} S''(x_{j+1}) \frac{(x_{j+1}-x)(x-x_j)}{x_{j+1}-x_j}$$

Weiter ist

$$S'(x_j) = \frac{f_{j+1}-f_j}{x_{j+1}-x_j} - \left( \frac{1}{3} S''(x_j) + \frac{1}{6} S''(x_{j+1}) \right) (x_{j+1}-x_j)$$

$$S'(x_{j+1}) = \frac{f_{j+1}-f_j}{x_{j+1}-x_j} + \left( \frac{1}{6} S''(x_j) + \frac{1}{3} S''(x_{j+1}) \right) (x_{j+1}-x_j)$$

Beweis: Da  $S''(x)$  auf  $[x_j, x_{j+1}]$  linear ist, gilt dort

$$S''(x) = S''(x_j) \frac{x_{j+1}-x}{x_{j+1}-x_j} + S''(x_{j+1}) \frac{x-x_j}{x_{j+1}-x_j}$$

Damit ist auf diesem Intervall

$$S(x) = \frac{1}{6} \frac{(x_{j+1}-x)^3}{x_{j+1}-x_j} S''(x_j) + \frac{1}{6} \frac{(x-x_j)^3}{x_{j+1}-x_j} S''(x_{j+1}) + c_j(x-x_j) + d_j(x_{j+1}-x)$$

mit gewissen Integrationskonstanten  $c_j, d_j$ .

Die Interpolationsforderung

$$S(x_j) = \frac{1}{6}S''(x_j)(x_{j+1} - x_j)^2 + d_j(x_{j+1} - x_j) = f_j$$

$$S(x_{j+1}) = \frac{1}{6}S''(x_{j+1})(x_{j+1} - x_j)^2 + c_j(x_{j+1} - x_j) = f_{j+1}$$

führen auf

$$c_j = \frac{f_{j+1}}{x_{j+1} - x_j} - \frac{1}{6}S''(x_{j+1})(x_{j+1} - x_j)$$

$$d_j = \frac{f_j}{x_{j+1} - x_j} - \frac{1}{6}S''(x_j)(x_{j+1} - x_j)$$

Damit ist für  $x \in [x_j, x_{j+1}]$

$$S'(x) = \frac{1}{2}S''(x_{j+1})\frac{(x-x_j)^2}{x_{j+1}-x_j} - \frac{1}{2}S''(x_j)\frac{(x_{j+1}-x)^2}{x_{j+1}-x_j} + \frac{f_{j+1}-f_j}{x_{j+1}-x_j} \\ + \left(\frac{1}{6}S''(x_j) - \frac{1}{6}S''(x_{j+1})\right)(x_{j+1} - x_j)$$

Wertet man  $S'$  an den Stellen  $x_j$  und  $x_{j+1}$  aus, so folgt die Behauptung.

Für die Funktionen selbst erhält man die Darstellung

$$S(x) = \frac{1}{6}S''(x_j)\left(\frac{(x_{j+1}-x)^2}{x_{j+1}-x_j} - (x_{j+1} - x_j)(x_{j+1} - x)\right) \\ + \frac{1}{6}S''(x_{j+1})\left(\frac{(x-x_j)^2}{x_{j+1}-x_j} - (x_{j+1} - x_j)(x - x_j)\right) \\ + f_{j+1}\frac{x-x_j}{x_{j+1}-x_j} + f_j\frac{x_{j+1}-x}{x_{j+1}-x_j}$$

für  $x \in [x_j, x_{j+1}]$ , die man um auslöschungseffekte zu vermeiden besser in der oben angegebenen Form auswertet.

□

### Satz 6.26

Erfüllt die zweimal stetig differenzierbare kubische Splinefunktion

$f$  die Interpolationsbedingungen  $S(x_j) = f_j, j = 0, \dots, N$ , so gilt für

$$j = 1, \dots, N-1 \quad \frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}}S''(x_{j+1}) + 2S''(x_j) + \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}}S''(x_{j-1})$$

$$= \frac{6}{x_{j+1}-x_{j-1}} \left( \frac{f_{j+1}-f_j}{x_{j+1}-x_j} - \frac{f_j-f_{j-1}}{x_j-x_{j-1}} \right) = 6f[x_{j-1}, x_j, x_{j+1}] \quad (2)$$

Beweis: Aus Satz 6.25 angewandt auf das Intervall  $[x_{j-1}, x_j]$ , folgt

$$S'(x_j) = \frac{f_j-f_{j-1}}{x_j-x_{j-1}} + \left(\frac{1}{6}S''(x_{j-1}) + \frac{1}{3}S''(x_j)\right)(x_j - x_{j-1})$$

und andererseits ( $[x_j, x_{j+1}]$ )

$$S'(x_j) = \frac{f_{j+1}-f_j}{x_{j+1}-x_j} - \left(\frac{1}{3}S''(x_j) + \frac{1}{6}S''(x_{j+1})\right)(x_{j+1} - x_j)$$

Daraus folgt

$$\left(\frac{1}{6}S''(x_{j-1}) + \frac{1}{3}S''(x_j)\right)(x_j - x_{j-1}) + \left(\frac{1}{3}S''(x_j) + \frac{1}{6}S''(x_{j+1})\right)(x_{j+1} - x_j) \\ = \frac{f_{j+1}-f_j}{x_{j+1}-x_j} - \frac{f_j-f_{j-1}}{x_j-x_{j-1}}$$

oder nach Multiplikation mit  $\frac{6}{x_{j+1}-x_j}$  die behauptung.

□

Um die Momente  $S''(x_j), j = 0, \dots, N$  und damit den Spline  $S$  bestimmen zu können, fehlen noch zwei Gleichungen. Diese Gleichungen lassen sich nur aus Zusatzbedingungen herleiten.

### Definition 6.27

Eine zweimal stetig differenzierbare kubische Splinefunktion  $S$ , deren zweite Ableitungen in den Randpunkten  $a = x_0$  und  $b = x_N$

(3)  $S''(x_0) = 0, S''(x_N) = 0$  erfüllt, heißt natürliche Splinefunktion.

### Satz 6.28

Die Gleichungen (2) aus Satz 6.26 und (3) bestimmen die Momente  $S''(x_0), \dots, S''(x_n)$  eindeutig. Dieses Momentengleichungssystem ist sehr gut konditioniert. Bzgl. der Maximumnorm als Vektornorm ist die Norm seiner Koeffizientenmatrix unabhängig von der Intervalleinteilung durch (3) und die Norm seiner Inversen durch 1 beschränkt.

Beweis: Sei  $M_0 = b_0, M_N = b_N$  und für  $j = 1, \dots, N - 1$

$$\frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}}M_{j+1} + 2N_j + \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}}M_{j-1} = b_j$$

$$\text{Sei weiter } M = \max_{j=0, \dots, N} |M_j|, B = \max_{j=0, \dots, N} |b_j|$$

Dann ist für  $j = 1, \dots, N - 1$

$$2|M_j| = \left| b_j - \frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}}M_{j+1} - \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}}M_{j-1} \right|$$

$$\leq B + \frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}}M + \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}}M = B + M$$

also

$$|M_j| \leq \frac{1}{2}B + \frac{1}{2}M$$

Da außerdem

$$|M_0| = \frac{1}{2}|b_0| + \frac{1}{2}|M_0| \leq \frac{1}{2}B + \frac{1}{2}M$$

$$|M_N| = \frac{1}{2}|b_N| + \frac{1}{2}|M_N| \leq \frac{1}{2}B + \frac{1}{2}M$$

gilt  $M \leq \frac{1}{2}B + \frac{1}{2}M$  oder  $M \leq B$ .

Damit ist das Gleichungssystem eindeutig lösbar und die Norm der inversen der Koeffizientenmatrix durch 1 beschränkt. Wegen

$$\frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}} + 2 + \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}} = 2$$

ist die Norm der Matrix gleich 3.

□

### Satz 6.29

Ist  $f : [a, b] \rightarrow \mathbb{R}$  eine zweimal stetig differenzierbare Funktion und  $S$  ein zweimal stetig differenzierbarer kubischer Spline mit  $S(x_j) = f(x_j)$  für  $j = 0, \dots, N$ , so gilt

$$\int_a^b |f''(x)|^2 dx = \int_a^b |S''(x)|^2 dx + \int_a^b |S''(x) - f''(x)|^2 dx - 2(S'(x) - f'(x))S''(x)|_{x_0}^{x_N}$$

Beweis: Zunächst gilt:

$$\int_a^b |S''(x)|^2 dx + \int_a^b |S''(x) - f''(x)|^2 dx = 2 \int_a^b (S''(x) - f''(x))S''(x) dx + \int_a^b |f''(x)|^2 dx$$

Weiter gilt:

$$\begin{aligned} \int_a^b (S''(x) - f''(x))S''(x) dx &= \sum_{j=0}^{N-1} \int_{x_j}^{x_{j+1}} (S''(x) - f''(x))S''(x) dx \\ &= \sum_{j=0}^{N-1} [(S'(x) - f'(x))S''(x)|_{x_j}^{x_{j+1}} - \int_{x_j}^{x_{j+1}} (S'(x) - f'(x))S'''(x) dx] \\ &= (S'(x) - f'(x))S''(x)|_{x_0}^{x_N} - \sum_{j=0}^{N-1} \int_{x_j}^{x_{j+1}} (S'(x) - f'(x))S'''(x) dx \end{aligned}$$

Da  $S'''(x)$  auf den Teilintervallen  $x_j \leq x \leq x_{j+1}$  konstant ist, gilt

$$\begin{aligned} \sum_{j=0}^{N-1} \int_{x_j}^{x_{j+1}} (S'(x) - f'(x))S'''(x) dx &= \sum_{j=0}^{N-1} S'''(\frac{x_j+x_{j+1}}{2}) \int_{x_j}^{x_{j+1}} (S'(x) - f'(x)) dx \\ &= \sum_{j=0}^{N-1} S'''(\frac{x_j+x_{j+1}}{2}) [S(x) - f(x)]_{x_j}^{x_{j+1}} \\ &= 0 \end{aligned}$$

Einsetzen ergibt die Behauptung.

□

Folgerung:

Für die Interpolierende natürliche kubische Splinefunktion  $S$  gilt:

$$\int_a^b |f''(x)|^2 dx = \int_a^b |S''(x)|^2 dx + \int_a^b |f''(x) - S''(x)|^2 dx \quad (*)$$

und damit für  $f \neq S$

$$\int_a^b |f''(x)|^2 dx > \int_a^b |S''(x)|^2 dx$$

Damit ist  $S$  die "glatteste" Funktion mit den Werten  $S(x_j) = f_j$ .

Bemerkung:

Der Satz 6.29 legt es nahe, statt der Zusatzbedingung  $S''(x_0) = 0$ ,  $S''(x_N) = 0$ . Die Bedingungen  $S'(x_0) = f'(x_0), S'(x_N) = f'(x_N)$  zu betrachten. Dabei bleiben die beschriebenen Optimalitätseigenschaften erhalten.

### Satz 6.30

Die Zusatzbedingungen  $S'(x_0) = f'(x_0), S'(x_N) = f'(x_N)$  sind äquivalent zu  $2S''(x_0) + S''(x_1) = \frac{6}{x_1 - x_0} \left( \frac{f_1 - f_0}{x_1 - x_0} - f'_0 \right)$ ,

$$S''(x_{N-1}) + 2S''(x_N) = \frac{6}{x_N - x_{N-1}} \left( f'_N - \frac{f_N - f_{N-1}}{x_N - x_{N-1}} \right),$$

Das Momentengleichungssystem ist auch mit diesen Zusatzbedingungen eindeutig lösbar und die Norm seiner Koeffizientenmatrix durch 3 mal die Norm von deren Inversen durch 1 beschränkt.

### Satz 6.31

Für jeden zweimal stetig differenzierbaren kubischen Spline  $S$ , der die zweimal stetig differenzierbare Funktion  $f$  in den Punkten  $x_0 < \dots < x_N$  interpoliert, gilt auf dem Intervall  $[x_j, x_{j+1}]$

$$|S(x) - f(x)| \leq (x_{j+1} - x_j) \int_{x_j}^{x_{j+1}} |S''(t) - f''(t)| dt$$

$$|S'(x) - f'(x)| \leq \int_{x_j}^{x_{j+1}} |S''(t) - f''(t)| dt$$

Falls  $|S''(t) - f''(t)| \leq k$  ist, folgt:

$$|S(x) - f(x)| \leq k(x_{j+1} - x_j)^2,$$

$$|S'(x) - f'(x)| \leq k(x_{j+1} - x_j)$$

Beweis: Wegen  $S(x_j) - f(x_j) = S(x_{j+1}) - f(x_{j+1}) = 0$  gibt es nach dem Satz von Rolle ein  $\eta \in (x_j, x_{j+1})$  mit  $S'(\eta) - f'(\eta) = 0$ . Für  $x_j \leq x \leq x_{j+1}$

$$\text{gilt daher } |S'(x) - f'(x)| = \left| \int_{\eta}^x (S''(t) - f''(t)) dt \right| \leq \int_{x_j}^{x_{j+1}} |S''(t) - f''(t)| dt$$

Damit ist wegen  $S(x_j) - f(x_j) = 0$

$$|S(x) - f(x)| = \left| \int_{x_j}^x (S'(t) - f'(t)) dt \right| \leq (x_{j+1} - x_j) \int_{x_j}^{x_{j+1}} |S''(t) - f''(t)| dx$$

□



**Satz 6.32**

Ist  $h$  die maximale Teilintervalllänge, so gilt für den zweimal stetig differenzierbaren kubischen Spline  $S$ , der die viermal stetig differenzierbare Funktion  $f$  in den Stützstellen  $x_j$  interpoliert und den Zusatzbedingungen  $S'(x_0) = f'(x_0), S'(x_N) = f'(x_N)$  genügt, die Abschätzung  $\|S'' - f''\|_\infty \leq \frac{3}{8}h^2\|f^{(4)}\|_\infty$

Beweisidee: Man betrachtet das lineare Gleichungssystem

$$2D_0 + D_1 = B_0$$

$$\frac{x_{j+1}-x_j}{x_{j+1}-x_{j-1}}D_{j+1} + 2D_j + \frac{x_j-x_{j-1}}{x_{j+1}-x_{j-1}}D_{j-1} = B_j, j = 1, \dots, N-1$$

$$D_{N-1} + 2D_N = B_N$$

für die Differenzen  $D_j := S''(x_j) - f''(x_j)$  mit entsprechenden rechten Seiten  $B_j$ .

Eine sorgfältige Analyse mit Hilfe des Taylorschen Satzes mit Integralrestglied zeigt:

$$\max_{j=0, \dots, N} |B_j| \leq \frac{1}{4}h^2\|f^{(4)}\|_\infty$$

Nach Satz 6.28 ist die Norm der Inversen der Koeffizientenmatrix des Gleichungssystems durch 1 beschränkt. Daher ist

$$\max_{j=0, \dots, N} |D_j| \leq \max_{j=0, \dots, N} |B_j|$$

Zusammengefasst ergibt sich

$$\max_{j=0, \dots, N} |S''(x_j) - f''(x_j)| \leq \frac{1}{4}h^2\|f^{(4)}\|_\infty \quad (*)$$

Wir betrachten nun die stückweise lineare Interpolation

$$(If'')(x) = \frac{x-x_j}{x_{j+1}-x_j}f''(x_{j+1}) + \frac{x_{j+1}-x}{x_{j+1}-x_j}f''(x_j), x_j \leq x \leq x_{j+1}$$

der zweiten Ableitung von  $f$ . Aus (\*) folgt dann  $\|S'' - If''\|_\infty = \max_{j=0, \dots, N} |S''(x_j) - f''(x_j)| \leq \frac{1}{4}h^2\|f^{(4)}\|_\infty$

Andererseits ist nach Satz 6.6

$$\|f'' - If''\|_\infty \leq \frac{1}{8}h^2\|f^{(4)}\|_\infty$$

und damit schließlich

$$\|S'' - f''\|_\infty \leq \left(\frac{1}{4}h^2 + \frac{1}{8}h^2\right)\|f^{(4)}\|_\infty = \frac{3}{8}h^2\|f^{(4)}\|_\infty$$

Folgerung:

Für viermal stetig differenzierbares  $f$  gilt

$$\|S - f\|_\infty = \mathcal{O}(h^4)$$

$$\|S' - f'\|_\infty = \mathcal{O}(h^3)$$

$$\|S'' - f''\|_\infty = \mathcal{O}(h^2)$$

**Satz 6.33**

Für den natürlichen interpolierenden Spline (mit den Zusatzbedingungen  $S''(x_0) = S''(x_N) = 0$ ) gilt für  $j = 0, \dots, N-1$

$$|S''(x_j) - f''(x_j)| \leq \frac{1}{4}h^2 \|f^{(4)}\|_\infty + \left(\frac{1}{2}\right)^j |f''(x_0)| + \left(\frac{1}{2}\right)^{N-j} |f''(x_N)|$$

Beweisidee: dann definiere die Splinefunktion  $S_L$  und  $S_R$  durch die Bedingungen  $S_L^{(2)}(x_0) = 1, S_L^{(2)}(x_N) = 0, S_L(x_j) = 0, j = 0, \dots, N$

$S_R^{(2)}(x_0) = 0, S_R^{(2)}(x_N) = 1, S_R(x_j) = 0, j = 0, \dots, N$  und setze weiter

$$\tilde{S}(x) = S(x) + f''(x_0)S_L(x) + f''(x_N)S_R(x)$$

Für den interpolierenden Spline  $\tilde{S}$  gilt wegen  $S''(x_0) = f''(x_0),$

$S''(x_N) = f''(x_N)$  analog zum Beweis von Satz 6.32

$$\max_{j=0, \dots, N} |\tilde{S}''(x_j) - f''(x_j)| \leq \frac{1}{4}h^2 \|f^{(4)}\|_\infty$$

Daraus folgt für  $j = 0, \dots, N$

$$|S''(x_j) - f''(x_j)| \leq \frac{1}{4}h^2 \|f^{(4)}\|_\infty + |f''(x_0)| |S_L^{(2)}(x_j)| + |f''(x_N)| |S_R^{(2)}(x_j)|$$

Wie man durch Analyse des Momentengleichungssystems zeigen kann, ist

$$|S_L^{(2)}(x_j)| \leq \left(\frac{1}{2}\right)^j, |S_R^{(2)}(x_j)| \leq \left(\frac{1}{2}\right)^{N-j}$$

□

## 7. Numerische Integration

Problem: näherungsweise Berechnung bestimmter Integrale.

### Beispiel 7.1

$$\frac{8}{\pi \ln 2} \int_0^1 \frac{\ln(1+x)}{1+x^2} dx = 1$$

### 7.1 Zusammengesetzte Quadraturformeln

Zu berechnen sei  $\int_a^b f(x) dx$ . Gegeben sei eine Zerlegung

$a = x_0 < x_1 < \dots < x_n = b$ . Dann gilt mit  $g_i(t) = f\left(\frac{x_i+x_{i+1}}{2} + t\frac{x_{i+1}-x_i}{2}\right)$

$$\begin{aligned} \int_a^b f(x) dx &= \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} f(x) dx = \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \int_{-1}^1 f\left(\frac{x_i+x_{i+1}}{2} + t\frac{x_{i+1}-x_i}{2}\right) dt \\ &= \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \int_{-1}^1 g_i(t) dt \end{aligned}$$

Idee: Ersetze die Integrale  $\int_{-1}^1 g_i(t) dt$  durch finite Ausdrücke

$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^n a_k g(t_k)$$

wobei die  $a_k$  von  $g$  unabhängige Gewichte und die  $t_k \in [-1, 1]$  von  $g$  unabhängige Knoten sind.

Endergebnis:

$$\int_a^b f(x) dx \rightarrow \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \sum_{k=0}^n a_k f\left(\frac{x_i+x_{i+1}}{2} + t\frac{x_{i+1}-x_i}{2}\right)$$

### Beispiel Rechteckregel

$$\int_{-1}^1 g(t) dt \rightarrow 2g(-1)$$

$$\int_a^b f(x) dx \rightarrow \sum_{i=0}^{n-1} (x_{i+1} - x_i) f(x_i)$$

### Mittelpunktsregel

$$\int_{-1}^1 g(t) dt \rightarrow 2g(0)$$

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

### Trapezregel

$$\int_{-1}^1 g(t) dt \rightarrow g(-1) + g(1)$$

$$\int_a^b f(x) dx \rightarrow \sum_{i=0}^{n-1} (x_{i+1} - x_i) \frac{f(x_i) + f(x_{i+1})}{2}$$

### Simpsonregel

$$\int_{-1}^1 g(t) dt \rightarrow \frac{1}{3}g(-1) + \frac{4}{3}g(0) + \frac{1}{3}g(1)$$

$$\int_a^b f(x) dx \rightarrow \sum_{i=0}^{n-1} \frac{(x_{i+1} - x_i)}{6} \left( f(x_i) + 4f\left(\frac{x_i + x_{i+1}}{2}\right) + f(x_{i+1}) \right)$$

### Beispiel 7.2

Integral aus Beispiel 7.1

$$x_i = \frac{i}{n}, i = 0, \dots, n$$

Konvergenz für  $n \rightarrow \infty$ ?

Beobachtung: Fehler Simpsonregel  $\mathcal{O}\left(\frac{1}{n^4}\right)$ ?

$\ll$  Fehler Mittelpunktsregel, Trapezregel  $\mathcal{O}\left(\frac{1}{n^2}\right)$ ?

$\ll$  Fehler Rechtecksregel  $\mathcal{O}\left(\frac{1}{n}\right)$ ?

### Satz 7.3

Die auf dem Referenzintervall vorgegebene Formel  $\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^n a_k g(t_k)$

sei exakt für alle Polynome vom Grad  $q - 1 \geq 0$ . Die Funktion  $f : [a, b] \rightarrow \mathbb{R}$

sei  $q$ -mal stetig differenzierbar und  $a = x_0 < \dots < x_n = b$  sei eine beliebige

Zerlegung des Intervalls  $[a, b]$ . Dann ist  $\left| \int_a^b f(x) dx - \sum_{i=0}^{n-1} \frac{x_{i+1} - x_i}{2} \sum_{k=0}^m a_k f\left(\frac{x_i + x_{i+1}}{2} + t_k \frac{x_{i+1} - x_i}{2}\right) \right|$

$$\leq C \sum_{k=0}^{n-1} (x_{i+1} - x_i)^q \int_{x_i}^{x_{i+1}} |f^{(q)}(x)| dx$$

mit einer von der Wahl der Formel für das Referenzintervall abhängigen

Konstante  $C$ .

Hilfssatz: Es gilt

$$\int_{-1}^1 g(t) dt - \sum_{k=0}^m a_k g(t_k) = \int_{-1}^1 K(s) g^{(q)}(s) ds$$

$$K(s) = \frac{1}{q!} (1-s)^q - \frac{1}{(q-1)!} \sum_{k=0}^m a_k (t_k - s)_t^{q-1}$$

mit der abgeschnittenen Potenz

$$(t-s)_t^{q-1} = \begin{cases} (t-s)^{q-1}, & \text{falls } s \leq t \\ 0, & \text{sonst} \end{cases}$$

Beweis: Nach dem Taylorschen Satz ist  $g(t) = \sum_{j=0}^{q-1} \frac{1}{j!} g^{(j)}(-1)(t+1)^j + \frac{1}{(q-1)!} \int_{-1}^t (t-s)^{q-1} g^{(q)}(s) ds$

(Beweis durch Induktion über  $q$  mit partieller Integration)

$$\text{also } g(t) = \sum_{j=0}^{q-1} \frac{1}{j!} g^{(j)}(-1)(t+1)^j + \frac{1}{(q-1)!} \int_{-1}^{+1} (t-s)^{q-1} g^{(q)}(s) ds$$

Da der polynomiale Anteil nach Voraussetzung exakt integriert wird folgt:

$$\begin{aligned} & \int_{-1}^1 g(t) dt - \sum_{k=0}^m a_k g(t_k) \\ &= \int_{-1}^1 \left( \frac{1}{(q-1)!} \int_{-1}^1 (t-s)^{q-1} g^{(q)}(s) ds \right) dt - \sum_{k=0}^m a_k \left( \frac{1}{(q-1)!} \int_{-1}^1 (t_k - s)_t^{q-1} g^{(q)}(s) ds \right) \\ &= \int_{-1}^1 \left( \frac{1}{(q-1)!} \int_{-1}^1 (t-s)^{q-1} dt \right) g^{(q)}(s) ds - \int_{-1}^1 \left( \sum_{k=0}^m a_k \frac{1}{(q-1)!} (t_k - s)_t^{q-1} \right) g^{(q)}(s) ds \\ &= \int_{-1}^1 K(s) g^{(q)}(s) ds \end{aligned}$$

mit dem Kern

$$\begin{aligned} K(s) &= \frac{1}{(q-1)!} \left( \int_{-1}^1 (t-s)_t^{q-1} dt - \sum_{k=0}^m a_k (t_k - s)_t^{q-1} \right) \\ &= \frac{1}{(q-1)!} \left( \frac{1}{q} (1-s)^q - \sum_{k=0}^m a_k (t_k - s)_t^{q-1} \right) \end{aligned}$$

□

Beweis von Satz 7.3:

Setzt man  $g_i(t) = f\left(\frac{x_i+x_{i+1}}{2} + t \frac{x_{i+1}-x_i}{2}\right)$ , gilt nach dem Hilfssatz

$$\begin{aligned} & \int_a^b f(x) dx - \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \sum_{k=0}^m a_k f\left(\frac{x_i+x_{i+1}}{2} + t_k \frac{x_{i+1}-x_i}{2}\right) \\ &= \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \left( \int_{-1}^1 g_i(t) dt - \sum_{k=0}^m a_k g_i(t_k) \right) \\ &= \sum_{i=0}^{n-1} \frac{x_{i+1}-x_i}{2} \int_{-1}^1 K(s) g_i^{(q)}(s) ds \end{aligned}$$

$$\begin{aligned}
\text{mit } C &:= \frac{1}{2^{q-1}} \max_{-1 \leq s \leq 1} |K(s)| \text{ ist } \left| \int_{-1}^1 K(s) g_i^{(q)}(s) ds \right| \leq 2^{q-1} C \int_{-1}^1 |g_i^{(q)}(s)| ds \\
&= 2^{q-1} C \int_{-1}^1 \left| f^{(q)} \left( \frac{x_i + x_{i+1}}{2} + s \frac{x_{i+1} - x_i}{2} \right) \right| \left( \frac{x_{i+1} - x_i}{2} \right)^q ds \\
&= 2^{q-1} C \left( \frac{x_{i+1} - x_i}{2} \right)^{q-1} \int_{x_i}^{x_{i+1}} |f^{(q)}(x)| dx
\end{aligned}$$

Summation über  $n$  ergibt die Behauptung.

□

äquidistante Zerlegungen:  $x_{i+1} - x_i = h$

$$\begin{aligned}
&\left| \int_a^b f(x) dx - \frac{h}{2} \sum_{i=0}^{n-1} \sum_{k=0}^m a_k f \left( \frac{x_i + x_{i+1}}{2} + \frac{h}{2} t_k \right) \right| \\
&\leq C \sum_{i=0}^{n-1} h^q \int_{x_i}^{x_{i+1}} |f^{(q)}(x)| dx \\
&= Ch^q \int_a^b |f^{(q)}(x)| dx
\end{aligned}$$

exakt für Polynome vom Grad  $q - 1$

$\Rightarrow$  Konvergenzordnung  $\mathcal{O}(h^q)$

Bemerkung:

Äquidistante Zerlegungen sind nicht immer sinnvoll!

$$\text{Ziel: } (x_{i+1} - x_i)^q \int_{x_i}^{x_{i+1}} |f^{(q)}(x)| dx \approx \frac{x_{i+1} - x_i}{b-a} \varepsilon$$

$$\Rightarrow |\text{Fehler}| \leq \varepsilon$$

## Konvergenz bei Riemann integrierbaren Funktionen

$$\begin{aligned}
&\int_a^b f(x) dx - \sum_{i=0}^{n-1} \frac{x_{i+1} - x_i}{2} \sum_{k=0}^m a_k f \left( \frac{x_i + x_{i+1}}{2} + t_k \frac{x_{i+1} - x_i}{2} \right) \\
&= \frac{1}{2} \sum_{k=0}^m a_k \underbrace{\left( \int_a^b f(x) dx - \sum_{i=0}^{n-1} (x_{i+1} - x_i) f \left( \frac{x_i + x_{i+1}}{2} + t_k \frac{x_{i+1} - x_i}{2} \right) \right)}_{\text{Riemannsche Zwischensumme}}
\end{aligned}$$

$\Rightarrow$  Konvergenz

## 7.2 Interpolatorische Quadraturformeln

Ziel: Konstruiere Formeln

$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g(t_k) \quad (*)$$

möglichst hoher polynomialer Exaktheit.

### Satz 7.4

Gegeben seien die  $m + 1$  von einander verschiedenen Knoten  $t_0, \dots, t_n \in [-1, 1]$ . Dann ist die Formel(\*) genau dann für alle Polynome

$$n\text{-ten Grades exakt, falls } a_k = \int_{-1}^1 \prod_{\substack{j=0 \\ j \neq k}}^m \frac{t-t_j}{t_k-t_j} dt$$

Beweis: Polynome  $m$ -ten Grades haben die Darstellung

$$P(t) = \sum_{k=0}^m l_k(t) P(t_k), l_k(t) = \prod_{\substack{j=0 \\ j \neq k}}^m \frac{t-t_j}{t_k-t_j}$$

□

### Beispiel 7.5

$$\int_{-1}^1 g(t) dt \rightarrow a_0 g(-1) + a_1 g(1)$$

$$a_0 = \int_{-1}^1 \frac{t-1}{-1-1} dt = 1, a_1 = \int_{-1}^1 \frac{t-(-1)}{1-(-1)} dt = 1$$

Bemerkung: Die Gewichte  $a_k$  lassen sich über das lineare Gleichungssystem

$$\sum_{k=0}^m t_k^j a_k = \int_{-1}^1 t^j dt, j = 0, \dots, m \text{ bestimmen.}$$

### Beispiel 7.6

Ist  $t_0 = -1, t_1 = 0, t_2 = 1$  ergibt sich

$$\begin{array}{rcl} a_0 & +a_1 & +a_2 = 2 \\ -a_0 & & +a_2 = 0 \\ a_0 & & +a_2 = \frac{2}{3} \end{array} \left| \begin{array}{l} g(t) = 1 \\ g(t) = t \\ g(t) = t^2 \end{array} \right.$$

und damit  $a_0 = \frac{1}{3}, a_1 = \frac{4}{3}, a_2 = \frac{1}{3}$ , also die Simpsonregel

$$\int_{-1}^1 g(t) dt \rightarrow \frac{1}{3}g(-1) + \frac{4}{3}g(0) + \frac{1}{3}g(1),$$

die auch für  $g(t) = t^3$  und damit für alle Polynome 3. Grades exakt ist.

### Beispiel 7.7

$$\frac{8}{\pi \ln 2} \int_0^1 \frac{\ln(1+x)}{1+x^2} dx (= 1)$$

$$x_i = \frac{i}{n}, i = 0, \dots, n$$

Beobachtung: Die Simpsonregel ist ungleich genauer als Mittelpunkt- oder Trapezregel.

### Newton-Cotes-Formeln

$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g\left(-1 + k \frac{2}{m}\right) \quad (*)$$

exakt für Polynome vom Grad  $m$ .

### Satz 7.8

Ist  $m$  gerade, so ist die Newton-Cotes-Formel (\*) sogar für alle Polynome vom Grad  $m + 1$  exakt.

Beweis: Aus Symmetriegründen ist mit  $m = 2l$

$$\sum_{k=0}^m a_k g\left(-1 + k \frac{2}{m}\right) = a_k g(0) + \sum_{k=1}^l a_{l+k} \left(g\left(-\frac{k}{l}\right) + g\left(\frac{k}{l}\right)\right)$$

Die Formel ist daher für alle ungeraden Funktionen  $g$  exakt.

□

Bemerkung: Für  $m \geq 7$  treten negative Gewichte  $a_k$  auf, was die Brauchbarkeit dieser Formeln dann stark einschränkt.

### Satz 7.9

Die Quadraturformel der Form

$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g(t_k)$$

kann maximal für Polynome  $(2m + 1)$ -ten Grades exakt sein.

Beweis: Das Polynom  $P(t) = \prod_{j=0}^m (t - t_j)^2$  vom Grad  $2m + 2$  kann durch diese

Formel nicht mehr exakt integriert werden, denn es gilt:

$$\int_{-1}^1 P(t) dt > 0 = \sum_{k=0}^m a_k P(t_k)$$

□

### Satz 7.10

Ist die Quadraturformel



$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g(t_k)$$

für alle Polynome vom Grad  $2m + 1$  exakt, so sind ihre Gewichte  $a_k$  positiv.

Beweis:  $a_j = \sum_{k=0}^m a_k \prod_{\substack{i=0 \\ i \neq j}}^m \left( \frac{t_k - t_i}{t_j - t_i} \right)^2$

## 7.3 Orthogonalpolynome

Gegebene: Skalarprodukt  $\langle f, g \rangle = \int_a^b f(x)g(x)w(x)dx$  auf dem Raum aller stetigen

Funktionen von  $[a, b]$  nach  $\mathbb{R}$  und einer festen Gewichtsfunktion  $\omega > 0$

### Definition 7.11

Die Polynome  $P_0, P_1, P_2, \dots$  bilden ein System von Orthogonalpolynomen, zu dem gegebenen Skalarprodukt, wenn  $P_k$  ein Polynom  $k$ -ten Grades und führendem Koeffizienten  $\neq 0$  ist und für  $l \neq k$   $\langle P_k, P_l \rangle = 0$  ist.

### Beispiel 7.12

Ist  $[a, b] = [-1, 1]$  und  $w(x) = 1$ , so erhält man bis auf Normierung die

Legendre-Polynome:

$$P_k(x) = \frac{1}{2^k k!} \left( \frac{d}{dx} \right)^k \left( (1 - x^2)^k \right)$$

### Satz 7.13

Zu dem gegebenen Skalarprodukt gibt es eindeutig bestimmte Orthogonalpolynome  $P_k$  mit führendem Koeffizienten 1. Sie genügen

$$\text{einer Dreitermrekursion } P_k(x) = (x - a_k)P_{k-1}(x) - b_k^2 P_{k-2}(x)$$

für  $k = 1, 2, 3, \dots$  mit  $P_{-1}(x) = 0, P_0(x) = 1$ .

Beweis: Für  $k \geq 1$  gilt  $xP_{k-1} = \sum_{j=0}^k c_j P_j$  mit den Koeffizienten

$$c_j = \frac{\langle xP_{k-1}, P_j \rangle}{\langle P_j, P_j \rangle} = \frac{\langle P_{k-1}, P_j \rangle}{\langle P_j, P_j \rangle}$$

Da  $xP_j$  für  $j \leq k - 3$  ein Polynom vom Grad  $\leq k - 2$  und damit  $c_j = 0$  ist,

$$\text{folgt } xP_{k-1} = c_k P_k + c_{k-1} P_{k-1} + c_{k-2} P_{k-2}$$

unter den gegebenen Normierungsbedingungen also

$$P_k = (x - c_{k-1})P_{k-1} - c_{k-2}P_{k-2}$$

Da sich  $xP_{k-2}$  und  $P_{k-1}$  wegen der Normierungsbedingung nur durch ein

Polynom vom Grad  $k - 2$  unterscheiden, ist

$$c_{k-2} = \frac{\langle P_{k-1}, xP_{k-2} \rangle}{\langle P_{k-2}, P_{k-2} \rangle} = \frac{\langle P_{k-1}, P_{k-1} \rangle}{\langle P_{k-2}, P_{k-2} \rangle} > 0$$

□

### Beispiel 7.14

Die Legendre-Polynome genügen der Rekursion  $P_0(x) = 1, P_1(x) = x$  und  $(k+1)P_{k+1}(x) = (2k+1)xP_k(x) - kP_{k-1}(x)$ .

### Satz 7.15

Das Orthogonalpolynom  $P_k$  aus Definition 7.11 besitzt  $k$  einfache Nullstellen im Innern von  $[a, b]$ .

Beweis: Wegen  $w(x) > 0$  für  $a < x < b$  sind  $\int_a^b P_k(x)w(x)dx = 0$  muss  $P_k$  für

$k \geq 1$  mindestens eine Nullstelle mit Vorzeichenwechsel in  $(a, b)$  besitzen.

Seien nun  $x_1, \dots, x_m$  die Nullstellen von  $P_k$  im Innern von  $[a, b]$  in aufsteigender Reihenfolge. Sei

$$Q(x) = \prod_{i=1}^m (x - x_i)$$

Dann ist  $P_k(x)Q(x)$  Vorzeichenkonstant und daher

$$\int_a^b P_k(x)Q(x)w(x)dx \neq 0$$

Dies ist aber nach Definition von  $P_k$  nur dann möglich, wenn  $m \geq k$  ist, d.h.  $m = k$  ist.

□

### interpolatorische Quadraturformel

$$\int_{-1}^1 g(t)dt \rightarrow \sum_{k=0}^m a_k g(t_k), a_k = \int_{-1}^1 \prod_{\substack{j=0 \\ j \neq k}}^m \frac{t-t_j}{t_k-t_j} dt$$

mindestens exakt für Polynome vom Grad  $m$  und maximal für Polynome vom Grad  $2m+1$ .

Frage: Lässt sich der Grad  $2m+1$  erreichen?

### Legendre-polynome $P_k, k = 0, 1, 2, \dots$ von Grad $k$

$$\int_{-1}^1 P_k(x)P_l(x)dx = \delta_{kl}, k, l = 0, 1, 2, \dots \text{ haben } k \text{ Nullstellen im Innern des}$$

Intervalls  $[-1, 1]$ .

## 7.4 Gaußsche Quadraturformeln

### Satz 7.16

Es gibt genau eine Quadraturformel  $\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g(t_k)$  die für alle Polynome vom Grad  $2m + 1$  exakt ist. Die Knoten  $t_k$  dieser Quadraturformel sind die Nullstellen des Legendre-Polynoms  $(m + 1)$ -ten Grades und liegen im Innern des Intervalls  $[-1, 1]$ . Die Gewichte  $a_k$  sind positiv.

Beweis: Jedes Polynom  $P$  vom Grad  $2m + 1$  lässt sich eindeutig in der Form

$$P(t) = Q(t) \prod_{j=0}^m (t - t_j) + R(t)$$

mit Polynomen  $Q$  und  $R$  vom Grad  $m$  darstellen. Es gilt dann:

$$\int_{-1}^1 P(t) dt = \int_{-1}^1 Q(t) \prod_{j=0}^m (t - t_j) dt + \int_{-1}^1 R(t) dt$$

$$\sum_{k=0}^m a_k P(t_k) = \sum_{k=0}^m a_k R(t_k)$$

Damit ist die gegen Quadraturformel genau dann für alle Polynome vom Grad  $2m + 1$  exakt, wenn sie

- 1) für alle Polynome  $m$ -ten Grades exakt ist, und wenn
- 2) für alle Polynome  $Q$  vom Grad  $m$

$$\int_{-1}^1 Q(t) \prod_{j=0}^m (t - t_j) dt = 0 \text{ ist.}$$

Die Bedingung 2) ist aber genau dann erfüllt, wenn die  $t_k$  die Nullstellen des Legendre Polynoms vom Grad  $m + 1$  sind. Die Positivität der  $a_k$  wurde bereits im Satz 7.10 bewiesen.

□

merke:  $m + 1$  Stützstellen

=exakt für Polynome vom Grad  $2m + 1$

=Fehlerordnung  $O(h^{2m+2})$

### Beispiel 7.17

Die Mittelpunkregel

$$\int_{-1}^1 g(t) dt \rightarrow 2g(0)$$

ist die Gaußsche Quadraturformel für den Fall  $m = 1$ .

### Beispiel 7.18

$$\frac{8}{\pi \ln 2} \int_0^1 \frac{\ln(1+x)}{1+x^2} dx = 1$$

Mit 2 Teilintervallen liefert die 4-Punkt Gaußformel der Fehlerordnung 10 bereits den auf 10 Stellen genauen Wert!

### Satz 7.19

Für die Gaußsche Quadraturformel

$$\int_{-1}^1 g(t) dt \rightarrow \sum_{k=0}^m a_k g(t_k)$$

vom Grad  $2m + 1$  gilt:

$$\left| \int_a^b f(x) dx - \frac{b-a}{2} \sum_{k=0}^m a_k f\left(\frac{a+b}{2} + t_k \frac{b-a}{2}\right) \right| \leq 2(b-a) E_{2m+1}(f)$$

Dabei ist  $E_{2m+1}(f)$  der Fehler der besten Approximation der Funktion  $f : [a, b] \rightarrow \mathbb{R}$  durch ein Polynom  $(2m + 1)$ -ten Grades.

Beweis: Für alle Polynome vom Grad  $2m + 1$  gilt:

$$\begin{aligned} & \left| \int_a^b f(x) dx - \frac{b-a}{2} \sum_{k=0}^m a_k f\left(\frac{a+b}{2} + t_k \frac{b-a}{2}\right) \right| \\ &= \left| \int_a^b (f - P)(x) dx - \frac{b-a}{2} \sum_{k=0}^m a_k (f - P)\left(\frac{a+b}{2} + t_k \frac{b-a}{2}\right) \right| \\ &\leq \int_a^b |f(x) - P(x)| dx + \frac{b-a}{2} \sum_{k=0}^m |a_k| \left| (f - P)\left(\frac{a+b}{2} + t_k \frac{b-a}{2}\right) \right| \\ &\leq \left( (b-a) + \frac{b-a}{2} \sum_{k=0}^m |a_k| \right) \max_{a \leq x \leq b} |f(x) - P(x)| \end{aligned}$$

und daher

$$\begin{aligned} & \left| \int_a^b f(x) dx - \frac{b-a}{2} \sum_{k=0}^m a_k f\left(\frac{a+b}{2} + t_k \frac{b-a}{2}\right) \right| \\ &\leq \left( 1 + \frac{1}{2} \sum_{k=0}^m |a_k| \right) (b-a) E_{2m+1}(f) \end{aligned}$$

Wegen der Positivität der  $a_k$  ist

$$1 + \frac{1}{2} \sum_{k=0}^m |a_k| = 1 + \frac{1}{2} \sum_{k=0}^m a_k = 1 + \frac{1}{2} \int_{-1}^1 1 dt = 2$$

□

## 8. Gewöhnliche Differentialgleichungen

Anfangwertproblem 1. Ordnung

Gegeben sei eine Funktion  $f : [a, b] \times \Omega \rightarrow \mathbb{R}^n$ ,  $\Omega \subseteq \mathbb{R}^n$  und ein Anfangswert  $y_0 \in \Omega$ . Gesucht ist eine Funktion  $y : [a, b] \rightarrow \Omega$  mit  $y'(x) = f(x, y(x))$  für  $a \leq x \leq b$

$$y(a) = y_0$$

Existenz und Eindeutigkeit (oder lokale Ex. und Eind.)

-> Analysis II

hier: Konstruktion von Näherungsverfahren

### Einschrittverfahren

Gesucht: Näherungswerte  $Y_k$  für die Lösungswerte  $y(x_k)$  auf dem Gitter

$$a = x_0 < x_1 < \dots < x_N = b, \quad (x_{k+1} = x_k + h_k)$$

### Euler-Verfahren

$$y(x_{k+1}) = y(x_k) + \int_{x_k}^{x_{k+1}} y'(x) dx = y(x_k) + \int_{x_k}^{x_{k+1}} f(x, y(x)) dx \approx y(x_k) + h_k f(x_k, y(x_k))$$

$$\frac{y(x_{k+1}) - y(x_k)}{h_k} \approx y'(x_k) = f(x_k, y(x_k))$$

$$Y_{k+1} = Y_k + h_k f(x_k, Y_k)$$

sehr grobe Approximation!

### verbessertes Euler-Verfahren

$$y(x_{k+1}) = y(x_k) + \int_{x_k}^{x_{k+1}} f(x, y(x)) dx \approx y(x_k) + h_k f(x_k + \frac{1}{2}h_k, y(x_k + \frac{1}{2}h_k))$$

$$\frac{y(x_{k+1}) - y(x_k)}{h_k} \approx y'(x_k + \frac{1}{2}h_k) = f(x_k + \frac{1}{2}h_k, y(x_k + \frac{1}{2}h_k))$$

$$f(x_k + \frac{1}{2}h_k, y(x_k + \frac{1}{2}h_k)) \approx f(x_k + \frac{1}{2}h_k, y(x_k)) + \frac{1}{2}h_k f(x_k, y(x_k))$$

$$Y_{k+1} = Y_k + h_k f(x_k + \frac{1}{2}h_k, Y_k + \frac{1}{2}h_k f(x_k, Y_k))$$

Einschrittverfahren:

$$Y_{k+1} = Y_k + h_k \phi(x_k, Y_k, h_k)$$

Probleme:

-In jedem Schritt wird ein neuer Fehler eingeschleppt

-Einmal vorhandene Fehler pflanzen sich fort

### Definition

Die Funktion  $y(x)$  sei also die Lösung des gegebenen Anfangswertproblems und  $p(x, y, h)$  die Verfahrensfunktion des Einschrittverfahrens. Dann heißt

$$\tau(x, h) := \frac{1}{h}(y(x+h) - y(x)) - \varphi(x, y(x), h)$$

der Abbruchfehler des Verfahrens in der Lösung  $y$ . Es ist

$$y(x+h) = y(x) + h\varphi(x, y(x), h) + h\tau(x, h)$$

### Beispiel Euler

$$\tau(x, h)|_i = \frac{1}{h}(y(x+h) - y(x)) - f(x, y(x))|_i$$

$$= \frac{1}{h}(y(x+h) - y(x)) - y'(x)|_i$$

$$= \frac{1}{2}y''(y_i)h$$

$$\Rightarrow \tau(h, x) = \mathcal{O}(h)$$

verbesserter Euler:

$$\tau(h, x) = \mathcal{O}(h^2)$$

Ziel:

globalen Verfahrensfehler  $\max \|Y_k - y(x_k)\|$  durch die Abbruchfehler  $\tau(x_k, h_k)$  abschätzen.

### Satz

Die Verfahrensfunktion des Einschrittverfahrens genüge der Lipschitzbedingung  $\|\varphi(x, y, h) - \varphi(x, z, h)\| \leq L\|y - z\|$  für alle  $y, z \in \mathbb{R}^n, h \in [0, H], x \in [a, b - h]$ .

Dann gilt die Fehlerabschätzung  $\|y(x_k) - Y_k\| \leq e^{L(x_k - a)} \sum_{j=0}^{k-1} h_j \|\tau(x_j, h_j)\|$ ,

$$k = 1, 2, \dots, N.$$

Beweis: Für  $k = 0, 1, \dots, N - 1$  gilt:

$$\|y(x_{k+1}) - Y_{k+1}\| = \|(y(x_k) + h_k\varphi(x_k, y(x_k), h_k) + h_k\tau(x_k, h_k)) - (Y_k + h_k\varphi(x_k, Y_k, h_k))\| \leq \|y(x_k) - Y_k\| + h_k\|\varphi(x_k, y(x_k), h_k) - \varphi(x_k, Y_k, h_k)\| + h_k\|\tau(x_k, h_k)\|$$

$$\leq \|y(x_k) - Y_k\| + h_kL\|y(x_k) - Y_k\| + h_k\|\tau(x_k, h_k)\|$$

$$\Rightarrow \|y(x_{k+1}) - Y_{k+1}\| \leq e^{h_kL}\|y(x_k) - Y_k\| + h_k\|\tau(x_k, h_k)\| \quad (*) \text{ wegen}$$

$1 + t \leq e^t$ . Aus (\*) folgt durch Induktion über  $k$  die Behauptung  $k = 0 \checkmark$

$$k \rightarrow k+1: \|y(x_{k+1}) - Y_{k+1}\| \leq e^{h_k L} e^{L(x_k - a)} \underbrace{\sum_{j=0}^{k-1} h_j \|\tau(x_j, h_j)\|}_{=e^{L(x_{k+1} - a)}} + h_k \|\tau(x_k, h_k)\| \quad \square$$

$$\leq e^{L(x_{k+1} - a)} \sum_{j=0}^k h_j \|\tau(x_k, h_j)\|$$

Folgerung: Wegen  $\sum_{j=0}^{k-1} h_j \|\tau(x_j, h_j)\| \leq (x_k - a) \max_{j=0, \dots, k-1} \|\tau(x_j, h_j)\|$  gilt unter der Annahme

$\|\tau(x, h)\| \leq Kh^q$  mit  $h = \max h_j$  die Abschätzung

$$\max_{k=0, \dots, N} \|y(x_k) - Y_k\| \leq (b-a)e^{L(b-a)} Kh^q = Ch^q$$

Der Fehler strebt also mit der Ordnung  $\mathcal{O}(h^q)$  gegen Null.

Konsistenz (=Verhalten des Abbruchfehlers)

+ Stabilität (Verfahrensfehler lässt sich durch Abbruchfehler abschätzen)

Konvergenz (mit der Ordnung des Abbruchfehlers)

Euler  $\mathcal{O}(h)$ , verb. Euler  $\mathcal{O}(h^2)$

### klassisches Runge-Kutta-Verfahren

$$\varphi(x, y, h) = \frac{1}{6}K_1 + \frac{1}{3}K_2 + \frac{1}{3}K_3 + \frac{1}{6}K_4$$

$$K_1 = f(x, y)$$

$$K_2 = f\left(x + \frac{h}{2}, y + \frac{h}{2}K_1\right)$$

$$K_3 = f\left(x + \frac{h}{2}, y + \frac{h}{2}K_2\right)$$

$$K_4 = f(x + h, y + hK_3)$$

Fehlerordnung:  $\tau(x, h) = \mathcal{O}(h^4)$

geht für  $y'(x) = f(x)$  in die Simpsonregel über.