

Numerik II

(Prof. Arnold – WS 2005/2006)
(by Georg Kusch)

TODO : Def. L_1 -Norm

gram-schmidtsches orth-verf
bulge-chasing verfahren
Horner-Schema zu Bem 9.18 ?

Wiederholung : Komplexe Zahlen

$z = a + ib = r \cdot e^{ij}$, $r = \text{Betrag von } z$, $j = \arg z$
(d.h. $a = r \cos j$ und $b = r \sin j$)

Rechenregeln :

$$(a_1, b_1) \pm (a_2, b_2) = (a_1 \pm a_2, b_1 \pm b_2)$$
$$(a_1, b_1) \cdot (a_2, b_2) = (a_1 a_2 - b_1 b_2, a_1 b_2 + a_2 b_1)$$

Dreiecksungleichung :

$$\forall z_1, z_2 \in \mathbb{C} \text{ gilt : } |z_1 + z_2| \leq |z_1| + |z_2|$$

Eulersche Beziehung :

$$e^{ix} = \cos x + i \sin x$$

!! In \mathbb{C} gibt es keine Ordnungsrelationen . **!!**

6.5 Kondition des Eigenwertproblems

geg.: $A \in \mathbb{C}^{n \times n}$

ges.: $\lambda \in \mathbb{C}$, $x \in \mathbb{R}^n \setminus \{0\}$ mit $Ax = \lambda x$

$$k_{abs} = \|I'(A)\| = \frac{\|x_0\| \cdot \|y_0\|}{|\langle x_0, y_0 \rangle|} = \frac{1}{|\cos(\angle(x_0, y_0))|}$$

mit $x_0 \in \mathbb{C}^n$ Eigenvektor der Matrix $A \in \mathbb{C}^{n \times n}$ und y_0 konjugierter Eigenvektor $\in \mathbb{C}^n$

$$k_{abs} = \frac{\|A\|}{|\lambda|} \|I'(A)\| \quad ???$$

Bemerkung 6.6 : Normale Matrizen

$A \in \mathbb{C}^{n \times n}$ heißt **normal** , wenn $AA^* = A^*A$ ($A^* = \bar{A}^T$)

Unter Verwendung der Schurschen Normalform zeigt man , dass A genau dann normal ist ,

wenn gilt : $A = U^* \Lambda U$

, mit einer unitären Matrix $U \in \mathbb{C}^{n \times n}$, $U^*U = I$, $\Lambda = \text{diag } I_i$ mit $I_i \in \mathbb{C}$

Ist $x_0 \in \mathbb{C}^n$ Eigenvektor von A zum einfachen Eigenwert λ_0 , so ist $U^* \Lambda U x_0 = \lambda_0 x_0$.

Folgerung 6.7 :

Das Eigenwertproblem für einfache Eigenwerte normaler Matrizen ist gut konditioniert mit $k_{abs} = 1$.

Beispiel 6.8 : schlecht konditionierte EWP

a)

$$\text{Sei } A = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \text{ und } B = \begin{pmatrix} 0 & 1 \\ d & 0 \end{pmatrix} \quad \text{mit } |d| \ll 1$$

Die Eigenwerte sind $\lambda_{1,2}(A) = 0$ und $\lambda_{1,2}(B) = \pm\sqrt{d}$.

$$\left(0 = \det(\lambda I - B) = \begin{vmatrix} \lambda & -1 \\ -d & \lambda \end{vmatrix} \right)$$

$$k_{abs} \geq \frac{|\lambda_1(B) - \lambda_1(A)|}{\|B - A\|_2} = \frac{\sqrt{d}}{d} = \frac{1}{\sqrt{d}} \rightarrow \infty \quad \text{für } d \rightarrow 0$$

b) (Wilkinson)

$$\text{Sei } P(\lambda) = \prod_{i=1}^{20} (\lambda - i) \in \Pi_{20} \quad (\lambda_i = i \text{ Nullstellen des Polynoms})$$

Diese Nullstellen sind gegenüber kleinen Störungen der Polynomkoeffizienten extrem empfindlich.

So hat z.B. $\tilde{P}(\lambda) = P(\lambda) - 2^{-23} \lambda^{19}$ die Nullstellen

$$\lambda \approx 16.73 \pm 2.8i \quad \text{und} \quad \lambda \approx 19.50 \pm 1.9i$$

6.2 Vektoriteration („power method“ / „von Mises“)

Algorithmus :

geg.: $A \in \mathbb{R}^{n \times n}$ und Anfangsnäherung $y^{(0)} \in \mathbb{R}^n$

$$\text{Setze } x^{(0)} := \frac{y^{(0)}}{\|y^{(0)}\|_2}, \quad k := 0$$

Schritt k :

$$y^{(k+1)} := Ax^{(k)}$$

$$x^{(k+1)} := \frac{y^{(k+1)}}{\|y^{(k+1)}\|_2}$$

$$\lambda^{(k+1)} := (x^{(k)})^T y^{(k+1)} \quad \text{Eigenwert}$$

$$k := k + 1$$

Eigenschaften :

Mittels vollständiger Induktion folgt: $x^{(k)} = \frac{A^k x^{(0)}}{\|A^k x^{(0)}\|_2}$

(konvergiert gegen den betragsgrößten Eigenvektor – siehe Information Retrieval HITS-Algorithmus)

Sei A diagonalisierbar mit Eigenwerten λ_i $|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$
 und den Eigenvektoren z_i $(i = 1, \dots, n)$.

Dann ist $x^{(0)} = \sum_{i=1}^n a_i z_i$ mit gewissen $a_i \in \mathbb{C}$.

Ist λ_1 dominanter Eigenwert von A , d.h. $|\lambda_1| > |\lambda_i|$ $(i = 2, \dots, n)$, so ist insbesondere $\lambda_1 \in \mathbb{R}$ und es gilt

$$x^{(k)} = \frac{a_1 \lambda_1^k z_1 + \sum_{i=2}^n a_i \lambda_i^k z_i}{\|a_1 \lambda_1^k z_1 + \sum_{i=2}^n a_i \lambda_i^k z_i\|_2} \quad \left(A^k \cdot \sum_i a_i z_i = \sum_i a_i (A^k z_i) = \sum_i a_i \lambda_i^k z_i \right)$$

$$= \operatorname{sgn}(a_1 \lambda_1^k) \frac{z_1 + r^{(k)}}{\|z_1 + r^{(k)}\|_2} \quad \text{mit } r^{(k)} := \sum_{i=2}^n \frac{a_i}{a_1} \left(\frac{\lambda_i}{\lambda_1} \right)^k z_i \rightarrow 0 \text{ für } k \rightarrow \infty$$

, wobei $a_1 \neq 0$ vorausgesetzt wird.

Ergebnis :

$x_\infty := \lim_{k \rightarrow \infty} x^{(k)}$ ist Eigenvektor zum Eigenwert λ_1

und $\lim_{k \rightarrow \infty} n^{(k)} = x_\infty^T A x_\infty = x_\infty^T \lambda_1 x_\infty = \lambda_1 \|x_\infty\|_2^2 = \lambda_1$

D.h., die Vektoriteration berechnet näherungsweise den betragsmäßig größten Eigenwert von A und eines zugehörigen Eigenvektors.

Die Normierung der $x^{(k)}$ verhindert einen Überlauf.

Satz 6.10 : Rayleigh-Quotient

Unter den Voraussetzungen von Bem. 6.9 gilt für eine symmetrische Matrix $A \in \mathbb{R}^{n \times n}$:

$$n^{(k+1)} = \lambda_1 \left(1 + O \left(\left(\frac{\lambda_2}{\lambda_1} \right)^{2k} \right) \right)$$

, d.h. die Rayleigh-Quotienten $n^{(k+1)} = \frac{(x^{(k)})^T A x^{(k)}}{(x^{(k)})^T x^{(k)}}$ konvergieren quadratisch gegen λ_1 .

Bemerkung 6.11 Inverse Vektoriteration (Wielandt)

geg.: $A \in \mathbb{R}^{n \times n}$, $\lambda \in \mathbb{C}$ sei Eigenwert von A

$B \in \mathbb{R}^{n \times n}$, $B := A - \bar{\lambda}I$, $\lambda - \bar{\lambda}$ Eigenwert von B („shift“ der Matrix A)

A regulär $\Rightarrow C := A^{-1}$, λ^{-1} Eigenwert von C

Die formale Anwendung der Vektoriteration auf B^{-1} führt zu einer Approximation des betragskleinsten Eigenwertes einer gegebenen Matrix.

praktisch:

$B := A - \bar{\lambda}I$ mit einer Näherung $\bar{\lambda}$ eines Eigenwertes λ_j von A

Eigenwerte von B : $\lambda_i = \bar{\lambda}$ ($i = 1, \dots, n$) „shift“

Eigenwerte von B^{-1} : $\frac{1}{\lambda_i - \bar{\lambda}}$ ($i = 1, \dots, n$)

A, B und B^{-1} haben die gleichen Eigenvektoren.

Für $\bar{\lambda} \approx \lambda_j$ ist $\left| \frac{1}{\lambda_j - \bar{\lambda}} \right| \gg \left| \frac{1}{\lambda_i - \bar{\lambda}} \right| \quad \forall i \neq j$

\Rightarrow schnelle Konvergenz der Vektoriteration

Implementierung:

$$y^{(k+1)} := B^{-1}x^{(k)} \quad | \cdot B$$

$$(A - \bar{\lambda}I)y^{(k+1)} = x^{(k)}$$

$$\underbrace{\bar{\lambda}^{(k+1)}}_{\approx \text{Eigenwert von } A} := \bar{\lambda} + \frac{1}{\underbrace{\lambda^{(k+1)}}_{\approx \text{Eigenwert von } B^{-1}}}$$

Aufwand:

- Pro Iterationsschritt Lösung eines Gleichungssystems (kubisch)

1x LR-Zerlegung (=LU-Zerlegung) von $A - \bar{\lambda}I$ $O\left(\frac{n^3}{3}\right)$

und pro Iterationsschritt

1x Vorwärtssubstitution $Lz = x^{(k)}$

1x Rückwärtssubstitution $Uy = z$ (bzw. $Ry = z$)

jeweils $O(n^2)$

Problem:

$A - \bar{\lambda}I$ ist sehr schlecht konditioniert für sehr gute Näherungen der Eigenwerte, da $A - \lambda_j I$ singulär ist.

Allerdings nicht die Richtung (bleibt genau), sondern nur der Betrag. Dieser wird jedoch im weiteren Verlauf der Berechnung nicht benötigt.

Konsequenz : LR-Zerlegung mit Spaltenpivotisierung für $A - \bar{I}I$
 Verschwindet ein Pivotelement , so wird es durch *eps* ersetzt.

Modifikation :

Ersetze \bar{I} für $k > 0$ durch die (bessere) Approximation $\bar{n}^{(k)}$ des Eigenwertes I_j .

$$(A - \bar{n}^{(k)}I)y^{(k+1)} = x^{(k)} \quad , \quad \bar{n}^{(k+1)} := \bar{n}^{(k)} + \frac{1}{n^{(k+1)}}$$

Im allgemeinen unattraktiv , wegen des hohen Zusatzaufwands zur neuen LR-Zerlegung von $A - \bar{n}^{(k)}I$ in **jedem** Iterationsschritt.

Idee :

Vor Beginn der Iteration die Matrix A durch Ähnlichkeitstransformationen auf einfachere Gestalt bringen.
 (Ähnlichkeitstransformationen lassen die Eigenwerte unverändert.)

Bemerkung 6.12 : Transformation auf Hessenberg-Form

Transformationsschritt :

Bestimme zu gegebener Matrix $A = (a_{ij}) \in \mathbb{R}^{n \times n}$ ($A = QR$ siehe Householder)

einen Vektor $\tilde{v}_1 \in \mathbb{R}^{n-1}$ so , dass

$$\left(I_{n-1} - 2 \frac{\tilde{v}_1 \tilde{v}_1^T}{\tilde{v}_1^T \tilde{v}_1} \right) \cdot \begin{pmatrix} a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{pmatrix} = \mathbf{a} \cdot \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

Mit $v_1 = \begin{pmatrix} 0 \\ \tilde{v}_1 \end{pmatrix} \in \mathbb{R}^n$ und $Q_1 := I - 2 \frac{v_1 v_1^T}{v_1^T v_1}$ ist $A^{(2)} := Q_1 A Q_1^T$,

denn die Rechts-Multiplikation mit Q_1^T lässt wegen $v_{1,1} = 0$ die erste Spalte unverändert und transformiert nur die restlichen Spalten.

Die Matrizen $A^{(2)}$ und $A^{(1)} := A$ sind zueinander ähnlich und haben insbesondere dieselben Eigenwerte.

allgemein :

Analog zur QR-Zerlegung von A bestimmt man sukzessive (nacheinander) $n - 2$ Spiegelungsvektoren

$$v_k = \begin{pmatrix} \mathbf{s}_k \\ \tilde{v}_k \end{pmatrix} \in \mathbb{R}^n \quad \text{mit} \quad \tilde{v}_k \in \mathbb{R}^{n-k} \quad \text{und die zugehörigen Spiegelungsmatrizen}$$

$$Q_k := I - 2 \frac{v_k v_k^T}{v_k^T v_k} \in \mathbb{R}^{n \times n} \quad \text{so , dass mit} \quad A^{(k+1)} := Q_k A^{(k)} Q_k^T \quad \text{die}$$

Matrix $\tilde{A} := A^{(n-2)}$ Hessenberg-Form besitzt.

Mittels vollständiger Induktion folgt , dass \tilde{A} und A zueinander ähnlich sind :

$$\tilde{A} := Q A Q^T \quad \text{mit} \quad Q = Q_{n-2} Q_{n-1} \dots Q_2 Q_1 .$$

Ist \tilde{x}_i Eigenvektor von \tilde{A} zum Eigenwert I_i , so ist $x_i := Q^T \tilde{x}_i = Q_1^T Q_2^T \dots Q_{n-2}^T \tilde{x}_i$ Eigenvektor von A zum Eigenwert I_i , denn

$$\begin{aligned} \tilde{A}\tilde{x}_i = I_i\tilde{x}_i &\Leftrightarrow QAQ^T\tilde{x}_i = I_i\tilde{x}_i && | \cdot Q^T \\ & \underbrace{AQ^T}_{x_i} \tilde{x}_i = I_i \underbrace{Q^T}_{x_i} \tilde{x}_i \end{aligned}$$

Hat \tilde{A} Hessenberg-Form, so auch $\tilde{A} - \bar{I}I$ für jedes $\bar{I} \in C$.

Die LR-Zerlegung von $\tilde{A} - \bar{I}I$ erfordert $O(n^2)$ Rechenoperationen.

Aufwand für die Hessenberg-Form: $O\left(\frac{4}{3}n^3\right)$.

Spezialfall: Ist A symmetrisch, so ist \tilde{A} Tridiagonalmatrix.

6.3 QR-Algorithmus

geg.: $A \in R^{n \times n}$, setze $A_0 := A$

Algorithmus:

- Berechne die QR-Zerlegung von $A_k = Q_k R_k$
- $A_{k+1} = R_k Q_k$

Wegen $A_{k+1} = R_k Q_k = Q_k^T A_k Q_k$ folgt mittels vollständiger Induktion, dass alle A_k ähnlich zu A sind.

Nachteil: pro Iterationsschritt QR-Zerlegung und Matrix-Matrix-Multiplikation

Alternative: Transformiere A zunächst auf Hessenberg-Form \tilde{A} und wende den QR-Algorithmus auf \tilde{A} an.

Vorteil: QR-Zerlegung mit $n-1$ Givens-Rotationen $\Omega_{n-1,n}, \dots, \Omega_{1,2}$

$$\Omega_{1,2} \tilde{A} = \Omega_{1,2} \begin{pmatrix} * & * & \dots & * \\ * & * & & \vdots \\ 0 & * & & \\ \vdots & & \ddots & * \\ 0 & \dots & 0 & * & * \end{pmatrix} = \begin{pmatrix} * & * & \dots & * \\ 0 & * & & \vdots \\ 0 & * & & \\ \vdots & & \ddots & * \\ 0 & \dots & 0 & * & * \end{pmatrix}$$

(Nur Veränderung der ersten zwei Zeilen.)

$$\text{usw. , d.h. } \Omega_{n-1,n} \dots \Omega_{1,2} \cdot \tilde{A} = R_0, \quad Q_0^T = \Omega_{n-1,n} \dots \Omega_{1,2}$$

$$\tilde{A}_1 = R_0 Q_0 = R_0 \Omega_{1,2}^T \dots \Omega_{n-1,n} \quad (\text{ebenfalls Hessenberg-Form})$$

Ergebnisse:

- Anwendung des QR-Algorithmus auf Hessenberg-Matrix liefert eine Folge von Hessenberg-Matrizen.
- Rechenaufwand: pro Iterationsschritt $O(n^2)$ für $2 \cdot (n-1)$ Givens-Rotationen

Ist A symmetrisch und tridiagonal, so reduziert sich der Rechenaufwand pro Iterationsschritt auf $O(n)$.

Satz 6.14 : Konvergenz des QR-Algorithmus (siehe 5.7 in Numerik1-Skript)

geg.: $A \in \mathbb{R}^{n \times n}$ symmetrisch und mit den Eigenwerten $\lambda_1, \dots, \lambda_n$, wobei $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$.

Für die Matrizen Q_k, R_k und A_k gilt :

- a) $\lim_{k \rightarrow \infty} Q_k = I$
- b) $\lim_{k \rightarrow \infty} Q_k = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) = \lim_{k \rightarrow \infty} A_k$
- c) $a_{ij}^{(k)} = O\left(\left|\frac{\lambda_i}{\lambda_j}\right|^k\right)$ für $i > j$

Beweis zu a) :

Es gilt $A^k = \underbrace{Q_1 \dots Q_k}_{P_k} \underbrace{R_k \dots R_1}_{S_k}$, denn für $k = 1$ ist dies trivial und für $k > 1$ folgt aus

$$A^k = P_k S_k \text{ wegen } A_{k+1} = Q_{k+1} R_{k+1} = Q_k^T \dots Q_1^T A Q_1 \dots Q_k = \underbrace{Q_1^T A Q_1}_{=A_2} \dots Q_k = P_k^{-1} A P_k \text{ auch}$$

$$A^{k+1} = A A^k = \underbrace{A P_k}_{P_k A_{k+1}} S_k = P_k Q_{k+1} R_{k+1} S_k = P_{k+1} S_{k+1}$$

P_k ist orthogonal, S_k ist obere Dreiecksmatrix.

$\Rightarrow A^k = P_k S_k$ ist (eine) QR-Zerlegung von A^k .

Bemerkung 6.15 : Allgemeine Konvergenzaussagen

- a) Der QR-Algorithmus konvergiert auch für mehrfache Eigenwerte.
- b) Falls unterschiedliche Eigenwerte denselben Betragswert haben, so konvergieren die A_k nicht gegen Λ , sondern es bleiben 2×2 Blöcke entlang der Hauptdiagonalen stehen (Jordan-Blöcke?). Man kann die Eigenwerte trotzdem aus diesen Blöcken ausrechnen.

Bemerkung 6.16 : Verbesserung des QR-Algorithmus durch Shift-Technik

a) **Modifikation :**

$$A_k - \tilde{I}_k I = Q_k R_k$$

QR-Zerlegung für $\tilde{I}_k \approx \lambda_i$

$$A_{k+1} = R_k Q_k + \tilde{I}_k I$$

Es gilt $A_{k+1} = Q_k^T (A_k - \tilde{I}_k I) Q_k + \tilde{I}_k I = Q_k^T A_k Q_k$

$\Rightarrow A_{k+1}$ und A_k sind einander ähnlich.

b) **Konvergenz**

Analog zu Satz 6.14 zeigt man unter Verwendung von $(A - \tilde{I}_k I) \dots (A - \tilde{I}_1 I) = Q_1 \dots Q_k R_k \dots R_1$

$$a_{ij}^{(k)} = O\left(\left|\frac{\lambda_i - \tilde{I}_1}{\lambda_j - \tilde{I}_1}\right| \dots \left|\frac{\lambda_i - \tilde{I}_k}{\lambda_j - \tilde{I}_k}\right|\right) \quad \forall i > j$$

Falls $\tilde{I}_l \approx \lambda_l$ ($l \geq 1$), so ergibt sich eine Konvergenzbeschleunigung.

c) Implementierung : Expliziter Shift (nach Wilkinson)

Ziel : $a_{n,n-1}^{(k+1)} \approx 0 \Rightarrow a_{n,n}^{(k+1)}$ ist eine sehr gute Näherung für Eigenwert von A

Wähle \tilde{I}_k als denjenigen Eigenwert von $\begin{pmatrix} a_{n-1,n-1}^{(k)} & a_{n-1,n}^{(k)} \\ a_{n,n-1}^{(k)} & a_{n,n}^{(k)} \end{pmatrix} \in R^{2 \times 2}$, der näher an $a_{n,n}^{(k)}$ liegt.

\Rightarrow Konvergenz von $a_{n,n-1}^{(k)} \rightarrow 0$ und $a_{n,n}^{(k)} \rightarrow I_i(A)$ quadratisch und häufig sogar kubisch.

Berechnung des Eigenwertes $I_i(A)$ in der Regel mit $O(n)$ Rechenoperationen.

Erweiterung auf nichtsymmetrische Matrizen mit konjugiert komplexen Eigenwerten :

Implizite Shift-Technik (nach Francis) – ermöglicht rein reelle Berechnungen

d) nächste Eigenwerte und Abbruchkriterium

Deflation : Ist $\max_{1 \leq j < n} |a_{n,j}^{(k)}| = O(eps)$, so ist $a_{n,n}^{(k)}$ eine gute Näherung für einen Eigenwert von A .

Der QR-Algorithmus wird mit $A = A_k(1 : (n-1), 1 : (n-1)) \in R^{(n-1) \times (n-1)}$ fortgesetzt.

(D.h. letzte Zeile und letzte Spalte weglassen.)

Falls $n = 2$, so werden die letzten beiden Eigenwerte explizit berechnet.

e) Gesamtaufwand für symmetrische Matrizen

- $O\left(\frac{4}{3}n^3\right)$ Rechenoperationen für Transformation auf Tridiagonal-Form

- $n \cdot O(n) = O(n^2)$ Rechenoperationen für QR-Algorithmus

f) Berechnung der Eigenvektoren

$$A = U\Lambda U^T \Rightarrow AU = U\Lambda$$

\Rightarrow Spaltenvektoren von U sind die Eigenvektoren von A

$$PAP^T = \tilde{A} \quad (\text{Transformation auf Tridiagonalgestalt})$$

$$\tilde{A} = A_0 = Q_0 Q_1 \dots Q_k A_{k+1} Q_k^T \dots Q_1^T Q_0^T$$

$$\Rightarrow U = P^T Q_0 Q_1 \dots Q_k$$

Beispiel 6.17 : QR-Algorithmus mit Shift-Technik

$$A = \begin{pmatrix} 12 & 1 & \dots & 0 \\ 1 & 9 & 1 & \vdots \\ & 1 & 6 & 1 \\ \vdots & & 1 & 3 & 1 \\ 0 & \dots & & 1 & 0 \end{pmatrix}$$

k	$a_{5,4}^{(k)}$	$a_{5,5}^{(k)}$	\tilde{I}_k
1	-4.5 E-3	0	-0.302776
2	1.1 E-10	-0.316869	-0.316876
3		-0.316876	-0.316876
4	9.1 E-23	-0.316876	-0.316876

$$\Rightarrow I_5 \approx -0.316876$$

Fortsetzung mit der Matrix $A_4(1:4, 1:4)$:

$$\Rightarrow a_{4,3}^{(6)} = -1.1 E - 18 \quad \Rightarrow I_4 \approx 2.98386$$

Fortsetzung mit der Matrix $A_6(1:3, 1:3)$:

$$\Rightarrow a_{3,2}^{(8)} = 1.8 E - 20 \quad \Rightarrow I_3 \approx 6.0$$

Explizite Berechnung der beiden letzten Eigenwerte aus $A_8(1:2, 1:2)$:

$$I_2 \approx 9.01614, \quad I_1 \approx 12.316876$$

Zum Vergleich der Konvergenzgeschwindigkeiten : QR-Algorithmus ohne Shift-Technik

$$a_{2,1}^{(12)} = 0.0337458$$

$$a_{3,2}^{(12)} = 0.011408$$

$$a_{3,3}^{(12)} = -0.316876$$

Konvergenzgeschwindigkeit entspricht wegen $\left| \frac{I_5}{I_4} \right| \approx 0.1$ den Abschätzungen von Satz 6.14.

6.4 Lancosz-Methoden

Bemerkung 6.18 : Eigenwertprobleme als Extremwertaufgaben

Zu einer symmetrischen Matrix $A \in R^{n \times n}$ und $x \in R^n \setminus \{0\}$ bezeichnet

$$\mathbf{m}(x) := \frac{\langle x, Ax \rangle}{\langle x, x \rangle} \text{ den Rayleigh-Quotienten.}$$

Es gilt : $\min_{x \neq 0} \mathbf{m}(x) = I_{\min}(A)$ und $\max_{x \neq 0} \mathbf{m}(x) = I_{\max}(A)$ mit dem kleinsten bzw. größtem Eigenwert von A .

Denn A ist diagonalisierbar : $A = U\Lambda U^T$ mit $\Lambda = \text{diag}(I_1, \dots, I_n)$, so dass

wegen $U^T A U = \Lambda$ für $y := U^T x$ folgt :

$$\langle y, y \rangle = \langle U^T x, U^T x \rangle = x^T U U^T x = \langle x, x \rangle \quad \text{und}$$

$$\langle x, Ax \rangle = \langle U y, A U y \rangle = \langle y, U^T A U y \rangle = \langle y, \Lambda y \rangle = \sum_{i=1}^n I_i y_i^2$$

Idee :

Approximiere $I_{\min}(A) = \min \{ \mathbf{m}(x) : x \in R^n, x \neq 0 \}$ durch $\min \{ \mathbf{m}(x) : x \in V_k, x \neq 0 \}$ mit einem niedrigdimensionalen Unterraum $V_k \subset R^n$. Analog $I_{\max}(A)$.

Ansatz :

$$\text{Krylovraum } V_k = V_k(x) = K_k(x, A) = \text{span}\{x, Ax, \dots, A^{k-1}x\}$$

Bemerkung 6.19 : Lanczos-Algorithmus

$$A = A^T \in \mathbb{R}^{n \times n}$$

Bestimmen einer Orthonormalbasis v_1, \dots, v_l von $V_l(x) = K_l(x, A)$ mittels Drei-Term-Rekursion.
(Effizienter als z.B. Gram-Schmidtsches -Verfahren)

$$v_0 := 0, \quad v_1 := \frac{x}{\|x\|_2}$$

for $k = 1 : (l-1)$

$$\mathbf{a}_k := \langle v_k, Av_k \rangle$$

$$w_{k+1} := Av_k - \mathbf{a}_k v_k - \mathbf{b}_k v_{k-1} \quad (\mathbf{b}_1 \text{ nicht nötig, da } v_0 = 0)$$

$$\mathbf{b}_{k+1} := \|w_{k+1}\|_2$$

$$v_{k+1} := \frac{w_{k+1}}{\mathbf{b}_{k+1}}, \text{ falls } \mathbf{b}_{k+1} \neq 0$$

Abbruch, falls $w_{k+1} = 0$

endfor

In Matrixschreibweise : $AQ_k = Q_k T_k$ mit

$$Q_k = [v_1, v_2, \dots, v_k] \in \mathbb{R}^{n \times k} \text{ und } T_k = \begin{pmatrix} \mathbf{a}_1 & \mathbf{b}_2 & & & & & \\ \mathbf{b}_2 & \mathbf{a}_2 & \mathbf{b}_3 & & & & \\ & \mathbf{b}_3 & \mathbf{a}_3 & \mathbf{b}_4 & & & \\ & & \ddots & \ddots & \ddots & & \\ & & & \mathbf{b}_{k-1} & \mathbf{a}_{k-1} & \mathbf{b}_k & \\ & & & & \mathbf{b}_k & \mathbf{a}_k & \end{pmatrix}$$

Es gilt : $T_k = Q_k^T A Q_k$, jedoch sind A und T_k einander **nicht ähnlich**. ($k < n$)

zu zeigen :

- $A^k x \in \text{span}\{v_1, \dots, v_{k+1}\}$
- $\langle w_{k+1}, v_j \rangle = 0 \quad (j = 1, \dots, k)$

Beweis :

zu a) Folgt aus $w_{k+1} = \mathbf{g} \cdot A^k x + r_k$ mit einem $\mathbf{g} \neq 0$ und einem $r_k \in K_k(x, A)$
Nachweis mittels vollständiger Induktion

zu b)

$$\langle w_{k+1}, v_j \rangle = \underbrace{\langle Av_k, v_j \rangle}_{= \langle v_k, Av_j \rangle \text{ wegen } A=A^T} - \mathbf{a}_k \underbrace{\langle v_k, v_j \rangle}_{=0 \text{ für } j < k} - \mathbf{b}_k \underbrace{\langle v_{k-1}, v_j \rangle}_{=0 \text{ für } j \neq k-1}$$

$$Av_j \in K_{j+1}(x, A) = \text{span}\{v_1, v_2, \dots, v_{j+1}\} \Rightarrow \text{Skalarprodukt} = 0 \text{ für } j < k-1$$

d.h., $\langle w_{k+1}, v_j \rangle = 0$ für $j = 1, 2, \dots, k-2$

Fall $j = k-1$:

$$\begin{aligned}
\langle w_{k+1}, v_{k-1} \rangle &= \left\langle v_k, \underbrace{w_k + a_{k-1}v_{k-1} + b_{k-1}v_{k-2}}_{\text{siehe obige Def. von } w_k} \right\rangle - \underbrace{b_k}_{=1} \langle v_{k-1}, v_{k-1} \rangle \\
&= \langle v_k, Av_{k-1} \rangle - b_k \\
&= \langle v_k, w_k \rangle - b_k \\
&= 0 \quad , \text{ da } w_k = b_k v_k \text{ und } \langle v_k, v_k \rangle = 1
\end{aligned}$$

Fall $j = k$:

$$\langle w_{k+1}, v_k \rangle = \langle v_k, Av_k \rangle - a_k \langle v_k, v_k \rangle - b_k \cdot 0 = \langle v_k, Av_k \rangle - a_k = 0$$

Anwendung :

Zu $u \in R^k$ betrachte $y = Q_k u \in K_k(x, A) \subset R^n$, dann ist

$$\langle y, y \rangle = \langle Q_k u, Q_k u \rangle = \langle u, Q_k^T Q_k u \rangle = \langle u, u \rangle \quad \text{und}$$

$$\langle y, Ay \rangle = \langle Q_k u, A Q_k u \rangle = \langle u, Q_k^T A Q_k u \rangle = \langle u, T_k u \rangle$$

Also ist

$$I_{\min}^{(k)} = \min_{\substack{y \in V_k(x) \\ y \neq 0}} \frac{\langle y, Ay \rangle}{\langle y, y \rangle} = \min_{\substack{u \in R^k \\ u \neq 0}} \frac{\langle u, T_k u \rangle}{\langle u, u \rangle} = I_{\min}(T_k)$$

Analog ist $I_{\max}^{(k)} = I_{\max}(T_k)$.

Wegen $V_k \subset V_{k+1}$ folgt weiterhin $I_{\min}^{(k+1)} \leq I_{\min}^{(k)}$ und $I_{\max}^{(k+1)} \geq I_{\max}^{(k)}$.

Satz 6.20 : Eigenwertberechnung mittels Lanczos-Verfahren

Sei A eine symmetrische Matrix mit den Eigenwerten $I_1 \leq \dots \leq I_n$ und zugehörigen orthonormalen Eigenvektoren h_1, \dots, h_n .

Aus dem Lanczos-Verfahren zum Startwert $x \neq 0$ ergibt sich eine ONB v_1, \dots, v_k von $V_k(x) = K_k(x, A)$, sowie die Tridiagonalmatrix $T_k \in R^{k \times k}$ mit den Eigenwerten $m_1 \leq \dots \leq m_k$.

Dann gilt :
$$I_n \geq m_k \geq I_n - \frac{(I_n - I_1) \tan^2(\angle(v_1, h_n))}{T_{k-1}^2 \left(1 + 2 \frac{I_n - I_{n-1}}{I_{n-1} - I_1} \right)}$$
 mit dem Tschebyscheff-Polynom T_{k-1} .

6.5 Singulärwertzerlegung

Satz 6.21

Sei $A \in R^{m \times n}$. Dann gibt es zwei orthogonale Matrizen $U \in R^{m \times m}$ und $V \in R^{n \times n}$, so dass $U^T A V = \Sigma = \text{diag}(s_1, \dots, s_p)$ mit $p = \min(m, n)$ und $s_1 \geq \dots \geq s_p \geq 0$.

Diese Zerlegung $U^T A V = \Sigma$ heißt Singulärwertzerlegung von A und die \mathbf{s}_i sind die Singulärwerte von A .

Beweis mittels vollständiger Induktion

zu zeigen :

Es gibt orthogonale Matrizen U_1, V_1 mit $U_1^T A V_1 = \begin{pmatrix} \mathbf{s}_1 & 0 \\ 0 & B \end{pmatrix}$, wobei $B \in \mathbb{R}^{(m-1) \times (n-1)}$ ist.

Sei $\mathbf{s}_1 := \|A\|_2 = \max_{\|v\|_2=1} \|Av\|_2$, dann gibt es Vektoren $v_1 \in \mathbb{R}^n, u_1 \in \mathbb{R}^m$ mit

$$\|v_1\|_2 = 1, \quad u_1 := \frac{1}{\mathbf{s}_1 A v_1}, \quad \|u_1\|_2 = 1$$

und $\{v_1\}, \{u_1\}$ können zu orthogonalen Matrizen $V_1 = [v_1, \dots, v_n], U_1 = [u_1, \dots, u_m]$ erweitert werden.

$$A_1 := U_1^T A V_1 = U_1^T [A v_1, A v_2, \dots, A v_n] = [\mathbf{s}_1 \cdot u_1, A v_2, \dots, A v_n] = \begin{pmatrix} \mathbf{s}_1 & w^T \\ 0 & B \end{pmatrix}$$

(wegen $U_1^T (\mathbf{s}_1 u_1) = \mathbf{s}_1$ und $U_1^T u_j = 0$ für $j = 2, \dots, m$)

mit einem $w \in \mathbb{R}^{n-1}$ und einer Matrix $B \in \mathbb{R}^{(m-1) \times (n-1)}$.

Wegen $\|C\|_2 = \sqrt{\mathbf{I}_{\max}(C^T C)}$ gilt

$$\|A_1\|_2^2 = \mathbf{I}_{\max}(A_1^T A_1) = \mathbf{I}_{\max}(V_1^T A^T U_1 U_1^T A V_1) = \mathbf{I}_{\max}(V_1^T (A^T A) V_1) = \mathbf{I}_{\max}(A^T A)$$

(wegen Ähnlichkeit der Matrizen)

$$= \|A\|_2^2 = \mathbf{s}_1^2$$

\Rightarrow

$$\mathbf{s}_1^2 = \|A_1\|_2^2 \geq \frac{\left\| A_1 \cdot \begin{pmatrix} \mathbf{s}_1 \\ w \end{pmatrix} \right\|_2^2}{\left\| \begin{pmatrix} \mathbf{s}_1 \\ w \end{pmatrix} \right\|_2^2} \quad (\text{Verträglichkeit von Matrix- und Vektornorm})$$

$$= \frac{\left\| \begin{pmatrix} \mathbf{s}_1^2 + w^T w \\ \vdots \end{pmatrix} \right\|_2^2}{\mathbf{s}_1^2 + w^T w} \geq \frac{(\mathbf{s}_1^2 + w^T w)^2}{\mathbf{s}_1^2 + w^T w} = \mathbf{s}_1^2 + w^T w$$

$$\Rightarrow w = 0, \quad U_1^T A V_1 = \begin{pmatrix} \mathbf{s}_1 & 0 \\ 0 & B \end{pmatrix}$$

Bemerkung 6.22 : Eigenschaften der Singulärwertzerlegung

Für die in Matrix 6.21 konstruierte Matrixzerlegung $U^T A V = \Sigma$ gilt :

- a) $Av_i = \mathbf{s}_i u_i$, $A^T u_i = \mathbf{s}_i v_i$ ($i = 1, \dots, p$) mit den Spaltenvektoren u_1, \dots, u_m von U und v_1, \dots, v_n von V .
- b) Es ist $\mathbf{s}_1 \geq \mathbf{s}_2 \geq \dots \geq \mathbf{s}_r > \mathbf{s}_{r+1} = \dots = \mathbf{s}_p = 0 \Leftrightarrow \text{rang}(A) = r$
 $\ker(A) = \text{span}\{v_{r+1}, \dots, v_n\}$
 $\text{im}(A) = \text{span}\{u_1, \dots, u_r\}$
- c) $\mathbf{s}_1 = \|A\|_2$
- d) $\|A\|_F^2 = \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 = \|\Sigma\|_F^2 = \sum_{i=1}^p \mathbf{s}_i^2$
- e) $\mathbf{k}_2(A) = \|A\|_2 \cdot \|A^{-1}\|_2 = \frac{\mathbf{s}_1}{\mathbf{s}_n}$ für reguläre $A \in R^{n \times n}$
- f) $A^T A$ hat die Eigenwerte $\mathbf{s}_1^2, \dots, \mathbf{s}_p^2$ zu den Eigenvektoren v_1, \dots, v_p
 $A^T A$ hat die Eigenwerte $\mathbf{s}_1^2, \dots, \mathbf{s}_p^2$ zu den Eigenvektoren u_1, \dots, u_p

weitere Eigenschaften :

- Gegeben sei eine Matrix $A \in R^{m \times n}$ und orthogonale Matrizen $P \in R^{m \times m}$, $Q \in R^{n \times n}$.
 Dann haben A und $B := PAQ$ die gleichen Singulärwerte.
- Sei $U^T AV = \Sigma = \text{diag}(\mathbf{s}_1, \dots, \mathbf{s}_r, 0, \dots, 0)$ die Singulärwertzerlegung einer Matrix $A \in R^{m \times n}$ mit $\text{rang}(A) = r$.

Dann ist die Pseudoinverse $A^+ \in R^{n \times m}$ gegeben durch $A^+ = V\Sigma^+U^T$

mit $\Sigma^+ = \text{diag}\left(\frac{1}{\mathbf{s}_1}, \dots, \frac{1}{\mathbf{s}_r}, 0, \dots, 0\right)$

(Penrose-Axiome $AA^+A = A$, $x = A^+b$)

Bemerkung : Bi-Diagonalisierung rechteckiger Matrizen

geg. : $A \in R^{m \times n}$, $m \geq n$

ges. : orthogonale Matrizen $P \in R^{m \times m}$, $Q \in R^{n \times n}$ mit

$$PAQ = \begin{pmatrix} B \in R^{n \times n} \\ \dots \\ 0_{(m-n) \times n} \end{pmatrix} = \begin{pmatrix} * & * & & & \\ & \ddots & \ddots & & \\ & & \ddots & * & \\ & & & \ddots & * \\ \dots & \dots & \dots & \dots & \dots \\ 0 & \dots & & 0 & \\ \vdots & & & \vdots & \\ 0 & \dots & & 0 & \end{pmatrix} , B = \text{Bi-Diagonalmatrix}$$

Vorgehensweise :

Wähle P_1 als Householder-Spiegelung , welche die erste Spalte von A in ein Vielfaches des ersten Einheitsvektors transformiert.

$$P_1 A = \begin{pmatrix} * & * & \dots & * \\ 0 & * & \dots & * \\ \vdots & \vdots & & \vdots \\ 0 & * & \dots & * \end{pmatrix}$$

Wähle Q_1 als Householder-Spiegelung so, dass $P_1 A Q_1 = \begin{pmatrix} * & * & 0 & \dots & 0 \\ 0 & * & * & \dots & * \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & * & * & \dots & * \end{pmatrix}$

(Die erste Komponente des Spiegelungsvektors in Q_1 wird $= 0$ gesetzt.)

\Rightarrow 1. Spalte von $P_1 A$ bleibt unverändert

Mittels vollständiger Induktion erhält man die Bi-Diagonalgestalt.

$$P A Q = \begin{pmatrix} B \\ 0 \end{pmatrix} = \underbrace{P_{n-1} \dots P_1}_{=: P} A \underbrace{Q_1 \dots Q_{n-2}}_{=: Q}$$

Singulärwertzerlegung (bidiagonaler Matrizen) - Algorithmus :

siehe „chasing“ / „bulge chasing“.

MATLAB-Befehl : `svd` (single value decomposition)

Idee : Anwendung des QR-Algorithmus auf $B^T B$ ohne explizite Berechnung von $B^T B$
(Anwendung von Givens-Drehungen)

TODO : genaue Beschreibung

Ergebnis :

Bi-Diagonalmatrix $B_1 = \tilde{\Omega}_{n-1,n} \dots \tilde{\Omega}_{1,2} B \Omega_{1,2}^T \dots \Omega_{n-1,n}^T$

Die Transformation $B^T B \rightarrow B_1^T B_1$ entspricht dem ersten Schritt des QR-Algorithmus angewendet auf $B^T B$, jedoch ohne explizite Berechnung von $B^T B$.

Nach Satz 6.14 konvergiert $B_k^T B_k$ gegen $\Lambda = \text{diag}(\mathbf{s}_1^2, \dots, \mathbf{s}_n^2) = \Sigma^T \Sigma$

$$B_k \xrightarrow[k \rightarrow \infty]{} \Sigma$$

Aufwand : für quadratische Matrizen ($m = n$)

- $\frac{4}{3} n^3 + O(n^2)$ Rechenoperationen für Reduktion auf Bi-Diagonalgestalt
- $O(n^2)$ Rechenoperationen für den modifizierten QR-Algorithmus

7. Interpolation

Problemstellung :

Gegeben seien diskrete Funktionswerte und ggf. Ableitungen einer Funktion $f : \mathbb{R} \rightarrow \mathbb{R}$ an den Stützstellen t_i ($i = 0, 1, \dots, n$)

Und zwar $f(t_i), f'(t_i), \dots, f^{(c_i)}(t_i)$ mit einem $c_i \geq 0$ (=Stützweite)

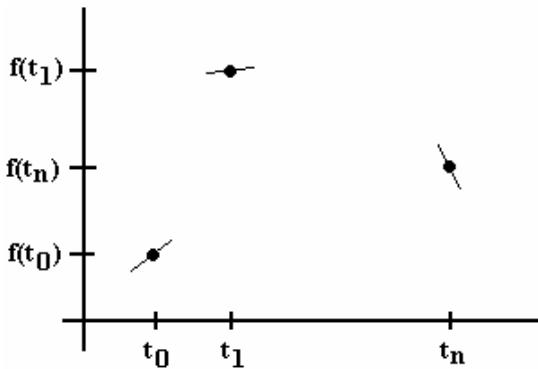


Abb : Beispiel für die zwei Fälle $c_i = 0$ und $c_i = 1$ (im ersten Fall z.B. quadratisches Polynom)

Gesucht : Eine Approximation von f durch eine mit wenig Aufwand und numerisch stabil auswertbare Funktion j , wie z.B.

- Polynome
- stückweise Polynome (bei sehr vielen Stützstellen Aufteilung)
- trigonometrische Funktionen
- Exponentialfunktionen
- rationale Funktionen

Interpolationsbedingungen : $j^{(j)}(t_i) = f^{(j)}(t_i)$, $\begin{pmatrix} j = 0, 1, \dots, c_i \\ i = 0, 1, \dots, n \end{pmatrix}$

7.1 Klassische Polynominterpolation

Satz 7.2 : Interpolationspolynome

Zu beliebigen $n + 1$ Stützpunkten (t_i, f_i) , $(i = 0, 1, \dots, n)$ mit $t_i \neq t_j \quad \forall i \neq j$ gibt es genau ein Polynom höchstens n -ten Grades

$$P(t) \in \Pi_n = \{p(t) : p(t) = a_0 + a_1 t + \dots + a_n t^n, \quad a_i \in R\} \quad \text{mit}$$

$$P(t_i) = f_i, \quad (i = 0, 1, \dots, n)$$

(Zusammenhang siehe Vandermonde Determinante)

Beweis :

a) Eindeutigkeit

Gilt für $P_1, P_2 \in \Pi_n \quad P_1(t_i) = f_i = P_2(t_i) \quad \text{für } i = 0, 1, \dots, n$, so hat das

Polynom $P_1 - P_2 \in \Pi_n$ wegen $P_1(t_i) - P_2(t_i) = 0$ mindestens $n + 1$ verschiedene Nullstellen.

$$\Rightarrow P_1 - P_2 = \text{Nullpolynom}$$

$$\Rightarrow P_1 = P_2$$

b) Existenz

Konstruiere Lagrange-Polynome $L_i \in \Pi_n \quad (i = 0, 1, \dots, n)$

$$\text{mit } L_i(t_k) = \begin{cases} 1 & \text{für } i = k \\ 0 & \text{sonst} \end{cases} = d_{ik}$$

$$\Rightarrow P(t) = \sum_{i=0}^n f_i L_i(t)$$

Lagrangesche Darstellung des Interpolationspolynoms :

$$L_i(t) = \prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - t_j}{t_i - t_j} = \frac{(t - t_0)(t - t_1) \dots (t - t_{i-1})(t - t_{i+1}) \dots (t - t_n)}{(t_i - t_0)(t_i - t_1) \dots (t_i - t_{i-1})(t_i - t_{i+1}) \dots (t_i - t_n)}$$

$$= \frac{w(t)}{(t - t_i) \cdot w'(t_i)} \quad \text{mit dem Knotenpolynom}$$

$$w(t) := \prod_{j=0}^n (t - t_j) = (t - t_0)(t - t_1) \dots (t - t_n) \quad \text{q.e.d.}$$

Hinweis :

$$\{L_0, L_1, \dots, L_n\} \text{ ist Orthonormalbasis bezüglich } \langle P, Q \rangle := \sum_{i=0}^n P(t_i) \cdot Q(t_i)$$

Bemerkung 7.3 : Interpolationspolynom

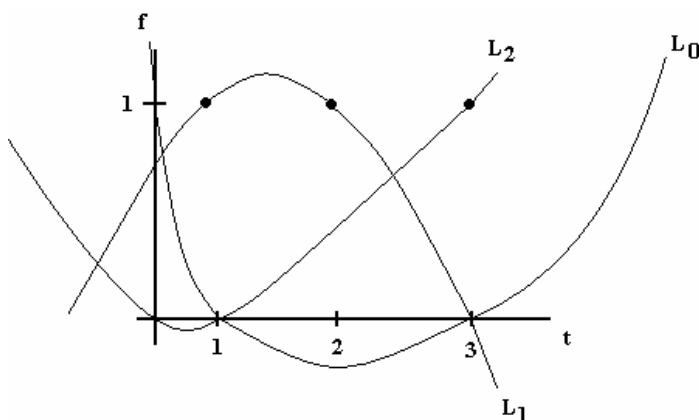
- a) Das (eindeutig bestimmte) Polynom P aus Satz 7.2 heißt Interpolationspolynom zu den Stützpunkten (t_i, f_i) .
- b) Als Abbildung $P: \underbrace{C[a, b]}_{\text{Menge der stetigen Fkt}} \rightarrow \Pi_n$ mit $a \leq t_i \leq b$ besitzt die klassische Polynominterpolation

$$\text{die Kondition } k_{abs} = \max_{t \in [a, b]} \sum_{i=0}^n |L_i(t)|.$$

Diese sogenannte Lebesgue-Konstante hängt von der Lage der Stützstellen t_0, \dots, t_n in $[a, b]$ ab und wird besonders klein für die Tschebyscheff-Knoten.

Beispiel 7.4 : Lagrange-Polynome

geg.: $(0,1), (1,3), (3,2)$



$$L_0(t) = \frac{(t-1)(t-3)}{(-1)(-3)}, \quad L_1(t) = \frac{t \cdot (t-3)}{1 \cdot (-2)}, \quad L_2(t) = \frac{t \cdot (t-1)}{3 \cdot 2}$$

$$P(t) = L_0(t) + 3 \cdot L_1(t) + 2 \cdot L_2(t)$$

Bemerkung 7.5 : Algorithmus von Aitken-Neville

geg.: (t_i, f_i) , $(i = 0, 1, \dots, n)$

ges.: $P(t)$ effizient

Bezeichnungen :

- Zu $k + 1 \leq n + 1$ seien paarweise Indizes i_j definiert mit $j = 0, 1, \dots, k$ und $i_j \in \{0, 1, \dots, n\}$.
- $P_{i_0}, P_{i_1}, \dots, P_{i_k} \in \Pi_k$ seien die Interpolationspolynome mit $P_{i_0}, \dots, P_{i_k}(t_i) = f_i$ für $i = i_0, \dots, i_k$

Eigenschaften :

a) $P_i(t) = f_i$

b)
$$P_{i_0}, \dots, P_{i_k}(t) = \frac{(t - t_{i_0})P_{i_1}, \dots, P_{i_k}(t) - (t - t_{i_k})P_{i_0}, \dots, P_{i_{k-1}}(t)}{t_{i_k} - t_{i_0}}$$
 ,

denn beide Ausdrücke nehmen für $t = t_{i_j}$ den Wert f_{i_j} an und sind Polynome aus Π_k .

Rechenschema : für n=3

	$k = 0$	$k = 1$	$k = 2$	$k = 3$
t_0	$f_0 \equiv P_0(t)$	$P_{01}(t) = \dots$	$P_{012}(t) = \dots$	$P_{0123}(t) = P(t)$
t_1	$f_1 \equiv P_1(t)$			
t_2	$f_2 \equiv P_2(t)$			
t_3	$f_3 \equiv P_3(t)$			

Beispiel : Stützpunkte $(0,1)$, $(1,3)$, $(3,2)$

	$k = 0$	$k = 1$	$k = 2$
0	1	$P_{01}(2) = \frac{(2-0) \cdot 3 - (2-1) \cdot 1}{1-0} = 5$	$P_{012}(2) = \frac{(2-0) \cdot \frac{5}{2} - (2-3) \cdot 5}{3-0} = \frac{10}{3}$
1	3		
2	2		

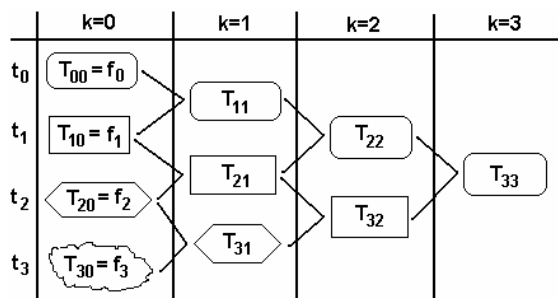
kurz :

$$P_{i,i+1,\dots,i+k} =: T_{i+k,k} \quad \Rightarrow T_{i,0} = f_i$$

$$P_{i-1,\dots,i} =: T_{i,k} \quad T_{i,k} = \frac{(t-t_{i-k})T_{i,k-1} - (t-t_i)T_{i-1,k-1}}{t_i - t_{i-k}}$$

$$T_{i,k} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\frac{t-t_{i-k}}{t-t_i} - 1} \quad , 1 \leq k \leq i \quad , i \geq 0$$

Speicherschema für k=3 :



(Gleiche Formen belegen gleichen Speicherplatz.)

Implementierung :

Bezeichnung : $T_j := T_{i,k}$ mit $j = i - k$

for $i = 0 : n$

$T_i := f_i$

for $j = i - 1 : -1 : 0$

$$T_j := T_{j+1} + \frac{T_{j+1} - T_j}{\left(\frac{t-t_j}{t-t_i} - 1\right)}$$

endfor

endfor

Bemerkung 7.6 : Dividierte Differenzen

Idee : $P_{i_0,i_1,\dots,i_{k-1}}$ und P_{i_0,i_1,\dots,i_k} haben jeweils denselben Funktionswert in $t_i, t_{i_1}, \dots, t_{i_{k-1}}$.

$\Rightarrow P_{i_0,i_1,\dots,i_k}(t) = P_{i_0,i_1,\dots,i_{k-1}}(t) + (t-t_{i_0})(t-t_{i_1})\dots(t-t_{i_{k-1}})f_{i_0,i_1,\dots,i_k}$ mit einem $f_{i_0,i_1,\dots,i_k} \in \mathbb{R}$.

$$f_{i_0,i_1,\dots,i_k} = \frac{f_{i_1,i_1,\dots,i_k} - f_{i_0,i_1,\dots,i_{k-1}}}{t_{i_k} - t_{i_0}}$$

Eigenschaften :

f_{i_0,i_1,\dots,i_k} ist invariant gegenüber Permutationen der Indizes i_0, i_1, \dots, i_k , da $P_{i_0,i_1,\dots,i_k}(t)$ durch die Interpolationsbedingungen $P(t_{i_j}) = f_{i_j}$, $(j = 0, 1, \dots, k)$

eindeutig und von der Reihenfolge der t_j unabhängig festgelegt ist.

Steigungsschema :

$$f_{i,i+1,\dots,i+k} = \frac{f_{i+1,\dots,i+k} - f_{i,\dots,i+k-1}}{t_{i+k} - t_i}$$

	$k = 0$	$k = 1$	$k = 2$
t_0	f_0	$f_{01} = \frac{f_1 - f_0}{t_1 - t_0}$	$f_{012} = \frac{f_{12} - f_{01}}{t_2 - t_0}$
t_1	f_1		
t_2	f_2	$f_{12} = \frac{f_2 - f_1}{t_2 - t_1}$	

Beispiel : Stützpunkte $(0,1)$, $(1,3)$, $(3,2)$

0	1	$\frac{3-1}{1-0} = 2$	$\frac{-\frac{1}{2}-2}{3-0} = -\frac{5}{6}$
1	3		
3	2	$\frac{2-3}{3-1} = -\frac{1}{2}$	

Implementierung :

```

for i = n : -1 : 0
    a_i = f_i
    for j = i + 1 : n
        a_j = (a_j - a_{j-1}) / (t_j - t_i)
    endfor
endfor

```

Bemerkung 7.7 : Newtonsche Darstellung des Interpolationspolynoms

Aus Bemerkung 7.6 erbit sich rekursiv die Newtonsche Darstellung des Interpolationspolynoms :

$$\begin{aligned}
 P(t) = P_{0,1,\dots,n}(t) &= P_{0,1,\dots,n-1}(t) + (t-t_0)(t-t_1)\dots(t-t_{n-1}) \cdot f_{0,1,\dots,n} \\
 &= \dots
 \end{aligned}$$

$$\begin{aligned}
&= f_0 + (t-t_0)f_{0,1} + (t-t_0)(t-t_1)f_{0,1,2} + \dots + (t-t_0)(t-t_1)\dots(t-t_{n-1})f_{0,1,\dots,n} \\
&= f_0 + (t-t_0)[f_{0,1} + (t-t_1)f_{0,1,2} + \dots + (t-t_{n-2})(f_{0,1,\dots,n-1} + (t-t_{n-1})f_{0,1,\dots,n})]
\end{aligned}$$

Beispiel : Stützpunkte $(0,1)$, $(1,3)$, $(3,2)$

$$\begin{aligned}
P(t) &= 1 + (t-0) \left(2 + (t-1) \cdot \left(-\frac{5}{6} \right) \right) \\
P(2) &= 1 + (2-0) \left(2 + (2-1) \left(-\frac{5}{6} \right) \right) = 1 + 2 \left(2 - \frac{5}{6} \right) = 1 + \frac{7}{3} = \frac{10}{3}
\end{aligned}$$

Implementierung : Newton-Horner-Schema

```

P := a_n
for i : n-1 : -1 : 0
    P := P · (t - t_i) + a_i
endfor

```

Satz 7.8 : Restglied der Polynominterpolation

Gegeben seien $f : C^{n+1}[a,b] \rightarrow R$ und das Interpolationspolynom $P(t)$ zu den Stützstellen $(t_i, f(t_i))$, $i = 0,1,\dots,n$ mit $n+1$ paarweise voneinander verschiedenen Stützstellen $t_i \in [a,b]$, $i = 0,1,\dots,n$

Zu jedem $\bar{t} \in [a,b]$ gibt es ein $\mathbf{t} \in [a,b]$ mit $f(\bar{t}) - P(\bar{t}) = \frac{\mathbf{w}(\bar{t})f^{(n+1)}(\mathbf{t})}{(n+1)!} = \text{Restglied}$,
mit $\mathbf{w}(t) := (t-t_0)(t-t_1)\dots(t-t_n)$ und $\mathbf{t} = \mathbf{t}(\bar{t})$.

Folgerung 7.9 : Eigenschaften der dividierten Differenzen

- a) Unter den Voraussetzungen von Satz 7.8 gilt für jedes $\bar{t} \in [a,b]$: $f(\bar{t}) - P(\bar{t}) = \mathbf{w}(\bar{t}) \cdot f[t_0, t_1, \dots, t_n, \bar{t}]$, wobei $f[t_0, t_1, \dots, t_n, \bar{t}]$ die dividierte Differenz $f_{0,1,\dots,n+1}$ zu den $n+2$ Stützpunkten $(t_i, f(t_i)), (\bar{t}, f(\bar{t}))$ bezeichnet.
- b) Es gilt $f[t_0, t_1, \dots, t_n] = \frac{f^{(n)}(\mathbf{t})}{n!}$ für ein $\mathbf{t} \in [a,b]$.
- c) Für $f \in \Pi_n$ ist $f[t_0, t_1, \dots, t_n, t_{n+1}, \dots, t_{n+k}] \equiv 0$, $k > 0$

Beispiel 7.10 : Restglied der Polynominterpolation

$$\begin{aligned}
f(t) &= \sin t, \quad t_i = i \cdot \frac{\pi}{10}, \quad i = 0,1,\dots,5, \quad n = 5 \\
\sin t - P(t) &= \frac{1}{720} (t-t_0)(t-t_1)\dots(t-t_5)(-\sin t)
\end{aligned}$$

$$\bar{t} \in \left[0, \frac{p}{2}\right]: \quad |\sin \bar{t} - P(\bar{t})| \leq \frac{1}{720} \left(\frac{p}{2}\right)^2 \left(\frac{2p}{5}\right)^2 \left(\frac{3p}{10}\right)^2 = \frac{p^6}{2 \cdot 10^5} \approx 0.005$$

beachte: Ausserhalb des Intervalls $\left(\min_i t_i, \max_i t_i\right)$ („Extrapolation“)
häufig sehr große Fehler $(|w(t)| \gg 1)$.

Betrachte eine Folge (Δ_m) von Intervallteilungen mit $\Delta_m = \{a = t_0^{(m)} < t_1^{(m)} < \dots < t_{n_m}^{(m)} = b\}$ mit zugehörigen Interpolationspolynomen $P(t; \Delta_m)$.

Satz von Faber :

Zu jeder solchen Folge existiert ein $f \in C[a, b]$ so, dass $P(t; \Delta_m)$ auf $[a, b]$ nicht gleichmäßig konvergiert.
(D.h. es gibt keine optimale Folge von Stützstellen.)

Beispiel : Funktion von Runge

$$f(t) = \frac{1}{1+t^2}, \quad [a, b] = [-5, 5]$$

$$t_i^{(m)} = -5 + i \cdot \frac{10}{m}, \quad i = 0, 1, \dots, m$$

Bemerkung 7.11 : Hermite-Interpolation

$$\text{Interpolationsbedingungen } P^{(j)}(t_i) = f^{(j)}(t_i), \quad \begin{cases} i = 0, 1, \dots, n \\ j = 0, 1, \dots, c_i \end{cases}$$

klassische Hermite-Interpolation :

$$P \in \Pi_r \text{ mit } r = \sum_{i=0}^n c_i + n$$

Lagrange-Darstellung :

Bestimme $L_{ij}(t) \in \Pi_r$ mit

$$L_{ij}^{(l)}(t_k) = \begin{cases} 1 & \text{falls } i = k \text{ und } j = l \\ 0 & \text{sonst} \end{cases}$$

\Rightarrow

$$P(t) = \sum_{i=0}^n \sum_{j=0}^{c_i} f^{(j)}(t_i) L_{ij}(t)$$

Newtonsche Darstellung :

Knoten mehrfach angeben, Steigung $f[t_i, t_i]$ ersetzen durch vorgegebene Daten $f'(t_i)$,

$f[t_i, t_i, t_i]$ ersetzen durch vorgegebene Daten $\frac{1}{2} f''(t_i)$, ... usw.

Beispiel 1:

Interpolationsbedingungen $P(t_i) = f_i$, $P'(t_i) = f_i'$, $P''(t_i) = f_i''$, $i = 0,1$
 Steigungsschema :

t_0	f_0	$f_{00} = f_0'$	$f_{000} = \frac{1}{2} f_0''$	$f_{0001} = \frac{f_{001} - f_{000}}{t_1 - t_0}$... usw.
t_0	f_0	$f_{00} = f_0'$		
t_0	f_0	$f_{01} = \frac{f_1 - f_0}{t_1 - t_0}$	$f_{001} = \frac{f_{01} - f_{00}}{t_1 - t_0}$	
t_1	f_1		$f_{011} = \frac{f_{11} - f_{01}}{t_1 - t_0}$	
t_1	f_1	$f_{11} = f_1'$		
t_1	f_1	$f_{11} = f_1'$	$f_{111} = \frac{1}{2} f_1''$	

$$P(t) = f_0 + (t-t_0) \cdot f_0' + (t-t_0)^2 \cdot \frac{1}{2} f_0'' + (t-t_0)^3 \cdot f_{0001} + (t-t_0)^3 (t-t_1) \cdot f_{00011} + (t-t_0)^3 (t-t_1)^2 \cdot f_{000111}$$

Beispiel 2: Taylorpolynom

Interpolationsbedingungen $P^{(k)}(t_0) = f^{(k)}(t_0)$, $k = 0,1,\dots,c_0$

$$P(t) = \sum_{k=0}^{c_0} \frac{1}{k!} f^{(k)}(t_0) (t-t_0)^k$$

Bemerkung 7.12: stückweise Hermite-Interpolation

Problem : Polynome hoher Ordnung neigen zu Oszillation

Kompromiss : Approximierende Funktionen, die nicht unendlich oft stetig differenzierbar sind

Intervallgitter : $\Delta := \{t_i : a = t_0 < t_1 < \dots < t_m = b\}$ auf $[a, b]$

Approximation von $f : [a, b] \rightarrow R$ durch ein \mathbf{j} aus dem Hermiteschen Funktionenraum

$$H_{\Delta}^{(n)} := \{ \mathbf{j} \in C^{n-1}[a, b] \mid \mathbf{j}|_{[t_i, t_{i+1}]} = \mathbf{p}_i|_{[t_i, t_{i+1}]} \text{ („eingeschränkt auf“)} \}$$

für ein $\mathbf{p}_i \in \Pi_{2n-1}$, $i = 0,1,\dots,m-1$.

praktisch:

Bestimmung der Hermite-Polynome $P_i(t)$, $i = 0,1,\dots,m-1$

mit $P_i^{(k)}(t_i) = f^{(k)}(t_i)$, $P_i^{(k)}(t_{i+1}) = f^{(k)}(t_{i+1})$, $k = 0,1,\dots,n-1$

$$j|_{[t_i, t_{i+1}]} := P_i|_{[t_i, t_{i+1}]}, \quad P_i \in \Pi_{2n-1}$$

Auswertung von j :

Bestimme i mit $t_i \leq t \leq t_{i+1}$

Spezialfall : äquidistantes Gitter $t_i := a + \frac{i(b-a)}{m}, \quad i = 0, 1, \dots, m$

$$i := \left\lfloor \frac{t-a}{h} \right\rfloor \quad \text{mit} \quad h := \frac{b-a}{m}$$

allgemein : binäre Suche

$$\underline{i} := 0, \quad \bar{i} := m;$$

do
{

$$i_* := \left\lfloor \frac{\underline{i} + \bar{i}}{2} \right\rfloor$$

if $t \geq t_{i_*}$ **then** $\underline{i} := i_*$

else $\bar{i} := i_*$;

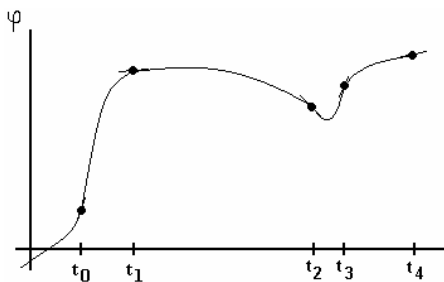
}

while $(\bar{i} - \underline{i} > 1)$;

$i := \underline{i}$;

Beispiel :

$$n = 2, \quad P_i \in \Pi_3, \quad j \in C^1[a, b]$$



(Polynome hoher Ordnung neigen zum Schwingen zwischen den Stützstellen.)

Restglied , Fehlerabschätzungen :

Approximation der Ableitungen von f :

$$\left| f^{(k)}(t) - P_i^{(k)}(t) \right| \leq \frac{|(t-t_i)(t-t_{i+1})|^{m-k}}{k!(2n-2k)!} \cdot (t_{i+1} - t_i)^k \cdot \max_{t \in [t_i, t_{i+1}]} |f^{(2n)}(t)|$$

$$\leq \frac{\|\Delta\|^{2n-k}}{2^{2n-2k} \cdot k!(2n-2k)!} \cdot \max_{t \in [a,b]} |f^{(2n)}(t)|$$

7.2 Spline-Interpolation

Definition 7.13 : Polynomialer Spline

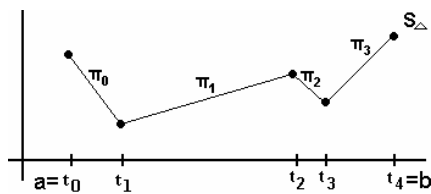
Sei $\Delta = \{t_0, \dots, t_{l+1}\}$ ein Gitter von $l+2$ paarweise voneinander verschiedenen Knoten t_i mit $a = t_0 < t_1 < \dots < t_{l+1} = b$.

Ein polynomialer Spline vom Grad $k-1$ (Ordnung k) bzgl. Δ ist eine Funktion $s_\Delta \in C^{k-2}[a,b]$, die auf jedem Teilintervall $[t_i, t_{i+1}]$ mit einem $p_i \in \Pi_{k-1}$ übereinstimmt, ($i = 0, 1, \dots, l$).

Bezeichnung: $S_{k,\Delta} :=$ Menge der Splines vom Grad $k-1$ bzgl. Δ

Beispiel 7.14 :

lineare Splines : Menge der stetigen, stückweise linearen Funktionen ($S_{2,\Delta}$)



kubische Splines : stückweise kubische Funktion, zweimal stetig differenzierbar ($S_{4,\Delta}$)

Bemerkung 7.15 : Eigenschaften polynomialer Splines

a) $S_{k,\Delta}$ ist ein Vektorraum

b) $\Pi_{k-1} \subset S_{k,\Delta}$, $((t-t_i)^{k-1})_+ \in S_{k,\Delta}$ mit $((t-t_i)^{k-1})_+ := \begin{cases} (t-t_i)^{k-1} & \text{falls } t \geq t_i \\ 0 & \text{sonst} \end{cases}$

c) Bedingungen zur eindeutigen Bestimmung eines kubischen Splines :

Interpolationsbedingungen : $s_\Delta(t_i) = f_i$, ($i = 0, 1, \dots, l+1$)

i) $s_\Delta''(a) = s_\Delta''(b) = 0$ (natürlicher kubischer Spline)

ii) $s_\Delta^{(k)}(a) = s_\Delta^{(k)}(b)$, $k = 0, 1, 2$ (periodischer kubischer Spline – nur sinnvoll, falls $f_0 = f_{l+1}$)

iii) $s_\Delta'(a) = f_a'$, $s_\Delta'(b) = f_b'$ (eingespannter kubischer Spline / vollständiger kub. Spline / kubischer Hermite-Spline)

kubischer Spline :

4 Parameter je Teilintervall

$\Rightarrow 4 \cdot (l+1)$ Koeffizienten

Interpolationsbedingungen = $l+2$

aus C^2 folgen $3l$ Bedingungen

$\Rightarrow \sum = 4(l+1) - 2$

Daher : 2 Bedingungen aus den obigen i)-iii) auswählen

d) Die Bedingungen i)-iii) führen jeweils zu kubischen Splines s_Δ , welche die

Halbnorm $\|s_\Delta\| := \left(\int_a^b (s''(t))^2 dt \right)^{\frac{1}{2}}$ über einer gewissen Klasse von Funktionen minimieren.
(Krümmung wird minimiert.)

Bezeichnung: $K^r[a, b] := \{f : f^{(k)} \text{ auf } [a, b] \text{ absolutstetig, } (k = 0, 1, \dots, r-1), f^{(r)} \in L^2[a, b]\}$

$$v \in L^2[a, b] \Leftrightarrow \int_a^b (v(t))^2 dt \text{ existiert} \quad (\text{siehe Lebesgue/Maßtheorie})$$

Definition Absolutstetigkeit:

$f : [a, b] \rightarrow \mathbb{R}$ heißt absolutstetig, falls es zu jedem $\epsilon > 0$ ein $\delta > 0$ gibt, so dass für beliebige a_i, b_i mit $a \leq a_1 < b_1 < a_2 < b_2 < \dots < b_n \leq b$ und $\sum_i |b_i - a_i| < \delta$ gilt:

$$\sum_i |f(b_i) - f(a_i)| < \epsilon.$$

Absolutstetige Funktionen sind fast überall stetig differenzierbar und es gilt:

$$f(t) - f(a) = \int_a^t f'(t) dt \quad \text{mit } t \in [a, b] \quad \text{sowie} \quad \int_a^b f(t) g'(t) dt = [f(t) g(t)]_a^b - \int_a^b f'(t) g(t) dt.$$

Lemma 7.16 : Identität von Holladay

Ist $f \in K^2[a, b]$, $\Delta = \{a = t_0 < t_1 < \dots < t_{l+1} = b\}$ und $s_\Delta \in S_{4, \Delta}$, so gilt:

$$\|f - s_\Delta\|^2 = \|f\|^2 - \|s_\Delta\|^2 - 2 \left([(f'(t) - s_\Delta'(t)) \cdot s_\Delta''(t)]_a^b - \sum_{i=0}^{l+1} [(f(t) - s_\Delta(t)) \cdot s_\Delta'''(t)]_{t_{i-1+0}}^{t_i-0} \right)$$

Satz 7.17 : Splines und Interpolierende minimaler Krümmung

Zu $\Delta = \{a = t_0 < t_1 < \dots < t_{l+1} = b\}$ seien Stützwerte $\{y_0, y_1, \dots, y_{l+1}\}$ und ein interpolierender Spline $s_\Delta \in S_{4, \Delta}$ gegeben.

- i) Ist $s_\Delta''(a) = s_\Delta''(b) = 0$, so gilt für jedes $f \in K^2[a, b]$ mit $f(t_i) = y_i$: $\|f\| \geq \|s_\Delta\|$.
- ii) Ist $s_\Delta^{(k)}(a) = s_\Delta^{(k)}(b)$, $(k = 0, 1, 2)$, so gilt für jedes $f \in K^2[a, b]$ mit $f(t_i) = y_i$: $\|f\| \geq \|s_\Delta\|$.
- iii) Ist $f \in K^2[a, b]$, $f(t_i) = y_i$ und $f'(a) = s_\Delta'(a)$, $f'(b) = s_\Delta'(b)$, so gilt: $\|f\| \geq \|s_\Delta\|$.

In jedem der drei Fälle ist s_{Δ} eindeutig bestimmt.

Bemerkung 7.18 : Bestimmung der Splinekoeffizienten (für kubische Splines)

$$S_{4,\Delta} : s_{\Delta}|_{[t_i, t_{i+1}]} = a_i + b_i(t-t_i) + \frac{1}{2}c_i(t-t_i)^2 + \frac{1}{6}d_i(t-t_i)^3 \quad \text{mit}$$

$$a_i = s_{\Delta}(t_i) \quad (\text{durch Interpolationsbedingungen vorgegeben})$$

$$b_i = s_{\Delta}'(t_i)$$

$$c_i = s_{\Delta}''(t_i)$$

$$d_i = s_{\Delta}'''(t_i)|_{t \rightarrow t_i+0}$$

$$\text{Aus } s_{\Delta}''(t)|_{[t_i, t_{i+1}]} \text{ ist linear folgt : } s_{\Delta}''(t) = c_{i+1} \frac{t-t_i}{h_i} - c_i \frac{t-t_{i+1}}{h_i} \quad \text{mit } h_i := t_{i+1} - t_i$$

\Rightarrow

$$s_{\Delta}(t) = c_{i+1} \frac{(t-t_i)^3}{6h_i} - c_i \frac{(t-t_{i+1})^3}{6h_i} + B_i(t-t_i) + A_i \quad \left(\iint \text{ von } s_{\Delta}'' \right)$$

Interpolationsbedingungen :

$$s_{\Delta}(t_i) = c_i \frac{h_i^2}{6} + A_i \stackrel{!}{=} f_i$$

$$s_{\Delta}(t_{i+1}) = c_{i+1} \frac{h_i^2}{6} + B_i h_i + A_i \stackrel{!}{=} f_{i+1}$$

\Rightarrow

$$A_i = f_i - \frac{h_i^2}{6} c_i$$

$$B_i = \frac{f_{i+1} - f_i}{h_i} - \frac{h_i}{6} (c_{i+1} - c_i)$$

Ausmultiplizieren von $(t-t_{i+1})^3 = \left((t-t_i) - \underbrace{(t_{i+1}-t_i)}_{=h_i} \right)^3$ und Koeffizientenvergleich ergibt :

$$a_i = f_i$$

$$b_i = \frac{f_{i+1} - f_i}{h_i} - \frac{h_i}{6} (2c_i + c_{i+1})$$

$$d_i = \frac{c_{i+1} - c_i}{h_i}$$

Bestimmung der c_i : aus den Stetigkeitsbedingungen für $s_{\Delta}'(t)$

$$s_{\Delta}'(t)|_{t_i+0} = b_i = \frac{f_{i+1} - f_i}{h_i} - \frac{h_i}{6} (2c_i + c_{i+1})$$

$$s_{\Delta}'(t)|_{t=t_0} = \left(b_{i-1} + h_{i-1}c_{i-1} + \frac{h_{i-1}^2}{2}d_{i-1} \right) = \frac{f_i - f_{i-1}}{h_{i-1}} - \frac{h_{i-1}}{6}(2c_{i-1} + c_i) + h_{i-1}c_{i-1} + \frac{h_{i-1}}{2}(c_i - c_{i-1})$$

$$= \frac{f_i - f_{i-1}}{h_{i-1}} + \frac{h_{i-1}}{6}c_{i-1} + \frac{h_i}{3}c_i$$

⇒

$$\frac{h_{i-1}}{6}c_{i-1} + \frac{h_{i-1} + h_i}{3}c_i + \frac{h_{i+1}}{6}c_{i+1} = \frac{f_{i+1} - f_i}{h_i} - \frac{f_i - f_{i-1}}{h_{i-1}}, \quad (i = 1, \dots, l)$$

mit $c_{l+1} := s_{\Delta}''(b)$

natürlicher kubischer Spline : $c_0 = c_{l+1} = 0$

$$\begin{pmatrix} \frac{h_0 + h_1}{3} & \frac{h_1}{6} & & & \\ \frac{h_1}{6} & \frac{h_1 + h_2}{3} & \frac{h_2}{6} & & \\ & & \ddots & \ddots & \\ & & & \frac{h_{l-1}}{6} & \frac{h_{l-1} + h_l}{3} \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \\ \vdots \\ c_l \end{pmatrix} = \begin{pmatrix} \frac{f_2 - f_1}{h_1} - \frac{f_1 - f_0}{h_0} \\ \frac{f_3 - f_2}{h_2} - \frac{f_2 - f_1}{h_1} \\ \vdots \\ \frac{f_{l+1} - f_l}{h_l} - \frac{f_l - f_{l-1}}{h_{l-1}} \end{pmatrix}$$

=MATLAB-Befehl "diff(f) "

Lösung des Gleichungssystems mittels Cholesky-Zerlegung für tridiagonale Matrizen.

- Fragen : - Eindeutige Lösbarkeit ?
 - Koeffizientenmatrix positiv definit ?

Lemma 7.19 : Ein Kriterium für die Regularität von Matrizen

Gilt für eine Matrix $A \in R^{n \times n}$ mit $a_{ii} = 2$ und $\sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \leq 1$, $(i = 1, \dots, n)$,

so ist A regulär und $\|A^{-1}\|_{\infty} \leq 1$.

Beweis :

a) Für Vektoren $w, z \in R^n$ mit $w = Az$ bestimmt man $i_0 \in \{1, \dots, n\}$ so,

$$\text{dass } |z_{i_0}| = \max_{1 \leq i \leq n} |z_i| = \|z\|_{\infty}.$$

$$\text{Dann ist } \|w\|_{\infty} \geq |w_{i_0}| = \left| \sum_{j=1}^n a_{i_0,j} \cdot z_j \right| \geq |a_{i_0,i_0} \cdot z_{i_0}| - \left| \sum_{\substack{j=1 \\ j \neq i_0}}^n a_{i_0,j} \cdot z_j \right|$$

$$\geq 2|z_{i_0}| - \sum_{\substack{j=1 \\ j \neq i_0}}^n \left(|a_{i_0,j}| \cdot \max_{1 \leq j \leq n} |z_j| \right)$$

$$\geq (2-1)\|z\|_{\infty} = \|z\|_{\infty} \quad \Rightarrow \text{Eindeutigkeit}$$

Satz 7.22 : Approximationsgüte der zweiten Ableitung

Mit den Bezeichnungen aus Bem. 7.21 gilt für den vollständigen (eingespannten) interpolierten kubischen Spline :

$$\|c_{\Delta} - f_{\Delta}''\|_{\infty} \leq \|r_{\Delta}\|_{\infty} \leq \frac{3}{4} M \cdot \|\Delta\|^2 ,$$

sofern $f \in C^4[a,b]$ und $|f^4(t)| \leq M$ mit $(t \in [a,b])$ ist.

Satz 7.23 : Approximationsgüte vollständiger kubischer Splines

Gegeben sei eine Funktion $f \in C^4[a,b]$ mit $\|f^{(4)}\|_{\infty} \leq M$, sowie ein Gitter $\Delta = \{a = t_0 < t_1 < \dots < t_{l+1} = b\}$

mit einer Konstanten $K \geq \frac{\max_{0 \leq i \leq l} h_i}{\min_{0 \leq i \leq l} h_i}$.

Dann gibt es zum vollständigen interpolierenden kubischen Spline s_{Δ} Konstanten C_0, C_1, C_2, C_3 die von Δ und K unabhängig sind , und für die gilt :

$$|f^{(k)}(t) - s_{\Delta}^{(k)}(t)| \leq C_k \cdot M \cdot K \cdot \|\Delta\|^{4-k} , (k = 0,1,2,3)$$

mit $t \in [a,b]$, sofern $s_{\Delta}^{(k)}$ in t definiert ist.

Folgerung 7.24 : Gleichmäßige Konvergenz interpolierender Splines

Zu einer Gitterfolge $(\Delta_m)_{m \geq 0}$ mit $\|\Delta_m\| \rightarrow 0$ und $\sup_{m,j} \frac{\|\Delta_m\|}{h_j^{(m)}} \leq K < \infty$ konvergieren die vollständigen

interpolierenden kubischen Splines $(s_{\Delta_m})_{m \geq 0}$ gleichmäßig auf $[a,b]$ gegen $f \in C^4[a,b]$, falls $\|\Delta_m\| \rightarrow 0$.

Ebenso konvergieren die Ableitungen von s_{Δ_m} gleichmäßig gegen die entsprechenden Ableitungen von f .

Bemerkung 7.25 : B-Splines

Statt der expliziten Berechnung der Splinekoeffizienten ist häufig eine Darstellung $s_{\Delta}(t) = \sum_{i=0}^{l+k-1} \mathbf{a}_i B_i(t)$ mit einer

geeigneten Basis $\{B_0, B_1, \dots, B_{l+k-2}, B_{l+k-1}\}$ von $s_{k,\Delta}$ mit Koeffizienten $\mathbf{a}_i \in R$ vorteilhaft.

Beispiel :

- approximierender Spline $s_{\Delta}(t_i) \approx f_i$, $(i = 0,1,\dots,m)$ mit $m \gg l$

$$\sum_{i=0}^m (s_{\Delta}(t_i) - f_i)^2 \rightarrow \min$$

$\Rightarrow \mathbf{a}_i$ als Lösung eines überbestimmten linearen Gleichungssystems. (vgl. Bem. 4.8)

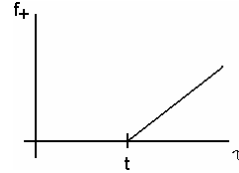
- Spline s_{Δ} mit zusätzlichen Forderungen wie Nichtnegativität , Monotonie , Konvexität , ...

Definition : B-Splines

Zu einer Knotenfolge $\dots, \mathbf{t}_{i-1}, \mathbf{t}_i, \mathbf{t}_{i+1}, \dots, \mathbf{t}_{i+r}, \mathbf{t}_{i+r+1}$ mit $\mathbf{t}_j \leq \mathbf{t}_{j+1}$, $(j \in \mathbb{Z})$ und $\mathbf{t}_i < \mathbf{t}_{i+r}$

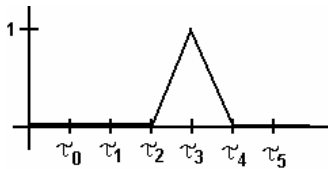
heisst $B_{i,r}(t) := (\mathbf{t}_{i+r} - \mathbf{t}_i) f_t^{r-1}[\mathbf{t}_i, \mathbf{t}_{i+1}, \dots, \mathbf{t}_{i+r}]$

mit $f_t^{r-1}(\mathbf{t}) = (f_t(\mathbf{t}))^{r-1}$, $f_t(\mathbf{t}) := (\mathbf{t} - t)_+ := \begin{cases} \mathbf{t} - t & \text{für } \mathbf{t} > t \\ 0 & \text{sonst} \end{cases}$



i-ter Spline der Ordnung k .

Beispiel :



$B_{3,2}(t)$ - „Hütchenfunktion“

$r = 1$:

$$f_t^{r-1}(\mathbf{t}) = (f_t(\mathbf{t}))^0 = \begin{cases} 1 & \text{für } \mathbf{t} > t \\ 0 & \text{sonst} \end{cases}$$

$$f_t^{r-1}[\mathbf{t}_i, \mathbf{t}_{i+1}] = \begin{cases} \frac{1-1}{\mathbf{t}_{i+1} - \mathbf{t}_i} = 0 & \text{für } t < \mathbf{t}_i \\ \frac{\mathbf{t}_{i+1} - \mathbf{t}_i}{1-0} & \text{für } \mathbf{t}_i \leq t \leq \mathbf{t}_{i+1} \\ \frac{0-0}{\mathbf{t}_{i+1} - \mathbf{t}_i} = 0 & \text{für } t > \mathbf{t}_{i+1} \end{cases}$$

Darstellung : $s_\Delta(t) = \sum_i \mathbf{a}_i B_i(t)$

Rekursive Berechnung :

Unter Verwendung der Produktregel für dividierte Differenzen Newtonsche Darstellung des Interpolationspolynoms von $f_t^{r-1}(\mathbf{t})$ zu Stützstellen $\mathbf{t}_i, \mathbf{t}_{i+1}, \dots, \mathbf{t}_{i+r}$:

$$\mathbf{p}_{r-1}(\mathbf{t}) = \sum_{s=i}^{i+r} \left(f_t^{r-1}[\mathbf{t}_i, \dots, \mathbf{t}_s] \cdot (\mathbf{t} - \mathbf{t}_i) \cdot \dots \cdot (\mathbf{t} - \mathbf{t}_{s-1}) \right) \in \Pi_r$$

Für $t \in [\mathbf{t}_i, \dots, \mathbf{t}_{i+r}]$ kann $f_t^{r-1}[\mathbf{t}_i, \dots, \mathbf{t}_{i+r}]$ numerisch stabil als konvexe Linearkombination von $f_t^{r-2}[\mathbf{t}_i, \dots, \mathbf{t}_{i+r-1}]$ und $f_t^{r-2}[\mathbf{t}_{i+1}, \dots, \mathbf{t}_{i+r}]$ berechnet werden.

Für $t < t_i$ und für $t_{i+r} < t$ fällt $f_t^{r-1}(t)$ in $[t_i, t_{i+r}]$ mit einem Polynom vom Grad $\leq r-1$ zusammen $\Rightarrow f_t^{r-1}[t_i, \dots, t_{i+r}] = 0$.

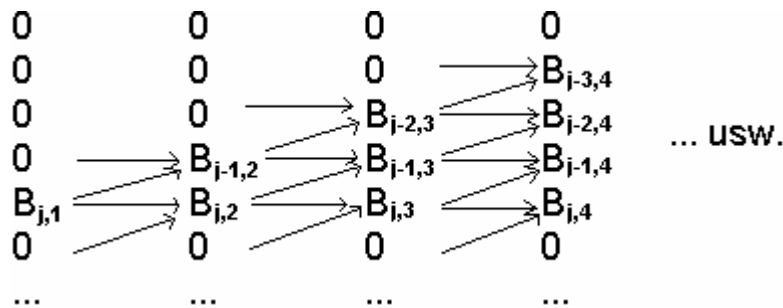
Ist $t_i < t_{i+r-1}$ und $t_{i+1} < t_{i+r}$, so folgt:

$$B_{i,r}(t) = \frac{t-t_i}{t_{i+r-1}-t_i} B_{i,r-1}(t) + \frac{t_{i+1}-t}{t_{i+r}-t_{i+1}} B_{i+1,r-1}(t). \quad (*)$$

Zu gegebenem $t \in R$ gibt es höchstens r B-Splines $B_{i,r}(t) \neq 0$.

Sei j so gewählt, dass $t \in [t_j, t_{j+1}] \Rightarrow B_{i,r}(t) = 0$, für $i \leq j-r$ oder $i \geq j+1$.

Schema:



Beispiel:

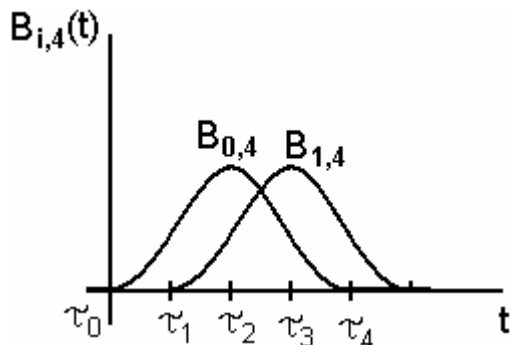
$$t_i = i, \quad (i = 0, 1, 2, \dots), \quad t = 3.5 \in [t_3, t_4], \quad j = 3$$

(siehe *)

$i \setminus r$	1	2	3	4
0	0	0	0	1/48
1	0	0	1/8	23/48
2	0	1/2	3/4	23/48
3	1	1/2	1/8	1/48
4	0	0	0	0

$$\text{z.B.: } B_{2,3}(3.5) = \frac{3.5-2}{4-2} B_{2,2}(3.5) + \frac{5-3.5}{5-3} B_{3,2}(3.5) = \frac{1.5}{2} \cdot \frac{1}{2} + \frac{1.5}{2} \cdot \frac{1}{2} = \frac{3}{8} + \frac{3}{8} = \frac{3}{4}$$

$$r = k = 4: \quad s_\Delta(3.5) = \sum_i a_i B_{i,4}(3.5) = \frac{a_0 + 23a_1 + 23a_2 + a_3}{48}$$



Eigenschaften :

- a) Nichtnegativität $B_{i,r}(t) \geq 0, t \in R$
- b) Lokaler Träger $\text{supp } B_{i,r} \subset [t_i, t_{i+r}]$
- c) $\sum_i B_{i,r}(t) \equiv 1, t \in R$
- d) Zum Gitter $\Delta = \{t_i : a = t_0 < t_1 < \dots < t_{l+1} = b\}$ betrachte die Knotenfolge
 - $t_{-k+1} = t_{-k+2} = \dots = t_{-1} = t_0 := t_0 = a$
 - $t_j = t_j, (j = 1, \dots, l)$
 - $t_{l+1} = t_{l+2} = \dots = t_{l+k} := t_{l+1} = b$
 und die B-Splines $B_{-k+1,k}(t), B_{-k+2,k}(t), \dots, B_{l,k}(t)$

Satz : (Curry , Schoenberg)

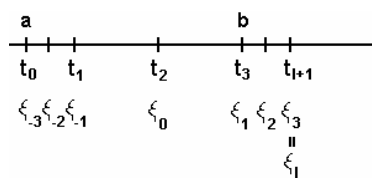
Zu $k = 4$ betrachte $l + k$ Interpolationsbedingungen $f_j = \sum_{i=-k+1}^l a_i B_{i,k}(x_j), (j = -k + 1, \dots, l)$
 \Rightarrow Lineares Gleichungssystem der Dimension $l + k$ mit Bandmatrix der Breite k zur Bestimmung von a_{-k+1}, \dots, a_l .

Elemente der Koeffizientenmatrix : $B_{i,k}(x_j)$ aus obigem Rechenschema ,
 zweckmäßig : $x_{-k+1} < x_{-k+2} < \dots < x_l$.

Satz : (Schoenberg , Whitney)

Der interpolierende Spline ist genau dann eindeutig bestimmt , wenn $B_{i,k}(x_j) \neq 0, (i = -k + 1, \dots, l)$.

praktisch : für $k = 4$



z.B. für gerades k :

$$x_{-k+1} := t_0, \quad x_j := t_{j+k/2}, \quad (j = -k/2 + 1, \dots, -k/2 + l) \quad d.h. t_1, \dots, t_l$$

zusätzliche Interpolationsbedingungen :

$$x_{-k+1+j} \in (t_0, t_1), \quad x_{l-j} \in (t_l, t_{l+1}), \quad (j = 1, \dots, k/2 - 1)$$

7.3 Trigonometrische Interpolation

Problemstellung : Interpolation periodische Daten durch trigonometrische Polynome

$N = 2n + 1$: (ungerade)

$$T_R^N := \left\{ \Phi_{2n+1}(t) := \frac{a_0}{2} + \sum_{j=1}^n (a_j \cos(jt) + b_j \sin(jt)) \right\} \quad , R \text{ bedeutet : } a_j, b_j \in R$$

$N = 2n$: (gerade)

$$T_R^N := \left\{ \Phi_{2n}(t) := \frac{a_0}{2} + \left(\sum_{j=1}^{n-1} (a_j \cos(jt) + b_j \sin(jt)) \right) + \frac{a_n}{2} \cos(nt) \right\}$$

Die komplexe Darstellung kommt ohne Fallunterscheidung aus :

$$T_C^N := \left\{ \Phi_N(t) := \sum_{j=0}^{N-1} c_j \cdot e^{ijt} \quad , c_j \in C \right\}$$

Die Ansatzfunktionen $\{1, \cos(jt), \sin(jt)\}$ bzw. $\{e^{ijt}\}$ sind linear unabhängig.

\Rightarrow Die Interpolationsaufgabe zu N Stützstellen (t_k, f_k) , $(k = 0, 1, \dots, N-1)$ ist eindeutig lösbar.

Typische Aufgabenstellung :

Digitale Datenverarbeitung , Datenerfassung im festen Takt

\Rightarrow äquidistante Knoten , normiert auf $t_k = k \cdot \frac{2p}{N}$

Bemerkung 7.27 : Trigonometrische Interpolation auf äquidistantem Gitter

Zu $t_k = k \cdot \frac{2p}{N}$ sind $w_k := e^{it_k} = e^{ik \frac{2p}{N}}$ die N-ten Einheitswurzeln $((w_k)^N = 1)$

Die Interpolationsbedingungen $\Phi_N(t_k) = f_k$, $(k = 0, \dots, N-1)$ bestimmen Φ_N mit :

$$\begin{pmatrix} 1 & w_0 & w_0^2 & \dots & w_0^{N-1} \\ 1 & w_1 & w_1^2 & \dots & w_1^{N-1} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & w_{N-1} & w_{N-1}^2 & \dots & w_{N-1}^{N-1} \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ c_1 \\ \vdots \\ c_{N-1} \end{pmatrix} = \begin{pmatrix} f_0 \\ f_1 \\ \vdots \\ f_{N-1} \end{pmatrix} \quad , (\det \neq 0)$$

Lemma 7.28 : Komplexe und reelle trigonometrische Polynome

Gilt für das (komplexe) trigonometrische Polynom $\Phi_N(t) = \sum_{j=0}^{N-1} c_j e^{ijt}$ mit $\Phi_n(t) \in R$, $t \in R$, so gilt :

$$\Phi_N \in T_R^N \text{ mit } a_j = 2\text{Re}(c_j) = c_j + c_{N-j} \text{ und } b_j = -2\text{Im}(c_j) = i(c_j - c_{N-j})$$

Lemma 7.29 : Eigenschaften der Basisfunktionen

- a) Für die N-ten Einheitswurzeln $w_j := e^{ij\frac{2p}{N}}$ gilt : $\sum_{j=0}^{N-1} w_j^k w_k^{-l} = N \cdot d_{k,l}$.
- b) Die Basisfunktionen $y_k(t) := e^{ikt}$ der (komplexen) trigonometrischen Interpolation sind orthonormal bezüglich des diskreten Skalarprodukts $\langle f, g \rangle := \frac{1}{N} \sum_{j=0}^{N-1} (f(t_j) \cdot \overline{g(t_j)})$, d.h. $\langle y_k, y_l \rangle = d_{k,l}$.

Satz 7.30 : Lösung der trigonometrischen Interpolationsaufgabe

Für äquidistante Stützstellen $t_k = k \frac{2p}{N}$ ist die Lösung der trigonometrischen Interpolationsaufgabe $\Phi_n(t_k) = f_k$ gegeben durch :

$$\Phi_n(t) = \sum_{j=0}^{N-1} c_j e^{ijt} \quad \text{mit} \quad c_j := \frac{1}{N} \sum_{k=0}^{N-1} f_k w_k^{-j}$$

Beweis :

$$\begin{aligned} \Phi_n(t_l) &= \sum_{j=0}^{N-1} c_j e^{ijt_l} = \sum_{j=0}^{N-1} c_j \underbrace{\left(e^{il\frac{2p}{N}} \right)^j}_{=w_l} && \text{- wegen } t_l = l \frac{2p}{N} \\ &= \sum_{j=0}^{N-1} \left(\frac{1}{N} \sum_{k=0}^{N-1} f_k w_k^{-j} \right) \cdot w_l^j && \text{- } c_j \text{ eingesetzt} \\ &= \sum_{j=0}^{N-1} f_k \cdot \left(\frac{1}{N} \sum_{k=0}^{N-1} w_k^{-j} w_l^j \right) \\ &= \sum_{j=0}^{N-1} f_k \cdot d_{k,l} && \text{- nach Lemma 7.29a) mit } w_k^{-j} = w_j^{-k}, w_l^j = w_j^l \\ &= f_l \quad , (l = 0, 1, \dots, N-1) \end{aligned}$$

Bemerkung 7.31 : Diskrete Fourier-Transformation

Für $2p$ - periodische Funktionen $f \in L^2(R)$ erhält man mit den Fourierkoeffizienten

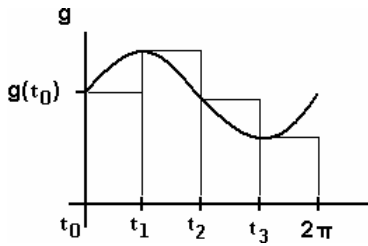
$$\hat{f}(j) := \frac{1}{2p} \int_0^{2p} f(t) e^{-ijt} dt \quad , j \in Z$$

die abgebrochene Fourier-Reihe $f_n(t) := \sum_{j=-n}^n \hat{f}(j) e^{ijt}$.

Setzt man für $N = 2n + 1$: $c_{-j} := c_{N-j}$, $j = 1, \dots, n$, so ist

$$\Phi_N(t_k) = \sum_{j=0}^{N-1} c_j e^{ijt_k} = \sum_{j=-n}^n c_j e^{ijt_k} .$$

Approximiert man $\int_0^{2p} g(t) dt \approx \frac{2p}{N} \sum_{k=0}^{N-1} g(t_k)$,



so ergibt sich: $\hat{f}(j) \approx \frac{1}{N} \sum_{k=0}^{N-1} f(t_k) e^{-ij t_k} = \frac{1}{N} \sum_{k=0}^{N-1} f_k \mathbf{w}_k^{-j} = c_j$ (letzter Schritt nach 7.30)

Die Abbildung $F_N : C^N \rightarrow C^N$ mit $(f_k)_{k=0}^{N-1} \mapsto (c_j)_{j=0}^{N-1}$ aus Satz 7.30 heißt Diskrete Fourier-Transformation.

Die Umkehrabbildung $F_N^{-1} : (c_j)_{j=0}^{N-1} \rightarrow (f_k)_{k=0}^{N-1}$ heißt inverse DFT oder Fourier-Synthese.

$$\left(f_k := \sum_{j=0}^{N-1} c_j \mathbf{w}_j^k, (k = 0, 1, \dots, N-1) \right)$$

Anwendung in der Signalverarbeitung :

Elimination hochfrequenter Komponenten, die häufig durch Messfehler verfälscht (Messrauschen) und darüber hinaus die praktische Anwendung oft nicht relevant ist.

„Abschneiden“ hoher Frequenzen bei Hinzunahme der ersten j_{\max} Frequenzen :

$$\tilde{f}_k := \sum_{j=0}^{j_{\max}} c_j \mathbf{w}_j^k + \sum_{j=1}^{j_{\max}} c_{N-j} \mathbf{w}_{N-j}^k, (k = 0, 1, \dots, N-1)$$

Ist das Eingangssignal mit der Frequenz f_{input} abgetastet, so enthält das Ausgangssignal alle Frequenzen

bis $\bar{f}_{output} := \frac{j_{\max}}{N} \cdot f_{input}$.

Matlab-Befehle : fft, ifft

Bemerkung 7.32 : Schnelle Fourier-Transformation (FFT)

Problem : Standard-Algorithmus zur Auswertung von F_N oder F_N^{-1} erfordert $O(n^2)$ Rechenoperationen. (Matrix-Vektor-Multiplikation)

Algorithmus : Cooley-Tuckey :

Sei $N = 2M$ gerade und $\mathbf{w} = e^{i \frac{2p}{N}}$.

Dann gilt für $\mathbf{a}_j = \sum_{k=0}^{N-1} f_k \mathbf{w}^{kj}$: $\mathbf{a}_{2l} = \sum_{k=0}^{M-1} g_k \mathbf{x}^{kl}$ und $\mathbf{a}_{2l+1} = \sum_{k=0}^{M-1} h_k \mathbf{x}^{kl}$

mit $M = \frac{N}{2}$ und $\mathbf{x} := \mathbf{w}^2$ und $g_k := f_k + f_{k+M}$ und $h_k := (f_k - f_{k+M}) \cdot \mathbf{w}^k$

Mit $2M$ Additionen und $2M$ Multiplikationen ($\mathbf{w}^k \rightarrow \mathbf{w}^{k+1} = \mathbf{w} \cdot \mathbf{w}^k$, $h_k = (\dots) \cdot \mathbf{w}^k$) wird die Berechnung von N Summen der Länge N zurückgeführt auf $2M = N$ Summen der Länge $M = \frac{N}{2}$.

Die rekursive Anwendung ist besonders einfach für $N = 2^p$

⇒

Der Aufwand zur Berechnung von $\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_{N-1}$ (Analyse oder Synthese)

beträgt $2N \log_2 N$ Multiplikationen.

Für $N = 8192$ ist $O(n^2) \approx 64 \cdot 10^6$ und $2N \log_2 N \approx 10^5$.

Speicherschema :

wünschenswert : Überspeicherung von (f_k) durch (\mathbf{a}_j) ohne aufwendigen Komponententausch , jedoch unter Beibehaltung der einfachen Termstruktur für \mathbf{a}_{2l} und \mathbf{a}_{2l+1} .

Problem : Einfaches Überspeichern tauscht Komponenten.

Lösung : "Bit-Reversal" – Reihenfolge der Bits umdrehen

k	k_{dual}	$N = 8 \rightarrow M = 4$	$N = 4 \rightarrow M = 2$	k_{dual}
0	000	0	0	000
1	001	2	4	100
2	010	4	2	010
3	011	6	6	110
4	100	1	1	001
5	101	3	5	101
6	110	5	3	011
7	111	7	7	111

Permutation $\mathbf{s} : \{0,1,\dots,N-1\} \rightarrow \{0,1,\dots,N-1\}$ mit $\mathbf{s} \left(\sum_{l=0}^{N-1} a_l 2^l \right) := \sum_{l=0}^{N-1} a_{N-1-l} 2^l$, $a_i \in \{0,1\}$

Einfache Implementierung durch Bit-Manipulationen.

Algorithmus 7.33 : Schnelle Fourier-Transformation (FFT)

Sei $N = 2^p$ und $\mathbf{w} = e^{i\frac{2\pi}{N}}$ oder $\mathbf{w} = e^{-i\frac{2\pi}{N}}$.

Eingabe: $f_0, f_1, \dots, f_{N-1} \in \mathbb{C}$

Ausgabe: $\mathbf{a}_0, \mathbf{a}_1, \dots, \mathbf{a}_{N-1} \in \mathbb{C}$ mit $\mathbf{a}_j := \sum_{k=0}^{N-1} f_k \mathbf{w}^{kj}$

Algorithmus:

```

 $N_{red} := n; \quad z := \mathbf{w};$ 
while  $N_{red} > 1$  do
   $M_{red} := N_{red} / 2;$ 
  for  $j = 0 : (N / N_{red} - 1)$  do
     $l := j \cdot N_{red};$ 
    for  $k = 0 : (M_{red} - 1)$  do <- g's und h's ausrechnen
       $a := f_{l+k} + f_{l+k+M_{red}};$ 
       $f_{l+k+M_{red}} := (f_{l+k} - f_{l+k+M_{red}}) \cdot z^k;$ 
       $f_{l+k} := a;$ 
    endfor;
  endfor;
   $N_{red} := M_{red};$ 
   $z := z^2;$ 
endwhile;

for  $k = 0 : (N - 1)$ 
   $\mathbf{a}_{s(k)} := f_k;$ 

```

Die FFT ist typisches Beispiel eines 'divide-and-conquer' – Algorithmus. (Ist in Signalprozessoren hardwaremäßig implementiert.)

Bem. 7.34 : Auswertung trigonometrischer Polynome

Auswertung von $\sum_{j=1}^n a_j \cos(jt)$, $\sum_{j=1}^n b_j \sin(jt)$?

Fall $t = t_k \Rightarrow$ iFFT (Synthese)

Fall $t \neq t_k$:

Algorithmus von Goertzel :

Für $t \neq r\pi$, ($r = 0, \pm 1, \pm 2, \dots$) berechnen sich

$$U_j := \frac{1}{\sin t} \sum_{k=j}^{n-1} y_k \cdot \sin((k-j+1) \cdot t) \quad , (j = 1, \dots, n-1)$$

mit $U_1 \cdot \sin t = \sum_{k=1}^{n-1} y_k \sin(kt)$ und $U_1 \cdot \cos t - U_2 = \sum_{k=1}^{n-1} y_k \cos(kt)$ rekursiv aus :

$$U_j := y_j + 2U_{j+1} \cdot \cos t - U_{j+2} \quad , (j = n-1, n-2, \dots, 1)$$

mit den zusätzlich initialisierten $U_{n+1} = U_n := 0$

Vorteil : Statt (n-1)mal muss der Winkelausdruck nur 1mal ausgewertet werden.

Nachteil : Obige Drei-Term-Rekursion ist numerisch instabil für $|t| \ll 1$:

Für $|t| \ll 1$ ist die Verwendung von $\cos t$ in der rekursiven Berechnung der U_j mittels der Drei-Term-Rekursion nachteilig. (vgl. Bsp. 3.32 Numerik I)

Relative Fehler $O(u)$ in $c = \cos t$ entsprechen einem absoluten Fehler $O(u)$ in t , hingegen entspricht einem relativen Fehler $O(u)$ in t ein absoluter Fehler $O(u \cdot |t|) \ll O(u)$ in t .

\Rightarrow Verwendung von $\cos t$ für $|t| \ll 1$ vermeiden !

\Rightarrow Verbesserung : gekoppelte Zwei-Term-Rekursion :

Algorithmus von Reinsch

Fall $\cos t > 0$:

$$U_{n+1} := 0$$

$$\Delta U_n := 0$$

for $j = n-1 : -1 : 0$

$$U_{j+1} := U_{j+2} + \Delta U_{j+1}$$

$$\Delta U_j := y_j + \Delta U_{j+1} - 4 \sin^2 \left(\frac{t}{2} \right) \cdot U_{j+1} \quad (*)$$

endfor

Fall $\cos t \leq 0$:

$$U_{n+1} := 0$$

$$\Delta U_n := 0$$

for $j = n-1 : -1 : 0$

$$U_{j+1} := -U_{j+2} + \Delta U_{j+1}$$

$$\Delta U_j := y_j - \Delta U_{j+1} + 4 \cos^2 \left(\frac{t}{2} \right) \cdot U_{j+1}$$

endfor

beachte : Der Winkelausdruck muss nur 1mal berechnet werden

\Rightarrow

$$\sum_{k=1}^{n-1} y_k \sin(kt) = U_1 \sin t$$

$$\sum_{k=1}^{n-1} y_k \cos(kt) = \Delta U_0 + 2 \cdot \left\{ \begin{array}{l} \sin^2 \left(\frac{t}{2} \right) \text{ für } \cos t > 0 \\ -\cos^2 \left(\frac{t}{2} \right) \text{ für } \cos t \leq 0 \end{array} \right\} \cdot U_1$$

Begründung für (*):

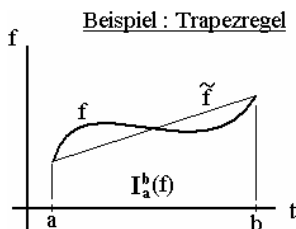
$$\begin{aligned}
 \Delta U_j - \Delta U_{j+1} &= U_j - U_{j+1} \\
 - (U_{j+1} - U_{j+2}) &= U_j + U_{j+2} - 2U_{j+1} \\
 &\stackrel{\text{Goertzel}}{=} y_j + 2(\cos t - 1) \cdot U_{j+1} \\
 &= y_j + 2 \left(\cos \left(\frac{t}{2} + \frac{t}{2} \right) - 1 \right) \cdot U_{j+1} \\
 &= y_j + 2 \left(\cos^2 \left(\frac{t}{2} \right) - \sin^2 \left(\frac{t}{2} \right) - 1 \right) \cdot U_{j+1} \\
 &= y_j + 2 \left(\left(1 - \sin^2 \left(\frac{t}{2} \right) \right) - \sin^2 \left(\frac{t}{2} \right) - 1 \right) \cdot U_{j+1} \\
 &= y_j + 2 \left(-2 \sin^2 \left(\frac{t}{2} \right) \right) \cdot U_{j+1} \\
 &= y_j - 4 \sin^2 \left(\frac{t}{2} \right) \cdot U_{j+1}
 \end{aligned}$$

8. Numerische Integration reeller Funktionen

Problemstellung:

geg.: $f: [a, b] \rightarrow \mathbb{R}$ stückweise stetig

ges.: Näherung für $I(f) := I_a^b(f) = \int_a^b f(t) dt$



Berechnung numerischer Integrale = Quadratur / Quadraturformel

Idee:

Ersetze f durch eine „einfachere“ Funktion \tilde{f} und verwende $I_a^b(f) \approx \int_a^b \tilde{f}(t) dt$,

wobei \tilde{f} eine Interpolation von f ist, durch Polynome, Splines, ...

Eigenschaften von $I(f)$:

- a) Linearität : $I(\mathbf{a}f + \mathbf{b}g) = \mathbf{a} \cdot I(f) + \mathbf{b} \cdot I(g)$
- b) Positivität : $f \geq 0 \Rightarrow I(f) \geq 0$
- c) Additivität : $I_a^c(f) + I_c^b(f) = I_a^b(f)$ für ein $c \in (a, b)$

Bem 8.2 : Kondition der Quadratur

Wegen $\left| \int_a^b (f + \mathbf{d}_f) dt - \int_a^b f dt \right| = \left| \int_a^b \mathbf{d}_f(t) dt \right| \leq \int_a^b |\mathbf{d}_f(t)| dt = \|\mathbf{d}_f\|_1$ hat die Auswertung

von $I(f)$ bzgl. der L_1 -Norm die Kondition $\mathbf{k}_{abs} = 1$ und $\mathbf{k}_{rel} = \frac{I(|f|)}{|I(f)|}$.

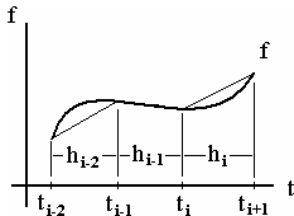
In Analogie zur Auslöschung bei Subtraktion zweier annähernd gleichgroßer Zahlen ist der kritische Fall durch (stark) oszillierende Funktionen f gegeben.

8.1 Newton-Cotes-Quadratur

Bem 8.3 : Zusammengesetzte Trapezregel bzw. Trapezsummen

Gitter $\Delta = \{a = t_0 < t_1 < \dots < t_n = b\}$ mit Schrittweiten $h_i := t_{i+1} - t_i$, $(i = 0, \dots, n - 1)$

$$I_a^b(f) \approx T^{(n)} := \sum_{i=0}^{n-1} h_i \cdot \frac{f(t_i) + f(t_{i+1})}{2}$$



Wegen $\min_{t_i \leq t \leq t_{i+1}} f(t) \leq \frac{f(t_i) + f(t_{i+1})}{2} \leq \max_{t_i \leq t \leq t_{i+1}} f(t)$ ist die Folge der $T^{(n)}$ nach oben und unten durch die Riemannschen Ober- bzw. Untersummen beschränkt.
 \Rightarrow Konvergenz für $f \in C[a, b]$

Definition 8.4 : Quadraturformeln

Gegeben seien $t_0, \dots, t_n \in [a, b]$ (=Knoten) und die zugehörigen Gewichte $w_0, \dots, w_n \in \mathbb{R}$ mit $\sum_{i=0}^n w_i = 1$.

Dann heisst $\tilde{I}(f) := (b - a) \cdot \sum_{i=0}^n w_i f(t_i)$ Quadraturformel zur numerischen Berechnung von $I_a^b(f)$.

Eine Quadraturformel heisst positiv, falls alle $w_i > 0$.

Bem 8.5 : Interpolation und Quadratur

Idee : Approximiere f durch ein Interpolationspolynom zu den Stützstellen $t_0, \dots, t_n \in [a, b]$.

Und zwar mit den Lagrangeschen Basispolynomen : $f(t) \approx \tilde{f}(t) := \sum_{i=0}^n f(t_i) L_i^{(n)}(t)$.

z.B. Trapezregel auf einem Teilintervall = lineare Interpolation

Ergebnis :

$$\tilde{I}(f) := I(\tilde{f}) = \int_a^b \tilde{f}(t) dt = \sum_{i=0}^n \left(f(t_i) \cdot \int_a^b L_i^{(n)}(t) dt \right) = \sum_{i=0}^n w_i f(t_i)$$

$$\text{mit } w_i := \frac{1}{b-a} \cdot \int_a^b L_i^{(n)}(t) dt$$

Eigenschaften :

- $\tilde{I}(\mathbf{p}) = I(\mathbf{p})$, $(\mathbf{p} \in \Pi_n)$

- Sind $n+1$ paarweise verschiedene Knoten $t_0, \dots, t_n \in [a, b]$ gegeben , so bestimmt die

Bedingung $\tilde{I}(\mathbf{p}) = I(\mathbf{p})$ eindeutig die Gewichte der Quadraturformel $(b-a) \cdot \sum_{i=0}^n w_i f(t_i)$, denn durch

Einsetzen der Lagrangeschen Basispolynome folgt :

$$\begin{aligned} \int_a^b L_j^{(n)}(t) dt &= I(L_j^{(n)}) = \tilde{I}(L_j^{(n)}) = (b-a) \cdot \sum_{i=0}^n w_i L_j^{(n)}(t_i) \\ &= (b-a) \cdot \sum_{i=0}^n w_i \mathbf{d}_{i,j} = (b-a) \cdot w_j \quad , (j = 0, \dots, n) \end{aligned}$$

Bem. 8.6 : Newton-Cotes-Formeln

Idee : Wähle Stützstellen t_0, \dots, t_n äquidistant in $[a, b]$.

\Rightarrow

Newton-Cotes-Formeln :

$$t_i = a + i \cdot \frac{b-a}{n} \quad , (i = 0, \dots, n)$$

$$\tilde{I}(f) = (b-a) \cdot \sum_{i=0}^n w_i f(t_i) = (b-a) \cdot \sum_{i=0}^n w_i f\left(a + i \cdot \frac{b-a}{n}\right)$$

Newton-Cotes-Gewichte :

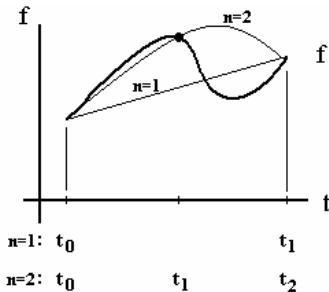
$$w_i = \frac{1}{b-a} \int_a^b \left(\prod_{\substack{j=0 \\ j \neq i}}^n \frac{t - (a + j \cdot h)}{(a + i \cdot h) - (a + j \cdot h)} \right) dt = \frac{1}{n} \int_0^n \prod_{\substack{j=0 \\ j \neq i}}^n \frac{s-j}{i-j} ds$$

$$\text{mit der Schrittweite } h := \frac{b-a}{n} \text{ und } s := \frac{t-a}{h} . \quad \left(ds = \frac{1}{h} dt = \frac{n}{b-a} dt \right)$$

Beispiele :

n	w_i	Name	Fehler
1	1/2 , 1/2	Trapezregel	$\frac{h^3}{12} f'''(\mathbf{t})$

2	1/6, 2/3, 1/6	Keplersche Fassregel	$\frac{h^5}{90} f^{(4)}(\mathbf{t})$
3	1/8, 3/8, 3/8, 1/8	Newtonsche 3/8-Regel	$\frac{3h^5}{80} f^{(4)}(\mathbf{t})$



n=1 : Trapez
n=2 : quadratisches Interpolationspolynom

Jeweils mit einem von f, a und b abhängigen $\mathbf{t} \in [a, b]$.

Die Positivität ist garantiert für $n \leq 7$; ab $n = 8$ auch negative Gewichte.

Lemma 8.7 : Mittelwertsatz in Integralform

Sei $g, h \in C[a, b]$.

Hat g auf $[a, b]$ stets ein und dasselbe Vorzeichen, d.h. $g(t) \geq 0$ oder $g(t) \leq 0$ für $t \in [a, b]$,

so gibt es ein $\mathbf{t} \in [a, b]$ mit $\int_a^b h(t)g(t) dt = h(\mathbf{t}) \cdot \int_a^b g(t) dt$.

Lemma 8.8 : Approximationsfehler der Trapezregel

Sei $f \in C^2[a, b]$ und $h := b - a$.

Dann gibt es ein $\mathbf{t} \in [a, b]$, so dass sich der Approximationsfehler der

Trapezregel $\tilde{I}(f) = T := \frac{b-a}{2}(f(a) + f(b))$ darstellen lässt als: $T - \int_a^b f(t) dt = \frac{h^3}{12} f''(\mathbf{t})$.

Lemma 8.9 : Approximationsfehler der Keplerschen Fassregel

Die Keplersche Fassregel $S = \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right)$ ist für beliebige

Polynome $p \in \Pi_3$ exakt.

Ist $f \in C^4[a, b]$, so lässt sich der Approximationsfehler darstellen als: $S - \int_a^b f(t) dt = \frac{h^5}{90} f^{(4)}(\mathbf{t})$,

mit $h := \frac{b-a}{2}$ und $\mathbf{t} \in [a, b]$.

Lemma 8.10 : Approximationsfehler der zusammengesetzten Trapezregel

Zu äquidistanten Stützstellen $t_i = a + i \cdot h$, ($i = 0, \dots, n$), $h := \frac{b-a}{n}$

sei $T_i := \frac{h}{2}(f(t_i) + f(t_{i+1}))$ die Trapezregel auf $[t_i, t_{i+1}]$.

Der Approximationsfehler $T(h) = \sum_{i=0}^{n-1} T_i = \frac{b-a}{n} \cdot \left(\frac{f(a) + f(b)}{2} + \sum_{i=1}^{n-1} f(a+i \cdot h) \right)$

lässt sich für $f \in C^2[a, b]$ darstellen als: $T(h) - \int_a^b f(t) dt = \frac{h^2}{12}(b-a)f''(\mathbf{t})$,

mit einem $\mathbf{t} \in [a, b]$.

8.2 Gauß-Christoffel-Quadratur

Bem 8.11 : Problemstellung :

Wie sind für vorgegebene $n \in \mathbb{N}$ die $n+1$ Knoten $t_0^{(n)}, t_1^{(n)}, \dots, t_n^{(n)}$ zu wählen, damit man

Gewichte $w_0^{(n)}, w_1^{(n)}, \dots, w_n^{(n)}$ für die Quadraturformel $\tilde{I}_n(f) := \sum_{i=0}^n (w_i^{(n)} \cdot f(t_i^{(n)}))$ angeben kann,

welche alle Polynome $p \in \Pi_N$ mit möglichst großem $N \in \mathbb{N}$ exakt integriert ($\tilde{I}_n(p) = I(p)$) ?

Hypothese :

$N = 2n+1$ möglich, da $2n+2$ freie Parameter $t_i^{(n)}, w_i^{(n)}$, $(i = 0, 1, \dots, n)$

Analytischer Rahmen :

$I(f) = \int_a^b w(t) f(t) dt$ gewichtete Integrale

$w(t) > 0$ so, dass die Folge der Momente $\mathbf{m}_k := \int_a^b w(t) \cdot t^k dt$ existiert und beschränkt bleibt.

(Zulässig ist auch $a = -\infty, b = +\infty$)

Beispiele :

$w(t)$	$\frac{1}{\sqrt{1-t^2}}$	e^{-t}	e^{-t^2}	1
Intervall $[a, b]$	$[-1, 1]$	$[0, \infty)$	$(-\infty, \infty)$	$[-1, 1]$

Lemma 8.12 : Orthogonalität von Knotenpolynomen

Ist für jedes $k \leq n$ die Quadraturformel \tilde{I}_k für alle Polynome $p \in \Pi_{2k+1}$ exakt, d.h. $\tilde{I}_k(p) = I(p)$, so sind die Polynome $\{P_0, P_1, \dots, P_n\}$ mit $P_k(t) := (t - t_0^{(k)}) \dots (t - t_1^{(k)}) \dots (t - t_k^{(k)}) \in \Pi_k$, $(k = 0, \dots, n)$

orthogonal bezüglich des von w induzierten Skalarprodukts $\langle f, g \rangle := \int_a^b w(t) f(t) g(t) dt$.

Lemma 8.13 : Eigenschaften von Orthogonalpolynomen

- a) Zu einem gegebenen Skalarprodukt $\langle f, g \rangle := \int_a^b w(t) f(t) g(t) dt$ gibt es genau dann eine Familie $\{P_0, P_1, \dots, P_n\}$ von Polynomen $P_k \in \Pi_k$ mit führenden Koeffizienten = 1 und $\langle P_j, P_k \rangle = d_{j,k} \cdot \langle P_k, P_k \rangle$.

Beweis : Gram-Schmidtsches Orthogonalisierungsverfahren angewendet auf $\{1, t, t^2, \dots, t^n\}$ (wobei die Normierung entfällt)

- b) Jedes dieser Polynome P_k hat genau k einfache Nullstellen in (a, b) .

Lemma 8.14 : Gewichte der Gauß-Quadraturformeln

Seien $t_0^{(n)}, \dots, t_n^{(n)}$ die Nullstellen des $(n + 1)$ -ten Orthogonalpolynoms P_{n+1} .

Dann gilt für jede Quadraturformel $\tilde{I}_n(f) = \sum_{i=0}^n w_i^{(n)} f(t_i^{(n)})$:

$$w_i^{(n)} = \int_a^b w(t) L_i^{(n)}(t) dt, \quad (i = 0, \dots, n) \Leftrightarrow \tilde{I}_n \text{ ist exakt auf } \Pi_n$$

$$\Leftrightarrow \tilde{I}_n \text{ ist exakt auf } \Pi_{2n+1}$$

Satz 8.15 : Gauß-Christoffel-Quadratur

- a) Es existieren eindeutig bestimmte Knoten $t_0^{(n)}, \dots, t_n^{(n)}$ und Gewichte $w_0^{(n)}, \dots, w_n^{(n)}$, so dass die Quadraturformel $\tilde{I}_n(f) = \sum_{i=0}^n w_i^{(n)} f(t_i^{(n)})$ Polynome bis zum Grad $2n + 1$ exakt integriert, d.h.

$$\tilde{I}_n(p) = \int_a^b w(t) p(t) dt, \quad p \in \Pi_{2n+1} \quad \text{- d.h. die obige Hypothese gilt}$$

- b) Die Knoten $t_0^{(n)}, \dots, t_n^{(n)}$ sind die Nullstellen des $(n + 1)$ -ten Orthogonalpolynoms P_{n+1} bzgl. der Gewichtsfunktion $w(t)$ und die Gewichte ergeben sich aus $w_i^{(n)} = \int_a^b w(t) L_i^{(n)}(t) dt$ mit den Lagrangeschen Basispolynomen $L_i^{(n)}(t), (i = 0, 1, \dots, n)$.

- c) Die Gewichte $w_i^{(n)}$ sind positiv, d.h. die Quadraturformel \tilde{I}_n ist positiv.

d) Es gilt: $w_i^{(n)} = \frac{\langle P_n, P_n \rangle}{P_{n+1}'(t_i^{(n)}) \cdot P_n(t_i^{(n)})}$

Def. 8.16 : Gauß-Christoffel-Quadratur

Die nach Satz 8.15 eindeutig bestimmten positiven Quadraturformeln \tilde{I}_n heißen Gauß-Christoffel-Formeln

zum Gewicht w .

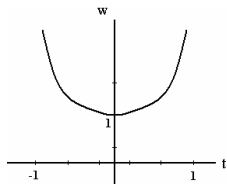
Satz 8.17 : Approximationsfehler der GCQ

Für jedes $f \in C^{2n+2}[a, b]$ lässt sich der Approximationsfehler der GCQ ausdrücken durch :

$$\int_a^b w(t) f(t) dt - \tilde{I}_n(f) = \frac{f^{(2n+2)}(\mathbf{t})}{(2n+2)!} \cdot \langle P_{n+1}, P_{n+1} \rangle \quad \text{für ein } \mathbf{t} \in [a, b].$$

Beispiel 8.18 : Gauß-Tschebyscheff-Quadratur

$$w(t) = \frac{1}{\sqrt{1-t^2}}, \quad t \in (-1, 1)$$



Die in Bemerkung 5.12.c) eingeführten Tschebyscheff-Polynome $T_k = \cos k \cdot \arccos t$, $t \in [-1, 1]$

sind orthogonal bzgl. $\langle f, g \rangle := \int_{-1}^1 \frac{f(t)g(t)}{\sqrt{1-t^2}} dt$, denn

$$\begin{aligned} \int_{-1}^1 \frac{T_k(t) \cdot T_j(t)}{\sqrt{1-t^2}} dt &= \int_p^0 \frac{-\sin x \cdot \cos(kx) \cdot \cos(jx)}{\sin x} dx \quad (t = \cos x \Rightarrow dt = -\sin x dx) \\ &= \int_0^p \cos(kx) \cdot \cos(jx) dx = \begin{cases} p & , \text{ falls } k = j = 0 \\ p/2 & , \text{ falls } k = j > 0 \\ 0 & \text{sonst} \end{cases} \end{aligned}$$

Die Tschebyscheff-Knoten $t_i^{(n)}$ ergeben sich als Nullstellen der Orthogonalpolynome

$$P_{n+1}(t) := \frac{1}{2^n} T_{n+1}(t) \quad \text{zu} \quad t_i^{(n)} = \cos\left(\frac{2i+1}{2n+2} p\right), \quad (i = 0, 1, \dots, n) \quad (*)$$

Gewichte für $n > 0$ nach Satz 8.15.d) :

$$\begin{aligned} w_i^{(n)} &= \frac{1}{\frac{1}{2^n} T_{n+1}'(t_i^{(n)}) \cdot \frac{1}{2^{n-1}} T_n(t_i^{(n)})} \cdot \int_{-1}^1 \frac{\left(\frac{1}{2^{n-1}} T_n(t)\right)^2}{\sqrt{1-t^2}} dt \\ &= \frac{2 \frac{p}{2}}{T_{n+1}'(t_i^{(n)}) \cdot T_n(t_i^{(n)})} \end{aligned}$$

Mit $T_{n+1}'(t) \cdot T_n(t) = -\sin(n+1) \cdot \arccos t \cdot (n+1) \cdot \frac{-1}{\sqrt{1-t^2}} \cdot \cos n \cdot \arccos t$ folgt :

$$= (n+1) \frac{\sin((n+1)x) \cdot \cos(nx)}{\sin x}, \quad (t = \cos x)$$

$$\begin{aligned}
&= (n+1) \frac{\sin((2n+1)x) - \cos((n+1)x) \cdot \sin(nx)}{\sin x} \\
&= (n+1) \frac{\sin((2n+2)x) \cdot \cos x - \cos((2n+2)x) \cdot \sin x - \cos((n+1)x) \cdot \sin(nx)}{\sin x}
\end{aligned}$$

Einsetzen von $x = \frac{2i+1}{2n+2} \mathbf{p}$ ergibt :

- $\sin((2n+2)x) = 0$
- $\cos((n+1)x) = \cos\left(\frac{2i+1}{2} \mathbf{p}\right) = 0$
- $\cos((2n+2)x) = \cos((2i+1)\mathbf{p}) = -1$

$$\Rightarrow w_i^{(n)} = \frac{\mathbf{p}}{n+1}$$

Ergebnis der Gauß-Tschebyscheff-Quadratur :

$$\tilde{I}_n(f) = \sum_{i=0}^n w_i^{(n)} f(t_i^{(n)}) = \frac{\mathbf{p}}{n+1} \cdot \sum_{i=0}^n f\left(\cos\left(\frac{2i+1}{2n+2} \mathbf{p}\right)\right) \approx \int_{-1}^1 \frac{f(t)}{\sqrt{1-t^2}} dt$$

Approximationsfehler :

$$\int_{-1}^1 \frac{f(t)}{\sqrt{1-t^2}} dt - \tilde{I}_n(f) = \frac{\mathbf{p}}{2^{2n+1} \cdot (2n+2)!} \cdot f^{(2n+2)}(t) \quad \text{für ein } t \in [-1,1].$$

Desweiteren :

$$\langle P_{n+1}, P_{n+1} \rangle = \frac{\mathbf{p}}{2} \cdot \left(\frac{1}{2^n}\right)^2$$

Bem. 8.19 : Typische Gauß-Christoffel-Quadraturformeln

$w(t)$	Intervall $[a,b]$	orthogonales Polynomsystem
$\frac{1}{\sqrt{1-t^2}}$	$[-1,1]$	Tschebyscheff-Polynome T_n
e^{-t}	$[0, \infty)$	Laguerre-Polynome L_n
e^{-t^2}	$(-\infty, \infty)$	Hermite-Polynome H_n
1	$[-1,1]$	Legendre-Polynome P_n

Vorteile :

Sehr hohe Approximationsordnung für sehr glatte Integranden schon bei geringer Knotenzahl , auch für singuläre Integranden und/oder unendliche Integrationsintervalle.

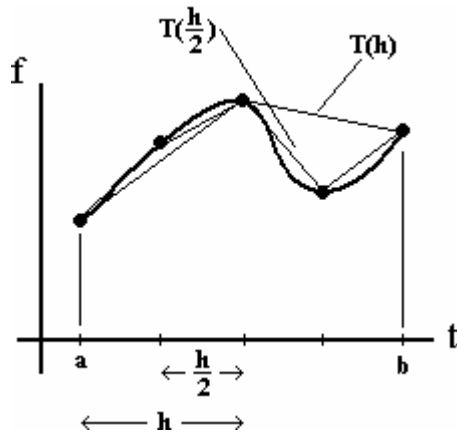
Nachteile :

- Im Allgemeinen aufwendige Berechnung der Knoten und Gewichte durch Berechnung der orthogonalen Polynome oder anschließende Nullstellen-Bestimmung.
- Erhöht man zur Verringerung des Approximationsfehlers die Zahl der Knoten , so sind i.A. sämtliche Funktionswerte $f(t_i^{(n)})$ neu zu berechnen.

8.3 Extrapolationsverfahren, Romberg-Quadratur

Beispiel 8.20: Simpsonregel

Berechnung zweier Näherungswerte $T(h)$ und $T\left(\frac{h}{2}\right)$ mit der zusammengesetzten Trapezregel.
(vgl. Bem. 8.3 und Lemma 8.10)



$$T(h) = \frac{h}{2} f(a) + h \cdot f(a+h) + \dots + \frac{h}{2} f(b)$$

$$T\left(\frac{h}{2}\right) = \frac{h}{4} f(a) + \frac{h}{2} f\left(a + \frac{h}{2}\right) + \frac{h}{2} f(a+h) + \dots + \frac{h}{4} f(b)$$

$$T(h) = I(f) + \frac{h^2}{12} (b-a) \cdot f''(t_h)$$

$$T\left(\frac{h}{2}\right) = I(f) + \frac{h^2}{48} (b-a) \cdot f''\left(t_{h/2}\right)$$

Idee:

Ist $f''(t_h) \approx f''\left(t_{h/2}\right)$, so ergibt sich eine sehr viel genauere Näherung mit der Simpsonregel:

$$S\left(\frac{h}{2}\right) := T\left(\frac{h}{2}\right) + \frac{1}{3} \left(T\left(\frac{h}{2}\right) - T(h) \right)$$

Es gilt:

$$S\left(\frac{h}{2}\right) = \frac{h}{6} f(a) + \frac{2h}{3} f\left(a + \frac{h}{2}\right) + \frac{h}{3} f(a+h) + \dots + \frac{h}{6} f(b)$$

$$= \sum_{i=0}^{n-1} S_i \quad \text{mit } S_i := \frac{h}{6} \left(f(a+ih) + 4 \cdot f\left(a + \left(i + \frac{1}{2}\right)h\right) + f(a+(i+1)h) \right)$$

$$S_i \approx \int_{a+ih}^{a+(i+1)h} f(t) dt \quad \text{Keplersche Fassregel}$$

$$S(h) \approx \int_a^b f(t) dt \quad \text{Simpsonregel}$$

Analog zu Lemma 8.10 zeigt man $S(h) - \int_a^b f(t) dt = \frac{h^4}{180} (b-a) \cdot f^{(4)}(t_s)$ mit einem $t_s \in (a,b)$.

\Rightarrow Fehler i.A. viel kleiner als bei $T\left(\frac{h}{2}\right)$, keine zusätzlichen Auswertungen von f erforderlich.

$S\left(\frac{h}{2}\right)$ ergibt sich durch Auswertung des linearen Interpolationspolynoms $\mathbf{p}(s)$

durch $(h, T(h))$ und $\left(\frac{h}{2}, T\left(\frac{h}{2}\right)\right)$ für $s = 0$. (Extrapolationsverfahren)

Satz 8.21 : Euler-MacLaurinsche Summenformel

Ist $g \in C^{2m+2}[0,1]$, so gilt :

$$\int_0^1 g(t) dt = \frac{g(0) + g(1)}{2} + \sum_{k=1}^m \frac{B_{2k}}{(2k)!} \cdot (g^{(2k-1)}(0) - g^{(2k-1)}(1)) - \frac{B_{2m+2}}{(2m+2)!} \cdot g^{(2m+2)}(t)$$

mit einem $t \in (0,1)$.

Hierbei bezeichnet B_k die Bernoullizahlen $B_k(0)$ mit

$$B_1(x) := x - \frac{1}{2} \quad \text{und} \quad B_{k+1}'(x) = (k+1) \cdot B_k(x) \quad , \quad (k \geq 1)$$

$$\text{z.B. } B_2 = \frac{1}{6} \quad , \quad B_4 = \frac{1}{30} \quad , \quad B_6 = \frac{1}{42} \quad , \quad B_8 = -\frac{1}{30} \quad , \quad (B_{2k} \sim (2k)! \text{ für } k \rightarrow \infty)$$

Folgerung 8.22 : h^2 -Entwicklung des Approximationsfehler der zusammengesetzten Trapezregel

Unterteilt man $[a,b]$ äquidistant in N Teilintervalle der Länge $h := \frac{b-a}{N}$, so ergibt die wiederholte

Anwendung von Satz 8.21

$$T(h) = \int_a^b f(t) dt + \sum_{k=1}^m \underbrace{h^{2k} \cdot \frac{B_{2k}}{(2k)!}}_{\text{Fehlersiehe } \frac{h^2}{12}, \text{ oben } \frac{B_2}{2}} \cdot (f^{(2k-1)}(b) - f^{(2k-1)}(a)) + h^{2m+2} \cdot \frac{B_{2m+2}}{(2m+2)!} \cdot f^{(2m+2)}(t)$$

mit einem $t \in (a,b)$

Beweis : Ausgehend von fortlaufender partieller Integration : $\int_0^1 g(t) dt = B_1(t) \cdot g(t) \Big|_0^1 - \int_0^1 B_1(t) \cdot g'(t) dt \dots$ usw.

Bem. 8.23 : Asymptotische Entwicklung

Ein Näherungsverfahren $T(h)$ zur Berechnung einer Größe t_0 besitzt eine asymptotische Entwicklung in h^p bis zur Ordnung $p \cdot m$, falls es Konstanten $C_p, C_{2p}, C_{3p}, \dots, C_{mp}$ gibt, so dass gilt:

$$T(h) = t_0 + C_p h^p + C_{2p} h^{2p} + \dots + C_{mp} h^{mp} + O(h^{(m+1)p}) \quad , (h \rightarrow 0)$$

Die zusammengesetzte Trapezregel hat für Funktionen $f \in C^{2m+2}[a, b]$ eine asymptotische Entwicklung in h^2 bis zur Ordnung $2m$. (vgl. Folgerung 8.22)

Bem. 8.24 : Extrapolationsverfahren

Ansatz: Auswertung von $T(h)$ für k verschiedene Schrittweiten $h = h_{i-k+1}, h = h_{i-k+2}, \dots, h = h_i$ und Bestimmung des Interpolationspolynoms $P_{ik} \in \Pi_{k-1}$ zu den k Stützpunkten

$$\left((h_{i-k+1}^p, T(h_{i-k+1}^p)), \dots, (h_i^p, T(h_i^p)) \right)$$

Näherungswert: $I(f) \approx T_{ik} := P_{ik}(0) \quad , (k = 1, \dots, i)$

Auswertung von $P_{ik}(0)$ mittels Neville-Schema (vgl. Bem. 7.5):

$$T_{i1} := T(h_i)$$

$$T_{ik} := T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\left(\frac{h_{i-k+1}}{h_i}\right)^p - 1} \quad , (k = 2, \dots, i)$$

Extrapolationstabelle:

$$\begin{array}{ccccccc} T_{11} & & & & & & \\ T_{21} & \rightarrow & T_{22} & & & & \\ \vdots & & & \rightarrow & & & \\ \vdots & & & & \ddots & & \\ T_{k1} & \rightarrow & T_{k2} & \rightarrow & & \rightarrow & T_{kk} \end{array}$$

Satz 8.25 : Approximationsfehler des Extrapolationsverfahrens

Sei $T(h)$ ein Verfahren mit einer asymptotischen Entwicklung in h^p bis zur Ordnung $m \cdot p$.

Für die Fehler der Extrapolationswerte T_{ik} zu den paarweise voneinander verschiedenen

Schrittweiten h_1, \dots, h_m gilt:

$$|T_{ik} - t_0| = |C_{kp}| \cdot h_{i-k+1}^p \cdot \dots \cdot h_{i-1}^p \cdot h_i^p + \sum_{j=i-k+1}^i O(h_j^{(k+1)p}) \quad , (k = 1, \dots, i \leq m \quad , \quad h_j \leq h \rightarrow 0)$$

Lemma 8.26 : Eigenschaften der Lagrangeschen Basispolynome

Für die Lagrangeschen Basispolynome $L_j^{(n)}$, $(j = 0, 1, \dots, n)$ zu den paarweise voneinander verschiedenen Stützstellen t_0, t_1, \dots, t_n gilt:

$$\sum_{j=0}^n L_j^{(n)}(0) \cdot t_j^r = \begin{cases} 1 & \text{für } r = 0 \\ 0 & \text{für } r = 1, \dots, n \\ (-1)^n t_0 t_1 \dots t_n & \text{für } r = n + 1 \end{cases}$$

Beweis:

$P(t) = t^r$ ist Interpolationspolynom zu den Stützpunkten $(t_0, t_0^r), (t_1, t_1^r), \dots, (t_n, t_n^r)$ für $r = 0, 1, \dots, n$.

$\Rightarrow P(t) = t^r = \sum_{j=0}^n (L_j^{(n)}(t) \cdot P(t_j)) = \sum_{j=0}^n (L_j^{(n)}(t) \cdot t_j^r)$ und damit die Behauptung für $r \leq n$ durch Einsetzen von $t = 0$.

Zu $r = n + 1$ betrachtet man $q(t) := t^{n+1} - \sum_{j=0}^n (L_j^{(n)}(t) \cdot t_j^{n+1}) \in \Pi_{n+1}$

Das Polynom q hat den führenden Koeffizienten 1 und Nullstellen t_0, t_1, \dots, t_n
 $\Rightarrow q(t) = (t - t_0)(t - t_1) \dots (t - t_n)$ und für Einsetzen von $t = 0$ gilt:

$$\sum_{j=0}^n (L_j^{(n)}(0) \cdot t_j^{n+1}) = -q(0) = (-1)^n t_0 \dots t_n$$

q.e.d.

Algorithmus 8.27 : Extrapolationsverfahren

geg.: Grundschriftweite H

Schritt 1: wähle Schrittweiten h_1, h_2, \dots mit $h_j = \frac{H}{n_j}$, $(n_{j+1} > n_j)$ und setze $i := 1$

Schritt 2: Berechne $T_{i1} = T(h_i)$

Schritt 3: Berechne T_{i2}, \dots, T_{ii} mit dem Schema von Neville:

$$T_{ik} = T_{i,k-1} + \frac{T_{i,k-1} - T_{i-1,k-1}}{\left(\frac{n_i}{n_{i-k+1}}\right)^p - 1}$$

Schritt 4: Abbruch, falls T_{ii} ausreichend genau oder $i \geq i_{\max}$
 Sonst $i := i + 1$ und gehe zu Schritt 2

Die Anwendung dieses Algorithmus auf die zusammengesetzte Trapezregel heißt **Romberg-Quadratur**.

Wie sind die Schrittweiten und Abbruchkriterien zu wählen?

Bem. 8.28 : Romberg-Quadratur

Aufwand: A_k ... Anzahl der zur Berechnung von T_{kk} benötigten Funktionsauswertungen

Aufwandsfolge $\tilde{A} = \{A_1, A_2, \dots\}$ zur Folge $F = \{n_1, n_2, \dots\}$ mit $0 < n_1 < n_2 < \dots$

Rombergfolge :

$$F_R = \{1, 2, 4, 8, 16, \dots\}, (n_i = 2^{i-1})$$

$$\tilde{A}_R = \{2, 3, 5, 9, 17, \dots\}, (A_i = n_i + 1)$$

$$T(h) = \underbrace{\frac{h}{2} \left(f(a) + 2 \sum_{j=1}^{n-1} f(a + 2jh) + f(b) \right)}_{\text{Summanden zu den geradzähligen Indizes}} + \underbrace{h \cdot \sum_{j=1}^n f(a + (2j-1)h)}_{\text{Summanden zu den ungeraden Indizes}}$$

$$= \frac{1}{2} T_{\substack{(2h) \\ \text{neu zu berechnen}}} + \underbrace{h \cdot \sum_{j=1}^n f(a + (2j-1)h)}_{\text{bleibt wie zuvor}} \quad \text{mit } h := \frac{H}{2n}$$

$$T_{i+1,1} = \frac{1}{2} T_{i1} + h_{i+1} \cdot \sum_{j=1}^{n_i} f(a + (2j-1) \cdot h_{i+1}) \quad \text{mit } h_i := \frac{H}{n_i} \text{ und } n_i = 2^{i-1}$$

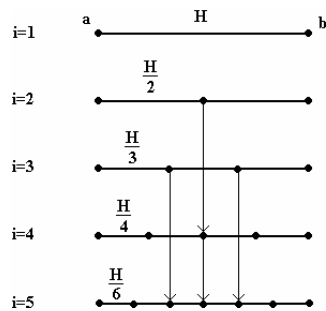
Es gilt : $T_{kk} = H \cdot \sum_{j=1}^{A_k} l_j f_j$ mit positiven Gewichten l_j

Bulirschfolge :

$$F_B = \{1, 2, 3, 4, 6, 8, 12, 16, 24, \dots\}$$

Doppeltes des Vorgängers bzw. :

$$n_i = \begin{cases} 1 & \text{falls } i = 1 \\ 2^k & \text{falls } i = 2k \\ 3 \cdot 2^{k-1} & \text{falls } i = 2k + 1 > 1 \end{cases}$$



$$\tilde{A}_B = \{2, 3, 5, 7, 9, 13, \dots\}$$

Der Aufwand wächst langsamer als für die Rombergfolge, jedoch treten für $k \geq 3$ negative Gewichte l_j in T_{kk} auf.

Harmonische Folge :

$$F_H = \{1, 2, 3, 4, 5, 6, \dots\}, (n_i = i)$$

$$A_H = \{2, 3, 5, 7, 11, 13, 19, \dots\}$$

Zur Quadratur ungeeignet, da bei der Berechnung von T_{i1} kaum auf bereits zuvor berechnete

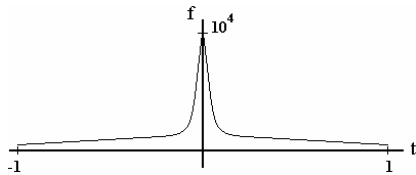
Funktionswerte von f zurückgegriffen werden kann.

Jedoch sehr gut geeignet für Probleme , bei denen A_k weitgehend unabhängig von n_1, n_2, \dots, n_{k-1} ist.

(Beispiel : Extrapolationsverfahren für Anfangswertprobleme gewöhnlicher Differentialgleichungen)

Beispiel. 8.29 : Nadelimpuls

$$I(f) = \int_{-1}^1 \frac{1}{10^{-4} + t^2} dt$$



Rombergfolge :

k	A_k	T_{kk}	$e_{kk} := \frac{ T_{kk} - I(f) }{ I(f) }$ (relativer Fehler)
1	2	1.998	0.99 (99%)
2	3	13334.0	42.0 (4200%)
3	5	2627.7	7.6
4	9	1551.9	4.0
...
8	129	293.01	0.061
...
13	4097	312.16	$1.1 \cdot 10^{-9}$

Dieselbe Fehlergenauigkeit erhält man mit TRAPEX (extrapolierende Trapezregel – siehe Deuflhard) mit 321 Funktions-Auswertungen.

Problem : Hoher Rechenaufwand der klassischen Romberg-Quadratur

Ausweg :

i) Zerlege $[-1,1]$ in Teilintervalle $[a_i, a_{i+1}]$ mit $a = a_0 < a_1 < \dots < a_N = b$

$$\Rightarrow I(f) = \sum_{j=1}^{N-1} I_{a_i}^{a_{i+1}}(f)$$

ii) Approximiere $I_{a_i}^{a_{i+1}}(f)$ mittels Romberg-Quadratur

$$T_{kk} \rightarrow T_{k(i),k(i)}$$

unterschiedliche Tiefe der Extrapolationsschemata in den einzelnen Teilintervallen

Adaptiver Algorithmus :

Bestimme a_i und $k(i)$ während der Berechnung automatisch so , dass bei möglichst niedrigem Gesamtrechenaufwand der Fehler unterhalb einer vorgesehenen Genauigkeitsschranke **TOL** bleibt.

Grundprinzip :

Gleichverteilung des Fehlers , d.h. Bestimmung der a_i so , dass jedes der N Teilintervalle einen annähernd gleich großen Beitrag zum Gesamtfehler leistet.

Fehlerschätzer :

Benötigt wird ein möglichst guter Schätzwert für $|T_{kk} - I(f)|$ bzw. für $\left| T_{k(i),k(i)} - \int_{a_i}^{a_{i+1}} f(t) dt \right|$
(, d.h. für jedes Intervall.)

ideal :

Obere und untere Schranke für den tatsächlichen Fehler , zumindest asymptotisch für $\frac{H}{n_k} \rightarrow 0$.

Aus Satz 8.25 folgt :

$$|T_{kk} - I(f)| = |C_{2k}| \cdot h_1^2 \cdot \dots \cdot h_{k-1}^2 h_k^2 + \sum_{j=1}^k O(h_j^{2(k+1)})$$

$$|T_{k,k-1} - I(f)| = |C_{2(k-1)}| \cdot h_2^2 \cdot \dots \cdot h_{k-1}^2 h_k^2 + \sum_{j=2}^k O(h_j^{2(k+1)})$$

(Der zweite Summand wird verschwindend klein und wird nicht betrachtet.)

Für ausreichend kleine Schranken ($h_1 \ll 1$) und Konstanten $|C_{2k}| \approx |C_{2(k-1)}|$ gilt :

$$|T_{kk} - I(f)| \ll |T_{k,k-1} - I(f)| \quad , \text{ also ist wegen}$$

$$|T_{k,k-1} - I(f)| - |T_{kk} - I(f)| \leq |T_{k,k-1} - T_{kk}| \leq |T_{k,k-1} - I(f)| + |T_{kk} - I(f)|$$

\Rightarrow

$$|T_{k,k-1} - T_{kk}| \text{ ist ein guter Schätzwert für } T_{k,k-1}.$$

Beachte :

Man schätzt unter Verwendung von $T_{k,k-1}$ und T_{kk} den Fehler von $T_{k,k-1}$, verwendet jedoch anschließend den im Allgemeinen genaueren Näherungswert T_{kk} für $I(f)$.

Algorithmus 8.27 Schritt 4 lautet somit : Abbruch, falls $|T_{i,i-1} - T_{ii}| \leq TOL$ oder $i \geq i_{\max}$

9. Numerische Lösung nichtlinearer Gleichungssysteme

geg. : $F : D \rightarrow R^n$ mit einer offenen Menge $D \subset R^n$

ges. : $\{x : F(x) = 0\}$

Spezialfälle :

a) F ist linear

$$F(x) = Ax - b$$

\Rightarrow Lösung von $Ax - b = 0$ (vgl. Kapitel 2 „Lösen von linearen Gleichungssystemen“)

b) $n = 1$

Nullstellen einer nichtlinearen reellen Funktion

Qualitativer Unterschied zu linearen Gleichungssystemen :

Eindeutigkeit der Lösung von $F(x) = 0$ kann im Allgemeinen nur lokal in einer Umgebung einer Nullstelle x^* von F garantiert werden , z.B. mit dem Satz über die Implizite Funktion.

Hinreichende Bedingung :

$$F = \begin{pmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{pmatrix} \text{ ist stetig differenzierbar und die Jacobimatrix } F_x = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1} & \dots & \frac{\partial f_n}{\partial x_n} \end{pmatrix} \text{ ist regulär.}$$

9.1 Nichtlineare Gleichungssysteme und Fixpunktiteration

Fixpunktgleichungen $\Phi(x) = x$ sind spezielle nichtlineare Gleichungssysteme $F(x) = 0$ mit $F(x) := \Phi(x) - x$.

Idee :

Iterative Bestimmung des Fixpunktes mittels Fixpunktiteration

$$x_{k+1} := \Phi(x_k) \quad , (k \geq 0)$$

für einen geeigneten Startwert $x_0 \in \mathbb{R}^n$.

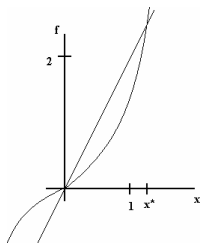
(Wie ist die Funktion Φ und der Startwert x_0 zu bestimmen ?)

Umgekehrt lässt sich jedes nichtlineare Gleichungssystem $F(x) = 0$ als Fixpunktgleichung $\Phi(x) = x$ schreiben.

z.B. mit $\Phi(x) := x - A^{-1} \cdot F(x)$ mit einer beliebigen regulären Matrix $A \in \mathbb{R}^{n \times n}$.

$$F(x) = 0 \iff \Phi(x) = x$$

Beispiel 9.3 : Fixpunktiteration zur Lösung nichtlinearer Gleichungen



Gesucht : $x^* \in \left(0, \frac{\pi}{2}\right)$ mit $f(x^*) = 0$.

Naheliegende Fixpunktgleichungen :

a) $x = \frac{1}{2} \tan x =: j_1(x)$

b) $x = \arctan(2x) =: j_2(x)$

k	$x_{k+1} = \mathbf{j}_1(x_k)$	$x_{k+1} = \mathbf{j}_2(x_k)$
0	1.2	1.2
1	1.2860	1.1760
2	1.70 > $p/2$	1.1687
...	divergent	...
6		1.1655
7		1.1655

Definition 9.4 : Kontrahierende Abbildung

Sei $D \subset \mathbb{R}^n$ eine offene Menge und $\Phi : \bar{D} \rightarrow \mathbb{R}^n$. (\bar{D} = Abschluss von D)

Die Abbildung Φ heißt kontrahierend auf \bar{D} , wenn es eine Konstante $\mathbf{a} \in [0,1)$ gibt mit

$$\|\Phi(x) - \Phi(y)\| \leq \mathbf{a} \cdot \|x - y\|, \quad (x, y \in \bar{D})$$

Lemma 9.5 : Stetig differenzierbare kontrahierende Abbildung

Sei $D \subset \mathbb{R}^n$ eine konvexe, offene Menge und $\Phi : \bar{D} \rightarrow \mathbb{R}^n$ stetig differenzierbar.

Existiert $\mathbf{a} := \sup_{x \in \bar{D}} \|\Phi_x(x)\|$ und ist $\mathbf{a} < 1$, so ist Φ kontrahierend auf \bar{D} .

Satz 9.6 : Banachscher Fixpunktsatz

Sei $E \subset \mathbb{R}^n$ kompakt und $\Phi : E \rightarrow E$ kontrahierend mit einer Lipschitz-Konstanten $\mathbf{a} < 1$.
Dann gilt :

- a) Φ hat genau einen Fixpunkt x^* in E , d.h. $\Phi(x^*) = x^*$.
- b) Für jeden Startwert $x_0 \in E$ konvergiert die Fixpunktiteration $x_{k+1} = \Phi(x_k)$ gegen x^* und es gilt :

i) $\|x^* - x_k\| \leq \frac{\mathbf{a}^k}{1 - \mathbf{a}} \cdot \|x_1 - x_0\|$ (a-priori-Fehlerschranke)

ii) $\|x^* - x_{k+1}\| \leq \frac{\mathbf{a}}{1 - \mathbf{a}} \cdot \|x_{k+1} - x_k\|$ (a-posteriori-Fehlerschranke)

Bem. 9.7 : praktische Aspekte der Fixpunktiteration

- a) Bestimmung einer Näherung für \mathbf{a}

a-priori : Für stetig differenzierbares Φ eine obere Schranke für $\|\Phi_x\|$ berechnen.

a-posteriori : $x_{k+2} - x_{k+1} = \Phi(x_{k+1}) - \Phi(x_k)$

$$\Rightarrow \mathbf{a} \approx \mathbf{a}_k \quad \text{mit} \quad \mathbf{a}_k := \frac{\|x_{k+2} - x_{k+1}\|}{\|x_{k+1} - x_k\|}$$

- b) Abbruch

Gegeben seien Fehlerschranken $ATOL$ und $RTOL$ für den absoluten und relativen Fehler.
Ein Schätzwert für \mathbf{a} sei bekannt.

- i) Anzahl der Iterationsschritte vorab festlegen :

$$x^* \approx x_k \quad \text{mit} \quad \mathbf{a}^k \leq \frac{(1-\mathbf{a}) \cdot (ATOL + RTOL \cdot \|x_0\|)}{\|x_1 - x_0\|}$$

$$\Rightarrow \quad k := 1 + \frac{\log\left(\frac{(1-\mathbf{a}) \cdot (ATOL + RTOL \cdot \|x_0\|)}{\|x_1 - x_0\|}\right)}{\log \mathbf{a}}$$

ii) Abbruchkriterium während der Iteration überprüfen :

$$x^* \approx x_{k+1} \quad , \text{ falls } \frac{\mathbf{a}_k}{1-\mathbf{a}_k} \cdot \|x_{k+1} - x_k\| \leq ATOL + RTOL \cdot \|x_k\| \quad \text{mit } \mathbf{a}_k \text{ wie oben}$$

Ist für ein $k_0 > 0$ $\mathbf{a}_{k_0} > 1$, so deutet dies auf Divergenz von (x_k)
 \Rightarrow (ergebnisloser) Abbruch des Verfahrens

c) Konvergenzbeschleunigung

Mittels Δ^2 -Methode von Aitken oder Methode von Steffensen

Definition 9.8 : Konvergenz , Konvergenzordnung

a) Ein Iterationsverfahren $x_{k+1} = \Psi_k(x_0, x_1, \dots, x_k)$ heißt konvergent , wenn die Folge der Iterierten (x_k) konvergiert.
 Anderenfalls heißt sie divergent.

b) Das Verfahren hat die Konvergenzordnung $p \geq 1$, wenn $x^* := \lim_{k \rightarrow \infty} x_k$ existiert und eine Konstante $C \geq 0$ ($p = 1 : C = \mathbf{a} \in [0,1)$) angegeben werden kann , so dass
 $\|x_{k+1} - x^*\| \leq C \cdot \|x_k - x^*\|^p$, ($k \geq k_0$).

Das Verfahren konvergiert linear , falls $p = 1$, quadratisch , falls $p = 2$

und superlinear , falls gilt : $\|x_{k+1} - x^*\| \leq C_k \cdot \|x_k - x^*\|$ mit $C_k \geq 0$ und $\lim_{k \rightarrow \infty} C_k = 0$.

9.2 Newton-Verfahren

geg.: $F : R^n \rightarrow R^n$ stetig differenzierbar

ges.: x^* mit $F(x^*) = 0$

Ansatz :

Linearisierung von F in einer Umgebung von x_k :

$$F_k(x) := F(x_k) + F_x(x_k) \cdot (x - x_k)$$

$F_k(x) \approx F(x)$ in einer (kleinen) Umgebung von x_k

Bestimme x_{k+1} als Nullstelle von F_k : $x_{k+1} := x_k + p_k$ mit $F_x(x_k) \cdot p_k = -F(x_k)$

Aufwand pro Iterationsschritt :

1 F - Auswertung

1 F_x - Auswertung

1 LGS lösen

Fixpunktdarstellung (formal) :

$$x_{k+1} = \Phi(x_k) \quad \text{mit} \quad \Phi(x) := x - (F_x(x))^{-1} \cdot F(x)$$

Affine Invarianz :

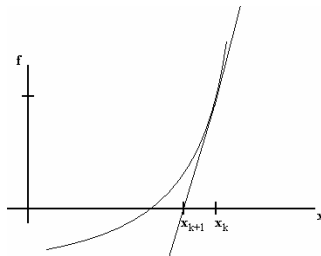
Für jede reguläre Matrix $A \in \mathbb{R}^{n \times n}$ ist $F(x) = 0 \Leftrightarrow A \cdot F(x) = 0$

$\Rightarrow \|F(x)\|$ ist als Abbruchkriterium ungeeignet, da z.B. mit 10^{-20} skaliert werden kann.

Das Newton-Verfahren ergibt (in exakter Arithmetik) für beide Gleichungssysteme dieselbe Folge (x_k) , es ist **affin-invariant** :

$$\begin{aligned} \Phi_A(x) &= x - \left(\frac{\partial}{\partial x} (A \cdot F(x)) \right)^{-1} \cdot A \cdot F(x) \\ &= x - \left(\frac{\partial}{\partial x} F(x) \right)^{-1} \cdot F(x) \\ &= \Phi(x) \end{aligned}$$

Geometrische Interpretation für n=1 :



x_{k+1} ist Schnittpunkt der im Punkt $(x_k, f(x_k))$ an den Graphen von f gelegten Tangente mit der Abszissenachse.

Beispiel 9.10 : Verfahren von Heron

Newton-Verfahren für $0 = f(x) := x^2 - a$ - d.h. Berechnung von \sqrt{a}

$$x_{k+1} = x_k - \frac{x_k^2 - a}{2x_k} = \frac{1}{2} \left(x_k + \frac{a}{x_k} \right)$$

praktisch :

binäre Gleitpunktdarstellung $a = m \cdot 2^p$ mit $0.5 < m \leq 1$

$$\sqrt[p]{a} = \begin{cases} 2^r \cdot \sqrt{m} & , \text{ falls } p = 2r \\ 2^r \cdot \sqrt{m} \cdot \sqrt[4]{1/2} & , \text{ falls } p = 2r - 1 \end{cases}$$

Begründung : Sehr schnelle Konvergenz des Newton-Verfahrens für $m \in (0.5, 1]$.

Beispiel :

$$m = 0.81 \quad , \quad x_0 = 1$$

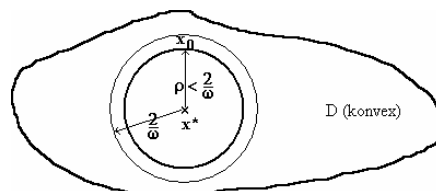
k	x_k
0	1.0000000000
1	0.9050000000
2	0.9000138122
3	0.9000000001

Satz 9.11 : Konvergenz des Newton-Verfahrens

Sei $D \subset \mathbb{R}^n$ offen und konvex , und $F : D \rightarrow \mathbb{R}^n$ stetig differenzierbar mit regulärer Jacobimatrix $F_x(x)$ für alle $x \in D$.

Für F_x gelte die Lipschitz-Bedingung $\left\| (F_x(x))^{-1} \cdot (F_x(x + s \cdot v) - F_x(x)) \cdot v \right\| \leq s \cdot \mathbf{w} \cdot \|v\|^2$ mit $s \in [0,1]$, $x \in D$, $v \in \mathbb{R}^n$ mit $x + v \in D$, Konstante $\mathbf{w} > 0$.

Hat $F(x) = 0$ eine Lösung $x^* \in D$ und gilt für den Startwert $x_0 \in D$: $\mathbf{r} := \|x^* - x_0\| < \frac{2}{\mathbf{w}}$, und ist $B_r(x^*) := \{x \in \mathbb{R}^n : \|x - x^*\| < r\} \subset D$, dann bleibt die durch das Newton-Verfahren definierte Folge (x_k) für $k > 0$ in $B_r(x^*)$ und konvergiert gegen x^* : $\lim_{k \rightarrow \infty} x_k = x^*$.



$$\text{Es gilt } \|x_{k+1} - x^*\| \leq \frac{\mathbf{w}}{2} \cdot \|x_k - x^*\|^2 \quad , (k \geq 0)$$

(D.h. in etwa Verdopplung der genauen Stellen pro Schritt.)

Darüber hinaus ist x^* einzige Lösung von $F(x) = 0$ in $B_{\frac{2}{\mathbf{w}}}(x^*)$.

Beweis :

Linearisiert man F um einen Punkt $x \in D$, so folgt aus dem Hauptsatz der Differential- und Integralrechnung folgende Abschätzung des Restgliedes für beliebige $y \in D$:

$$F(y) - (F(x) + F_x(x) \cdot (y - x))$$

$$\begin{aligned} \text{Führe Funktion } \Psi \text{ ein mit : } \Psi(s) &:= F(x + s(y - x)) \\ (F(y) &= \Psi(1) , F(x) = \Psi(0)) \end{aligned}$$

$$\begin{aligned}
&= \int_0^1 \Psi'(s) ds \\
&= \int_0^1 (F_x(x + s(y-x)) - F_x(x)) \cdot (y-x) ds \\
&\Rightarrow \left\| (F_x(x))^{-1} \cdot (F(y) - F(x) - F_x(x) \cdot (y-x)) \right\| \leq \int_0^1 s \cdot \mathbf{w} \cdot \|y-x\|^2 ds = \frac{\mathbf{w}}{2} \cdot \|y-x\|^2 \\
&\Rightarrow x_{k+1} - x^* = x_k - (F_x(x_k))^{-1} \cdot F(x_k) - x^* \\
&\quad = (F_x(x_k))^{-1} \cdot \left(\underbrace{F(x^*) - F(x_k) - F_x(x_k)}_{=0} \cdot (x^* - x_k) \right) \\
&\|x_{k+1} - x^*\| \leq \frac{\mathbf{w}}{2} \cdot \|x_k - x^*\|^2
\end{aligned}$$

Mittels vollständiger Induktion folgt hieraus $x_k \in B_r(x^*)$, ($k > 0$), denn

$$\text{aus } 0 < \|x_k - x^*\| \leq \mathbf{r} \text{ ergibt sich } \|x_{k+1} - x^*\| \leq \underbrace{\frac{\mathbf{w}}{2} \cdot \|x_k - x^*\|}_{\leq \mathbf{r} \cdot \frac{\mathbf{w}}{2} < 1} \cdot \underbrace{\|x_k - x^*\|}_{\leq \mathbf{r}} < \mathbf{r}.$$

Weiterhin ist $\|x_{k+1} - x^*\| \leq \mathbf{a} \cdot \|x_k - x^*\|$ mit $\mathbf{a} := \mathbf{r} \cdot \frac{\mathbf{w}}{2} < 1$, also konvergiert (x_k) gegen x^* .

Ist x^{**} eine Nullstelle von F in $B_{\frac{2}{\mathbf{w}}}(x^*)$, dann folgt oben mit $y := x^{**}$, $x := x^*$:

$$\begin{aligned}
\|x^{**} - x^*\| &= \left\| (F_x(x^*))^{-1} \cdot \left(\underbrace{F(x^{**}) - F(x^*) - F_x(x^*) \cdot (x^{**} - x^*)}_{=0} \right) \right\| \\
&\leq \frac{\mathbf{w}}{2} \cdot \underbrace{\|x^{**} - x^*\|}_{< \frac{2}{\mathbf{w}}} \cdot \|x^* - x^*\| < \|x^{**} - x^*\|
\end{aligned}$$

, also folgt $x^* = x^{**}$ und damit die lokale Eindeutigkeit der Lösung von $F(x) = 0$. **q.e.d.**

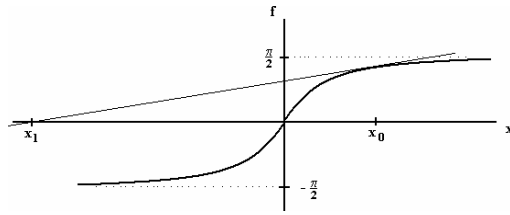
Bem. 9.12 : Gedämpftes Newton-Verfahren

Beispiel :

ges.: Nullstelle von $f(x) = \arctan x$

klassisches Newton-Verfahren :

$$x_{k+1} = x_k + p_k \quad \text{mit} \quad p_k := -\frac{f(x_k)}{f'(x_k)}$$



z.B. mit dem Startwert $x_0 = 10$ versagt das Verfahren (divergiert) : $\lim_{k \rightarrow \infty} f'(x_k) = 0$

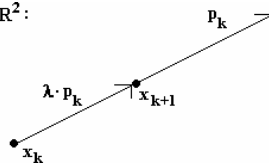
Schlussfolgerung : Das Newton-Verfahren konvergiert lokal sehr schnell , hat jedoch of schlechte globale Konvergenz.

Idee : Modifiziere das Newton-Verfahren so , dass $\left\| (F_x(x_k))^{-1} \cdot F(x_{k+1}) \right\| \leq \left\| (F_x(x_k))^{-1} \cdot F(x_k) \right\|$.

Betrachte zu p_k mit $F_x(x_k) \cdot p_k = -F(x_k)$ Vektoren $x_{k+1}(\mathbf{I}) := x_k + \mathbf{I} \cdot p_k$

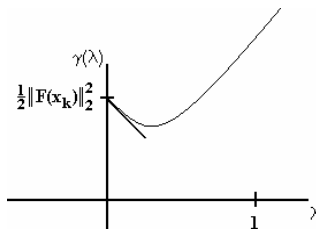
und $\mathbf{g}(\mathbf{I}) := \frac{1}{2} \cdot \|F(x_k + \mathbf{I} \cdot p_k)\|_2^2 \cdot \left(= \frac{1}{2} F^T F \right)$

z.B. im \mathbb{R}^2 :



Wegen $\mathbf{g}'(\mathbf{I}) = \left(F_x^T(x_k + \mathbf{I} \cdot p_k) \cdot F(x_k + \mathbf{I} \cdot p_k) \right)^T \cdot \underbrace{p_k}_{=-(F_x(x_k))^{-1} \cdot F(x_k)}$

ist $\mathbf{g}'(0) = -F^T(x_k) \cdot F(x_k) = -\|F(x_k)\|_2^2 < 0$



Strategie :

Verringere beginnend mit $\mathbf{I}_k^{(0)} := 1$ die Schrittweite \mathbf{I}_k so lange ,

$\left(\mathbf{I}_k^{(l+1)} := c \cdot \mathbf{I}_k^{(l)} \right)$ mit $c \in (0,1)$, z.B. $c = 0.5$

bis $\left\| (F_x(x_k))^{-1} \cdot F(x_k + \mathbf{I}_k^{(l)} \cdot p_k) \right\|_2 \leq \underbrace{\left\| (F_x(x_k))^{-1} \cdot F(x_k) \right\|}_{= -p_k}$

und setze $\mathbf{I}_k := \mathbf{I}_k^{(l)}$, $x_{k+1} := x_k + \mathbf{I}_k p_k$.

Zusatzaufwand pro Reduktionsschritt :

1 F - Auswertung und jeweils 1 Vorwärts- und Rückwärtssubstitution

Varianten :

1.) Setze $\mathbf{I}_k^{(0)} := \min\left(\frac{\mathbf{I}_{k-1}}{d}, 1\right)$ mit einer Konstanten $d < 1$, z.B. $d=1/2$

2.) Forderung $\left\| (F_x(x_k))^{-1} \cdot F(x_k + \mathbf{I}_k^{(0)} \cdot p_k) \right\|_2 \leq \left(1 - \frac{\mathbf{I}_k}{2}\right) \cdot \|p_k\|_2$, um möglichst „optimale“

Dämpfungsfaktoren \mathbf{I}_k zu erreichen

praktische Erfahrung :

Gegenüber dem klassischen Newton-Verfahren verbessert sich das Konvergenzverhalten oft erheblich.

Lokal quadratische Konvergenz jedoch nur , falls $\mathbf{I}_k = 1$, ($k \geq k_0$) .

(D.h. \mathbf{I} also wieder auf 1 gesetzt wird.)

Bem. 9.13 : Vereinfachtes Newton-Verfahren

Der große numerische Aufwand für häufige Neuberechnung und LR-Zerlegung der Jacobimatrix lässt sich verringern , indem die Jacobimatrix über möglichst viele Iterationsschritte konstant gehalten wird :

$$x_{k+l} := x_{k+l-1} + p_{k+l-1} \quad \text{mit} \quad F_x(x_k) \cdot p_{k+l-1} = -F(x_{k+l-1})$$

(konvergiert linear ! , nicht quadratisch !)

Konvergenz :

$$x_{k+l} = \Phi_k(x_{k+l-1}) \quad \text{mit} \quad \Phi_k(x) = x - (F_x(x_k))^{-1} \cdot F(x)$$

$$\frac{\partial}{\partial x} \Phi_k(x) = I - (F_x(x_k))^{-1} \cdot F_x(x) = (F_x(x_k))^{-1} \cdot (F_x(x_k) - F_x(x))$$

Ist $(F_x(x_k))^{-1} \cdot F_x(x)$ Lipschitz-stetig mit der Konstanten L , so ist

$$\left\| \frac{\partial}{\partial x} \Phi_k(x) \right\| = \left\| (F_x(x_k))^{-1} \cdot (F_x(x_k) - F_x(x)) \right\| \leq L \cdot \|x - x_k\| \quad \text{und}$$

$$\mathbf{a} := \sup \left\{ \left\| \frac{\partial}{\partial x} \Phi_k(x) \right\| : \|x - x_k\| \leq \Delta_0 \right\} \leq L \cdot \Delta_0 .$$

Also ist Φ_k kontraktiv , falls Δ_0 hinreichend klein.

Neuberechnung der Jacobimatrix :

$$\text{Schätzwert für Kontraktionskonstante } \mathbf{a}_{k+l} := \frac{\|x_{k+l} - x_{k+l-1}\|}{\|x_{k+l-1} - x_{k+l-2}\|} \quad (\text{Bem. 9.7.a})$$

Neuberechnung der Jacobimatrix $F_x(x_{k+l})$, falls $\mathbf{a}_{k+l} \geq \mathbf{a}_0$. (typische Werte : $\mathbf{a}_0 \in [0.25, 0.5]$)

Algorithmus :

geg.: Startwert x_0 , $\mathbf{a}_0 \in (0,1)$

Schritt 0 : $k := 0$, $\mathbf{a} := 1$

Schritt 1 : if $\mathbf{a} > \mathbf{a}_0$ then $-J := F_x(x_k)$
- LR-Zerlegung von J
- $l := 0, \mathbf{a} := 0$

Schritt 2 : Funktionsauswertung $F := F(x_k)$

Schritt 3 : Bestimme p_k mit $J \cdot p_k = -F$ mittels Vorwärts- und Rückwärtssubstitution

Schritt 4 : $x := x_k + p_k$
 $l := l + 1$
if $l > 1$ then $\mathbf{a} := \frac{\|p_k\|}{\|p_{k-1}\|}$

Schritt 5 : if $(l = 1 \text{ or } \mathbf{a} < 1)$ then $-x_{k+1} := x$
 $-k := k + 1$
if Konvergenz then stop
else goto Schritt 1

Erweiterung : Kombination des vereinfachten mit dem gedämpften Newton-Verfahren

Bem. 9.14 : Berechnung der Jacobimatrix

a) Analytische Differentiation :

Empfehlenswert ist die Anwendung mathematischer Hilfsprogramme (Mathematica, Matlab) und Export der Ableitungen als C- oder Fortran-Quelltext.

Vorteil : - analytisch exakter Ausdruck zur Berechnung der Ableitungen

Nachteile : - große Bearbeitungszeiten und hoher Speicherbedarf für kompliziertere Probleme
- Der analytische Ausdruck für $F(x)$ muss in der Syntax des Hilfsprogramms formuliert werden.

b) Algorithmische Differentiation (AD) :

Ausgangspunkt : C- oder Fortran-Unterprogramm zur Auswertung einer Funktion $F : R^n \rightarrow R^m$
(input : x , output : $F(x)$)

Werkzeuge : z.B. ADIFOR , ADOL-C

Ergebnis : Zusätzliches C- oder Fortran-Unterprogramm zur Auswertung der Jacobimatrix $F_x : R^n \rightarrow R^{m \times n}$ (input : x , output : $F_x(x)$)

Vorteil : - analytisch exakte Berechnungsvorschrift

Nachteile : - sehr großer Speicherbedarf für sehr große Probleme
- Unterstützt wird in der Regel ausschließlich der Sprachstandard.

c) Differenzenapproximation :

$$\left(\frac{\partial F}{\partial x}(x) \right)_{\bullet j} = \frac{F(x + \Delta \cdot e_j) - F(x)}{\Delta} \quad (\text{d.h. alle Elemente der } j\text{-ten Spalte})$$

Aufwand zur Berechnung von F_x für $F : R^n \rightarrow R^m$: $n+1$ Auswertungen von F

Wahl von Δ : Approximationsfehler $O(\Delta)$ und Verstärkung der Rundungsfehler $O\left(\frac{eps}{\Delta}\right)$

$$\Psi(\Delta) := \Delta + \frac{eps}{\Delta} \text{ wird minimal für } \Delta = \sqrt{eps}$$

$$\Rightarrow \Delta := \sqrt{eps} \cdot \max\{|x_j|, \sqrt[4]{eps}\}$$

(Wobei sich in der Praxis $\Delta := \sqrt{eps} \cdot \max\{\sqrt{|x_j|}, \sqrt[4]{eps}\}$ besser verhält.)

praktisch : Wähle Δ so , dass Δ und x_j gleiches Vorzeichen haben.

Alternativen : - Differenzenquotienten 2. Ordnung :

$$\frac{F(x + \Delta \cdot e_j) - F(x - \Delta \cdot e_j)}{2\Delta}$$

(geringerer Fehler , jedoch Verdopplung des Aufwands)

- Extrapolationsverfahren :

Hiermit ist prinzipiell eine sehr genaue Berechnung der Ableitungen möglich , für das Newton-Verfahren jedoch zu aufwendig.

Vorteil : - sehr einfache Implementierung

Nachteile : - hoher Aufwand für $n \gg 1$
- Sehr großer Fehler , falls F nicht genau ausgewertet werden kann.

typisches Beispiel : Auswertung von F beinhaltet selbst ein iteratives Verfahren (nested iterations)

d) Approximation der Jacobimatrix :

Für große Probleme kann es vorteilhaft sein , betragsmäßig sehr kleine Elemente der Jacobimatrix zu vernachlässigen.

\Rightarrow Reduktion der Anzahl der Nicht-Null-Elemente

e) Matrixfreie Algorithmen :

Löst man im Newton-Verfahren die linearen Gleichungssysteme $J \cdot p_k = -F(x_k)$ nicht exakt , sondern iterativ , so kann im CG-Verfahren (und entsprechenden Verfahren) für nichtsymmetrische Matrizen die Berechnung von J vermieden werden , weil nicht J selbst , sondern ausschließlich Matrix-Vektor-Produkte $J \cdot v$ auszuwerten sind :

$$J \cdot v = \frac{F(x + \Delta \cdot v) - F(x)}{\Delta}$$

Zusammenfassung :

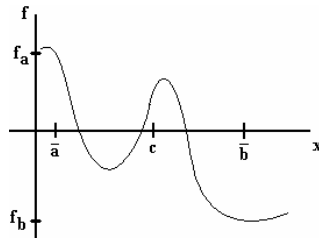
- Newton-Verfahren kann divergieren.
- Gedämpftes Newton-Verfahren konvergiert nur nahe der Nullstelle ($I \approx 1$) quadratisch , sonst linear.
- Vereinfachtes Newton-Verfahren : für hochdimensionale Funktionen
Da sich in der Nähe der Nullstelle die Ableitung kaum mehr ändert , braucht man die Jacobimatrix/Ableitungen nicht in jedem Schritt neu berechnen.

9.3 Nullstellen reeller Funktionen

Bem 9.15 : Bisektionsverfahren

geg. : $f \in C[\bar{a}, \bar{b}]$ mit Vorzeichenwechsel in $[\bar{a}, \bar{b}] \Rightarrow f(\bar{a}) \cdot f(\bar{b}) < 0$

ges. : $x^* \in (\bar{a}, \bar{b})$ mit $f(x^*) = 0$, Existenz von x^* folgt aus Zwischenwertsatz



Algorithmus :

Initialisierung : $a := \bar{a}$, $b := \bar{b}$, $f_a := f(a)$, $f_b := f(b)$

repeat

$$c := \frac{a+b}{2} \quad , \quad f_c := f(c)$$

if $f_a \cdot f_c < 0$ **then** $b := c$, $f_b := f_c$
else $a := c$, $f_a := f_c$

until $(b - a \leq TOL)$ **OR** $(f_c = 0)$

$$x^* \approx x_{bisec} := \begin{cases} c & \text{falls } f_c = 0 \\ \frac{a+b}{2} & \text{sonst} \end{cases}$$

Eigenschaften :

- einfache Implementierung
- stets konvergent (Bei Rechnung in exakter Arithmetik gilt : $f(x^*) = 0$ für ein $x^* \in [a, b]$.)
- In Gleitpunktarithmetik muss $[a, b]$ keine Einschließung der Nullstelle ergeben (Alternative : Intervallarithmetik)
- Anzahl der Iterationsschritte : $k \geq 1 + \log_2 \left(\frac{\bar{b} - \bar{a}}{TOL} \right)$

Bem 9.16 : Regula falsi

Idee : Einschließung der Nullstelle wie im Bisektionsverfahren , aber Berücksichtigung des Funktionsverlaufs in $[a, b]$.

Algorithmus :

Bestimme c als Nullstelle des (linearen) Interpolationspolynoms zu (a, f_a) , (b, f_b) :

$$f_a + \frac{c-a}{b-a} \cdot (f_b - f_a) = 0$$

\Rightarrow

$$c = \frac{a \cdot f_b - b \cdot f_a}{f_b - f_a}$$

Konvergenz :

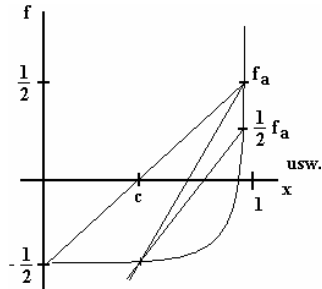
- Häufig schneller als im Bisektionsverfahren.

- In exakter Arithmetik gilt : Intervall $[a,b]$ ist eine Einschließung der Nullstelle x^* von f .

Problem:

Langsame Konvergenz , wenn während der Iteration eines der beiden Intervallenden unverändert bleibt.

Beispiel : $f(x) = x^{10} - \frac{1}{2}$, $x^* = \sqrt[10]{\frac{1}{2}} \approx 1$, $[\bar{a}, \bar{b}] = [0,1]$



Verbesserung:

Bleibt a über 2 Iterationsschritte unverändert , so setze $f_a := \frac{1}{2} f_a$. (analog für b)

verbesserter Algorithmus:

Initialisierung : $a := \bar{a}$, $b := \bar{b}$, $f_a := f(a)$, $f_b := f(b)$, $i := 0$

repeat

$$c = \frac{a \cdot f_b - b \cdot f_a}{f_b - f_a} , \quad f_c := f(c)$$

if $f_a \cdot f_c < 0$ **then**

$$b := c , \quad f_b := f_c$$

if $i = -1$ **then** $f_a := \frac{1}{2} f_a$

$$i := -1$$

else

$$a := c , \quad f_a := f_c$$

if $i = 1$ **then** $f_b := \frac{1}{2} f_b$

$$i := 1$$

until $(b - a \leq TOL)$ **OR** $(f_c = 0)$

$$x^* \approx x_{reg} := \begin{cases} c , & \text{falls } f_c = 0 \\ \frac{a+b}{2} & \text{sonst} \end{cases}$$

Bem 9.17 : Sekantenverfahren

Idee : Nutze (ähnlich wie im Newton-Verfahren) Informationen über f in einer (kleinen) Umgebung von x_k , vermeide jedoch die Auswertung von f' .

Vorteil : Schnellere Konvergenz als regula falsi, falls „gute“ Startwerte bekannt sind.

Nachteil : Keine Einschließung der Nullstelle.

Verfahren :

Bestimme zu gegebenen $(x_{k-1}, f_{k-1}), (x_k, f_k)$, (mit $f_k = f(x_k)$) das lineare Interpolationspolynom und wähle x_{k+1} als dessen Nullstelle :

$$f_k + \frac{x_{k+1} - x_k}{x_{k-1} - x_k} \cdot (f_{k-1} - f_k) = 0 \quad (*)$$

\Rightarrow

$$x_{k+1} := \frac{x_{k-1} \cdot f_k - x_k \cdot f_{k-1}}{f_k - f_{k-1}}$$

Alternative Interpretation : Inverse Interpolation

Bestimme das Interpolationspolynom p zu den Stützpunkten (f_{k-i}, x_{k-i}) , ($i = 0, 1, \dots, l$) und wähle $x_{k+1} := p(0)$.

Die Existenz von p lässt sich nachweisen für paarweise voneinander verschiedene x_{k-i} mit $|x_{k-i} - x^*| \ll 1$, ($i = 0, 1, \dots, l$), sofern x^* einfache Nullstelle von f ist.

praktische : $l = 2$ oder 3 oder 4

Konvergenz :

$$x_{k+1} = x_k - \left(f_k \cdot \frac{x_{k-1} - x_k}{f_{k-1} - f_k} \right) \quad (\text{nach Umstellen von } (*))$$

$$= x_k - (f_k - f(x^*)) \cdot \frac{1}{\frac{f_{k-1} - f_k}{x_{k-1} - x_k}}$$

$$= x_k - f[x_k, x^*] \cdot (x_k - x^*)$$

$$= \frac{1}{f[x_{k-1}, x_k]}$$

\Rightarrow

$$x_{k+1} - x^* = (x_k - x^*) \cdot \left(1 - \frac{f[x_k, x^*]}{f[x_{k-1}, x_k]} \right)$$

$$= (x_k - x^*) \cdot (x_{k-1} - x^*) \cdot \frac{f[x_{k-1}, x_k, x^*]}{f[x_{k-1}, x_k]}$$

Ist x^* einfache Nullstelle von $f \in C^2$ und x_{k-1}, x_k in einer hinreichend kleinen Umgebung von x^* , so

gibt es eine Konstante $M > 0$ mit $\left| \frac{f[x_{k-1}, x_k, x^*]}{f[x_{k-1}, x_k]} \right| \leq M$,

denn $f[x_{k-1}, x_k] = f'(x_{k-1} + \mathbf{J} \cdot (x_l - x_{k-1}))$ für ein $\mathbf{J} \in (0,1)$

und $f[x_{k-1}, x_k, x^*] = \frac{1}{2} f''(\mathbf{x})$ für ein $\mathbf{x} \in [\mathbf{a}, \mathbf{b}]$ mit $x_{k-1}, x_k, \mathbf{x} \in [\mathbf{a}, \mathbf{b}]$.

Mit $e_k := M \cdot |x_k - x^*|$ folgt $e_{k+1} \leq e_k \cdot e_{k-1}$ und hieraus mittels vollständiger

Induktion $e_k \leq K^{q^k}$, ($k \geq 0$) mit $K := \max\{e_0, \sqrt[q]{e_1}\}$ und $q = \frac{1}{2}(1 + \sqrt{5})$, denn $q^2 = q + 1$

und $K^{q^k} \cdot K^{q^{k-1}} = K^{q^k + q^{k-1}} = K^{q^{k-1} \cdot (q+1)} = K^{q^{k-1} \cdot q^2} = K^{q^{k+1}}$.

Wählt man Startwerte x_0, x_1 so, dass $K < 1$, so folgt die Konvergenz des Sekantenverfahrens.

Die Ordnung beträgt mindestens $q = \frac{1}{2}(1 + \sqrt{5}) \approx 1.618$

Bem 9.18 : Nullstellen von Polynomen

Enger Zusammenhang zwischen Polynomnullstellen und Eigenwerten von Matrizen.
(Die Eigenwertberechnung ist numerisch oft günstiger.)

Newton-Verfahren zur Berechnung einer Nullstelle von $p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0$

Auswertung von $p(x_k), p'(x_k)$ mittels Horner-Schema

Zur Wahl der Startwerte vgl. Literatur *Stoer*, Abschnitt 5.5.

Hat man eine Nullstelle \mathbf{x}_1 (näherungsweise) bestimmt, so wendet man zur Bestimmung weiterer Nullstellen das

Newton-Verfahren auf $p_1(x) := \frac{p(x)}{x - \mathbf{x}_1}$ an. (*)

Numerisch stabile Auswertung von $\frac{p_1(x_k)}{p_1'(x_k)}$ mittels Trick von Maehly :

$$x_{k+1} = x_k - \frac{p_1(x_k)}{p_1'(x_k)} = x_k - \frac{p(x_k)}{p'(x_k) - \frac{p(x_k)}{x_k - \mathbf{x}_1}}$$

- wegen der Ableitung von (*): $p_1'(x) = \frac{p'(x) \cdot (x - \mathbf{x}_1) - p(x)}{(x - \mathbf{x}_1)^2}$

9.4 : Minimierungsprobleme

Bem 9.19 : Newton-Verfahren für Minimierungsprobleme

geg. : $f : \mathbb{R}^n \rightarrow \mathbb{R}$, f stetig differenzierbar

ges. : x^* mit $f(x^*) = \min_{x \in \mathbb{R}^n} f(x)$ (= freies Minimierungsproblem)

Notwendige Bedingung für lokales Extremum : $\nabla f(x^*) = 0$

Hinweis : Gradient von f : $\nabla f(x) = \begin{pmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_n}(x) \end{pmatrix}$

Ist x^* einzige Nullstelle von ∇f und wächst $f(x)$ für $\|x\| \rightarrow \infty$ unbeschränkt, d.h. gibt es zu jedem $C > 0$ ein $M > 0$ mit $f(x) \geq C$ für alle $x \in \mathbb{R}^n$ mit $\|x\| \geq M$, so ist $\nabla f(x^*) = 0$ notwendig und hinreichend für eine globale Minimalstelle von f .

Berechnung von x^* durch Anwendung des (gedämpften) Newton-Verfahrens auf $F(x) = 0$ mit $F(x) := \nabla f(x)$:

$$x_{k+1} = x_k - \mathbf{I}_k \cdot (\nabla^2 f(x_k))^{-1} \cdot \nabla f(x_k)$$

mit $\nabla^2 f(x) = \text{Hessematrix} = \left(\frac{\partial^2 f}{\partial x_i \partial x_j}(x) \right)_{i,j=1}^n \in \mathbb{R}^{n \times n}$.

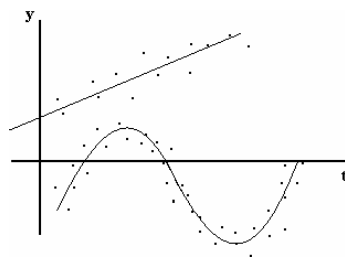
Bestimmung des Dämpfungsparameters \mathbf{I}_k so, dass $f(x_{k+1}) < f(x_k)$.

Praktisch sind häufig beschränkte Minimierungsprobleme anzutreffen.

D.h. Menge $E \subset \mathbb{R}^n$ und gesucht : $x^* \in E$ mit $f(x^*) = \min \{f(x) : x \in E\}$

Bem. 9.20 : Nichtlineare Ausgleichsprobleme

geg. : Messdaten (t_k, y_k) , $(k = 1, \dots, m)$ denen ein (angenommener) funktionaler Zusammenhang $y = \mathbf{j}(t; x_1, x_2, \dots, x_n)$ zu Grunde liegt.



linear : $y = a \cdot t + b$

nichtlinear : $y = a \cdot \sin(\mathbf{w} \cdot t + \mathbf{j}) + b$

ges. : Parameter x_1, x_2, \dots, x_n so, dass $\mathbf{j}(t_k; x_1, x_2, \dots, x_n)$ möglichst gut y_k approximiert :

$$\|F(x)\|_2 \rightarrow \min \quad \text{mit } F(x_1, x_2, \dots, x_n) := \begin{pmatrix} y_1 - \mathbf{j}(t_1; x_1, x_2, \dots, x_n) \\ \vdots \\ y_m - \mathbf{j}(t_m; x_1, x_2, \dots, x_n) \end{pmatrix}$$

„Ausgleichsrechnung“

Spezialfall :

Ist \mathbf{j} linear in x_1, x_2, \dots, x_n , so ergibt sich x_1, x_2, \dots, x_n als Lösung eines überbestimmten linearen Gleichungssystems. Numerische Lösung mittels QR-Zerlegung.

allgemein :

Mit $f(x) := \frac{1}{2} \|F(x)\|_2^2$ ergibt sich in Bem 9.19

$$\begin{aligned} \nabla f(x) &= (F_x(x))^T \cdot F(x) \quad \text{und} \\ \nabla^2 f(x) &= (F_x(x))^T \cdot F_x(x) + (F_{xx}(x))^T \cdot F(x). \end{aligned}$$

Hat $F_x(x)$ Vollrang und ist $\|F(x^*)\| \ll 1$, so ist $(F_x(x))^T \cdot F_x(x)$ symmetrisch und positiv definit und $\nabla^2 f(x)$ regulär in einer Umgebung von $x^* \Rightarrow$ Newton-Verfahren anwendbar

Problem : Auswertung von $F_{xx}(x)$ kann sehr aufwendig sein

Bem. 9.12 : Gauß-Newton-Verfahren

geg. : $F : R^n \rightarrow R^m$, F stetig differenzierbar

ges. : $x^* \in R^n$ mit $\|F(x^*)\|_2 = \min_x \|F(x)\|_2$

Ansatz :

Linearisierung von F in x_k

$$F(x) \approx F_k(x) := F(x_k) + F_x(x_k) \cdot (x - x_k)$$

Bestimmung von \bar{x} als Lösung des linearen Ausgleichsproblems

$$\begin{aligned} \min_x \frac{1}{2} \|F_k(x)\|_2^2 &= \min_x \frac{1}{2} \|F(x_k) + F_x(x_k) \cdot (x - x_k)\|_2^2 \\ &= \min_x \frac{1}{2} \|Ax - b\|_2^2 \end{aligned}$$

mit $A := F_x(x_k) \in R^{m \times n}$, $b := F_x(x_k) \cdot x_k - F(x_k)$

gedämpftes Newton-Verfahren :

$$p_k := \bar{x} - x_k, \quad x_{k+1} = x_k + \mathbf{I}_k p_k$$

Normalgleichungen :

$$(F_x(x_k))^T \cdot F_x(x_k) \cdot (\bar{x} - x_k) = -(F_x(x_k))^T \cdot F(x_k)$$

Im Vergleich zum Newton-Verfahren entfällt die aufwendige Auswertung von $(F_{xx}(x_k))^T \cdot F(x_k)$

in $\nabla^2 f(x_k)$.

Ist $\|F(x_k)\| \ll 1$, so ist die Konvergenz jedoch ähnlich schnell wie im Newton-Verfahren

für $f(x) := \frac{1}{2} \|F(x)\|_2^2$.

Häufig Konvergenzprobleme, falls $\|F(x^*)\| = O(1)$.

Praktisch ergibt sich für derartige Probleme jedoch ohnehin eine sehr schlechte Übereinstimmung von y_k und $\mathbf{j}(t_k; x_1, x_2, \dots, x_n)$.

10. Numerische Lösung von Differentialgleichungen

Matlab-Befehl zum Lösen von DGL : ode45

Bem 10.1 : Systeme gewöhnlicher DGL

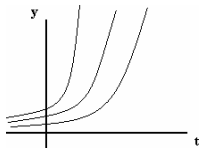
geg. : $f : [0, T] \times \mathbb{R}^{n_y} \rightarrow \mathbb{R}^{n_y}$ System gew. DGL 1. Ordnung - $y'(t) = f(t, y(t))$ (*)

ges. : Lösung $y : [0, T] \rightarrow \mathbb{R}^{n_y}$, $y \in C^1$, und für jedes $t \in [0, T]$ soll (*) erfüllt sein.

Beispiel :

$$y' = y, \quad n_y = 1$$

$$\text{D.h. } y(t) = c \cdot e^t, \quad c \in \mathbb{R}$$



Spezialfall :

$$f(t, y) = f(t)$$

$$y(t) = y(0) + \int_0^t f(t) dt$$

(Die Menge der Quadraturprobleme (Integration) ist Untermenge der Anfangswertprobleme von DGL.)

klassisch : $y(t_n)$ für „viele“ $t_n \in [0, T]$

besser : „stetige Lösungsdarstellung“ (dense output)
 $\tilde{y}(t) \in C[0, T]$ mit $\tilde{y}(t) \approx y(t)$, ($t \in [0, T]$)

Systeme höherer Ordnung : Äquivalenz zu Systemen 1. Ordnung (siehe Analysis III)

Nichtautonome Systeme :

Häufig betrachtet man nur autonom Systeme $y'(t) = f(t, y(t))$, denn diese sind äquivalent zu $Y'(t) = F(Y(t))$ mit $Y(t) = \begin{pmatrix} t \\ y(t) \end{pmatrix} \in \mathbb{R}^{n_y+1}$, $F(Y) = \begin{pmatrix} 1 \\ f(t, y) \end{pmatrix}$

- Anfangswertprobleme:

Vorgabe von n Anfangsbedingungen $y(0) = y_0 \in \mathbb{R}^{n_y}$ zu (*)

Satz von Picard-Lindelöf:

Ist f stetig und Lipschitz-stetig bzgl. y , so hat das AWP zu (*) eine eindeutige Lösung.

$$y^{(k+1)}(t) = y_0 + \int_0^t f(\tau, y^{(k)}(\tau)) d\tau, (k \geq 0), \quad y^{(0)}(t) \equiv y_0, (t \in [0, T])$$

praktisch: sehr langsame Konvergenz \Rightarrow für praktische Rechnungen ungeeignet (jedoch gut parallelisierbar)

- Randwertprobleme:

Vorgabe von insgesamt n_y Bedingungen an die Randwerte $y(0)$ und $y(T)$:

$$r(y(0), y(T)) = 0 \quad \text{mit } r: \mathbb{R}^{n_y} \times \mathbb{R}^{n_y} \rightarrow \mathbb{R}^{n_y}$$

Beispiel:

$$x'' + x = 0 \quad (*)$$

$$y = \begin{pmatrix} x \\ x' \end{pmatrix}, \quad y_1' = y_2, \quad y_2' = -y_1$$

$$y(t) = \begin{pmatrix} y_1(t) \\ y_2(t) \end{pmatrix}, \quad f(t, y) = \begin{pmatrix} y_2 \\ -y_1 \end{pmatrix}$$

AWP:

$$x(0) = y_1(0) = x_0$$

$$x'(0) = y_2(0) = x_0'$$

RWP:

$$[0, T] = [0, \mathbf{p}]$$

Lösung von (*): $x(t) = c_1 \sin t + c_2 \cos t$ mit $c_1, c_2 \in \mathbb{R}$

a) $x(0) = x(\mathbf{p}) = 0 = y_1(0) = y_1(\mathbf{p})$
 $\Rightarrow c_2 = 0$

D.h. $x(t) = c_1 \sin t$ mit $c_1 \in \mathbb{R}$

b) $x(0) = y_1(0) = 0, \quad x(\mathbf{p}) = y_1(\mathbf{p}) = 1$

Aus $x(0) = 0$ folgt $c_2 = 0$

Aus $x(\mathbf{p}) = 1$ folgt $c_2 = -1$

D.h. widersprüchliche Bedingungen \Rightarrow nicht lösbar

Im Allgemeinen nur lokale Existenz- und Eindeutigkeitsaussagen in einer Umgebung einer (isolierten) Lösung.

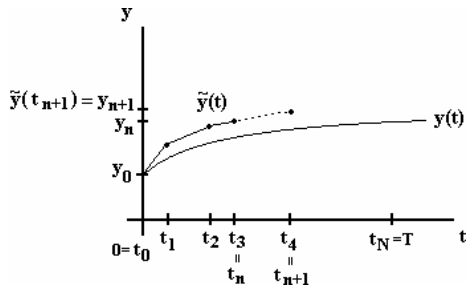
Bem 10.2 : Diskretisierungsverfahren für gewöhnliche DGL

Def. Zeitgitter: $\Delta := \{t_0, t_1, \dots, t_N : 0 = t_0 < t_1 < \dots < t_N = T\}$

Explizites Eulerverfahren :

Sei Näherungslösung $\tilde{y}(t)$ gegeben auf $[0, t_n]$ und $y_n := \tilde{y}(t_n)$.

Gesucht : Stetige Fortsetzung auf $[0, t_{n+1}]$ mit $t_{n+1} - t_n =: h_n > 0$.



$$t \in [t_n, t_{n+1}]: \quad y(t) = y(t_n) + \int_{t_n}^t y'(t) dt \stackrel{(*)}{=} y(t_n) + \int_{t_n}^t f(t, y(t)) dt$$

$$\approx \tilde{y}(t_n) + \int_{t_n}^t f(t_n, \tilde{y}(t_n)) dt = y_n + (t - t_n) \cdot f(t_n, y_n)$$

$$y(t_{n+1}) \approx \tilde{y}(t_{n+1}) := y_{n+1} := y_n + \underbrace{h_n}_{=t_{n+1}-t_n} \cdot f(t_n, y_n)$$

Verfahrensvorschrift des expliziten Eulerverfahrens :

$$\frac{y_{n+1} - y_n}{h_n} = f(t_n, y_n)$$

Stetige Lösungsdarstellung (Polygonzugverfahren) :

$$y(t_n + \mathbf{q} \cdot h_n) \approx \tilde{y}(t_n + \mathbf{q} \cdot h_n) := y_n + \mathbf{q} \cdot h_n \cdot f(t_n, y_n) = y_n + \mathbf{q} \cdot (y_{n+1} - y_n) \quad , \quad (\mathbf{q} \in [0,1])$$

Trapezregel :

$$\tilde{y}(t_{n+1}) = \tilde{y}(t_n) + \frac{h_n}{2} \cdot (f(t_n, \tilde{y}(t_n)) + f(t_{n+1}, \tilde{y}(t_{n+1})))$$

$$\frac{y_{n+1} - y_n}{h_n} = \frac{1}{2} \cdot (f(t_n, y_n) + f(t_{n+1}, y_{n+1}))$$

- Nichtlineares Gleichungssystem zur Bestimmung von y_{n+1}
- Zur Lösung von Anfangswertproblemen nur geeignet , falls h_n wesentlich größer sein kann als im Eulerverfahren.
- Vorteilhaft für RWP , da symmetrisches Verfahren

Definition 10.3 : Einschrittverfahren (ESV)

Zur numerischen Lösung des AWP $y(0) = y_0$ zu (*) wird die Lösung $y(t)$ approximiert durch eine Gitterfunktion $y_\Delta : \Delta \rightarrow \mathbb{R}^{n_y}$ mit $y_n := y_\Delta(t_n) = \tilde{y}(t_n) \approx y(t_n)$.

Einschrittverfahren haben die Form $y_{n+1} = y_n + h_n \cdot \mathbf{j}(t_n, y_n; h_n, f)$ mit der Verfahrensfunktion \mathbf{j} .

Euler : $\mathbf{j}(t_n, y_n; h_n, f) = f(t_n, y_n)$

Trapezregel : \mathbf{j} ergibt sich als Lösung eines nichtlinearen Gleichungssystems

Bem. 10.4 : Explizite Runge-Kutta-Verfahren

S-stufiges explizites Runge-Kutta-Verfahren

$$y_{n+1} = y_n + h_n \cdot \sum_{j=1}^S b_j \cdot f(t_n + c_j \cdot h_n, Y_{nj})$$

mit Stufenvektoren $Y_{ni} = y_n + h_n \cdot \sum_{j=1}^{i-1} a_{ij} \cdot f(t_n + c_j \cdot h_n, Y_{nj})$, $(i = 1, \dots, S)$

$$Y_{n1} = y_n$$

$$Y_{n1}' := f(t_n + c_1 h_n, Y_{n1})$$

$$Y_{n2} = y_n + h_n \cdot a_{21} \cdot f(t_n + c_1 h_n, Y_{n1}) = y_n + h_n \cdot a_{21} \cdot Y_{n1}'$$

$$Y_{n2}' := f(t_n + c_2 h_n, Y_{n2})$$

...

$$y_{n+1} = y_n + h_n \cdot \sum_{j=1}^S b_j \cdot Y_{nj}'$$

mit :

b_j - Gewichte

c_i - Knoten

a_{ij} - Runge-Kutta-Parameter

Butcher-Schema :

c_2	a_{21}				
c_3	a_{31}	a_{32}			
\vdots	\vdots		\ddots		
c_S	a_{S1}	a_{S2}	\dots	$a_{S,S-1}$	
	b_1	b_2	\dots	b_{S-1}	b_S

Spezialfall :

$$y'(t) = 1$$

$$Y_{ni} = y_n + h_n \cdot \sum_{j=1}^{i-1} a_{ij} \cdot 1$$

Stufenvektoren $Y_{ni} \approx y(t_n + c_i h_n)$

$$y(t_n + c_i h_n) = y(t_n) + c_i h_n$$

$$c_i = \sum_{j=1}^{i-1} a_{ij}, \quad (i = 1, \dots, S)$$

stetige Lösungsdarstellung :

$$y(t_n + \mathbf{q} \cdot h_n) \approx y_n + h_n \cdot \sum_{j=1}^S b_j(\mathbf{q}) \cdot Y_{nj}$$

mit $b_j(0) = 0$ und $b_j(1) = b_j$ für $j = 1, \dots, S$

Beispiel :

Arensdorf-Orbit (siehe evt. Folien)

- a) Explizites Eulerverfahren $h = 5.0_{E-4}$ sehr schlechte Näherung (unbrauchbar)
 $4 \cdot 10^4$ Integrationssschritte und $4 \cdot 10^4$ Funktionsauswertungen
- b) Explizites Eulerverfahren $h = 2.0_{E-5}$ akzeptable Näherung
 10^6 Integrationssschritte und 10^6 Funktionsauswertungen
- c) Runge-Kutta $h = 2.0_{E-3}$ akzeptable Näherung
 $4 \cdot 10^4$ Integrationssschritte und $4 \cdot 10^4$ Funktionsauswertungen
- d) DOPRI5 (Dormand/Prince) sehr gute Näherung
75 Integrationssschritte und 535 Funktionsauswertungen

zum klassischen Runge-Kutta-Verfahren :

$S = 4$ (d.h. 4stufig)

$$Y_{n1} = y_n$$

$$Y_{n2} = y_n + \frac{h_n}{2} \cdot f(t_n, y_n)$$

$$Y_{n3} = y_n + \frac{h_n}{2} \cdot f\left(t_n + \frac{h_n}{2}, Y_{n2}\right)$$

$$Y_{n4} = y_n + \frac{h_n}{2} \cdot f\left(t_n + \frac{h_n}{2}, Y_{n3}\right)$$

$$y_{n+1} = y_n + \frac{h_n}{6} \cdot \left(f(t_n, y_n) + 2 \cdot f\left(t_n + \frac{h_n}{2}, Y_{n2}\right) + 2 \cdot f\left(t_n + \frac{h_n}{2}, Y_{n3}\right) + f\left(t_n + \frac{h_n}{2}, Y_{n4}\right) \right)$$

Butcher-Schema :

$1/2$	$1/2$			
$1/2$	0	$1/2$		
1	0	10	1	

	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$
--	---------------	---------------	---------------	---------------

Bem 10.5 : Diskretisierungsfehler

a) **Lokaler Fehler des ESV (local error) :**

$$y(t_{n+1}) - (y(t_n) + h_n \cdot \mathbf{j}(t_n, y(t_n); h_n, f))$$

Beispiel : Explizites Euler-Verfahren

$$y_{n+1} = y_n + h_n \cdot f(t_n, y_n)$$

$$(le)_n := y(t_{n+1}) - \left(y(t_n) + h_n \cdot \underbrace{f(t_n, y(t_n))}_{=y'(t_n)} \right)$$

$$y(t_{n+1}) = y(t_n + h_n) = y(t_n) + h_n \cdot y'(t_n) + \mathcal{O}(h_n^2)$$

$$\Rightarrow (le)_n = \mathcal{O}(h_n^2)$$

Runge-Kutta : $(le)_n = \mathcal{O}(h_n^5)$

DOPRI5 : $(le)_n = \mathcal{O}(h_n^6)$

b) **Globaler Fehler :**

$$\mathbf{e}_N := \|y_N - y(t_N)\|$$

$$\|y_{n+1} - y(t_{n+1})\| = \left\| \begin{aligned} &y_n + h_n \cdot \mathbf{j}(t_n, y_n; h_n, f) - (y(t_n) + h_n \cdot \mathbf{j}(t_n, y(t_n); h_n, f)) \\ &+ \underbrace{(y(t_n) + h_n \cdot \mathbf{j}(t_n, y(t_n); h_n, f)) - y(t_{n+1})}_{=-(le)_n} \end{aligned} \right\|$$

wegen Dreiecksungleichung folgt

$$\leq \|y_n + h_n \cdot \mathbf{j}(t_n, y_n; h_n, f) - (y(t_n) + h_n \cdot \mathbf{j}(t_n, y(t_n); h_n, f))\| + \|(le)_n\|$$

$$\leq \|y_n - y(t_n)\| + h_n \cdot \|\mathbf{j}(t_n, y_n; h_n, f) - \mathbf{j}(t_n, y(t_n); h_n, f)\| + \|(le)_n\|$$

$$\leq \|y_n - y(t_n)\| + h_n \cdot L \cdot \|y_n - y(t_n)\| + \|(le)_n\| \quad (\text{Lipschitz-Konstante } L)$$

$$= (1 + h_n \cdot L) \cdot \|y_n - y(t_n)\| + \|(le)_n\|$$

mit $\mathbf{e}_n := \|y_n - y(t_n)\|$ folgt die Rekursion :

$$\mathbf{e}_{n+1} \leq (1 + h_n \cdot L) \cdot \mathbf{e}_n + \|(le)_n\|$$

$$\leq e^{h_n \cdot L} \cdot \mathbf{e}_n + \|(le)_n\|$$

$$\mathbf{e}_n \leq e^{h_{n-1} \cdot L} \cdot \mathbf{e}_{n-1} + \|(le)_{n-1}\|$$

Einsetzen :

\Rightarrow

$$\begin{aligned} \mathbf{e}_{n+1} &\leq e^{h_n \cdot L} \cdot e^{h_{n-1} \cdot L} \cdot \mathbf{e}_{n-1} + \|(le)_n\| + e^{h_n \cdot L} + \|(le)_{n-1}\| \\ &\leq \dots \\ &\leq e^{(h_n + h_{n-1} + \dots + h_0) \cdot L} \cdot \mathbf{e}_0 + \|(le)_n\| + e^{h_n \cdot L} + \|(le)_{n-1}\| \\ &\quad + e^{(h_n + h_{n-1}) \cdot L} \cdot \|(le)_{n-2}\| + \dots + e^{(h_n + h_{n-1} + \dots + h_1) \cdot L} \cdot \|(le)_0\| \end{aligned}$$

wegen $h_0 + h_1 + \dots + h_n = t_{n+1} - t_0 \leq T$ folgt

$$\mathbf{e}_{n+1} \leq e^{LT} \cdot \sum_{i=0}^n \|(le)_i\|$$

Für $h_n = h = \text{const}$ ist $\mathbf{e}_{n+1} \leq e^{LT} \cdot \sum_{i=0}^n C \cdot h^2 = e^{LT} \cdot n \cdot C \cdot h \cdot h$ und

wegen $n \cdot h < T$ ist $\mathbf{e}_{n+1} \leq C \cdot e^{LT} \cdot T \cdot h$.

(=Beschreibung des globalen Fehlers des Euler-Verfahrens.)

D.h. Halbierung der Schrittweite entspricht einer Halbierung des Fehlers.

Für Runge-Kutta : $\mathbf{e}_{n+1} \leq C \cdot e^{LT} \cdot T \cdot h^4$

Für DOPRI5 : $\mathbf{e}_{n+1} \leq C \cdot e^{LT} \cdot T \cdot h^5$ (Halbierung der Schrittweite bedeutet $\frac{1}{32}$ Fehler)