

---

# Approximation

Roland Pulch

Skript zur Vorlesung im Sommersemester 2024

Institut für Mathematik und Informatik  
Universität Greifswald

Inhalt:

1. Approximation mit Polynomen und Splines
2. Approximation in normierten Räumen
3. Parameterbestimmung

Literatur:

H.R. Schwarz, N. Köckler: Numerische Mathematik. (8. Aufl.) Vieweg+Teubner 2011. (Kapitel 3)

R. Schaback, H. Wendland: Numerische Mathematik. (5. Aufl.) Springer 2005. (Kapitel 12)

R.A. DeVore, G.G. Lorentz: Constructive Approximation. Springer 1993.

O. Christensen, K.L. Christensen: Approximation Theory: From Taylor Polynomials to Wavelets. Birkhäuser 2005.

---

# Inhaltsverzeichnis

<b>1</b>	<b>Approximation mit Polynomen und Splines</b>	<b>3</b>
1.1	Interpolation mit Polynomen . . . . .	3
1.2	Approximation mit Polynomen . . . . .	7
1.3	Interpolation mit Splines . . . . .	13
1.4	Ausgleichsspline . . . . .	21
<b>2</b>	<b>Approximation in normierten Räumen</b>	<b>34</b>
2.1	Allgemeine Approximationstheorie . . . . .	34
2.2	Fourier-Reihen . . . . .	51
	<b>Literaturverzeichnis</b>	<b>59</b>

# 1 Approximation mit Polynomen und Splines

Wir betrachten das folgende Problem. Dabei bezeichnet  $C[a, b]$  die Menge aller stetigen Funktionen  $f : [a, b] \rightarrow \mathbb{R}$ , wobei stets  $a < b$  vorausgesetzt wird.

**Aufgabenstellung:** Gegeben sei eine Funktion  $f \in C[a, b]$ . Aus einer Menge  $\mathcal{G} \subset C[a, b]$  von Funktionen einfacher Gestalt soll ein  $g \in \mathcal{G}$  gefunden werden, so dass  $\|f - g\|_\infty$  klein wird.

Diese Aufgabenstellung ist allgemein formuliert, wodurch eine Approximation  $g$  noch nicht eindeutig definiert ist. Unter gewissen Voraussetzungen an die Menge  $\mathcal{G}$  existiert eine Bestapproximation  $\hat{g}$ , d.h. es gilt dann

$$\|f - \hat{g}\|_\infty \leq \|f - g\|_\infty \quad \text{für alle } g \in \mathcal{G}.$$

In diesem Kapitel wird die folgende Strategie angewendet: Es seien genau  $n + 1$  Knoten  $a \leq x_0 < x_1 < x_2 < \dots < x_{n-1} < x_n \leq b$  ausgewählt. An diesen Stellen erfolgen Auswertungen der Funktion  $f$ . Gesucht ist dann eine Funktion  $g \in \mathcal{G}$  derart, dass die Abweichungen

$$|f(x_j) - g(x_j)| \quad \text{für } j = 0, 1, \dots, n \quad (1.1)$$

möglichst klein werden. Auch hier ist die Approximation  $g$  noch nicht eindeutig bestimmt.

In manchen Anwendungen ist die Funktion  $f$  sogar unbekannt und es liegt nur eine endliche Anzahl von Funktionsauswertungen vor, z.B. aus Messungen. In diesem Fall muss das obige Konzept verwendet werden.

Eine naheliegende Idee ist, die Funktion  $f$  mit einer Funktion  $g \in \mathcal{G}$  an den Stützpunkten zu interpolieren, sofern dies möglich ist. Dann gilt entsprechend  $f(x_j) = g(x_j)$  für alle  $j = 0, 1, \dots, n$  und die Differenzen (1.1) werden alle zu null. Die Hoffnung ist, dass dann  $g$  auch nahe bei  $f$  außerhalb der Knoten liegt und somit  $\|f - g\|_\infty$  klein ausfällt.

## 1.1 Interpolation mit Polynomen

Die Interpolation kann mit Polynomen erfolgen. Wir bezeichnen die Menge aller Polynome vom Grad höchstens  $n$  mit  $\mathcal{P}_n$ . Ein Wertepaar  $(x_j, y_j)$  heißt Stützpunkt und  $x_j$  darin die Stützstelle. Es gilt der folgende Satz für das Interpolationsproblem.

**Satz 1.1** *Zu Stützpunkten  $(x_j, y_j)$  für  $j = 0, 1, \dots, n$  mit  $x_j \neq x_i$  für  $j \neq i$  existiert ein eindeutiges Polynom  $p \in \mathcal{P}_n$  mit  $p(x_j) = y_j$  für  $j = 0, 1, \dots, n$ .*

Beweis: siehe [19], Satz 2.1.1.1.

Das gesuchte Polynom  $p \in \mathcal{P}_n$  besitzt die Gestalt

$$p(x) = \alpha_0 + \alpha_1 x + \alpha_2 x^2 + \cdots + \alpha_{n-1} x^{n-1} + \alpha_n x^n$$

mit a priori unbekanntem reellen Koeffizienten  $\alpha_0, \dots, \alpha_n$ . Daraus ersieht man die Monom-Basis

$$\mathcal{P}_n = \text{span} \{1, x^1, x^2, \dots, x^{n-1}, x^n\}.$$

Für theoretische Untersuchungen bei der Polynominterpolation sind die Lagrange-Polynome

$$L_i(x) := \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} \quad \text{für } i = 0, 1, \dots, n$$

als Basis geeignet, d.h. es gilt  $\mathcal{P}_n = \text{span}\{L_0, \dots, L_n\}$ . Die Lagrange-Polynome zeichnen sich durch die Eigenschaft

$$L_i(x_j) = \begin{cases} 1 & \text{für } i = j \\ 0 & \text{für } i \neq j \end{cases}$$

aus. Es gilt  $\text{grad}(L_i) = n$  für alle  $i$ . Die Konstruktion dieser Basis liefert die einfache Darstellung

$$p(x) = \sum_{i=0}^n y_i L_i(x)$$

des Interpolationspolynoms, d.h. die Koeffizienten sind aus den Stützpunkten bekannt. Ein entscheidender Nachteil ist jedoch, dass sich bei hinzufügen einer neuen Stützstelle  $(x_{n+1}, y_{n+1})$  zu den bisherigen Stützpunkten dann alle Basispolynome ändern.

In der Praxis werden daher die Newton-Polynome eingesetzt. Sie sind definiert durch

$$N_i(x) := \prod_{j=0}^{i-1} (x - x_j) = (x - x_0)(x - x_1) \cdots (x - x_{i-1})$$

für  $i = 1, \dots, n$  und  $N_0(x) := 1$ . Es gilt  $\mathcal{P}_n = \text{span}\{N_0, \dots, N_n\}$  sowie  $\text{grad}(N_i) = i$  für alle  $i$ . Das Interpolationspolynom besitzt somit die Darstellung

$$p(x) = \sum_{i=0}^n \beta_i N_i(x).$$

Die gesuchten Koeffizienten  $\beta_i$  können mit der Methode der Dividierten Differenzen berechnet werden. Der Rechenaufwand ist dabei proportional zu  $n^2$ . Die Auswertungen des Polynoms erfolgen dann mit einer Modifikation des Horner-Schemas für die Newton-Polynome (Aufwand proportional zu  $n$ ). Alternativ kann

das Interpolationspolynom direkt ausgewertet werden mit dem Algorithmus von Aitken-Neville, wobei der Rechenaufwand wieder proportional zu  $n^2$  ausfällt. Näheres hierzu findet man in [17, 19].

Wird zu den gegebenen Stützpunkten  $(x_j, y_j)$  für  $j = 0, \dots, n$  ein weiterer Stützpunkt  $(x_{n+1}, y_{n+1})$  hinzugefügt, dann ändern sich die Koeffizienten  $\beta_0, \dots, \beta_n$  nicht und es muss nur der neue Koeffizient  $\beta_{n+1}$  berechnet werden.

Gegeben sei eine Folge  $(\Delta_m)_{m \in \mathbb{N}}$  von Mengen aus Stützstellen

$$\Delta_m = \{x_0^{(m)}, x_1^{(m)}, \dots, x_{n_m}^{(m)}\}$$

mit  $a \leq x_0^{(m)} < x_1^{(m)} < \dots < x_{n_m}^{(m)} \leq b$  für alle  $m$ . Wir definieren

$$\rho(\Delta_m) := \max \left\{ \left| x_1^{(m)} - a \right|, \left| x_2^{(m)} - x_1^{(m)} \right|, \dots, \left| x_{n_m}^{(m)} - x_{n_m-1}^{(m)} \right|, \left| b - x_{n_m}^{(m)} \right| \right\}$$

als Kennzahl für die Feinheit der Zerlegung. Gilt  $\rho(\Delta_m) \rightarrow 0$ , dann geht die Anzahl  $n_m$  der Stützpunkte notwendigerweise gegen unendlich.

Beispiele für Folgen von Stützstellen sind ( $n_m = m$ ):

i) *äquidistante Knoten*:

$$h := \frac{b-a}{m}, \quad x_j^{(m)} := a + jh \quad \text{für } j = 0, 1, \dots, m.$$

ii) *Tschebycheff-Knoten*:

$$\text{in } [-1, 1]: \quad \xi_j^{(m)} := \cos \left( \frac{2j+1}{2(m+1)} \pi \right) \quad \text{für } j = 0, 1, \dots, m,$$

$$\text{in } [a, b]: \quad x_j^{(m)} := a + \frac{b-a}{2} (\xi_j^{(m)} + 1) \quad \text{für } j = 0, 1, \dots, m.$$

Beide Folgen erfüllen  $\rho(\Delta_m) \rightarrow 0$ .

Die Frage ist nun, für welche Folgen  $(\Delta_m)_{m \in \mathbb{N}}$  die zugehörige Folge  $(p_m)_{m \in \mathbb{N}}$  der Interpolationspolynome gleichmäßig gegen  $f$  konvergiert.

**Satz 1.2 (Marcinkiewicz 1939)** *Zu jeder stetigen Funktion  $f : [a, b] \rightarrow \mathbb{R}$  existiert eine Folge  $(\Delta_m)_{m \in \mathbb{N}}$ , so dass die zugehörigen Interpolationspolynome gleichmäßig gegen  $f$  konvergieren.*

Beweis: siehe [14].

Leider liefert der Beweis dieses Satzes kein Konstruktionsverfahren für die Stützstellen. Zudem gilt das folgende negative Resultat bezüglich der Approximation mit einer fest vorgegebenen Verfeinerung.

**Satz 1.3 (Faber 1914)** *Zu jeder festen Folge  $(\Delta_m)_{m \in \mathbb{N}}$  existiert eine stetige Funktion  $f : [a, b] \rightarrow \mathbb{R}$ , so dass die zugehörigen Interpolationspolynome nicht gleichmäßig gegen  $f$  konvergieren.*

Beweis: siehe [3].

Unter stärkeren Voraussetzungen an die zu approximierende Funktion erhalten wir jedoch das gewünschte Verhalten.

**Satz 1.4** *Sei  $f : \mathbb{C} \rightarrow \mathbb{C}$  analytisch und  $f|_{[a,b]}$  reellwertig. Dann konvergieren die Interpolationspolynome zu jeder Folge  $(\Delta_m)_{m \in \mathbb{N}}$  mit  $n_m \rightarrow \infty$  gleichmäßig gegen  $f$ .*

Beweis: siehe [10], Kapitel 5, Abschnitt 4.3.

Im Fall der Tschebyscheff-Knoten läßt sich die Voraussetzung noch wesentlich abschwächen.

**Satz 1.5** *Ist  $f : [a, b] \rightarrow \mathbb{R}$  eine Lipschitz-stetige Funktion, dann konvergiert die Folge der Interpolationspolynome zu den Tschebycheff-Knoten gleichmäßig gegen  $f$ .*

Je glatter die Funktion ist, desto schneller erfolgt die Konvergenz.

**Satz 1.6** *Für  $f \in C^{k+1}[a, b]$  mit  $k > 1$  ist die Konvergenzgeschwindigkeit bei den Tschebycheff-Knoten gekennzeichnet durch  $\|f - p_m\|_\infty = o\left(\frac{\log m}{m^k}\right)$  für  $m \rightarrow \infty$ .*

Hinreichend für die Lipschitz-Stetigkeit als Voraussetzung in Satz 1.5 ist bereits  $f \in C^1[a, b]$ . In Satz 1.6 zeigt die Abschätzung  $\log_{10} m \leq m$  für alle  $m$  die obere Schranke  $o(m^{-(k-1)})$ . Somit ist das Ziel für eine wichtige Teilmenge der stetigen Funktionen bereits erreicht. Ein Nachteil dieser Konstruktion ist, dass die Tschebycheff-Knoten nicht ineinander geschachtelt sind. Eine Erhöhung der Knotenanzahl erfordert somit eine komplette Neuberechnung des Interpolationspolynoms.

**Beispiel von Runge:** Die Funktion

$$f : [-a, a] \rightarrow \mathbb{R}, \quad f(x) := \frac{1}{1+x^2}$$

wird betrachtet für  $a \geq 4$ , z.B.  $a = 5$ . Diese Funktion ist nicht analytisch in  $\mathbb{C}$ , jedoch Lipschitz-stetig in  $[-a, a]$ . Für äquidistante Knoten konvergieren die Interpolationspolynome nicht gleichmäßig gegen  $f$ , sondern das Interpolationspolynom

geht sogar gegen unendlich am Rand. Für die Tschebycheff-Knoten konvergieren die Interpolationspolynome gleichmäßig gegen  $f$ .

Es verbleibt die Frage nach einer Konstruktion der Stützstellen im Fall einer beliebigen stetigen Funktion. Als Ausblick betrachten wir das folgende rekursive Verfahren zur Konstruktion des Interpolationspolynoms.

### „Gieriger Algorithmus“:

Sei  $x_0 := \arg \max_{x \in [a,b]} |f(x)|$ .

Setze  $p_0 := f(x_0)$  und  $\Delta_0 := \{x_0\}$ .

für  $n = 1, 2, 3, \dots$

bestimme  $x_n := \arg \max_{x \in [a,b]} |f(x) - p_{n-1}(x)|$

setze  $\Delta_n := \Delta_{n-1} \cup \{x_n\}$

bilde  $p_n \in \mathcal{P}_n$  aus den Stützstellen  $\Delta_n$

Dieser Algorithmus ist „gierig“, da er die Stelle mit dem größten Fehler als neue Stützstelle verwendet und damit der Fehler an dieser Stelle im nächsten Interpolationspolynom zu null wird. Ein Vorteil ist, dass die Stützstellenmengen ineinander geschachtelt sind, d.h.  $p_{n+1}$  kann mit wenig Aufwand aus  $p_n$  bestimmt werden bei Verwendung der Newton-Basis. Das Verfahren vermeidet die Einschränkung aus dem Satz von Faber, da die Stützstellen in Abhängigkeit von  $f$  gewählt werden. Eine Konvergenzanalyse dieses Verfahrens ist noch offen.

## 1.2 Approximation mit Polynomen

In diesem Abschnitt betrachten wir sowohl die Approximation von kontinuierlichen Funktionen als auch diskreten Daten durch Polynome.

### Approximation kontinuierlicher Funktionen

Bezüglich der Approximation mit Polynomen ohne Interpolationsbedingungen gilt der allgemeine Satz.

**Satz 1.7 (Weierstraß)** *Zu jeder Funktion  $f \in C[a, b]$  existiert eine Folge  $(p_n)_{n \in \mathbb{N}}$  von Polynomen, welche gleichmäßig gegen  $f$  konvergiert.*

Satz 1.2 impliziert bereits diese Aussage. Der Beweis des Satzes 1.7 kann jedoch auch konstruktiv erfolgen über die Bernstein-Polynome.

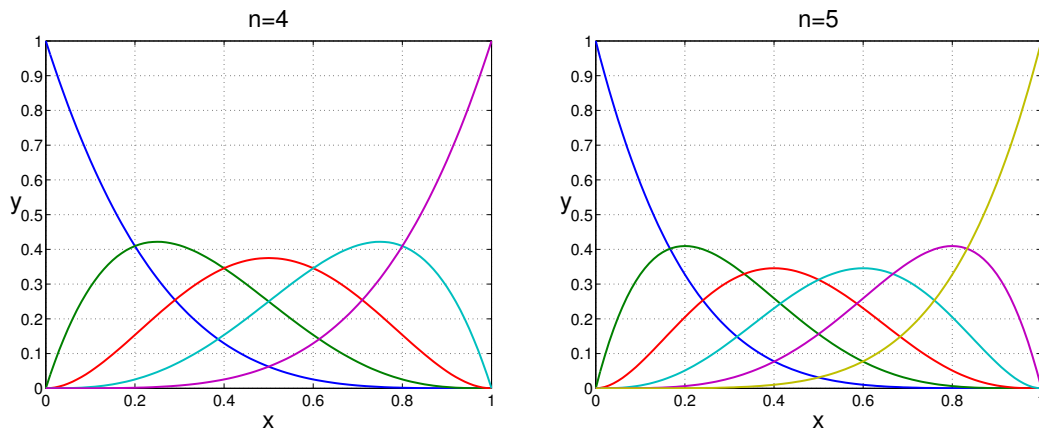


Abbildung 1: Beispiele für Bernstein-Polynome im Intervall  $[0, 1]$ .

**Def. 1.1** Für ein Intervall  $[a, b]$  lauten die Bernstein-Polynome vom Grad  $n$

$$B_{i,n}^{[a,b]}(x) := \frac{1}{(b-a)^n} \binom{n}{i} (x-a)^i (b-x)^{n-i}$$

für  $i = 0, 1, \dots, n$ .

Es gilt  $\mathcal{P}_n = \text{span}\{B_{0,n}^{[a,b]}, \dots, B_{n,n}^{[a,b]}\}$  und  $\text{grad}(B_{i,n}^{[a,b]}) = n$  für alle  $n$ .

Weitere Eigenschaften der Bernstein-Polynome:

- (i) Positivität:  $B_{i,n}^{[a,b]}(x) > 0$  für alle  $i = 0, 1, \dots, n$  und  $x \in (a, b)$ ,
- (ii) Zerlegung der Eins:

$$\sum_{i=0}^n B_{i,n}(x) = 1 \quad \text{für alle } x \in [a, b]. \quad (1.2)$$

Die Eigenschaft (i) ist offensichtlich, während die Eigenschaft (ii) aus dem Binomischen Lehrsatz folgt. Eine Folgerung aus (i) und (ii) ist  $B_{i,n}^{[a,b]}(x) \leq 1$  für alle  $i = 0, 1, \dots, n$  und  $x \in [a, b]$ . O.E.d.A. verwenden wir im folgenden das Intervall  $[0, 1]$ , wodurch sich die Gestalt der Bernstein-Polynome vereinfacht zu

$$B_{i,n}(x) = \binom{n}{i} x^i (1-x)^{n-i} \quad \text{für } i = 0, 1, \dots, n.$$

Diese Polynome lassen sich wie folgt zur Approximation einsetzen.



**Def. 1.2** Zu einer stetigen Funktion  $f : [0, 1] \rightarrow \mathbb{R}$  heißt

$$B_n(x) := \sum_{i=0}^n f\left(\frac{i}{n}\right) B_{i,n}(x) \quad (1.3)$$

das  $n$ -te Bernstein-Polynom.

Wir benötigen noch folgende Hilfsaussage.

**Lemma 1.1** Es gilt die Abschätzung

$$\sum_{i=0}^n \left(x - \frac{i}{n}\right)^2 \binom{n}{i} x^i (1-x)^{n-i} \leq \frac{1}{4n}$$

für  $0 \leq x \leq 1$ .

Beweis:

Man benötigt hierfür (1.2) und die beiden Eigenschaften

$$\begin{aligned} \sum_{i=0}^n \frac{i}{n} \binom{n}{i} x^i (1-x)^{n-i} &= x, \\ \sum_{i=0}^n \frac{i(n-i)}{n^2} \binom{n}{i} x^i (1-x)^{n-i} &= \frac{n-1}{n} x(1-x). \end{aligned}$$

Wir rechnen nach

$$\begin{aligned} \sum_{i=0}^n \frac{i}{n} \binom{n}{i} x^i (1-x)^{n-i} &= \sum_{i=1}^n \frac{i n!}{n i! (n-i)!} x^i (1-x)^{n-i} = \sum_{i=1}^n \binom{n-1}{i-1} x^i (1-x)^{n-i} \\ &= x \sum_{i=0}^{n-1} \binom{n-1}{i} x^i (1-x)^{n-1-i} = x(x + (1-x))^{n-1} = x \end{aligned}$$

und

$$\begin{aligned} \sum_{i=0}^n \frac{i(n-i)}{n^2} \binom{n}{i} x^i (1-x)^{n-i} &= \sum_{i=1}^{n-1} \frac{i(n-i)n!}{n^2 i! (n-i)!} x^i (1-x)^{n-i} \\ &= \frac{n-1}{n} x(1-x) \sum_{i=1}^{n-1} \frac{(n-2)!}{(i-1)! (n-1-i)!} x^{i-1} (1-x)^{n-i-1} \\ &= \frac{n-1}{n} x(1-x) \sum_{i=0}^{n-2} \frac{(n-2)!}{i! (n-2-i)!} x^i (1-x)^{n-i-2} \\ &= \frac{n-1}{n} x(1-x)(x + (1-x))^{n-2} = \frac{n-1}{n} x(1-x). \end{aligned}$$

Sei als Abkürzung  $R_i := \binom{n}{i} x^i (1-x)^{n-i}$ . Dadurch erhalten wir

$$\begin{aligned}
 \sum_{i=0}^n \left(x - \frac{i}{n}\right)^2 R_i &= \sum_{i=0}^n \left(x^2 - \frac{2i}{n}x + \frac{i^2}{n^2}\right) R_i \\
 &= x^2 \sum_{i=0}^n R_i - 2x \sum_{i=0}^n \frac{i}{n} R_i + \sum_{i=0}^n \frac{i^2}{n^2} R_i \\
 &= x^2 - 2x \cdot x - \sum_{i=0}^n \frac{i(n-i)}{n^2} R_i + \sum_{i=0}^n \frac{i}{n} R_i \\
 &= -x^2 - \left(1 - \frac{1}{n}\right) x(1-x) + x = \frac{1}{n} x(1-x).
 \end{aligned}$$

Man kann schließlich abschätzen  $x(1-x) \leq \frac{1}{4}$  für  $x \in [0, 1]$ . □

Es ergibt sich dann das gewünschte Verhalten.

**Satz 1.8** *Für eine stetige Funktion  $f : [0, 1] \rightarrow \mathbb{R}$  konvergiert die Folge der Bernstein-Polynome gleichmäßig gegen  $f$ .*

Beweis:

Sei  $\varepsilon > 0$ . Da  $f$  auf  $[0, 1]$  gleichmäßig stetig ist gibt es ein  $\delta > 0$ , so dass für alle  $x, y \in [0, 1]$  gilt

$$|x - y| < \delta \quad \Rightarrow \quad |f(x) - f(y)| < \varepsilon.$$

Wegen (1.2) und (1.3) gilt

$$|f(x) - B_n(x)| = |f(x) \cdot 1 - B_n(x)| \leq \sum_{i=0}^n |f(x) - f(\frac{i}{n})| \binom{n}{i} x^i (1-x)^{n-i}.$$

Wir spalten die Summe in zwei Teile auf gemäß

$$A_n := \left\{ i : 0 \leq i \leq n, \left| x - \frac{i}{n} \right| < \delta \right\}, \quad A'_n := \left\{ i : 0 \leq i \leq n, \left| x - \frac{i}{n} \right| \geq \delta \right\},$$

d.h.  $A_n \cap A'_n = \emptyset$  und  $A_n \cup A'_n = \{0, 1, \dots, n\}$ . Mit der gleichmäßigen Stetigkeit gilt

$$|f(x) - f(\frac{i}{n})| < \varepsilon \quad \text{für } i \in A_n$$

und desweiteren die grobe Abschätzung mit Dreiecksungleichung

$$|f(x) - f(\frac{i}{n})| \leq |f(x)| + |f(\frac{i}{n})| \leq 2\|f\|_\infty \quad \text{für } i \in A'_n.$$

Sei  $R_i := \binom{n}{i} x^i (1-x)^{n-i}$ . Damit können wir abschätzen unter Verwendung von Lem-

ma 1.1

$$\begin{aligned}
\sum_{i=0}^n |f(x) - f(\frac{i}{n})| R_i &= \sum_{i \in A_n} |f(x) - f(\frac{i}{n})| R_i + \sum_{i \in A'_n} |f(x) - f(\frac{i}{n})| R_i \\
&< \varepsilon \sum_{i \in A_n} R_i + 2\|f\|_\infty \sum_{i \in A'_n} R_i \\
&\leq \varepsilon \sum_{i=0}^n R_i + \frac{2\|f\|_\infty}{\delta^2} \sum_{i=0}^n (x - \frac{i}{n})^2 R_i \\
&\leq \varepsilon + \frac{\|f\|_\infty}{2n\delta^2}.
\end{aligned}$$

Obwohl  $A_n, A'_n$  von  $x$  abhängen ist die obere Schranke nun gleichmäßig für alle  $x$ . Da  $n$  beliebig war, wählen wir  $n_0$  derart, dass  $\|f\|_\infty / (2n_0\delta^2) < \varepsilon$  erfüllt ist. Für alle  $x \in [0, 1]$  und alle  $n \geq n_0$  gilt dann  $|f(x) - B_n(x)| < 2\varepsilon$  und die gleichmäßige Konvergenz ist gezeigt.  $\square$

Als Folgerung aus Satz 1.8 ergibt sich Satz 1.7. Zur Konvergenzgeschwindigkeit der Approximation gilt die asymptotische Formel.

**Satz 1.9 (Voronovskaya 1932)** *Sei die Funktion  $f : [0, 1] \rightarrow \mathbb{R}$  beschränkt und  $x_0 \in [0, 1]$ , so dass  $f$  in einer Umgebung von  $x_0$  differenzierbar ist und  $f''(x_0)$  existiert. Dann gilt*

$$\lim_{n \rightarrow \infty} n [f(x_0) - B_n(x_0)] = \frac{1}{2} x_0 (1 - x_0) f''(x_0).$$

Beweis: siehe Kapitel 10, Theorem 3.1 in [2].

Für ein stetiges  $f$  mit  $f''(x_0) \neq 0$  in einem Punkt  $x_0$  liegt damit Konvergenz in  $x_0$  mit der Konvergenzgeschwindigkeit  $\mathcal{O}(\frac{1}{n})$  vor. Die Geschwindigkeit der gleichmäßigen Konvergenz ist damit höchstens so schnell wie  $\mathcal{O}(\frac{1}{n})$ , d.h. relativ langsam im Vergleich zu beispielsweise der Splineinterpolation. Zudem bewirkt eine erhöhte Glattheit von  $f$  keine Konvergenzbeschleunigung.

Für nur Lipschitz-stetige Funktionen ergibt sich eine langsamere Konvergenz.

**Satz 1.10** *Für eine Lipschitz-stetige Funktion  $f : [0, 1] \rightarrow \mathbb{R}$  mit Lipschitz-Konstante  $L > 0$  gilt*

$$\|B_n - f\|_\infty \leq \frac{L}{2} \cdot \frac{1}{\sqrt{n}}.$$

Beweis: siehe [10], Kapitel 4, Abschnitt 2.5.

Für die Ableitungen ergibt sich jedoch ein erstaunlich starkes Resultat.

**Satz 1.11** Falls  $f \in C^k[0, 1]$ , dann konvergieren die  $i$ -ten Ableitungen der Bernstein-Approximationen  $B_n$  gleichmäßig gegen die Ableitungen  $f^{(i)}$  für jedes  $i \leq k$ .

Beweis: siehe [2], Kapitel 10, Theorem 2.1.

## Approximation diskreter Daten

Statt einer kontinuierlichen Funktion werden nun nur endlich viele Stützpunkte angenähert.

**Aufgabenstellung:** Gegeben Stützpunkte  $(x_i, y_i)$  für  $i = 0, 1, \dots, m$  mit  $x_j \neq x_i$  für  $j \neq i$ . Gesucht ist ein Polynom  $p \in \mathcal{P}_n$  mit  $n < m$ , so dass die Zielfunktion

$$J := \sum_{i=0}^m (y_i - p(x_i))^2 \quad (1.4)$$

minimal wird.

Wir setzen das gesuchte Polynom in der Monom-Basis an, d.h.

$$p(x) = \alpha_0 + \alpha_1 x + \dots + \alpha_{n-1} x^{n-1} + \alpha_n x^n.$$

Gesucht ist daher

$$\min_{\alpha_0, \dots, \alpha_n} \sum_{i=0}^m (y_i - (\alpha_0 + \alpha_1 x_i + \dots + \alpha_{n-1} x_i^{n-1} + \alpha_n x_i^n))^2.$$

Im Spezialfall  $n = m$  können für die Koeffizienten  $\alpha_0, \dots, \alpha_n$  das Interpolationspolynom eingesetzt werden. Dadurch wird die Zielfunktion (1.4) null und somit liegt bereits das eindeutige Minimum vor.

Bei dieser Approximationsaufgabe ist ein Polynom von niedrigem Grad  $n$  erwünscht, welches dennoch eine Menge von Stützpunkten mit hohem  $m$  gut approximiert. Daher wird die Frage der Konvergenz für  $n \rightarrow \infty$  nicht gestellt. Für  $n \geq m$  liegt das Interpolationspolynom vor und somit gelten die Aussagen aus Abschnitt 1.1.

Die obige Aufgabenstellung ist bei der Approximation von kontinuierlichen Funktionen  $f$  sinnvoll, wenn die Funktionsauswertungen  $y_i := f(x_i) + \varepsilon_i$  mit Fehlern  $\varepsilon_i$  (z.B. Messfehlern) behaftet sind. Eine Interpolation der Werte würde damit unnötigerweise den Fehler mit erfassen.

Wir definieren die Vektoren

$$\theta := (\alpha_0, \dots, \alpha_n)^\top \in \mathbb{R}^{n+1}, \quad y := (y_0, \dots, y_m)^\top \in \mathbb{R}^{m+1}$$

und die Matrix

$$\Phi := \begin{pmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \cdots & x_m^n \end{pmatrix} \in \mathbb{R}^{(m+1) \times (n+1)}. \quad (1.5)$$

Die Aufgabenstellung lässt sich daher als lineares Ausgleichsproblem schreiben, d.h.

$$\min_{\theta \in \mathbb{R}^{n+1}} \|y - \Phi\theta\|_2^2$$

mit der Euklidischen Norm  $\|\cdot\|_2$ . Die Matrix  $\Phi$  besitzt vollen Spaltenrang, da die Stützstellen paarweise verschieden sind. Somit existiert eine eindeutige Lösung des Ausgleichsproblems. Die Lösung  $\theta$  kann über eine  $QR$ -Zerlegung von  $\Phi$  berechnet werden, siehe [19]. Die  $QR$ -Zerlegung erfolgt mit entweder Householder-Transformationen oder Givens-Rotationen.

Ein Hinzufügen von weiteren Stützpunkten  $(x_i, f(x_i))$  bedeutet hier die Vergrößerung der Matrix  $\Phi$  um weitere Zeilen. Wurde eine Matrix  $\Phi$  bereits transformiert, dann brauchen die Transformationen nur auf den neu hinzugekommenen Teil angewendet zu werden.

### 1.3 Interpolation mit Splines

Eine naheliegende Idee ist es, statt mit Polynomen die Interpolation mit stückweise polynomialen Funktionen durchzuführen. Sei eine Zerlegung  $\Delta$  der Form  $a = x_0 < x_1 < \cdots < x_n = b$  gegeben. Es bezeichnet

$$\rho(\Delta) := \max_{j=1, \dots, n} |x_j - x_{j-1}|$$

die Feinheit der Zerlegung. Man beachte, dass die beiden Randpunkte hier Stützstellen sind.

**Def. 1.3** Die Menge der Splines vom Grad  $k$  zu einer Zerlegung  $\Delta$  von  $[a, b]$  ist definiert durch

$$\mathcal{S}_k(\Delta) := \{s \in C^{k-1}[x_0, x_n] : s|_{[x_{j-1}, x_j]} \in \mathcal{P}_k \text{ für } j = 1, \dots, n\}.$$

In der Praxis werden meist nur Splines vom Grad  $k = 1, 2, 3, 4, 5$  eingesetzt, wobei der Fall  $k = 3$  am häufigsten auftritt.

## Lineare Splines

Im Fall  $k = 1$  entstehen als Splines stetige Streckenzüge. Der lineare interpolierende Spline interpoliert dann eine Funktion  $f$  an den Stützpunkten  $(x_j, f(x_j))$  für  $j = 0, 1, \dots, n$ . Es folgt die Formel

$$s(x) = \frac{x_j - x}{x_j - x_{j-1}} f(x_{j-1}) + \frac{x - x_{j-1}}{x_j - x_{j-1}} f(x_j) \quad \text{für } x \in [x_{j-1}, x_j].$$

Die Stetigkeit von  $s$  ist mit dieser Konstruktion sichergestellt. Folgende Approximationsgüte kann gezeigt werden.

**Satz 1.12** Sei  $x_j = a + jh$  für  $j = 0, 1, \dots, n$  mit  $h := \frac{b-a}{n}$ . Ist  $f \in C^2[a, b]$ , dann gilt für den linearen interpolierenden Spline  $s$

$$\|f - s\|_\infty \leq \frac{1}{8} \|f''\|_\infty h^2 \quad \text{und} \quad \|f' - s'\|_\infty \leq \frac{1}{2} \|f''\|_\infty h,$$

wobei die Ableitungen definiert sind über  $s'(x_0) = s'(x_{0+})$ ,  $s'(x_n) = s'(x_{n-})$ ,  $s'(x_j) = \frac{1}{2}(s'(x_{j-}) + s'(x_{j+}))$  für  $j = 1, \dots, n-1$ .

Beweis: siehe [22].

Falls  $f$  nur stetig auf  $[a, b]$  ist, dann kann man die gleichmäßige Konvergenz des linearen interpolierenden Splines noch durch die gleichmäßige Stetigkeit von  $f$  nachweisen.

**Satz 1.13** Für eine Folge  $(\Delta_n)_{n \in \mathbb{N}}$  von Zerlegungen mit  $\rho(\Delta_n) \rightarrow 0$  konvergieren die linearen interpolierenden Splines  $s$  gleichmäßig gegen  $f \in C[a, b]$ .

Beweis:

Sei  $\varepsilon > 0$ . Da  $f$  gleichmäßig stetig ist gibt es ein  $\delta > 0$ , so dass

$$|x - x'| < \delta \quad \Rightarrow \quad |f(x) - f(x')| < \varepsilon$$

für alle  $x, x' \in [a, b]$ . Es gibt ein  $n_0$ , so dass  $\rho(\Delta_n) < \delta$  für alle  $n \geq n_0$ . Für festes  $x \in [a, b]$  existiert jeweils ein von  $n$  abhängiges  $j \in \{1, \dots, n\}$  mit  $x \in [x_{j-1}, x_j]$ . Dadurch gilt

$$s(x) - f(x) = \underbrace{\frac{x_j - x}{x_j - x_{j-1}}}_{=: \eta_j} f(x_{j-1}) + \underbrace{\frac{x - x_{j-1}}{x_j - x_{j-1}}}_{=: \xi_j} f(x_j) - f(x).$$

Es gilt stets  $\eta_j + \xi_j = 1$  und  $\eta_j, \xi_j \geq 0$ . Daher können wir abschätzen

$$\begin{aligned} |s(x) - f(x)| &= |\eta_j f(x_{j-1}) + \xi_j f(x_j) - (\eta_j + \xi_j) f(x)| \\ &\leq \eta_j |f(x_{j-1}) - f(x)| + \xi_j |f(x_j) - f(x)| \\ &\leq \eta_j \varepsilon + \xi_j \varepsilon = \varepsilon \end{aligned}$$

wegen  $|x - x_j| < \delta$  und  $|x - x_{j-1}| < \delta$ . Somit ist die gleichmäßige Konvergenz gezeigt.  $\square$

Die Approximation mit linearen Splines ist in der Praxis jedoch meistens ungeeignet, da diese Funktionen nicht stetig differenzierbar sind.

### Konstruktion kubischer Splines

Wir verwenden im folgenden stets kubische Splines, d.h.  $k = 3$ . Dadurch folgt die stückweise Darstellung

$$s(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3 \quad \text{für } x_i \leq x \leq x_{i+1}. \quad (1.6)$$

Damit gilt für  $i = 0, 1, \dots, n - 1$

$$s(x_i) = a_i, \quad s'(x_i) = b_i, \quad s''(x_i) = 2c_i, \quad s'''(x_i) = 6d_i.$$

Die  $4n$  Koeffizienten  $a_i, b_i, c_i, d_i$  stellen Freiheitsgrade dar. Die Glattheitsbedingung  $s \in C^2[x_0, x_n]$  ergibt jedoch bereits  $3(n - 1)$  Bedingungen.

Seien Stützpunkte  $(x_j, y_j)$  für  $j = 0, 1, \dots, n$  gegeben. Ein kubischer interpolierender Spline  $s$  erfüllt die Eigenschaft

$$s(x_j) = y_j \quad \text{für } j = 0, 1, \dots, n. \quad (1.7)$$

Somit werden  $n + 1$  Bedingungen gestellt. Mit dieser Forderung ist  $s$  noch nicht eindeutig bestimmt. Wir stellen die natürlichen Randbedingungen

$$s''(x_0) = s''(x_n) = 0, \quad (1.8)$$

wodurch insgesamt  $4n$  Bedingungen vorliegen. Der kubische interpolierende Spline ist damit eindeutig festgelegt wie wir nachher zeigen.

Die Krümmung einer Funktion  $f \in C^2$  wird definiert als

$$\kappa(x) := \frac{f''(x)}{\sqrt{1 + f'(x)^2}^3} \approx f''(x)$$

für kleine  $|f'(x)|$ . Wir definieren die Gesamtkrümmung

$$J(f) := \|f''\|_{L^2}^2 = \int_{x_0}^{x_n} (f''(x))^2 \, dx. \quad (1.9)$$

Für den kubischen Spline  $s \in C^2[x_0, x_n]$  mit der Eigenschaft (1.7) und den Randbedingungen (1.8) folgt  $J(s) \leq J(f)$  für alle  $f \in C^2[x_0, x_n]$  mit gleicher Interpolationseigenschaft, siehe [19]. Somit minimiert der kubische interpolierende Spline mit natürlichen Randbedingungen die Gesamtkrümmung (1.9).

Nun berechnen wir die Koeffizienten in (1.6). Wir definieren die Schrittweiten  $h_i := x_i - x_{i-1}$  und die Werte

$$M_i := s''(x_i) \quad \text{für } i = 0, 1, \dots, n-1, n.$$

Die Randbedingungen (1.8) ergeben  $M_0 = M_n = 0$ . Die Eigenschaft  $s^{(4)} \equiv 0$  und die Stetigkeit von  $s''$  führt auf

$$s''(x) = M_i \frac{x_{i+1} - x}{h_{i+1}} + M_{i+1} \frac{x - x_i}{h_{i+1}} \quad \text{für } x_i \leq x \leq x_{i+1}. \quad (1.10)$$

Integration liefert

$$\begin{aligned} s'(x) &= -M_i \frac{(x_{i+1}-x)^2}{2h_{i+1}} + M_{i+1} \frac{(x-x_i)^2}{2h_{i+1}} + B_i \\ s(x) &= M_i \frac{(x_{i+1}-x)^3}{6h_{i+1}} + M_{i+1} \frac{(x-x_i)^3}{6h_{i+1}} + B_i(x-x_i) + A_i \end{aligned}$$

mit Konstanten  $A_i, B_i$ . Die Interpolationsbedingungen (1.7) ergeben

$$M_i \frac{h_{i+1}^2}{6} + A_i = y_i, \quad M_{i+1} \frac{h_{i+1}^2}{6} + B_i h_{i+1} + A_i = y_{i+1}$$

für  $i = 0, 1, \dots, n-1$  und stellen gleichzeitig die Stetigkeit von  $s$  sicher. Die Integrationskonstanten können nun bestimmt werden

$$A_i = y_i - M_i \frac{h_{i+1}^2}{6}, \quad B_i = \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{h_{i+1}}{6} (M_{i+1} - M_i).$$

Damit folgen die Konstanten in der Darstellung (1.6) als

$$\begin{aligned} a_i &= y_i \\ b_i &= \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{h_{i+1}}{6} (2M_i + M_{i+1}) \\ c_i &= \frac{1}{2} M_i \\ d_i &= \frac{1}{6h_{i+1}} (M_{i+1} - M_i). \end{aligned}$$

Es verbleibt also nur das Problem die Koeffizienten  $M_i$  zu bestimmen. Die erste Ableitung des Splines kann dargestellt werden als

$$s'(x) = -M_i \frac{(x_{i+1} - x)^2}{2h_{i+1}} + M_{i+1} \frac{(x - x_i)^2}{2h_{i+1}} + \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{h_{i+1}}{6} (M_{i+1} - M_i)$$

für  $x_i \leq x \leq x_{i+1}$ . Als links- und rechtsseitige Grenzwerte folgen

$$\begin{aligned} s'(x_i-) &= \frac{y_i - y_{i-1}}{h_i} + \frac{h_i}{3} M_i + \frac{h_i}{6} M_{i-1}, \\ s'(x_i+) &= \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{h_{i+1}}{3} M_i - \frac{h_{i+1}}{6} M_{i+1}. \end{aligned}$$

Die Stetigkeit von  $s'$  erzeugt die Bedingungen

$$h_i M_{i-1} + 2(h_i + h_{i+1}) M_i + h_{i+1} M_{i+1} = 6 \left( \frac{y_{i+1} - y_i}{h_{i+1}} - \frac{y_i - y_{i-1}}{h_i} \right)$$





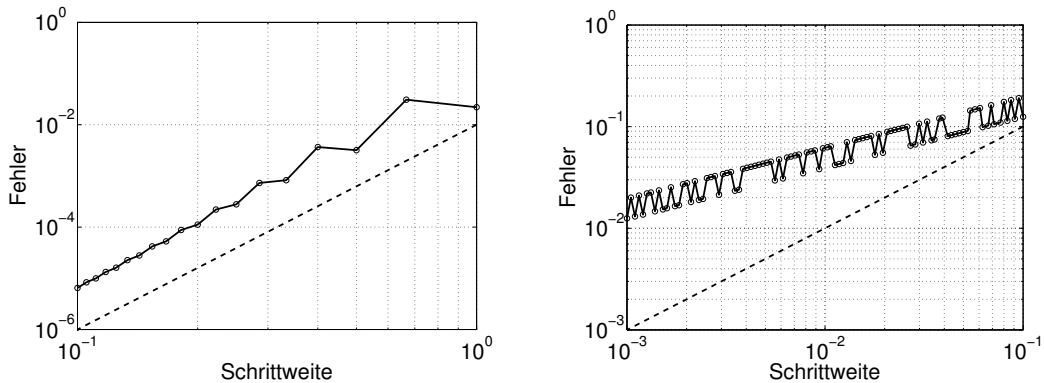


Abbildung 2: Approximationsfehler in Maximumnorm für kubischen interpolierenden Spline bei  $f(x) = \frac{1}{1+x^2}$  in  $[-5, 5]$  zusammen mit Gerade für  $H^4$  (links) und bei  $f(x) = \sqrt{|x|}$  in  $[-1, 1]$  zusammen mit Gerade für  $H$  (rechts) in doppelt-logarithmischer Skala.

Beweis: siehe [15].

Für den kubischen interpolierenden Spline mit den natürlichen Randbedingungen (1.8) kann eine Konstante  $C = \frac{3}{4}$  nachgewiesen werden falls  $f''(a) = f''(b) = 0$  gilt. Diese Randwerte können erzeugt werden, indem zur ursprünglichen Funktion  $f$  ein Polynom (höchstens) dritten Grades addiert wird.

Im Fall von äquidistanten Stützstellen gilt  $H := h_{\max} = h_{\min}$  und es folgt  $c = 1$ . Wir erhalten

$$\|s - f\|_{\infty} \leq \|f^{(4)}\|_{\infty} H^4,$$

d.h. eine Konvergenzgeschwindigkeit  $\mathcal{O}(\frac{1}{n^4})$ . Im Gegensatz zur Polynominterpolation ist die Wahl äquidistanter Knoten bei der Splineinterpolation günstig.

Abbildung 2 zeigt den Approximationsfehler in der Maximumnorm für äquidistante Stützstellen beim Beispiel von Runge, wo die Funktion beliebig oft differenzierbar ist, und bei einer Wurzelfunktion, die an einer Stelle nicht differenzierbar ist. Entsprechend erkennen wir bei der hinreichend glatten Funktion eine Konvergenz von vierter Ordnung, obwohl die Voraussetzung  $f''(a) = f''(b) = 0$  nicht erfüllt ist. Bei der nicht-glatten Funktion liegt die Konvergenz auch vor, jedoch mit einer Geschwindigkeit langsamer als lineare Konvergenz.

Desweiteren wird in Abbildung 3 die Konvergenzgeschwindigkeit des kubischen interpolierenden Splines bei äquidistanten Knoten verglichen mit der Polynominterpolation bei Tschebycheff-Knoten. Beim Beispiel von Runge ist die Funktion beliebig oft differenzierbar, wodurch die totale Variation aller Ableitungen

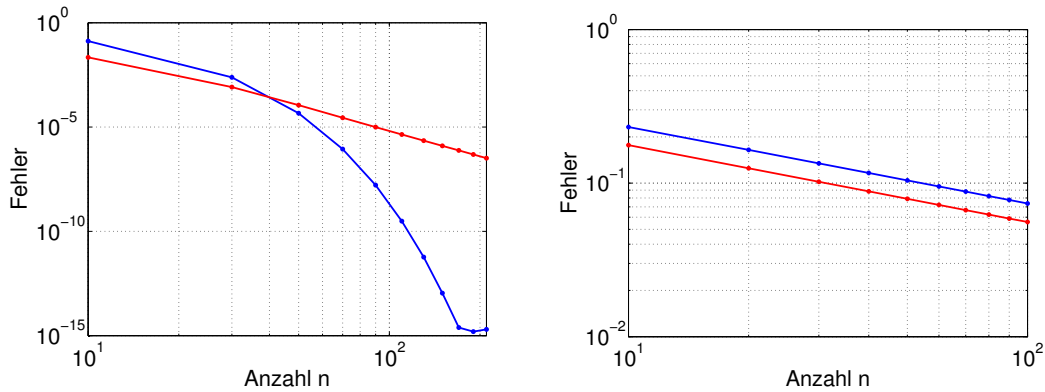


Abbildung 3: Approximationsfehler bei  $f(x) = \frac{1}{1+x^2}$  in  $[-5, 5]$  (links) und bei  $f(x) = \sqrt{|x|}$  in  $[-1, 1]$  (rechts) für kubischen interpolierenden Spline (rot) und für Polynominterpolation mit Tschebycheff-Knoten (blau) in doppelt-logarithmischer Skala.

existiert. Nach Satz 1.6 wird die Konvergenzgeschwindigkeit bei Tschebycheff-Knoten beliebig schnell. Dies ist auch erkennbar bis zu einem Fehler in der Größenordnung der Maschinengenauigkeit ( $\varepsilon_0 \approx 10^{-16}$ ). Bei der Wurzelfunktion ist für die Tschebycheff-Knoten die Voraussetzung für die Konvergenz aus Satz 1.5 nicht erfüllt. Trotzdem beobachten wir eine Konvergenz mit der gleichen Geschwindigkeit wie für den Spline.

Es verbleibt die Frage, ob bei nur stetiger Funktion  $f$  auf  $[a, b]$  der kubisch interpolierende Spline gleichmäßig gegen  $f$  konvergiert falls  $\rho(\Delta_n) \rightarrow 0$ . Hierzu gibt es nur Resultate für den periodischen Spline.

**Satz 1.15** Sei  $f \in C[0, 1]$  und  $f(0) = f(1)$ . Zu einer Folge von Zerlegungen  $(\Delta_n)_{n \in \mathbb{N}}$  des Intervalls  $[0, 1]$  mit  $\rho(\Delta_n) \rightarrow 0$  definiert man  $K_n := \frac{h_{\max}}{h_{\min}}$ . Falls die Menge aller  $K_n$  beschränkt bleibt, dann konvergiert die Folge der kubischen interpolierenden periodischen Splines gleichmäßig gegen  $f$ .

Beweis: siehe [18].

Offensichtlich ist die Voraussetzung für den Satz 1.15 für äquidistante Stützstellen gegeben, da dann  $K_n = 1$  für alle  $n$  gilt. Die Beschränktheit der  $K_n$  bedeutet, dass die Gitterfolge nicht entartet, d.h.  $h_{\min}$  kann nicht beliebig klein im Vergleich zu  $h_{\max}$  werden.

### Konvergenzaussagen bei Splines allgemeinen Grades

Wir betrachten jetzt Splines von einem beliebigen Grad  $2m-1$  für ein  $m \geq 2$ . Eine

zu approximierende Funktion  $f \in C^m[a, b]$  sei gegeben. Es gilt  $\dim(\mathcal{S}_{2m-1}(\Delta)) = n + 2m - 1$ . Dann gibt es bei interpolierenden Splines  $2m - 2$  Freiheitsgrade, die mit einem der folgenden Randbedingungstypen festgelegt werden.

(i) *natürliche Randbedingungen*

$$s^{(j)}(a) = s^{(j)}(b) = 0 \quad \text{für } j = m, \dots, 2m - 2.$$

(ii) *vollständige Randbedingungen*

$$s^{(j)}(a) = f^{(j)}(a) \quad \text{und} \quad s^{(j)}(b) = f^{(j)}(b) \quad \text{für } j = 1, \dots, m - 1.$$

(iii) *periodische Randbedingungen*

$$s^{(j)}(a) = s^{(j)}(b) \quad \text{für } j = 1, \dots, 2m - 2$$

unter der Voraussetzung  $f^{(j)}(a) = f^{(j)}(b)$  für  $j = 0, \dots, m - 1$ .

Man kann zeigen, dass diese Bedingungen jeweils einen interpolierenden Spline eindeutig festlegen. Für die Konvergenz der Approximation gilt dann das folgende Resultat.

**Satz 1.16** *Sei  $f \in C^m[a, b]$  und  $s \in \mathcal{S}_{2m-1}(\Delta)$  der interpolierende Spline eines der Typen (i), (ii), (iii). Dann gilt die Abschätzung*

$$\|f^{(j)} - s^{(j)}\|_{\infty} \leq \frac{m!}{\sqrt{m}} \frac{1}{j!} \|f^{(m)}\|_{L_2} \rho(\Delta)^{m-j-\frac{1}{2}}$$

für beliebiges  $m \geq 2$  und  $j = 0, 1, \dots, m - 1$ .

Beweis: siehe [10], Kapitel 6, Abschnitt 5.4.

Man beachte, dass hier die  $L_2$ -Integralnorm von  $f^{(m)}$  auftritt, welche sich aber durch die Maximumnorm abschätzen läßt über

$$\|g\|_{L_2} = \left( \int_a^b g(x)^2 dx \right)^{\frac{1}{2}} \leq \sqrt{b-a} \|g\|_{\infty}.$$

Im Fall  $m = 2$  der kubischen Splines folgen die nächsten beiden Abschätzungen mit  $h_{\max} = \rho(\Delta)$

$$\begin{aligned} \|f - s\|_{\infty} &\leq \sqrt{2} \|f''\|_{L_2} h_{\max}^{\frac{3}{2}}, \\ \|f' - s'\|_{\infty} &\leq \sqrt{2} \|f''\|_{L_2} h_{\max}^{\frac{1}{2}}, \end{aligned}$$

Die gleichmäßige Konvergenz ist somit für die Funktion und ihre erste Ableitung garantiert. Im Gegensatz zu den Resultaten aus dem vorhergehenden Abschnitt brauchen hier keine Forderungen an die minimale Schrittweite  $h_{\min}$  gestellt zu werden.

Die Konvergenzraten aus Satz 1.16 sind nicht optimal als Preis dafür, dass die Aussagen für beliebigen Grad  $2m - 1$  gelten. Daher lassen sich für konkrete Gerade oft bessere Konvergenzgeschwindigkeiten ohne höhere Glattheitsforderungen finden. Optimale Schranken und kleinstmögliche Konstanten für kubische interpolierende Splines mit vollständigen Randbedingungen unter der Annahme  $f \in C^4$  finden sich in [9].

## 1.4 Ausgleichsspline

Die folgende Vorgehensweise folgt im wesentlichen der Beschreibung in [16]. Seien die Stützpunkte  $(x_j, y_j)$  für  $j = 0, 1, \dots, n$  gegeben mit  $x_0 < x_1 < \dots < x_n$ . Dabei kann  $n$  groß sein. Falls die Werte  $y_j$  (und/oder  $x_j$ ) fehlerbehaftet sind beispielsweise durch Messfehler, dann ist eine Interpolation der Stützpunkte nicht sinnvoll. Daher approximieren wir die Stützpunkte mit einer Funktion  $f \in C^2$ . Abbildung 4 verdeutlicht diesen Ansatz.

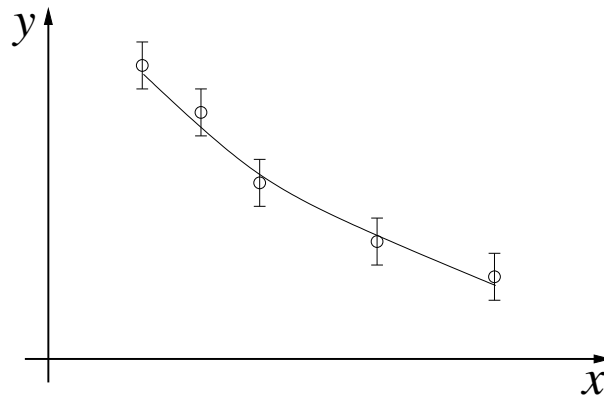


Abbildung 4: Interpolation von Daten mit Messfehlern.

### Minimierung mit Nebenbedingung

Da wir Oszillationen in der Approximation  $f$  vermeiden möchten, stellen wir die Minimierungsaufgabe

$$\min_{f \in C^2[x_0, x_n]} J(f) \quad \text{mit} \quad J(f) := \int_{x_0}^{x_n} (f''(x))^2 dx \quad (1.12)$$

unter der Nebenbedingung

$$\sum_{i=0}^n \left( \frac{f(x_i) - y_i}{w_i} \right)^2 \leq S \quad (1.13)$$

mit Gewichten  $w_i > 0$  und dem Glättungsparameter  $S \geq 0$ . Die Gewichte  $w_i$  werden in Abhängigkeit von der Größe der Messfehler gewählt. Ein sinnvoller Bereich für den Glättungsparameter ist  $S \in [N - \sqrt{2N}, N + \sqrt{2N}]$  mit  $N := n + 1$ .

Für  $S = 0$  folgt  $f(x_i) = y_i$  für alle  $i$  und die optimale Funktion  $f$  ist gerade der interpolierende kubische Spline mit natürlichen Randbedingungen. Ohne die Nebenbedingung (1.13) ergibt sich das Minimum der Aufgabe (1.12) aus  $f'' \equiv 0$  ( $J(f) = 0$ ), d.h. eine Gerade  $f(x) = \alpha x + \beta$ . Diesen Fall kann man als  $S \rightarrow \infty$  interpretieren. Die bestapproximierende Gerade bezüglich der Stützpunkte und der Gewichte ist die Ausgleichsgerade. Die Ausgleichsgerade ist Lösung des linearen Ausgleichsproblems

$$\min_{z \in \mathbb{R}^2} \|D^{-1}(Az - y)\|_2^2 \quad (1.14)$$

mit

$$D = \begin{pmatrix} w_0 & & & \\ & w_1 & & \\ & & \ddots & \\ & & & w_n \end{pmatrix}, \quad A = \begin{pmatrix} x_0 & 1 \\ x_1 & 1 \\ \vdots & \vdots \\ x_n & 1 \end{pmatrix}, \quad z = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}, \quad y = \begin{pmatrix} y_0 \\ y_1 \\ \vdots \\ y_n \end{pmatrix}.$$

Sei  $f^*(x) = \alpha^*x + \beta^*$  die eindeutige Lösung des Ausgleichsproblems (1.14). Wir erhalten für das Residuum

$$S_0 := \sum_{i=0}^n \left( \frac{\alpha^*x_i + \beta^* - y_i}{w_i} \right)^2.$$

Für  $S \in [0, S_0]$  existiert eine eindeutige Lösung der Minimierungsaufgabe (1.12) mit Nebenbedingung (1.13). Für  $S > S_0$  ist die Lösung nicht eindeutig und wir wählen die Ausgleichsgerade  $f^*$  in diesem Fall. Man beachte im folgenden, dass jede Gerade auch einen kubischen Spline mit natürlichen Randbedingungen darstellt.

Die Ungleichung (1.13) kann in eine Gleichung umgeformt werden durch Einführung einer Variablen  $z \in \mathbb{R}$

$$\sum_{i=0}^n \left( \frac{f(x_i) - y_i}{w_i} \right)^2 + z^2 - S = 0.$$

Wir koppeln diese Bedingung an die Minimierung (1.12) über einen Lagrange-Parameter  $p$ . Es folgt das Funktional

$$\tilde{J}(f, p, z) = \int_{x_0}^{x_n} (f''(x))^2 dx + 2p \left( \sum_{i=0}^n \left( \frac{f(x_i) - y_i}{w_i} \right)^2 + z^2 - S \right).$$

Eine Variationsrechnung kann nun durchgeführt werden basierend auf diesem Funktional, siehe [8].

Sei  $s$  die Lösung der Aufgabe (1.12),(1.13). Wir betrachten ein stückweise kubisches Polynom

$$s(x) = a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3 \quad \text{für } x \in [x_i, x_{i+1}] \quad (1.15)$$

mit  $i = 0, 1, \dots, n - 1$ . An den Stützstellen sind die links- und rechtsseitigen Grenzwerte dann

$$s^{(k)}(x_i+) = \lim_{h \rightarrow 0, h > 0} s^{(k)}(x_i + h), \quad s^{(k)}(x_i-) = \lim_{h \rightarrow 0, h > 0} s^{(k)}(x_i - h).$$

Wegen  $s \in C^2$  erhalten wir die Bedingungen

$$s(x_i+) = s(x_i-), \quad s'(x_i+) = s'(x_i-), \quad s''(x_i+) = s''(x_i-) \quad (1.16)$$

für  $i = 1, \dots, n - 1$ . Wir fordern die Randbedingungen

$$s''(x_0) = s''(x_n) = 0. \quad (1.17)$$

Desweiteren wird die Sprungbedingung

$$s'''(x_i-) - s'''(x_i+) = 2p \frac{s(x_i) - y_i}{w_i^2} \quad \text{für } i = 0, 1, \dots, n \quad (1.18)$$

mit  $s'''(x_0-) = s'''(x_n+) = 0$  gestellt. Wenn der Lagrange-Parameter  $p$  gegeben ist, dann liefern (1.16),(1.17),(1.18) hier  $4n$  Gleichungen für die  $4n$  unbekanntenen Koeffizienten in (1.15).

## Berechnung der Koeffizienten

Mit den Schrittweiten  $h_i := x_{i+1} - x_i$  erhalten wir

$$\begin{aligned} s(x) &= a_i + b_i(x - x_i) + c_i(x - x_i)^2 + d_i(x - x_i)^3, \\ s'(x) &= b_i + 2c_i(x - x_i) + 3d_i(x - x_i)^2, \\ s''(x) &= 2c_i + 6d_i(x - x_i), \\ s'''(x) &= 6d_i \end{aligned}$$

für  $x \in [x_i, x_{i+1}]$  und  $i = 0, 1, \dots, n - 1$ . Es folgt:

- Die dritte Bedingung aus (1.16) und die Randbedingungen (1.17) liefern

$$2c_i + 6d_i h_i = 2c_{i+1} \quad \Rightarrow \quad d_i = \frac{c_{i+1} - c_i}{3h_i}$$

für  $i = 0, 1, \dots, n-1$  mit  $c_0 = c_n = 0$ .

- Die erste Bedingung aus (1.16) führt auf

$$a_i + b_i h_i + c_i h_i^2 + d_i h_i^3 = a_{i+1} \quad \Rightarrow \quad b_i = \frac{a_{i+1} - a_i}{h_i} - c_i h_i - d_i h_i^2$$

für  $i = 0, 1, \dots, n-1$ .

- Die zweite Bedingung aus (1.16) ergibt

$$b_i + 2c_i h_i + 3d_i h_i^2 = b_{i+1}.$$

Mit den Beziehungen zwischen  $b_i$  und  $d_i$  erreichen wir

$$\frac{h_{i-1}}{3} c_{i-1} + \frac{2}{3} (h_{i-1} + h_i) c_i + \frac{h_i}{3} c_{i+1} = \frac{1}{h_{i-1}} a_{i-1} - \left( \frac{1}{h_{i-1}} + \frac{1}{h_i} \right) a_i + \frac{1}{h_i} a_{i+1}$$

für  $i = 1, \dots, n-1$ .

Wir definieren die Vektoren  $c := (c_1, \dots, c_{n-1})^\top$  und  $a := (a_0, \dots, a_n)^\top$  sowie die Matrizen  $T \in \mathbb{R}^{(n-1) \times (n-1)}$  und  $Q \in \mathbb{R}^{(n+1) \times (n-1)}$

$$T := \begin{pmatrix} t_{11} & t_{12} & & 0 \\ t_{21} & \ddots & \ddots & \\ & \ddots & \ddots & t_{n-2, n-1} \\ 0 & & t_{n-1, n-2} & t_{n-1, n-1} \end{pmatrix}, \quad Q := \begin{pmatrix} q_{11} & & & 0 \\ q_{21} & \ddots & & \\ q_{31} & \ddots & q_{n-1, n-1} & \\ & \ddots & q_{n, n-1} & \\ 0 & & q_{n+1, n-1} & \end{pmatrix}$$

mit  $t_{ii} = \frac{2}{3}(h_{i-1} + h_i)$ ,  $t_{i+1, i} = t_{i, i+1} = \frac{1}{3}h_i$ ,  $q_{ii} = \frac{1}{h_{i-1}}$ ,  $q_{i+1, i} = -\frac{1}{h_{i-1}} - \frac{1}{h_i}$ ,  $q_{i+2, i} = \frac{1}{h_i}$ . Die Tridiagonalmatrix  $T$  ist (bis auf ein Vielfaches) identisch mit der Matrix (1.11) aus der kubischen Splineinterpolation. Es folgt das Gleichungssystem

$$Tc = Q^\top a. \quad (1.19)$$

Daher koppelt  $T$  die Koeffizienten  $c_i$  aus der zweiten Ableitung mit einem Differenzenschema aus den Werten  $a_i$ .

Die Sprungbedingung (1.18) ergibt direkt

$$6(d_{i-1} - d_i) = 2p \frac{a_i - y_i}{w_i^2} \quad (1.20)$$



für  $i = 0, \dots, n$  mit  $d_{-1} := d_n := 0$ . Ersetzen der  $d_i$  durch die  $c_i$  liefert

$$Qc = pD^{-2}(y - a).$$

Durch die Gleichungen

$$Q^\top(D^2Qc) = Q^\top(p(y - a)) = pQ^\top y - pQ^\top a = pQ^\top y - pTc$$

kann die Berechnung von  $a$  und  $c$  aus (1.19),(1.20) entkoppelt werden zu

$$(Q^\top D^2Q + pT)c = pQ^\top y, \quad p \neq 0 : \quad a = y - \frac{1}{p}D^2Qc. \quad (1.21)$$

Die Matrix  $Q^\top D^2Q + pT$  ist symmetrisch und positiv definit für  $p \geq 0$  wegen

$$x^\top(Q^\top D^2Q + pT)x = \underbrace{\|DQx\|_2^2}_{>0} + p\underbrace{x^\top Tx}_{>0} > 0 \quad \text{für jedes } x \neq 0.$$

Man beachte, dass  $T$  positiv definit ist und  $Q$  vollen Rang besitzt. Desweiteren ist  $Q^\top D^2Q + pT$  eine Bandmatrix mit Breite 5. Für gegebenen Parameter  $p$  können wir sukzessive  $c_i, a_i, d_i, b_i$  berechnen.

### Bestimmung des Lagrange-Parameters

Die Aufgabe ist somit auf die Identifizierung des Lagrange-Parameters  $p$  zurückgeführt. Die Funktion  $s$  erfüllt die Bedingung

$$\sum_{i=0}^n \left( \frac{s(x_i) - y_i}{w_i} \right)^2 + z^2 - S = 0.$$

Die Summe besitzt die Darstellung

$$\sum_{i=0}^n \left( \frac{s(x_i) - y_i}{w_i} \right)^2 = \|D^{-1}(y - a)\|_2^2.$$

Mit

$$D^{-1}(y - a) = \frac{1}{p}DQc = DQ(Q^\top D^2Q + pT)^{-1}Q^\top y,$$

was auch für  $p = 0$  gilt wegen der Stetigkeit, folgt

$$F(p)^2 = S - z^2$$

mit

$$F(p) := \|DQ(Q^\top D^2Q + pT)^{-1}Q^\top y\|_2. \quad (1.22)$$

Es kann gezeigt werden, dass  $F(p)^2$  eine streng monoton fallende und konvexe Funktion für  $p \geq 0$  ist. Also ist auch  $F(p)$  injektiv.

Deshalb erhalten wir eine Verbindung zwischen  $p$  und  $z$ . Wir fordern

$$p \geq 0 \quad \text{und} \quad pz = 0. \quad (1.23)$$

Somit treten zwei Fälle auf:

- 1. Fall:  $p = 0$

Die Gleichung (1.21) liefert direkt  $c = 0$ . Es folgt  $d_i = 0$  sowie  $b_i = b_{i+1}$  und  $a_{i+1} = a_i + h_i b_i$  für  $i = 0, 1, \dots, n-1$ . Somit ist die Funktion  $s \in C^2$  eine Gerade. Es gilt  $J(s) = 0$  für (1.9). Im Unterfall  $S = S_0$  ist dann die Ausgleichsgerade die eindeutige Lösung des Minimierungsproblems mit Nebenbedingung. Für jede andere Gerade ist das Residuum des Ausgleichsproblems größer als  $S_0$ . In (1.14) gilt  $Az = a$ . Mit  $F(p)^2 = \|D^{-1}(y - a)\|_2^2$  folgt  $F(0)^2 = S_0$ . Im Unterfall  $S > S_0$  ist die Lösung des Minimierungsproblems mit Nebenbedingung nicht eindeutig. Dann wähle die Ausgleichsgerade.

- 2. Fall:  $p > 0$

Bedingung (1.23) impliziert  $z = 0$ . Der Lagrange-Parameter ist Lösung der nichtlinearen Gleichung

$$F(p)^2 - S = 0.$$

Da  $F^2$  eine konvexe Funktion ist, liegt im Newton-Verfahren sogar globale Konvergenz vor. Man kann als Startwert beispielsweise  $p^{(0)} = 0$  setzen.

Es folgt

$$(Q^\top D^2 Q + pT)^{-1} \stackrel{p \gg 1}{\approx} \frac{1}{p} T^{-1} \xrightarrow{p \rightarrow \infty} 0 (\in \mathbb{R}^{(n-1) \times (n-1)})$$

und damit

$$\lim_{p \rightarrow \infty} F(p) = 0 (\in \mathbb{R}).$$

Dementsprechend besitzt die nichtlineare Gleichung eine Lösung für jeden Glättungsparameter  $0 < S < S_0$ . Die Funktion  $F(p)$  ist in Abbildung 5 dargestellt.

Die Newton-Iteration für die äquivalente Gleichung  $F(p) - \sqrt{S} = 0$  lautet

$$p^{(j+1)} = p^{(j)} - \frac{F(p^{(j)}) - \sqrt{S}}{F'(p^{(j)})} = p^{(j)} - \frac{F(p^{(j)})^2 - F(p^{(j)})\sqrt{S}}{F(p^{(j)})F'(p^{(j)})}. \quad (1.24)$$

Man kann zeigen, dass

$$F(p)F'(p) = pu^\top T(Q^\top D^2 Q + pT)^{-1}Tu - u^\top Tu$$

mit  $u = p^{-1}c = (Q^\top D^2 Q + pT)^{-1}Q^\top y$ . Es folgt  $F(p) = \|DQu(p)\|$ .

Wir erkennen, dass entweder  $s$  zur Ausgleichsgeraden aus (1.14) resultiert ( $p = 0$ ) oder die Ungleichung (1.13) wird zur Gleichung ( $p > 0, z = 0$ ).

Wir zeigen noch nachträglich die Monotonie der Funktion (1.22).

**Lemma 1.2** *In der Funktion  $F$  aus (1.22) sei  $Q^\top y \neq 0$ . Dann ist  $F^2$  streng monoton fallend und konvex für alle  $p \geq 0$ .*

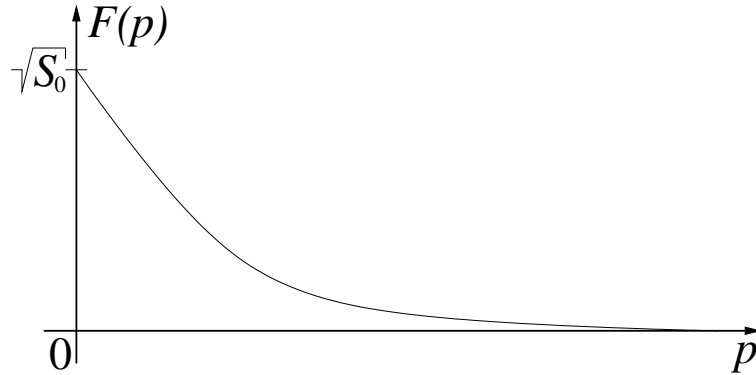


Abbildung 5: Funktion  $F(p)$  in Abhängigkeit vom Lagrange-Parameter  $p$ .

Beweis:

Die Matrix  $T$  ist symmetrisch und positiv definit. Dadurch gilt  $T = U^\top \hat{D} U$  mit Diagonalmatrix  $\hat{D} = \text{diag}(\hat{d}_1, \dots, \hat{d}_{n-1})$  und orthogonaler Matrix  $U$ . Für

$$T^{\frac{1}{2}} := U^\top \hat{D}^{\frac{1}{2}} U \quad \text{mit} \quad \hat{D}^{\frac{1}{2}} = \text{diag} \left( \sqrt{\hat{d}_1}, \dots, \sqrt{\hat{d}_{n-1}} \right)$$

gilt dann  $T^{\frac{1}{2}} T^{\frac{1}{2}} = T$ . Es ist auch  $T^{\frac{1}{2}}$  symmetrisch und positiv definit. Dadurch erhalten wir

$$Q^\top D^2 Q + pT = T^{\frac{1}{2}} (T^{-\frac{1}{2}} Q^\top D^2 Q T^{-\frac{1}{2}} + pI) T^{\frac{1}{2}}$$

sowie

$$(Q^\top D^2 Q + pT)^{-1} = T^{-\frac{1}{2}} (T^{-\frac{1}{2}} Q^\top D^2 Q T^{-\frac{1}{2}} + pI)^{-1} T^{-\frac{1}{2}}.$$

Es bezeichne  $S := T^{-\frac{1}{2}} Q^\top D^2 Q T^{-\frac{1}{2}}$ . Die Matrix  $S$  ist symmetrisch und positiv definit, da  $Q^\top D^2 Q$  positiv definit ist und  $T^{-\frac{1}{2}}$  symmetrisch sowie regulär ist. Es seien  $\lambda_1, \dots, \lambda_{n-1} > 0$  die Eigenwerte von  $S$  und  $v_1, \dots, v_{n-1}$  eine zugehörige Orthonormalbasis aus Eigenvektoren.

Wir berechnen

$$\begin{aligned} F(p)^2 &= (DQT^{-\frac{1}{2}}(S + pI)^{-1}T^{-\frac{1}{2}}Q^\top y)^\top (DQT^{-\frac{1}{2}}(S + pI)^{-1}T^{-\frac{1}{2}}Q^\top y) \\ &= y^\top QT^{-\frac{1}{2}}(S + pI)^{-1}T^{-\frac{1}{2}}Q^\top D^2QT^{-\frac{1}{2}}(S + pI)^{-1}T^{-\frac{1}{2}}Q^\top y \\ &= z^\top (S + pI)^{-1}S(S + pI)^{-1}z \end{aligned}$$

mit  $z := T^{-\frac{1}{2}}Q^\top y$ . Da  $T^{-\frac{1}{2}}$  regulär ist, gilt nach Voraussetzung  $z \neq 0$ .

In der Basisdarstellung der Eigenvektoren ist  $z = \alpha_1 v_1 + \dots + \alpha_{n-1} v_{n-1}$  mit  $\alpha_j \neq 0$  für mindestens ein  $j$ . Die Matrix  $(S + pI)^{-1}S(S + pI)^{-1}$  besitzt die gleichen Eigenvektoren wie  $S$  und die zugehörigen Eigenwerte lauten

$$\mu_i := \frac{\lambda_i}{(\lambda_i + p)^2} \quad \text{für } i = 1, \dots, n-1.$$

Dadurch gilt

$$F(p)^2 = \sum_{i=1}^{n-1} \frac{\alpha_i^2 \lambda_i}{(\lambda_i + p)^2} = \sum_{\alpha_i \neq 0} \frac{\alpha_i^2 \lambda_i}{(\lambda_i + p)^2}.$$

Wir differenzieren

$$\frac{d}{dp} F(p)^2 = -2 \sum_{\alpha_i \neq 0} \frac{\alpha_i^2 \lambda_i}{(\lambda_i + p)^3}$$

und

$$\frac{d^2}{dp^2} F(p)^2 = 6 \sum_{\alpha_i \neq 0} \frac{\alpha_i^2 \lambda_i}{(\lambda_i + p)^4}$$

Offensichtlich gilt  $(F^2)' < 0$  und  $(F^2)'' > 0$  für alle  $p \geq 0$ . □

### Algorithmus

Zur Berechnung des Ausgleichssplines wird hier die einzelne nichtlineare Gleichung  $F(p) - \sqrt{S} = 0$  mit der Newton-Iteration gelöst. Die Auswertung von (1.24) erfolgt in den Schritten:

1. Berechnung der Cholesky-Zerlegung

$$R^\top R = Q^\top D^2 Q + pT.$$

2. Löse  $R^\top Ru = Q^\top y$  mit Vorwärtssubstitution  $R^\top r = Q^\top y$  und Rückwärtssubstitution  $Ru = r$ .
3. Berechne  $F := DQu$ ,  $\bar{F} := F^\top F$ ,  $f := Tu$ ,  $\bar{f} := u^\top f$ .
4. Löse  $R^\top v = f$  und bestimme  $g := v^\top v$ .

Nun wird die Newton-Iteration (1.24) zu

$$p^{(j+1)} = p^{(j)} - \frac{\bar{F} - \sqrt{S\bar{F}}}{p^{(j)}g - \bar{f}}.$$

Die Iteration kann durchgeführt werden bis die Lösung auf Maschinengenauigkeit vorliegt. Eine geeignete Abbruchbedingung lautet  $F(p^{(j)})^2 \leq S$ .

Dieser Algorithmus ist im Fall  $S = 0$  nicht durchführbar, da kein entsprechender Lagrange-Parameter existiert ( $p \rightarrow +\infty$  für  $S \rightarrow 0$ ).

### Verifikation der Optimalität

In den vorhergehenden Abschnitten haben wir eine kubische Splinefunktion identifiziert, welche die Bedingung (1.13) erfüllt. Jetzt ist noch zu zeigen, dass diese Funktion die Minimierungsaufgabe löst.

**Satz 1.17** *Der Ausgleichsspline  $s$  definiert durch die Bedingungen (1.15), (1.16), (1.17), (1.18), (1.23) stellt ein Minimum laut (1.12) unter der Nebenbedingung (1.13) dar, d.h. für alle Funktionen  $f \in C^2[x_0, x_n]$  mit*

$$\sum_{i=0}^n \left( \frac{f(x_i) - y_i}{w_i} \right)^2 \leq S, \quad (1.25)$$

folgt

$$\int_{x_0}^{x_n} f''(x)^2 dx \geq \int_{x_0}^{x_n} s''(x)^2 dx.$$

Beweis:

Wir erhalten

$$\begin{aligned} \int_{x_0}^{x_n} f''(x)^2 dx &= \int_{x_0}^{x_n} (f''(x) - s''(x))^2 + 2(f''(x) - s''(x))s''(x) + s''(x)^2 dx \\ &\geq 2 \int_{x_0}^{x_n} (f''(x) - s''(x))s''(x) dx + \int_{x_0}^{x_n} s''(x)^2 dx. \end{aligned}$$

Nachzuweisen ist somit

$$\int_{x_0}^{x_n} (f''(x) - s''(x))s''(x) dx \geq 0 \quad (1.26)$$

für alle  $f \in C^2[x_0, x_n]$  welche (1.25) erfüllen. Partielle Integration liefert

$$\begin{aligned} &\int_{x_0}^{x_n} (f''(x) - s''(x))s''(x) dx \\ &= \sum_{i=0}^{n-1} \int_{x_i}^{x_{i+1}} (f''(x) - s''(x))s''(x) dx \\ &= \sum_{i=0}^{n-1} [(f' - s')s'']_{x_i}^{x_{i+1}} - [(f - s)s''']_{x_i}^{x_{i+1}} + \int_{x_i}^{x_{i+1}} (f(x) - s(x))s^{(4)}(x) dx. \end{aligned}$$

Wegen  $f, s \in C^2[x_0, x_n]$  und den Randbedingungen (1.17) für  $s$  verschwindet der erste Term. Da  $s$  ein stückweise kubisches Polynom ist, siehe (1.15), gilt  $s^{(4)} \equiv 0$  und der dritte Term verschwindet. Mit  $s'''(x_0-) = s'''(x_n+) = 0$  und der

Sprungbedingung (1.18) wird der zweite Term weiter umgeformt

$$\begin{aligned}
& \int_{x_0}^{x_n} (f''(x) - s''(x))s''(x) \, dx \\
&= -\sum_{i=0}^n (f(x_i) - s(x_i))(s'''(x_i-) - s'''(x_i+)) \\
&= -\sum_{i=0}^n (f(x_i) - s(x_i))2p \frac{s(x_i) - y_i}{w_i^2} \\
&= -2p \sum_{i=0}^n \frac{(f(x_i) - y_i)(s(x_i) - y_i)}{w_i^2} - \frac{(s(x_i) - y_i)^2}{w_i^2} \\
&= 2p \left( S - z^2 - \sum_{i=0}^n \frac{(f(x_i) - y_i)(s(x_i) - y_i)}{w_i^2} \right).
\end{aligned}$$

Für  $p = 0$  ist die Ungleichung (1.26) erfüllt. Für  $p > 0$  liefert die Bedingung (1.23) dann  $z = 0$ . Wir haben zu zeigen, dass

$$\sum_{i=0}^n \frac{(f(x_i) - y_i)(s(x_i) - y_i)}{w_i^2} \leq S.$$

Mit der Definition  $u := D^{-1}(f - y)$  und  $v := D^{-1}(s - y)$  entspricht diese Behauptung  $u^\top v \leq S$ . Die Nebenbedingung (1.13) führt direkt auf  $\|u\|_2^2 \leq S$ ,  $\|v\|_2^2 \leq S$ . Mit der Cauchy-Schwarzschen Ungleichung schätzen wir ab

$$|u^\top v| \leq \|u\|_2 \|v\|_2 \leq \sqrt{S} \cdot \sqrt{S} = S.$$

Somit ist die Eigenschaft (1.26) gezeigt. □

### Beispiel 1:

Wir verwenden eine Menge aus fünf Stützpunkten. Die Gewichte werden auf  $w_i = 1$  für alle  $i$  gesetzt. Abbildung 6 zeigt die resultierenden Ausgleichssplines für verschiedene Parameter  $S$ . Zum einen ergibt sich eine Funktion ähnlich zum kubischen interpolierenden Spline im Fall  $S \approx 0$ . Zum anderen entsteht die Ausgleichsgerade für hohe Werte  $S$ .

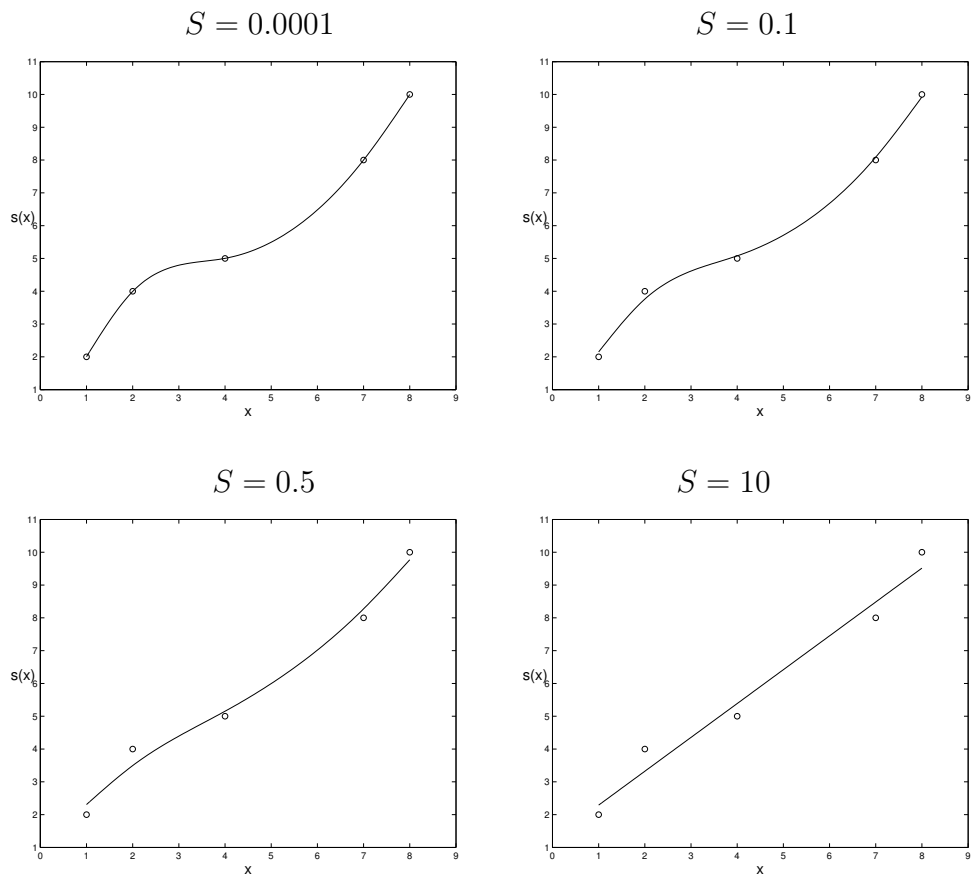


Abbildung 6: Ausgleichsspline für identische Stützpunkte und unterschiedliche Glättungsparameter  $S$ .

### Beispiel 2:

Sei  $x_i = a + ih$  mit  $h = \frac{b-a}{n}$  für  $i = 0, 1, \dots, n$ . Wir betrachten die Punkte

$$y_i = g(x_i) + r_i \quad \text{für } i = 0, 1, \dots, n$$

mit der vorgegebenen Funktion

$$g(x) = 1 + \sin(2\pi x).$$

Die Werte  $r_i$  stellen Zufallszahlen aus einer Gleichverteilung in  $[-0.1, 0.1]$  dar. Der Erwartungswert und die Varianz ergeben sich daher zu

$$\mathbb{E}(r_i) = 0, \quad \text{Var}(r_i) = \sigma^2 := \frac{0.2^2}{12} = \frac{1}{300} \quad \text{für } i = 0, 1, \dots, n.$$

Alternativ können in der Praxis auch paarweise verschiedene Varianzen auftreten.

Zu den Daten  $(x_i, y_i)$  kann man versuchen eine Funktion der Form

$$f(x) = \alpha + \beta \sin(\omega x + \varphi)$$

mit a priori unbekanntem Parametern  $\alpha, \beta, \omega, \varphi \in \mathbb{R}$  anzupassen. Falls nur  $\alpha$  und  $\beta$  unbekannt sind, dann folgt ein lineares Ausgleichsproblem. Jedoch ergibt sich ein nichtlineares Gleichungssystem falls alle Parameter als unbekannt angenommen werden. Zudem erfordert die Auswahl des Ansatzes für  $f$  eine Untersuchung der Gestalt bzw. Struktur der Daten  $(x_i, y_i)$ .

Im Gegensatz hierzu wenden wir den Ansatz des Ausgleichssplines an, welcher unabhängig von der Gestalt der Eingabedaten ist. Wir setzen  $n = 100$ ,  $a = 0$ ,  $b = 1.5$  und wählen die Gewichte  $w_i = \sigma$  für alle  $i$ . Abbildung 7 stellt die entstehenden Ausgleichssplines für verschiedene Parameter  $S$  dar. Die Wahl  $S = n + 1$  liefert eine gute Approximation der zugrunde liegenden Funktion  $g$ . Viel größere Werte  $S$  führen auf eine schlechte Approximation sowohl der Daten als auch der Funktion  $g$ , da zu stark geglättet wird. Viel kleinere Werte  $S$  erzeugen unerwünschte Oszillationen im Spline, wobei die Stützpunkte zwar besser approximiert werden, jedoch die Funktion  $g$  schlechter erfasst wird.



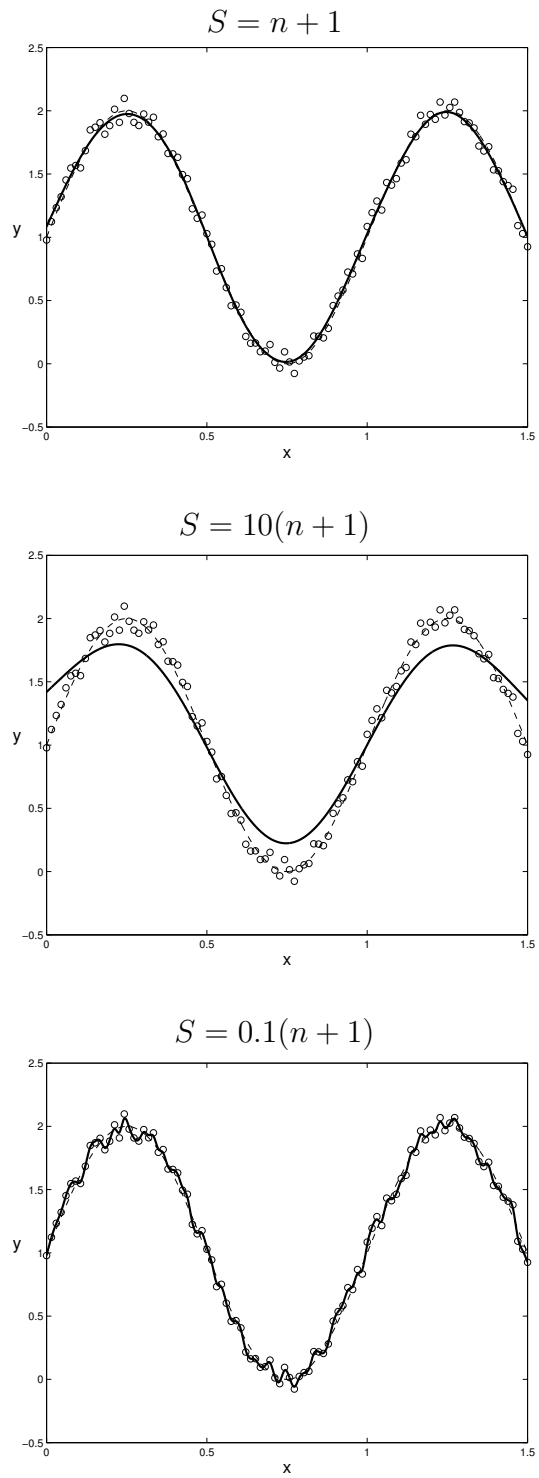


Abbildung 7: Ausgleichsspline (—) für identische Stützpunkte ( $\circ$ ) und verschiedene Parameter  $S$  zusammen mit der ursprünglichen Funktion  $g$  (- - -).

## 2 Approximation in normierten Räumen

In diesem Kapitel betrachten wir eine allgemeine Approximationstheorie in normierten Vektorräumen. Als günstigsten Fall erweisen sich dabei Hilbert-Räume. Die Ansätze zur Approximation werden konkret bei Fourier-Reihen und Wavelets eingesetzt.

### 2.1 Allgemeine Approximationstheorie

Zunächst wiederholen wir einige Begriffe aus der Topologie.

#### Definitionen aus der Topologie

Wir betrachten einen reellen normierten Vektorraum  $(V, \|\cdot\|)$ , welcher eine beliebige Dimension haben kann. Die Norm induziert eine Metrik  $d$  über

$$d(v_1, v_2) := \|v_1 - v_2\| \quad \text{für } v_1, v_2 \in V.$$

Zu einem  $u \in V$  und  $\varepsilon > 0$  definieren wir die Kugel

$$B_\varepsilon(u) := \{v \in V : \|v - u\| < \varepsilon\}.$$

Für eine Folge  $(v_i)_{i \in \mathbb{N}} \subset V$  schreiben wir

$$\hat{v} = \lim_{i \rightarrow \infty} v_i \quad \text{falls} \quad \lim_{i \rightarrow \infty} \|v_i - \hat{v}\| = 0.$$

Damit folgen die nachstehenden Begriffe.

**Def. 2.1** Sei  $(V, \|\cdot\|)$  ein normierter Vektorraum.

(i) Eine Menge  $M \subseteq V$  heißt offen, wenn

$$\forall v \in M \exists \varepsilon > 0 : B_\varepsilon(v) \subset M.$$

(ii) Eine Menge  $M \subseteq V$  heißt abgeschlossen, wenn  $V \setminus M$  offen ist.

(iii) Eine Menge  $M \subseteq V$  heißt beschränkt, wenn ein  $\varepsilon > 0$  existiert mit  $M \subset B_\varepsilon(0)$ .

(iv) Eine Menge  $M \subseteq V$  heißt vollständig, wenn jede Cauchy-Folge aus  $M$  einen Grenzwert in  $M$  besitzt.

Ist der ganze Raum  $V$  vollständig, dann nennt man  $(V, \|\cdot\|)$  einen Banach-Raum. Von besonderer Bedeutung sind auch kompakte Mengen, welche dadurch gekennzeichnet sind, dass jede offene Überdeckung eine endliche Teilüberdeckung besitzt.

**Def. 2.2** Eine Menge  $K \subset V$  heißt kompakt, wenn aus

$$K \subset \bigcup_{i \in I} O_i$$

mit offenen Mengen  $O_i$  die Existenz endlich vieler Indizes  $i_1, i_2, \dots, i_m$  mit

$$K \subset O_{i_1} \cup O_{i_2} \cup \dots \cup O_{i_m}$$

folgt.

Eine kompakte Menge ist stets beschränkt und abgeschlossen. Die Umkehrung gilt meistens nicht. Jedoch ist im  $\mathbb{R}^n$  jede beschränkte und abgeschlossene Menge auch kompakt nach dem Satz von Heine-Borel. Eine wichtige Eigenschaft kompakter Mengen ist der Satz von Bolzano-Weierstraß.

**Satz 2.1** Sei  $K \subset V$  eine kompakte Teilmenge. Dann besitzt jede Folge aus  $K$  eine konvergente Teilfolge mit Grenzwert in  $K$ .

Beweis: siehe [6], Satz 9.

Man beachte, dass diese Aussage auch ohne die Vollständigkeit von  $V$  gilt.

**Satz 2.2** Ein beliebiger Vektorraum  $V$  mit endlicher Dimension  $n$  ist isomorph zu  $\mathbb{R}^n$ .

Beweis: siehe [5], Abschnitt 2.2.4, Korollar 2.

Aus diesem Satz folgt sofort, dass ein endlichdimensionaler Vektorraum vollständig ist, da  $\mathbb{R}^n$  für jedes  $n$  vollständig ist. Insbesondere gilt die Aussage für endlichdimensionale Untervektorräume  $U \subset V$  eines möglicherweise unendlichdimensionalen Vektorraums  $V$ .

**Satz 2.3** Ein endlichdimensionaler Untervektorraum  $U \subseteq V$  eines normierten Vektorraums  $(V, \|\cdot\|)$  ist abgeschlossen.

Beweis:

Wir zeigen, dass  $V \setminus U$  offen ist. Sei  $v \in V \setminus U$ . Angenommen, es gelte für kein  $\varepsilon > 0$  die Inklusion  $B_\varepsilon(v) \subset V \setminus U$ . Dann enthielte jede Kugel um  $v$  Elemente aus  $U$  und somit gäbe es eine Folge aus  $U$ , die gegen  $v$  konvergiert. Da nach Satz 2.2  $U$  isomorph zu einem  $\mathbb{R}^n$  ist, folgt wegen der Vollständigkeit von  $\mathbb{R}^n$  dann für den Grenzwert  $v \in U$ . Durch diesen Widerspruch gibt es ein  $\varepsilon > 0$  mit  $B_\varepsilon(v) \subset V \setminus U$ .  $\square$

Besonders günstig erweisen sich Vektorräume  $V$  mit einem Skalarprodukt  $\langle \cdot, \cdot \rangle$ , d.h. einer positiv definiten, symmetrischen Bilinearform. Das Skalarprodukt erzeugt eine Norm durch

$$\|v\| = \sqrt{\langle v, v \rangle} \quad \text{für } v \in H.$$

Somit liegt auch ein normierter Vektorraum vor. Ist  $(V, \langle \cdot, \cdot \rangle)$  vollständig, dann liegt ein Hilbert-Raum vor. Ist  $(V, \langle \cdot, \cdot \rangle)$  nicht vollständig, so spricht man von einem Prä-Hilbert-Raum.

## Strikte Normen

Desweiteren benötigen wir den Konvexitätsbegriff.

**Def. 2.3** Sei  $(V, \|\cdot\|)$  ein normierter Vektorraum.

(i) Eine Menge  $M \subseteq V$  heißt konvex, falls gilt

$$v_1, v_2 \in M \quad \Rightarrow \quad \lambda v_1 + (1 - \lambda)v_2 \in M$$

für alle  $\lambda \in [0, 1]$ .

(ii) Eine Menge  $M \subseteq V$  heißt streng konvex, falls mit  $v_1, v_2 \in M$  und  $v_1 \neq v_2$  folgt

$$\forall \lambda \in (0, 1) \exists \varepsilon > 0 : B_\varepsilon(\lambda v_1 + (1 - \lambda)v_2) \subset M.$$

Für normierte Räume betrachten wir damit die folgende Verschärfung.

**Def. 2.4** Ein normierter Vektorraum  $(V, \|\cdot\|)$  heißt strikt normiert, falls die abgeschlossene Einheitskugel  $\{v \in V : \|v\| \leq 1\}$  streng konvex ist.

Äquivalent zur strikten Normiertheit ist die Bedingung

$$\|v_1\|, \|v_2\| \leq 1, \quad v_1 \neq v_2 \quad \Rightarrow \quad \|\lambda v_1 + (1 - \lambda)v_2\| < 1.$$

Gilt  $\|v_1\| \leq 1$  und  $\|v_2\| < 1$  sowie  $v_1 \neq v_2$ , dann folgt

$$\|\lambda v_1 + (1 - \lambda)v_2\| \leq |\lambda| \cdot \|v_1\| + |1 - \lambda| \cdot \|v_2\| < |\lambda| + |1 - \lambda| = 1$$

für  $0 < \lambda < 1$ , wodurch die Bedingung aus der strengen Konvexität erfüllt ist. Somit braucht nur der Fall  $\|v_1\| = \|v_2\| = 1$  betrachtet zu werden. Notwendig und hinreichend für die strikte Normiertheit ist daher das Kriterium

$$\|v_1\| = \|v_2\| = 1, \quad v_1 \neq v_2 \quad \Rightarrow \quad \|\lambda v_1 + (1 - \lambda)v_2\| < 1.$$

Auf dem  $\mathbb{R}^n$  haben wir als übliche Normen

$$\|x\|_p = \sqrt[p]{\sum_{j=1}^n |x_j|^p} \quad (2.1)$$

für  $p \geq 1$  mit der Maximumnorm

$$\|x\|_\infty = \max_{j=1, \dots, n} |x_j|$$

als Grenzfall. Die Normen (2.1) sind genau dann strikt, wenn  $1 < p < \infty$  gilt.

Der folgende Satz liefert eine Aussage für (Prä-)Hilbert-Räume.

**Satz 2.4** *Wird in einem normierten Vektorraum  $(V, \|\cdot\|)$  die Norm von einem Skalarprodukt induziert, dann ist der Raum strikt normiert.*

Beweis:

Sei  $\|v_1\| = \|v_2\| = 1$  und  $v_1 \neq v_2$ . Wir müssen

$$\|(1 - \lambda)v_1 + \lambda v_2\| < 1 \quad \text{für } \lambda \in (0, 1)$$

zeigen. Jeder Fall  $0 < \lambda \leq \frac{1}{2}$  kann dargestellt werden als

$$(1 - \lambda)v_1 + \lambda v_2 = \frac{1}{2} \left( \underbrace{v_1}_{=: \tilde{v}_1} + \underbrace{(v_1 + 2\lambda(v_2 - v_1))}_{=: \tilde{v}_2} \right).$$

Es gilt  $\|\tilde{v}_1\| = 1$  sowie

$$\|\tilde{v}_2\| = \|(1 - 2\lambda)v_1 + 2\lambda v_2\| \leq |1 - 2\lambda| \cdot \|v_1\| + |2\lambda| \cdot \|v_2\| = |1 - 2\lambda| + |2\lambda| = 1.$$

Gilt  $\|\tilde{v}_2\| < 1$  dann ist die obige Bedingung sofort erfüllt. Es genügt somit den Fall  $\lambda = \frac{1}{2}$  zu betrachten.

Da die Norm von einem Skalarprodukt erzeugt wird, gilt die Parallelogrammgleichung

$$\|v_1 + v_2\|^2 + \|v_1 - v_2\|^2 = 2(\|v_1\|^2 + \|v_2\|^2) \quad \text{für alle } v_1, v_2.$$

Sei  $v_1 \neq v_2$  und  $\|v_1\| = \|v_2\| = 1$ . Die Parallelogrammgleichung liefert

$$\begin{aligned} \left\| \frac{1}{2}v_1 + \frac{1}{2}v_2 \right\|^2 &= \frac{1}{4}\|v_1 + v_2\|^2 = \frac{1}{4}(2\|v_1\|^2 + 2\|v_2\|^2 - \|v_1 - v_2\|^2) \\ &< \frac{1}{2}\|v_1\|^2 + \frac{1}{2}\|v_2\|^2 = 1. \end{aligned}$$

Somit folgt  $\|\frac{1}{2}v_1 + \frac{1}{2}v_2\| < 1$ . □

## Bestapproximationen

Die folgende Definition gilt für beliebige Teilmengen.

**Def. 2.5** Sei  $(V, \|\cdot\|)$  ein normierter Vektorraum und  $M \subset V$  mit  $M \neq \emptyset$  sowie  $v \in V$  fest gegeben.

(i) Die reelle Zahl

$$E_M(v) = \inf_{u \in M} \|u - v\|$$

wird als Minimalabstand des Vektors  $v$  zur Teilmenge  $M$  bezeichnet.

(ii) Eine Folge  $(u_i)_{i \in \mathbb{N}} \subset M$  heißt  $M$ -Minimalfolge an  $v$ , wenn

$$\lim_{i \rightarrow \infty} \|u_i - v\| = E_M(v).$$

(iii) Man nennt einen Vektor  $\hat{u} \in M$  eine Bestapproximation in  $M$  oder auch  $M$ -Proximum falls

$$\|\hat{u} - v\| = \inf_{u \in M} \|u - v\|.$$

Der Minimalabstand existiert immer und ist eindeutig. Die Existenz einer Minimalfolge ergibt sich aus der Eigenschaft des Infimums. Eine naheliegende Fragestellung ist nun, für welche Teilmengen die Existenz und Eindeutigkeit einer Bestapproximation garantiert werden kann.

## Beispiele:

- Wir betrachten  $(\mathbb{R}^2, \|\cdot\|_2)$  und  $M = \{x \in \mathbb{R}^2 : \|x\|_2 \leq 1\}$ . Zu  $v \notin M$  gibt es ein eindeutiges  $M$ -Proximum, nämlich

$$\hat{u} = \frac{1}{\|v\|_2} v.$$

Es folgt  $E_M(v) = \|v\|_2 - 1$ . Dass  $\hat{u}$  ein  $M$ -Proximum ist erkennt man wie folgt: Sei  $\tilde{u} \in \mathbb{R}^2$  mit  $\|\tilde{u} - v\|_2 < \|\hat{u} - v\|_2$ . Es folgt

$$\|v\|_2 \leq \|\tilde{u} - v\|_2 + \|\tilde{u}\|_2 < \|\hat{u} - v\|_2 + \|\tilde{u}\|_2 = \|v\|_2 - 1 + \|\tilde{u}\|_2.$$

Daraus folgt sofort  $\|\tilde{u}\|_2 > 1$  und somit  $\tilde{u} \notin M$ .

2. Wir verwenden  $(C[0, 1], \|\cdot\|_\infty)$  und  $M = \{e^{\beta x} : \beta > 0\}$  sowie  $v \equiv \frac{1}{2}$ . Es gilt

$$\|e^{\beta x} - \frac{1}{2}\|_\infty = e^\beta - \frac{1}{2} > \frac{1}{2} \quad \text{für alle } \beta > 0.$$

Für  $\beta \rightarrow 0$  sieht man  $E_M(v) = \frac{1}{2}$ . Somit existiert kein  $M$ -Proximum von  $v$ . Eine Minimalfolge ist durch jede Funktionenfolge in  $M$  mit  $\beta \rightarrow 0$  gegeben.

Folgendes Lemma liefert ein hinreichendes Kriterium für die Existenz einer Bestapproximation.

**Lemma 2.1** *Sei  $M \subset V$  mit  $M \neq \emptyset$  für einen normierten Raum  $(V, \|\cdot\|)$  und  $(u_i)_{i \in \mathbb{N}}$  eine  $M$ -Minimalfolge an  $v \in V$ . Besitzt die Folge einen Häufungspunkt  $\hat{u} \in M$ , dann ist  $\hat{u}$  ein  $M$ -Proximum an  $v$ .*

Beweis:

Ist  $\hat{u} \in M$  ein Häufungspunkt der Folge, dann gibt es eine Teilfolge  $(u_{i_k})_{k \in \mathbb{N}}$ , die gegen  $\hat{u}$  konvergiert, d.h.

$$\lim_{k \rightarrow \infty} \|u_{i_k} - \hat{u}\| = 0.$$

Wir setzen an

$$\|\hat{u} - v\| \leq \|\hat{u} - u_{i_k}\| + \|u_{i_k} - v\| \quad \text{für alle } k$$

und somit

$$\|\hat{u} - v\| \leq \lim_{k \rightarrow \infty} \|\hat{u} - u_{i_k}\| + \|u_{i_k} - v\| = 0 + E_M(v) = E_M(v).$$

Wegen der Abschätzung  $E_M(v) \leq \|u - v\|$  für alle  $u$  folgt die Behauptung.  $\square$

**Satz 2.5** *Ist  $K \subset V$  mit  $K \neq \emptyset$  eine kompakte Teilmenge in einem normierten Raum  $(V, \|\cdot\|)$ , dann existiert für jedes  $v \in V$  ein  $K$ -Proximum an  $v$ .*

Beweis:

Es existiert zu jeder nichtleeren Teilmenge  $K \subset V$  eine  $K$ -Minimalfolge. Als Folge in einer kompakten Teilmenge existiert nach Satz 2.1 eine konvergente Teilfolge mit Grenzwert in  $K$ . Da der Grenzwert ein Häufungspunkt der  $K$ -Minimalfolge ist, stellt dieser Grenzwert laut Lemma 2.1 ein  $K$ -Proximum dar.  $\square$

**Lemma 2.2** *Ist  $M \subset V$  für einen normierten Vektorraum  $(V, \|\cdot\|)$ , dann ist jede  $M$ -Minimalfolge (an beliebiges  $v \in V$ ) beschränkt.*

Beweis:

Sei  $(u_i)_{i \in \mathbb{N}}$  eine  $M$ -Minimalfolge für  $v \in V$ . Dann gibt es ein  $n^* \in \mathbb{N}$  mit

$$E_M(v) \leq \|u_i - v\| \leq E_M(v) + 1$$

für alle  $i \geq n^*$ . Es folgt

$$\|u_i\| \leq \|v\| + \|u_i - v\| \leq E_M(v) + 1 + \|v\| =: K_1$$

für  $i \geq n^*$ . Mit  $K_2 = \max\{\|u_i\| : i < n^*\}$  erhalten wir

$$\|u_i\| \leq \max\{K_1, K_2\} \quad \text{für alle } i \in \mathbb{N}.$$

Damit ist die Folge beschränkt. □

Das nächste Resultat nennt man auch den Fundamentalsatz der Approximationstheorie in normierten Vektorräumen.

**Satz 2.6** *Ist  $U \subseteq V$  ein endlichdimensionaler Untervektorraum eines normierten Vektorraums  $(V, \|\cdot\|)$ , dann existiert für jedes gegebene  $v \in V$  ein  $U$ -Proximum an  $v$ .*

Beweis:

Sei  $(u_i)_{i \in \mathbb{N}}$  eine  $U$ -Minimalfolge für  $v \in V$ . Nach Lemma 2.2 ist diese Folge beschränkt in  $V$ . Somit ist die Folge trivialerweise auch beschränkt in  $U$ , d.h.  $u_i \in K$  für alle  $i$  mit  $K = \{u \in U : \|u\| \leq c\}$  bei hinreichend hohem  $c$ . Die Menge  $K$  ist als abgeschlossene und beschränkte Teilmenge eines endlichdimensionalen Vektorraums kompakt in  $U$ . Satz 2.1 zeigt die Existenz einer konvergenten Teilfolge mit Grenzwert in  $K$ . Da der Grenzwert ein Häufungspunkt der  $U$ -Minimalfolge ist, stellt dieser Grenzwert laut Lemma 2.1 wieder ein  $U$ -Proximum dar. □

Man beachte, dass in diesem Beweis der Satz 2.5 nicht direkt angewendet wurde, denn dafür müsste man zuerst nachweisen, dass  $K$  auch kompakt in  $V$  und nicht nur in  $U$  ist.



**Beispiel:** Wir betrachten den normierten Raum  $(C[a, b], \|\cdot\|_\infty)$ . Es bezeichnet  $\mathcal{P}_n$  die Menge der Polynome mit Grad höchstens  $n$ . Somit ist  $\mathcal{P}_n \subset C[a, b]$  ein Unterraum der Dimension  $n + 1$ . Satz 2.6 garantiert die Existenz einer Bestapproximation in  $\mathcal{P}_n$  zu beliebigem  $f \in C[a, b]$ .

Nachdem wir hinreichende Bedingungen für die Existenz einer Bestapproximation gefunden haben wenden wir uns nun der Eindeutigkeitsfrage zu.

Eine strenge Konvexität impliziert die Eindeutigkeit, jedoch nicht die Existenz.

**Satz 2.7** *Ist  $M \subseteq V$  mit  $M \neq \emptyset$  eine streng konvexe Teilmenge eines normierten Vektorraums  $(V, \|\cdot\|)$ , dann existiert höchstens ein  $M$ -Proximum an  $v \in V$ .*

Beweis:

Angenommen es existiert mindestens ein Proximum. Es seien  $\hat{u}_1$  und  $\hat{u}_2$  beide  $M$ -Proxima an  $v$  sowie  $\hat{u}_1 \neq \hat{u}_2$ . Dann folgt

$$E_M(v) \leq \|v - \frac{1}{2}(\hat{u}_1 + \hat{u}_2)\| \leq \frac{1}{2}\|v - \hat{u}_1\| + \frac{1}{2}\|v - \hat{u}_2\| = \frac{1}{2}E_M(v) + \frac{1}{2}E_M(v) = E_M(v).$$

Dadurch haben wir  $\|v - \frac{1}{2}(\hat{u}_1 + \hat{u}_2)\| = E_M(v)$ .

Weil  $M$  streng konvex ist gibt es für  $\lambda = \frac{1}{2}$  Werte  $\mu$  mit

$$\tilde{u} = \frac{1}{2}(\hat{u}_1 + \hat{u}_2) + \mu(v - \frac{1}{2}(\hat{u}_1 + \hat{u}_2)) \in M$$

für alle  $|\mu| < \mu_0$  mit  $\mu_0 > 0$  hinreichend klein. Sei  $0 < \hat{\mu} < 1$  einer dieser Werte. Wir erhalten

$$\|\tilde{u} - v\| = \|\frac{1}{2}(1 - \hat{\mu})(\hat{u}_1 + \hat{u}_2) - (1 - \hat{\mu})v\| = |1 - \hat{\mu}| \cdot \|\frac{1}{2}(\hat{u}_1 + \hat{u}_2) - v\| = (1 - \hat{\mu})E_M(v).$$

Also folgt  $\|\tilde{u} - v\| < E_M(v)$  im Widerspruch zu  $E_M(v) \leq \|u - v\|$  für beliebiges  $u \in M$ .  $\square$

Die nächste Schlußfolgerung ergibt sich direkt aus Satz 2.5 und Satz 2.7.

**Korollar 2.1** *Ist  $K \subset V$  mit  $K \neq \emptyset$  eine kompakte und steng konvexe Teilmenge eines normierten Vektorraums  $(V, \|\cdot\|)$ , dann existiert zu jedem  $v \in V$  ein eindeutiges  $K$ -Proximum an  $v$ .*

Nach Satz 2.6 existiert für einen endlichdimensionalen Untervektorraum die Bestapproximation. In strikt normierten Räumen erhalten wir auch die Eindeutigkeit.

**Lemma 2.3** *Ein normierter Vektorraum  $(V, \|\cdot\|)$  ist genau dann strikt normiert, wenn für beliebige  $v, w \in V \setminus \{0\}$  gilt*

$$\|v\| + \|w\| = \|v + w\| \quad \Rightarrow \quad \exists \alpha \in \mathbb{R} : w = \alpha v.$$

Beweis:

Wir zeigen nur die wichtigere Implikation. Für die Umkehrung siehe [15], Theorem 15.18.

Sei  $(V, \|\cdot\|)$  strikt normiert. Wir verwenden  $v, w \in V \setminus \{0\}$  mit der Eigenschaft

$$\|v\| + \|w\| = \|v + w\|.$$

O.E.d.A. sei  $\|v\| \leq \|w\|$ . Es folgt mit der Dreiecksungleichung

$$\left\| \frac{v}{\|v\|} + \frac{w}{\|v\|} \right\| \leq \left\| \frac{v}{\|v\|} + \frac{w}{\|w\|} \right\| + \left\| \frac{w}{\|v\|} - \frac{w}{\|w\|} \right\|$$

und damit

$$\begin{aligned} \left\| \frac{v}{\|v\|} + \frac{w}{\|w\|} \right\| &\geq \left\| \frac{v}{\|v\|} + \frac{w}{\|v\|} \right\| - \left\| \frac{w}{\|v\|} - \frac{w}{\|w\|} \right\| \\ &= \frac{\|v + w\|}{\|v\|} - \left( \frac{1}{\|v\|} - \frac{1}{\|w\|} \right) \|w\| \\ &= \frac{\|v\| + \|w\|}{\|v\|} - \left( \frac{\|w\|}{\|v\|} - \frac{\|w\|}{\|w\|} \right) = 2. \end{aligned}$$

Also gilt

$$\left\| \frac{1}{2} \cdot \frac{v}{\|v\|} + \frac{1}{2} \cdot \frac{w}{\|w\|} \right\| \geq 1.$$

Wegen der strikten Normiertheit ist diese Ungleichung nur möglich falls  $\frac{v}{\|v\|} = \frac{w}{\|w\|}$ . Somit erhalten wir  $w = \alpha v$  mit  $\alpha = \frac{\|w\|}{\|v\|}$ .  $\square$

**Satz 2.8** *Sei  $(V, \|\cdot\|)$  ein strikt normierter Vektorraum und  $U \subseteq V$  ein Untervektorraum. Dann existiert höchstens ein  $U$ -Proximum zu einem beliebigen  $v \in V$ .*

Beweis:

Gilt  $v \in U$ , dann ist das  $U$ -Proximum eindeutig  $v$  selbst. Sei daher  $v \notin U$ . Sind  $\hat{u}_1$  und  $\hat{u}_2$  beide  $U$ -Proxima, dann folgt

$$E_U(v) \leq \|v - \frac{1}{2}(\hat{u}_1 + \hat{u}_2)\| \leq \frac{1}{2}\|v - \hat{u}_1\| + \frac{1}{2}\|v - \hat{u}_2\| = \frac{1}{2}E_U(v) + \frac{1}{2}E_U(v) = E_U(v).$$

Somit gilt

$$\left\| \frac{1}{2}(v - \hat{u}_1) + \frac{1}{2}(v - \hat{u}_2) \right\| = \left\| v - \frac{1}{2}(\hat{u}_1 + \hat{u}_2) \right\| = \left\| \frac{1}{2}(v - \hat{u}_1) \right\| + \left\| \frac{1}{2}(v - \hat{u}_2) \right\|.$$

Da der Raum strikt normiert ist können wir Lemma 2.3 anwenden und erhalten ein  $\alpha \in \mathbb{R}$ , so dass

$$\frac{1}{2}(v - \hat{u}_1) = \alpha \frac{1}{2}(v - \hat{u}_2) \quad \text{und damit} \quad (1 - \alpha)v = \hat{u}_1 - \alpha \hat{u}_2.$$

Weil  $v \notin U$  angenommen ist, erhalten wir  $1 - \alpha = 0$ , d.h.  $\alpha = 1$ . Also folgt  $\hat{u}_1 = \hat{u}_2$ .  $\square$

Als Folgerung aus Satz 2.6 und Satz 2.8 erhalten wir die nächste Aussage.

**Korollar 2.2** *Ist  $U \subseteq V$  ein endlichdimensionaler Untervektorraum eines strikt normierten Vektorraums  $(V, \|\cdot\|)$ , dann existiert für jedes gegebene  $v \in V$  ein eindeutiges  $U$ -Proximum an  $v$ .*

Falls die Eindeutigkeit nicht gegeben ist, dann liefert folgender Satz eine Information zur Struktur der Bestapproximationen.

**Satz 2.9** *Ist  $M \subset V$  mit  $M \neq \emptyset$  eine konvexe Teilmenge eines normierten Vektorraums  $(V, \|\cdot\|)$ , dann ist für jedes  $v \in V$  die Menge der  $M$ -Proxima an  $v$  konvex.*

Beweis:

Seien  $u_1$  und  $u_2$  beide  $M$ -Proxima an  $v$  und  $\lambda \in [0, 1]$ . Die Konvexität garantiert hier  $(1 - \lambda)u_1 + \lambda u_2 \in M$ . Es folgt

$$\begin{aligned} \|(1 - \lambda)u_1 + \lambda u_2 - v\| &= \|(1 - \lambda)(u_1 - v) + \lambda(u_2 - v)\| \\ &\leq |1 - \lambda| \cdot \|u_1 - v\| + |\lambda| \cdot \|u_2 - v\| \\ &= (1 - \lambda)E_M(v) + \lambda E_M(v) = E_M(v). \end{aligned}$$

Wegen  $E_M(v) \leq \|u - v\|$  für alle  $u \in M$  folgt  $E_M(v) = \|(1 - \lambda)u_1 + \lambda u_2 - v\|$ , d.h.  $(1 - \lambda)u_1 + \lambda u_2$  ist ebenfalls ein  $M$ -Proximum an  $v$ .  $\square$

Gibt es also zwei verschiedene Bestapproximationen an  $v$ , so existiert bereits ein Kontinuum aus Bestapproximationen zu  $v$ . Satz 2.9 gilt insbesondere für (unendlichdimensionale) Untervektorräume.

## Bestapproximation in Hilbert-Räumen

In Vektorräumen mit Skalarprodukt lassen sich die Bestapproximationen zu Untervektorräumen leicht charakterisieren und bestimmen. Laut Satz 2.4 ist ein (Prä-)Hilbert-Raum strikt normiert und das Proximum existiert jeweils und ist eindeutig.

**Def. 2.6** *Ist  $V$  ein Vektorraum mit Skalarprodukt  $\langle \cdot, \cdot \rangle$  und  $M \subseteq V$ , dann ist das orthogonale Komplement von  $M$  gegeben durch*

$$M^\perp = \{v \in V : \langle v, u \rangle = 0 \text{ für alle } u \in M\}.$$

Diese Definition wird insbesondere für Untervektorräume eingesetzt.

**Satz 2.10** *Sei  $(V, \langle \cdot, \cdot \rangle)$  ein (Prä-)Hilbert-Raum und  $U \subseteq V$  ein Untervektorraum. Ein Element  $\hat{u} \in U$  ist genau dann ein  $U$ -Proximum an  $v$ , wenn  $\hat{u} - v \in U^\perp$  gilt.*

Beweis:

Sei  $\hat{u} \in U$  und  $\hat{u} - v \in U^\perp$ . Für beliebiges  $u \in U$  folgern wir

$$\begin{aligned} \|u - v\|^2 &= \|\hat{u} - v + u - \hat{u}\|^2 \\ &= \|\hat{u} - v\|^2 + 2 \underbrace{\langle \hat{u} - v, u - \hat{u} \rangle}_{=0} + \|u - \hat{u}\|^2 \geq \|\hat{u} - v\|^2. \end{aligned}$$

Somit ist  $\hat{u}$  ein  $U$ -Proximum an  $v$ .

Sei nun  $\hat{u} \in U$  und  $\hat{u} - v \notin U^\perp$ . Dadurch existiert ein  $w \in U \setminus \{0\}$  mit  $\langle \hat{u} - v, w \rangle \neq 0$ . Mit  $t \in \mathbb{R}$  erhalten wir

$$\|\underbrace{\hat{u} + tw}_{\in U} - v\|^2 = \|\hat{u} - v\|^2 + 2t\langle \hat{u} - v, w \rangle + t^2\|w\|^2 =: f(t),$$

d.h. ein Polynom zweiten Grades in  $t$ . Differentiation zeigt

$$f'(t) = 2\langle \hat{u} - v, w \rangle + 2t\|w\|^2.$$

Ein Minimum von  $f$  liegt daher bei  $t^* = \frac{\langle \hat{u} - v, w \rangle}{\|w\|^2} \neq 0$  vor. Also ist  $\hat{u} + t^*w$  eine bessere Approximation und somit  $\hat{u}$  kein  $U$ -Proximum an  $v$ .  $\square$

Der nächste Satz zeigt die Konstruktion der Bestapproximation.

**Satz 2.11** Sei  $(V, \langle \cdot, \cdot \rangle)$  ein (Prä-)Hilbert-Raum und  $U \subseteq V$  ein endlichdimensionaler Untervektorraum mit gegebener Basis  $\{u_1, \dots, u_n\}$ . Es ist  $\hat{u} \in U$  mit

$$\hat{u} = \sum_{k=1}^n \alpha_k u_k$$

genau dann das  $U$ -Proximum an  $v \in V$ , wenn die Normalgleichungen

$$\sum_{k=1}^n \langle u_k, u_j \rangle \alpha_k = \langle v, u_j \rangle \quad \text{für } j = 1, 2, \dots, n \quad (2.2)$$

erfüllt sind.

Beweis:

Die Normalgleichungen lassen sich äquivalent umschreiben in

$$\left\langle \left( \sum_{k=1}^n \alpha_k u_k \right) - v, u_j \right\rangle = 0 \quad \text{für } j = 1, 2, \dots, n.$$

Nach Satz 2.10 ist  $\hat{u}$  genau dann ein  $U$ -Proximum an  $v$ , wenn  $\hat{u} - v \in U^\perp$ . Ist demnach  $\hat{u} - v \in U^\perp$ , dann gilt  $\langle \hat{u} - v, u_j \rangle = 0$  für alle  $j$  und die Normalgleichungen sind erfüllt. Seien umgekehrt die Normalgleichungen erfüllt und  $\tilde{u} \in U$  beliebig mit

$$\tilde{u} = \sum_{j=1}^n \beta_j u_j.$$

Daraus folgt

$$\langle \hat{u} - v, \tilde{u} \rangle = \left\langle \hat{u} - v, \sum_{j=1}^n \beta_j u_j \right\rangle = \sum_{j=1}^n \beta_j \underbrace{\langle \hat{u} - v, u_j \rangle}_{=0} = 0,$$

wodurch  $\hat{u} - v \in U^\perp$  nachgewiesen ist.  $\square$

Die Normalgleichungen (2.2) bilden ein lineares Gleichungssystem  $Ax = b$  für die unbekanntenen Koeffizienten  $\alpha_1, \dots, \alpha_n$  der Bestapproximation. Die beteiligte Matrix

$$A = \begin{pmatrix} \langle u_1, u_1 \rangle & \cdots & \langle u_1, u_n \rangle \\ \vdots & & \vdots \\ \langle u_n, u_1 \rangle & \cdots & \langle u_n, u_n \rangle \end{pmatrix}$$

nennt man Gramsche Matrix. Diese Matrix ist offensichtlich symmetrisch. Sie ist auch positiv definit wegen

$$x^\top Ax = \sum_{i,j=1}^n x_i x_j \langle u_i, u_j \rangle = \left\langle \sum_{i=1}^n x_i u_i, \sum_{j=1}^n x_j u_j \right\rangle = \left\| \sum_{j=1}^n x_j u_j \right\|^2 > 0$$

falls  $x \neq 0$ . Das lineare Gleichungssystem kann somit durch die Cholesky-Zerlegung gelöst werden.

Besonders günstig sind hier Orthonormalbasen des endlichdimensionalen Teilraums, d.h.

$$\langle u_i, u_j \rangle = \begin{cases} 0 & \text{für } i \neq j, \\ 1 & \text{für } i = j. \end{cases}$$

In diesem Fall reduziert sich die Gramsche Matrix zur Einheitsmatrix und die Koeffizientendarstellung der Bestapproximation wird direkt

$$\alpha_k = \langle v, u_k \rangle \quad \text{für } k = 1, 2, \dots, n.$$

Ist eine beliebige Basis von  $U$  gegeben, so lässt sich daraus eine Orthonormalbasis mit dem Verfahren von Gram-Schmidt konstruieren, siehe [5] Abschnitt 5.4.9.

### Approximationsprinzip von Korovkin

Wir betrachten in diesem Abschnitt den normierten Raum  $(C[a, b], \|\cdot\|_\infty)$  und setzen als Abkürzung  $I = [a, b]$  für  $a < b$ . Diskutiert werden lineare Operatoren auf diesem Vektorraum, wobei nur stetig Operatoren zugelassen sind. Die Stetigkeit ist äquivalent zur Beschränktheit des Operators.

**Def. 2.7** Eine Abbildung  $L : C(I) \rightarrow C(I)$  heißt

(i) linear, falls

$$L(\alpha f + \beta g) = \alpha Lf + \beta Lg$$

für alle  $f, g \in C(I)$  und alle  $\alpha, \beta \in \mathbb{R}$  gilt.

(ii) monoton, wenn die Folgerung

$$f \leq g \quad \Rightarrow \quad Lf \leq Lg \quad \text{für } f, g \in C(I)$$

gilt, wobei die Ungleichungen punktweise für alle  $x \in I$  gelten.

(iii) positiv, falls

$$0 \leq f \quad \Rightarrow \quad 0 \leq Lf \quad \text{für } f \in C(I).$$

(iv) beschränkt, falls

$$\sup \{ \|Lf\|_\infty : f \in C(I), \|f\|_\infty \leq 1 \} < \infty.$$

Man kann leicht zeigen, dass für eine lineare Abbildung dann Monotonie und Positivität äquivalent sind.

### Beispiele:

#### 1. Polynominterpolationsoperator

Zu einer Menge von Stützstellen  $a \leq x_0 < x_1 < \dots < x_n \leq b$  lautet der Operator der Polynominterpolation

$$(P_n f)(x) = \sum_{j=0}^n f(x_j) \cdot \prod_{k=0, k \neq j}^n \frac{x - x_k}{x_j - x_k}.$$

Dieser Operator ist linear und beschränkt. Jedoch ist der Operator nicht positiv für z.B. alle  $n \geq 3$ . Man definiere die stetige Funktion  $f \geq 0$  durch

$$f(x) = \left| \prod_{j=1}^n x - x_j \right|.$$

Das Interpolationspolynom  $P_n f$  besitzt dann die Nullstellen  $x_1, \dots, x_n$  und ist von null verschieden wegen  $(P_n f)(x_0) > 0$ . Es folgt, dass alle Nullstellen einfach sind und daher ein Vorzeichenwechsel stattfinden muss. Also nimmt  $P_n f$  auch negative Werte in  $[a, b]$  an.

#### 2. Bernsteinoperator

Auf  $I = [0, 1]$  definiert sich dieser Operator durch

$$(B_n f)(x) = \sum_{i=0}^n f\left(\frac{i}{n}\right) \cdot \binom{n}{i} x^i (1-x)^{n-i},$$

vergleiche (1.3). Der Operator basiert auf den Bernstein-Polynomen  $B_{i,n}$ . Für diese gilt  $0 \leq B_{i,n}(x) \leq 1$  für alle  $x \in I$ . Somit ist der Operator linear, beschränkt und positiv.

Um das Prinzip von Korovkin anzuwenden brauchen wir noch folgenden Begriff. Dabei sei  $e_1 \in C(I)$  die Funktion, welche konstant eins ist.

**Def. 2.8** *Eine Menge  $Q \subset C(I)$  mit  $Q = \{f_1, \dots, f_K\}$  und  $e_1 \in Q$  heißt Testmenge, wenn es eine Funktion  $p \in C(I \times I)$  gibt mit den Eigenschaften:*

(i) *Es existieren Funktionen  $a_1, \dots, a_K \in C(I)$  mit*

$$p(t, x) = \sum_{k=1}^K a_k(t) f_k(x).$$

(ii)  $p(t, x) \geq 0$  für alle  $(t, x) \in I \times I$ .

(iii)  $p(t, t) = 0$  für alle  $t \in I$ .

Hilfreich sind noch die nächsten beiden Begriffe.

**Def. 2.9** Zu  $g \in C(I \times I)$  lautet die Nullstellenmenge

$$Z(g) = \{(t, x) \in I \times I : g(t, x) = 0\}.$$

Zu  $f \in C(I)$  ergibt sich die zugehörige Differenzfunktion

$$d_f(t, x) = f(x) - f(t).$$

Offensichtlich gilt  $d_f(t, t) = f(t) - f(t) = 0$  und somit

$$(t, t) \in Z(d_f) \quad \text{für alle } t \in I.$$

Der folgende Satz stellt eine Verallgemeinerung eines Approximationssatzes von P.P. Korovkin aus dem Jahre 1953 dar.

**Satz 2.12** Sei  $(L_n)_{n \in \mathbb{N}}$  mit  $L_n : C(I) \rightarrow C(I)$  eine Folge monotoner linearer Operatoren und sei  $Q$  eine Testmenge mit zugehöriger Funktion  $p$ . Gilt

$$\lim_{n \rightarrow \infty} \|L_n f - f\|_\infty = 0 \quad \text{für alle } f \in Q,$$

dann folgt sogar

$$\lim_{n \rightarrow \infty} \|L_n f - f\|_\infty = 0 \quad \text{für alle } f \in C(I),$$

die die Bedingung  $Z(p) \subseteq Z(d_f)$  erfüllen.

Beweis:

(i) Wir zeigen zuerst, dass die Bedingung

$$\lim_{n \rightarrow \infty} \max_{t \in I} |(L_n d_f(t, \cdot))(t)| = 0 \tag{2.3}$$

hinreichend ist für

$$\lim_{n \rightarrow \infty} \|L_n f - f\|_\infty = 0.$$

Wegen  $d_f(t, \cdot) = f - f(t)e_1$  folgt  $f = f(t)e_1 + d_f(t, \cdot)$  und weiter

$$f - L_n f = f - f(t)L_n e_1 - L_n d_f(t, \cdot).$$



Für  $t \in I$  erhalten wir die gleichmäßige Abschätzung

$$\begin{aligned} |f(t) - (L_n f)(t)| &\leq |f(t)e_1 - f(t)(L_n e_1)(t)| + |(L_n d_f(t, \cdot))(t)| \\ &\leq \|f\|_\infty \|e_1 - L_n e_1\|_\infty + \max_{t \in I} |(L_n d_f(t, \cdot))(t)|. \end{aligned}$$

Da  $e_1 \in Q$  gilt  $\|e_1 - L_n e_1\|_\infty \rightarrow 0$  und die Bedingung (2.3) liefert  $\|f - L_n f\|_\infty \rightarrow 0$ .

(ii) Als zweites weisen wir nach, dass die Bedingung (2.3) für alle  $f \in C(I)$  mit der Eigenschaft  $Z(p) \subseteq Z(d_f)$  erfüllt ist.

Die Differenzfunktion ist stetig in  $x$  und  $t$ . Zu jedem  $\varepsilon > 0$  gibt es eine offene Umgebung  $U$  von  $Z(d_f)$  mit

$$|d_f(t, x)| < \varepsilon \quad \text{für alle } (t, x) \in U.$$

Für die Diagonale gilt

$$\{(t, x) \in I \times I : t = x\} \subseteq Z(d_f)$$

bei beliebigem  $f$ . Die Bedingung  $Z(p) \subseteq Z(d_f)$  impliziert

$$p(t, x) > 0 \quad \text{für alle } (t, x) \in U^C = (I \times I) \setminus U.$$

Ist  $U^C = \emptyset$ , dann folgt  $|d_f(t, x)| < \varepsilon$  für alle  $(t, x) \in I \times I$ . Anderenfalls ist  $U^C \neq \emptyset$  abgeschlossen und somit kompakt in  $I \times I$ . Dadurch existiert das Minimum

$$M = \min_{(t, x) \in U^C} p(t, x) > 0.$$

Wir folgern

$$|d_f(t, x)| \leq \|d_f\|_\infty \leq \|d_f\|_\infty \frac{p(t, x)}{M} \quad \text{für alle } (t, x) \in U^C$$

und somit

$$|d_f(t, x)| \leq \frac{\|d_f\|_\infty}{M} p(t, x) + \varepsilon \quad \text{für alle } (t, x) \in I \times I.$$

Anwendung des monotonen Operators  $L_n$  bezüglich  $x$  bei festem  $t$  zeigt

$$|(L_n d_f(t, \cdot))(x)| \leq \frac{\|d_f\|_\infty}{M} (L_n p(t, \cdot))(x) + \varepsilon (L_n e_1)(x)$$

und mit  $x = t$

$$|(L_n d_f(t, \cdot))(t)| \leq \frac{\|d_f\|_\infty}{M} \max_{t \in I} (L_n p(t, \cdot))(t) + \varepsilon \|L_n e_1\|_\infty$$

gleichmäßig für alle  $t \in I$ . Wegen  $p(t, t) = 0$  für alle  $t$  gilt

$$\sum_{k=1}^K a_k(t) f_k(t) = 0 \quad \text{für alle } t \in I.$$

Andererseits ist

$$L_n p(t, \cdot) = L_n \left( \sum_{k=1}^K a_k(t) f_k(\cdot) \right) = \sum_{k=1}^K a_k(t) (L_n f_k)$$

und damit

$$(L_n p(t, \cdot))(t) = \sum_{k=1}^K a_k(t) (L_n f_k)(t) - \sum_{k=1}^K a_k(t) f_k(t) = \sum_{k=1}^K a_k(t) [(L_n f_k)(t) - f_k(t)].$$

Die Konvergenz der Folge  $(L_n)_{n \in \mathbb{N}}$  auf  $\text{span}(Q)$  impliziert

$$\lim_{n \rightarrow \infty} \max_{t \in I} (L_n p(t, \cdot))(t) = 0.$$

Desweiteren gilt

$$\|L_n e_1\|_\infty \leq \|L_n e_1 - e_1\|_\infty + \|e_1\|_\infty = \|L_n e_1 - e_1\|_\infty + 1$$

und mit  $e_1 \in Q$  folgt die Konvergenz der Operatoren. Also ist die betrachtete Folge  $(\|L_n e_1\|_\infty)_{n \in \mathbb{N}}$  beschränkt. Wir folgern

$$\lim_{n \rightarrow \infty} \max_{t \in I} |(L_n d_f(t, \cdot))(t)| = 0,$$

wodurch der Beweis abgeschlossen ist.  $\square$

### **Beispiel:** Bernstein-Operatoren

Wir folgern mit dem Prinzip von Korovkin nun die gleichmäßige Konvergenz der Bernstein-Polynome  $B_n$  aus (1.3) gegen die zu approximierende Funktion  $f$  in  $I = [0, 1]$ . Als Testmenge verwenden wir  $Q = \{e_1, e_2, e_3\}$  mit

$$e_1 = 1, \quad e_2 = x, \quad e_3 = x^2.$$

Wir definieren

$$p(t, x) = (t - x)^2 = t^2 - 2tx + x^2$$

und verifizieren

$$(i) \quad p(t, x) = a_1(t)e_1(x) + a_2(t)e_2(x) + a_3(t)e_3(x) = a_1(t) + a_2(t)x + a_3(t)x^2 \text{ mit} \\ a_1(t) = t^2, \quad a_2(t) = -2t, \quad a_3(t) = 1,$$

$$(ii) \quad p(t, x) = (t - x)^2 \geq 0,$$

$$(iii) \quad p(t, t) = (t - t)^2 = 0.$$

Desweiteren gilt  $p(t, x) > 0$  für  $x \neq t$ , wodurch  $Z(p) \subseteq Z(d_f)$  für beliebiges  $f$  folgt.

Man kann zeigen, dass gilt

$$\lim_{n \rightarrow \infty} \|L_n e_k - e_k\|_\infty = 0 \quad \text{für } k = 1, 2, 3.$$

Für  $e_1$  folgt dies direkt aus der Zerlegung der Eins in (1.2) und für  $e_2, e_3$  aus den ersten beiden Formeln im Beweis von Lemma 1.1. Somit sind die Voraussetzungen für Satz 2.12 erfüllt und es folgt die Konvergenz der Approximationen für beliebiges  $f \in C(I)$ . Dies stellt einen alternativen Beweis von Satz 1.8 dar.

## 2.2 Fourier-Reihen

In diesem Abschnitt betrachten wir periodische Funktionen mit fester Periode  $L > 0$ . O.E.d.A. sei  $L = 2\pi$ , wodurch wir uns auf das kompakte Intervall  $[-\pi, +\pi]$  zurückziehen können.

### Trigonometrische Polynome

**Def. 2.10** *Wir bezeichnen die Menge der stetigen bzw. quadratintegriblen periodischen Funktionen mit Periode  $2\pi$  durch*

$$C_p = \{f \in C(\mathbb{R}) : f(x + k2\pi) = f(x) \text{ für alle } k \in \mathbb{Z} \text{ und } x \in \mathbb{R}\},$$

$$L_p^2 = \{f : \mathbb{R} \rightarrow \mathbb{R} : f|_{[-\pi, +\pi]} \in L^2([-\pi, +\pi]), \\ f(\cdot + k2\pi) = f(\cdot) \text{ für alle } k \in \mathbb{Z}\}.$$

Ist  $f \in C[-\pi, +\pi]$  mit  $f(-\pi) = f(\pi)$ , dann kann  $f$  erweitert werden zu einer Funktion in  $C_p$ . Jedes  $f \in L^2[-\pi, +\pi]$  kann zu einer Funktion in  $L_p^2$  fortgesetzt werden. Umgekehrt kann jedes  $f \in L_p^2$  auf den Raum  $L^2[-\pi, +\pi]$  eingeschränkt werden. Es gilt  $C_p \subset L_p^2$ .

**Def. 2.11** *Ein reelles trigonometrisches Polynom ist eine Funktion der Gestalt*

$$f(x) = \frac{a_0}{2} + \sum_{k=1}^n a_k \cos(kx) + b_k \sin(kx) \quad (2.4)$$

mit Koeffizienten  $a_0, \dots, a_n \in \mathbb{R}$  und  $b_1, \dots, b_n \in \mathbb{R}$ . Es bezeichnet  $n$  den Grad des Polynoms ( $a_n \neq 0$  oder  $b_n \neq 0$ ). Sei  $\mathcal{T}_n$  die Menge aller trigonometrischen Polynome mit Grad höchstens  $n$ .

Offensichtlich ist jedes trigonometrische Polynom in  $C_p$ . Die trigonometrischen Polynome  $\mathcal{T}_n$  bilden einen Untervektorraum, d.h.  $\mathcal{T}_n \subset C_p \subset L_p^2$ .

Auf  $L^2[-\pi, +\pi]$  verwenden wir das modifizierte Skalarprodukt

$$\langle f, g \rangle = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x)g(x) \, dx \quad \text{für } f, g \in L^2[-\pi, +\pi]. \quad (2.5)$$

Die von diesem Skalarprodukt erzeugte Norm ist äquivalent zur üblichen  $L^2$ -Norm und wir bezeichnen sie mit  $\|\cdot\|_{L^2}$ . Das System

$$\left\{ \frac{1}{\sqrt{2}}, \sin(x), \cos(x), \sin(2x), \cos(2x), \dots, \sin(nx), \cos(nx) \right\}.$$

ist eine Orthonormalbasis von  $\mathcal{T}_n$  bezüglich des Skalarprodukts (2.5).

**Def. 2.12** Zu  $f \in L^2[-\pi, +\pi]$  lauten die zugehörigen Fourier-Koeffizienten

$$\begin{aligned} a_k &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(kx) \, dx & \text{für } k = 0, 1, 2, \dots, \\ b_k &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(kx) \, dx & \text{für } k = 1, 2, \dots \end{aligned}$$

Die Fourier-Koeffizienten implizieren die Operatoren  $S_n$  mit

$$(S_n f)(x) = \frac{a_0}{2} + \sum_{k=1}^n a_k \cos(kx) + b_k \sin(kx).$$

Die Eigenschaft  $f \in L^2[-\pi, +\pi]$  garantiert die Existenz der Fourierkoeffizienten als reelle Zahlen. Es gilt

$$a_k = \langle f(x), \cos(kx) \rangle \quad \text{für } k \geq 0, \quad b_k = \langle f(x), \sin(kx) \rangle \quad \text{für } k \geq 1.$$

Speziell für  $k = 0$  haben wir  $a_0 = \sqrt{2} \langle f(x), \frac{1}{\sqrt{2}} \rangle$  und es ist  $\frac{a_0}{2} = \langle f(x), \frac{1}{\sqrt{2}} \rangle \frac{1}{\sqrt{2}}$ .

## Konvergenz der Fourier-Reihe

Mit Satz 2.11 folgt, dass  $S_n f$  die Bestapproximation von  $f$  in  $\mathcal{T}_n$  bezüglich der Norm  $\|\cdot\|_{L^2}$  ist. Daraus ergibt sich die Frage, ob diese Bestapproximationen auch gegen die Funktion konvergieren.

**Satz 2.13** Für  $f \in L^2[-\pi, +\pi]$  konvergiert die Fourier-Reihe im quadratischen Mittel gegen  $f$ , d.h. es gilt

$$\lim_{n \rightarrow \infty} \|f - S_n f\|_{L^2} = 0.$$

Beweis: siehe [23], Satz V4.9.

Somit konvergiert die Fourier-Reihe im quadratischen Mittel gegen die zu approximierende Funktion. Eine zentrale Frage der Approximationstheorie war, ob für eine stetige Funktion die Fourier-Reihe auch gleichmäßig konvergiert.

**Satz 2.14** Sei  $f \in C_p$  eine stückweise stetig differenzierbare Funktion, d.h. es gibt eine Unterteilung  $-\pi = x_0 < x_1 < \dots < x_m = \pi$ , so dass  $f|_{(x_{j-1}, x_j)}$  für  $j = 1, \dots, m$  stetig differenzierbar in  $[x_{j-1}, x_j]$  ist. Dann konvergiert die Fourier-Reihe gleichmäßig gegen  $f$ .

Beweis:

Es bezeichne  $g_j : [x_{j-1}, x_j] \rightarrow \mathbb{R}$  die erste Ableitung von  $f|_{(x_{j-1}, x_j)}$ . Dann sei  $g : \mathbb{R} \rightarrow \mathbb{R}$  die periodische Funktion, die auf  $(x_{j-1}, x_j)$  mit  $g_j$  für alle  $j$  übereinstimmt. Die Fourier-Koeffizienten von  $g$  erfüllen die Besselsche Ungleichung

$$\frac{|a_0(g)|^2}{4} + \sum_{k=1}^{\infty} |a_k(g)|^2 + |b_k(g)|^2 \leq \|g\|^2.$$

Für  $k \neq 0$  zeigt partielle Integration

$$\int_{x_{j-1}}^{x_j} f(x) \cos(kx) \, dx = \left[ \frac{1}{k} f(x) \sin(kx) \right]_{x=x_{j-1}}^{x=x_j} - \frac{1}{k} \int_{x_{j-1}}^{x_j} g(x) \sin(kx) \, dx$$

und damit weil sich die Randterme gegenseitig aufheben

$$a_k(f) = \frac{1}{\pi} \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(x) \cos(kx) \, dx = -\frac{1}{k\pi} \sum_{j=1}^m \int_{x_{j-1}}^{x_j} g(x) \sin(kx) \, dx = -\frac{1}{k} b_k(g).$$

Analog folgt  $b_k(f) = \frac{1}{k} a_k(g)$ . Man beachte, dass  $g \in L^2[-\pi, +\pi]$  ist.

Für alle  $v, w \in \mathbb{R}$  gilt  $vw \leq \frac{1}{2}(v^2 + w^2)$ . Wir erhalten

$$|a_k(f)| = \frac{|b_k(g)|}{k} \leq \frac{1}{2} \left( \frac{1}{k^2} + |b_k(g)|^2 \right), \quad |b_k(f)| = \frac{|a_k(g)|}{k} \leq \frac{1}{2} \left( \frac{1}{k^2} + |a_k(g)|^2 \right).$$

Da die Reihen

$$\sum_{k=1}^{\infty} \frac{1}{k^2}, \quad \sum_{k=1}^{\infty} |a_k(g)|^2, \quad \sum_{k=1}^{\infty} |b_k(g)|^2$$

wegen der Besselschen Ungleichung alle konvergent sind, folgt

$$\frac{|a_0(f)|}{2} + \sum_{k=1}^{\infty} |a_k(f)| + |b_k(f)| < \infty.$$

Die Fourier-Reihe von  $f$  konvergiert wegen  $|\sin(kx)| \leq 1$ ,  $|\cos(kx)| \leq 1$  für alle  $k$  damit absolut und gleichmäßig. Es existiert also eine stetige Funktion  $\tilde{f}$  als Grenzwert. Aus gleichmäßiger Konvergenz folgt auch die Konvergenz im quadratischen Mittel. Die Fourier-Reihe von  $f$  konvergiert jedoch auch gegen  $f$  selbst nach Satz 2.13. Da der Grenzwert eindeutig ist folgt  $f = \tilde{f}$  in  $L^2[-\pi, +\pi]$ . Da  $f$  und  $\tilde{f}$  stetig sind folgt daraus  $f \equiv \tilde{f}$ .  $\square$

Man beachte, dass bei einer stückweise stetigen Funktion hier vorausgesetzt ist, dass an den Unstetigkeitsstellen der links- und rechtsseitige Grenzwert jeweils existiert.

Desweiteren hängt die Konvergenzgeschwindigkeit von der Glattheit der Funktion ab. Für das Abklingverhalten der Fourier-Koeffizienten gibt der folgende Satz eine Auskunft.

**Satz 2.15** *Ist  $f \in C_p$  eine  $m$ -mal stetig differenzierbare Funktion ( $m \geq 1$ ), dann gibt es eine Konstante  $C_f > 0$  und ein  $n^* \in \mathbb{N}$  mit*

$$\max\{|a_n|, |b_n|\} \leq C_f \cdot \frac{1}{n^m} \quad \text{für } n \geq n^*$$

und desweiteren gilt

$$\|f - S_n f\|_{\infty} \leq K \|f^{(m)}\|_{\infty} \frac{\ln n}{n^m} \quad \text{für } n \geq 2$$

mit einer Konstanten  $K > 0$  und der Maximumnorm  $\|\cdot\|_{\infty}$ .

Der Beweis der ersten Aussage erfolgt analog zum Beweis von Satz 2.14 mit partieller Integration. Zum Beweis der zweiten Aussage siehe [20].

Jedoch ist die Stetigkeit allein nicht hinreichend für die gleichmäßige Konvergenz. Dazu betrachten wir ein Gegenbeispiel. Die Reihe

$$f(x) = \sum_{k=1}^{\infty} \frac{\sin(2^{k^3} x)}{k^2} \quad \text{für } 0 \leq x \leq \pi$$

konvergiert absolut und gleichmäßig. Der Grenzwert  $f$  existiert damit und ist stetig. Es gilt  $f(0) = f(\pi) = 0$ . Wir erweitern diese Funktion zu  $\tilde{f} : [-\pi, \pi] \rightarrow \mathbb{R}$

durch  $\tilde{f}(x) = f(|x|)$ , d.h.  $\tilde{f}$  ist eine gerade Funktion. Die Fourier-Reihe von  $\tilde{f}$  besteht dann nur aus Cosinus-Termen. Man kann zeigen, dass diese Fourier-Reihe bei  $x = 0$  divergiert, siehe [4].

## Approximation mit Fejér-Operatoren

Erfreulicherweise ergibt sich die gleichmäßige Konvergenz mit trigonometrischen Polynomen für eine leichte Modifikation der Fourier-Summen. Dabei wird aus den Fourier-Summen ein arithmetisches Mittel konstruiert mittels der Cesàro-Summation.

**Def. 2.13** Zu  $f \in L^2[-\pi, +\pi]$  wird mit den Operatoren  $S_n$  der Fourier-Summen der Fejér-Operator

$$F_n f = \frac{1}{n} \sum_{j=0}^{n-1} S_j f$$

für  $n \in \mathbb{N}$  gebildet.

Desweiteren existieren für beide Operatoren Kerne.

**Def. 2.14** Der  $n$ -te Dirichlet-Kern ist die Funktion

$$K_n^D(x) = \frac{\sin((2n+1)\frac{x}{2})}{\sin(\frac{x}{2})}$$

und der  $n$ -te Fejér-Kern ist die Funktion

$$K_n^F(x) = \frac{1}{n} \left( \frac{\sin(n\frac{x}{2})}{\sin(\frac{x}{2})} \right)^2.$$

Abbildung 8 zeigt ein Beispiel für diese Funktionen. Jetzt können wir die Operatoren durch Integrale mit den Kernen darstellen. Man beachte, dass diese Repräsentation nur in der Theorie von Interesse ist. Für die numerische Auswertung der Operatoren werden diese Formeln nicht eingesetzt.

**Satz 2.16** Für  $f \in C_p$  gilt

$$(S_n f)(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) K_n^D(t-x) dt$$

und

$$(F_n f)(x) = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) K_n^F(t-x) dt$$

für alle  $x \in \mathbb{R}$  und alle  $n \geq 0$ .

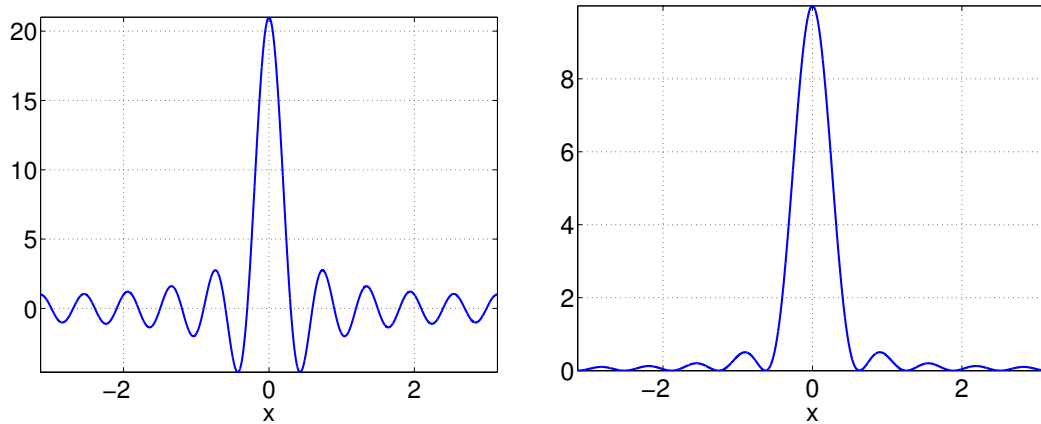


Abbildung 8: Dirichlet-Kern (links) und Fejér-Kern (rechts) für  $n = 10$  in  $[-\pi, +\pi]$ .

Beweis:

Wir erhalten die Hilfsformel

$$\frac{1}{2} + \sum_{k=1}^n \cos(kx) = \frac{1}{2} \cdot \frac{\sin((n + \frac{1}{2})x)}{\sin(\frac{x}{2})} \quad \text{für alle } n$$

aus

$$\begin{aligned} \frac{1}{2} + \sum_{k=1}^n \cos(kx) &= \frac{1}{2} \sum_{k=-n}^n \cos(kx) &= \frac{1}{2} \sum_{k=-n}^n e^{ikx} \\ &= \frac{1}{2} e^{-inx} \sum_{k=0}^{2n} e^{ikx} &= \frac{1}{2} e^{-inx} \frac{e^{i(2n+1)x} - 1}{e^{ix} - 1} \\ &= \frac{1}{2} \frac{e^{i(n+\frac{1}{2})x} - e^{-i(n+\frac{1}{2})x}}{e^{i\frac{x}{2}} - e^{-i\frac{x}{2}}} &= \frac{1}{2} \frac{\sin((n + \frac{1}{2})x)}{\sin(\frac{x}{2})}, \end{aligned}$$

wobei die Formeln für Sinus und Cosinus im Komplexen sowie der Wert der endlichen geometrischen Reihe verwendet wurden.

Mit den Additionstheoremen für trigonometrische Funktionen, der Periodizität



und der obigen Formel folgt

$$\begin{aligned}
(S_n f)(x) &= \frac{a_0}{2} + \sum_{k=1}^n a_k \cos(kx) + b_k \sin(kx) \\
&= \frac{1}{\pi} \int_{-\pi}^{\pi} f(t) \left[ \frac{1}{2} + \sum_{k=1}^n \cos(kt) \cos(kx) + \sin(kt) \sin(kx) \right] dt \\
&= \frac{1}{\pi} \int_{-\pi}^{\pi} f(t) \left[ \frac{1}{2} + \sum_{k=1}^n \cos(k(t-x)) \right] dt \\
&= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) \frac{\sin((n+\frac{1}{2})(t-x))}{\sin(\frac{t-x}{2})} dt.
\end{aligned}$$

Das Additionstheorem  $\sin(\alpha) \sin(\beta) = \frac{1}{2}(\cos(\alpha - \beta) - \cos(\alpha + \beta))$  zeigt uns

$$\begin{aligned}
\sum_{j=0}^{n-1} \sin((j+\frac{1}{2})x) \sin(\frac{x}{2}) &= \frac{1}{2} \sum_{j=0}^{n-1} \cos(jx) - \cos((j+1)x) \\
&= \frac{1}{2} [1 - \cos(nx)] = \sin^2(\frac{nx}{2}).
\end{aligned}$$

Aus Def. 2.13 und der obigen Formel für  $S_n$  folgt

$$\begin{aligned}
(F_n f)(x) &= \frac{1}{\pi} \int_{-\pi}^{\pi} f(t) \frac{1}{n} \left[ \sum_{j=0}^{n-1} \frac{\sin((j+\frac{1}{2})(t-x))}{2 \sin(\frac{t-x}{2})} \right] dt, \\
&= \frac{1}{2\pi} \int_{-\pi}^{\pi} f(t) \frac{1}{n} \frac{\sin^2(\frac{n(t-x)}{2})}{\sin^2(\frac{t-x}{2})} dt,
\end{aligned}$$

wodurch die Aussage gezeigt ist. □

Der Fejér-Operator ist offensichtlich linear. Mit der Darstellung des Operators aus Satz 2.16 folgt, dass der Operator auch positiv ist, denn der Kern ist eine nichtnegative Funktion. Dadurch können wir Satz 2.12 anwenden, wobei wir uns auf den Raum  $C[-\pi, +\pi]$  zurückziehen. Wir definieren eine Testmenge  $Q : \{e_1, e_2, e_3\} \in C_p$  mit

$$e_1 = 1, \quad e_2 = \cos(x), \quad e_3 = \sin(x)$$

sowie die Funktion

$$p(t, x) = 1 - \cos(t-x) = 1 - \cos(t) \cos(x) - \sin(t) \sin(x).$$

Wir verifizieren

- (i)  $p(t, x) = a_1(t)e_1(x) + a_2(t)e_2(x) + a_3(t)e_3$   
mit  $a_1(t) = 1$ ,  $a_2(t) = -\cos(t)$ ,  $a_3(t) = -\sin(t)$ ,
- (ii)  $p(t, x) \geq 0$  wegen  $-1 \leq \cos(t - x) \leq 1$ ,
- (iii)  $p(t, t) = 1 - \cos(0) = 0$ .

Die Nullstellenmenge von  $p$  in  $[-\pi, +\pi]^2$  ist

$$Z(p) = D \cup \{(-\pi, +\pi), (+\pi, -\pi)\}.$$

mit der Diagonalen  $D = \{(t, t) : t \in [-\pi, +\pi]\}$ . Für die Differenzfunktion gilt immer  $D \subseteq Z(d_f)$ . Da  $f$  periodisch ist, d.h.  $f(-\pi) = f(\pi)$ , folgt hier  $\{(-\pi, +\pi), (+\pi, -\pi)\} \subset Z(d_f)$ . Die Eigenschaft  $Z(p) \subseteq Z(d_f)$  ist somit für alle relevanten Funktionen  $f$  erfüllt. Es verbleibt zu zeigen, dass die Fejér-Operatoren für jede Funktion aus der Testmenge  $Q$  gleichmäßig konvergieren. Dies folgt aus

$$S_n e_1 = 1 \text{ für } n \geq 0, \quad S_n e_2 = \cos(x) \text{ für } n \geq 1, \quad S_n e_3 = \sin(x) \text{ für } n \geq 1,$$

denn wir erhalten

$$F_n e_1 = 1, \quad F_n e_2 = \frac{n-1}{n} \cos(x), \quad F_n e_3 = \frac{n-1}{n} \sin(x) \quad \text{für } n \geq 1$$

und somit

$$\lim_{n \rightarrow \infty} \|e_i - F_n e_i\|_\infty = 0 \quad \text{für } i = 1, 2, 3.$$

Wir notieren als Schlussfolgerung:

**Satz 2.17** *Ist  $f \in C_p$ , dann gilt mit den Fejér-Operatoren*

$$\lim_{n \rightarrow \infty} \|f - F_n f\|_\infty = 0,$$

wobei  $\|\cdot\|_\infty$  die Maximumnorm auf  $[-\pi, +\pi]$  bezeichnet.

Die nachfolgende direkte Konsequenz daraus nennt man auch den Weierstraßschen Approximationssatz für periodische Funktionen.

**Korollar 2.3** *Zu jeder Funktion  $f \in C_p$  existiert eine Folge von trigonometrischen Polynomen, welche gleichmäßig gegen  $f$  konvergiert.*

## Literatur

- [1] O. Christensen, K.L. Christensen: Approximation Theory: From Taylor Polynomials to Wavelets. (2. Aufl.) Birkhäuser, Boston, 2005.
- [2] R.A. DeVore, G.G. Lorentz: Constructive Approximation. Springer, Berlin, 1993.
- [3] G. Faber: Über die interpolatorische Darstellung stetiger Funktionen. Jber. Deutsch. Math.-Verein. 23 (1914), 192–210.
- [4] L. Fejér: Beispiele stetiger Funktionen mit divergenter Fourierreihe. Journal für Mathematik. Bd. 137, Heft 1 (1909).
- [5] G. Fischer: Lineare Algebra. (18. Aufl.) Springer, Wiesbaden, 2014.
- [6] O. Forster: Analysis 2. (10. Aufl.) Springer, Wiesbaden, 2013.
- [7] D. Meschede (Hrsg.): Gerthsen Physik. (23. Aufl.) Springer, Berlin 2006.
- [8] M. Günther, A. Jüngel: Finanzderivate mit MATLAB. (2. Aufl.) Vieweg + Teubner, Wiesbaden, 2010.
- [9] C.A. Hall: On error bounds for spline interpolation. J. Approximation Theory 1 (1968), 209–218.
- [10] G. Hämmerlin, K.-H. Hoffmann: Numerische Mathematik. (4. Aufl.) Springer, Berlin, 1994.
- [11] H.-P. Halvorsen: System Identification and Estimation in LabVIEW. Tutorial, Telemark University College, Norwegen, 2011.
- [12] G. Hübner: Stochastik. (5. Aufl.) Vieweg + Teubner, Wiesbaden, 2009.
- [13] A.K. Louis, P. Maaß, A. Rieder: Wavelets. Teubner, Stuttgart, 1994.

- [14] J. Marcinkiewicz: Sur l'interpolation d'opérations. C. R. Acad. des Sciences 208 (1939), 1272–1273.
- [15] R. Plato: Numerische Mathematik kompakt. (4. Aufl.) Vieweg + Teubner, Wiesbaden, 2010.
- [16] C. Reinsch: Smoothing by spline functions. Numer. Math. 10 (1967), 177–183.
- [17] H.R. Schwarz, N. Köckler: Numerische Mathematik. (8. Aufl.) Vieweg + Teubner, Wiesbaden, 2011.
- [18] A. Sharma, A. Meir: Degree of approximation of spline interpolation. J. Math. Mech. 15 (1966), 749–768.
- [19] J. Stoer: Numerische Mathematik 1. (9. Aufl.) Springer, Berlin, 2005.
- [20] G. Steidl, M. Tasche: Schnelle Fouriertransformation – Theorie und Anwendungen. Lehrbriefe der Fern Universität Hagen, 1996.
- [21] W. Törning, P. Spellucci: Numerische Mathematik für Ingenieure und Physiker, Band 1, Numerische Methoden der Algebra. (2. Aufl.) Springer, Berlin, 1988.
- [22] W. Törning, P. Spellucci: Numerische Mathematik für Ingenieure und Physiker, Band 2, Numerische Methoden der Analysis. (2. Aufl.) Springer, Berlin, 1990.
- [23] D. Werner: Funktionalanalysis. (7. Aufl.) Springer, Berlin, 2011.