

TECHNISCHE UNIVERSITÄT DRESDEN

Skript:

Numerik partieller Differentialgleichungen

Verfasser

Franziska Kühn

Daten

Prof. Dr. Hans-Görg Roos  
Wintersemester 2012/13  
Hauptstudium

# Inhaltsverzeichnis

<b>Einführung</b>	<b>2</b>
<b>1 Schwache Lösungen partieller Differentialgleichungen</b>	<b>4</b>
<b>2 Variationsgleichungen</b>	<b>11</b>
<b>3 Ritz-Galerkin-Verfahren</b>	<b>17</b>
<b>4 Finite-Elemente-Methode</b>	<b>20</b>
4.1 1D . . . . .	20
4.2 2D . . . . .	22
4.3 Einfache FE-Diskretisierung auf Standardgittern . . . . .	27
<b>5 Konvergenzanalyse</b>	<b>29</b>
5.1 Differenzenverfahren . . . . .	29
5.2 Finite Elemente . . . . .	31
<b>6 Nichtkonforme Aspekte</b>	<b>43</b>
<b>7 Weitere Verfahren</b>	<b>47</b>
7.1 Finite-Volumen-Methode . . . . .	47
7.2 Discontinuous Galerkin-Verfahren für elliptische Probleme . . . . .	50
<b>8 Parabolische Randwertaufgaben</b>	<b>55</b>

# Einführung

In vielen Anwendungen: partielle Differentialgleichungen 2. Ordnung. Beispiele:

- (i). Poisson-Gleichung:  $-\Delta u = f$  (stationäre Gleichung, d.h. unabhängig von  $t$ )
- (ii). Wärmeleitungsgleichung:  $u_t - \Delta u = f$  (instationäre Gleichung)
- (iii). Schwingungsgleichung (Wellengleichung):  $u_{tt} - \Delta u = f$
- (iv). Black-Scholes-Gleichung ( $\rightarrow$  Finanzderivate):

$$V_t + \frac{1}{2}\sigma^2 \cdot s^2 \cdot V_{ss} + (r - \delta) \cdot s \cdot V_s - r \cdot V = 0$$

- (v). Strömungsmechanik: inkompressible Strömungen (Navier-Stokes-Gleichung)

$$\begin{aligned}u_t + u \cdot \nabla u + \nabla p - \nu \cdot \Delta u &= f \\ \operatorname{div} u &= 0\end{aligned}$$

wobei  $u$  Geschwindigkeitsvektor,  $p$  Druck (nichtlineare Gleichung)

Vorlesung	Skript
$C_0^\infty(\Omega)$	$C_c^\infty(\Omega)$
$V^*$	$V'$ (Dualraum)
$h_K$	$\operatorname{diam} K$ (Durchmesser)

Tabelle 1: abweichende Notationen

# 1

## Schwache Lösungen partieller Differentialgleichungen

Wozu schwache Formulierung partieller Differentialgleichungen? „Klassische“ Lösungen partieller Differentialgleichungen: Betrachte  $-\Delta u = f$  in  $\Omega \subseteq \mathbb{R}^n$  offen,  $u = \varphi$  auf  $\partial\Omega$  mit  $\varphi, f$  stetig (Randwertaufgabe). Suche  $u \in C^2(\Omega) \cap C(\bar{\Omega})$ . Stellt man keine weiteren Bedingungen an  $\Omega$ , kann man die Existenz klassischer Lösungen nicht beweisen.

**Beispiel**  $\Delta u = 0$ ,  $u|_{\partial\Omega} = \varphi$ ,  $\Omega$  erfüllt „barrier property“, dann folgt die Existenz der klassischen Lösung (hinreichend für „barrier property“ ist z.B.  $\Omega$  strikt konvex).

**Beispiel**

$$-u'' + c \cdot u = f \qquad u(0) = u(1) = 0 \qquad (*)$$

Multipliziere mit  $v$  (mit  $v(0) = v(1) = 0$ ) und  $\int_0^1$ , dann

$$\begin{aligned} & - \int_0^1 u'' \cdot v + \int_0^1 c \cdot u \cdot v = \int_0^1 f \cdot v \\ \stackrel{\text{part. Int}}{\Rightarrow} & - \underbrace{[u' \cdot v]_0^1}_0 + \int_0^1 u' \cdot v' + \int_0^1 c \cdot u \cdot v = \int_0^1 f \cdot v \end{aligned}$$

$u$  heißt schwache (verallgemeinerte) Lösung von (\*), falls die letzte Gleichung für alle  $v$  mit  $v(1) = v(0) = 0$  gilt.

Aus welchem Raum muss man  $u, v$  wählen? Die schwache Ableitung von  $u$  muss quadratisch integrierbar sein! Neuen Ableitungsbegriff...

**Definition** Sei  $\Omega \subseteq \mathbb{R}^n$  und  $p \in [1, \infty)$ .

$$L^p(\Omega) := \left\{ f : \Omega \rightarrow \bar{\mathbb{R}}; f \text{ mb, } \int_{\Omega} |f|^p < \infty \right\}$$

Norm:

$$\|f\|_{L^p} := \left( \int_{\Omega} |f|^p \right)^{\frac{1}{p}}$$

Für  $p = \infty$ :

$$L^p(\Omega) := \{ f : \Omega \rightarrow \bar{\mathbb{R}}; f \text{ mb, } \text{esssup } f < \infty \}$$

Mit diesen Normen ist  $(L^p, \|\cdot\|_{L^p})$  für  $p \in [1, \infty]$  Banachraum, für  $p = 2$  Hilbertraum mit Skalarprodukt  $(f_1, f_2) := \int_{\Omega} f_1 \cdot f_2$ .

**Bemerkung** In Prä-Hilberträumen gilt

$$|(f_1, f_2)| \leq \|f_1\| \cdot \|f_2\|$$

(Cauchy-Schwarz-Ungleichung)

Sei  $u \in C^1(\bar{\Omega})$ ,  $v \in C^\infty(\Omega)$ . Dann gilt mittels partieller Integration:

$$\int_{\Omega} \frac{\partial u}{\partial x_j} \cdot v = \int_{\partial\Omega} u \cdot v \cdot \cos\langle n, e^j \rangle - \int_{\Omega} u \cdot \frac{\partial v}{\partial x_j} \quad (1.1)$$

mit  $n$  Normaleneinheitsvektor bzgl.  $\partial\Omega$  und  $e^i$  der  $i$ -te Einheitsvektor. Speziell für  $v \in C_c^\infty(\Omega)$ :

$$\int_{\Omega} \frac{\partial u}{\partial x_i} \cdot v = - \int_{\Omega} u \cdot \frac{\partial v}{\partial x_i}$$

**Definition** Sei  $u : \Omega \rightarrow \bar{\mathbb{R}}$  messbar. Ist  $w : \Omega \rightarrow \bar{\mathbb{R}}$  messbar und gilt

$$\int_{\Omega} w \cdot v = - \int_{\Omega} u \cdot \frac{\partial v}{\partial x_i} \quad (v \in C_c^\infty(\Omega))$$

heißt  $w$  schwache Ableitung von  $u$  nach  $x_i$ .

**Beispiel** Sei  $\Omega := (0, 1)$  und  $a \in (0, 1)$ ,

$$f(x) := \begin{cases} f_1(x) & x \in (0, a) \\ f_2(x) & x \in (a, 1) \end{cases}$$

mit  $f_1, f_2 \in C^1(\Omega)$ . Ist  $f$  schwach differenzierbar?

Beweis: Sei  $v \in C_c^\infty(0, 1)$ , dann

$$\begin{aligned} \int_0^1 f \cdot v' &= \int_0^a f_1 \cdot v' + \int_a^1 f_2 \cdot v' \\ &= (f_1(a) - f_2(a)) \cdot v(a) + \int_0^a f_1' \cdot v + \int_a^1 f_2' \cdot v \\ &= (f_1(a) - f_2(a)) \cdot v(a) + \int_0^1 w \cdot v \end{aligned}$$

für  $w := f_1 \cdot 1_{(0,a)} + f_2 \cdot 1_{(a,1)}$ . Also ist  $f$  schwach differenzierbar genau dann, wenn  $f_1(a) = f_2(a)$ .  $\square$

**Bemerkung** Analog definiert man schwache Ableitungen höherer Ordnung, dazu Multiindexschreibweise

$$\partial^\alpha u := \frac{\partial^{|\alpha|}}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} u \quad |\alpha| := \sum_{j=1}^n \alpha_j$$

für  $\alpha \in \mathbb{N}_0^n$ .

**Definition** (i). Sei  $u : \Omega \rightarrow \mathbb{R}^n$  messbar.  $w : \Omega \rightarrow \mathbb{R}^n$  messbar heißt verallgemeinerte Ableitung zum Multiindex  $\alpha$ , falls

$$\int_{\Omega} w \cdot v = (-1)^{|\alpha|} \cdot \int_{\Omega} u \cdot \partial^\alpha v \quad (v \in C_c^\infty(\Omega))$$

(ii). Sobolevraum:

$$H^1(\Omega) := \{u \in L^2(\Omega); \forall j = 1, \dots, n : \partial_j u \in L^2\}$$

mit Norm

$$\|u\|_1 := \left( \int_{\Omega} u^2 + \int_{\Omega} (\nabla u)^2 \right)^{\frac{1}{2}}$$

und Skalarprodukt

$$(u, v)_1 := (u, v) + (\nabla u, \nabla v) = \int_{\Omega} u \cdot v + \sum_{j=1}^n \int_{\Omega} \partial_j u \cdot \partial_j v$$

**Bemerkung** (i).  $H^1(\Omega)$  ist ein Hilbertraum.

(ii). Allgemeiner:

$$W_p^\ell(\Omega) := \{u \in L^p(\Omega); \forall \alpha \in \mathbb{N}_0^n, |\alpha| \leq \ell : \partial^\alpha u \in L^p(\Omega)\}$$

mit Norm

$$\|u\|_{W_p^\ell} := \left( \int_\Omega \sum_{|\alpha| \leq \ell} |\partial^\alpha u|^p \right)^{\frac{1}{p}}$$

für  $p \in [1, \infty]$ ,  $\ell \in \mathbb{N}$ . Für  $p = 2$ :  $W_2^\ell =: H^\ell$ .  $W_p^\ell(\Omega)$  sind Banachräume, für  $p = 2$  Hilberträume (Satz 4.4, PDE I).

**Beispiel** Sei  $\Omega \subseteq \mathbb{R}^n$  offen und beschränkt. Gegeben sei  $C^1(\bar{\Omega})$ . Mögliche Normen:

- (i).  $\|v\| := \sup_{\bar{\Omega}} |v| + \max_{\bar{\Omega}} |\nabla v|$ . Ist Banachraum, aber kein Hilbertraum.  
 (ii).  $\|v\|^2 := \int_\Omega v^2 + \int_\Omega (\nabla v)^2$ . Ist Prä-Hilbertraum, aber nicht vollständig. Vervollständigung ist  $H^1(\Omega)$ .

**Bemerkung** Alternativer Zugang zu Sobolev-Räumen (Neyers, Serrin, 1964): Für  $1 \leq p < \infty$  ist  $C^\infty(\Omega) \cap W_p^\ell(\Omega)$  dicht in  $W_p^\ell(\Omega)$ . Wichtig z.B. für Beweis von Ungleichungen in  $W_p^\ell$ , es genügt diese für Funktionen aus  $C^\infty(\Omega) \cap W_p^\ell(\Omega)$  zu zeigen.

**Beispiel** (Schwache Formulierung einer Randwertaufgabe (Neumann-Bedingung)) Sei  $c \in L^\infty(\Omega)$ ,  $f \in L^2(\Omega)$

$$-\Delta u + c \cdot u = f \qquad \left. \frac{\partial u}{\partial n} \right|_{\partial\Omega} = 0$$

wobei  $n$  Normaleneinheitsvektor. Multipliziere mit Testfunktion  $v$  und Integration gibt:

$$\begin{aligned} & - \int_\Omega \Delta u \cdot v + \int_\Omega c \cdot u \cdot v = \int_\Omega f \cdot v \\ \stackrel{\text{p.I.}}{\Rightarrow} & \underbrace{- \int_{\partial\Omega} \frac{\partial u}{\partial n} \cdot v}_0 + \int_\Omega \nabla u \cdot \nabla v + \int_\Omega c \cdot u \cdot v = \int_\Omega f \cdot v \end{aligned}$$

Gesucht ist also  $H^1(\Omega)$ , sodass

$$\int_\Omega \nabla u \cdot \nabla v + \int_\Omega c \cdot u \cdot v = \int_\Omega f \cdot v \quad (v \in H^1(\Omega))$$

Eigenschaften von  $H^1(\Omega)$ :

- $n = 1$ :  $\Omega = (0, 1)$ , dann gilt wegen  $(a \cdot 1 + b \cdot 1)^2 \leq 2(a^2 + b^2)$  und Jensen-Ungleichung:

$$\begin{aligned} v(x) &= v(y) + \int_y^x v'(t) dt \\ \Rightarrow v^2(x) &\leq 2 \cdot \left( v^2(y) + \int_0^1 v'^2(t) dt \right) \end{aligned}$$

Bestimme  $\int_0^1 dy$ :

$$v^2(x) \leq 2\|v\|_1^2 \tag{1.2}$$

Also ist jede  $H^1$ -Funktion beschränkt. Weiterhin:

$$|v(x) - v(y)| \leq \sqrt{x-y} \cdot \|v\|_1$$

Für  $v \in H^1$  approximiere durch  $C^\infty$ -Funktionen - aus obiger Gleichung folgt, dass diese gleichgradig stetig sind und somit folgt die Stetigkeit von  $v$  (Satz von Arzelà-Ascoli).

**Bemerkung** Man kann zeigen:  $v \in H^1 \Leftrightarrow v$  absolutstetig

- $n \geq 2$ :  $H^1$ -Funktionen sind nicht notwendigerweise beschränkt.

**Beispiel**  $n = 2$ ,  $\Omega := B(0, 1) \setminus \{0\}$  und  $u(x, y) := \ln \ln \frac{2}{\sqrt{x^2 + y^2}}$ . Dann ist  $u$  offenbar nicht beschränkt, aber es gilt  $u \in H^1(\Omega)$ .

$$\begin{aligned} \int_{\Omega} u_x^2 + u_y^2 &= \int_{\Omega} (u_r^2 + r^2 \cdot \underbrace{u_{\varphi}^2}_0) \cdot r \\ &= 2\pi \cdot \int_{r=0}^1 \left( \ln^2 \frac{2}{r} \cdot r \right)^{-1} dr \\ \xrightarrow{\ln \frac{2}{r} = t} &= -2\pi \cdot \int_{-\infty}^{\ln 2} \frac{1}{t^2} dt < \infty \end{aligned}$$

wegen

$$u_r = \frac{1}{\ln \frac{2}{r}} \cdot \left( -\frac{1}{r} \right)$$

(Die 2 ist als Faktor in  $u$  fundamental für die Existenz des Integrals!)

**Beispiel**  $n = 3$ ,  $\Omega := B(0, 1) \setminus \{0\}$ ,  $u(x, y, z) := r^{-\alpha}$  mit  $r := \sqrt{x^2 + y^2 + z^2}$ ,  $0 < \alpha < \frac{1}{2}$ . Dann

$$\int_{\Omega} (u_r)^2 = c \cdot \int_{\Omega} r^{-2\alpha-2} \cdot r^2 = c \cdot \int_{\Omega} r^{-2\alpha} < \infty$$

**Bemerkung** Für  $n \geq 2$ ,  $u \in H^1$ : Werte in Punkten existieren nicht notwendigerweise (Äquivalenzklassen!). Konsequenz: normale Interpolierende von  $H^1$ -Funktionen gibt es nicht.

### 1.1 Beispiel (Dirichletsche Randwertaufgabe)

$$-\Delta u = f \qquad u|_{\partial\Omega} = \varphi$$

Oft genügt es eine Funktion  $u_0$  zu finden mit  $u_0|_{\partial\Omega} = \varphi$ . Setze dann  $u := u_0 + v$  wobei  $v$  die Randwertaufgabe

$$-\Delta u = g \qquad u|_{\partial\Omega} = 0$$

(homogene Dirichlet-Bedingung) löst. Problem: Besitzt eine  $H^1$ -Funktion „Randwerte“? Ja. Damit dann auch schwache Formulierung des Problems möglich.

### Definition

$$W_{p,0}^{\ell}(\Omega) := \overline{C_c^{\infty}(\Omega)}^{W_p^{\ell}} \qquad H_0^p(\Omega) := W_{2,0}^p(\Omega)$$

**Beispiel** (Spuren von  $H^1$ -Funktionen, Spezialfall)  $\Omega$  Rechteck, hier speziell  $\Omega := [0, 1] \times [0, 1]$ ,  $v \in C^1(\Omega)$  mit

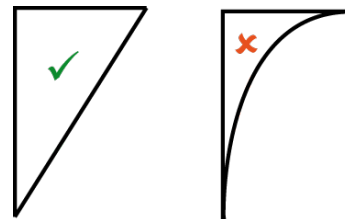
$$v(x, 0) = v(x, y) + \int_y^0 v_y(x, t) dt$$

Analog zu (1.2):

$$\begin{aligned} v^2(x, 0) &\leq c \cdot \left( v^2(x, y) + \int_0^1 v_y^2(x, t) dt \right) \\ \Rightarrow \int_0^1 v^2(x, 0) dx &\leq c \cdot \left( \|v\|_0^2 + \iint_{\Omega} v_y^2(x, y) d(x, y) \right) \\ &\leq c \cdot \|v\|_1 \end{aligned}$$

Allgemeineres Spurlemma gilt nur für Gebiete mit „gutartigen“ Rändern.

Wir setzen im folgenden generell voraus:  $\Omega$  beschränkt, offen, zusammenhängend. Weiterhin besitze  $\Omega$  einen Lipschitz-Rand, d.h.  $\forall x \in \partial\Omega \exists U \in \mathcal{U}_x : \partial\Omega \cap U$  ist Graph einer Lipschitz-stetigen Funktion.



**1.2 Lemma (Spurlemma)**

Sei  $\Omega \subseteq \mathbb{R}^n$  ein beschränktes Gebiet mit Lipschitz-Rand. Dann ist die Abbildung  $C^1(\bar{\Omega}) \cap W_p^1(\Omega) \rightarrow L^p(\partial\Omega), f \mapsto f|_{\partial\Omega}$  linear und beschränkt, d.h. es existiert  $c > 0$  mit

$$\|u\|_{L^p} \leq c \cdot \|u\|_{W_p^1(\Omega)} \quad (u \in C^1(\bar{\Omega}) \cap W_p^1(\Omega))$$

Beweis: Ohne Beweis. (Literatur: Braess: Finite Elemente) □

**1.3 Folgerung**

Es existiert eine Abbildung  $\gamma : W_p^1(\Omega) \rightarrow L^p(\partial\Omega)$ , sodass

$$\|\gamma(u)\|_{L^p(\partial\Omega)} \leq c \cdot \|u\|_{W_p^1} \quad (u \in W_p^1(\Omega))$$

mit  $\gamma(u) = u|_{\partial\Omega}$  für  $u \in C^1(\bar{\Omega}) \cap W_p^1(\Omega)$ .  $\gamma(u)$  heißt Spur ("Randwert") von  $u$ .

Beweis: Folgt aus Satz von Hahn-Banach. □

**Bemerkung** Ist jede  $L^2$ -Funktion auf dem Rand Spur einer  $H^1$ -Funktion? Nein.

**Definition** Die Menge aller  $L^2$ -Funktionen auf  $\partial\Omega$ , die Randwerte von einer  $H^1$ -Funktion sind, heißt  $H^{\frac{1}{2}}(\partial\Omega)$ .

**Beispiel** (Fortsetzung von Beispiel 1.1) (i). Homogen:

$$-\Delta u = f \quad u|_{\partial\Omega} = 0$$

Schwache Formulierung: Finde  $u \in H_0^1(\Omega)$  mit

$$\forall v \in H_0^1(\Omega) : (\nabla u, \nabla v) = (f, v)$$

(ii). Allgemein:

$$-\Delta u = f \quad u|_{\partial\Omega} = \varphi$$

mit  $\varphi \in H^{\frac{1}{2}}(\partial\Omega)$ .

**Bemerkung** Notation:  $H^{-1}(\Omega) := H_0^1(\Omega)'$ . Bekannt:  $H^{-1}(\Omega)$  ist Banachraum.

**Beispiel** (Fortsetzung) Sinnvoll ist z.B. die Aufgabe

$$(\nabla u, \nabla v) = g(v) \quad (v \in H_0^1(\Omega))$$

für  $g \in H^{-1}(\Omega)$ .

**1.4 Satz** (Friedrichs'sche Ungleichung, 1D)

Sei  $\Omega = (0, 1)$ . Dann gilt:

$$\exists c > 0 : \forall v \in H_0^1(\Omega) : \int_{\Omega} v^2 \leq c \cdot \int_{\Omega} v'^2$$



Beweis: Es gilt  $v(x) = \int_0^x v'(t) dt$  wegen  $v(0) = 0$ . Damit folgt

$$\begin{aligned} v^2(x) &\leq \int_0^1 v'(t)^2 dt \\ \Rightarrow \int_0^1 v^2(x) dx &\leq \int_0^1 v'(t)^2 dt \end{aligned} \quad \square$$

**1.5 Satz** (Friedrichs'sche Ungleichung)

$\Omega \subseteq \mathbb{R}^n$  sei enthalten in einem Würfel der Kantenlänge  $\ell$ . Dann gilt:

$$\forall v \in H_0^1(\Omega) : \|v\|_0 \leq \ell \cdot \left( \int_{\Omega} (\nabla v)^2 \right)^{\frac{1}{2}}$$

Beweis: Analog zum eindimensionalen Fall. □

**1.6 Folgerung**

Für  $v \in H^1(\Omega)$  sei

$$|v|_1 := \|\nabla v\|_0 = \left( \int_{\Omega} (\nabla v)^2 \right)^{\frac{1}{2}}$$

$(H_0^1(\Omega), |\cdot|_1)$  ist ein normierter Raum. Die Norm ist äquivalent zur  $H^1$ -Norm.

Beweis: Es gilt:

$$\|v\|_1^2 = \int_{\Omega} v^2 + \int_{\Omega} (\nabla v)^2 = \|v\|_0^2 + |v|_1^2$$

Zu zeigen: Es existieren  $c_1, c_2 > 0$  mit

$$c_1 \cdot \|v\|_1^2 \leq |v|_1^2 \leq c_2 \cdot \|v\|_1^2$$

Zweite Ungleichung ist klar. Erste Ungleichung:

$$\begin{aligned} c_1 \cdot \|v\|_1^2 \leq |v|_1^2 &\Leftrightarrow c_1 \cdot (\|v\|_0^2 + |v|_1^2) \leq |v|_1^2 \\ &\Leftrightarrow c_1 \cdot \|v\|_0^2 \leq (1 - c_1) \cdot |v|_1^2 \\ &\Leftrightarrow \|v\|_0^2 \leq \frac{1 - c_1}{c_1} \cdot |v|_1^2 \end{aligned}$$

Damit folgt die Behauptung aus Satz 1.5 □

**1.7 Satz** (Ungleichung von Poincaré-Friedrichs-Typ)

Sei  $\Omega \subseteq \mathbb{R}^n$  ein beschränktes Gebiet mit Lipschitz-Rand. Sei  $\Omega_1 \subseteq \Omega$  und  $\Omega_2 \subseteq \partial\Omega$  mit  $\lambda^n(\Omega_1) > 0$ ,  $\lambda^{n-1}(\Omega_2) > 0$ . Dann gilt für  $u \in H^1(\Omega)$ :

$$\begin{aligned} \|u\|_0^2 &\leq c \cdot \left( |u|_{1,\Omega}^2 + \left( \int_{\Omega_1} u \right)^2 \right) \\ \|u\|_0^2 &\leq c \cdot \left( |u|_{1,2}^2 + \left( \int_{\Omega_2} u \right)^2 \right) \end{aligned}$$

**Definition** Seien  $U, V$  normierte Räume.  $U$  heißt stetig eingebettet in  $V : \Leftrightarrow U \subseteq V, \|u\|_V \leq c \cdot \|u\|_U$  für  $u \in U$ . Notation:  $U \hookrightarrow V$ .

**Beispiel**  $H^1(\Omega)$  ist stetig eingebettet in  $L^2(\Omega)$ .

### 1.8 Satz (Einbettungssatz)

Sei  $\Omega \subseteq \mathbb{R}^n$  ein beschränktes Gebiet mit Lipschitz-Rand. Seien  $0 \leq j \leq k$  und  $1 \leq p < \infty$  und  $0 < \beta < 1$ . Dann gilt für  $k - j - \beta > \frac{n}{p}$ :  $W_p^k(\Omega) \hookrightarrow C^{j,\beta}(\bar{\Omega})$  wobei

$$C^{j,\beta}(\bar{\Omega}) := \{f \in C^j(\bar{\Omega}); \forall \alpha \in \mathbb{N}_0^n, |\alpha| = j : \partial^\alpha f \text{ } \beta\text{-Hölder-stetig}\}$$

**Bemerkung** (Spezialfälle vom Einbettungssatz) (i).  $p = 2, j = 0$ : Ist  $k - \beta > \frac{n}{2}$ , so gilt  $H^k \hookrightarrow C^{0,\beta}(\bar{\Omega})$ .

- (1)  $n = 1: k = 1, \beta < \frac{1}{2}$  ist zulässig.
- (2)  $n = 2: k = 2, \beta \in (0, 1)$  beliebig ist zulässig. Allerdings für  $k = 1$  keine Aussage (bereits gezeigt:  $H^1$ -Funktionen sind nicht notwendigerweise stetig).
- (3)  $n = 3: k = 2, \beta < \frac{1}{2}$  zulässig.

# 2

## Variationsgleichungen

Sei  $V$  ein Hilbertraum und  $V'$  der Dualraum zu  $V$ . Sei  $a : V \times V \rightarrow \mathbb{K}$  bilinear. Sei  $f$  ein lineares Funktional auf  $V$ . Variationsgleichung: Gesucht ist  $u \in V$  mit

$$\forall v \in V : a(u, v) = f(v) \quad (2.1)$$

Äquivalente Formulierung des Problems: Man kann den Operator

$$A : V \rightarrow V', v \mapsto \langle Au, v \rangle := a(u, v)$$

definieren. Dann ist mit  $f \in V'$  die Gleichung (2.1) äquivalent zu  $Au = f$  in  $V'$ . Bei uns:  $V$  ist stets ein Sobolevraum oder Teilraum eines Sobolevraums.

### Beispiel

$$-\Delta u + c \cdot u = f$$

dann

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} c \cdot u \cdot v$$

mit

$$v \in \begin{cases} H_0^1(\Omega) & \text{für } u|_{\partial\Omega} = 0 \\ H^1(\Omega) & \text{für } \frac{\partial u}{\partial n} \Big|_{\partial\Omega} = 0 \end{cases}$$

### Beispiel (Gemischte Randbedingungen)

$$-\Delta u + c \cdot u = f \quad u|_{\Gamma_1} = 0 \quad \frac{\partial u}{\partial n} + p \cdot u \Big|_{\Gamma_2} = q$$

mit  $\Gamma_1 \cup \Gamma_2 = \partial\Omega$ .

$$\begin{aligned} & - \int_{\Omega} \Delta u \cdot v + \int_{\Omega} c \cdot u \cdot v = \int_{\Omega} f \cdot v \\ \stackrel{\text{p.I.}}{\Rightarrow} & - \int_{\partial\Omega} \frac{\partial u}{\partial n} \cdot v + \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} c \cdot u \cdot v = \int_{\Omega} f \cdot v \end{aligned}$$

Weiterhin gilt

$$\int_{\partial\Omega} \frac{\partial u}{\partial n} \cdot v = \underbrace{\int_{\Gamma_1} \frac{\partial u}{\partial n} \cdot v}_{0 \text{ (} v \in V)} + \int_{\Gamma_2} \underbrace{\frac{\partial u}{\partial n}}_{(q-p \cdot u)} \cdot v$$

Definiere

$$\begin{aligned} f(v) &:= \int_{\Omega} f \cdot v + \int_{\Gamma_2} q \cdot v \\ V &:= \{v \in H^1; v|_{\Gamma_1} = 0\} \\ a(u, v) &:= \int_{\Gamma_2} p \cdot u \cdot v + \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} c \cdot u \cdot v \end{aligned} \quad (2.2)$$

**Bemerkung** (i). Unterscheidung:

- (1) wesentliche Randbedingungen: Gehen in Definition des Raumes  $V$  ein.
  - (2) natürliche Randbedingungen: Beeinflussen die Bilinearform.
- (ii). Gegebenes DGL-Problem und zugehörige Variationsgleichung sind nicht äquivalent. Ist aber die Lösung der Variationsgleichung ausreichend glatt, so ist eine „Rückrechnung“ möglich, siehe folgendes Beispiel.

**Beispiel (Gemischte Randbedingungen, Fortsetzung)** Sei  $u$  eine hinreichend glatte Lösung von

$$\forall v \in V : a(u, v) = f(v) \quad (2.3)$$

mit  $f, V, a$  wie in (2.2).

- (i). Die Variationsgleichung muss insbesondere für  $v \in H_c^1(\Omega)$  erfüllt sein. Für  $v \in H_0^1(\Omega)$  gilt

$$\begin{aligned} \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} c \cdot u \cdot v &= \int_{\Omega} f \cdot v \\ \stackrel{\text{p.I.}}{\Rightarrow} \int_{\Omega} (-\Delta u + c \cdot u - f) \cdot v &= 0 \\ \Rightarrow -\Delta u + c \cdot u - f &= 0 \end{aligned}$$

- (ii). Aus (i) folgt (mit dem üblichen Vorgehen), dass

$$-\underbrace{\int_{\partial\Omega} \frac{\partial u}{\partial n} \cdot v}_{\stackrel{v \in V}{=} \int_{\Gamma_2} \frac{\partial u}{\partial n} \cdot v} + \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} c \cdot u \cdot v = \int_{\Omega} f \cdot v$$

und somit folgt aus (2.3)

$$\begin{aligned} \int_{\Gamma_2} \left( \frac{\partial u}{\partial n} - (q - p \cdot u) \right) \cdot v &= 0 \\ \Rightarrow \frac{\partial u}{\partial n} + p \cdot u \Big|_{\Gamma_2} &= q \end{aligned}$$

**Beispiel (Stokes-Problem, 2D)**

$$\begin{aligned} -\Delta u_1 + \frac{\partial p}{\partial x_1} &= f_1 \quad \text{in } \Omega \\ -\Delta u_2 + \frac{\partial p}{\partial x_2} &= f_2 \quad \text{in } \Omega \\ \operatorname{div} u := \frac{\partial u_1}{\partial x_1} + \frac{\partial u_2}{\partial x_2} &= 0 \quad \text{in } \Omega \\ u_2|_{\partial\Omega} = u_1|_{\partial\Omega} &= 0 \end{aligned}$$

wobei  $u = (u_1, u_2)$  Geschwindigkeitsvektor und  $p$  Druck.

- (i). Eine Variante:  $u \in V := \{v \in H_c^1(\Omega) \times H_c^1(\Omega); \operatorname{div} v = 0\}$  Erste Gleichung  $\cdot v_1$ , integriere, partielle Integration. Analog zweite Gleichung. Summation gibt (beachte  $\operatorname{div} v = 0$ )

$$\int_{\Omega} \nabla u_1 \cdot \nabla v_1 + \int_{\Omega} \nabla u_2 \cdot \nabla v_2 + 0 = \int_{\Omega} f_1 \cdot v_1 + \int_{\Omega} f_2 \cdot v_2$$

d.h.  $p$  wurde eliminiert.

- (ii). Andere Varianten: Gleichzeitige Bestimmung von  $p$  und  $v$ .

**Bemerkung** Nächstes Ziel: Existenz und Eindeutigkeit von Lösungen von Variationsgleichungen ( $\rightarrow$  Lax-Milgram-Lemma).

**Definition** Sei  $V$  ein Hilbertraum,  $a : V \times V \rightarrow \mathbb{K}$  bilinear.

- (i).  $a$  heißt beschränkt  $:\Leftrightarrow \exists C > 0 \forall v, w \in V : |a(v, w)| \leq C \cdot \|v\| \cdot \|w\|$ .
- (ii).  $a$  heißt  $V$ -elliptisch (koerzitiv)  $:\Leftrightarrow \exists \alpha > 0 \forall v \in V : a(v, v) \geq \alpha \cdot \|v\|^2$ .
- (iii).  $a$  heißt positiv  $:\Leftrightarrow \forall v \in V \setminus \{0\} : a(v, v) > 0$ .
- (iv).  $a$  heißt symmetrisch  $:\Leftrightarrow \forall v, w \in V : a(v, w) = a(w, v)$ .

### 2.1 Satz (Eindeutigkeit der Lösung)

Sei  $V$  ein Hilbertraum. Seien  $u_1, u_2$  Lösungen von

$$\forall v \in V : a(u_i, v) = f(v)$$

mit  $a : V \times V \rightarrow \mathbb{K}$  bilinear und positiv. Dann gilt  $u_1 = u_2$ .

Beweis: Offenbar gilt dann

$$\forall v \in V : a(u_1 - u_2, v) = 0$$

Wähle  $v := u_1 - u_2$ , dann folgt  $a(u_1 - u_2, u_1 - u_2) = 0$ . Wegen  $a$  positiv folgt die Eindeutigkeit der Lösung.  $\square$

**Bemerkung** Untersuche zunächst den symmetrischen Fall. Dazu wichtig: Riesz'scher Darstellungssatz.

### 2.2 Satz (Existenz von Lösungen (symmetrischer Fall))

Sei  $V$  ein Hilbertraum und  $a$  eine symmetrische  $V$ -elliptische Bilinearform auf  $V$ . Dann besitzt die Variationsgleichung

$$\forall v \in V : a(u, v) = f(v)$$

eine Lösung.

Beweis:  $p(v) := \sqrt{a(v, v)}$  ist eine Norm auf  $V$  („energetische Norm“) und  $a(\cdot, \cdot)$  ein Skalarprodukt. Aus dem Darstellungssatz von Riesz folgt: Es existiert  $u \in V$  mit

$$\forall v \in V : a(u, v) = f(v)$$

$\square$

### 2.3 Lemma (Zusammenhang zu Variationsproblem)

Sei  $V$  ein normierter Raum,  $u \in V$  und  $a$  eine symmetrische positive Bilinearform auf  $V$ . Dann äquivalent:

- (i).  $u$  minimiert das Funktional  $J(v) := \frac{1}{2}a(v, v) - f(v)$  (Variationsproblem)
- (ii).  $\forall v \in V : a(u, v) = f(v)$

Beweis:  $\bullet$  (ii)  $\Rightarrow$  (i):

Es gilt

$$\begin{aligned}
 J(w) &= \underbrace{\frac{1}{2}a(u, u) - f(u)}_{J(u)} + \underbrace{a(u, w - u) - f(w - u)}_{\stackrel{(ii)}{=} 0} + \underbrace{\frac{1}{2}a(w - u, w - u)}_{> 0} \\
 &> J(u)
 \end{aligned}$$

für alle  $w \neq u$ .

- (i)  $\Rightarrow$  (ii):

Nutze Standardtechnik in der Variationsrechnung. Man betrachte  $\mathbb{R} \ni t \mapsto J(u + t \cdot v)$  mit  $v \in V$ . Notwendige Optimalitätsbedingung:

$$\left. \frac{dJ(u + t \cdot v)}{dt} \right|_{t=0} = 0$$

und wegen

$$\begin{aligned} \frac{dJ(u + t \cdot v)}{dt} &= \frac{d}{dt} \left( \frac{1}{2} t^2 \cdot a(v, v) - t \cdot f(v) + \frac{1}{2} a(u, u) + a(u, v) \cdot t - f(u) \right) \\ &= a(u, v) + t \cdot a(v, v) - f(v) \end{aligned}$$

folgt somit die Behauptung. □

**Bemerkung** Dirichlet studierte das Variationsproblem

$$J(v) := \frac{1}{2} \int_{\Omega} (\nabla v)^2 - \int_{\Omega} f \cdot v$$

für  $v \in C^1(\Omega) \cap C(\bar{\Omega})$ ,  $v|_{\partial\Omega} = 0$ . Er bewies, dass eine Lösung existiert. Weierstraß zeigte später (zu Recht), dass der Beweis falsch ist. Studiert man die Gleichung auf  $V := H_c^1(\Omega)$  (Hilbertraum!), so ist die Existenz der Lösung mit Lemma 2.3 und Satz 2.2 klar.

#### 2.4 Lemma (Lax-Milgram)

Sei  $V$  ein Hilbertraum und  $a : V \times V \rightarrow \mathbb{K}$  eine beschränkte  $V$ -elliptische Bilinearform. Dann existiert eine eindeutige Lösung  $u \in V$  der Variationsgleichung

$$\forall v \in V : a(u, v) = f(v)$$

Beweis: Sei  $z \in V$  definiert durch

$$(z, v) = (y, v) - r \cdot (a(y, v) - f(v))$$

für  $r \in \mathbb{R}$ ,  $v \in V$  (Existenz und Eindeutigkeit: Darstellungssatz von Riesz). Betrachte die Abbildung  $y \mapsto T_r(y) := z$ .

Idee: Nachweis der Existenz eines Fixpunktes der Abbildung  $T_r$  (der Fixpunkt ist dann offenbar Lösung der Variationsgleichung). Dazu Anwendung des Banach'schen Fixpunktsatzes, d.h. es ist zu zeigen, dass  $T_r$  eine Kontraktion ist. Betrachte Hilfsproblem  $V \ni p \mapsto S_p \in V$  mit

$$\forall v \in V : (S_p, v) = a(p, v)$$

Dann gilt  $(T_r y, v) = (y - r S y + r \cdot g, v)$  wobei  $f(v) = (g, v)$  (Darstellungssatz von Riesz). Damit:

$$\begin{aligned} \|T_r y - T_r w\|^2 &= (T_r y - T_r w, T_r y - T_r w) \\ &= (y - w - r \cdot S(y - w), y - w - r \cdot S(y - w)) \\ &= \|y - w\|^2 - 2(y - w, S(y - w)) + r^2 \cdot \|S(y - w)\|^2 \end{aligned}$$

Eigenschaften von  $S$ :

- (i). Setze  $v = S_p$  in Definition von  $S$ , dann  $\|S p\|^2 = a(p, S p) \leq C \cdot \|p\| \cdot \|S p\|$ , also  $\|S p\| \leq C \cdot \|p\|$ .
- (ii). Setze  $v = p$  in Definition von  $S$ , dann  $(S p, p) = a(p, p) \geq \alpha \cdot \|p\|^2$  mit  $\alpha > 0$ , da  $V$ -elliptisch.

Unter Nutzung der Abschätzungen folgt

$$\|T_r y - T_r w\|^2 \leq \|y - w\|^2 \cdot \underbrace{(1 - 2r \cdot \alpha + r^2 \cdot C^2)}_{=: \beta}$$

Für  $r > 0$  hinreichend klein folgt  $\beta < 1$  und daher die Behauptung (genauer:  $r \in (0, \frac{2\alpha}{C^2})$ ). □

## 2.5 Folgerung (Stabilität)

Sei  $V$  ein Hilbertraum und  $a : V \times V \rightarrow \mathbb{K}$  eine beschränkte  $V$ -elliptische Bilinearform. Dann gilt für die Lösung  $u$  der Variationsgleichung

$$\forall v \in V : a(u, v) = f(v),$$

dass

$$\|u\| \leq \frac{\|f\|}{\alpha}$$

Beweis: Setze  $v := u$ , dann folgt

$$\alpha \cdot \|u\|^2 \leq a(u, u) = f(u) \leq \|f\| \cdot \|u\|$$

und somit folgt die Behauptung.  $\square$

**Bemerkung** Oft ist es relativ einfach im konkreten Fall die Beschränktheit der Bilinearform nachzuweisen. Nachweis der  $V$ -Elliptizität meist problematischer.

### Beispiel

$$-\Delta u|_{\Omega} = f \qquad u|_{\partial\Omega} = 0$$

Bilinearform ist gegeben durch

$$a(u, v) := (\nabla u, \nabla v)$$

für  $V = H_0^1(\Omega)$ . Aus Folgerung 1.6 bekannt:  $(H_0^1(\Omega), |\cdot|_1)$  ist normierter Raum, also

$$\begin{aligned} |(\nabla u, \nabla v)| &\stackrel{\text{CSU}}{\leq} |u|_1 \cdot |v|_1 \\ (\nabla u, \nabla u) &= |u|_1^2 \end{aligned}$$

### Beispiel (Konvektions-Diffusions-Reaktions-Gleichung)

$$-\Delta u + b \cdot \nabla u + c \cdot u = f \qquad u|_{\partial\Omega} = 0$$

mit  $c - \frac{1}{2} \operatorname{div} b \geq 0$ . Warum diese Voraussetzung? Die zugehörige Bilinearform ist gegeben durch

$$a(u, v) := (\nabla u, \nabla v) + (b \cdot \nabla u + c \cdot u, v)$$

und  $V := (H_0^1(\Omega), \|\cdot\|_1)$ . Nachweis der  $V$ -Elliptizität:

$$a(v, v) = (\nabla v, \nabla v) + (b \cdot \nabla v, v) + (c \cdot v, v)$$

Partielle Integration (siehe (1.1)):

$$\begin{aligned} \int_{\Omega} b_1 \cdot \partial_x v \cdot v &= - \int_{\Omega} v \cdot \partial_x (b_1 \cdot v) = - \int_{\Omega} \partial_x b_1 \cdot v^2 - \int_{\Omega} v \cdot b_1 \cdot \partial_x v \\ \Rightarrow \int_{\Omega} b_1 \cdot \partial_x v \cdot v &= - \frac{1}{2} \int_{\Omega} \partial_x b_1 \cdot v^2 \end{aligned}$$

Somit folgt  $(b \cdot \nabla v, v) = -\frac{1}{2} (\operatorname{div} b \cdot v, v)$ . Damit

$$a(v, v) = (\nabla v, \nabla v) + \underbrace{\left( c - \frac{1}{2} \operatorname{div} b \right) \cdot v, v}_{\geq 0}$$

und somit die  $V$ -Elliptizität.

**Beispiel**

$$-\Delta u + c \cdot u = f \qquad \frac{\partial u}{\partial n} \Big|_{\partial\Omega} = 0$$

Dann

$$a(u, v) := (\nabla u, \nabla v) + (c \cdot u, v)$$

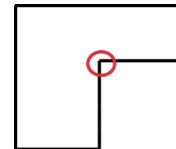
und  $V = (H^1(\Omega), \|\cdot\|_1)$ . Untersuchung der  $V$ -Elliptizität:

- (i). Gilt  $\|c\|_{\infty, \Omega} > 0$ , dann ist  $a$   $V$ -elliptisch.
- (ii).  $c = 0$ : Betrachte den Raum  $W := \{v \in H^1; \int_{\partial\Omega} v = 0\}$ . Dann kann  $W$ -Elliptizität nachgewiesen werden. (Beachte, dass im klassischen Sinne die Lösung nicht eindeutig ist!)

**Bemerkung** Später bei Finite-Elemente-Fehler-Analyse reicht  $u \in H^1(\Omega)$  nicht aus, um Fehlerabschätzungen zu beweisen (für  $u \in H^1$  nur Konvergenz beweisbar). Wünschenswert daher:  $u \in H^k(\Omega)$  mit  $k \geq 2$ . Dazu folgende Aussagen:

- (i). Sind die Daten einer elliptischen Randwertaufgabe 2. Ordnung hinreichend glatt und der Rand von  $\Omega$  hinreichend glatt, so gilt:  $f \in H^\ell \Rightarrow u \in H^{\ell+2}$  (Shift-Theorem).
- (ii). Betrachtet wird ein Dirichlet-Problem für eine elliptische Randwertaufgabe 2. Ordnung mit glatten Koeffizienten und  $\Omega$  konvex. Dann:  $f \in L^2(\Omega) \Rightarrow u \in H^2(\Omega)$ ,  $\|u\| \leq C \cdot \|f\|_{L^2}$ .

(Literatur: Grisvard: Elliptic problems in nonsmooth domains) Unangenehm sind also nichtkonvexe Gebiete, z.B. L-shaped.



**Beispiel** (i). Die Gleichung  $\Delta u = 0$  hat die exakte Lösung

$$u(r, \varphi) = r^{\frac{\pi}{\omega}} \cdot \sin\left(\frac{\pi}{\omega} \cdot \varphi\right)$$

Offenbar gilt dann

$$u|_{\varphi=0} = 0 \qquad u|_{\varphi=\omega} = 0 \qquad u|_{r=1} = \sin\left(\frac{\pi}{\omega} \cdot \varphi\right)$$

Falls  $\omega > \pi$  gilt  $u \notin H^2(\Omega)$ . Vgl. Übung 1, Aufgabe 6

- (ii). Bei gemischten Randwertproblem hat man noch mehr Probleme mit der Regularität. Die Funktion

$$u(r, \varphi) := r^{\frac{\pi}{2\omega}} \cdot \sin\left(\frac{\pi}{2\omega} \cdot \varphi\right)$$

löst  $\Delta u = 0$ . Dann Randwerte

$$u|_{\varphi=0} = 0 \qquad \frac{\partial u}{\partial n} \Big|_{\varphi=\omega} = 0 \qquad u|_{r=1} = \sin\left(\frac{\pi}{2\omega} \cdot \varphi\right)$$

Für  $\omega > \frac{\pi}{2}$  gilt  $u \notin H^2$ .



# 3

## Ritz-Galerkin-Verfahren

### Verfahren (Ritz-Verfahren)

Ziel: Finde  $v_0 \in V$  mit

$$J(v_0) = \min_{v \in V} J(v)$$

Idee: Betrachte Unterraum  $V_n \subseteq V$  mit  $\dim V_n < \infty$ . Finde  $v_n \in V_n$  mit

$$J(v_n) = \min_{v \in V_n} J(v)$$

### Verfahren (Galerkin-Verfahren)

Ziel: Finde  $u \in V$  mit

$$\forall v \in V : a(u, v) = f(v) \quad (3.1)$$

(stetiges Problem). Idee: Betrachte Unterraum  $V_n \subseteq V$  mit  $\dim V_n < \infty$ . Finde  $u_n \in V_n$  mit

$$\forall v \in V_n : a(u_n, v) = f(v) \quad (3.2)$$

(diskretes Problem).

**Bemerkung** (i). Sei  $a : V \times V \rightarrow \mathbb{K}$  eine symmetrische Bilinearform und  $J(v) := \frac{1}{2}a(v, v) - f(v)$ . Dann folgt aus Satz 2.3, dass das Ritz-Verfahren gleich dem Galerkin-Verfahren ist.

(ii). Beide Verfahren gehören zur Klasse der Projektionsverfahren. Weitere Projektionsverfahren sind beispielsweise die „Methode des gewichteten Residuums“ und Kollokation. Keine Projektionsverfahren sind zum Beispiel Differenzenverfahren und Varianten der FEM-Methode.

### 3.1 Lemma (Galerkin-Orthogonalität)

Es gilt

$$\forall v_n \in V_n : a(u - u_n, v_n) = 0$$

wobei  $u$  Lösung des stetigen und  $u_n$  Lösung des diskreten Problems.

Beweis: Setze  $v = v_n$  und bestimme (3.1)-(3.2). □

### 3.2 Satz (Cea-Lemma)

Sei  $V$  ein Hilbertraum und  $a : V \times V \rightarrow \mathbb{K}$  eine beschränkte  $V$ -elliptische Bilinearform. Dann gilt

$$\|u - u_n\| \leq \frac{M}{\alpha} \cdot \inf_{v_n \in V_n} \|u - v_n\|$$

Beweis: Es gilt für beliebige  $v_n \in V_n$ :

$$\begin{aligned} \alpha \cdot \|u - u_n\|^2 &\leq a(u - u_n, u - u_n) = a(u - u_n, u - v_n) + \underbrace{a(u - u_n, v_n - u_n)}_{\stackrel{3.1}{=} 0} \\ &\leq C \cdot \|u - u_n\| \cdot \|u - v_n\| \\ \Rightarrow \|u - u_n\| &\leq \frac{M}{\alpha} \cdot \inf_{v_n \in V_n} \|u - v_n\| \end{aligned}$$

□

**Bemerkung** Der Fehler des Galerkin-Verfahrens ist also bis auf eine multiplikative Konstante gleich dem Bestapproximationsfehler in  $V_n$ . Deshalb nennt man das Galerkin-Verfahren quasioptimal.

**Bemerkung** Praktisch benutzt man eine Basis des endlich-dimensionalen Raumes  $V_n$ , um (3.2) zu lösen. Sei  $\{\varphi_1, \dots, \varphi_N\}$  eine Basis von  $V_n$ . Dann existieren  $u_1, \dots, u_N$  mit

$$u_n = \sum_{i=1}^N u_i \cdot \varphi_i$$

Somit in (3.2) für  $j = 1, \dots, N$ :

$$\begin{aligned} a\left(\sum_{i=1}^N u_i \cdot \varphi_i, \varphi_j\right) &= f(\varphi_j) \\ \Leftrightarrow \sum_{i=1}^N a(\varphi_i, \varphi_j) \cdot u_i &= f(\varphi_j) \end{aligned}$$

(Auf Grund der Linearität ist dieses Gleichungssystem sogar äquivalent zu (3.2).) Damit  $N$  Gleichungen für  $N$  Unbekannte. Die Koeffizientenmatrix heißt oft Steifigkeitsmatrix.

**Beispiel**

$$-u'' = f \qquad u(0) = u(1) = 0$$

$V_n$  werde aufgespannt durch  $\varphi_1(x) := x \cdot (1 - x)$ .

(i). Ritz-Verfahren:

$$J(v) := \frac{1}{2} \int_0^1 v'^2 - \int_0^1 f \cdot v$$

Für  $u_h = u_1 \cdot \varphi_1$  gilt

$$J(u_h) = \frac{1}{2} \int_0^1 (u_1 \cdot \varphi_1')^2 - \int_0^1 f \cdot u_1 \cdot \varphi_1$$

Notwendige Bedingung:  $\frac{dJ}{du_1} = 0$ .

(ii). Galerkin-Verfahren:

$$a(u, v) := (u', v') \qquad V = H_0^1(\Omega)$$

Für  $u_h = u_1 \cdot \varphi = 1$  gilt

$$a(u_n, \varphi_1) = u_1 \cdot \int_0^1 \varphi_1' \cdot \varphi_1' = \int_0^1 f \cdot \varphi_1$$

(iii). Methode des gewichteten Residuums: Betrachte die Gleichung  $Lu = f$  (hier:  $Lu := -u''$ ). Dann löse

$$\int_0^1 (Lu_n - f) \cdot \varphi_j = 0 \quad (j = 1, \dots, N)$$

wobei  $\{\varphi_1, \dots, \varphi_N\}$  Basis von  $V_n$ . Hier im Beispiel erhält man

$$\int_0^1 (-u_1 \cdot \varphi_1'' - f) \cdot \varphi_1 = 0$$

Ist in dem Fall identisch mit Galerkin-Verfahren (folgt aus partieller Integration).

(iv). Kollokation: Löse

$$\forall j = 1, \dots, N : (Lu_n - f)(\xi_j) = 0$$

$\xi_j$  heißen Kollakationsstellen (gegeben). In obigem Beispiel: Wähle  $\xi_1 = \frac{1}{2}$ . Forderung:

$$(-u_1 \cdot \varphi_1'' - f) \left( \frac{1}{2} \right) = 0$$

#### Beispiel

$$-u'' = f \qquad u(0) = u'(1) = 0$$

Schwache Formulierung des Problems ist gegeben durch

$$a(u, v) := (u', v') \qquad V := \{v \in H^1; v(0) = 0\}$$

Wähle  $V_n := \text{lin} \left\{ \frac{x^k}{k}; k = 1, \dots, N \right\}$ . Dann gilt  $V_n \subseteq V$ . Elemente der Steifigkeitsmatrix:

$$a_{ij} = (\varphi_i', \varphi_j') = \int_0^1 x^{i+j-2} = \frac{1}{i+j-1}$$

(Hilbert-Matrix, hat schlechte Kondition!). Also: Wahl von  $V_n$  wichtig.

#### Beispiel

$$-u'' = f \qquad u(0) = u(1) = 0$$

Sei  $V_n := \text{lin}\{\sin(j \cdot \pi \cdot x); j = 1, \dots, N\}$ , dann gilt  $V_n \subseteq V$ . Elemente der Steifigkeitsmatrix:

$$a_{ij} = \pi^2 \cdot \int_0^1 \cos(i \cdot \pi \cdot x) \cdot \cos(j \cdot \pi \cdot x) = 0$$

für  $i \neq j$ , d.h. Steifigkeitsmatrix ist Diagonalmatrix.

Ist im Allgemeinen unrealistisch. Idee: Wähle Räume  $V_n$  derart, dass  $(\varphi_i, \varphi_j) = 0$  für möglichst viele  $i, j$ . Also: Basisfunktionen mit lokalem Träger ( $\rightarrow$  Splines  $\rightarrow$  FEM). Alternative Idee: orthogonale Polynome benutzen ( $\rightarrow$  spektrale Methoden) .

# 4

## Finite-Elemente-Methode

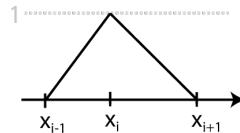
- FEM = Galerkin + Splines.
- Elliptische Randwertaufgaben 2. Ordnung:  $V = H^1(\Omega)$  oder Teilraum von  $H^1(\Omega)$ . Deshalb: Zur FEM benötigt man dann nur stetige Splines.
- Elliptische Randwertaufgaben 4. Ordnung:  $V = H^2(\Omega)$  oder Teilraum von  $H^2(\Omega)$ . Man benötigt dann  $C^1$ -Splines.

### 4.1 1D

Randwertaufgaben für Differentialgleichungen 2. Ordnung:  $C^0$ -Splines.

#### 4.1.1 Stückweise lineare Splines

Sei  $\Pi = \{0 = x_0 < \dots < x_N = 1\}$  eine Partition von  $[0, 1]$  und  $h_i := x_i - x_{i-1}$ . Definiere  $V_h$  als den Raum der stückweise linearen stetigen Funktionen auf dem Gitter (evtl. noch homogene Dirichlet-Bedingungen berücksichtigen). Seien  $\{\varphi_1, \dots, \varphi_n\}$  die Basisfunktionen von  $V_h$  mit  $\varphi_i(x_j) := \delta_{ij}$  („Hutfunktionen“).



#### Beispiel

$$-u'' = f \quad u(0) = u(1) = 0 \quad (4.1)$$

Schwache Formulierung:

$$\forall v_h \in V_h : (u'_h, v'_h) = (f, v_h)$$

wobei  $V_h := \text{lin}\{\varphi_1, \dots, \varphi_{N-1}\}$ . Es gilt  $u_h = \sum_{i=1}^{N-1} u_i \cdot \varphi_i$  mit  $u_j = u_h(x_j)$ . (Näherungswerte für  $u$  in Gitterpunkten) Koeffizientenmatrix:

$$(\varphi'_i, \varphi'_j) = \begin{cases} 0 & |i-j| \geq 2 \\ \frac{1}{h_i^2} \cdot h_i + \frac{1}{h_{i+1}^2} \cdot h_{i+1} = \frac{1}{h_i} + \frac{1}{h_{i+1}} & i = j \\ -\frac{1}{h_{i+1}^2} \cdot h_{i+1} & i = j - 1 \\ -\frac{1}{h_i^2} \cdot h_i & i = j + 1 \end{cases}$$

Damit tridiagonales Gleichungssystem:

$$-\frac{1}{h_i} u_{i-1} + \left( \frac{1}{h_i} + \frac{1}{h_{i+1}} \right) \cdot u_i - \frac{1}{h_{i+1}} \cdot u_{i+1} = \int_{x_{i-1}}^{x_{i+1}} f \cdot \varphi_i \quad (4.2)$$

### Verfahren (Differenzenverfahren)

Betrachte die gegebene Randwertaufgabe  $Lu = f$  mit gewissen Randbedingungen im Gebiet  $\Omega$ . Man nehme ein Gitter mit (endlich vielen) Gitterpunkten in  $\Omega$ . Nun wird die Differentialgleichung in den Gitterpunkten (in  $\Omega$ ) betrachtet und die Ableitungen werden durch Differenzenquotienten (gebildet mit Hilfe von Werten in den Gitterpunkten) approximiert. Auf diese Weise erhält man ein Gleichungssystem für die Näherungswerte von  $u$  in den Gitterpunkten.

**Beispiel** Wende Differenzenverfahren auf (4.1) an. Es gilt

$$-u''(x_i) \approx -2 \frac{1}{h_i + h_{i+1}} \cdot \left( \frac{u_{i+1} - u_i}{h_{i+1}} - \frac{u_i - u_{i-1}}{h_i} \right)$$

Benutze dazu Taylor-Entwicklung:

$$\begin{aligned} \frac{u(x_i + h_{i+1}) - u(x_i)}{h_{i+1}} &= u'(x_i) + \frac{u''(x_i)}{2} \cdot h_{i+1} + \frac{u^{(3)}(x_i)}{6} \cdot h_{i+1}^2 + o(h_{i+1}^3) \\ \frac{u(x_i) - u(x_i - h)}{h_i} &= u'(x_i) - \frac{u''(x_i)}{2} \cdot h_i + \frac{u^{(3)}(x_i)}{6} \cdot h_i^2 + o(h_i^3) \end{aligned}$$

(Konsistenzanalyse). Damit Differenzenverfahren gegeben durch

$$-\frac{2}{h_i + h_{i+1}} \cdot \left( \frac{u_{i+1} - u_i}{h_{i+1}} - \frac{u_i - u_{i-1}}{h_i} \right) = f(x_i) \quad (4.3)$$

für  $i = 1, \dots, N - 1$  mit  $u_0 = u_N := 0$ . Offenbar ist dieses Gleichungssystem gleich (4.2) bis auf einen Skalierungsfaktor, d.h. lineare finite Elemente und Standard-Differenzenverfahren (1D) stimmen in diesem Beispiel fast überein.

**Definition** (Konsistenz von Differenzenverfahren) Sei  $L_h u_h = f_h$  das Gleichungssystem für den Vektor der Näherungswerte in den Gitterpunkten gemäß dem Differenzenverfahren angewendet auf  $Lu = f$ . Sei  $r_h u$  der Vektor der Werte von  $u$  in den Gitterpunkten.

- (i). Das Verfahren heißt konsistent  $:\Leftrightarrow \|L_h r_h u - f_h\| \rightarrow 0$  für  $h \rightarrow 0$  wobei  $h$  maximaler Abstand der Gitterpunkten.
- (ii). Das Verfahren heißt konsistent von der Ordnung  $p$   $:\Leftrightarrow \|L_h r_h u - f_h\| = \mathcal{O}(h^p)$ .

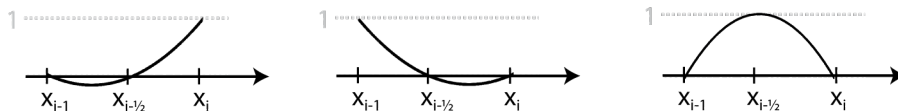
**Beispiel** Untersuche die Konsistenz (bzgl. der Maximumnorm) des obigen Beispiels. Aus Taylorentwicklung folgt

$$\begin{aligned} [L_h r_h u - f_h]_i &= -\frac{2}{h_i + h_{i+1}} \cdot \left( \frac{u_{i+1} - u_i}{h_{i+1}} - \frac{u_i - u_{i-1}}{h_i} \right) - f(x_i) \\ &= \underbrace{u''(x_i) - f(x_i)}_0 + \frac{u^{(3)}(x_i)}{3} \cdot (h_{i+1} - h_i) + \mathcal{O}(h^2) \end{aligned}$$

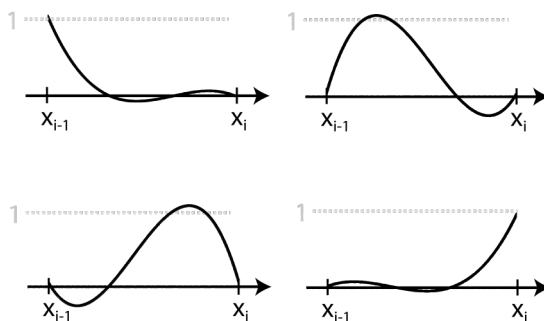
falls  $u^{(4)}$  beschränkt ist. Damit: Die Konsistenzordnung, gemessen in der Maximumnorm, ist 2 auf einem äquidistanten Gitter, anderenfalls 1.

#### 4.1.2 Finite Elemente höherer Ordnung

- quadratische Finite Elemente: Sei  $\Pi = \{x_0 < \dots < x_N\}$  eine Partition und  $x_{i-\frac{1}{2}} := \frac{x_i + x_{i-1}}{2}$ . Werden in der Praxis am häufigsten eingesetzt.

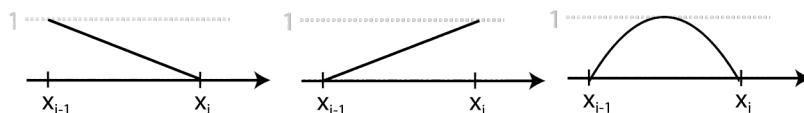


- Kubische Finite Elemente:



d.h. Definition der Basisfunktionen nach Lagrange-Prinzip (nodale Basisfunktionen). Die Gitterpunkten heißen auch Knoten.

- Alternative: hierarchische Basen. Zum Beispiel für quadratische Finite Elemente: Nutze weiterhin affin-lineare Funktionen, aber ergänze die Basis um eine quadratische Funktion.

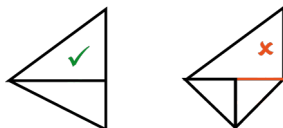


## 4.2 2D

In diesem Abschnitt sei  $\Omega \subseteq \mathbb{R}^2$  und polygonal.

**Definition** Eine Zerlegung  $\mathcal{T}_h$  von  $\Omega$  in polygonale Teilgebiete  $(K_i)_{i \in I}$  heißt zulässig, wenn

- $\bigcup_{i \in I} K_i = \bar{\Omega}$
- Besteht  $K_i \cap K_j$  aus einem Punkt, so ist dieser Eckpunkt von  $K_i$  und  $K_j$ .
- Besteht  $K_i \cap K_j$  aus mehr als einem Punkt, so ist  $K_i \cap K_j$  eine Kante von  $K_i$  und  $K_j$ .



**Definition** Gegeben sei eine zulässige Zerlegung  $\mathcal{T}_h$  in  $\Omega$  und  $K$  ein Element von  $\mathcal{T}_h$ . Ein finites Element ist ein Tripel  $(K, P_K, \Sigma_K)$  wobei

- $P_K$  ein endlich-dimensionaler Teilraum der stetigen Funktionen (oft: Polygone).
- $\Sigma_K$  ist eine Menge von linear unabhängigen Funktionalen auf  $P_K$ . Jedes Element  $p \in P_K$  wird durch die Vorgabe der Funktionswerte dieser linearen Funktionalen eindeutig bestimmt („Unisolvenz“).

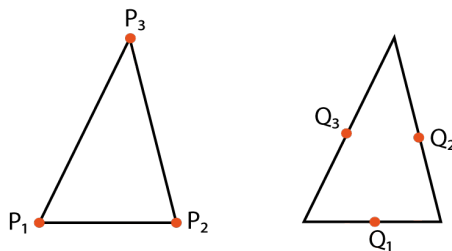
**Beispiel (Unisolvenz)** (i). Sei  $K \subseteq \mathbb{R}^2$  ein Dreieck und  $P_K := \{(x, y) \mapsto a + b \cdot x + c \cdot y; a, b, c \in \mathbb{R}\}$ . Dann offenbar  $\dim P_K = 3$  und

$$\Sigma_K^1 := \{f \mapsto f(P_i) = \delta_{P_i}(f); i = 1, 2, 3\}$$

ist unisolvent wobei  $P_1, P_2, P_3$  die Ecken des Dreiecks bezeichnen. Ebenfalls unisolvent ist

$$\Sigma_K^2 := \{f \mapsto f(Q_i) = \delta_{Q_i}(f); i = 1, 2, 3\}$$

wobei  $Q_1, Q_2, Q_3$  die Seitenhalbierenden bezeichnen.



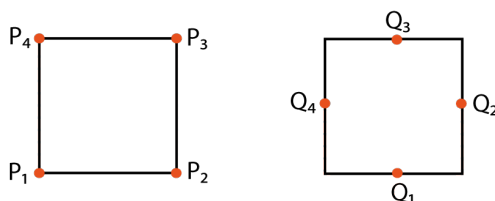
- (ii). Sei  $K \subseteq \mathbb{R}^2$  ein Rechteck und  $P_K := \{(x, y) \mapsto a + b \cdot x + c \cdot y + d \cdot x \cdot y; a, b, c, d \in \mathbb{R}\}$ . Es gilt  $\dim P_K = 4$ . Sei

$$\Sigma_K^1 := \{\delta_{Q_i}; i = 1, \dots, 4\}$$

wobei  $Q_i$  die Seitenmittelpunkte bezeichnen. Ist nicht unisolvent, denn  $x \cdot y|_{Q_i} = 0$ . Unisolvent ist dagegen

$$\Sigma_K^2 := \{\delta_{P_i}; i = 1, \dots, 4\}$$

wobei  $P_i$  die Eckpunkte bezeichnen.



Möchte man die Freiheitsgrade in den Punkten  $Q_i$  erhalten, kann man stattdessen  $P_K := \{(x, y) \mapsto a + b \cdot x + c \cdot y + d \cdot (x^2 - y^2); a, b, c, d \in \mathbb{R}\}$ . Damit ist  $(K, P_K, \Sigma_K^1)$  unisolvent: Sei  $K = [-1, 1]^2$ . Für  $x = 0, y = \pm 1$  bzw.  $y = 0, x = \pm 1$  erhält man das Gleichungssystem

$$\begin{pmatrix} 1 & 0 & 1 & -1 \\ 1 & 0 & -1 & -1 \\ 1 & 1 & 0 & 1 \\ 1 & -1 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \\ c \\ d \end{pmatrix} = 0$$

Dieses System hat nur die triviale Lösung.

**Definition** Gegeben sei eine zulässige Zerlegung  $\mathcal{T}_h$  in  $\Omega$ . Der Finite-Elemente-Raum  $V_h$  ist definiert durch

$$V_h := \{v; \forall K \in \mathcal{T}_h : v_h|_K \in P_K\}$$

Gilt  $V_h \subseteq V$ , so heißt der Raum konform.

**Bemerkung** Für uns hauptsächlich interessant:  $V \subseteq H^1(\Omega)$ . Bekannt: Für stückweise glatte Funktionen, die global stetig sind, liegen in  $H^1$ . Also zu prüfen: globale Stetigkeit.

**Beispiel** (i). Dreieck mit  $\Sigma_K^1$ : Auf der gemeinsamen Kante zweier Dreiecke ist die Funktion eine lineare Funktion einer Variablen. Diese ist also eindeutig bestimmt durch die Funktionswerte der Endpunkte. Also konform.

- (ii). Dreieck mit  $\Sigma_K^2$ : Der Raum ist nicht konform, da die lineare Funktion auf der gemeinsamen Kante nicht durch einen Funktionswert eindeutig bestimmt ist.

Analog für Rechteck.

**Beispiel** (Die wichtigsten konformen Elemente für elliptische Randwertaufgaben zweiter Ordnung) (i). Rechteckelemente:  $Q_k$ -Elemente

- $k = 1$ : Sei

$$P_K := \{(x, y) \mapsto a + b \cdot x + c \cdot y + d \cdot x \cdot y; a, b, c, d \in \mathbb{R}\}$$

$$= \text{lin}\{(x, y) \mapsto x^{\alpha_1} \cdot y^{\alpha_2}; 0 \leq \alpha_i \leq 1\}$$

$$\Sigma_K := \{\delta_{P_i}; i = 1, \dots, 4\}$$

(siehe vorheriges Beispiel).  $Q_1 := (K, P_K, \Sigma_K)$ .  $P_K$  hat nodale Basis

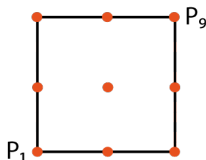
$$\{(x, y) \mapsto (1 - x) \cdot (1 - y), (x, y) \mapsto (1 - x) \cdot y, (x, y) \mapsto (1 - y) \cdot x, (x, y) \mapsto x \cdot y\}$$

- $k = 2$ : Sei

$$P_K := \text{lin}\{(x, y) \mapsto x^{\alpha_1} \cdot y^{\alpha_2}; 0 \leq \alpha_i \leq 2\}$$

$$\Sigma_K := \{\delta_{P_i}; i = 1, \dots, 9\}$$

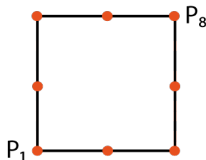
dann  $Q_2 := (K, P_K, \Sigma_K)$ . Die nodalen Basisfunktionen sind die Produkte der eindimensionalen nodalen Basen (Lagrange-Basispolynome).



- $k = 3$ : Verwende bikubische Funktionen, dann  $\dim P_K = 16$ .

Weitere finite Elemente:

- Manchmal versucht man „innere“ Freiheitsgrade zu vermeiden:



8 Freiheitsgrade, also muss  $\dim P_K = 8$  gelten. Wähle  $P_K$  als die lineare Hülle der nodalen Basisfunktionen von  $Q_2$  bis auf  $(x, y) \mapsto (1 - x^2) \cdot (1 - y^2)$  (diese ist auf dem Rand des Quadrates 0).

- Ein weiteres biquadratisches Element (spielt wichtige Rolle bei Superkonvergenz): Sei

$$P_K := \text{lin}\{(x, y) \mapsto x^{\alpha_1} \cdot y^{\alpha_2}; 0 \leq \alpha_i \leq 2\}$$

$$\Sigma_K := \{\delta_{P_i}; i = 1, \dots, 4\} \cup \left\{ f \mapsto \int_{K_i} f; i = 1, \dots, 4 \right\} \cup \left\{ f \mapsto \int_K f \right\}$$

wobei  $\int_{K_i} f$  das (Weg)Integral von  $f$  über die Kante  $K_i$  bezeichnet und  $P_i$  die Eckpunkte.

- (ii). Dreieckelemente:  $P_k$ -Elemente. Die Basisfunktionen beschreibt man mit Hilfe baryzentrischer Koordinaten  $\lambda_i$  (s. Bemerkung).

- $k = 1$ :

$$P_K := \text{lin}\{(x, y) \mapsto x^{\alpha_1} \cdot y^{\alpha_2}; 0 \leq \alpha_1 + \alpha_2 \leq 1\}$$

$$\Sigma_K := \{\delta_{P_i}; i = 1, 2, 3, \}$$

wobei  $P_i$  die Eckpunkte bezeichnen. Die Basisfunktionen sind  $\lambda_1, \lambda_2, \lambda_3$ . Offenbar gilt  $\lambda_i(a_\ell) = \delta_{i\ell}$ .

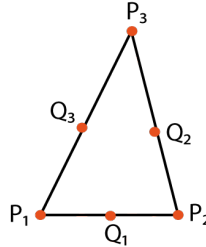


- $k = 2$ :

$$P_K := \text{lin}\{(x, y) \mapsto x^{\alpha_1} \cdot y^{\alpha_2}; 0 \leq \alpha_1 + \alpha_2 \leq 2\}$$

$$\Sigma_K := \{\delta_{P_i}; i = 1, 2, 3, \} \cup \{\delta_{Q_i}; i = 1, 2, 3\}$$

wobei  $Q_i$  die Seitenmittelpunkte der Kanten bezeichnen.



Basisfunktionen: Funktionen, die im Eckpunkt  $P_i$  1 sind und in allen anderen ausgewählten Punkten 0, sind gegeben durch

$$2\lambda_i \cdot \left(\lambda_i - \frac{1}{2}\right)$$

Funktionen, die in  $Q_i$  1 ist und in allen anderen Punkten 0, sind gegeben durch

$$4\lambda_1 \cdot \lambda_2 \quad 4\lambda_2 \cdot \lambda_3 \quad 4 \cdot \lambda_1 \cdot \lambda_3$$

Diese Funktionen bilden zusammen eine Basis.

- $k = 3$ : Sei  $P_K$  die Menge der kubischen Funktionen auf  $K$ , dann  $\dim P_K = 10$ . Sei  $\Sigma_K := \{\delta_{P_i}; i = 1, \dots, 10\}$ . Basisfunktionen:

$$\frac{1}{2}\lambda_i \cdot (3\lambda_i - 1) \cdot (3\lambda_i - 2) \quad i = 1, 2, 3$$

$$\frac{9}{2} \cdot \lambda_i \cdot \lambda_j \cdot (3\lambda_i - 1) \quad i, j \in \{1, 2, 3\}, i \neq j$$

$$27 \cdot \lambda_1 \cdot \lambda_2 \cdot \lambda_3$$

Modifikation: Serendipity-Element  $\tilde{P}_3$ . Sei  $P_K$  die Menge aller kubischen Polynome der Form

$$\sum_i p(a_i) \cdot \frac{1}{2}\lambda_i \cdot (3\lambda_i - 1) \cdot (3\lambda_i - 2) + \sum_{i,j} p(a_{ij}) \cdot \frac{9}{2} \cdot \lambda_i \cdot \lambda_j \cdot (3\lambda_i - 1) + p(a_{123}) \cdot \lambda_1 \cdot \lambda_2 \cdot \lambda_3$$

mit

$$p(a_{123}) = -\frac{1}{6} \sum_i p(a_i) + \frac{1}{4} \sum_{i,j} (p(a_{ij}) + p(a_{ijj})) \quad (*)$$

(damit  $\dim P_K = 9$ ). Warum (\*)? Man möchte, dass  $P_2 \subseteq \tilde{P}_3$ .

Nachweis, dass (\*) für quadratische Funktionen erfüllt ist:

- (1) Bestimme quadratisches Interpolationspolynom  $p$  durch  $a_1, a_{112}, a_2$ . Sei  $a_{12} := \frac{a_1+a_2}{2}$ . Berechne  $p(a_{12})$ .
- (2) Bestimme quadratisches Interpolationspolynom  $p$  durch  $a_1, a_{122}, a_2$ . Sei  $a_{12} := \frac{a_1+a_2}{2}$ . Berechne  $p(a_{12})$ .
- (3) Gleichsetzen der Gleichungen für  $p(a_{12}) = q(a_{12})$  liefert

$$p(a_{12}) = -\frac{1}{16} \cdot (p(a_1) + p(a_2)) + \frac{9}{16} \cdot (p(a_{112}) + p(a_{122}))$$

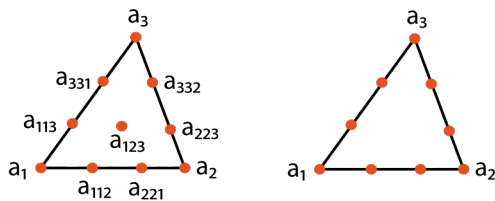


Abbildung 4.1: Links:  $P_3$ -Element, rechts:  $\tilde{P}_3$ -Element

- (4) Bestimmung der quadratischen Interpolierenden  $q$  aus  $p(a_i), p(a_{ij})$ . Berechne  $q(a_{123})$ , das ergibt (\*).

**Bemerkung (Baryzentrische Koordinaten)** Sei  $x = \sum_{i=1}^3 \lambda_i \cdot a^i \in \mathbb{R}^2$  wobei  $\lambda_i \in \mathbb{R}$  ( $i \in \{1, 2, 3\}$ ) mit  $\sum_{i=1}^3 \lambda_i = 1$ .  $\lambda_1, \lambda_2, \lambda_3$  heißen baryzentrische Koordinaten. Analog in  $\mathbb{R}^d$ :  $a^1, \dots, a^{d+1}$  mögen ein nichtentartetes Simplex bilden. Durch

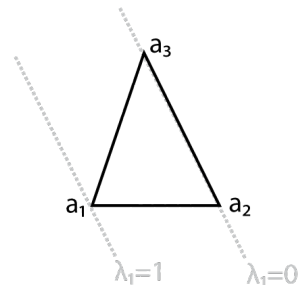
$$x = \sum_{i=1}^{d+1} \lambda_i \cdot a^i \qquad \sum_{i=1}^{d+1} \lambda_i = 1$$

sind eindeutig  $\lambda_1, \dots, \lambda_{d+1}$  bestimmt. Für  $0 \leq \lambda_i \leq 1$  erhält man die konvexe Hülle der gegebenen Punkte. Geometrische Interpretation in  $\mathbb{R}^2$ :

$$\begin{aligned} x_1 &= a_1^1 \cdot \lambda_1 + a_1^2 \cdot \lambda_2 + a_1^3 \cdot \lambda_3 \\ x_2 &= a_2^1 \cdot \lambda_1 + a_2^2 \cdot \lambda_2 + a_2^3 \cdot \lambda_3 \\ 1 &= \lambda_1 + \lambda_2 + \lambda_3 \end{aligned}$$

Cramer'sche Regel:

$$\lambda_1 = \frac{\begin{vmatrix} x_1 & a_1^2 & a_1^3 \\ x_2 & a_2^2 & a_2^3 \\ 1 & 1 & 1 \end{vmatrix}}{\begin{vmatrix} a_1^1 & a_1^2 & a_1^3 \\ a_2^1 & a_2^2 & a_2^3 \\ 1 & 1 & 1 \end{vmatrix}} = \frac{\text{Flächeninhalt } \Delta(x, a_2, a_3)}{\text{Flächeninhalt } K}$$



$\lambda_1$  ist „linear“ in  $x_1, x_2$ . Analog für  $\lambda_2, \lambda_3$ .

**Bemerkung ( $C^1$ -Elemente)**  $C^1$ -Elemente werden für die konforme FEM-Diskretisierung von Randwertaufgaben 4. Ordnung benötigt.

- (i). Rechteck: Es gibt relativ einfache  $C^1$ -Elemente, s. Übung. (In 1D: Kubische Splines haben 4 Freiheitsgrade. Definiere

$$\varphi_i(x_j) := \delta_{ij} \qquad \varphi_i'(x_j) = 0$$

und

$$\psi_i(x_j) = 0 \qquad \psi_i'(x_j) = \delta_{ij}$$

dann bilden die Funktionen jeweils eine Basis.)

- (ii). Dreieck: T10-Element: Sei  $P_K$  die Menge der kubischen Funktionen (also  $\dim P_K = 10$ ). Die Freiheitsgrade seien die Funktionswerte und Ableitungen in den Ecken sowie der Funktionswert im Schwerpunkt. Ist kein  $C^1$ -Element!

Zenisek: Für  $C^1$ -Elemente muss die Dimension des Polynomraumes  $P_K$  mindestens 18 sein. Bekanntes  $C^1$ -Element ist das Argyris-Element. Dabei ist  $P_K$  die Menge der Polynome 5. Grades ( $\dim P_K = 21$ ). 18 Freiheitsgrade in Ecken, dazu Normalenableitungen in Seitenhalbierenden.

Praktisch verwendet man  $C^1$ -Elemente man nicht-konforme Elemente oder gemischte Methoden.

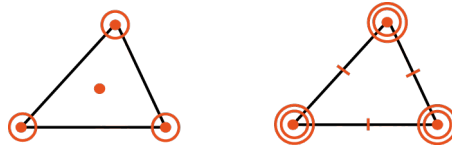


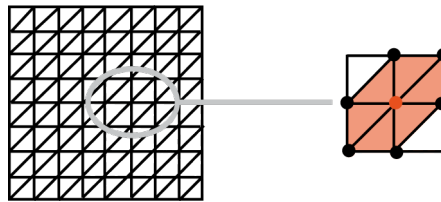
Abbildung 4.2: Links: T10-Element, rechts: Argyris-Element

### 4.3 Einfache FE-Diskretisierung auf Standardgittern

#### Beispiel

$$-\Delta u = f \quad u|_{\partial\Omega} = 0$$

wobei  $\Omega := (0, 1)^2$ . Verwende Friedrichs-Keller-Triangulation und diskretisiere mit linearen FE.



Zu jedem inneren Gitterpunkt gehört eine Basisfunktion. Ansatz:

$$u_h = \sum_{i,j} u_{ij} \cdot \varphi_{ij}$$

wobei  $u_{ij}$  Näherungswerte der Lösung in den Gitterpunkten. Koeffizientenmatrix besteht aus Elementen der Form:

$$\int_{\Omega} \nabla \varphi_{ij} \cdot \nabla \varphi_{\ell k}$$

Typisches Vorgehen bei der Berechnung:

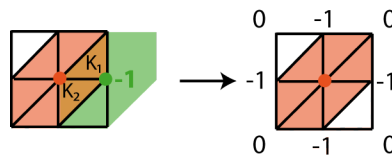
$$\int \nabla \varphi_{ij} \cdot \nabla \varphi_{i+1,j} = \int_{K_1 \cup K_2} \nabla \varphi_{ij} \cdot \nabla \varphi_{i+1,j}$$

mit

$$\begin{aligned} \varphi_{ij}|_{K_1} &= \frac{1}{h} \cdot (h - (x - x_i)) \\ \varphi_{i+1,j}|_{K_1} &= \frac{1}{h} \cdot (x - x_i) \end{aligned}$$

Damit folgt beispielsweise

$$\begin{aligned} \int_{K_1} \nabla \varphi_{ij} \cdot \nabla \varphi_{i+1,j} &= \int_{K_1} -\frac{1}{h^2} = \frac{h^2}{2} \cdot \left(-\frac{1}{h^2}\right) \\ &= -\frac{1}{2} \end{aligned}$$



Analoges Vorgehen gibt

$$\int \nabla \varphi_{ij} \cdot \nabla \varphi_{i+1,j} = -1$$

Ergebnis:

$$-u_{i-1,j} - u_{i+1,j} - u_{i,j+1} - u_{i,j-1} + 4u_{ij} = \int_{\text{spt } \varphi_{ij}} f \cdot \varphi_{ij}$$

Praktische Generierung des diskreten Problems: siehe zweiter Teil der Vorlesung.

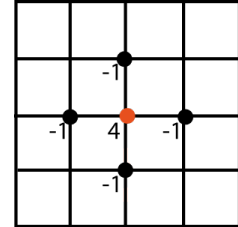
**Bemerkung** (Vergleich mit Differenzenverfahren)

Verwende äquidistantes Gitter mit Gitterbreite  $h$ . Diskretisierung für  $-u''$  (in 1D, siehe (4.3)):

$$\frac{-u_{i+1} + 2u_i - 2u_{i+1}}{h^2}$$

Diskretisierung für  $-\Delta u$ :

$$\frac{u_{i+1,j} + 2u_{ij} - u_{i-1,j}}{h^2} + \frac{u_{i,j+1} + 2u_{ij} - u_{i,j-1}}{h^2} = f(x_i, y_j) =: f_{ij}$$



(Fünf-Punkte-Differenzenstern). Ist also „quasi gleich“ der obigen FEM-Diskretisierung mit linearen FE (bis auf Skalierung). Hat Konsistenzordnung 2.

**Beispiel**

$$-\Delta u = f \quad u|_{\partial\Omega} = 0$$

wobei  $\Omega := (0,1)^2$ . Verwende äquidistantes Gitter und diskretisiere mit bilinearen FE. Ergebnis ist ein 9-Punkte-Stern:

$$\frac{1}{3h^2} \cdot \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix} \quad (4.4)$$

Vergleich mit Differenzenverfahren: Welche 9-Punkt-Differenzenverfahren sind sinnvoll? Dazu zunächst Untersuchung der Konsistenzordnung (in Supremumnorm). Allgemeiner 9-Punkte-Stern:

$$\begin{aligned} &\alpha_1 u(x-h, y+h) + \alpha_2 u(x, y+h) + \alpha_3 u(x+h, y+h) \\ &+ \alpha_4 u(x-h, y) + \alpha_5 u(x, y) + \alpha_6 u(x+h, y) \\ &+ \alpha_7 u(x-h, y-h) + \alpha_8 u(x, y-h) + \alpha_9 u(x+h, y-h) \end{aligned}$$

Taylorentwicklung gibt Bedingung für Konsistenzordnung 2. Ergebnis:

$$\frac{1}{c(h)} \begin{bmatrix} \delta & \nu & \delta \\ \nu & \mu & \nu \\ \delta & \nu & \delta \end{bmatrix} \quad \mu + 4\nu + \delta = 0 \quad \nu + 2\delta = -1$$

Speziell:

- (i).  $\delta = 0$ : 5-Punkte-Stern ( $\nu = -1, \mu = 4$ )
- (ii).  $\nu = \delta = -\frac{1}{3}, \mu = \frac{8}{3}$ : 9-Punkte-Stern aus (4.4)

# 5

## Konvergenzanalyse

### 5.1 Differenzenverfahren

Betrachte das stetige Problem  $Lu = f$  mit Diskretisierung  $L_h u_h = f_h$ . Zu untersuchende Eigenschaften:

- Konsistenz der Ordnung  $p$ :

$$\|L_h r_h u - f_h\| \leq c_K \cdot h^p$$

- Stabilität des diskreten Problems: Für alle  $v_h, w_h \in V_h$  gilt

$$\|v_h - w_h\| \leq C_S \cdot \|L_h v_h - L_h w_h\|$$

Ist  $L_h$  linear, ist dies äquivalent zu  $\|w_h\| \leq C_S \cdot \|L_h w_h\|$ . Ist  $L_h$  zusätzlich bijektiv, dann äquivalent zu  $\|L_h^{-1}\| \leq C_S$ .

#### 5.1 Satz

Stabilität und Konsistenz der Ordnung  $p$  impliziert

$$\|u_h - r_h u\| \leq C_K \cdot C_S \cdot h^p$$

d.h. Konvergenz der Ordnung  $p$ .

Beweis:

$$\|u_h - r_h u\| \leq C_S \cdot \underbrace{\|L_h u_h - L_h r_h u\|}_{f_h} \leq C_S \cdot C_K \cdot h^p \quad \square$$

**Bemerkung** (i). Schwierig zu zeigen ist meist die Stabilität. Die Konsistenz kann mit Hilfe der Taylorentwicklung gezeigt werden.

(ii). Übliche Vektornormen:

- Supremumnorm  $\|\cdot\|_\infty$ : Hilfsmittel sind invers-monotone Matrizen.
- euklidische Norm  $|\cdot|_2$  (diskrete  $L^2$ -Norm): Hilfsmittel ist diskrete Fouriertransformation. Nachteil: konstante Koeffizienten notwendig.
- diskrete  $H^1$ -Norm: s. finite Elemente

**Definition** Sei  $A \in \mathbb{R}^{n \times n}$ .

- Existiert  $A^{-1}$  und gilt  $A^{-1} \geq 0$  (d.h. alle Matrixelemente  $\geq 0$ ), so heißt  $A$  invers-monoton.
- Es gelte  $a_{ij} \leq 0$  für  $i \neq j$  und  $A^{-1} \geq 0$ . Dann heißt  $A$   $M$ -Matrix.
- $A$  heißt (streng) diagonaldominant  $\Leftrightarrow |a_{ii}| > \sum_{j \neq i} |a_{ij}|$ .
- $A$  heißt schwach diagonaldominant  $\Leftrightarrow |a_{ii}| \geq \sum_{j \neq i} |a_{ij}|$  und „>“ gelte für mindestens ein  $i \in \{1, \dots, n\}$ .

(v).  $A$  heißt *reduzibel*  $\Leftrightarrow \exists N_1, N_2 \subsetneq \{1, \dots, n\}, N_1 \cup N_2 = \{1, \dots, n\} \forall i \in N_1, j \in N_2 : a_{ij} = 0$ .

**Bemerkung** Betrachte das Gleichungssystem  $A \cdot x = b, b \geq 0$ . Ist  $A$  *invers-monoton*, dann folgt  $x \geq 0$ .

## 5.2 Satz

Sei  $A \in \mathbb{R}^{n \times n}$  mit  $a_{ij} \leq 0$  für  $i \neq j$ . Dann gilt:

- (i).  $A$  *diagonaldominant*  $\Rightarrow A$  ist  $M$ -Matrix.
- (ii).  $A$  *schwach diagonaldominant* und *irreduzibel*  $\Rightarrow A$  ist  $M$ -Matrix.
- (iii).  $M$ -Kriterium:  $A$  ist  $M$ -Matrix  $\Leftrightarrow \exists e \in \mathbb{R}^n, e > 0 : A \cdot e > 0$ . In diesem Fall gilt

$$\|A^{-1}\|_\infty \leq \frac{\|e\|_\infty}{\min(A \cdot e)}$$

Beweis: (i). siehe Übung 4, Aufgabe 1

(ii). siehe Übung 4, Aufgabe 1

(iii). Beweis der Normabschätzung: Aus  $A^{-1} \geq 0$  folgt

$$\|A^{-1}\|_\infty = \|A^{-1} \cdot (1, \dots, 1)^T\|_\infty$$

(links: Zeilensummennorm). Aus  $A \cdot e \geq \min(A \cdot e) \cdot (1, \dots, 1)^T$  folgt durch Anwendungen von  $A^{-1}$ :

$$A^{-1} \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \leq \frac{e}{\min(A \cdot e)} \quad \square$$

**Beispiel** (i). Betrachte Differenzenverfahren auf  $-\Delta u + u = f$  auf äquidistantem Gitter. Differenzenverfahren führt zu dem Differenzenstern

$$\frac{1}{h^2} \cdot \begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Diese Matrix ist eine  $M$ -Matrix nach 5.2(i).

(ii). Die Matrix

$$a_{ij} := \begin{cases} 2 & i = j \\ -1 & |i - j| = 1 \\ 0 & \text{sonst} \end{cases}$$

erfüllt die Bedingungen aus 5.2(ii) und ist daher eine  $M$ -Matrix. Die Matrix  $B := \frac{1}{h^2} \cdot A$  ist Diskretisierung von  $-u''$ ,  $u(0) = u(1) = 0$  auf äquidistantem Gitter.

Idee: Finde Funktion  $x \mapsto e(x)$  mit  $e(0) = e(1) = 0, e > 0$  in  $(0, 1), -e'' > 0$ .  $x \mapsto e(x) := x \cdot (1 - x)$  erfüllt diese Bedingungen. Der Vektor  $e := r_h e(x)$  erfüllt  $e > 0$  und es gilt  $Ae = (2, \dots, 2)$  (Polynome zweiten Grades werden durch Diskretisierung exakt approximiert). Aus 5.2(iii) folgt:

$$\|B^{-1}\| \leq \frac{1}{2} = \frac{1}{8}$$

**Bemerkung** Nachteil der  $M$ -Matrizen: Theorie ist nur für gewisse kompakte Differenzensterne anwendbar.

## 5.2 Finite Elemente

**Beispiel** Schon gezeigt: Diskretisierung von  $-\Delta u$  mit Friedrichs-Keller-Triangulierung und linearen Elementen führt zu

$$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix}$$

Frage: Erhält man bei einem beliebigem Gitter eine  $M$ -Matrix?

Nutze nodale Basis  $(\varphi_j)_j$ . Für einen Gitterpunkt  $i$  sei  $\Lambda_i$  die Menge der benachbarten Gitterpunkte. Es ergibt sich folgende  $i$ -te Gleichung:

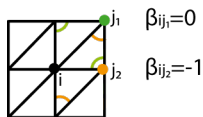
$$\alpha_i \cdot u_i + \sum_{j \in \Lambda_i} \beta_{ij} \cdot u_j = \int_{\Omega} f \cdot \varphi_i$$

wobei

$$0 = \alpha_i + \sum_{j \in \Lambda_i} \beta_{ij}$$

$$\beta_{ij} = -\frac{1}{2} \cdot (\cot \gamma_{ij}^1 + \cot \gamma_{ij}^2)$$

Angewendet auf Friedrichs-Keller-Triangulierung ergibt sich der bereits bekannte Differenzenstern:



Gilt  $\beta_{ij} \leq 0$  für  $i \neq j$ ? Wegen

$$\cot \gamma_{ij}^1 + \cot \gamma_{ij}^2 = \frac{\cos \gamma_{ij}^1}{\sin \gamma_{ij}^1} + \frac{\cos \gamma_{ij}^2}{\sin \gamma_{ij}^2} = \frac{\sin(\gamma_{ij}^1 + \gamma_{ij}^2)}{\sin \gamma_{ij}^1 \cdot \sin \gamma_{ij}^2}$$

ist  $\gamma_{ij}^1 + \gamma_{ij}^2 \leq \pi$  eine hinreichende Bedingung. (Es ist nicht nötig, dass alle Winkel  $\leq \frac{\pi}{2}$ .) Dies impliziert (bis auf Sonderfälle), dass die Steifigkeitsmatrix eine  $M$ -Matrix ist.

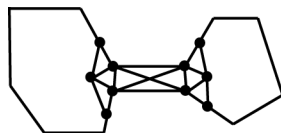


Abbildung 5.1: Beispiel für reduzierbare Steifigkeitsmatrix

**Bemerkung** (i). Für quadratische finite Elemente (und Elemente höherer Ordnung) kann man im Allgemeinen inverse Monotonie der Steifigkeitsmatrix nicht sichern.

(ii). Fehlerabschätzungen in Supremumnorm für lineare Elemente, basierend auf  $M$ -Matrix-Eigenschaften: s. Ciarlet. Resultierende Abschätzungen sind nicht optimal.

**Bemerkung** (zu allgemeinen elliptischen Randwertaufgaben) Sei  $\Omega := (0, 1)^2$  und betrachte die Randwertaufgabe

$$-\Delta u + b \cdot \nabla u + c \cdot u = f \quad u|_{\partial\Omega} = 0$$

Differenzenverfahren auf äquidistantem Gitter führt zu Differenzenstern

$$A := \frac{1}{h^2} \cdot \begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix} + \sum_{j=1}^2 \frac{b_j}{2h} \cdot \begin{bmatrix} 0 & 0 & 0 \\ -1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & c & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Ist  $A$  eine  $M$ -Matrix? Antwort: Ja, für hinreichend kleine  $h$  (denn dann dominiert der erste Term).

Singulär gestörter Fall:

$$-\varepsilon \cdot \Delta u + b \cdot \nabla u + c \cdot u = f \qquad u|_{\partial\Omega} = 0$$

mit  $0 < \varepsilon \ll 1$ . Dann funktioniert die obige Argumentation nicht mehr. Idee: Benutze einseitige Differenzenquotienten.

Upwind-Differenzenverfahren: Seien o.B.d.A.  $b_1, b_2 > 0$ .

$$\frac{\varepsilon}{h^2} \cdot \begin{bmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{bmatrix} + \frac{b_1}{2h} \cdot \begin{bmatrix} 0 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \frac{b_2}{2h} \cdot \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -1 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & 0 \\ 0 & c & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Ist  $M$ -Matrix. Nachteil: Methode 1. Ordnung.

Wie erzeugt man stabile Verfahren im singulär gestörten Fall mit finiten Elementen? Schwierig...

(Konvergenztheorie für finite Elemente)

Ausgang: Cea-Lemma

$$\|u - u_h\| \leq C \cdot \inf_{v_h \in V_h} \|u - v_h\|$$

Abschätzung des Approximationsfehler ist schwierig, deshalb nutzt man

$$\inf_{v_h \in V_h} \|u - v_h\| \leq \|u - \Pi u\|$$

wobei  $\Pi u \in V_h$  die Projektion von  $u$  in  $V_h$ . Zunächst ist es üblich  $\Pi u = u^I$  zu wählen wobei  $u^I$  die „Interpolierende“ von  $u$  in  $V_h$  (dazu wird (bei Lagrange-Elementen)  $u \in H^2$  vorausgesetzt).

**Beispiel** (Pobrowolski, 1990) Sei  $\Omega := (0, 1)^2$ ,

$$-\Delta u = f \qquad u|_{\partial\Omega} = 0$$

Verwende gleichmäßige Zerlegung in Quadrate und bilineare Elemente. Sei

$$c_I := \sup \frac{\|\nabla(u - u^I)\|_0}{\|u\|_2} \cdot h^{-1}$$

$$c := \sup \frac{\|\nabla(u - u_h)\|_0}{\|u\|_2} \cdot h^{-1}$$

wobei  $u^I$  die Interpolierende von  $u$  und  $u_h$  die FE-Lösung. Man kann zeigen:

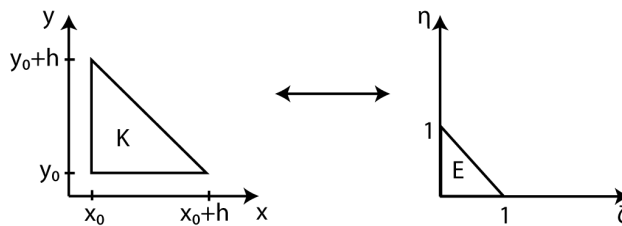
$$\frac{1}{\pi} = c_I = c$$

Wie schätzt man den Projektionsfehler (oft: Interpolationsfehler) ab?

- (i). Transformation des Integrals über ein Element in ein Integral über einen sogenannten Referenzelement.
- (ii). Abschätzung des Fehlers auf dem Referenzelement (oft mit Bramble-Hilbert-Lemma)
- (iii). Rücktransformation

**Beispiel** Betrachte gleichschenkliges Dreieck  $K$ . Sei  $u \in H^2$ . Abzuschätzen sei  $|u - u^I|_{1,K}$  bei linearen finiten Elementen. Betrachte Einheitsdreieck  $E$  als Referenzelement.





Dann ist die Transformation gegeben durch

$$\xi = \frac{x - x_0}{h} \qquad \eta = \frac{y - y_0}{h}$$

Ausführen der obigen 3 Schritte:

(i). Ziel: Abschätzen von  $\int_K ((u - u^I)_x)^2$ . Es gilt

$$\int_K ((u - u^I)_x)^2 = \underbrace{2|K|}_{\text{Funktionaldet.}} \cdot \int ((u - u^I)_\xi)^2 \cdot \frac{1}{h^2}$$

(ii). Abschätzung auf  $E$  (hier wieder  $x \hat{=} \xi$ ,  $y \hat{=} \eta$ ): Hier Abschätzung nicht mit Bramble-Hilbert-Lemma, sondern elementar.

$$\begin{aligned} \partial_x(u - u^I) &= \partial_x u - (u(1, 0) - u(0, 0)) \\ &= \int_0^1 (\partial_x u(x, y) - \partial_x u(\xi, y) + \partial_x u(\xi, y) - \partial_x u(\xi, 0)) d\xi \\ &= \int_0^1 \left( \int_\xi^x \partial_x^2 u(\mu, y) d\mu + \int_0^y \partial_x \partial_y u(\xi, r) dr \right) d\xi \end{aligned}$$

Mit Cauchy-Schwarz-Ungleichung folgt:

$$\int_E (\partial_x(u - u^I))^2 \leq C \cdot \int_E (u_{xx}^2 + u_{xy}^2)$$

(iii). Rücktransformation:

$$\int_E u_{\xi\xi}^2 + u_{\xi\eta}^2 = \frac{1}{2|K|} \cdot \int_K (u_{xx}^2 + u_{xy}^2) \cdot (h^2)^2$$

Es folgt:

$$\int_K ((u - u^I)_x)^2 \leq c \cdot h^2 \cdot \int_K u_{xx}^2 + u_{xy}^2$$

Analog:  $y$ -Ableitung. Summation über  $K$  gibt

$$|u - u^I|_1 \leq c \cdot h \cdot |u|_2$$

(Interpolationsfehler für lineare finite Elemente)

### 5.3 Satz (Bramble-Hilbert-Lemma)

Sei  $q : H^{k+1}(B) \rightarrow \mathbb{R}$  ein sublineares beschränktes Funktional (sublinear: subadditiv  $q(u + v) \leq q(u) + q(v)$  und absolut homogen  $q(\lambda \cdot v) = |\lambda| \cdot q(v)$ ) und es gelte  $q(w) = 0$  für  $w \in P_k$  (Polynome  $k$ -ten Grades). Dann existiert  $C > 0$  mit

$$|q(v)| \leq C \cdot |v|_{k+1}$$

Beweis: (i). Sei  $v \in H^{k+1}(B)$ . Zeige: Es existiert  $w \in P_k$  mit  $\int_B D^\alpha(v+w) = 0$  für alle  $|\alpha| \leq k$ .  
 (Für  $k = 2$ : Dann zu zeigen: Es existiert  $w \in P_2$  mit

$$\int_B D^\alpha(v+w) = 0$$

für  $|\alpha| \leq 2$ . Setze

$$w(x_1, x_2) := c_0 + c_1 \cdot x_1 + c_2 \cdot x_2 + c_3 \cdot x_1^2 + c_4 \cdot x_1 \cdot x_2 + c_5 \cdot x_2^2$$

Wähle beispielsweise  $\alpha = (2, 0)$ , dann gibt  $\int w_{x_1, x_1} = -\int v_{x_1, x_1}$  die Konstante  $c_3$ . Analog für  $\alpha = (1, 1)$ ,  $\alpha = (0, 2)$ . Betrachte dann  $\alpha = (1, 0)$ ,  $(0, 1)$ ,  $(0, 0)$ .

(ii). Genutzt wird die Poincaré-Ungleichung

$$\|u\|_{k+1}^2 \leq c \cdot \left( |u|_{k+1}^2 + \sum_{|\alpha| \leq k} \int_B |D^\alpha u|^2 \right)$$

Dann folgt

$$\|v+w\|_{k+1}^2 \leq C \cdot \|v+w\|_{k+1}^2$$

(iii). Es gilt wegen der Sublinearität:

$$\begin{aligned} |q(v)| &\leq |q(v+w)| + \underbrace{|q(w)|}_0 \leq C \cdot \|v+w\|_{k+1} \\ &\stackrel{(ii)}{\leq} c \cdot \|v+w\|_{k+1} \stackrel{w \in P_k}{=} c \cdot |v|_{k+1} \end{aligned} \quad \square$$

**Bemerkung** Im letzten Beispiel:  $k = 1$ ,  $q(v) = |v - v^I|_{1,E}$  für  $v \in H^2$ .

Untersucht werden jetzt Lagrange-Elemente, zudem sogenannte affin-äquivalente Familien.

**Definition** (Äquivalenz von Lagrange-Elementen) Sei  $(K, P, \Sigma)$  ein finites Element und  $(\hat{K}, \hat{P}, \hat{\Sigma})$  ein Referenzelement. Es existiere eine bijektive Abbildung  $F: \hat{K} \rightarrow K$  mit

- (i).  $P = \{p: K \rightarrow \mathbb{R}; p \circ F \in \hat{P}\}$
- (ii).  $\{F(\hat{a}); \hat{a} \in \hat{K} \text{ erzeugt Freiheitsgrade auf } \hat{K}\} = \{a; a \in K \text{ erzeugt Freiheitsgrade auf } K\}$

Ist  $F$  zusätzlich affin-linear, dann handelt es sich um eine affin-äquivalente Familie.

**Beispiel** (i).  $P_k$ -Elemente: Zwei Dreiecke  $K, \hat{K}$ . Dann existiert eine bijektive affin-lineare Abbildung  $F$  wie in der letzten Definition.

(ii).  $Q_k$ -Elemente: Bildet man ein Quadrat unter einer affinen Abbildung ab, dann erhält man ein Parallelogramm. Triangulierung in Parallelogrammen (mit Affin-Äquivalenz) also möglich. Im Allgemeinen für Rechtecke nicht.

#### 5.4 Lemma

Sei  $K \subseteq \mathbb{R}^d$  ein Dreieck und  $K'$  ein Referenzelement,  $K \ni x = \varphi(p) := B \cdot p + b$  mit  $B$  regulär. Ist  $u \in H^\ell(K)$ , so gilt  $v := u \circ \varphi \in H^\ell(K')$  und

$$\begin{aligned} |v|_{\ell, K'} &\leq c \cdot \|B\|^\ell \cdot (\det B)^{\frac{1}{2}} \cdot |u|_{\ell, K} \\ |u|_{\ell, K} &\leq c \cdot \|B^{-1}\|^\ell \cdot (\det B)^{\frac{1}{2}} \cdot |v|_{\ell, K'} \end{aligned}$$

Beweis: (für  $\ell = 1$ ) Nach Definition gilt

$$|v|_{1,K'}^2 = \int_{K'} \sum_i \left( \frac{\partial v}{\partial p_i} \right)^2$$

Wegen

$$\begin{aligned} \frac{\partial v}{\partial p_j} &= \sum_i \frac{\partial u}{\partial x_i} \cdot \frac{\partial x_i}{\partial p_j} \\ \Rightarrow \left| \frac{\partial v}{\partial p_j} \right| &\leq \max_i \left| \frac{\partial u}{\partial x_i} \right| \cdot \underbrace{\sum_i \left| \frac{\partial x_i}{\partial p_j} \right|}_{b_{ij}} \\ &\leq c \cdot \max_i \left| \frac{\partial u}{\partial x_i} \right| \cdot \|B\| \end{aligned}$$

(Beachte dazu, dass alle Normen auf  $\mathbb{R}^{d \times d}$  äquivalent sind.) Damit folgt aus Transformationssatz:

$$\begin{aligned} |v|_{1,K'}^2 &\leq c \cdot \|B\|^2 \cdot (\det B)^{-1} \cdot \int_K \left( \max_i \left| \frac{\partial u}{\partial x_i} \right| \right)^2 \\ &\leq C \cdot \|B\|^2 \cdot (\det B)^{-1} \cdot \int_K \sum_i \left( \frac{\partial u}{\partial x_i} \right)^2 \quad \square \end{aligned}$$

**Bemerkung** Verallgemeinerung:

$$|v|_{\ell,p,K'} \leq C \cdot \|B\|^\ell \cdot (\det B)^{-\frac{1}{p}} \cdot |u|_{\ell,p,K}$$

für  $p \in [1, \infty]$ .

### 5.5 Lemma

Sei  $K \subseteq \mathbb{R}^d$  ein Dreieck und  $K'$  ein Referenzelement. Es seien  $\varrho, R > 0$  und  $x_1, x_2 \in \mathbb{R}^d$  mit  $B(x_1, \varrho) \subseteq K \subseteq B(x_2, R)$ . Dann existieren  $c_1, c_2 > 0$  mit

$$\|B\| \leq c_1 \cdot R \qquad \|B^{-1}\| \leq \frac{c_2}{\varrho}$$

Beweis: Für  $K'$  existieren  $\varrho', R'$  und  $p_1, p_2$  wie für  $K$ . Folglich gilt  $p_1 + p \in K'$  für alle  $p \in B(0, \varrho')$ . Sei  $x_0 := B \cdot p_1 + b \in K$  und  $x := B \cdot (p_1 + p) + b \in K$ . Dann gilt  $\|x - x_0\| \leq 2R$ . Nun ist

$$\|B\| = \frac{1}{\varrho'} \cdot \sup_{\|p\|=\varrho'} \|Bp\| = \frac{1}{\varrho'} \sup_{\|p\|=\varrho'} \|x - x_0\| \leq \frac{2}{\varrho'} \cdot R$$

Analog: 2. Abschätzung. □

### 5.6 Satz

Gegeben sei eine affine Familie von Lagrange-FE über einer Zerlegung  $Z$  von  $\Omega$  sowie eine Projektion  $\Pi_{Z,k} : H^{k+1}(\Omega) \rightarrow P_{Z,k}$  (stückweise Polynome vom Grad  $k$  über Zerlegung). Ziel: Abschätzung des Projektionsfehlers  $|u - \Pi_{Z,k}u|_{r,K}$ . Dann gilt

$$|u - \Pi_{Z,k}u|_{r,K} \leq c \cdot \frac{(\text{diam } K)^{k+1}}{\varrho_K^r} \cdot |u|_{k+1,K}$$

falls  $K$  eine Kugel vom Radius  $\varrho$  enthält.

Beweis: (i). Transformation: Nach Lemma 5.4 gilt

$$|u - \Pi_{Z,k}u|_{r,K} \leq c \cdot \|B^{-1}\|^r \cdot (\det B)^{\frac{1}{2}} \cdot \|v - \Pi_{K',k}v\|_{r,K'}$$

(ii). Bramble-Hilbert-Lemma:

$$|u - \Pi_{Z,k} u|_{r,K} \leq c \cdot \|B^{-1}\|^r \cdot (\det B)^{\frac{1}{2}} \cdot |v|_{k+1,K'}$$

(iii). Rücktransformation: Lemma 5.4 gibt

$$|u - \Pi_{Z,k} u|_{r,K} \leq c \cdot \|B^{-1}\|^r \cdot \|B\|^{k+1} \cdot |u|_{k+1,K}$$

$K$  enthalte eine Kugel mit Radius  $\varrho_K$ . Dann gilt mit Lemma 5.5:

$$|u - \Pi_{Z,k} u|_{r,K} \leq c \cdot \frac{(\text{diam } K)^{k+1}}{\varrho_K^r} \cdot |u|_{k+1,K} \quad \square$$

**Definition** Man nennt eine Familie von Zerlegungen shape-regulär (quasi-uniform), wenn ein  $c > 0$  existiert, sodass für alle  $K$  dieser Familie die Abschätzung  $\frac{\text{diam } K}{\varrho_K} \leq c$  gilt.

**Bemerkung** (i). Quasi-uniform wird manchmal anders definiert.

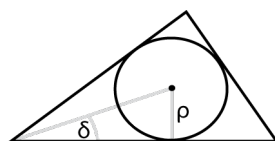


Abbildung 5.2: Shape-Regularität bedeutet anschaulich: Für die Familie der Zerlegungen darf  $\delta$  nicht gegen Null gehen.

(ii). In  $2D$ , Dreiecke: Shape-regulär  $\Leftrightarrow$  Für alle Winkel  $\alpha$  gilt  $\alpha \geq \alpha_0$  für ein  $\alpha_0 > 0$  (Minimalwinkelbedingung).

(iii). Für eine shape-reguläre affine Familie von Lagrange-FE gilt also nach Theorem 5.6:

$$|u - \Pi_{Z,k} u|_{r,K} \leq c \cdot (\text{diam } K)^{k+1-r} \cdot |u|_{k+1,K}$$

Nächstes Ziel: Globale Abschätzung.

### 5.7 Korollar

Gegeben sei eine affine Familie von Lagrange-FE über einer Zerlegung  $Z$  von  $\Omega$  sowie eine Projektion  $\Pi_{Z,k} : H^{k+1}(\Omega) \rightarrow P_{Z,k} \cap H^r(\Omega)$ . Dann gilt

$$|u - \Pi_{Z,k} u|_{r,\Omega} \leq c \cdot h^{k+1-r} \cdot |u|_{r+1,\Omega}$$

mit  $h := \max_K \text{diam } K$ .

**Beispiel**  $d \in \{2, 3\}$ ,  $u \in H^2(\Omega)$ ,  $P_k$ -Elemente

(i). lineare Interpolation ( $k = 1$ ):  $u^I$  sei lineare Interpolierende von  $u$  (existiert wegen  $d \in \{2, 3\}$ ,  $u \in H^2$ ). Dann gilt  $u^I \in H^1$ . Für  $r = 0$  bzw.  $r = 1$  in Korollar 5.7 folgt

$$\begin{aligned} \|u - u^I\|_0 &\leq c \cdot h^2 \cdot |u|_2 \\ |u - u^I|_1 &\leq c \cdot h \cdot |u|_2 \end{aligned}$$

(Die Abschätzungen sind bzgl.  $h$ -Potenzen optimal.) Es gilt auch:

$$\|u - u^I\|_\infty \leq c \cdot h^2 \cdot |u|_{2,\infty}$$

(ii). Interpolation durch Polynome vom Grad  $k$ : Sei dazu  $u \in H^{k+1}(\Omega)$ . Wegen  $u^I \in H^1$  kann in Korollar 5.7 nur  $r = 0$  oder  $r = 1$  gewählt werden.

$$\begin{aligned} \|u - u^I\|_0 &\leq c \cdot h^{k+1} \cdot |u|_{k+1} \\ |u - u^I|_1 &\leq c \cdot h^k \cdot |u|_{k+1} \end{aligned}$$

**5.8 Satz (A-priori-Abschätzung)**

Betrachte elliptische Randwertaufgabe 2. Ordnung, sodass die Voraussetzungen des Lax-Milgram-Lemmas erfüllt sind ( $V \subseteq H^1$ ). Benutze Cea-Lemma, um 5.7 für FEM-Abschätzungen zu verwenden.

$$\|u - u_h\|_1 \leq C \cdot \|u - u^I\|_1$$

Weitere Voraussetzungen:

- zulässige shape-reguläre Zerlegung
- $P_k$ -Elemente (2D/3D)
- $u \in H^{k+1}(\Omega)$

Dann:

$$\|u - u_h\|_1 \leq c \cdot h^k \cdot |u|_{k+1}$$

**Bemerkung** (i). Problem: Im Allgemeinen sind  $C$  sowie  $|u|_{k+1}$  unbekannt.

- (ii). Kann man für den FE-Fehler analoge Abschätzungen für andere Normen beweisen? Problem: Cea-Lemma nur für  $\|\cdot\|_1$ -Norm anwendbar.

**5.9 Satz (Abschätzung des  $L^2$ -Fehlers)**

Es seien die Voraussetzungen aus Satz 5.8 erfüllt. Betrachte das Hilfsproblem

$$a(v, w) = (g, v) \quad (v \in V) \tag{5.1}$$

für  $g \in L^2$  („duales Problem“). (5.1) besitzt nach Lax-Milgram-Lemma eine eindeutige Lösung  $w$ . Setze nun zusätzlich voraus, dass  $w \in H^2(\Omega)$  und  $\|w(g)\|_2 \leq c \cdot \|g\|_0$ . Dann gilt

$$\|u - u_h\|_0 \leq C \cdot h^{k+1} \cdot |u|_{k+1}$$

Beweis: Dualitätstrick (oder Nitschetrick). Wähle  $g := u - u_h$  und  $v := u - u_h$ , dann

$$\begin{aligned} \|u - u_h\|_0^2 &= (g, v) = a(u - u_h, w) = a(u - u_h, w - w_h) \\ &\leq M \cdot \|u - u_h\|_1 \cdot \|w - w_h\|_1 \end{aligned}$$

für beliebige  $w_h \in V_h$  (Galerkinorthogonalität). Aus Theorem 5.8 folgt für  $w_h := w^I$ :

$$\begin{aligned} \|u - u_h\|_0^2 &\leq C \cdot h^k \cdot |u|_{k+1} \cdot h \cdot \|w\|_2 \\ &\leq c \cdot C \cdot h^{k+1} \cdot |u|_{k+1} \cdot \|u - u_h\|_0 \end{aligned} \quad \square$$

**Bemerkung** (i). Randwertaufgabe  $-\Delta w = g$ ,  $w|_{\partial\Omega} = 0$  mit  $\Omega$  konvex erfüllt die Voraussetzungen des Dualitätstricks.

- (ii). Randwertaufgabe  $-\varepsilon \cdot \Delta u + b \cdot \nabla u + c \cdot u = f$ ,  $u|_{\partial\Omega} = 0$  (singulär gestörtes Problem).  $c$  ist dann gegeben durch  $c = c' \cdot \varepsilon^{-\frac{3}{2}}$ , d.h. Nitsche-Trick ist nicht anwendbar.

(iii). Analoge Abschätzungen gelten auch für  $Q_k$ -Elemente.

**Bemerkung (Numerische Konvergenzrate)** Sei  $u - u_h = e_h$ . Annahme:  $e_h \leq C \cdot h^\beta$ . Frage: Beobachtet man praktisch eine Konvergenzrate  $\beta$ ?

Heuristische Überlegung: Angenommen  $e_h = C \cdot h^\beta$ . Dann

$$e_{\frac{h}{2}} = c \cdot \left(\frac{h}{2}\right)^\beta$$

(Verfeinerung der Schrittweite). Dann folgt

$$\begin{aligned} \ln e_h &= \ln C + \beta \cdot \ln h \\ \ln e_{\frac{h}{2}} &= \ln C + \beta \cdot (\ln h - \ln 2) \\ \Rightarrow \beta &= \frac{\ln e_h - \ln e_{\frac{h}{2}}}{\ln 2} \end{aligned}$$

$\beta$  heißt numerische Konvergenzrate.

**Beispiel** (i).  $-\Delta u = f$  in  $(0, 1)^2$ ,  $u|_{\partial\Omega} = 0$  Sei

$$u(x) := \sin(\pi \cdot x) \cdot \sin(\pi \cdot y)$$

(Idee: Gebe exakte Lösung  $u$  vor und berechne daraus  $f$ .) Berechnung der numerischen (experimentellen) Konvergenzraten bestätigt die theoretischen Ergebnisse.

(ii).  $-\Delta u = 0$ ,  $\Omega := \{(r, \varphi); 0 \leq r < 1, 0 < \varphi < \frac{3\pi}{2}\}$  (nicht konvex!), exakte Lösung

$$u(x) := r^{\frac{2}{3}} \cdot \sin\left(\frac{2}{3}\varphi\right) \notin H^2$$

Numerische Konvergenzrate  $\approx \frac{2}{3}$  für  $P_1$ -,  $P_2$ -Elemente.

### 5.10 Lemma (Inverse Ungleichung)

Sei  $d = 2$  und benutze lineare finite Elemente. Sei  $\frac{\max_K \text{diam } K}{\text{diam } K} \leq C$  auf der Familie von shape-regulären Triangulationen (d.h. eine uniforme Zerlegung). Dann gilt

$$\|v_h\|_{\infty} \leq \frac{C}{h} \cdot \|v_h\|_0 \quad (v_h \in V_h)$$

**Bemerkung** (i). Abschätzung kann nur gelten, weil  $V_h$  endlich-dimensional ist (d.h. alle Normen sind äquivalent). Aber: Konstante abhängig von Raumdimension!

(ii). Manchmal wird die Bedingung  $\frac{\max_K \text{diam } K}{\text{diam } K} \leq C$  auch quasi-uniform genannt.

(iii). Bei vorausgesetzter Uniformität sind lokale Verfeinerungen nicht zulässig. Starke Voraussetzung!

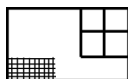


Abbildung 5.3: Lokale Verfeinerungen führen zur Verletzung der Uniformität.

Beweis:

Verwende eine Quadraturformel, die für quadratische Funktionen exakt ist.

Dann gilt

$$\int_K v_h^2 = \frac{\lambda(K)}{60} \cdot \left( 3 \sum v_{i,K}^2 + 8 \cdot \sum v_{ij,K}^2 + 27 \cdot v_{ijk,K}^2 \right)$$

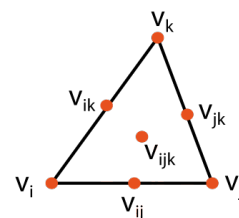
Da  $v_h$  linear ist, liegt das Maximum von  $v_h$  in einem Eckpunkt. Damit

$$\int_K v_h^2 \geq C \cdot \lambda(K) \cdot \max_K |v_h^2|$$

Da Triangulation shape-regulär ist, folgt

$$\|v_h\|_{\infty, K} \leq \frac{c}{\text{diam } K} \cdot \|v_h\|_{0, K}$$

(lokale inverse Ungleichung, bisher wurde Uniformität noch nicht benötigt). Summation und Uniformität gibt Behauptung.  $\square$



**5.11 Satz** ( $L^\infty$ -Abschätzung für lineare finite Elemente)

Die Familie der Zerlegungen sei uniform und shape-regulär. Weiterhin gelte  $u \in W_{2,\infty}$  und der Nitsche-Trick sei anwendbar. Dann gilt

$$\|u - u_h\|_\infty \leq C \cdot h$$

Beweis: Aus Lemma 5.10 und Theorem 5.7 folgt

$$\begin{aligned} \|u - u_h\|_\infty &\leq \|u - u^I\|_\infty + \|u^I - u_h\|_\infty \\ &\leq C \cdot h^2 \cdot |u|_{2,\infty} + \frac{C}{h} \cdot \|u^I - u_h\|_0 \\ &\leq C \cdot h^2 \cdot |u|_{2,\infty} + \frac{C}{h} + (\|u^I - u\|_0 + \|u - u_h\|_0) \\ &\leq C \cdot (h^2 \cdot |u|_{2,\infty} + h \cdot |u|_2) \end{aligned} \quad \square$$

**Bemerkung** Ist diese Abschätzung optimal?  $\|u - u^I\|_\infty = O(h^2) \dots$ ? Praxis zeigt: Konvergenzordnung 1 ist nicht optimal für  $u \in W_{2,\infty}$ . Es gilt

$$\|u - u_h\|_\infty \leq C \cdot h^2 \cdot |\ln h|$$

Allgemein: Die Anwendung globaler inverser Ungleichungen führt oft zu nicht-optimalen Abschätzungen.

**Bemerkung** (Versuch einer  $L^\infty$ -Abschätzung für lineare finite Elemente (2D)) Einbettungssatz gibt  $H^2 \hookrightarrow L^\infty$ , also

$$\|u_h - u\|_{\infty, K'}^2 \leq C \cdot \|u_h - u\|_{2, K'}^2$$

wobei  $K'$  das Referenzelement bezeichne. Damit folgt

$$\begin{aligned} \|u_h - u\|_{\infty, K} &\leq c \cdot \|u_h - u\|_{\infty, K'} \\ \Rightarrow \max_K \|u_h - u\|_{\infty, K}^2 &\leq c^2 \cdot \|u_h - u\|_{\infty, K'}^2 \leq c^2 \cdot C \cdot \|u_h - u\|_{2, K'}^2 \\ &= C' \cdot (\|u_h - u\|_{2, K'}^2 + |u_h - u|_{1, K'}^2 + \|u_h - u\|_0^2) \end{aligned}$$

Transformation auf  $K$ :

$$\|u_h - u\|_\infty^2 \leq \frac{C'}{\lambda(K)} \cdot (h^4 \cdot |u_h - u|_{2, K}^2 + h^2 \cdot |u_h - u|_{1, K}^2 + \|u_h - u\|_{0, K}^2)$$

Summation über alle  $K$  gibt

$$\begin{aligned} \|u_h - u\|_\infty^2 \cdot \lambda(\Omega) &\leq C' \cdot (h^4 \cdot |u_h - u|_2^2 + h^2 \cdot |u_h - u|_1^2 + \|u_h - u\|_0^2) \\ &\leq C' \cdot (h^4 \cdot |u|_2^2 + h^4 \cdot |u|_2^2 + h^4 \cdot |u|_2^2) \\ \Rightarrow \|u_h - u\|_{\infty, K} &\leq C \cdot h^2 \cdot |u|_2 \end{aligned}$$

Diese Aussage ist falsch! (Veröffentlicht 1986 in „Numerische Mathematik“). Problem:

$$\|u_h - u\|_{\infty, K} \leq c(K) \cdot \|u_h - u\|_{\infty, K'}$$

d.h. Konstante hängt von  $K$  ab. Globale Abschätzung (wie oben verwendet) i.A. nicht möglich.

**Bemerkung** Optimale  $L^\infty$ -Abschätzungen nutzen die Greensche Funktion (oder eine Approximation). Idee: Sei  $x_0 \in \Omega$ .  $G$  löse

$$\forall v \in V : a(v, G) = v(x_0)$$

Dann folgt

$$\begin{aligned} (u - u_h)(x_0) &= a(u - u_h, G) = a(u - u_h, G - G_h) \\ &\leq M \cdot \|u - u_h\|_1 \cdot \|G - G_h\|_1 \end{aligned}$$

Hoffnung:  $\|u - u_h\|_1$  gibt  $h^k$  und  $\|G - G_h\|_1$  gibt  $h$ . Idee nur so für  $n = 1$  realisierbar ( $v \in H^1$  besitzt für  $n \geq 2$  keine Werte in Punkten). Modifikation: Suche stattdessen  $G \in W^{1,1}$ , sodass

$$\forall v \in W^{1,\infty} : a(v, G) = v(x_0)$$

Neue Existenztheorie nötig. Dann

$$|(u - u_h)(x_0)| \leq M \cdot \|u - u_h\|_{W^{1,\infty}} \cdot \|G - G_h\|_{W^{1,1}}$$

Auf uniformen Gittern kann man dann beweisen, dass

$$\|u_h - u\|_\infty \leq \begin{cases} h^2 \cdot |\ln h| & P_1 - \text{Elemente} \\ h^{k+1} & P_k - \text{Elemente, } k \geq 2 \end{cases}$$

für  $u \in W^{k+1,\infty}$

### 5.12 Satz (Allgemeine inverse Ungleichungen)

$v$  sei eine Funktion aus einem FE-Raum und Zerlegung sei shape-regulär.

$$\|v\|_{\ell,p,K} \leq C \cdot (\text{diam } K)^{m-\ell+n \cdot (\frac{1}{p}-\frac{1}{q})} \cdot \|v\|_{m,q,K}$$

wobei  $n$  Raumdimension,  $\ell \geq m$ ,  $p, q \in [1, \infty]$ .

Beweis: (i).  $m = 0$ :

$$\begin{aligned} |v|_{\ell,p,K} &\leq c \cdot (\text{diam } K)^{-\ell} \cdot h^{\frac{n}{p}} \cdot |v|_{\ell,p,K'} \\ \Rightarrow \|v\|_{\ell,p,K} &\leq c \cdot (\text{diam } K)^{-\ell} \cdot h^{\frac{n}{p}} \cdot \|v\|_{\ell,p,K'} \\ &\leq c \cdot (\text{diam } K)^{-\ell} \cdot h^{\frac{n}{p}} \cdot \|v\|_{0,q,K'} \end{aligned}$$

(Im letzten Schritt: Normäquivalenz auf endlich-dimensionalen Räumen.) Rücktransformation gibt

$$\|v\|_{\ell,p,K} \leq C \cdot (\text{diam } K)^{-\ell} \cdot h^{\frac{n}{p}-\frac{n}{q}} \cdot \|v\|_{0,q,K}$$

(ii).  $m \geq 1$ : Sei zunächst  $|\alpha| < \ell - m$ , dann

$$\begin{aligned} \|\partial^\alpha v\|_{0,p,K} &\leq \|v\|_{\ell-m,p,K} \stackrel{(i)}{\leq} C \cdot (\text{diam } K)^{m-\ell+n \cdot (\frac{1}{p}-\frac{1}{q})} \cdot \|v\|_{0,q,K} \\ &\leq C \cdot (\text{diam } K)^{m-\ell+n \cdot (\frac{1}{p}-\frac{1}{q})} \cdot \|v\|_{m,q,K} \end{aligned}$$

Für  $\ell - m \leq |\alpha| \leq \ell$ : Dann  $\alpha = \beta + \gamma$  mit  $|\beta| = \ell - m$  (dann  $|\gamma| \leq m$ ),

$$\begin{aligned} \|\partial^\alpha v\|_{0,p,K} &= \|\partial^\beta (\partial^\gamma v)\|_{0,p,K} \leq \|\partial^\gamma v\|_{\ell-m,p,K} \\ &\stackrel{(i)}{\leq} C \cdot (\text{diam } K)^{m-\ell+n \cdot (\frac{1}{p}-\frac{1}{q})} \cdot \|\partial^\gamma v\|_{0,q,K} \\ &\leq C \cdot (\text{diam } K)^{m-\ell+n \cdot (\frac{1}{p}-\frac{1}{q})} \cdot \|v\|_{m,q,K} \end{aligned} \quad \square$$

**Beispiel (Zur Winkelbedingung)** Betrachte Transformation auf Referenzelement  $E$ :

$$\xi = \frac{x - x_0}{h_1} \qquad \eta = \frac{y - y_0}{h_2}$$

Sei  $w := u - u^I$ .

(i). Interpolationsfehler (in  $L^2$ -Norm):

$$\begin{aligned} \int_K w^2 &= D \cdot \int_E w^2 \stackrel{5.3}{\leq} c \cdot D \cdot |w|_{2,E}^2 \\ &\leq c \cdot \int_K (h_1^4 u_{xx}^2 + h_1^2 \cdot h_2^2 \cdot u_{xy}^2 + h_2^4 \cdot u_{yy}^2) \end{aligned}$$

(anisotrope Interpolationsfehlerabschätzung).



(ii). Interpolationsfehler (in  $H^1$ -Seminorm): Naiver Versuch:

$$\int_K w_y^2 \leq c \cdot D \cdot \int_E w_\eta^2 \cdot \frac{1}{h_2^2} \stackrel{5.3}{\leq} c \cdot \int_K \frac{1}{h_2^2} \cdot (h_1^4 u_{xx}^2 + h_1^2 \cdot h_2^2 \cdot u_{xy}^2 + h_2^4 \cdot u_{yy}^2)$$

Der Term  $\frac{h_1^4}{h_2^2}$  ist störend (kein beliebiges Verhältnis von  $h_1$  zu  $h_2$ , falls sich dieser „schön“ verhalten soll). Statt Bramble-Hilbert-Lemma muss also andere Abschätzung auf Referenzelement verwendet werden.

Neuer Versuch: Es gilt zum Beispiel

$$\frac{\partial}{\partial y}(u - u^I) = \int_0^1 \left( \int_\eta^y \partial_{22} u(x, \nu) d\nu + \int_0^x \partial_{12} u(\nu, \eta) d\nu \right) d\eta$$

(vgl. auch Beispiel vor Satz 5.3). Damit

$$\begin{aligned} \int_K w_y^2 &\leq c \cdot D \cdot \int_E w_\eta^2 \cdot \frac{1}{h_2^2} \leq c \cdot D \cdot \frac{1}{h_2^2} \cdot \int_E w_{\xi,\eta}^2 + w_{\eta,\eta}^2 \\ &\leq c \cdot \left( \int_K h_1^2 \cdot u_{xy}^2 \cdot h_2^2 \cdot u_{yy}^2 \right) \end{aligned}$$

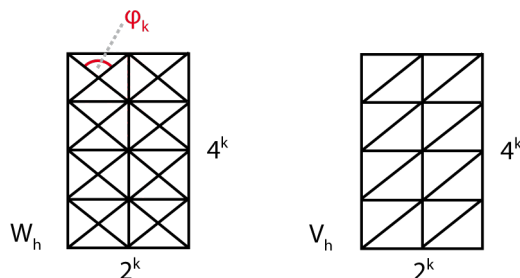
(anisotrope Abschätzung).

Damit: Spitze Winkel verursachen (für lineare finite Elemente) keine größeren Probleme. Schon 1976 wurde bewiesen, dass

$$\inf_{v_h \in V_h} \|u - v_h\| \leq c \cdot h \cdot |u|_2$$

(lineare finite Elemente, 2D), falls die Maximalwinkelbedingung gilt (d.h.  $\exists \alpha_0 < \pi : \alpha \leq \alpha_0$ ).

**Beispiel** (Maximalwinkelbedingung nicht notwendig für Konvergenz von linearen FEM (Krzek u.a. 2012))

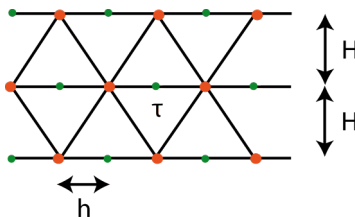


Es gilt  $\varphi_k \rightarrow \pi$  für  $k \rightarrow \infty$ , d.h. die Maximalwinkelbedingung in  $W_h$  ist verletzt. Wegen  $W_h \subseteq V_h$  gilt

$$\inf_{v_h \in W_h} \|u - v_h\| \leq \inf_{v_h \in V_h} \|u - v_h\| \leq c \cdot h \cdot |u|_2$$

(denn in  $V_h$  ist die Maximalwinkelbedingung erfüllt.)

**Beispiel** (Nichtkonvergenz einer FEM (Babuska, Aziz 1976))



Betrachte lineare finite Elemente. Sei  $H = h_0 \cdot h^\beta$  mit  $h_0 \in \mathbb{R}$ . Annahme:

$$\min_{v_h \in V_h} \|u - v_h\|_1 =: c_1 \rightarrow 0 \quad (h \rightarrow 0)$$

dann Konvergenz der FEM gegen  $u$ , nach Cea-Lemma. Für  $u := x^2$

$$\int_{\tau} \left( \frac{\partial}{\partial y} v_h \right)^2 \leq c_1^2$$

$$\Rightarrow \left| \frac{\partial v_h}{\partial y} \right| \leq \frac{c_2 \cdot c_1}{(H \cdot h)^{\frac{1}{2}}} = c_2 \cdot c_1 \cdot h^{-\frac{1}{2} - \frac{\beta}{2}}$$

Sei  $\eta := u - v_h$  und

$$\varphi(x, y) := \frac{1}{2} \cdot (\eta(x + h, y) + \eta(x - h, y)) - \eta(x, y)$$

$$\psi(x, y) := \frac{1}{2} \cdot (v_h(x + h, y) + v_h(x - h, y)) - v_h(x, y)$$

(Differenzenquotienten!) Damit folgt  $\varphi = h^2 - \psi$  (denn quadratische Funktionen werden exakt diskretisiert). Da  $v_h$  linear ist, gilt  $\psi(\text{grüne Punkte}) = 0$ . Deshalb folgt mittels Taylorentwicklung:

$$|\psi(i \cdot h, y)| \leq H \cdot c_3 \cdot c_1 \cdot h^{-\frac{1+\beta}{2}} = c_4 \cdot c_1 \cdot h^{\frac{\beta-1}{2}}$$

Wähle  $\beta \geq 5$ , dann

$$\varphi(i \cdot h, y) \geq c_5 \cdot h^2$$

mit  $c_5 > 0$ . Seien  $x, y \in [-1, 1]$ . Setze

$$z_i := \int_{-1}^1 \eta(i \cdot h, y) dy$$

$$\vartheta_i := \int_{-1}^1 \varphi(i \cdot h, y) dy$$

Offenbar gilt dann

$$\vartheta_i \geq c_6 \cdot h^2 \qquad z_{i-1} - 2z_i + z_{i+1} = 2\vartheta_i$$

Aus  $\|\eta\|_1 \leq c_1$  folgt  $|z_i| \leq C \cdot c_1$  (es gilt:  $u \in H^1 \Rightarrow u \in L^\infty \otimes L^1$ , „anisotroper Einbettungssatz“).  
Aus

$$\Delta z_i := -z_{i-1} + 2z_i - z_{i+1} = -2\vartheta_i \leq -2c_6 \cdot h^2$$

Anwendung des diskreten Maximumprinzips mit Schrankenfunktion

$$w_i := \alpha \cdot c_1 + c_5 \cdot ((i \cdot h)^2 - 1)$$

dann  $\Delta w_i = -2c_5 \cdot h^2 \geq \Delta z_i$  für geeignetes  $c_5$  und  $w_i|_{\partial\Omega} \geq z_i|_{\partial\Omega}$  durch entsprechende Wahl von  $\alpha$ . Maximumprinzip gibt  $w_i \geq z_i$ . Speziell gilt  $z_0 \leq w_0 = \alpha \cdot c_1 - c_3$ . Widerspruch zu  $|z_0| \rightarrow 0$  für  $h \rightarrow 0$ .

Numerisches Experiment ergab:

- $\beta = 2$ : Konvergenz gegen falsche Lösung ( $\neq 0$ )
- $\beta > 2$ : Konvergenz gegen Null

Beweis offen

# 6

## Nichtkonforme Aspekte

- Ausgangsproblem: Variationsgleichung

$$a(u, v) = f(v) \quad (v \in V)$$

- Konforme FEM:

$$a(u_h, v_h) = f(v_h) \quad (v_h \in V_h \subseteq V)$$

Bei „richtiger“ nicht-konformer FEM:  $V_h \not\subseteq V$ . Andere mögliche Modifikation von konformen FEM:

$$a_h(u_h, v_h) = f_h(v_h) \quad (v_h \in V_h \subseteq V)$$

**Beispiel** Sei  $\Omega$  polygonal. Betrachte Randwertaufgabe

$$-\Delta u = f \qquad u|_{\partial\Omega} = 0 \qquad (6.1)$$

Nutze Dreieckszerlegung mit linearen finiten Elementen, die Freiheitsgrade seien die Funktionswerte in den Seitenmittelpunkten (Crouzeix-Raviart-Element). Dann gilt  $V_h \not\subseteq H^1$ .

Warum nutzt man dieses Element?

- Oft im Zusammenhang mit Stokes/Navier-Stokes.
- Basisfunktionen sind  $L^2$ -orthogonal (dazu: Quadraturformel mit diesen Knoten ist exakt für  $P_2$ , wegen  $\varphi_i \cdot \varphi_j \in P_2$  folgt damit Behauptung).

Schwache Formulierung von (6.1) ist

$$(\nabla u, \nabla v) = (f, v)$$

Diskretes Problem ist zum Beispiel (beachte  $V_h \not\subseteq H^1$ , damit globale Integrale nicht mehr definiert):

$$\sum_K \int_K \nabla u_h \cdot \nabla v_h = (f, v_h)$$

Benötigen neue Theorie für nichtkonforme Methoden.

**Beispiel** (Bekannte nichtkonforme Elemente) (i).  $Q_1^{\text{rot}}$  (Rannacher, Turek): Basisfunktionen  $\{1, x, y, x^2 - y^2\}$ , Freiheitsgrade sind Integralmittel über Kanten

(ii). Wilson-Element:  $P_2$ -Elemente mit den Freiheitsgraden

- Funktionswerte in den Eckpunkten
- $\frac{h_1}{h_2} \int_K p_{xx}$  und  $\frac{h_2}{h_1} \int_K p_{yy}$

wobei  $h_1, h_2$  die Kantenlängen des Rechtecks bezeichnen.

(iii). Versuch einer Verallgemeinerung des Crouzeix-Raviart-Elements: Nutze Funktionswerte in Gaußpunkten als Freiheitsgrade und  $P_2$ -Elemente. Ist nicht unisolvent!



Abbildung 6.1: Links: Versuchte Verallgemeinerung, rechts: Morley-Element

- (iv). Morley-Element:  $P_2$ -Elemente mit Freiheitsgraden als Funktionswerte in Ecken und Normalableitungen in den Seitenmittelpunkten. Offenbar  $V_h \not\subseteq H^2$ , sogar  $V_h \not\subseteq H^1$ . Wird insbesondere für Probleme vierter Ordnung benutzt.

**6.1 Satz** (2. Lemma von Strang)

Betrachte diskretes Problem

$$a_h(u_h, v_h) = f_h(v_h) \quad (v_h \in V_h \not\subseteq V) \quad (6.2)$$

Sei  $Z_h := V + V_h := \{z; z = v + v_h, v \in V, v_h \in V_h\}$ . Weiterhin gelte:

- (i).  $a_h$  sei eine auf  $Z_h$  definierte symmetrische Bilinearform.
- (ii).  $\|z_h\|_h := \sqrt{a_h(z_h, z_h)}$  sei eine Norm auf  $V_h$ .
- (iii).  $|a_h(z_h, w_h)| \leq M_0 \cdot \|z_h\|_h \cdot \|w_h\|_h$  für  $w_h, z_h \in Z_h$ .

Dann gilt

$$\|u - u_h\|_h \leq \underbrace{M \cdot \inf_{z_h \in V_h} \|u - z_h\|_h}_{\text{Approximationsfehler}} + \underbrace{\sup_{w_h \in V_h} \frac{|f_h(w_h) - a_h(u, w_h)|}{\|w_h\|_h}}_{\text{Konsistenzfehler}}$$

Beweis: Sei  $u_h$  die Lösung von (6.2). Sei  $z_h \in V_h$  beliebig. Dann

$$\begin{aligned} \|u_h - z_h\|_h^2 &= a(u_h - z_h, \underbrace{u_h - z_h}_{=: w_h}) \stackrel{(6.2)}{=} a_h(u - z_h, w_h) + f_h(w_h) - a_h(u, w_h) \\ &\stackrel{(iii)}{\leq} M_0 \cdot \|u - z_h\|_h \cdot \|w_h\|_h + f_h(w_h) - a_h(u, w_h) \end{aligned}$$

Damit folgt

$$\|u - u_h\|_h \leq \|u - z_h\|_h + \|u_h - z_h\|_h \leq M \cdot \|u - z_h\|_h + \frac{f_h(w_h) - a_h(u, w_h)}{\|w_h\|_h} \quad \square$$

**Beispiel** (Poisson-Gleichung mit Crouzeix-Raviart-Element) Sei  $\Omega$  polygonal und konvex. Betrachte die Randwertaufgabe

$$-\Delta u = f \quad u|_{\partial\Omega} = 0$$

Diskretisierung mit Crouzeix-Raviart-Elementen gibt FEM-Raum  $V_h$  (dabei: liegt Seite auf Rand, so setze Funktionswert in Seitenmittelpunkt gleich Null). Diskretes Problem:

$$a_h(u_h, v_h) := \sum_K (\nabla u_h, \nabla v_h)_K = (f, v_h)$$

Damit

$$\|v_h\|_h^2 := \sum_K (\nabla v_h, \nabla v_h)_K$$

Ist dies eine Norm auf  $V_h$ ?

Beweis: Zu untersuchen: Folgt aus  $\|v_h\|_h = 0$ , dass  $v_h = 0$ ? (Andere Norm-Eigenschaften folgen leicht.) Aus  $\|v_h\|_h = 0$  folgt, dass  $v_h$  stückweise konstant. Da je zwei Dreiecke eine gemeinsame Seite haben (und auf dem Seitenmittelpunkt der Funktionswert übereinstimmt), folgt, dass  $v_h$  global konstant ist. Damit  $v_h = 0$  wegen  $v_h = 0$  auf dem Rand.  $\square$

Damit folgt: Theorem 6.1 ist anwendbar.

- (i). Approximationsfehler: Da der konforme FEM-Raum ein Teilraum des nicht-konformen FEM-Raum ist, gilt

$$\inf_{z_h \in V_h} \|u - z_h\|_h \leq C \cdot h \cdot |u|_2$$

- (ii). Konsistenzfehler: Mittels partieller Integration folgt

$$\begin{aligned} F &:= f_h(w_h) - a_h(u, w_h) = \sum_K ((f, w_h)_K - (\nabla u, \nabla w_h)_K) \\ &= \sum_K \int_{\partial K} \frac{\partial u}{\partial n} \cdot w_h \end{aligned}$$

Seien  $e_{ij}$  die Kanten von  $K_i$ , dann

$$\begin{aligned} F &= \sum_i \sum_{j=1}^3 \int_{e_{ij}} \frac{\partial u}{\partial n} \cdot w_h \\ &= \sum_i \sum_{j=1}^3 \int_{e_{ij}} \frac{\partial u}{\partial n} \cdot (w_h - \bar{w}_h) \end{aligned}$$

wobei  $\bar{w}_h$  Konstante auf  $e_{ij}$  und gleich Null auf Randkanten (jede Nichtrandkante tritt 2 Mal auf...)

$$F = \sum_i \sum_{j=1}^3 \int_{e_{ij}} \left( \frac{\partial u}{\partial n} - \alpha \right) \cdot (w_h - \bar{w}_h)$$

mit  $\alpha$  Konstante. Funktioniert, wenn  $\int_{e_{ij}} w_h = \int_{e_{ij}} \bar{w}_h$  (damit Wahl von  $\bar{w}_h$  festgelegt, ist wohldefiniert). Wähle  $\alpha := \frac{\partial u^I}{\partial n}$ .

Abschätzung von  $\left\| \frac{\partial u}{\partial n} - \frac{\partial u^I}{\partial n} \right\|$  notwendig (analog dann  $\|w_h - \bar{w}_h\|_{0,e}$ ). Hilfsmittel: multiplikative Spur-Ungleichung.

$$\|v\|_{0,e}^2 \leq c \cdot \left( \|v\|_{0,K} \cdot |v|_{1,K} + \frac{1}{\text{diam } K} \cdot \|v\|_{0,K}^2 \right) \quad (6.3)$$

Damit erhält man

$$\begin{aligned} \left\| \frac{\partial u}{\partial n} - \frac{\partial u^I}{\partial n} \right\|_{0,e} &\leq c \cdot \left( \left\| \frac{\partial u}{\partial n} - \frac{\partial u^I}{\partial n} \right\|_{0,K} \cdot |u|_{2,K} + \frac{1}{\text{diam } K} \cdot \left\| \frac{\partial u}{\partial n} - \frac{\partial u^I}{\partial n} \right\|_{0,K}^2 \right) \\ &\leq c \cdot \text{diam } K \cdot |u|_{2,K}^2 \end{aligned}$$

Somit gilt:

$$\|u - u_h\|_h \leq C \cdot h \cdot |u|_2$$

## 6.2 Satz (1. Lemma von Strang)

Betrachte diskretes Problem

$$a_h(u_h, v_h) = f_h(v_h) \quad (v_h \in V_h \subseteq V) \quad (6.4)$$

(d.h.  $a, f$  werden modifiziert, aber  $V_h \subseteq V$ ).  $a_h$  sei  $V_h$ -elliptisch und beschränkt. Dann

$$\|u - u_h\| \leq M \cdot \inf_{v_h \in V_h} \|u - v_h\| + \sup_{w_h \in V_h} \left( \frac{|a(v_h, w_h) - a_h(v_h, w_h)|}{\|u - u_h\|} + \frac{|f_h(w_h) - f(w_h)|}{\|u - u_h\|} \right)$$

Beweis: Es sei  $u_h$  die Lösung von (6.4). Sei  $v_h \in V_h, w_h := u_h - v_h$ . Dann folgt aus der  $V_h$ -Elliptizität:

$$\begin{aligned} \alpha \cdot \|u_h - v_h\|^2 &\leq a_h(u_h - v_h, w_h) \\ &= a(u - v_h, w_h) + a(v_h, w_h) - a_h(v_h, w_h) + f_h(w_h) - f(w_h) \end{aligned}$$

Damit:

$$\alpha \cdot \|u_h - v_h\|^2 \leq M \cdot \|u - v_h\| \cdot \|w_h\| + a(v_h, w_h) - a_h(v_h, w_h) + f_h(w_h) - f(w_h)$$

Aus der Dreiecksungleichung folgt die Behauptung.  $\square$

**Bemerkung** Warum betrachtet man diskrete Probleme der Form (6.4)? Standard-Anwendungen:

- Quadraturformel zur Berechnung der Integrale
- Stabilisierungsmethoden (z.B. Stromliniendiffusion)
- Analysis von Finite-Volumen-Methoden

**Beispiel** (Stromliniendiffusion, 1978)

$$-\Delta u + b \cdot \nabla u + c \cdot u = f \qquad u|_{\partial\Omega} = 0$$

Bei dominanter Konvektion ist das Galerkin-Verfahren wenig geeignet. Schwache Formulierung bei Galerkin-Verfahren ist gegeben durch

$$a_G(u, v) := (\nabla u, \nabla v) + (b \cdot \nabla u + c \cdot u, v) = (f, v) \qquad (v \in H_0^1(\Omega))$$

Stromliniendiffusion:

$$\begin{aligned} a_h(u_h, v_h) &:= a_G(u_h, v_h) + \sum_K (-\Delta u_h + b \cdot \nabla u_h + c \cdot u_h, b \cdot \nabla v_h) \cdot \delta_K \\ f_h(v_h) &:= f(v_h) + \sum_K (v, b \cdot \nabla v_h) \cdot \delta_K \end{aligned}$$

weiterhin gelte  $V_h \subseteq V$ . Vorteil dieser Formulierung:

$$a_h(u, v_h) = (f_h, v_h) \qquad (v_h \in V_h)$$

Ist  $a_h$   $V_h$ -elliptisch? Ja, aber nicht trivial.

# 7

## Weitere Verfahren

### 7.1 Finite-Volumen-Methode

**Bemerkung** (Zur Geschichte) • Samarskij 1977 („Integralbilanzmethode“)

- Heinrich 1987 („Finite difference methods on irregular networks“)
- Bank, Rose 1987 („Box-Methode“)

**Beispiel** (1D)

$$\frac{d}{dx} \left( k(x) \cdot \frac{d}{dx} u(x) \right) + q(x) \cdot u(x) = f(x) \quad u(0) = u(1) = 0$$

Integriere Gleichung über „Kontrollvolumen“  $(x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}})$ :

$$\left[ -k(x) \cdot \frac{d}{dx} u(x) \right]_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} + \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} q(x) \cdot u(x) dx = \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f(x) dx \quad (7.1)$$

(„Erhaltungsgleichung“). Die Gleichung (7.1) ist Ausgangspunkt für Finite-Volumen-Methode. Annahme:  $u|_{[x_{i-\frac{1}{2}}, x_{i+\frac{1}{2}}]} = u_i$  (d.h.  $u$  wird jetzt auf dem Kontrollvolumen stückweise konstant approximiert). Approximation von (7.1):

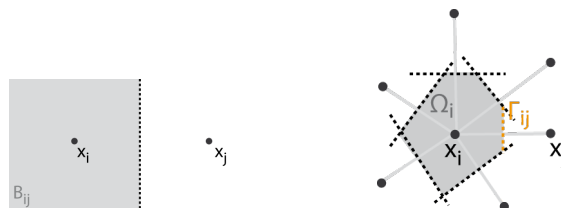
$$-k_{i+\frac{1}{2}} \cdot \frac{u_{i+1} - u_i}{h_{i+1}} + k_{i-\frac{1}{2}} \cdot \frac{u_i - u_{i-1}}{h_i} + u_i \cdot \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} q(x) dx = \int_{x_{i-\frac{1}{2}}}^{x_{i+\frac{1}{2}}} f(x) dx$$

(Entspricht einem bekannten Differenzenverfahren.)

**Definition** (Voronoi-Boxen 2D) Oft benutzt man als Kontrollbereiche sogenannte Voronoi-Boxen: Gegeben seien Punkte  $x_i$  in der Ebene. Sei

$$B_{ij} := \{x \in \mathbb{R}^2; |x - x_i| < |x - x_j|\}$$

Dann heißt  $\Omega_i := \bigcap_{j \neq i} B_{ij}$  Voronoi-Box zu  $x_i$ .



**Bemerkung** (i). Für  $\Omega \subseteq \mathbb{R}^2$  polygonal seien  $x_i \in \bar{\Omega}$  geeignet fixiert. Definiert Zerlegung in Voronoi-Boxen.

(ii). Notation:  $j \in \Lambda_i \Leftrightarrow x_j$  definiert ein Randstück von  $\Omega_i$ . Dieses Randstück wird mit  $\Gamma_{ij}$  bezeichnet.

**Beispiel** Sei  $\Omega \subseteq \mathbb{R}^2$  polygonal. Betrachte Randwertaufgabe

$$-\Delta u + c \cdot u = f \quad u|_{\partial\Omega} = 0$$

wobei  $c \geq 0$ . Es seien (hinreichend viele)  $x_i \in \bar{\Omega}$  gegeben, die eine Zerlegung von  $\Omega$  in Voronoi-Boxen  $\Omega_i$  definieren. (Die Voronoi-Boxen entsprechen nun den „Kontrollbereichen“.)

$$-\int_{\partial\Omega_i} \frac{\partial u}{\partial n} + \int_{\Omega_i} c \cdot u = \int_{\Omega_i} f$$

(„Bilanzgleichung“). Approximiere  $u$  stückweise konstant, d.h.  $u|_{\Omega_i} = u_i$ . Damit

$$-\sum_{j \in \Lambda_i} m_{ij} \cdot \frac{u_j - u_i}{\ell_{ij}} + u_i \cdot \int_{\Omega_i} c = \int_{\Omega_i} f$$

wobei

$$|\Gamma_{ij}| =: m_{ij} \quad |x_i - x_j| =: \ell_{ij}$$

wobei

$$|\Gamma_{ij}| =: m_{ij} \quad |x_i - x_j| =: \ell_{ij}$$

Für  $m_i := |\Omega_i|$  erhält man

$$-\sum_{j \in \Lambda_i} m_{ij} \cdot \frac{u_j - u_i}{\ell_{ij}} + u_i \cdot m_i \cdot c_i = f_i \cdot c_i \quad (7.2)$$

**Bemerkung** (Primale/duale Vernetzung)

Gegeben sei ein Dreiecksgitter (einer FEM), genannt primales Gitter. Dann heißt das Gitter, das durch die entsprechenden Voronoi-Boxen gegeben ist, das duale Gitter. (Zum Teil werden auch andere Boxen genutzt, z.B. Donald-Boxen.)



Welche Beziehungen gibt es zwischen linearen finiten Elementen auf primalem Gitter und der Finite-Volumen-Methode auf dem dualen Gitter? Aus dem letzten Beispiel folgt, dass für  $-\Delta u$ : lineare finite Elemente = Finite-Volumen-Methode mit Voronoi-Boxen.

### 7.1 Lemma

Für eine gegebene Funktion  $\xi$  bezeichne  $\bar{\xi}$  die stückweise konstante Funktion auf  $\Omega_i$  mit  $\bar{\xi}|_{\Omega_i} = \xi(x_i)$ . Es gilt

$$\|w_h - \bar{w}_h\| \leq c \cdot h \cdot |w_h|_1 \quad (v_h \in V_h)$$

Beweis: Es gilt

$$(w_h - \bar{w}_h)(x) = \nabla w_h \cdot (x - x_i) \\ \Rightarrow |w_h - \bar{w}_h| \leq c \cdot h \cdot |\nabla w_h|$$

**Beispiel** (Interpretation der FVM als lineare FEM) Sei  $V_h$  der Raum der linearen finiten Elemente über der primalen Triangulation,  $u_h(x_i) := u_i$ . Summiere über  $i$  in (7.2):

$$\underbrace{\sum_i v_i \cdot \left( -\sum_{j \in \Lambda_i} m_{ij} \cdot \frac{u_j - u_i}{\ell_{ij}} + m_i \cdot c_i \cdot u_i \right)}_{=: a_h(u_h, v_h)} = \underbrace{\sum_i m_i \cdot f_i \cdot v_i}_{=: f_h(v_h)}$$

(Diese Gleichung ist zu (7.2) äquivalent, betrachte dazu  $v_i := \varphi_i$  wobei  $(\varphi_i)_i$  die nodale Basis bezeichnet.)

Fehleranalyse der Finite-Volumen-Methode mit Hilfe des Strang-Lemmas 6.2:



(i). Ist  $a_h$   $V_h$ -elliptisch? Es gilt

$$a_h(v_h, v_h) = \sum_i v_i \cdot \left( - \sum_{j \in \Lambda_i} m_{ij} \cdot \frac{v_j - v_i}{\ell_{ij}} \right) + \underbrace{\sum_i m_i \cdot c_i \cdot v_i^2}_{\geq 0} \geq (\nabla v_h, \nabla v_h)$$

(Beachte dazu, dass nach obiger Bemerkung die FVM-Diskretisierung von  $-\Delta u$  mit den linearen finiten Elementen übereinstimmt.)

(ii). Abschätzung des Konsistenzfehlers: Beachte dazu, dass

$$\sum_i v_i \cdot \left( - \sum_{j \in \Lambda_i} m_{ij} \cdot \frac{u_j - u_i}{\ell_{ij}} \right) = (\nabla u_h, \nabla v_h)$$

(FVM = lineare FEM für  $-\Delta u$ ).

$$\begin{aligned} |(f, w_h) - (\bar{f}, \bar{w}_h)| &\leq |(f, w_h - \bar{w}_h)| + |(f - \bar{f}, \bar{w}_h)| \\ &\stackrel{7.1}{\leq} c \cdot h \cdot |w|_1 + c \cdot \underbrace{\|\bar{w}_h\|_0}_{\leq \|w_h\|_0 + \|w_h - \bar{w}_h\|_0} \\ &\leq c \cdot h \cdot \|w_h\|_1 \end{aligned}$$

Analoge Abschätzung für  $|(c \cdot v_h, w_h) - (\bar{c} \cdot \bar{v}_h, \bar{w}_h)|$

Damit folgt unter Anwendung von Theorem 6.2:

$$\|u - u_h\|_1 \leq c \cdot h \cdot |u|_2$$

(unter denselben Voraussetzungen wie bei linearen finiten Elementen).

**Bemerkung** (Konvektions-Diffusions-Gleichungen)

$$-\Delta u + b \cdot \nabla u + c \cdot u = f \quad u|_{\partial\Omega} = 0$$

Idee: Diskretisiere  $-\Delta u$  mit linearen finiten Elementen (gibt M-Matrix) und diskretisiere die Terme  $b \cdot \nabla u + c \cdot u$  mit Finite-Volumen-Methode. Warum ist es sinnvoll den Konvektionsterm mit FVM zu behandeln?

$$\begin{aligned} (b \cdot \nabla u, v) &= (\operatorname{div}(b \cdot u), v) - (\operatorname{div} b, u \cdot v) \\ &\approx (\operatorname{div}(b \cdot u), \bar{v}) - (\operatorname{div} b, \bar{u} \cdot \bar{v}) \\ &= \sum_i \sum_{j \in \Lambda_i} v_i \cdot \int_{\Gamma_{ij}} u(b \cdot n_{ij}) - \sum_i \sum_{j \in \Lambda_i} u_i \cdot v_i \cdot \int_{\Gamma_{ij}} b \cdot n_{ij} \end{aligned}$$

wobei  $n_{ij}$  der äußere Normalenvektor auf  $\Gamma_{ij}$ . Weitere Approximationen:

- Es sei  $\int_{\Gamma_{ij}} b_i \cdot n_{ij} \approx \beta_{ij}$ , zum Beispiel  $\beta_{ij} := |\Gamma_{ij}| \cdot b(Q_{ij}) \cdot n_{ij}$  wobei  $Q_{ij}$  der Seitenmittelpunkt von  $\Gamma_{ij}$ .
- $u \approx \lambda_{ij} \cdot u_i + (1 - \lambda_{ij}) \cdot u_j$ . Upwind-Strategie:

$$\lambda_{ij} := \begin{cases} 1 & \beta_{ij} > 0 \\ 0 & \beta_{ij} < 0 \end{cases}$$

(Damit: Matrix mit positiven Hauptdiagonalelementen und nichtpositiven Nebendiagonalelementen. Mit FEM nicht so einfach erreichbar.)

Somit:

$$b_h(u_h, v_h) := \sum_i \sum_{j \in \Lambda_i} \beta_{ij} \cdot (1 - \lambda_{ij}) \cdot (u_j - u_i) \cdot v_i$$

## 7.2 Discontinuous Galerkin-Verfahren für elliptische Probleme

**Bemerkung** (Zur Geschichte) • 1973 Discontinuous Galerkin-Verfahren für Transportprobleme  
 $(b \cdot \nabla u + c \cdot u = f)$

• 1998 Discontinuous Galerkin-Verfahren für elliptische Probleme (Oden, Baumann, Babuska)  
 Es gibt zwei verschiedene Zugänge zum Discontinuous Galerkin-Verfahren:

- (i). primal (Startpunkt: elliptisches Problem)
- (ii). dual (Ausgang: gemische Formulierung)

Hier Beschränkung auf primalen Zugang. Die beiden Zugänge sind äquivalent.

**Bemerkung** Betrachte im Folgenden das Problem

$$-\Delta u + c \cdot u = f \qquad u|_{\partial\Omega} = 0 \qquad (7.3)$$

Bezeichnungen:

- (i).  $\bar{\Omega} = \bigcup_{K \in \mathcal{T}} K$  Zerlegung in polygonale Elemente
- (ii).  $\mathcal{E}$  sei die Menge aller Ränder der Elemente von  $\mathcal{T}$  und  $\mathcal{E}_{\text{int}}$  die Menge aller Ränder, die nicht auf dem Rand liegen.
- (iii). Sei  $v$  eine stückweise glatte Funktion und  $e \in \mathcal{E}_{\text{int}}$ , dann definiere

$$[v]_e := v|_{\partial K \cap \Omega} - v|_{\partial K' \cap \Omega}$$

$$\langle v \rangle_e := \frac{1}{2} \cdot (v|_{\partial K \cap \Omega} + v|_{\partial K' \cap \Omega})$$

für  $K, K' \in \mathcal{T}$  („Sprung“ und „Mittelwert“; offenbar abhängig von der konkret gewählten Nummerierung).

(Herleitung der schwachen Formulierung)

Sei  $v$  eine stückweise glatte Funktion. Aus (7.3):

$$\sum_{K \in \mathcal{T}} \int_K (-\Delta u) \cdot v + \sum_{K \in \mathcal{T}} \int_K c \cdot u \cdot v = \sum_{K \in \mathcal{T}} \int_K f \cdot v$$

Partielle Integration gibt:

$$\sum_{K \in \mathcal{T}} \int_K (-\Delta u) \cdot v = \sum_K \int_K \nabla u \cdot \nabla v - \int_{\Gamma} (\nabla u) \cdot \mu \cdot v - \sum_K \int_{\partial K} (\nabla u) \cdot \nu \cdot v$$

wobei  $\Gamma := \partial\Omega$  und  $\mu, \nu$  die entsprechenden Normalenvektoren bezeichnen. Nutze für die Umformung des letzten Summanden folgende Gleichheit

$$a_+ \cdot b_+ - a_- \cdot b_- = \frac{a_+ + a_-}{2} \cdot (b_+ - b_-) + (a_+ - a_-) \cdot \frac{b_+ + b_-}{2}$$

Damit ergibt sich insgesamt

$$\sum_K (\nabla u, \nabla v)_K - \int_{\Gamma} (\nabla u) \cdot \mu \cdot v - \sum_{e \in \mathcal{E}_{\text{int}}} \int_e \langle (\nabla u) \cdot \nu \rangle_e \cdot [v]_e + (c \cdot u, v) = (f, v)$$

Addiere auf der linken Seite die Terme

$$\pm \int_{\Gamma} (\nabla v) \cdot \mu \cdot u \pm \sum_{e \in \mathcal{E}_{\text{int}}} \int_e \langle (\nabla v) \cdot \nu \rangle_e \cdot [u]_e$$

mit dem Ziel (im Fall  $-$ ) die entstehende Bilinearform zu symmetrisieren (beide sind 0, beachte  $u|_{\partial\Omega} = 0$  und Stetigkeit von  $u$ ). Addiere außerdem die Terme

$$\int_{\Gamma} \sigma \cdot u \cdot v + \int_{\mathcal{E}_{\text{int}}} \sigma \cdot [u] \cdot [v]$$

für geeignetes  $\sigma \geq 0$  (Strafterme, die  $V_h$ -Elliptizität sichern sollen). Definiere linke Seite als  $B_{\pm}(u, v)$ .  $V_h$  sei der FE-Raum der stückweise glatten (i.a. polynomialen) Funktionen. Damit diskretes Problem

$$B_{\pm}(u_h, v_h) = (f, v_h) \quad (v_h \in V_h) \quad (7.4)$$

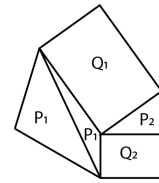
Damit Variante des Discontinuous-Galerkin erhalten:

- Vorzeichen  $-$ : symmetric interior penalties (SIP)
- Vorzeichen  $+$ : nonsymmetric interior penalties (NIP)

Ermöglicht auch „schwache“ Einbeziehung von Randbedingungen.

**Bemerkung** (Vorteile des Discontinuous Galerkin-Verfahren)

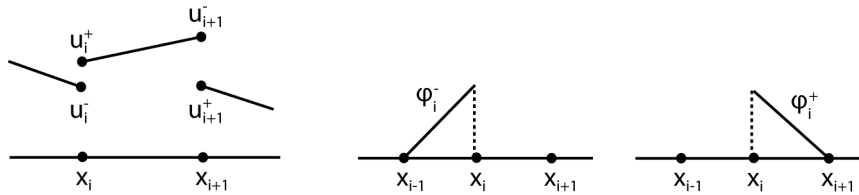
- (i). Flexibilität
- (ii). Gitter müssen nicht notwendig zulässig sein.



**Beispiel** Sei  $\Omega = (0, 1)$ . Betrachte Randwertaufgabe

$$-u'' = f \quad u(0) = u(1) = 0$$

Verwende lineare Elemente auf äquidistantem Gitter, setze  $\sigma := \frac{1}{h}$ .



Es ergibt sich das Gleichungssystem

$$\begin{aligned} \frac{1}{h} \cdot \left( \frac{1}{2}u_{i-1}^+ - u_i^- + 2u_i^+ - u_{i+1}^- - \frac{1}{2}u_{i+1}^+ \right) &= \int_{x_{i-1}}^{x_i} f \cdot \varphi_i^- \\ \frac{1}{h} \cdot \left( -\frac{1}{2}u_{i-1}^- - u_i^+ + 2u_i^- - u_{i+1}^+ + \frac{1}{2}u_{i+1}^- \right) &= \int_{x_i}^{x_{i+1}} f \cdot \varphi_i^+ \end{aligned}$$

Für  $u_i^+ = u_i^-$  ergibt sich insbesondere

$$\frac{1}{h} \cdot (-u_{i-1} + 2u_i - u_{i+1}) = \dots$$

d.h. das stetige Galerkin-Verfahren.

**Bemerkung** Im Folgenden bezeichne

$$\|v\|_{dG}^2 := \sum_{K \in \mathcal{T}} (|v|_{1,K}^2 + \|v\|_{0,K}^2) + \int_{\Gamma} \sigma \cdot v^2 + \int_{\mathcal{E}_{\text{int}}} \sigma \cdot [v]^2$$

(ist abhängig von Zerlegung!). Ist  $B_{\pm}$  bzgl. dieser Norm  $V_h$ -elliptisch?

**7.2 Lemma** (i).  $B_+$  ist  $V_h$ -elliptisch für alle  $\sigma \geq 0$ .

(ii).  $B_-$  ist für  $\sigma := \frac{\sigma_0}{h}$  mit  $\sigma_0$  hinreichend groß  $V_h$ -elliptisch.

Beweis: (i). Es gilt

$$B_+(v_h, v_h) = \|v_h\|_{dG}^2$$

d.h.  $V_h$ -Elliptizität für alle  $\sigma \geq 0$ .

(ii). Nach Definition gilt

$$\begin{aligned} B_-(v_h, v_h) &= \sum_K |v_h|_{1,K}^2 + (c \cdot v_h, v_h) + \int_{\Gamma} \sigma \cdot v_h^2 + \int_{\mathcal{E}_{\text{int}}} \sigma \cdot [v_h]^2 \\ &\quad - 2 \int_{\Gamma} (\nabla v_h) \cdot \mu \cdot v_h - 2 \sum_{e \in \mathcal{E}_{\text{int}}} \int_e \langle (\nabla v_h) \cdot \nu \rangle_e \cdot [v_h]_e \\ &\stackrel{!}{\geq} C \cdot \|v_h\|_{dG}^2 \end{aligned}$$

Schätze dazu die Terme in der zweiten Zeile nach oben ab. Hier nur Abschätzung des zweiten Terms:

$$\begin{aligned} \int_e \langle (\nabla v_h) \cdot \nu \rangle_e \cdot [v_h] &= \int_e \frac{1}{\sigma^{\frac{1}{2}}} \cdot \langle (\nabla v_h) \cdot \nu \rangle_e \cdot \sigma^{\frac{1}{2}} \cdot [v_h]_e \\ &\stackrel{\text{CSU}}{\leq} \left( \int_e \frac{1}{\sigma} \cdot \langle (\nabla v_h) \cdot \nu \rangle^2 \right)^{\frac{1}{2}} \cdot \left( \int_e \sigma \cdot [v_h]^2 \right)^{\frac{1}{2}} \\ &\leq \frac{\gamma}{2} \int_e \frac{1}{\sigma} \cdot \langle (\nabla v_h) \cdot \nu \rangle^2 + \frac{1}{2\gamma} \int_e \sigma \cdot [v_h]^2 \\ &\stackrel{*}{\leq} \frac{c \cdot \gamma}{2h \cdot \sigma} \cdot \int_K (\nabla v_h)^2 + \frac{1}{2\gamma} \int_e \sigma \cdot [v_h]^2 \end{aligned}$$

Dabei in (\*) genutzt:

$$\int_e \langle (\nabla v_h) \cdot \nu \rangle^2 \leq \frac{c}{h} \cdot \int_K (\nabla v_h)^2$$

(folgt aus multiplikativer Spurgleichung (6.3)). Wähle nun  $\sigma := \sigma_0 \cdot \frac{1}{h}$ , dann

$$-\frac{c \cdot \gamma}{2h \cdot \sigma} = -\frac{c \cdot \gamma}{2\sigma_0}$$

Der Term wird schön für  $\sigma_0$  hinreichend groß. □

### 7.3 Satz (Fehlerabschätzung für Discontinuous Galerkin)

Gegeben sei eine quasi-uniforme  $\Delta$ -Zerlegung. Verwende (unstetige) lineare Elemente. Weiterhin sei  $\sigma := \frac{\sigma_0}{h}$  mit  $\sigma_0$  hinreichend groß. Für  $u \in H^2$  gilt in (7.3)

$$\|u - u_h\|_{dG} \leq c \cdot h \cdot |u|_2$$

Beweis:  $\pi u \in V_h$  sei die stückweise  $L^2$ -Projektion von  $u$  in  $V_h$  (d.h.  $(\pi u, v)_K = (u, v)_K$  für  $v \in V_h$ ) und  $u^I$  die stetige lineare Interpolierende. Abschätzungen für  $u - \Pi u$  verlaufen analog zu Abschätzungen für  $u - u^I$ . Es sei

$$\eta := u - \pi u \qquad \xi := \pi u - u_h \quad (\in V_h)$$

Es gilt

$$\begin{aligned} \alpha \cdot \|\xi\|_{dG}^2 &\leq B(\xi, \xi) = B(\pi u - u_h, \xi) = B(\pi u - u, \xi) + \underbrace{B(u - u_h, \xi)}_0 \\ &\Rightarrow \alpha \cdot \|\xi\|_{dG}^2 \leq -B(\eta, \xi) \end{aligned}$$

wegen  $B(u - u_h, v_h) = f(v_h)$ . Ziel: Schätze die Beträge der 8 Summanden in  $B(\eta, \xi)$  durch  $g(\eta) \cdot \|\xi\|_{dG}$  ab. Die folgenden Abschätzungen werden als bekannt vorausgesetzt:

$$\begin{aligned} |\eta|_{1,K} &\leq c \cdot h \cdot |u|_{2,K} \\ \|\eta\|_{0,K} &\leq c \cdot h^2 \cdot |u|_{2,K} \\ \|\eta\|_{0,e} &\leq c \cdot h^{\frac{3}{2}} \cdot |u|_{2,K} \\ |\eta|_{1,e} &\leq c \cdot h^{\frac{1}{2}} \cdot |u|_{2,K} \end{aligned}$$

(Die dritte Ungleichung folgt aus den ersten beiden mit Spurgleichung (6.3)).

(i). Bei 4 Summanden: Cauchy-Schwarz.

$$\sum_K (\nabla \eta, \nabla \xi)_K \quad (\eta, \xi) \quad \int_{\Gamma} \sigma \cdot \eta \cdot \xi \quad \int_{\mathcal{E}_{\text{int}}} \sigma \cdot [\eta] \cdot [\xi]$$

(Für die letzten beiden Integrale nutze  $\sigma = \sigma^{\frac{1}{2}} \cdot \sigma^{\frac{1}{2}}$  und beachte  $\sigma = \frac{\sigma_0}{h}$ .)

(ii). Es gilt

$$\begin{aligned} \left| \int_e (\nabla \eta \cdot \nu) \cdot \xi \right| &\leq \left| \int_e (\nabla \eta \cdot \nu) \cdot \xi \cdot \frac{\sigma^{\frac{1}{2}}}{\sigma^{\frac{1}{2}}} \right| \\ &\stackrel{\text{CSU}}{\leq} |\eta|_{1,e} \cdot \|\xi\|_{dG} \leq c \cdot h^{\frac{1}{2}} \cdot h^{\frac{1}{2}} \cdot \|\xi\|_{dG} \end{aligned}$$

Analoge Abschätzung auf „innerem Rand“.

(iii). Es gilt

$$\begin{aligned} \left| \int_e (\nabla \xi \cdot \nu) \cdot \eta \right| &\stackrel{\text{CSU}}{\leq} \left( \frac{1}{h} \cdot \int_e \eta^2 \right)^{\frac{1}{2}} \cdot \left( h \cdot \int_e (\nabla \xi \cdot \nu)^2 \right)^{\frac{1}{2}} \\ &\leq c \cdot \left( \frac{1}{h} \cdot \int_e \eta^2 \right)^{\frac{1}{2}} \cdot \left( \int_K (\nabla \xi)^2 \right)^{\frac{1}{2}} \\ &\leq c \cdot h \cdot \|\xi\|_{dG} \quad \square \end{aligned}$$

**Beispiel** (Reine Konvektionsprobleme)

$$b \cdot \nabla u + c \cdot u = f \quad u|_{\Gamma_-} = g \quad (7.5)$$

wobei  $\Gamma_- := \{x \in \Gamma := \partial\Omega; b \cdot n < 0\}$  „Einströmrand“.

1D (siehe auch Ern, Guermond: Theory and practice of finite elements): Betrachte das Problem

$$u' = f \text{ in } (0, 1) \quad u(0) = 0$$

Sei  $X := \{v \in H^1; v(0) = 0\}$ . Für  $u \in X, v \in L^2$  definiere  $a(u, v) := \int_0^1 u' \cdot v$ . Dann ist die schwache Formulierung gegeben durch

$$a(u, v) = (f, v) \quad (v \in W := L^2(0, 1))$$

gegeben. Problem: Lösungsraum  $\neq$  Raum der Testfunktionen, d.h. Verallgemeinerung des Lax-Milgram-Lemmas und des Cea-Lemmas notwendig. Für lineare finite Elemente erhält man dann

$$\|u - u_h\|_0 \leq C \cdot h \quad |u - u_h|_1 \leq C$$

Zudem oszilliert die numerische Lösung. Alternativen zu Galerkin:

- Stromliniendiffusion
- Discontinuous Galerkin (für lineare Elemente erhält man  $\|u - u_h\|_0 \leq c \cdot h^{\frac{3}{2}}$  und Oszillationen verschwinden)

# 8

## Parabolische Randwertaufgaben

**Beispiel** (Wärmeleitung)

$$u_t - \Delta u = f \text{ in } \Omega \times (0, T) \qquad u(0, \cdot) = u_0$$

mit Randbedingungen auf  $\partial\Omega$ . Schwache Formulierung ...?

Annahme: homogene Randbedingung  $u|_{\partial\Omega} = 0$ . Interpretation: Gesucht wird Abbildung  $t \mapsto u(t) \in V := H_0^1(\Omega)$ . Für einen normierten Raum  $(V, \|\cdot\|_V)$  sei

$$L^2(0, T; V) := \left\{ u : [0, T] \rightarrow V \text{ mb; } \int_0^T \|u(t)\|_V^2 < \infty \right\}$$

Nun versucht man eine schwache Formulierung anzugeben mit  $u \in L^2(0, T; V)$ ,  $u' \in L^2(0, T; V')$ . Nächstes Problem: Ist  $u(0)$  überhaupt definiert?

**Definition** Es gelte  $V \subseteq H \subseteq V'$ . Weiterhin sei  $V$  separabel, reflexiv,  $H$  ein separabler Hilbert-Raum und  $V$  sei dicht in  $H$  mit  $\|v\|_H \leq c \cdot \|v\|$ . Dann heißt  $(V, H, V')$  ein Evolutionstripel.

### 8.1 Lemma

Bilden  $(V, H, V')$  ein Evolutionstripel, dann gilt:

$$u \in L^2(0, T; V), u' \in L^2(0, T; V') \Rightarrow u : [0, T] \rightarrow H \text{ stetig}$$

Beweis: s. Zeidler: Nichtlineare Funktionalanalysis. □

### 8.2 Lemma

$$\langle u'(t), v \rangle = (u'(t), v)_H$$

**Beispiel** (Fortsetzung) Setze im letzten Beispiel  $V := H_0^1(\Omega)$ ,  $H := L^2(\Omega)$ ,  $V' = H^{-1}(\Omega)$ . Damit schwache Formulierung: Finde  $u \in L^2(0, T; V)$ ,  $u' \in L^2(0, T; V')$  mit

$$\forall v \in V : \underbrace{\langle u'(t), v \rangle}_{\stackrel{8.2}{=} (u'(t), v)} + a(u(t), v) = (f(t), v) \qquad u(0, \cdot) = u_0 \in H \qquad (8.1)$$

wobei  $a$  die „übliche“ Bilinearform zum elliptischen Operator. Man kann zeigen: (8.1) besitzt eine eindeutige Lösung.

### 8.3 Satz (A priori-Abschätzung)

Es sei  $(V, H, V')$  ein Evolutionstripel. Die Bilinearform  $a$  in (8.1) sei  $V$ -elliptisch. Dann gilt

$$\|u(t)\|_H^2 + \alpha \cdot \int_0^t \|u\|_V^2 \leq \|u(0)\|_H^2 + \frac{1}{\alpha} \cdot \int_0^t \|f(s)\|_H^2 ds$$

Beweis: Setze in (8.1)  $v := u(t)$ . Dann

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} (u(t), u(t)) + a(u(t), u(t)) &= (f(t), u(t)) \\ \Rightarrow \frac{1}{2} \frac{d}{dt} \|u(t)\|_H^2 + \alpha \cdot \|u\|_V^2 &\leq \frac{1}{\alpha} \cdot \frac{\|f(t)\|_H^2}{2} + \alpha \cdot \frac{\|u(t)\|_H^2}{2} \end{aligned}$$

Aus  $\|u\|_V \leq \|u\|_H$  folgt mittels Integration:

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} \int_0^t \|u(s)\|_H^2 ds + \frac{\alpha}{2} \cdot \int_0^t \|u\|_V^2 ds &\leq \frac{1}{2\alpha} \cdot \int_0^t \|f(s)\|_H^2 ds \\ \Rightarrow \|u(t)\|_H^2 + \alpha \cdot \int_0^t \|u\|_V^2 ds &\leq \|u(0)\|_H^2 + \frac{1}{\alpha} \cdot \int_0^t \|f(s)\|_H^2 ds \quad \square \end{aligned}$$

#### 8.4 Satz (A priori-Abschätzung)

Es sei  $(V, H, V')$  ein Evolutionstripel. Die Bilinearform  $a$  in (8.1) sei  $H$ -elliptisch. Dann gilt

$$\|u(t)\|_H^2 \leq e^{-\alpha t} \cdot \|u(0)\|_H^2 + \frac{1}{\alpha} \cdot e^{-\alpha t} \cdot \int_0^t e^{\alpha s} \cdot \|f(s)\|_H^2 ds$$

Beweis: Setze in (8.1)  $v := u(t)$ . Dann

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} (u(t), u(t)) + a(u(t), u(t)) &= (f(t), u(t)) \\ \Rightarrow \frac{1}{2} \frac{d}{dt} \|u(t)\|_H^2 + \alpha \cdot \|u(t)\|_H^2 &\leq \frac{1}{\alpha} \cdot \frac{\|f(t)\|_H^2}{2} + \alpha \cdot \frac{\|u(t)\|_H^2}{2} \\ \Rightarrow \frac{d}{dt} \|u(t)\|_H^2 + \alpha \|u(t)\|_H^2 &\leq \frac{1}{\alpha} \cdot \|f(t)\|_H^2 \\ \Rightarrow \|u(t)\|_H^2 &\leq e^{-\alpha t} \cdot \|u(0)\|_H^2 + \frac{1}{\alpha} \cdot e^{-\alpha t} \cdot \int_0^t e^{\alpha s} \cdot \|f(s)\|_H^2 ds \end{aligned}$$

Dabei benutzt:

$$\frac{d}{dt} (e^{\alpha t} \cdot \|u(t)\|_H^2) \leq \frac{1}{\alpha} \cdot e^{\alpha t} \cdot \|f\|_H^2 \quad \square$$

**Bemerkung** Theorem 8.4 sagt: Der Einfluss der Anfangsbedingungen wird exponentiell gedämpft.

(Semidiskretisierung: Elemente bzgl. der räumlichen Variablen)

Sei  $V_h \subseteq V$  ein FE-Raum. Gesucht ist  $u_h(t)$  mit

$$\forall v_h \in V_h : \left( \frac{d}{dt} u_h(t), v_h \right) + a(u_h(t), v_h) = (f(t), v_h) \quad (8.2)$$

Es sei  $\{\varphi_i\}_{i \in I}$  Basis von  $V_h$ , dann

$$u_h(t) = \sum_{i \in I} u_i(t) \cdot \varphi_i$$

Einsetzen in (8.2) gibt

$$\begin{aligned} \sum_{i \in I} (\varphi_i, \varphi_j) \frac{d}{dt} u_i + a \left( \sum_{i \in I} u_i(t) \cdot \varphi_i, \varphi_j \right) &= (f(t), \varphi_j) \\ \Leftrightarrow \sum_{i \in I} (\varphi_i, \varphi_j) \frac{d}{dt} u_i + \sum_{i \in I} a(\varphi_i, \varphi_j) \cdot u_i(t) &= (f(t), \varphi_j) \end{aligned}$$

Sei  $I := \{1, \dots, N\}$  und  $U(t) := (u_1(t), \dots, u_N(t))^T$ . Dann wird die letzte Gleichung zu

$$D_h U'(t) + A_h U(t) = F(t) \quad U(0) = U_0$$



(d.h. System von Differentialgleichungen erster Ordnung mit Anfangsbedingungen).

Wie wählt man  $U_0$ ? Verschiedene Möglichkeiten, u.a.

- (i).  $U_0 := u_0^I$  (Interpolierende im FE-Raum)
- (ii).  $U_0$  als  $L^2$ -Projektion von  $u_0$

**Beispiel**

$$u_t - u_{xx} = f$$

auf  $(0, 1)$  mit homogenen Randbedingungen und entsprechenden Anfangsbedingungen. Nutze lineare finite Elemente und äquidistantes Gitter. Dann gilt

$$D_h = \frac{h}{6} \cdot (d_{ij})_{i,j} \qquad A_h := \frac{1}{h} \cdot (a_{ij})_{i,j}$$

mit

$$d_{ij} := \begin{cases} 4 & i = j \\ 1 & |i - j| = 1 \\ 0 & \text{sonst} \end{cases} \qquad a_{ij} := \begin{cases} 2 & i = j \\ -1 & |i - j| = 1 \\ 0 & \text{sonst} \end{cases}$$

Nächstes Ziel: Untersuchung der Steifigkeit. Dazu: Was sind die Eigenwerte von  $A_h$ ? Man kann die Eigenwerte explizit angeben, dabei stellt man fest, dass Eigenwerte  $\lambda_1 \approx h$ ,  $\lambda_* \approx \frac{1}{h}$  existieren (d.h.  $\text{cond}A_h \approx h^{-2}$ ). Damit folgt, dass es sich i.A. um ein steifes System handelt. Favorisiert werden dann A-stabile Methoden zur Zeitdiskretisierung.

**8.5 Lemma** (Fehlerabschätzung für Semidiskretisierung)

Betrachtet werde das Problem (8.1) mit Diskretisierung (8.2) (mit hinreichend schönem  $a$ ). Weiterhin sei  $a$  nicht abhängig von  $t$ . Sei  $\Theta := R_h u - u_h$ . Für lineare finite Elemente gilt dann

$$\|\Theta\|_H^2 + \alpha \cdot \int_0^t \|\Theta\|_V^2 \leq C(u) \cdot h^4$$

für  $U_0 := R_h u$ .

Beweis: Es gilt  $u_h - u = u_h - R_h u + R_h u - u$  wobei  $R_h$  die Projektion in  $V_h$  bezeichnet. Wähle hier speziell die Ritz-Projektion (d.h.  $R_h u \in V_h$  mit  $a(R_h u, v_h) = a(u, v_h)$  für  $v_h \in V_h$ ). Dann folgt insbesondere

$$a(u - R_h u, v_h) = 0 \qquad (v_h \in V_h) \tag{8.3}$$

(Galerkin-Orthogonalität). Der Fehler  $R_h u - u$  der Ritz-Projektion kann somit genau so abgeschätzt werden wie ein FE-Diskretisierungsfehler.

Noch abzuschätzen:  $u_h - R_h u$ . Setze

$$u_h - R_h u =: \Theta \qquad R_h u - u =: \varrho$$

Dann gilt

$$\left( \frac{d}{dt} \Theta, v_h \right) + a(\Theta, v_h) = (f, v_h) - \left( \frac{d}{dt} R_h u, v_h \right) - \underbrace{a(R_h u, v_h)}_{\stackrel{(8.3)}{=} a(u, v_h)}$$

Aus (8.1) bekannt:

$$a(u, v_h) = (f, v_h) - \left( \frac{d}{dt} u, v_h \right)$$

Damit:

$$\left(\frac{d}{dt}\Theta, v_h\right) + a(\theta, v_h) = \left(\frac{d}{dt}u - \frac{d}{dt}R_h u, v_h\right) = -\left(\frac{d}{dt}\varrho, v_h\right)$$

d.h. Gleichung für  $\Theta$  gefunden. Als Anfangsbedingungen hat man  $\Theta(0) = U_0 - R_h u_0$ . Verschiedene Wahlmöglichkeiten:

- (i). Eleganteste Wahl (für die Fehlerabschätzung) ist  $U_0 := R_h u_0$ .
- (ii). Alternativ für  $U_0 := u_0^I$ :

$$\Theta(0) = u_0^I - R_h u_0 = (u_0^I - u_0) + (u_0 - R_h u_0)$$

Für lineare finite Elemente würde dann  $|\Theta(0)| \leq C \cdot h$  folgen, wenn  $u_0 \in H^2$ .

Setze hier  $U_0 := R_h u_0$ . Nächster Schritt: Nutzung der a priori-Abschätzung 8.3. Dazu benötigen wir eine Abschätzung von  $\|\frac{d}{dt}\varrho\|_H$ . Da  $a$  nach Voraussetzung nicht von  $t$  abhängt, gilt

$$a\left(\frac{d}{dt}R_h u, v_h\right) = a\left(\frac{d}{dt}u, v_h\right)$$

d.h.  $\|\frac{d}{dt}\varrho\|_H$  verhält sich wie  $\|\varrho\|_H$ . □

### 8.6 Satz

Betrachtet werde das Problem (8.1) mit Diskretisierung (8.2) (mit hinreichend schönem  $a$ ). Weiterhin sei  $a$  nicht abhängig von  $t$ . Für lineare finite Elemente gilt dann

$$\|u - u_h\|_H \leq C(u) \cdot h^2$$

für  $U_0 := R_h u$ .

**Bemerkung** Es gilt  $\|u - R_h u\|_V \leq c(u) \cdot h$ .