

# Theorie und Numerik von Variationsungleichungen

Mitschrift  
von M. Dücker  
unter Mitarbeit von  
C. Fasel, J. Frohne und I. Cherlenyak

zu einer Vorlesung von Prof. F.-T. Suttmeier

Fachbereich 6 - Mathematik  
der Universität Siegen

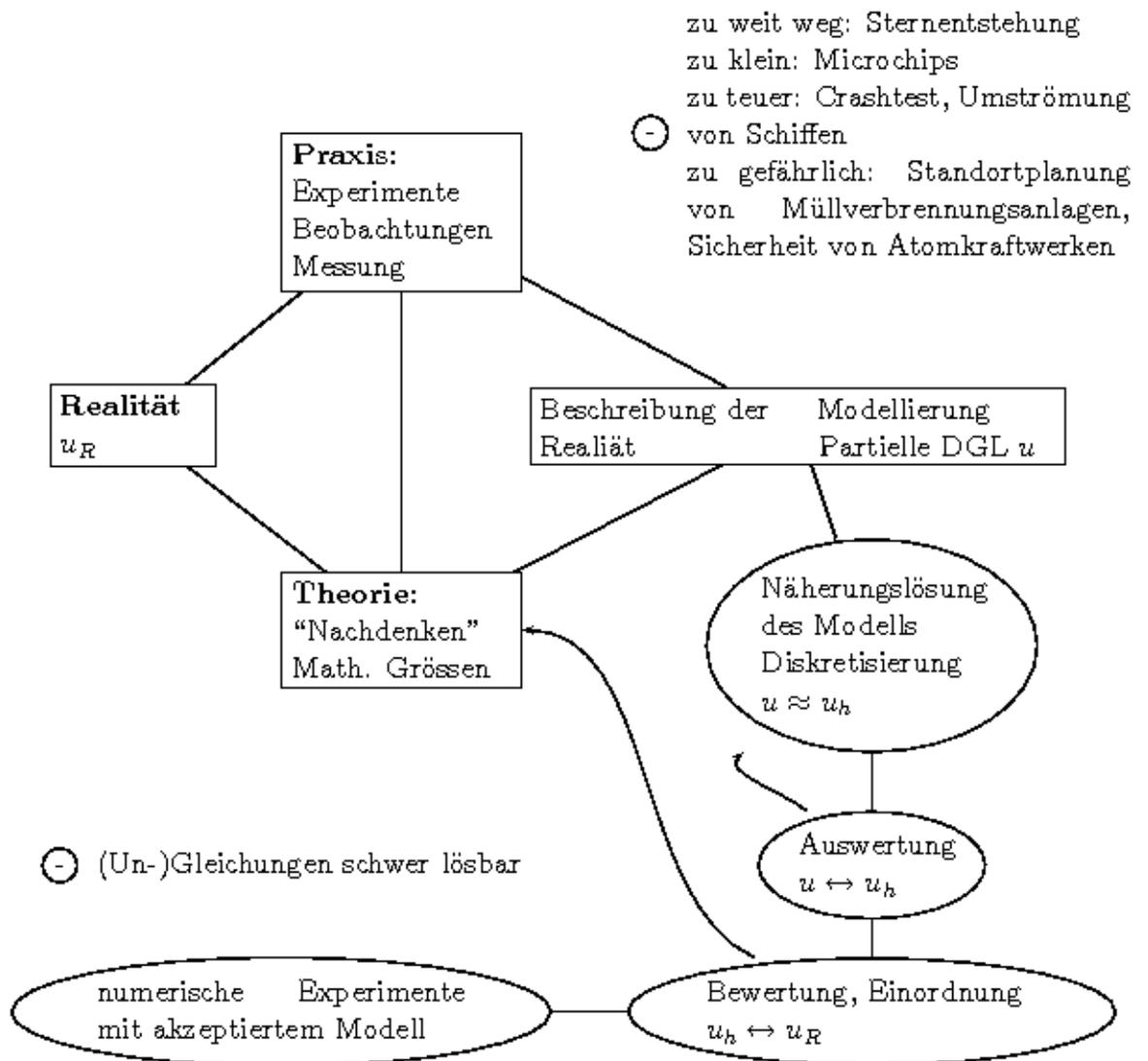
WS 2003/04

# Inhaltsverzeichnis

<b>1</b>	<b>Numerische Simulation: Eine Übersicht</b>	<b>2</b>
<b>2</b>	<b>Einleitung zur FEM (Finite Element Methode)</b>	<b>6</b>
2.1	Modell Beispiel . . . . .	6
2.2	Klassische und variationelle Formulierung . . . . .	7
2.3	Näherungslösung, Ritz-Galerkin-Verfahren . . . . .	9
2.3.1	Erste Fehlerabschätzung, Galerkin-Eigenschaft . . . . .	10
2.4	Einfache Finite Elemente . . . . .	11
2.4.1	Lineare Finite Elemente . . . . .	11
2.4.2	Interpolationsfehler . . . . .	12
2.4.3	Energiefehler-Abschätzung . . . . .	15
2.5	Variationsungleichungen . . . . .	16
2.5.1	Minimumsuche 1D . . . . .	16
2.5.2	Minimierung auf konvexer Menge $K \subset \mathbb{R}^N$ . . . . .	17
2.5.3	Minimierung auf $K \subset V$ . . . . .	17
2.6	A posteriori Fehlerschätzer . . . . .	19
2.7	Referenzelement, Gebietstransformation . . . . .	21
2.8	Rechentechnische Betrachtungen . . . . .	23
<b>3</b>	<b>FEM für elliptische Probleme</b>	<b>24</b>
3.1	Poisson Problem . . . . .	24
3.2	Natürliche und wesentliche Randbedingung . . . . .	26
3.3	Sobolev-Räume . . . . .	27
3.4	Abstrakte Formulierung . . . . .	29
3.5	Diskretisierung . . . . .	31
3.6	Variationsungleichungen . . . . .	32
3.7	Lineare Funktionale . . . . .	36
3.8	Interpolation . . . . .	38
<b>4</b>	<b>Minimierungsalgorithmen, iterative Methoden</b>	<b>41</b>
4.1	Positiv definite Matrizen . . . . .	41
4.2	Abstiegsverfahren . . . . .	42
4.3	Gradientenverfahren . . . . .	43
4.4	Projiziertes Gradientenverfahren . . . . .	46
4.5	Konjugiertes Gradientenverfahren (cg) . . . . .	47
4.6	Vorkonditionierung . . . . .	56

<b>5</b>	<b>Adaptivität</b>	<b>58</b>
5.1	Laplace-Problem . . . . .	58
5.2	Hindernisproblem . . . . .	60
5.3	Hindernisproblem in Lagrange-Formulierung . . . . .	60
5.4	Sattelpunktsuche . . . . .	63
5.5	Dualitätsargument . . . . .	66
5.5.1	A Priori Abschätzung . . . . .	67
5.5.2	A posteriori Abschätzung . . . . .	68
<b>6</b>	<b>Parabolische Probleme</b>	<b>70</b>
6.1	Parabolische Variationsungleichungen . . . . .	71
<b>7</b>	<b>Sattelpunktprobleme</b>	<b>73</b>
7.1	Hilfsmittel aus der Funktionalanalysis . . . . .	75
7.1.1	Adjungierte Operatoren . . . . .	75
7.1.2	Abstrakter Existenzsatz . . . . .	76
7.1.3	Abstrakter Konvergenzsatz . . . . .	77
7.2	Die Inf-sup-Bedingung . . . . .	78
7.3	Gemischte Finite-Element-Methoden . . . . .	79
7.4	Diskrete Sattelpunktprobleme . . . . .	80
7.5	Laplace-Gleichung als gemischtes Problem . . . . .	82
7.5.1	Primal-gemischte variationelle Formulierung . . . . .	82
7.5.2	Dual-gemischte Formulierung . . . . .	83
<b>8</b>	<b>Themengebiete</b>	<b>85</b>

# 1 Numerische Simulation: Eine Übersicht



**Beispiel 1.1** Auslenkung eines Drahtes.

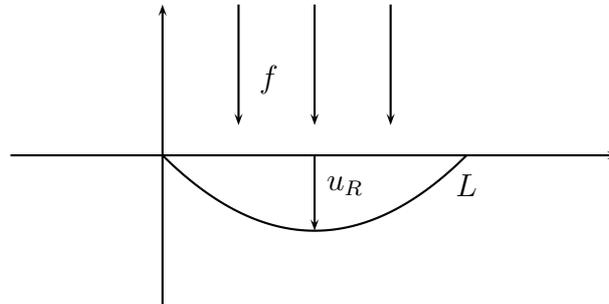


Abbildung 1.1: Draht der sich durch  $f$  verbiegt.

**Praxis:** Je grösser  $f$ , desto grösser ist die Auslenkung  $u_R$ .

**Theorie:** Längenänderung  $\Delta l$  ist proportional zur "elastischen Energie":

$$\Delta l \sim \frac{1}{2} \int_I (\partial_x u)^2 dx$$

**Modellierung**

$$-\partial_x^2 u = f, \quad u(0) = u(L) = 0$$

**Diskretisierung**

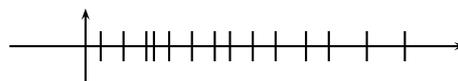


Abbildung 1.2: Diskretisierung

Wähle  $n$  bel. Punkte  $x_i \in I$ ,  $i = 1, \dots, n$ .

**Bezeichnung:**  $u_i = u(x_i)$ .

Approximation von  $\partial_x^2 u$  :

$$\partial_x^2 u(x_i) \approx \frac{u_{i-1} + 2u_i + u_{i+1}}{h^2}$$

Man erhält ein System von Gleichungen

$$-u_{i-1} + 2u_i - u_{i+1} = h^2 f_i, \quad i = 1, \dots, n.$$

Beachte die Randbedingungen:  $u(0) = u(L) = 0$ .

Mit den Bezeichnungen:  $x, b \in \mathbb{R}^n$ :

$$x = \begin{pmatrix} u_1 \\ \vdots \\ u_n \end{pmatrix}, \quad b = \begin{pmatrix} h^2 f_1 \\ \vdots \\ h^2 f_n \end{pmatrix}$$

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

$$A = \begin{pmatrix} 2 & -1 & & & & \\ -1 & 2 & -1 & & & \\ & & \ddots & \ddots & \ddots & \\ & & & -1 & 2 & -1 \\ & & & & -1 & 2 \end{pmatrix}$$

Löse das Gleichungssystem  $Ax = b$ .

Bestimme durch lineare Interpolation:  $x \rightarrow u_h$ .

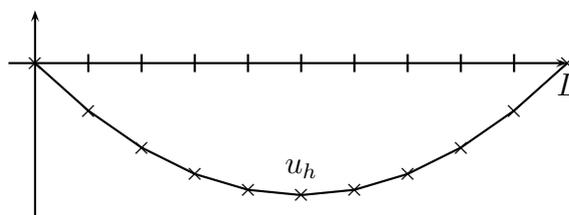


Abbildung 1.3: Interpolation

**Auswertung der Rechnung:** Ist  $\|u - u_h\|$  klein genug?

Bewertung in Bezug auf  $u_R$ :

$$\begin{aligned} \|u_R - u_h\| &= \|u_R - u + u - u_h\| \\ &\leq \underbrace{\|u_R - u\|}_{\text{Modellfehler}} + \underbrace{\|u - u_h\|}_{\text{Diskretisierungsfehler}} \end{aligned}$$

**Verbesserung des Modells:**

$$-\partial_x^2 u = f \rightarrow -\mu \cdot \partial_x^2 u = f, \quad \mu \text{ Elastizitätskonstante}$$

Akzeptiertes Modell z.B.:  $-0.75 \cdot \partial_x^2 u = f$

Numerische Experimente:

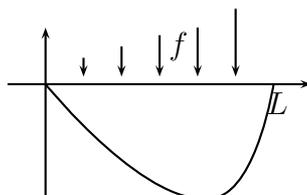
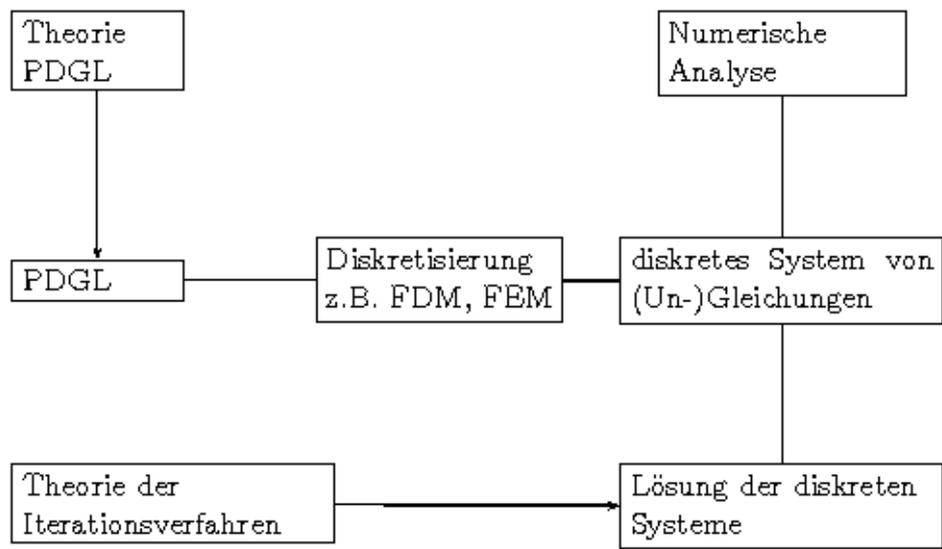


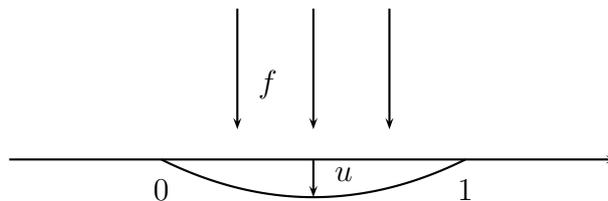
Abbildung 1.4: Numerische Experimente



# 2 Einleitung zur FEM (Finite Element Methode)

## 2.1 Modell Beispiel

### Elastischer Draht



**Beschreibung der Auslenkung  $u(x)$ :**

Betrachte die Längenänderung:  $\Delta l = \int_I \sqrt{1 + (\partial_x u)^2} dx - \int_I 1 dx$ .

Benutze Taylor-Entwicklung:  $\sqrt{1+y} \approx \sqrt{1+0} + \frac{1}{2\sqrt{1+0}}(y-0)$

$$\Delta l \approx \frac{1}{2} \int_I (\partial_x u)^2 dx$$

**Physik, Theorie:** Elastische Energie ist proportional zu  $\Delta l$ :

$$U_E \sim \Delta l$$

Durch Einwirken der Kraft  $f$  besitzt der Draht eine potentielle Energie:

$$U_f = - \int_I f \cdot u dx \quad (\text{Arbeit} = \text{Kraft} \cdot \text{Weg})$$

Stabile Gleichgewichtslage ist charakterisiert dadurch, dass die Gesamtenergie minimal wird.

$$U(u) = U_E + U_f = \frac{1}{2} \int_I (\partial_x u)^2 dx - \int_I f \cdot u dx$$

Also:

**Gesucht ist**  $u \in V$ , **so dass**

$$U(u) \leq U(v) \quad \forall v \in V \quad (M)$$

$V$  ist der Raum der Vergleichsfunktionen,  $V$  ist festgelegt durch:

1. Alle stetigen Funktionen mit Nullrandwerten.
2. Funktionen haben stückweise stetige, beschränkte 1. Ableitungen (Integrale müssen Sinn machen).

## 2.2 Klassische und variationelle Formulierung

### Variationsrechnung

Wähle  $\varphi \in V$  beliebig, aber fest.  
 Betrachte  $v = u + \varepsilon\varphi$  für  $\varepsilon \in \mathbb{R}$ .

Aufgrund der Minimaleigenschaft

$$U(u) \leq U(v) = U(u + \varepsilon\varphi) \quad \forall \varepsilon$$

folgt die notwendige Bedingung:

$$\begin{aligned} & \frac{d}{d\varepsilon} U(u + \varepsilon\varphi)|_{\varepsilon=0} = 0 \\ \Leftrightarrow & \frac{d}{d\varepsilon} \left[ \frac{1}{2} \int_I (\partial_x(u + \varepsilon\varphi))^2 dx - \int_I f \cdot (u + \varepsilon\varphi) dx \right]_{\varepsilon=0} = 0 \\ \Leftrightarrow & \left[ \frac{1}{2} \int_I 2(\partial_x u + \varepsilon \partial_x \varphi) \partial_x \varphi dx - \int_I f \cdot \varphi dx \right]_{\varepsilon=0} = 0 \\ \Leftrightarrow & \int_I \partial_x u \partial_x \varphi dx - \int_I f \cdot \varphi dx = 0 \end{aligned}$$

Das Funktional  $U(\cdot)$  ist konvex und daher ist die Bedingung

$$\int_I \partial_x u \partial_x \varphi dx = \int_I f \varphi dx \quad \forall \varphi \in V \quad (V)$$

auch hinreichend.

Man spricht bei (V) von der variationellen (oder schwachen) Formulierung.  
 Obige Rechnung zeigt:

**Satz 2.1**  $(M) \Rightarrow (V)$ .

Umgekehrt gilt:

**Satz 2.2**  $(V) \Rightarrow (M)$ .

**Beweis:** Schreibweise:  $(v, w) = \int_I v(x)w(x) dx$  für stückweise stetige beschränkte Funktionen.

1.  $u$  sei Lösung von  $(V)$ .
2. Wähle  $v \in V$  und setze  $w = v - u$ .
3. Somit gilt:  $v = w + u$  und  $w \in V$ .

$$\begin{aligned}
 U(v) = U(u + w) &= \frac{1}{2}(\partial_x u + \partial_x w, \partial_x u + \partial_x w) - (f, u + w) \\
 &= \frac{1}{2}(\partial_x u, \partial_x u) - (f, u) + \underbrace{(\partial_x u, \partial_x w) - (f, w)}_{=0 \text{ (wegen) } (V)} + \\
 &\quad + \underbrace{\frac{1}{2}(\partial_x w, \partial_x w)}_{\geq 0} \\
 &\geq \frac{1}{2}(\partial_x u, \partial_x u) - (f, u) \\
 &= U(u)
 \end{aligned}$$

Also gilt:

$$U(u) \leq U(v) \quad \forall v \in V.$$

□

**Satz 2.3** (Klassische Formulierung)

Sei  $u$  Lösung von  $(V)$ . Zusätzlich existiere  $\partial_x^2 u$  und sei stetig. Dann gilt:

$$-\partial_x^2 u = f \quad u(0) = u(1) = 0 \quad (D)$$

**Beweis:**

$$\int_0^1 \partial_x u \partial_x v dx - \int_0^1 f \cdot v dx = 0 \quad \forall v \in V$$

Partielle Integration liefert:

$$[\partial_x u \cdot v]_0^1 - \int_0^1 \partial_x^2 u \cdot v dx - \int_0^1 f \cdot v dx = 0$$

Beachte  $v(0) = v(1) = 0$ :

$$\int_0^1 (-\partial_x^2 u - f) v \, dx = 0 \quad \forall v \in V$$

$$\Rightarrow -\partial_x^2 u = f \text{ auf } [0, 1]$$

□

Verwende gleiche Rechentechnik (also partielle Integration) um zu zeigen:

**Satz 2.4**  $(D) \Rightarrow (V)$ .

**Zusammenfassung:**

$$(D) \Rightarrow (V) \Leftrightarrow (M)$$

## 2.3 Näherungslösung, Ritz-Galerkin-Verfahren

**Idee:** Approximation des Raumes  $V$  durch einen endlich dimensionalen Teilraum  $V^N$  mit  $\dim V^N = N$ .

Betrachte die schwache Formulierung  $(V)$

$$\int_I \partial_x u^N \partial_x \varphi^N \, dx = \int_I f \cdot \varphi^N \, dx \quad \forall \varphi \in V^N$$

Schreibweise:

$$a(v, w) = (\partial_x v, \partial_x w)$$

Also gilt kompakt:

$$a(u^N, \varphi^N) = (f, \varphi^N)$$

Wähle Basis für  $V^N$ :

$$V^N = \langle \varphi_1, \dots, \varphi_N \rangle$$

Darstellung für  $\varphi \in V^N$  :

$$\varphi = \sum_{j=1}^N v_j \varphi_j, \quad v_j \in \mathbb{R}$$

Es genügt

$$a(u^N, \varphi_i) = (f, \varphi_i) \quad \forall i = 1, \dots, N \quad (2.1)$$

zu erfüllen:

$$a\left(u^N, \sum_{i=1}^N v_i \varphi_i\right) - \left(f, \sum_{i=1}^N v_i \varphi_i\right) = \sum_{i=1}^N v_i \underbrace{\left(a(u^N, \varphi_i) - (f, \varphi_i)\right)}_{=0} = 0 \quad (2.2)$$

Wie berechnet man  $u^N$ ?

Einsetzen von  $u^N = \sum_j u_j \varphi_j$ ,  $u_j \in \mathbb{R}$  in (2.1):

$$\begin{aligned}
 & a(u^N, \varphi_i) = (f, \varphi_i) \\
 \Leftrightarrow & a\left(\sum_j u_j \varphi_j, \varphi_i\right) = (f, \varphi_i) \\
 \Leftrightarrow & \sum_j u_j a(\varphi_j, \varphi_i) = (f, \varphi_i) \\
 \Leftrightarrow & \sum_j u_j \left( \int_I \partial_x \varphi_j \partial_x \varphi_i dx \right) = (f, \varphi_i) \quad i = 1, \dots, N
 \end{aligned}$$

Die Bestimmung  $u^N \in V^N$  wird somit auf das Lösen eines linearen Gleichungssystems zurückgeführt:

$$Ax = b$$

Mit  $x^T = (u_1, \dots, u_N)$ ,  $b^T = (b_1, \dots, b_N)$ ,  $b_i = \int_I f \cdot \varphi_i dx$

$$A_{ij} = \int_I \partial_x \varphi_j \partial_x \varphi_i dx, \quad A \in \mathbb{R}^{N \times N}.$$

### 2.3.1 Erste Fehlerabschätzung, Galerkin-Eigenschaft

$$\begin{aligned}
 (\partial_x u, \partial_x \varphi) &= (f, \varphi) & \forall \varphi \in V \\
 (\partial_x u^N, \partial_x \varphi) &= (f, \varphi) & \forall \varphi \in V^N
 \end{aligned}$$

Weitere Voraussetzung:  $V^N \subset V$

Für ein  $\varphi \in V^N$  gilt:

$$\begin{aligned}
 (\partial_x u, \partial_x \varphi) &= (f, \varphi) \\
 -(\partial_x u^N, \partial_x \varphi) &= (f, \varphi) \\
 \hline
 (\partial_x u - \partial_x u^N, \partial_x \varphi) &= 0 \quad \text{Galerkin-Eigenschaft}
 \end{aligned}$$

Abschätzung des Fehlers:

$$\begin{aligned}
 \|\partial_x u - \partial_x u^N\|^2 &= (\partial_x u - \partial_x u^N, \partial_x u - \partial_x \varphi + \partial_x \varphi - \partial_x u^N) \\
 &= (\partial_x u - \partial_x u^N, \partial_x u - \partial_x \varphi) + (\partial_x u - \partial_x u^N, \underbrace{\partial_x \varphi - \partial_x u^N}_{\in V^N})
 \end{aligned}$$

Falls  $\varphi \in V^N$ , dann verschwindet der 2. Term. Es bleibt:

$$\|\partial_x u - \partial_x u^N\|^2 \leq \|\partial_x u - \partial_x u^N\| \|\partial_x u - \partial_x \varphi\|$$

Nach Division:

$$\|\partial_x u - \partial_x u^N\| \leq \|\partial_x u - \partial_x \varphi\|$$

Gehe über zu:

$$\|\partial_x u - \partial_x u^N\| \leq \inf_{\varphi \in V^N} \|\partial_x u - \partial_x \varphi\|$$

## 2.4 Einfache Finite Elemente

Allgemeine Überlegung:

- Lösung  $u$  sollte gut approximierbar sein.
- Leichte Berechnung der Einträge von  $A$ .
- $A$  sollte für die Numerik gute Eigenschaften haben. z.B. dünnbesetzt, moderate Kondition.

### 2.4.1 Lineare Finite Elemente

Idee:  $V^N$  besteht aus stückweise linearen Funktionen die global stetig sind.

Berechnung von  $A$ :

Basiswahl:

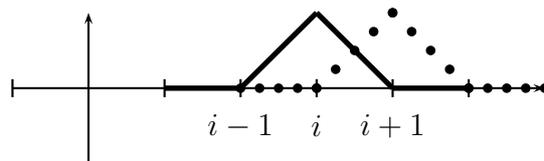


Abbildung 2.1: Basisfunktionen

Basisfunktionen  $\varphi_i$ :

$$\varphi_i(x_i) = 1, \quad \varphi_i(x_{i+1}) = 0, \quad \varphi_i \text{ linear auf } (x_i, x_{i+1})$$

$$\partial_x \varphi_i = -\frac{1}{h} \text{ auf } (x_i, x_{i+1})$$

$$\begin{aligned} \int_I \partial_x \varphi_i \partial_x \varphi_i \, dx &= \int_{x_{i-1}}^{x_{i+1}} \partial_x \varphi_i \partial_x \varphi_i \, dx \\ &= 2 \int_{x_i}^{x_{i+1}} \partial_x \varphi_i \partial_x \varphi_i \, dx \\ &= 2 \frac{1}{h^2} (x_{i+1} - x_i) = \frac{2}{h} = A_{ii} \\ \int_I \partial_x \varphi_i \partial_x \varphi_{i+1} \, dx &= \int_{x_i}^{x_{i+1}} \partial_x \varphi_i \partial_x \varphi_{i+1} \, dx \\ &= (x_{i+1} - x_i) \left( -\frac{1}{h} \right) \left( \frac{1}{h} \right) = -\frac{1}{h} = A_{i,i+1} = A_{i+1,i} \end{aligned}$$

$$A = \frac{1}{h} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix}$$

Rechte Seite:  $b_i = \int_I f \varphi_i dx$

Falls  $f$  konstant ist:  $b_i = f \int_{x_{i-1}}^{x_{i+1}} \varphi_i dx = hf_i$

## 2.4.2 Interpolationsfehler

**Motivation:**

$$\|\partial_x(u - u_h)\| \leq \inf_{\varphi \in V_h} \|\partial_x(u - \varphi)\| \leq \|\partial_x(u - I_h u)\|$$

$I_h u$  bezeichne die lineare Interpolierende von  $u$ .

**Satz 2.5** Auf einem Teilintervall  $T$ ,  $T = (a_1, a_2)$ ,  $h_T = a_2 - a_1$  der Zerlegung des Rechengebietes  $I \subset \mathbb{R}$  gilt:

$$\|v - I_h v\|_{L_\infty(T)} \leq ch_T^2 \|\partial_x^2 v\|_{L_\infty(T)} \quad (2.3)$$

$$\|\partial_x(v - I_h v)\|_{L_\infty(T)} \leq ch_T \|\partial_x^2 v\|_{L_\infty(T)} \quad (2.4)$$

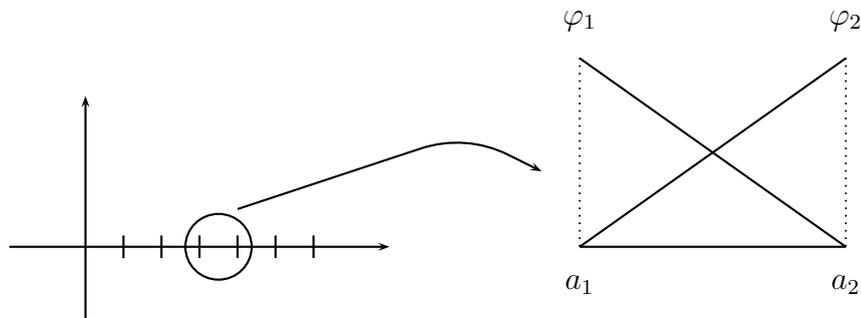


Abbildung 2.2: Skizze

**Beweis: Vorbereitungen:**

Die ‘‘Hutfunktionen’’  $\varphi_1, \varphi_2$  bestimmen die Basis für  $P_1(T)$ . Allgemein gilt für  $w \in P_1(T)$ :

$$w(x) = \sum_{i=1}^2 w(a_i) \varphi_i(x), \quad x \in T$$

Also

$$I_h v(x) = \sum_{i=1}^2 v(a_i) \varphi_i(x), \quad x \in T. \quad (2.5)$$

Betrachte die Taylor-Entwicklung von  $v$  um  $x$  in  $a_i$ :

$$v(a_i) = v(x) + \partial_x v(x) (a_i - x) + \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \quad (2.6)$$

Zu zeigen ist (2.3). Einsetzen von (2.6) in (2.5):

$$I_h v(x) = \sum_{i=1}^2 \left( v(x) + \partial_x v(x) (a_i - x) + \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \right) \varphi_i(x)$$

Umsortieren ergibt:

$$\begin{aligned} I_h v(x) &= v(x) \sum_{i=1}^2 \varphi_i(x) + \sum_{i=1}^2 \partial_x v(x) (a_i - x) \varphi_i(x) + \\ &\quad + \sum_{i=1}^2 \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \varphi_i(x) \end{aligned} \quad (2.7)$$

Wir zeigen später:

$$\sum_{i=1}^2 \varphi_i(x) = 1, \quad \sum_{i=1}^2 \partial_x v(x) (a_i - x) \varphi_i(x) = 0.$$

Was bleibt zunächst?

$$|I_h v(x) - v(x)| = \left| \sum_{i=1}^2 \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \varphi_i(x) \right|$$

Wegen  $\varphi_i(x) \leq 1$  und  $(a_i - x) \leq h_T$  gilt:

$$\left| I_h v(x) - v(x) \right| \leq \max_{\xi \in T} |\partial_x^2 v(\xi)| \cdot h_T^2$$

Bleibt zu zeigen:  $\sum_{i=1}^2 \varphi_i(x) = 1$ :

Betrachte für  $v(x) = 1$

$$\Rightarrow \partial_x v(x) = \partial_x^2 v(x) = 0$$

$I_h v = 1$ . Einsetzen in (2.7):

$$I_h v = 1 = 1 \cdot \sum_{i=1}^2 \varphi_i(x) \quad (\square)$$

Bleibt zu zeigen:  $\sum_{i=1}^2 \partial_x v(x) (a_i - x) \varphi_i(x) = 0$ .

Sei  $v$  gegeben. Setze für festes  $x$ :  $d = \partial_x v(x)$ .

Ansatz:

$$\begin{aligned} w(x) &= d \cdot x \\ I_h w &= w \\ \partial_x w &= d \\ \partial_x^2 w &= 0 \end{aligned}$$

Einsetzen in (2.7):

$$w = w \cdot 1 + \sum_{i=1}^2 d \cdot (a_i - x) \varphi_i(x) + 0$$

(□)

Zu zeigen (2.4):

$$\|\partial_x(v - I_h v)\|_{L^\infty(T)} \leq ch_T \|\partial_x^2 v\|_{L^\infty(T)}$$

Betrachte:

$$\partial_x I_h v(x) = \sum_{i=1}^2 v(a_i) \partial_x \varphi_i(x)$$

Einsetzen der Taylor-Entwicklung für  $v(a_i)$ :

$$\begin{aligned} \partial_x I_h v(x) &= v(x) \sum_{i=1}^2 \partial_x \varphi_i(x) + \sum_{i=1}^2 \partial_x v(x) (a_i - x) \partial_x \varphi_i(x) \\ &\quad + \sum_{i=1}^2 \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \partial_x \varphi_i(x) \end{aligned}$$

Nun gilt:  $\partial_x \varphi_1 = -\frac{1}{h_T}$ ,  $\partial_x \varphi_2 = \frac{1}{h_T}$

Damit

$$\sum_{i=1}^2 \partial_x \varphi_i(x) = 0$$

Weiterhin:

$$\begin{aligned} \sum_{i=1}^2 \partial_x v(x) (a_i - x) \partial_x \varphi_i(x) &= \partial_x v(x) (a_1 - x) \left(-\frac{1}{h_T}\right) + \\ &\quad + \partial_x v(x) (a_2 - x) \frac{1}{h_T} \\ &= \partial_x v(x) \frac{(a_2 - a_1)}{h_T} \\ &= \partial_x v(x) \end{aligned}$$

Es bleibt also:

$$\begin{aligned} |\partial_x I_h v - \partial_x v(x)| &= \left| \sum_{i=1}^2 \frac{1}{2} \partial_x^2 v(\xi_i) (a_i - x)^2 \partial_x \varphi_i(x) \right| \\ &\leq \max_{\xi \in T} |\partial_x^2 v(\xi)| h_T^2 \frac{1}{h_T} \end{aligned}$$

□

### 2.4.3 Energiefehler-Abschätzung

$$\begin{aligned} \|\partial_x u - \partial_x I_h u\|_I^2 &= \sum_T \int_T (\partial_x u - \partial_x I_h u)^2 dx \\ &\leq \sum_T \int_T h_T^2 \|\partial_x^2 u(x)\|_{L^\infty(T)}^2 dx \\ &\leq \sum_T h_T^2 \|\partial_x^2 u(x)\|_{L^\infty(T)}^2 \int_I 1 dx \\ &\leq h_{max}^2 \|\partial_x^2 u(x)\|_{L^\infty(I)}^2 \mu(I) \end{aligned}$$

mit  $\mu(I) = \int_I 1 dx$ . D.h es gilt:

$$\|\partial_x u - \partial_x I_h u\|_I \leq ch_{max} \|\partial_x^2 u(x)\|_{L^\infty(I)}$$

Wir hatten:

$$\|\partial_x u - \partial_x u_h\| \leq \|\partial_x u - \partial_x I_h u\| \leq ch_{max} \|\partial_x^2 u(x)\|_{L^\infty(I)}$$

**Satz 2.6** (Interpolationsfehler auf  $I$ .)

Auf  $I \subset \mathbb{R}$  gilt bei gegebener Zerlegung in  $T \subset I$  mit maximaler Grösse  $h$ :

$$\|\partial_x^i (v - I_h v)\| \leq ch^{2-i} \|\partial_x^2 v\|_{L^\infty(I)}, \quad i = 0, 1$$

**Beweis:**

$$\begin{aligned} \|\partial_x^i (v - I_h v)\|^2 &= \sum_T \int_T (\partial_x^i (v - I_h v))^2 dx \\ &\leq \sum_T \int_T ch_T^{2(2-i)} \|\partial_x^2 v\|_{L^\infty(T)}^2 dx \\ &\leq \sum_T ch_T^{2(2-i)} \|\partial_x^2 v\|_{L^\infty(T)}^2 \int_T 1 dx \\ &\leq ch^{2(2-i)} \|\partial_x^2 v\|_{L^\infty(I)}^2 \mu(I) \end{aligned}$$

□

**Satz 2.7** (Energiefehler)

Auf  $I \subset \mathbb{R}$  gilt bei gegebener Zerlegung in Teilintervalle  $T$  mit maximaler Grösse  $h$ :

$$\|\partial_x(u - u_h)\| \leq ch \|\partial_x^2 u\|_{L^\infty(I)}$$

**Beweis:** Beachte:

$$\begin{aligned} \|\partial_x(u - u_h)\| &\leq \inf_{\varphi \in V_h} \|\partial_x(u - \varphi)\| \\ &\leq \|\partial_x(u - I_h u)\| \\ &\leq ch \|\partial_x^2 u\|_{L^\infty(I)} \end{aligned}$$

□

## 2.5 Variationsungleichungen

### 2.5.1 Minimumsuche 1D

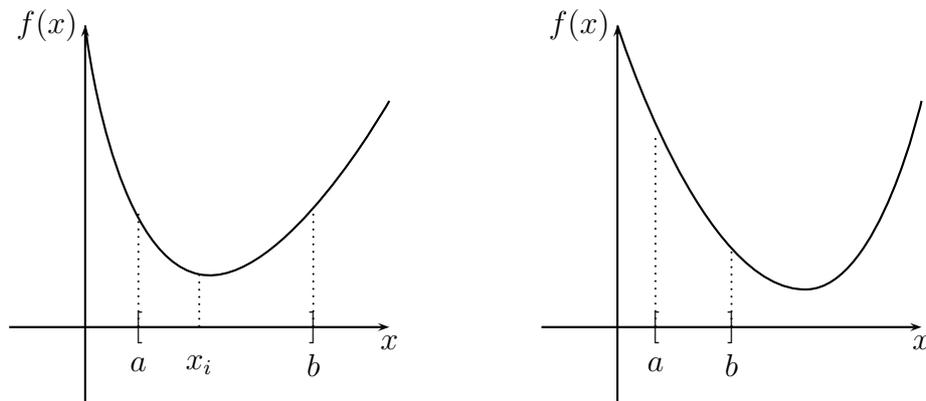


Abbildung 2.3: Fallunterscheidung

Weitere Voraussetzung:  $f$  stetig differenzierbar.

Fallunterscheidung auch für Minima am Rand:

$$\begin{aligned} f(a) &\leq f(x) \quad \forall x \rightarrow f'(a) \geq 0 \\ f(b) &\leq f(x) \quad \forall x \rightarrow f'(b) \leq 0 \\ f(x_i) &\leq f(x) \quad \forall x \rightarrow f'(x_i) = 0 \end{aligned}$$

Kompakte Darstellung der notwendigen Bedingung für eine Minimalstelle  $x_0$ :

$$f'(x_0) \cdot (x - x_0) \geq 0 \quad \forall x$$

### 2.5.2 Minimierung auf konvexer Menge $K \subset \mathbb{R}^N$

Für  $f : K \rightarrow \mathbb{R}$ . Gesucht ist  $x_0$ :

$$f(x_0) \leq f(x) \quad \forall x \in K$$

Betrachte  $F(\varepsilon) = f(x_0 + \varepsilon(x - x_0))$

Aus dem 1-dimensionalen Fall ist bekannt:

$$\begin{aligned} \frac{d}{d\varepsilon} F(\varepsilon)|_{\varepsilon=0} \cdot \varepsilon &\geq 0 \quad \forall \varepsilon \geq 0 \\ \Rightarrow \nabla f(x_0) \cdot (x - x_0) \cdot \varepsilon &\geq 0 \quad \forall \varepsilon \geq 0 \\ \Rightarrow \nabla f(x_0) \cdot (x - x_0) &\geq 0 \quad \forall x \in K \end{aligned}$$

### 2.5.3 Minimierung auf $K \subset V$

$$K = \{v \in V \mid v \geq g\}$$

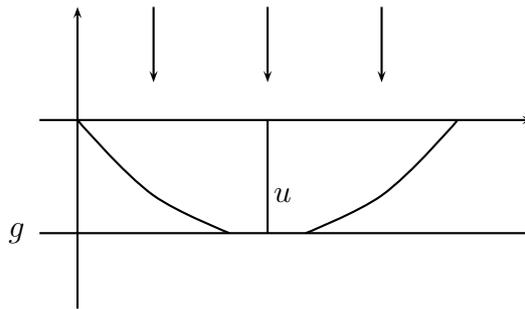


Abbildung 2.4: Skizze

**Bemerkung 2.1**  $K$  ist konvex:  $v_1, v_2 \in K, \alpha \in (0, 1)$ .

$$\alpha v_1 + (1 - \alpha)v_2 \geq \alpha g + (1 - \alpha)g = g$$

Welches Funktional wird minimiert?

$$U(v) = \frac{1}{2} \int_I (\partial_x v)^2 dx - \int_I f v dx$$

Betrachte:  $F(\varepsilon) = U(u + \varepsilon(v - u))$ .

Aus dem 1-dimensionalen Fall:

$$\frac{d}{d\varepsilon} F(\varepsilon)|_{\varepsilon=0} \cdot \varepsilon \geq 0$$

Ausrechnen liefert:

$$\begin{aligned} & \int_I \partial_x \left( u + \varepsilon(v - u) \right) \partial_x(v - u) dx - \int_I f(v - u) dx \Big|_{\varepsilon=0} \cdot \varepsilon \geq 0 \\ \Rightarrow & \left( \int_I \partial_x u \partial_x(v - u) dx - \int_I f(v - u) dx \right) \cdot \varepsilon \geq 0 \quad \forall \varepsilon \geq 0 \end{aligned}$$

Zusammengefasst:

$$\left( \partial_x u, \partial_x(v - u) \right) \geq (f, v - u) \quad \forall v \in K$$

**Bemerkung 2.2**  $\left( \partial_x u, \partial_x(v - u) \right) \geq (f, v - u)$  ist der Prototyp einer elliptischen Variationsungleichung 1. Art.

### Die Diskretisierung mit linearen Finiten Elementen

$$a(u, \varphi - u) \geq (f, \varphi - u) \quad \forall \varphi \in K \quad (2.8)$$

$$a(u_h, \varphi - u_h) \geq (f, \varphi - u_h) \quad \forall \varphi \in K_h = V_h \cap K \quad (2.9)$$

Insbesondere gilt:  $K_h \subset K$ .

#### Satz 2.8 (Energiefehler)

Voraussetzungen wie in Satz 2.7. Dann gilt:

$$\|\partial_x(u - u_h)\| \leq O(h)$$

**Beweis:** Bezeichnung  $u_i = I_h u$ . Ausgangspunkt:

$$\begin{aligned} a(u - u_h, u - u_h) &= a(u - u_h, u - u_i + u_i - u_h) \\ &= a(u - u_h, u - u_i) + a(u - u_h, u_i - u_h) \end{aligned} \quad (2.10)$$

Bei Variationsgleichungen war der 2. Term 0. Hier:

$$\begin{aligned} a(u - u_h, u_i - u_h) &= (f, u_i - u_h) - a(u_h, u_i - u_h) \quad \text{Term 1} \\ &\quad + a(u, u_i - u) - (f, u_i - u) \\ &\quad + a(u, u - u_h) - (f, u - u_h) \quad \text{Term 2} \end{aligned}$$

Term 1  $\leq 0$  wegen Test mit  $\varphi = u_i$  in (2.9).

Term 2  $\leq 0$  wegen Test mit  $\varphi = u_h$  in (2.8).

$$\begin{aligned} a(u, u_h - u) &\geq (f, u_h - u) \\ \Leftrightarrow a(u, u_h - u) - (f, u_h - u) &\geq 0 \\ \Leftrightarrow a(u, u - u_h) - (f, u - u_h) &\leq 0 \end{aligned}$$

Nun weiter in (2.10):

$$\begin{aligned} \dots &\leq \|\partial_x(u - u_h)\| \|\partial_x(u - u_i)\| + a(u, u_i - u) - (f, u_i - u) \\ &\leq \frac{1}{2} \|\partial_x(u - u_h)\|^2 + \frac{1}{2} \|\partial_x(u - u_i)\|^2 - \int_I (\partial_x^2 u)(u_i - u) dx - (f, u_i - u) \\ &\leq \frac{1}{2} \|\partial_x(u - u_h)\|^2 + \frac{1}{2} \|\partial_x(u - u_i)\|^2 + \|\partial_x^2 u\| \|u - u_i\| + \|f\| \|u - u_i\| \end{aligned}$$

$$\begin{aligned} \Rightarrow \frac{1}{2} \|\partial_x(u - u_h)\|^2 &\leq O(h^2) + (\|\partial_x^2 u\| + \|f\|) \|u - u_i\| \\ &\leq O(h^2) + O(h^2) \end{aligned}$$

□

## 2.6 A posteriori Fehlerschätzer

**Bisher:**

$$\|\partial_x(u - u_h)\| \leq ch \|\partial_x^2 u\|_{L^\infty(I)}$$

mit unbekannter Lösung  $u$ .

**Ziel:**

$$\|\partial_x(u - u_h)\| \leq \eta(u_h, f).$$

**Satz 2.9** Auf einem Intervall  $T = (x_i, x_{i+1})$  gilt für  $v \in V$  und  $v(x_i) = 0$ :

$$\|v\|_{L^2(T)} \leq h \|\partial_x v\|_{L^2(T)}$$

**Beweis:**  $y \in (x_i, x_{i+1}]$  :

$$\begin{aligned} v(y) &= \int_{x_i}^y \partial_x v(x) dx \\ &\leq \left( \int_{x_i}^y 1^2 dx \right)^{1/2} \cdot \left( \int_{x_i}^y (\partial_x v)^2 dx \right)^{1/2} \\ &\leq \sqrt{h} \|\partial_x v\|_{L^2(T)} \end{aligned}$$

Quadrieren liefert:

$$v^2(y) \leq h \|\partial_x v\|_{L^2(T)}^2$$

Integrieren liefert:

$$\begin{aligned} \int_T v^2 dx &\leq \int_T h \|\partial_x v\|^2 dy \\ &\leq h^2 \|\partial_x v\|^2 \end{aligned}$$

□

**Satz 2.10** (Stabilität der Interpolation)

Auf  $I \subset \mathbb{R}$  gilt bei gegebener Zerlegung in Zellen  $T$ :

$$\|\partial_x I_h v\|_{L^2(T)} \leq \|\partial_x v\|_{L^2(T)}$$

**Beweis:**  $y \in (x_i, x_{i+1})$

$$\begin{aligned}
 \partial_y I_h v(y) &= \frac{v(x_{i+1}) - v(x_i)}{h} \\
 &= \frac{1}{h} \int_{x_i}^{x_{i+1}} \partial_x v(x) \, dx \\
 &\leq \frac{1}{h} \left( \int_{x_i}^{x_{i+1}} 1^2 \, dx \right)^{1/2} \left( \int_{x_i}^{x_{i+1}} (\partial_x v)^2 \, dx \right)^{1/2} \\
 &\leq \frac{1}{\sqrt{h}} \|\partial_x v\|_{L^2(T)}
 \end{aligned}$$

Quadrieren und Integrieren liefert:

$$\begin{aligned}
 \int \partial_x I_h v(y)^2 \, dy &\leq \int_T \frac{1}{h} \|\partial_x v\|^2 \, dy \\
 \|\partial_x I_h v\|_{L^2(T)} &\leq \|\partial_x v\|_{L^2(T)}
 \end{aligned}$$

□

**Satz 2.11** (Energiefehlerschätzer)

Auf dem Rechengebiet  $I$  mit Zerlegung in  $T$  gilt:

$$\|\partial_x(u - u_h)\|_{L^2(I)} \leq c \left( \sum_T h_T^2 \rho_T^2 \right)^{1/2}$$

mit  $\rho_T = 2 \|f + \partial_x^2 u_h\|_{L^2(T)}$ .

**Beweis:** Schreibweisen:  $e = u - u_h$ ,  $e_i = I_h e$

$$\begin{aligned}
 \|\partial_x(u - u_h)\|_{L^2(I)}^2 &= (\partial_x u - \partial_x u_h, \partial_x e - \partial_x e_i) \\
 &= (f, e - e_i) - (\partial_x u_h, \partial_x e - \partial_x e_i)_I \\
 &= (f, e - e_i) - \sum_T (\partial_x u_h, \partial_x e - \partial_x e_i)_T \\
 &= (f, e - e_i) - \sum_T \left( (-\partial_x^2 u_h, e - e_i)_T + \underbrace{[\partial_x u_h \cdot (e - e_i)]_{x_i}^{x_{i+1}}}_{=0} \right) \\
 &= \sum_T (f + \partial_x^2 u_h, e - e_i)_T
 \end{aligned}$$

Betrachte:

$$\begin{aligned}
 (f + \partial_x^2 u_h, e - e_i)_T &\leq \|f + \partial_x^2 u_h\|_{L^2(T)} \|e - e_i\|_{L^2(T)} \\
 &\leq \|f + \partial_x^2 u_h\|_{L^2(T)} h_T \|\partial_x(e - e_i)\|_{L^2(T)} \\
 &\leq \|f + \partial_x^2 u_h\|_{L^2(T)} h_T \left( \|\partial_x e\|_{L^2(T)} + \underbrace{\|\partial_x e_i\|_{L^2(T)}}_{\leq \|\partial_x e\|_T} \right) \\
 &\leq \underbrace{\|f + \partial_x^2 u_h\|_{L^2(T)}}_{\rho_T} 2 h_T \|\partial_x e\|_{L^2(T)}
 \end{aligned}$$

Einsammeln:

$$\begin{aligned}
 \|\partial_x(u - u_h)\|^2 &\leq \sum_T h_T \rho_T \|\partial_x e\|_{L^2(T)}^2 \\
 &\leq \left( \sum_T h_T^2 \rho_T^2 \right)^{1/2} \underbrace{\left( \sum_T \|\partial_x e\|_{L^2(T)}^2 \right)^{1/2}}_{\|\partial_x e\|_I}
 \end{aligned}$$

$$\Rightarrow \|\partial_x(u - u_h)\| \leq \left( \sum_T h_T^2 \rho_T^2 \right)^{1/2} \quad \square$$

## 2.7 Referenzelement, Gebietstransformation

**Ziel:** Alle Rechnungen auf einem Referenzelement (z.B. Einheitsintervall).

- + Basisfunktionen nur einmal ausrechnen.
- + (numerische) Integrationsformeln werden nur auf dem Referenzelement benötigt.

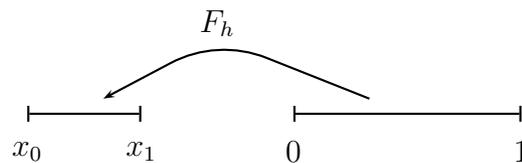


Abbildung 2.5: Gebietstransformation.

Vorbereitungen für die Substitutionsregel:

$$\begin{aligned}
 F_h : T_1 &\rightarrow T_h \\
 \xi &\mapsto x = x_0 + (1 - x_0)\xi
 \end{aligned}$$

$$\frac{d}{dx} : 1 = (x_1 - x_0) \frac{d\xi}{dx} \Rightarrow dx = \underbrace{(x_1 - x_0)}_{=:J} d\xi$$

$$F_h^{-1} : T_h \rightarrow T_1$$

$$\xi = \frac{x - x_0}{x_1 - x_0}, \quad \xi_x = \frac{d\xi}{dx} = \frac{1}{x_1 - x_0}$$

Eine Basisfunktion  $\varphi_i^h$  auf  $T_h$  wird wie folgt angesetzt:

$$\varphi_i^h(x) := \varphi_i^1(F_h^{-1}(x)) = \varphi_i^1(\xi)$$

Allgemein:

$$\begin{aligned} u(x) &= v(F_h^{-1}(x)) \\ \partial_x u(x) &= \partial_\xi v(\xi) \partial_x F_h^{-1}(x) \\ &= \partial_\xi v(\xi) \xi_x \end{aligned}$$

### Beispiel 2.1

$$\begin{aligned} \int_{T_h} f(x) \varphi_i^h(x) dx &= \int_{T_1} f(F_h(\xi)) \varphi_i^1(\xi) J d\xi \\ \int_{T_h} \partial_x \varphi_i^h(x) \partial_x \varphi_j^h(x) dx &= \int_{T_1} (\partial_x \varphi_i^1) \xi_x (\partial_\xi \varphi_j^1) \xi_x J d\xi \end{aligned}$$

Numerische Integration:

Allgemein:

$$\int_{T_1} g(\xi) d\xi \approx \sum_{k=1}^q w_k g(\xi_k)$$

mit Integrationsgewichten  $w_k$  und Stützstellen  $\xi_k$ ,  $k = 1, \dots, q$

**Beispiel 2.2** Für  $q = 2$ :  $\xi_1 = 0$ ,  $\xi_2 = 1$ ,  $w_1 = \frac{1}{2} = w_2$

Bei unserem Beispiel:

$$\int_{T_h} f(x) \varphi_i^h dx \approx \sum_{k=1}^q w_k f(F_h(\xi_k)) \varphi_i^1(\xi_k) J$$

Entsprechend:

$$\int_{T_h} \partial_x \varphi_j^h(x) \partial_x \varphi_i^h(x) dx \approx \sum_{k=1}^q w_k \left( \partial_\xi \varphi_j^1(\xi_k) \cdot \xi_x \right) \cdot \left( \partial_\xi \varphi_i^1(\xi_k) \cdot \xi_x \right) J$$

## 2.8 Rechentechnische Betrachtungen

$$A_{ij} = \int_I \partial_x \varphi_j \partial_x \varphi_i \, dx = \sum_T \int_T \partial_x \varphi_j \partial_x \varphi_i \, dx$$

for  $i = 1$  to  $n$

for  $j = 1$  to  $n$

$$A_{ij} = \int_I \partial_x \varphi_j \partial_x \varphi_i \, dx$$

**Praxis:** Summation vertauschen

forall  $T$

for  $i = 1$  to  $n$

for  $j = 1$  to  $n$

$$A_{ij} += \int_T \partial_x \varphi_j \partial_x \varphi_i \, dx$$

Was passiert auf einer Zelle  $T$ ?

for  $k = 1$  to  $q$

Berechne Basisfunktionen auf  $T^1$  in  $\xi_k$

for  $i = 1$  to  $\text{local}_n$

for  $j = 1$  to  $\text{local}_n$

$$A_{ij} += w_k \partial_{\xi} \varphi_j^1 \xi_x \partial_{\xi} \varphi_i^1 \xi_x J$$

# 3 FEM für elliptische Probleme

## 3.1 Poisson Problem

Rechengebiet  $\Omega \subset \mathbb{R}^2$ , beschränkt (meistens  $(0, 1)^2$ ).

Betrachte die Aufgabe:

$$\min \frac{1}{2} \int_{\Omega} (\nabla u)^2 dx - \int_{\Omega} f \cdot u dx \quad (M)$$

mit  $u = u(x)$ ,  $f(x) = f$ ,  $x = (x_1, x_2) \in \Omega$ ,  $f, u : \Omega \rightarrow \mathbb{R}$ ,  $\nabla u = (\partial_{x_1} u, \partial_{x_2} u)$

Suche die Lösung von (M) in  $V = \{\varphi \mid \varphi \text{ ist stetig auf } \Omega, \partial_{x_1} \varphi, \partial_{x_2} \varphi \text{ sind stückweise stetig, beschränkt, } \varphi = 0 \text{ auf } \partial\Omega\}$

**Satz 3.1** (Greensche Formel)

Für hinreichend glatte Funktionen  $v, w$  gilt:

$$\int_{\Omega} \nabla v \nabla w dx = - \int_{\Omega} v \Delta w dx + \int_{\partial\Omega=\Gamma} v \cdot \partial_n w d\Gamma$$

mit  $\Delta = \partial_{x_1}^2 + \partial_{x_2}^2$ ,  $\partial_n w = \nabla w \cdot n$  und  $n$  bezeichne die nach aussen gerichtete normierte Normale von  $\Gamma = \partial\Omega$

**Beweis:** Benutze den Divergenzatz in 2D für vektorwertige Funktionen.  $\square$

### Klassisches Poisson-Problem

$$\begin{aligned} -\Delta u &= f \quad \text{auf } \Omega & (D) \\ u &= 0 \quad \text{auf } \Gamma = \partial\Omega \end{aligned}$$

**Satz 3.2** Analog zum 1D-Fall gilt:

Für hinreichend glatte Lösung  $u$  gilt:

$$(D) \Rightarrow (M).$$

**Beweis:** Benutze die Greensche-Formel  $\square$

Betrachte die variationelle Formulierung

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

mit  $a(v, w) = (\nabla v, \nabla w)$ ,  $v, w \in V$  und  $(v, w) := \int_{\Omega} v w dx$   $v, w \in V$ .

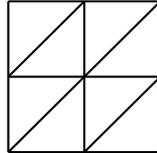
Analog zum 1D-Fall gilt:

**Satz 3.3**  $(V) \Leftrightarrow (M)$ .

**Beweis:** Variationsrechnung. □

### Finite Elemente

Triangulierung  $\mathbb{T}_h$ ,  $\Omega$  polygonal



Gitterparameter  $h = \max_{T \in \mathbb{T}_h} \text{diam}(T)$  mit  $\text{diam}(T) =$  längste Seite von  $T$ .

Diskreter Teilraum  $V_h \subset V$ :

$$V_h = \{\varphi \mid \varphi \in V, \varphi|_T \text{ ist linear für } T \in \mathbb{T}_h\}$$

Diskretisierung:

$$u_h \in V_h : a(u_h, \varphi) = (f, \varphi) \quad \forall \varphi \in V_h$$

**Satz 3.4** (Galerkin-Eigenschaft)

Es gilt

$$a(u - u_h, \varphi) = 0 \quad \forall \varphi \in V_h$$

**Beweis:**

$$\begin{aligned} a(u, \varphi) &= (f, \varphi) \quad \forall \varphi \in V \\ -a(u_h, \varphi) &= (f, \varphi) \quad \forall \varphi \in V_h \\ \hline a(u - u_h, \varphi) &= 0 \quad \forall \varphi \in V_h \end{aligned}$$

□

**Satz 3.5**  $u, u_h$  Lösungen von  $(V)$  und  $(V_h)$ , dann gilt:

$$\|\nabla(u - u_h)\| \leq \|\nabla(u - I_h u)\|$$

mit  $\|w\| = \sqrt{\int_{\Omega} w^2 dx}$  und  $I_h w$  ist definiert durch:

$w \in V : I_h w \in V_h$  und  $I_h w(x_i) = w(x_i)$  wobei  $x_i$  die Ecken aller  $T \in \mathbb{T}_h$  durchläuft.

**Satz 3.6**  $u, u_h$  Lösungen von  $(V), (V_h)$ :

$$\|\nabla(u - u_h)\| \leq ch$$

## 3.2 Natürliche und wesentliche Randbedingung

### Beispiel 3.1 Neumann-Problem

#### Klassisch

$$\begin{aligned} -\Delta u + u &= f & \text{auf } \Omega & \quad (D) \\ \frac{\partial u}{\partial n} &= g & \text{auf } \Gamma = \partial\Omega & \end{aligned}$$

$$g := g(x_1, x_2), \quad g : \mathbb{R}^2 \rightarrow \mathbb{R}$$

**Definition 3.1**  $\frac{\partial u}{\partial n} = g$  heisst Neumann-Bedingung.

$u = u_0$  heisst Dirichlet-Bedingung.

#### Variationelle Formulierung

$$a(u, \varphi) = (f, \varphi) + \int_{\Gamma} g\varphi \, d\Gamma \quad \forall \varphi \in V \quad (V)$$

$V = \{\varphi \mid \varphi \text{ ist stetig, } \partial_{x_i}\varphi \text{ stückweise stetig und beschränkt}\}$

$$\begin{aligned} a(u, \varphi) &= \int_{\Omega} \nabla u \nabla \varphi \, dx + \int_{\Omega} u\varphi \, dx \\ (f, \varphi) &= \int_{\Omega} f\varphi \, dx \end{aligned}$$

#### Minimum-Problem

$$u \in V : \min_{\varphi \in V} \frac{1}{2} a(\varphi, \varphi) - (f, \varphi) - \int_{\Gamma} g\varphi \, d\Gamma \quad (M)$$

**Satz 3.7**  $(D) \Rightarrow (V)$ .

#### Beweis:

1.  $-\Delta u + u = f$
2.  $(-\Delta u, \varphi) + (u, \varphi) = (f, \varphi) \quad \forall \varphi \in V$
3. Mit der Greenschen Formel:

$$(f, \varphi) = \int_{\Omega} \nabla u \nabla \varphi \, dx - \int_{\Gamma} \frac{\partial u}{\partial n} \varphi \, d\Gamma + \int_{\Omega} u\varphi \, dx$$

4. Benutze die Randbedingung  $\frac{\partial u}{\partial n} = g$

$$\int_{\Omega} \nabla u \nabla \varphi \, dx + \int_{\Omega} u \varphi \, dx = (f, \varphi) + \int_{\Gamma} g \varphi \, d\Gamma \quad \square$$

**Satz 3.8** Lösung  $u$  von (V) sei hinreichend glatt. Dann gilt (V)  $\Rightarrow$  (D).

**Beweis:** Greensche Formel:

$$\begin{aligned} (f, \varphi) + \int_{\Gamma} g \varphi \, d\Gamma &= a(u, \varphi) = \int_{\Gamma} \frac{\partial u}{\partial n} \varphi \, d\Gamma + \int_{\Omega} (-\Delta u + u) \varphi \, dx \\ \Leftrightarrow \int_{\Omega} (-\Delta u + u - f) \varphi \, dx + \int_{\Gamma} \left( \frac{\partial u}{\partial n} - g \right) \varphi \, d\Gamma &= 0 \quad \forall \varphi \in V \quad (3.1) \end{aligned}$$

Insbesondere gilt (3.1) für  $\bar{\varphi} \in V$  mit der zusätzlichen Bedingung  $\bar{\varphi} = 0$  auf  $\Gamma$ . Also gilt:

$$\int_{\Omega} (-\Delta u + u - f) \bar{\varphi} \, dx = 0$$

d.h.  $-\Delta u + u - f = 0$  auf  $\Omega$ . Somit reduziert sich (3.1) zu:

$$\int_{\Gamma} \left( \frac{\partial u}{\partial n} - g \right) \varphi \, d\Gamma = 0 \quad \forall \varphi \in V$$

Standardvariationsargument:  $\rightarrow \frac{\partial u}{\partial n} = g \quad \square$

**Bemerkung 3.1** Die Randbedingung  $\frac{\partial u}{\partial n} = g$  taucht in der variationellen Formulierung (V) nicht explizit auf. Sie heisst daher auch **natürliche Randbedingung**. Im Gegensatz dazu muss die Bedingung  $u = u_0$  explizit bei der Formulierung berücksichtigt werden. Sie heisst auch **wesentliche Randbedingung**.

### 3.3 Sobolev-Räume

**Bezeichnungen:**

Gebiet  $\Omega \subset \mathbb{R}^n$ , offen, stückweise glatter Rand.

$L_2(\Omega)$ : Menge aller Funktionen, deren Quadrat lebesgueintegrierbar ist.

Skalarprodukt:  $(v, w)_0 := (v, w)_{L_2} = \int_{\Omega} v(x)w(x) \, dx$

$L_2(\Omega)$  ist ein Hilbertraum mit Norm  $\|v\|_0 = \sqrt{(v, v)_0}$ .

**Definition 3.2** (Schwache Ableitung)

$u \in L_2(\Omega)$  besitzt in  $L_2(\Omega)$  die schwache Ableitung  $v = \partial^\alpha u$ , falls  $v \in L_2(\Omega)$  und  $(\varphi, v) = (-1)^{|\alpha|} (\partial^\alpha \varphi, u)_0 \quad \forall \varphi \in C_0^\infty(\Omega)$ .

Multiindex  $\alpha = (\alpha_1, \dots, \alpha_n) \quad \alpha_i \in \mathbb{N}_0, \quad |\alpha| := \sum \alpha_i$

$$\partial^\alpha = \partial_{x_1}^{\alpha_1} \partial_{x_2}^{\alpha_2} \dots \partial_{x_n}^{\alpha_n}$$

**Beispiel 3.2** In  $\mathbb{R}^2$ :  $\alpha = (1, 1): \partial^\alpha u = \partial_{x_1} \partial_{x_2} u$

$$(\varphi, \partial_{x_i} u) = -(\partial_{x_i} \varphi, u)$$

$\varphi \in C_0^\infty$ :  $\varphi \in C^\infty$  mit  $\text{supp } \varphi = \{x \in \Omega \mid \varphi(x_i) \neq 0\}$  ist kompakt in  $\Omega$  enthalten.

**Definition 3.3** (Sobolev Räume)

Sei  $m \in \mathbb{N}_0$ .

$H^m(\Omega) = \{u \in L_2(\Omega) \mid u \text{ besitzt schwache Ableitungen } \partial^\alpha u \text{ für alle } |\alpha| \leq m\}$

In  $H^m(\Omega)$  wird durch

$$(u, v)_m := \sum_{|\alpha| \leq m} (\partial^\alpha u, \partial^\alpha v)_0$$

ein Skalarprodukt definiert. Die zugehörige Norm ist

$$\|u\|_m = \sqrt{(u, u)_m} = \sqrt{\sum_{|\alpha| \leq m} \|\partial^\alpha u\|_{L_2(\Omega)}^2}$$

Man betrachtet auch:

$$|u|_m = \sqrt{\sum_{|\alpha|=m} \|\partial^\alpha u\|_0^2}$$

**Bemerkung 3.2** Mit  $\|\cdot\|_m$  ist  $H^m(\Omega)$  vollständig.

**Satz 3.9** Sei  $m \in \mathbb{N}_0$ . Dann ist  $C^\infty(\Omega) \cap H^m(\Omega)$  dicht in  $H^m(\Omega)$ .

**Definition 3.4** (Verallgemeinerung von Nullrandbedingungen)

Die Vervollständigung von  $C_0^\infty(\Omega)$  bezüglich der Sobolev-Norm  $\|\cdot\|_m$  wird mit  $H_0^m(\Omega)$  bezeichnet.

**Beispiel 3.3**

$$\begin{aligned} -\Delta u &= f && \text{auf } \Omega \\ u &= 0 && \text{auf } \Gamma \end{aligned}$$

Geeignete Wahl von  $V$ :  $V = H_0^1(\Omega)$

**Satz 3.10** (Poincaré-Ungleichung)

Sei  $\Omega$  in einem  $n$ -dimensionalen Würfel der Kantenlänge  $s$  enthalten, dann:

$$\|v\|_0 \leq s|v|_1 \quad \forall v \in H_0^1(\Omega)$$

**Beweis:** Übertragung der Idee aus Aufgabe 3.1

$$\|\partial^1 v\|_0^2 \leq \|v\|_0^2 + \|\partial^1 v\|_0^2 \leq s^2 \|\partial^1 v\|_0^2 + \|\partial^1 v\|_0^2$$

□

### 3.4 Abstrakte Formulierung

Abstrakter Rahmen:

$V$  ein Hilbertraum.  
 $(\cdot, \cdot)_V$  das zugehörige Skalarprodukt  
 $\|\cdot\|_V = \sqrt{(\cdot, \cdot)_V}$  die zugehörige Norm.  
 $a(\cdot, \cdot)$  eine Bilinearform auf  $V \times V$ .  
 $L(\cdot)$  eine Linearform auf  $V$ .

#### Variationelle Formulierung

$$u \in V : a(u, \varphi) = L(\varphi) \quad \forall \varphi \in V \quad (V) \quad (3.2)$$

#### Voraussetzungen:

V.i)  $a(\cdot, \cdot)$  ist symmetrisch.

V.ii)  $a(\cdot, \cdot)$  ist stetig:

$$|a(v, w)| \leq c \|v\|_V \|w\|_V \quad \forall v, w \in V, c > 0$$

V.iii)  $a(\cdot, \cdot)$  ist  $V$ -elliptisch:

$$a(v, v) \geq \alpha \|v\|_V^2 \quad \forall v \in V, \alpha > 0$$

V.iv)  $L(\cdot)$  ist stetig:

$$|L(v)| \leq \Lambda \|v\|_V \quad \forall v \in V, \Lambda > 0$$

#### Satz 3.11 (Existenzsatz)

Angenommen, es gelten die Bedingungen V.i)-V.iv):

$\exists!$  Lösung  $u \in V$  von (V) mit der Stabilitätsabschätzung

$$\|u\|_V \leq \frac{\Lambda}{\alpha}.$$

#### Beweis:

##### 1. Eindeutigkeit

Annahme  $\exists u_1, u_2 \in V$ ,  $u_1, u_2$  lösen (V).

$$\begin{aligned} a(u_1, \varphi) &= L(\varphi) \quad \forall \varphi \in V \\ -a(u_2, \varphi) &= L(\varphi) \quad \forall \varphi \in V \\ \hline a(u_1 - u_2, \varphi) &= 0 \quad \forall \varphi \in V \end{aligned}$$

Wähle  $\varphi = u_1 - u_2$  :

$$a(u_1 - u_2, u_1 - u_2) = 0$$

Benutze V.iii):

$$\alpha \|u_1 - u_2\|_V^2 \leq a(u_1 - u_2, u_1 - u_2) = 0$$

## 2. Existenz

Idee: Reduziere  $(V)$  auf ein Fixpunktproblem

Rieszscher Darstellungssatz (für Hilberträume)

$\exists A \in \mathcal{L}(V, V) =:$  Lineare Abbildungen, stetig, von  $V$  nach  $V$  und ein  $l \in V$ , so dass

$$a(u, v) = (Au, v)_V \quad \forall u, v \in V$$

und

$$L(v) = (l, v)_V \quad \forall v \in V$$

Betrachte  $(V)$ :

$\forall \varphi \in V$  gilt:

$$\begin{aligned} a(u, \varphi) - L(\varphi) &= 0 \\ \Leftrightarrow (Au - l, \varphi)_V &= 0 \\ \Leftrightarrow (-\rho(Au - l), \varphi)_V &= 0 \quad \forall \rho > 0 \\ \Leftrightarrow (u - \rho(Au - l) - u, \varphi)_V &= 0 \\ \Leftrightarrow u &= -\rho(Au - l) \quad \forall \rho > 0 \end{aligned}$$

Betrachte  $W_\rho : V \rightarrow V$  mit  $W_\rho(v) = v - \rho(Av - l)$ .

Abschätzung von  $\|W_\rho(v_1) - W_\rho(v_2)\|_V^2$  :

$$\begin{aligned} \|W_\rho(v_1) - W_\rho(v_2)\|_V^2 &= \|v_1 - \rho(Av_1 - l) - v_2 + \rho(Av_2 - l)\|_V^2 \\ &= \left( v_1 - \rho Av_1 - (v_2 - \rho Av_2), v_1 - v_2 - \rho(Av_1 - Av_2) \right)_V \\ &= \left( v_1 - v_2, v_1 - v_2 \right)_V - 2\rho \left( A(v_1 - v_2), v_1 - v_2 \right)_V \\ &\quad + \rho^2 \left( A(v_1 - v_2), A(v_1 - v_2) \right)_V \\ &= \|v_1 - v_2\|_V^2 - 2\rho a(v_1 - v_2, v_1 - v_2) + \rho^2 \|A(v_1 - v_2)\|_V^2 \end{aligned}$$

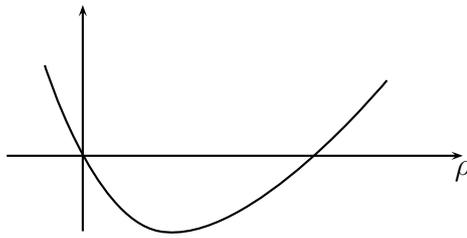
Benutze  $V$ -elliptisch:

$$\begin{aligned} \dots &\leq \|v_1 - v_2\|_V^2 - 2\rho \alpha \|v_1 - v_2\|_V^2 + \rho^2 \|A\|_V^2 \|v_1 - v_2\|_V^2 \\ &= (1 - 2\rho \alpha + \rho^2 \|A\|^2) \|v_1 - v_2\|_V^2 \end{aligned}$$

Kurvendiskussion für  $(1 - 2\rho \alpha + \rho^2 \|A\|^2)$ :

Bedingung dafür, dass  $W_\rho$  eine Kontraktionsabbildung ist:

$$1 - 2\rho \alpha + \rho^2 \|A\|^2 < 1, \text{ d.h. } p(\rho) = -\rho(2\alpha - \|A\|^2 \rho) < 0$$



$$\begin{aligned} \Rightarrow \rho > 0 \text{ und } \rho < \frac{2\alpha}{\|A\|^2} \\ \Rightarrow p(\rho) < 0 \\ \Rightarrow \text{Kontraktionsbedingung} \end{aligned}$$

Also: Mit dieser Wahl ist  $W_\rho(v) = v - \rho(Av - l)$  eine strikte Kontraktionsabbildung

→ Es gibt einen Fixpunkt.

→ Es gibt eine Lösung von (V).

### 3. Stabilitätsabschätzung

Wähle  $\varphi = u$  in (V) und benutze  $V$ -elliptisch (V.iii)) und die Stetigkeit von  $L$  (V.iv)):

$$\begin{aligned} \alpha \|u\|_V^2 \leq a(u, u) = L(u) \leq \Lambda \|u\|_V \\ \Leftrightarrow \|u\|_V \leq \frac{\Lambda}{\alpha} \end{aligned}$$

□

### Abstraktes Minimierungsproblem

Finde  $u \in V$ , so dass

$$F(u) \leq F(\varphi) \quad \forall \varphi \in V \quad (M)$$

gilt, mit  $F(\varphi) = \frac{1}{2} a(\varphi, \varphi) - L(\varphi)$ .

**Satz 3.12** (V)  $\Leftrightarrow$  (M).

## 3.5 Diskretisierung

$$\begin{aligned} u \in V : \quad a(u, \varphi) &= L(\varphi) \quad \forall \varphi \in V \\ u_h \in V_h : \quad a(u_h, \varphi) &= L(\varphi) \quad \forall \varphi \in V_h \subset V \end{aligned}$$

$$V_h = \langle \varphi_1, \dots, \varphi_n \rangle$$

$$\varphi \in V_h : \varphi = \sum_{i=1}^N \alpha_i \varphi_i, \quad \alpha_i \in \mathbb{R}$$

$$u_h \in V_h : u_h = \sum_{j=1}^N x_j \varphi_j, \quad x_j \in \mathbb{R}$$

$$a(u_h, \varphi_i) = L(\varphi_i), \quad i = 1, \dots, N$$

$$\rightarrow \sum_{j=1}^N a(\varphi_j, \varphi_i) x_j = L(\varphi_i), \quad i = 1, \dots, N$$

$\rightarrow$  Matrixform  $Ax = b$ ,  $A \in \mathbb{R}^{N \times N}$ ,  $x, b \in \mathbb{R}^N$

$$A_{ij} = a(\varphi_j, \varphi_i), \quad b_i = L(\varphi_i)$$

**Satz 3.13** *Unter den Voraussetzungen V.i),... V.iv) ist A symmetrisch und positiv definit.*

**Satz 3.14** *Es gelte V.i),..., V.iv), dann gilt:*

$$\|u_h\|_V \leq \frac{\Lambda}{\alpha}$$

**Satz 3.15** *Für den Diskretisierungsfehler gilt:*

$$\|u - u_h\|_V \leq \frac{c}{\alpha} \|u - \varphi\| \quad \forall \varphi \in V_h$$

## 3.6 Variationsungleichungen

Mit obigen Bezeichnungen betrachte das Problem

$$a(u, \varphi - u) \geq L(\varphi - u) \quad \forall \varphi \in K \subset V$$

$K$  ist abgeschlossen und konvex. Man spricht von einer elliptischen Variationsungleichung 1. Art.

**Lemma 3.1** *Sei  $K \subset V$  abgeschlossen und konvex. Dann gilt:*

$$\forall x \in V \exists! y \in K, \text{ so dass } \|x - y\| = \inf_{\varphi \in K} \|x - \varphi\|$$

*Der Punkt  $y$  heisst Projektion von  $x$  auf  $K$  :  $y = P_K(x)$ .*

**Beweis:**

1. "Es gibt ein  $y$ "

Sei  $\varphi_k$  eine Minimalfolge, d.h.

$$\lim_{k \rightarrow \infty} \|\varphi_k - x\| = d = \inf_{\varphi \in K} \|\varphi - x\|$$

Durch Ausmultiplizieren verifiziert man

$$\|\varphi_k - \varphi_l\|^2 = 2\|x - \varphi_k\|^2 + 2\|x - \varphi_l\|^2 - 4 \left\| x - \frac{1}{2}(\varphi_k + \varphi_l) \right\|^2.$$

Da  $K$  konvex:  $\frac{1}{2}(\varphi_k + \varphi_l) \in K$  ist  $d^2 \leq \left\| x - \frac{1}{2}(\varphi_k + \varphi_l) \right\|^2$ .

Zusammen:

$$\|\varphi_k - \varphi_l\|^2 \leq 2 \underbrace{\|x - \varphi_k\|^2}_{\rightarrow d^2} + 2 \underbrace{\|x - \varphi_l\|^2}_{\rightarrow d^2} - 4d^2$$

Und damit gilt:

$$\lim_{k, l \rightarrow \infty} \|\varphi_k - \varphi_l\| = 0$$

Da  $V$  vollständig ist und  $K$  abgeschlossen

$$\exists y \in K \text{ mit } \lim_{k \rightarrow \infty} \varphi_k = y.$$

Wegen der Stetigkeit der Norm gilt:

$$\|x - y\| = \lim_{k \rightarrow \infty} \|x - \varphi_k\| = d$$

2. "Eindeutigkeit von  $y$ "

Seien  $y_1, y_2 \in K$  mit

$$\|x - y_1\| = \|x - y_2\| = \inf_{\varphi \in K} \|x - \varphi\|.$$

Analog zu 1.:

$$\begin{aligned} \|y_1 - y_2\|^2 &\leq 2\|x - y_1\|^2 + 2\|x - y_2\|^2 - 4 \left\| x - \frac{1}{2}(y_1 + y_2) \right\|^2 \\ &\leq 2d^2 + 2d^2 - 4d^2 = 0 \end{aligned}$$

$$\rightarrow \|y_1 - y_2\|^2 \leq 0$$

□

**Satz 3.16** Sei  $K \subset V$  abgeschlossen und konvex, dann gilt:  $y = P_K(x)$  genau dann, wenn gilt:

$$y \in K : (y - x, \varphi - y) \geq 0 \quad \forall \varphi \in K$$

**Beweis:**

“ $\Rightarrow$ ”  $x \in V$  und  $y = P_K(x) \in K$

$K$  ist konvex  $\Rightarrow (1-t)y + t\varphi = y + t(\varphi - y) \in K, \quad 0 \leq t \leq 1$

Betrachte

$$\begin{aligned}\Phi(t) &= \|x - y - t(\varphi - y)\|^2 \\ &= \|x - y\|^2 - 2t(x - y, \varphi - y) + t^2 \|\varphi - y\|^2\end{aligned}$$

$\Phi(t)$  nimmt bei  $t = 0$  das Minimum an.

$$\begin{aligned}\Rightarrow \quad & \Phi'(0) = 0 \\ \Leftrightarrow & -2(x - y, \varphi - y) \geq 0 \\ \Leftrightarrow & (x - y, \varphi - y) \leq 0 \\ \Leftrightarrow & (y - x, \varphi - y) \geq 0\end{aligned}$$

“ $\Leftarrow$ ” Sei  $\varphi \in K$  beliebig, aber fest.

$$\begin{aligned}0 &\leq (y - x, \varphi - y) \\ &= (y - x, (\varphi - x) + (x - y)) \\ &= (y - x, x - y) + (y - x, \varphi - x) \\ &= -\|x - y\|^2 + (y - x, \varphi - x) \\ \Leftrightarrow \|x - y\|^2 &\leq (y - x, \varphi - x) \\ &\leq \|y - x\| \|\varphi - x\| \\ \Leftrightarrow \|x - y\| &\leq \|\varphi - x\| \quad \forall \varphi \in K\end{aligned}$$

□

**Korollar 3.1** Sei  $K \subset V$  abgeschlossen und konvex. Dann ist  $P_K$  nicht-expansiv, d.h. es gilt:

$$\|P_K(x) - P_K(x')\| \leq \|x - x'\| \quad \forall x, x' \in V$$

**Beweis:**

Gegeben seien  $x, x' \in V, y = P_K(x), y' = P_K(x')$ .

$$\begin{aligned}y \in K &: (y, \varphi - y) \geq (x, \varphi - y) \quad \forall \varphi \in K \\ y' \in K &: (y', \varphi - y') \geq (x', \varphi - y') \quad \forall \varphi \in K\end{aligned}$$

1. Ungleichung:  $\varphi = y'$

2. Ungleichung:  $\varphi = y$

Addiere beide Ungleichungen:

$$\begin{aligned}\|y - y'\|^2 &= (y - y', y - y') \\ &\leq (x - x', y - y') \\ &\leq \|x - x'\| \|y - y'\|\end{aligned}$$

$$\Leftrightarrow \|y - y'\| \leq \|x - x'\|$$

□

**Satz 3.17** Das Problem  $u \in K$

$$a(u, \varphi - u) \geq L(\varphi - u) \quad \forall \varphi \in K$$

hat eine eindeutige Lösung in  $V$ .

**Beweis:**

1. Eindeutigkeit:

$$\begin{aligned} a(u_1, \varphi - u_1) &\geq L(\varphi - u_1) & \forall \varphi \in K, u_1 \in K \\ a(u_2, \varphi - u_2) &\geq L(\varphi - u_2) & \forall \varphi \in K, u_2 \in K \end{aligned}$$

Testen mit  $\varphi = u_2$  bzw.  $\varphi = u_1$  und Addition:

$$\alpha \|u_1 - u_2\|^2 \leq a(u_2 - u_1, u_2 - u_1) \leq 0$$

$$\Rightarrow \|u_1 - u_2\| \leq 0$$

2. Existenz:

Benutze Riesz'schen Darstellungssatz

$$\begin{aligned} a(u, v) &= (Au, v) & \forall u, v \in V \\ L(v) &= (l, v) & \forall v \in V \\ (Au, \varphi - u) &\geq (l, \varphi - u) & \forall \varphi \in K \\ \left( -(Au - l), \varphi - u \right) &\leq 0 & \forall \varphi \in K \\ \left( (u - \rho(Au - l)) - u, \varphi - u \right) &\leq 0 & \forall \varphi \in K \end{aligned}$$

Dies ist äquivalent zu

$$u = P_K(u - \rho(Au - l))$$

Betrachte  $W_\rho : V \rightarrow V$

$$W_\rho(v) = P_K(v - \rho(Av - l))$$

Seien  $v_1, v_2 \in V$ :

$$\|W_\rho(v_1) - W_\rho(v_2)\|^2 \leq \|v_2 - v_1\|^2 + \rho^2 \|A(v_2 - v_1)\|^2 - 2\rho\alpha a(v_2 - v_1, v_2 - v_1)$$

Schliesslich:

$$\|W_\rho(v_1) - W_\rho(v_2)\|^2 \leq (1 - 2\rho\alpha + \rho^2\|A\|^2) \|v_2 - v_1\|^2$$

Polynomdivision liefert:

$W_\rho$  ist eine Kontraktion, falls  $0 < \rho < \frac{2\alpha}{\|A\|^2}$  gilt.

Kontraktion  $\rightarrow$  Fixpunkt  $\rightarrow$  Fixpunkt ist Lösung. □

### 3.7 Lineare Funktionale

**Bezeichnung:**  $X, Y$  normierte  $\mathbb{R}$ -Vektorräume.

Wir untersuchen “lineare Operatoren”.

$T : X \rightarrow Y$ , d.h. lineare, stetige Abbildungen von  $X \rightarrow Y$

**Lemma 3.2** *Ist  $T : X \rightarrow Y$  linear, so sind äquivalent:*

1.  $T$  ist stetig.
2.  $T$  ist stetig in  $x_0$  für ein  $x_0 \in X$ .
3.  $\sup_{\|x\|_X \leq 1} \|Tx\|_Y < \infty$ .
4.  $\exists c > 0$  mit  $\|Tx\|_Y \leq c\|x\|_X \quad \forall x \in X$ .

**Definition 3.5**  $L(X, Y) := \{T : X \rightarrow Y \mid T \text{ ist stetig und linear}\}$  “stetige (oder beschränkte) Operatoren”.

Operatornorm von  $T$ :

$$\|T\|_{L(X, Y)} := \sup_{\|x\|_X \leq 1} \|Tx\|_Y$$

$\|T\|_{L(X, Y)}$  ist die kleinste Zahl mit  $\|Tx\|_Y \leq \|T\|_{L(X, Y)} \|x\|_X$ .

**Definition 3.6**

1.  $X' := L(X, \mathbb{R})$  ist der “Dualraum” von  $X$ . Elemente von  $X'$  heißen “Lineare Funktionale”.
2. Für  $T \in L(X, Y)$  ist  $N(T) := \{x \in X \mid Tx = 0\}$  der Nullraum von  $T$ .

**Bemerkung 3.3**  $N(T)$  ist ein abgeschlossener Unterraum.

1. *Abgeschlossen:*

Betrachte  $x_k \rightarrow x$  für  $k \rightarrow \infty$ ,  $x_k \in N(T)$ ,  $x \in X$ .

$$\lim_{k \rightarrow \infty} T(x_k) = 0 = T(x)$$

$$\rightarrow T(x) = 0 \rightarrow x \in N(T)$$

2. *Unterraum:*

$$x_1, x_2 \in N(T)$$

$$T(x_1 + x_2) = T(x_1) + T(x_2) = 0 + 0 = 0$$

□

**Satz 3.18** (Rieszscher Darstellungssatz)

$X$  ein Hilbertraum. Betrachte:  $J : X \rightarrow X'$ ,  $x \mapsto (\cdot, x)_X$

Aussage:  $J$  ist ein linearer Isomorphismus

**Beweis:**

1.  $J$  ist linear.

$$J(x_1) = (\cdot, x_1)_X, \quad J(x_2) = (\cdot, x_2)_X$$

$$J(x_1) + J(x_2) = (\cdot, x_1)_X + (\cdot, x_2)_X = (\cdot, x_1 + x_2) = J(x_1 + x_2)$$

2.  $J(x) \in X'$ .

$$|J(x)(y)| = (y, x)_X \leq \|x\| \|y\|$$

$$\sup_{\|y\| \leq 1} |J(x)(y)| \leq \|x\| < \infty$$

$$\text{Also } \|J(x)\|_{X'} \leq \|x\|_X.$$

3.  $J$  ist injektiv.

$$\text{Betrachte } \left| J(x) \frac{x}{\|x\|} \right| = \frac{(x, x)}{\|x\|} = \|x\|.$$

Also

$$\|J(x)\|_{X'} \geq \|x\|$$

d.h.

$$x \neq 0 \Rightarrow J(x) \neq 0.$$

Also: "Nur die Null geht auf die Null".

**Randbemerkung:** Die Abbildung  $J$  ist eine Isometrie. Aus 1. und 3. folgt:  $\|J(x)\|_X = \|x\|_X$

4.  $J$  ist surjektiv:

**Beweisstruktur:** Konstruiere zu gegebenem  $0 \neq x'_0 \in X'$  ein  $w \in X$  mit  $x'_0(x) = (x, w)_X \quad \forall x \in X$ .

$P$  bezeichne die Projektion auf den abgeschlossenen Unterraum  $N(x'_0)$ .

Wähle  $e \in X$  mit  $x'_0(e) = 1$ .

Setze  $x_0 = e - Pe$ .

Es gilt:

$$x'_0(x_0) = x'_0(e) - x'_0(Pe) = 1 - 0 = 1$$

Also insbesondere  $x_0 \neq 0$ .

Wiederholung: Projektion:  $(Px - x, \varphi - Px) \geq 0 \quad \forall \varphi \in K$

Hier:  $(\tilde{y} - Pe, x_0)_X = (\tilde{y} - Pe, e - Pe) \leq 0 \quad \forall \tilde{y} \in N(x'_0)$   
 $y, Pe \in N(x'_0) \Rightarrow \tilde{y} = y + Pe, \tilde{y} = -y + Pe \in N(x'_0)$ .

Also

$$(y, x_0)_X \leq 0 \text{ und } (-y, x_0)_X \leq 0 \\ \Rightarrow (y, x_0) = 0 \quad \forall y \in N(x'_0)$$

$\forall x \in X$  ist  $x = x - x'_0(x) x_0 + x'_0(x) x_0$  wegen

$$x'_0(x - x'_0(x) x_0) = x'_0(x) - x'_0(x) \underbrace{x'_0(x_0)}_{=1} = 0$$

$$(x, x_0)_X = \left( \underbrace{x - x'_0(x) x_0}_{\in N(x'_0)} + x'_0(x) x_0, x_0 \right) \\ = (x'_0(x) x_0, x_0) \\ = x'_0(x) \|x_0\|^2$$

D.h.

$$x'_0(x) = \left( x, \frac{x_0}{\|x_0\|^2} \right)_X \\ = J \left( \frac{x_0}{\|x_0\|^2} \right) (x)$$

□

$L(v) = (l, v), a(u, v) = (Au, v) \quad a(v, w) \leq c\|v\|\|w\|$   
 $a(x, \cdot) \rightarrow$  ist lineares beschränktes Funktional  
 Riesz:  $a(x, \cdot) = (\tilde{x}, \cdot) = (Ax, \cdot)$

### 3.8 Interpolation

Wiederholung:

$$\|u - u_h\|_V \leq \frac{c}{\alpha} \|u - \varphi\|_V \quad \forall \varphi \in V_h \\ \leq \frac{c}{\alpha} \|u - I_h u\|_V$$

Interpolation in 2D für lineare Funktionen.

**Bezeichnungen:**

$$\begin{aligned} h_T &= \text{diam}(T) \text{ (längste Seite).} \\ \rho_T &= \text{Inkreisradius.} \\ h &= \max_{T \in \mathbb{T}} h_T. \end{aligned}$$

Wir betrachten Familien von Triangulierungen  $\mathbb{T}_h$ , für die unabhängig von  $h$  gilt:

$$\frac{\rho_T}{h_T} \geq \beta > 0 \quad \forall T \in \mathbb{T}_h$$

(“Die Dreiecke dürfen nicht zu dünn werden”.)

Seien  $v_i, i = 1, \dots, N$  die Knoten von  $\mathbb{T}_h$ .

Für  $u \in C^0(\bar{\Omega})$  definiere

$$I_h u(v_i) = u(v_i) \quad i = 1, \dots, N$$

und  $I_h u$  sei stückweise linear.

**Satz 3.19** Sei  $T \in \mathbb{T}_h$  ein Dreieck mit den Knoten  $a_i, i = 1, 2, 3$ .

Sei  $v \in C^0(T)$ . Die Interpolierende  $I_h v \in P_1(T)$  sei definiert durch

$$I_h v(a_i) = v(a_i), \quad i = 1, 2, 3.$$

Dann gilt:

1.  $\|v - I_h v\|_{L_\infty(T)} \leq 2(h_T)^2 \max_{|\alpha|=2} \|D^\alpha v\|_{L_\infty(T)}.$
2.  $\max_{|\alpha|=1} \|D^\alpha(v - I_h v)\|_{L_\infty(T)} \leq 6 \frac{h_T^2}{\rho_T} \max_{|\alpha|=2} \|D^\alpha v\|_{L_\infty(T)}$   
wobei gilt:  $\|v\|_{L_\infty(T)} = \max_{x \in T} |v(x)|.$

**Beweis:** Verläuft analog zu 1D-Interpolation mit längeren Taylorentwicklungen wegen 2D. □

**Satz 3.20** Unter den Voraussetzungen des vorherigen Satzes kann man zeigen:

1.  $\|v - I_h v\|_{L_2(T)} \leq c h_T^2 |v|_{H^2(T)}$
2.  $|v - I_h v|_{H^1(T)} \leq c \frac{h_T^2}{\rho_T} |v|_{H^2(T)}$

**Satz 3.21** Unter obigen Voraussetzungen gilt:

1.  $\|v - I_h v\|_{L_2(\Omega)} \leq c h^2 |v|_{H^2(\Omega)}$

$$2. |v - I_h v|_{H^1(\Omega)} \leq c \frac{h}{\beta} |v|_{H^2(\Omega)}$$

**Satz 3.22** (*Höhere Polynomansätze*)

$I_h v \in P_r(T)$  mit  $r \geq 1$ . Es gilt:

$$1. \|v - I_h v\|_{L_2(\Omega)} \leq ch^{r+1} |v|_{H^{r+1}(\Omega)}$$

$$2. |v - I_h v|_{H^1(\Omega)} \leq ch^r |v|_{H^{r+1}(\Omega)}$$

**Satz 3.23** (*Fehlende Regularität*)

$1 \leq s \leq r + 1$

$$1. \|v - I_h v\|_{L_2(\Omega)} \leq ch^s |v|_{H^s(\Omega)}$$

$$2. |v - I_h v|_{H^1(\Omega)} \leq ch^{s-1} |v|_{H^s(\Omega)}$$

# 4 Minimierungsalgorithmen, iterative Methoden

$Ax = b$ ,  $A \in \mathbb{R}^{n \times n}$  symmetrisch, positiv definit  $x, b \in \mathbb{R}^n$ .

Betrachte

$$f(x) = \frac{1}{2} x^T Ax - b^T x.$$

Minimalstelle  $\bar{x}$  von  $f(x)$  erfüllt  $A\bar{x} - b = 0$ .

## 4.1 Positiv definite Matrizen

**Bemerkung 4.1** Sei  $\|\cdot\|$  eine Norm auf  $\mathbb{C}^n$ ,  $A \in M(n, n) := \mathbb{C}^{n \times n}$ .  
 $A$  sei regulär, dann definiert

$$\|x\|_A = \|Ax\|, \quad x \in \mathbb{C}^n$$

ebenfalls eine Norm.

**Definition 4.1** Sei  $(\cdot, \cdot)$  das euklidische Skalarprodukt auf  $\mathbb{C}^n$ . Dann heisst  $A \in M(n, n)$  positiv definit, wenn  $A = A^H$  und  $(Ax, x) > 0 \quad \forall x \in \mathbb{C}^n, x \neq 0$ .

**Bemerkung 4.2**  $A \in M(n, n)$  mit  $(Ax, x) > 0 \quad \forall x \in \mathbb{C}^n, x \neq 0 \Leftrightarrow A = A^H$   
und alle Eigenwerte sind positiv.

Es gibt die Darstellung  $A = TDT^H$  mit  $T$  unitär,  $D$  diagonal.

**Definition 4.2** Sei  $A^{1/2} := TD^{1/2}T^H$ , dann heisst  $\|x\|_A := \|A^{1/2}x\|_2$ ,  
 $\|\cdot\|_2 = \sqrt{(\cdot, \cdot)}$  die Energienorm.

**Bemerkung 4.3** Für das Skalarprodukt  $(x, y)_A := (Ax, y)$ ,  $x, y \in \mathbb{C}^n$  gilt:

$$\|x\|_A = \sqrt{(Ax, x)}$$

**Bemerkung 4.4**  $A$  positiv definit  $\Leftrightarrow A^{-1}$  positiv definit.

**Definition 4.3** Die Kondition von  $A \in M(n, n)$ ,  $A$  regulär, ist definiert durch:

$$\text{cond}(A) = \|A\| \|A^{-1}\|$$

**Bemerkung 4.5** Die zugrundeliegende Vektornorm sei  $\|\cdot\|_2$ . Für  $A \in M(n, n)$ ,  $A$  positiv definit, werde  $\text{cond}(A)$  bestimmt in der zugeordneten Matrixnorm. Dann:

$$\text{cond}(A) = \frac{\lambda_{\max}}{\lambda_{\min}}$$

$\lambda_{\max}$  : grösster Eigenwert von  $A$ ,  $\lambda_{\min}$  : kleinster Eigenwert von  $A$ .

**Lemma 4.1** Sei  $A$  eine positiv definite Matrix mit Spektralkondition  $\kappa$ . Dann gilt für jeden Vektor  $x \neq 0$ :

$$\frac{(x^T A x)(x^T A^{-1} x)}{(x^T x)^2} \leq \kappa \quad (4.1)$$

**Beweis:** Anordnung der Eigenwerte:  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . Betrachte die Situation nach unitärer Transformation im Raum der Eigenvektoren.

Dann schreibt sich die linke Seite von (4.1):

$$\frac{\left(\sum_{i=1}^n \lambda_i x_i^2\right) \left(\sum_{i=1}^n \lambda_i^{-1} x_i^2\right)}{\left(\sum_{j=1}^n x_j^2\right)^2} \quad (4.2)$$

Substitution:

$$z_i = \frac{x_i^2}{\sum_{j=1}^n x_j^2}.$$

Es gilt:

$$\sum_{i=1}^n z_i = 1.$$

Einsetzen in (4.2) liefert:

$$(4.2) = \underbrace{\left(\sum_{i=1}^n \lambda_i z_i\right)}_{\leq \lambda_n} \underbrace{\left(\sum_{i=1}^n \lambda_i^{-1} z_i\right)}_{\leq \lambda_1^{-1}} \leq \frac{\lambda_n}{\lambda_1} = \kappa \quad (4.3)$$

wegen  $\sum_{i=1}^n \lambda_i z_i \leq \sum_{i=1}^n \lambda_n z_i = \lambda_n$ . Analog  $\sum_{i=1}^n \lambda_i^{-1} z_i \leq \sum_{i=1}^n \frac{1}{\lambda_1} z_i = \frac{1}{\lambda_1}$  □

## 4.2 Abstiegsverfahren

**Aufgabe:**  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  stetig differenzierbar.

Gesucht:  $\tilde{x} \in \mathbb{R}^n$ , so dass

$$f(\tilde{x}) \leq f(x) \quad \forall x \in \mathbb{R}^n$$

**Lemma 4.2** *Unter den obigen Voraussetzungen sei  $d = -\nabla f(x) \neq 0$ . Dann gilt:  $f(x + td) < f(x)$  für hinreichend kleines  $t > 0$ .*

**Beweis:** Betrachte die Richtungsableitung

$$\lim_{t \rightarrow 0} \frac{f(x + td) - f(x)}{t} = \nabla f(x)^T \cdot d < 0 \text{ wegen } d = -\nabla f(x).$$

Also gilt

$$\frac{f(x + td) - f(x)}{t} < 0$$

für hinreichend kleines  $t$ .

$\Rightarrow f(x + td) - f(x) < 0$  wegen  $t > 0$ . □

**Algorithmus 4.1** *Für  $k = 0, 1, \dots$*

1. *Berechne  $d_k = -\nabla f(x_k)$ .*
2. *Liniensuche: Man sucht für  $f$  das Minimum auf der Linie  $\{x + td_k\}$  für  $t > 0$ .*

## 4.3 Gradientenverfahren

Spezielle Aufgabe:

$$f(x) = \frac{1}{2}x^T Ax - b^T x$$

$A$  positiv definit.

Minimum falls

$$Ax - b = 0 \Leftrightarrow Ax = b.$$

Benutze Algorithmus 4.1. Hier speziell:

1.  $d_k = b - Ax_k$
2.  $t = \frac{d_k^T d_k}{d_k^T A d_k}$

Rechnung:

1.  $\nabla f(x) = Ax - b$
2. Minimumsuche für quadratisches  $f$ :

$$\begin{aligned}
 f(x_k + td_k) &= \min \\
 \Rightarrow \quad \partial_t f(x_k + td_k) &= 0 \\
 \Rightarrow \quad \nabla f(x_k + td_k) \cdot d_k &= 0 \\
 \Leftrightarrow \quad (A(x_k + td_k) - b) \cdot d_k &= 0 \\
 \Leftrightarrow \quad (\underbrace{(Ax_k - b)}_{=-d_k} + tAd_k) \cdot d_k &= 0 \\
 \Leftrightarrow \quad (-d_k + tAd_k) \cdot d_k &= 0 \\
 \Leftrightarrow \quad t(Ad_k)d_k &= d_k^T d_k \\
 \Leftrightarrow \quad t &= \frac{d_k^T d_k}{d_k^T Ad_k}
 \end{aligned}$$

□

**Lemma 4.3** Mit  $f(x) = \frac{1}{2}x^T Ax - b^T x$  gilt:

$$f(x) - f(\tilde{x}) = \frac{1}{2}\|x - \tilde{x}\|_A^2$$

wobei gilt:  $A\tilde{x} = b$

**Beweis:**

1. "Linke Seite"

$$\begin{aligned}
 f(x) - f(\tilde{x}) &= \frac{1}{2}x^T Ax - b^T x - \left( \frac{1}{2}\tilde{x} \underbrace{A\tilde{x}}_{=b} - b^T \tilde{x} \right) \\
 &= \frac{1}{2}x^T Ax - b^T x + \frac{1}{2}b^T \tilde{x}
 \end{aligned}$$

2. "Rechte Seite"

$$\begin{aligned}
 \frac{1}{2}\|x - \tilde{x}\|_A^2 &= \frac{1}{2}(x - \tilde{x})^T A(x - \tilde{x}) \\
 &= \frac{1}{2}(x - \tilde{x})^T (Ax - b) \\
 &= \frac{1}{2}x^T Ax - \frac{1}{2}\underbrace{\tilde{x}^T A}_{=b^T} x - \frac{1}{2}x^T b + \frac{1}{2}\tilde{x}^T b
 \end{aligned}$$

Vergleiche 1. und 2. → "1.=2."

□

**Lemma 4.4** Sei  $\tilde{x}$  Lösung von  $Ax = b$ , dann gilt:

$$\frac{\|x_k - \tilde{x}\|_A^2}{d_k^T A^{-1} d_k} = 1$$

**Beweis:** Es gilt:

$$d_k = b - Ax_k = A(\tilde{x} - x_k) \Leftrightarrow -A^{-1}d_k = x_k - \tilde{x}$$

Betrachte nun

$$\begin{aligned} \|x_k - \tilde{x}\|_A^2 &= (x_k - \tilde{x})^T A (x_k - \tilde{x}) \\ &= (d_k^T A^{-1}) A (A^{-1}d_k) \\ &= d_k^T A^{-1}d_k \end{aligned}$$

Nach Division folgt die Behauptung.  $\square$

**Satz 4.1** Sei  $\tilde{x} \in \mathbb{R}^n$ , so dass  $A\tilde{x} = b$ .

Für den Iterationsfehler nach  $(k+1)$ -Schritten des Gradientenverfahrens gilt:

$$\|x_{k+1} - \tilde{x}\|_A^2 \leq \|x_k - \tilde{x}\|_A^2 \left(1 - \frac{1}{\kappa}\right)$$

mit  $\kappa = \text{cond}(A)$

**Beweis:** Laut Algorithmus:

1.  $d_k = b - Ax_k$
2.  $t_k = \frac{d_k^T d_k}{d_k^T A d_k}$

Einsetzen in

$$\begin{aligned} f(x_{k+1}) &= f(x_k + t_k d_k) \\ &= \frac{1}{2} (x_k + t_k d_k)^T A (x_k + t_k d_k) - b^T (x_k + t_k d_k) \\ &= f(x_k) + t_k \underbrace{d_k^T (Ax_k - b)}_{=-d_k} + \frac{1}{2} t_k^2 d_k^T A d_k \\ &= f(x_k) - \frac{1}{2} \frac{(d_k^T d_k)^2}{d_k^T A d_k}. \end{aligned}$$

$$\Leftrightarrow f(x_{k+1}) - f(\tilde{x}) = f(x_k) - f(\tilde{x}) - \frac{1}{2} \frac{(d_k^T d_k)^2}{d_k^T A d_k}$$

Benutze Lemma 4.3:

$$\frac{1}{2} \|x_{k+1} - \tilde{x}\|_A^2 = \frac{1}{2} \|x_k - \tilde{x}\|_A^2 - \frac{1}{2} \frac{(d_k^T d_k)^2}{d_k^T A d_k} \cdot 1$$

Benutze Lemma 4.4:

$$\begin{aligned} \|x_{k+1} - \tilde{x}\|_A^1 &= \|x_k - \tilde{x}\|_A^2 - \frac{(d_k^T d_k)^2 \|x_k - \tilde{x}\|_A^1}{(d_k^T A d_k) (d_k^T A^{-1} d_k)} \\ &\leq \|x_k - \tilde{x}\|_A^2 \left(1 - \frac{(d_k^T d_k)^2}{(d_k^T A d_k) (d_k^T A^{-1} d_k)}\right) \end{aligned}$$

Benutze Lemma 4.1:

$$\|x_{k+1} - \tilde{x}\|_A^2 \leq \|x_k - \tilde{x}\|_A^2 \left(1 - \frac{1}{\kappa}\right)$$

□

## 4.4 Projiziertes Gradientenverfahren

Aufgabe: Finde  $u \in K \subset \mathbb{R}^n$ , so dass

$$\min = \frac{1}{2}u^T Au - f^T u =: J(u)$$

$A \in M(n, n)$  positiv definit,  $f \in \mathbb{R}^n$ ,  $K := \{x \in \mathbb{R}^n \mid x(i) \geq 0, i = 1, \dots, n\}$

**Algorithmus 4.2**  $P_K$  bezeichne die Projektion auf  $K$ .

1. *Initialisierung:* Wähle  $u_0 \in K$ .
2. *Iteration:* for  $k = 0, 1, \dots$

$$u_{k+1} = P_K\left(u_k - \alpha_k J'(u_k)\right), \quad \alpha_k > 0$$

mit  $J'(x) = Ax - f$ .

Schritt 2. zerfällt wie folgt:

$$\begin{aligned} u_{k+1/2} &= u_k + \alpha_k(f - Au_k) \quad (\text{Gradientenschritt}) \\ u_{k+1} &= P_K(u_{k+1/2}) \\ u_{k+1}(i) &= \max\left(0, u_{k+1/2}(i)\right) \end{aligned}$$

**Satz 4.2**  $\exists \alpha, \beta > 0$ , so dass mit  $\alpha \leq \alpha_k \leq \beta$  der Algorithmus 4.2 gegen die Lösung  $u$  konvergiert.

**Beweis:**

1. Wir zeigen:  $u = P_K\left(u - \alpha_k J'(u)\right)$

Aus Abschnitt 2.5 folgt:

$$(J'(u), \varphi - u) \geq 0 \quad \forall \varphi \in K$$

Mit  $\alpha_k > 0$  gilt:

$$\begin{aligned} & \left(\alpha_k J'(u), \varphi - u\right) \geq 0 \\ \Leftrightarrow & \left(u - u + \alpha_k J'(u), \varphi - u\right) \geq 0 \\ \Leftrightarrow & \left(u - (u - \alpha_k J'(u)), \varphi - u\right) \geq 0 \quad \forall \varphi \in K \\ \Leftrightarrow & u = P_K\left(u - \alpha_k J'(u)\right) \end{aligned}$$

2. Wir rechnen:

$$\begin{aligned} \|u_{k+1} - u\| &= \left\| P_K(u_k - \alpha_k J'(u_k)) - P_K(u - \alpha_k J'(u)) \right\| \\ &\leq \left\| u_k - u - \alpha_k (J'(u_k) - J'(u)) \right\| \end{aligned}$$

Quadrieren liefert:

$$\|u_{k+1} - u\|^2 \leq \|u_k - u\|^2 - 2\alpha_k (u_k - u, J'(u_k) - J'(u)) + \alpha_k^2 \|J'(u_k) - J'(u)\|^2$$

Nun gilt:

$$J'(u_k) - J'(u) = (Au_k - f) - (Au - f) = A(u_k - u)$$

Weiterhin ist  $A$  positiv definit:  $(u_k - u)^T A(u_k - u) \geq \lambda_{\min} \|u_k - u\|^2$

Einsetzen liefert:

$$\begin{aligned} \|u_{k+1} - u\|^2 &\leq \|u_k - u\|^2 - 2\alpha_k \lambda_{\min} \|u_k - u\|^2 + \alpha_k^2 \|A\|^2 \|u_k - u\|^2 \\ &= \|u_k - u\|^2 (1 - 2\alpha_k \lambda_{\min} + \alpha_k^2 \|A\|^2) \end{aligned}$$

Faktordiskussion  $\Rightarrow \alpha_k > 0$  und  $\alpha_k < \frac{2\lambda_{\min}}{\|A\|^2}$  □

## 4.5 Konjugiertes Gradientenverfahren (cg)

Idee bisher:  $x_{i+1} = x_i + \alpha_i d_i$ ,

$$f(x_{i+1}) = \min_{z \in \text{span}[d_i]} f(x_i + z) \text{ "eindimensionale Minimierung"}$$

Verbesserung:

$$f(x_{i+1}) = \min_{z \in \text{span}[d_{i-1}, g_i]} f(x_i + z)$$

$d_{i-1} = x_i - x_{i-1}$  "Richtung der letzten Korrektur"

$$g_i = Ax_i - b$$

**Definition 4.4** Sei  $A \in M(n, n)$  positiv definit.

Zwei Vektoren  $x, y \in \mathbb{R}^n$  heissen konjugiert oder  $A$ -orthogonal falls

$$x^T A y = 0$$

ist.

**Bemerkung 4.6**  $x_1, \dots, x_k \in \mathbb{R}^n$  paarweise konjugiert  $\Rightarrow x_1, \dots, x_k$  sind linear unabhängig.

**Rechnung**

$$\begin{aligned} \sum_{i=1}^k \alpha_i x_i &= 0 \quad | \cdot (Ax_j)^T \quad j = 1, \dots, k \\ \Leftrightarrow \sum_{i=1}^k \alpha_i x_j^T Ax_i &= 0 \\ \Leftrightarrow \underbrace{\alpha_i x_j^T Ax_j}_{>0} &= 0 \Rightarrow \alpha_j = 0, \quad j = 1, \dots, k \end{aligned}$$

Also ist der Nullvektor nur trivial darstellbar. □

**Lemma 4.5** (Konjugierte Richtungen sind gut)

Seien  $d_0, \dots, d_{n-1}$  konjugierte Richtungen. Weiterhin:  $\tilde{x} = A^{-1}b$  (die Lösung).  
Dann gilt:

$$\tilde{x} = \sum_{i=0}^{n-1} \alpha_i d_i, \quad \alpha_i = \frac{d_i^T b}{d_i^T A d_i}$$

(Lösung ist direkt hinschreibbar.)

**Beweis:** Ansatz:

$$\begin{aligned} \tilde{x} &= \sum_{k=0}^{n-1} \alpha_k d_k \quad | \cdot (Ad_i)^T, \quad i = 0, \dots, n-1 \\ \Leftrightarrow d_i^T A \tilde{x} &= \sum_{k=0}^{n-1} d_i^T A \alpha_k d_k \\ &= \alpha_i d_i^T A d_i \\ \Leftrightarrow \alpha_i &= \frac{d_i^T A \tilde{x}}{d_i^T A d_i} = \frac{d_i^T b}{d_i^T A d_i} \end{aligned}$$

□

**Lemma 4.6** (Hilfssatz über konjugierte Richtungen)

Seien  $d_0, \dots, d_{n-1}$  konjugierte Richtungen:

Für jedes  $x_0 \in \mathbb{R}^n$  liefert die durch

$$x_{i+1} = x_i + \alpha_i d_i, \quad \alpha_i = \frac{-g_i^T d_i}{d_i^T A d_i}, \quad g_i = Ax_i - b$$

für  $i \geq 0$  erzeugte Folge nach (höchstens)  $n$ -Schritten die Lösung  $x_n = A^{-1}b$ .

**Beweis:** Betrachte  $A(\tilde{x} - x_0) = (b - Ax_0)$  Benutze Lemma 4.5:

$$(\tilde{x} - x_0) = \sum_{i=0}^{n-1} \alpha_i d_i, \quad \alpha_i = \frac{d_i^T (b - Ax_0)}{d_i^T A d_i}$$

Bleibt zu zeigen:

$$-\frac{g_i^T d_i}{d_i^T A d_i} = \frac{d_i^T (b - Ax_0)}{d_i^T A d_i}$$

Rechnung:

$$\begin{aligned}\alpha_i &= \frac{-d_i^T (Ax_0 - b)}{d_i^T Ad_i} \\ &= \frac{-d_i^T (Ax_0 - Ax_i + Ax_i - b)}{d_i^T Ad_i} \\ \Leftrightarrow \alpha_i &= \frac{-d_i^T (Ax_i - b)}{d_i^T Ad_i} - \frac{d_i^T (Ax_0 - x_i)}{d_i^T Ad_i}\end{aligned}$$

Laut Algorithmus:

$$x_i = x_0 + \sum_{j=0}^{j<i} \alpha_j d_j \quad | \cdot (Ad_i)^T$$

$d_i$  ist konjugiert zu allen  $d_j$ :

$$\begin{aligned}(x_i - x_0) &= \sum_{j=0}^{j<i} \alpha_j d_j \\ \Leftrightarrow d_i^T A(x_i - x_0) &= \sum_{j=0}^{j<i} \alpha_j d_i^T Ad_j = 0\end{aligned}$$

Insgesamt also  $\alpha_i = \frac{-d_i^T g_i}{d_i^T Ad_i}$  □

**Korollar 4.1** *Unter den Voraussetzungen von Lemma 4.6 minimiert die  $k$ -te Iterierte  $x_k$  die Funktion  $f$  in  $x_0 + V_k$  mit  $V_k = \text{span}[d_0, \dots, d_{k-1}]$ . Insbesondere gilt:  $d_i^T g_k = 0$  für  $i < k$*

**Beweis:**

1. Es genügt  $d_i^T g_k = 0$ ,  $i < k$  zu zeigen

$$\begin{aligned}f(x_k) &= \min_{\alpha_i} f\left(x_0 + \sum_{i=0}^{i<k} \alpha_i d_i\right) \\ \Leftrightarrow \frac{\partial}{\partial \alpha_i} f(x_k) &= 0 \\ \Leftrightarrow \nabla f(x_k)^T \cdot d_i &= 0 \\ \Leftrightarrow (Ax_k - b)^T \cdot d_i &= 0 \\ \Leftrightarrow g_k^T d_i &= 0\end{aligned}$$

2.

$$\begin{aligned}
 0 &\stackrel{!}{=} d_k^T g_{k+1} \\
 &= d_k^T (Ax_k - b) \\
 &= d_k^T \left( A \left( x_k - \frac{g_k^T d_k}{d_k^T A d_k} d_k \right) - b \right) \\
 &= d_k^T (Ax_k - b) - \frac{d_k^T A d_k}{d_k^T A d_k} g_k^T d_k \\
 &= d_k^T g_k - g_k^T d_k \\
 &= 0
 \end{aligned}$$

3. Induktion (zz  $d_i^T g_k = 0$ ,  $i < k$ )

IA:  $k = 1$ :  $d_0^T g_1 = 0$  erfüllt wegen 2.

IV:  $d_i^T g_{k-1} = 0$ ,  $i < k - 1$

IS: Aufgrund des Algorithmus in Lemma 4.6:

$$\begin{array}{rcl}
 x_k - x_{k-1} &= & \alpha_{k-1} d_{k-1} \quad | \cdot A \\
 A(x_k - x_{k-1}) &= & \alpha_{k-1} A d_{k-1} \\
 Ax_k - b - (Ax_{k-1} - b) &= & \alpha_{k-1} A d_{k-1} \\
 g_k - g_{k-1} &= & \alpha_{k-1} A d_{k-1} \quad | \cdot d_i^T
 \end{array}$$

Also gilt für  $i < k - 1$

$$d_i^T (g_k - g_{k-1}) = 0.$$

Benutze Induktionsannahme  $d_i^T g_{k-1} = 0$ ,  $i < k - 1$ .

$$d_i g_k = 0 \quad \text{für } i < k - 1$$

Für  $i = k - 1$  benutze 2. und somit

$$d_i g_k = 0 \quad \text{für } i < k.$$

□

### Algorithmus 4.3 *cg-Verfahren*

1. *Initialisierung*:  $x_0 \in \mathbb{R}^n$  als Startwert.

Setze  $d_0 = -g_0 = b - Ax_0$ .

2. *Iteration* über  $k = 0, 1, \dots$

$$\begin{aligned}
 \alpha_k &= \frac{-g_k^T d_k}{d_k^T A d_k} \\
 x_{k+1} &= x_k + \alpha_k d_k \\
 g_{k+1} &= g_k + \alpha_k A d_k \\
 \beta_k &= \frac{g_{k+1}^T A d_k}{d_k^T A d_k} \\
 d_{k+1} &= -g_{k+1} + \beta_k d_k
 \end{aligned}$$

**Satz 4.3** (Eigenschaften des cg-Verfahrens)

Solange  $g_{k+1} \neq 0$  gelten folgende Aussagen:

1.  $d_{k-1} \neq 0$ .
2.  $V_k := \text{span}[g_0, Ag_0, A^2g_0, \dots, A^{k-1}g_0]$  "Krylov-Raum".
3.  $d_0, \dots, d_{k-1}$  sind paarweise konjugiert.
4. Es ist  $f(x_k) = \min_{z \in V_k} f(x_0 + z)$ .

**Beweis:** Induktion:

IA:  $k = 1$  ✓

IV: Satz 4.3 gelte für  $k$ .

IS: Zunächst  $g_k = g_{k-1} + \alpha_{k-1}Ad_{k-1}$ .

Wegen  $\text{span}[g_0, Ag_0, \dots, A^{k-1}g_0] = \text{span}[d_0, \dots, d_{k-1}]$  gibt es die Darstellung

$$d_{k-1} = \sum_{j=0}^{k-1} \gamma_j A^j g_0.$$

Also

$$g_k = g_{k-1} + \alpha_{k-1} \sum_{j=0}^{k-1} \gamma_j A^j g_0$$

und damit  $\text{span}[g_0, \dots, g_k] \subset V_{k+1}$ .

Nach Annahme:  $d_0, \dots, d_{k-1}$  konjugiert und wegen Optimalität von  $x_k$ :

$$d_i^T g_k = 0 \quad i < k \quad (\text{wegen Korollar 4.2})$$

Falls  $g_k \neq 0 \Rightarrow g_k$  linear unabhängig von  $(d_0, \dots, d_{k-1})$

$\Rightarrow g_k \notin V_k$ .

Also ist  $\text{span}[g_0, \dots, g_k]$  ein  $(k+1)$ -dimensionaler Raum und kein echter Unterraum von  $V_{k+1}$ .

Zusammen mit  $\text{span}[g_0, \dots, g_k] \subset V_{k+1}$  folgt die Gleichheit

$$\text{span}[g_0, Ag_0, \dots, A^k g_0] = \text{span}[g_0, \dots, g_k].$$

Zeige nun:  $V_{k+1} = \text{span}[d_0, \dots, d_k]$

Betrachte dazu:

$$\begin{aligned} g_k + d_k &= g_k - g_k + \beta_{k-1}d_{k-1} \\ \Leftrightarrow g_k + d_k &= \beta_{k-1}d_{k-1} \end{aligned}$$

Also  $g_k + d_k \in V_k$ .

$$\begin{aligned} \text{span}[g_0, \dots, g_{k-1}, g_k] &= \text{span}[g_0, \dots, g_{k-1}, d_k] \\ &= \text{span}[d_0, \dots, d_{k-1}, d_k] \end{aligned}$$

Also  $V_{k+1} = \text{span}[d_0, \dots, d_k]$ .

$d_k \neq 0$  wegen  $g_k + d_k \in V_k$ .

Zeige nun:  $d_0, \dots, d_k$  sind paarweise konjugiert.

Algorithmus:  $d_k = -g_k + \beta_{k-1}d_{k-1} \mid \cdot (Ad_i)^T$

$$\Rightarrow d_i^T Ad_k = -d_i^T Ag_k + \beta_{k-1}d_i^T Ad_{k-1}$$

a) Fall  $i < k - 1$ :

Nach Annahme:  $\beta_{k-1}d_i^T Ad_{k-1} = 0$ .

Weiterhin ist  $d_i \in V_{k-1} \Rightarrow Ad_i \in V_k$ .

$$\Rightarrow Ad_i = \sum_{j=0}^{k-1} \delta_j d_j$$

Und damit

$$\begin{aligned} d_i^T Ag_k &= (Ad_i)^T g_k \\ &= \sum_{j=0}^{k-1} \delta_j \underbrace{d_j^T g_k}_{=0} \end{aligned}$$

$$d_i^T Ag_k = 0.$$

b) Fall  $i = k - 1$ :

Wegen Wahl von  $\beta_{k-1} = \frac{g_k^T Ad_{k-1}}{d_{k-1}^T Ad_{k-1}}$  (Algorithmus) gilt:

$$\begin{aligned} d_{k-1}^T Ad_k &= -d_{k-1}^T Ag_k + \frac{g_k^T Ad_{k-1}}{d_{k-1}^T Ad_{k-1}} d_{k-1}^T Ad_{k-1} \\ &= 0 \end{aligned}$$

Zu 4. Minimaleigenschaft: Anwendung von Korollar 4.2 □

**Satz 4.4** *Unter all diesen Verfahren liefert das cg-Verfahren den kleinsten Fehler  $\|x_k - \tilde{x}\|_A$ .*

**Vorbereitung:** Sei  $p \in P_k$  (Polynome mit maximalem Grad  $k$ ),  $z \in \mathbb{R}$ ,

$p(z) = \sum_{i=0}^k \alpha_i z^i$ . Dann definiert man für  $A \in M(n, n)$ :

$$p(A) = \sum_{i=0}^k \alpha_i A^i$$

**Satz 4.5** Es gebe ein Polynom  $p \in P_k$  mit  $p(0) = 1$  und  $|p(z)| \leq r \forall z \in \sigma(A)$ .  
 $\sigma(A) =$  Menge aller Eigenwerte von  $A$ .

Dann gilt für das cg-Verfahren mit beliebigem  $x_0 \in \mathbb{R}^n$ :

$$\|x_k - \tilde{x}\|_A \leq r \|x_0 - \tilde{x}\|_A$$

**Beweis:**

1. Darstellung von  $(y - \tilde{x})$ ,  $y \in x_0 + V_k$ .

$$\text{Setze } q(z) = \frac{p(z) - 1}{z} = \sum_{i=1}^k \gamma_i z^{(i-1)}.$$

$$\rightarrow q(A) = \sum_{i=1}^k \gamma_i A^{(i-1)}$$

$$y = x_0 + \left( \sum_{i=1}^k \gamma_i A^{(i-1)} \right) g_0 \in x_0 + V_k$$

$$y = x_0 + q(A)g_0$$

Betrachte

$$\begin{aligned} y - \tilde{x} &= x_0 - \tilde{x} + y - x_0 \\ &= x_0 - \tilde{x} + q(A) g_0 \\ &= (x_0 - \tilde{x}) + \left( p(A) - 1 \right) A^{-1} A (x_0 - \tilde{x}) \end{aligned}$$

$$\Leftrightarrow y - \tilde{x} = p(A)(x_0 - \tilde{x})$$

2. Darstellung von  $\|y - \tilde{x}\|_A^2$  und  $\|x_0 - \tilde{x}\|_A^2$ .

Sei  $\{z_j\}_{j=1}^n$  ein ON-System aus Eigenvektoren zu  $A$ , d.h.  $Az_j = \lambda_j z_j$ .

Entwickle nun  $x_0 - \tilde{x} = \sum_{j=1}^n c_j z_j$ . Dann gilt:

$$\begin{aligned} y - \tilde{x} &= p(A)(x_0 - \tilde{x}) \\ &= \sum_{j=1}^n c_j p(A) z_j \\ y - \tilde{x} &= \sum_{j=1}^n c_j p(\lambda_j) z_j \end{aligned}$$

Berechne

$$\begin{aligned} \|x_0 - \tilde{x}\|_A^2 &= (x_0 - \tilde{x})^T A (x_0 - \tilde{x}) \\ &= \left( \sum_{j=1}^n c_j z_j \right)^T \left( \sum_{j=1}^n c_j \underbrace{Az_j}_{\lambda_j z_j} \right) \\ &= \sum_{j=1}^n \lambda_j |c_j|^2 \end{aligned}$$

Und entsprechend

$$\|y - \tilde{x}\|_A^2 = \sum_{j=1}^n \lambda_j |c_j p(\lambda_j)|^2$$

Abschätzung:

$$\begin{aligned} \|y - \tilde{x}\|_A^2 &= \sum_{j=1}^n \lambda_j |c_j p(\lambda_j)|^2 \\ &\leq r^2 \sum_{j=1}^n \lambda_j |c_j|^2 \\ &= r^2 \|x_0 - \tilde{x}\|_A^2 \end{aligned}$$

3. Abschluss:

$$\begin{aligned} \|y - \tilde{x}\|_A &\leq r \|x_0 - \tilde{x}\|_A, \quad y \in x_0 + V_k \\ \|x_k - \tilde{x}\|_A &\leq \|y - \tilde{x}\|_A \leq r \|x_0 - \tilde{x}\|_A \end{aligned}$$

□

**Bemerkung 4.7** (*Tschebyscheff-Polynome*)

Übliche Definition:  $T_k(x) := \cos(k \arccos x)$

Äquivalent dazu:  $T_k(x) = \frac{1}{2} \left( (x + \sqrt{x^2 - 1})^k + (x - \sqrt{x^2 - 1})^k \right)$

Eigenschaften:  $T_k(1) = 1$ ,  $|T_k(x)| \leq 1$  für  $-1 \leq x \leq 1$

Für cg-Analyse mache speziellen Ansatz:

$$p(z) := \frac{T_k \left( \frac{(b+a) - 2z}{b-a} \right)}{T_k \left( \frac{b+1}{b-a} \right)} \quad \text{mit } 0 < a < b$$

**Bemerke:**

1. Es gilt:  $p(0) = 1$ .
2. Für  $z \in [a, b]$  gilt:  $-1 \leq \frac{(b+a) - 2z}{b-a} \leq 1$   
 $\Rightarrow \left| T_k \left( \frac{(b+a) - 2z}{b-a} \right) \right| \leq 1$ .

Ab jetzt:  $a = \lambda_{\min}$ ,  $b = \lambda_{\max}$ ,  $\kappa = \frac{b}{a}$ .

**Satz 4.6** (Konvergenz des cg-Verfahrens)

Für das cg-Verfahren gilt mit beliebigem  $x_0 \in \mathbb{R}^n$ :

$$\|x_k - \tilde{x}\|_A \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - \tilde{x}\|_A$$

**Beweis:** Benutze spezielles  $p(z)$  und Satz 4.4. Damit:

$$\begin{aligned} \|x_k - \tilde{x}\|_A &\leq \max_{z \in [a,b]} p(z) \|x_0 - \tilde{x}\|_A \\ &\leq \max_{z \in [a,b]} \frac{T_k \left( \frac{b+a-2z}{b-a} \right)}{T_k \left( \frac{b+a}{b-a} \right)} \|x_0 - \tilde{x}\|_A \\ &\leq \max_{z \in [a,b]} \frac{1}{T_k \left( \frac{b+a}{b-a} \right)} \|x_0 - \tilde{x}\|_A \end{aligned}$$

Betrachte

$$T_k \left( \frac{(b+a)\frac{1}{a}}{(b-a)\frac{1}{a}} \right) = T_k \left( \frac{\kappa+1}{\kappa-1} \right).$$

Für  $z \geq 1$  gilt aufgrund der äquivalenten Definition

$$T_k(z) \geq \frac{1}{2} (z + \sqrt{z^2 - 1})^k.$$

$$\Rightarrow \|x_k - \tilde{x}\|_A \leq \frac{1}{\frac{1}{2} \left( \left( \frac{\kappa+1}{\kappa-1} \right) + \sqrt{\left( \frac{\kappa+1}{\kappa-1} \right)^2 - 1} \right)^k} \|x_0 - \tilde{x}\|_A$$

Nun gilt  $\kappa - 1 = (\sqrt{\kappa} + 1)(\sqrt{\kappa} - 1)$  und damit

$$\frac{\kappa+1}{\kappa-1} + \sqrt{\frac{(\kappa+1)^2 - (\kappa-1)^2}{(\kappa-1)^2}} = \frac{\kappa+1 + \sqrt{4\kappa}}{\kappa-1} = \frac{\sqrt{\kappa}+1}{\sqrt{\kappa}-1}.$$

Zusammen also:

$$\|x_k - \tilde{x}\|_A \leq 2 \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - \tilde{x}\|_A$$

□

## 4.6 Vorkonditionierung

### Grundidee

1. Transformiere  $Ax = b \rightarrow \tilde{A}\tilde{x} = \tilde{b}$  mit  $\text{cond}(\tilde{A}) < \text{cond}(A)$ .
2. Löse  $\tilde{A}\tilde{x} = \tilde{b}$ .
3. Rücktransformation  $\tilde{x} \rightarrow x$ .

1. Transformation

Sei  $C$  (symmetrisch) positiv definit gewählt mit  $C = HH^T$  (z.B. wegen Cholesky).

Betrachte nun  $Ax = b \Leftrightarrow H^{-1}AH^{-T}H^Tx = H^{-1}b$  wobei  $H^{-T} = (H^T)^{-1}$ .

Schreibe  $\tilde{A} = H^{-1}AH^{-T}$ ,  $\tilde{x} = H^Tx$ ,  $\tilde{b} = H^{-1}b$  und damit  $\tilde{A}\tilde{x} = \tilde{b}$ .

2. Zur Wahl von  $C$ :

Betrachte die Ähnlichkeitstransformation

$$H^{-T}\tilde{A}H^T = H^{-T}(H^{-1}AH^{-T})H^T = C^{-1}A.$$

Somit ist  $\tilde{A}$  ähnlich zu  $C^{-1}A$ .

Falls  $C = A$ , dann liegt Ähnlichkeit zu  $I$  vor, das hiesse  $\text{cond}(\tilde{A}) = 1$ .

Aber  $C$  sollte etwas mit  $A$  zu tun haben.

Einfachstes praktikables Beispiel:  $C = \text{diag}(A)$ .

Somit lautet der "ad-hoc"-Algorithmus:

Wende das cg-Verfahren an auf  $\tilde{A}\tilde{x} = \tilde{b}$ .

Besserer Algorithmus benutzt die Ähnlichkeit  $\tilde{A} \sim C^{-1}A$ .

### Algorithmus 4.4 Vorkonditioniertes cg-Verfahren

Startwert sei  $x_0 \in \mathbb{R}^n$ .

Setze  $g_0 = Ax_0 - b$ ,  $d_0 = -h_0 = -C^{-1}g_0$ .

Iteration  $k = 0, 1, 2, \dots$

$$\begin{aligned} \alpha_k &= \frac{g_k^T h_k}{d_k^T A d_k} \\ x_{k+1} &= x_k + \alpha_k d_k \\ g_{k+1} &= g_k + \alpha_k A d_k \\ h_{k+1} &= C^{-1} g_{k+1} \quad \text{Vorkonditionierung} \\ \beta_k &= \frac{g_{k+1}^T h_{k+1}}{g_k^T h_k} \\ d_{k+1} &= -h_{k+1} + \beta_k d_k \end{aligned}$$

**Satz 4.7** Für das vorkonditionierte cg-Verfahren gilt:

$$\|x_k - \tilde{x}\|_A \leq \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^k \|x_0 - \tilde{x}\|_A$$

mit  $\kappa = \text{cond}(C^{-1}A)$ .

Extremfall  $C = A$ : Nach Initialisierung schon fertig.

# 5 Adaptivität

## 5.1 Laplace-Problem

$V = H_0^1(\Omega)$ , Triangulierung  $\mathbb{T}_h$  aus Dreiecken,  $V_h \subset V$ , linearer Ansatz

$$\begin{aligned} u \in V \quad (\nabla u, \nabla \varphi) &= (f, \varphi) \quad \forall \varphi \in V \\ u_h \in V_h \quad (\nabla u_h, \nabla \varphi) &= (f, \varphi) \quad \forall \varphi \in V_h \end{aligned} \quad (5.1)$$

**Satz 5.1** (Energiefehlerschätzer)

Für den Diskretisierungsfehler  $e = u - u_h$  in (5.1) gilt die a-posteriori Abschätzung

$$\|\nabla e\|^2 \leq C \sum_T h_T^2 \rho_{1,T}^2 + h_T \rho_{2,T}^2$$

mit

$$\begin{aligned} \rho_{1,T} &= \|f + \Delta u_h\|_T \quad (\text{Gebietsresiduum}) \\ \rho_{2,T} &= \sum_{j=1}^3 \int_{\partial T_j} \frac{1}{2} [\partial_n u_h]_{T_j} d\Gamma \end{aligned}$$

und  $x \in \partial T_j$ :  $[\partial_n u_h](x) := \partial_n u_h(x - \varepsilon u) - \partial_n u_h(x + \varepsilon u)$  für  $\varepsilon \rightarrow 0$ .  
("Sprungterme")

**Beweis:**

1. Rechnung

$$\begin{aligned} \|\nabla e\|^2 &= (\nabla u - \nabla u_h, \nabla e) \\ &= (\nabla u - \nabla u_h, \nabla e - \nabla I_h e) \quad \text{Galerkin-Orthogonalität} \\ &= (f, e - I_h e) - \sum_T \left( \nabla u_h, \nabla (e - I_h e) \right)_T \\ &= (f, e - I_h e) - \sum_T \left( (-\Delta u_h, e - I_h e)_T + \int_{\partial T} \partial_n u_h (e - I_h e) d\Gamma \right) \\ &= \sum_T (f + \Delta u_h, e - I_h e)_T - \sum_T \sum_{j=1}^3 \int_{\partial T_j} \frac{1}{2} [\partial_n u_h] (e - I_h e) d\Gamma \end{aligned}$$

2. Wiederholung

Interpolation: Satz 3.23 (linearer Ansatz)

$$\|v - I_h v\|_{L_2(T)} \leq ch_T |v|_{H^1(T)}$$

Spursatz:

$$\int_{\Gamma_i} v^2 d\Gamma \leq \frac{2}{r} \|v\|_0^2 + r |v|_{H^1}^2$$

Hier:  $r \sim h_T$

3. Abschätzung

$$\|\nabla e\|^2 \leq \sum_T \|f + \Delta u_h\| ch_T \|\nabla e\|_T + \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T} \|e - I_h e\|_{\partial T}$$

Anwendung des Spursatzes auf  $\|e - I_h e\|_{\partial T}^2$ :

$$\begin{aligned} \|e - I_h e\|_{\partial T}^2 &\leq \frac{2}{h_T} \|e - I_h e\|_T^2 + h_T \|\nabla(e - I_h e)\|_T^2 \\ &\leq \frac{2}{h_T} ch_T^2 \|\nabla e\|_T^2 + h_T \|\nabla e\|_T^2 \\ &\leq ch_T \|\nabla e\|_T^2 \end{aligned}$$

Somit also:

$$\begin{aligned} \|\nabla e\|^2 &\leq \sum_T \|f + \Delta u_h\|_T h_T \|\nabla e\|_T + \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T} c\sqrt{h_T} \|\nabla e\|_T \\ &\leq c \left( \sum_T h_T^2 \|f + \Delta u_h\|_T^2 \right)^{1/2} + \left( \sum_T \|\nabla e\|_T^2 \right)^{1/2} \\ &\quad + c \left( \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T}^2 h_T \right)^{1/2} \left( \sum_T \|\nabla e\|_T^2 \right)^{1/2} \end{aligned}$$

$$\|\nabla e\| \leq c \left( \sum_T h_T^2 \|f + \Delta u_h\|_T^2 \right)^{1/2} + c \left( \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T}^2 h_T \right)^{1/2}$$

$$\|\nabla e\|^2 \leq c \sum_T h_T^2 \|f + \Delta u_h\|_T^2 + c \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T}^2 h_T$$

Wobei benutzt wird:  $ab \leq \frac{1}{2}a^2 + \frac{1}{2}b^2$  □

## 5.2 Hindernisproblem

Voraussetzungen wie in 5.1

Zusätzlich:  $K = \{\varphi \in V \mid \varphi \geq 0\}$

$$\begin{aligned} u \in K &: (\nabla u, \nabla(\varphi - u)) \geq (f, \varphi - u) \quad \forall \varphi \in K \\ u_h \in K_h &: (\nabla u_h, \nabla(\varphi - u_h)) \geq (f, \varphi - u_h) \quad \forall \varphi \in K_h := K \cap V_h \end{aligned} \quad (5.2)$$

**Lemma 5.1** Sei  $a(v, w) := (\nabla v, \nabla w)$  und  $e_i = I_h e = I_h(u - u_h)$ . Dann gilt:

$$a(u - u_h, e_i) \leq a(u, e_i - e) - f(e_i - e)$$

**Beweis:**

$$\begin{aligned} a(u - u_h, e_i) &= (f, e_i) - a(u_h, e_i) + a(u, e_i - e) - (f, e_i - e) + a(u, e) - (f, e) \\ &= (f, u_i - u_h) - a(u_h, u_i - u_h) \leq 0 \\ &= a(u, u - u_h) - (f, u - u_h) \\ &= -a(u, u_h - u) + (f, u_h - u) \leq 0 \end{aligned}$$

Also  $a(u - u_h, e_i) \leq a(u, e_i - e) - (f, e_i - e)$ . □

**Satz 5.2** Für den Diskretisierungsfehler  $e = u - u_h$  in (5.2) gilt die a-posteriori-Abschätzung

$$\|\nabla e\|^2 \leq c \sum_T (h_T^2 \rho_{1,T}^2 + h_T \rho_{2,T}^2)$$

mit Bezeichnungen wie in Satz 5.1.

**Beweis:**

$$\begin{aligned} (\nabla u - \nabla u_h, \nabla e) &= (\nabla u - \nabla u_h, \nabla(e - e_i)) + (\nabla u - \nabla u_h, \nabla e_i) \\ &\leq (\nabla u - \nabla u_h, \nabla(e - e_i)) + (\nabla u, \nabla(e_i - e)) - (f, e_i - e) \\ &= -(\nabla u_h, \nabla(e - e_i)) + (f, e - e_i) \\ &= (f, e - e_i) - (\nabla u_h, \nabla(e - e_i)) \end{aligned}$$

Jetzt geht's weiter wie im Beweis zu Satz 5.1. □

## 5.3 Hindernisproblem in Lagrange-Formulierung

Definiere zunächst das Lagrange-Funktional

$$\mathcal{L}(\varphi, w) = \frac{1}{2} a(\varphi, \varphi) - (f, \varphi) + \int_{\Omega} w(\psi - \varphi).$$

$\psi$  bezeichne das Hindernis bei der Bedingung  $u \geq \psi$ ,  $\psi : \Omega \rightarrow \mathbb{R}$ .  
 $\varphi \in V = H_0^1(\Omega)$ ,  $w \in \Lambda = \{q \in L^2 \mid q \geq 0\}$ .

**Definition 5.1**  $(u, \lambda) \in V \times \Lambda$  heisst *Sattelpunkt* von  $\mathcal{L}$ , falls für  $(\varphi, w) \in V \times \Lambda$  gilt:

$$\mathcal{L}(u, w) \leq \mathcal{L}(u, \lambda) \leq \mathcal{L}(\varphi, \lambda)$$

*Minimum bzgl  $u$ , Maximum bzgl  $\lambda$ .*

**Lemma 5.2** Sei  $(u, \lambda) \in V \times \Lambda$  *Sattelpunkt* von  $\mathcal{L}$ , dann

$$\begin{aligned} a(u, \varphi) - (\lambda, \varphi) &= (f, \varphi) \quad \forall \varphi \in V \\ (u, w - \lambda) &\geq (\psi, w - \lambda) \quad \forall w \in \Lambda. \end{aligned}$$

**Beweis:**

1.

$$\begin{aligned} \frac{1}{2} a(u, u) - (f, u) + \int_{\Omega} \lambda(\psi - u) &\leq \frac{1}{2} a(\varphi, \varphi) - (f, \varphi) + \int_{\Omega} \lambda(\psi - \varphi) \\ \Rightarrow \frac{1}{2} a(u, u) - (g, u) &\leq \frac{1}{2} a(\varphi, \varphi) - (g, \varphi) \quad \forall \varphi \in V. \end{aligned}$$

Über Variationsrechnung:

$$\begin{aligned} a(u, \varphi) &= (g, \varphi) \quad \forall \varphi \in V \\ &= (f + \lambda, \varphi) \end{aligned}$$

$$\Leftrightarrow a(u, \varphi) - (\lambda, \varphi) = (f, \varphi) \quad \forall \varphi \in V$$

2.

$$\begin{aligned} \int_{\Omega} w(\psi - u) &\leq \int_{\Omega} \lambda(\psi - u) \\ \Leftrightarrow \int_{\Omega} (w - \lambda)u &\geq \int_{\Omega} (w - \lambda)\psi \end{aligned}$$

□

**Lemma 5.3** Sei  $(u, \lambda)$  *Sattelpunkt* von  $\mathcal{L}$ , dann ist  $u$  Lösung von (5.2).

**Beweis:**

1.  $u$  fest,  $w$  variiert.

$$\int_{\Omega} w(\psi - u) \leq \int_{\Omega} \lambda(\psi - u)$$

Wähle  $w := y + \lambda$  mit  $y \geq 0$ , d.h.  $w \geq 0$ .

Einsetzen liefert:

$$\int_{\Omega} y(\psi - u) \leq 0$$

$$\Rightarrow \psi - u \leq 0 \Leftrightarrow u \geq \psi \Leftrightarrow u \in K$$

2. Testen mit  $w = 0$  und  $w = 2\lambda$ .

a)  $0 \leq \int_{\Omega} \lambda(\psi - u).$

b)  $\int_{\Omega} 2\lambda(\psi - u) \leq \int_{\Omega} \lambda(\psi - u)$

$\Rightarrow \int_{\Omega} \lambda(\psi - u) \leq 0$

$\Rightarrow \int_{\Omega} \lambda(\psi - u) = 0$  (Komplementarität)

3.  $\lambda$  fest,  $\varphi$  variiert mit  $\varphi \geq \psi$ .

$$\frac{1}{2} a(u, u) - (f, u) + \int_{\Omega} \lambda(\psi - u) \leq \frac{1}{2} a(\varphi, \varphi) - (f, \varphi) + \int_{\Omega} \lambda(\psi - \varphi)$$

$$\Rightarrow \frac{1}{2} a(u, u) - (f, u) \leq \frac{1}{2} a(\varphi, \varphi) - (f, \varphi) \quad \forall \varphi \geq \psi$$

Letzteres ist äquivalent zu (5.2)

$$a(u, \varphi - u) \geq (f, \varphi - u)$$

□

Sattelpunkt  $(u, \lambda)$  ist charakterisiert durch:

$$\begin{aligned} a(u, \varphi) - (\lambda, \varphi) &= (f, \varphi) \\ (u, w - \lambda) &\geq (\psi, w - \lambda) \end{aligned} \tag{5.3}$$

### Diskretisierung

Wähle  $V_h \subset V$ ,  $L_h \subset L_2$ ,  $\Lambda_h = L_h \cap \Lambda$ .

**Beispiel 5.1**  $V_h$  stückweise linear, stetig.

$\Lambda_h$  zellweise konstant.

Diskrete Version zu (5.3)

$$\begin{aligned} a(u_h, \varphi) - (\lambda_h, \varphi) &= (f, \varphi) & \forall \varphi \in V_h \\ (u_h, w - \lambda_h) &\geq (\psi, w - \lambda_h) & \forall w \in \Lambda_h \end{aligned} \tag{5.4}$$

Definiere:  $\delta := \max\{0, \psi - u_h\}$

**Lemma 5.4** *Es gilt:*

$$\left( \lambda_h - \lambda, u - (u_h + \delta) \right) - \left( \lambda_h, \psi - (u_h + \delta) \right) \geq 0.$$

**Beweis:**

$$\begin{aligned}
 \underbrace{(\lambda_h - \lambda, u - \psi)}_{\geq 0} &+ (\lambda_h - \lambda, \psi - (u_h + \delta)) - (\lambda_h, \psi - (u_h + \delta)) \\
 &\geq (-\lambda, \psi - (u_h + \delta)) \\
 &= \int_{\Omega} \underbrace{(-\lambda)}_{\leq 0} \underbrace{(\psi - (u_h + \delta))}_{\leq 0} \\
 &\geq 0
 \end{aligned}$$

□

**Orthogonalität**

Teste mit  $\varphi_h \in V_h$  in (5.3) und (5.4) und “(5.3)–(5.4)”.

$$a(u - u_h, \varphi_h) - (\lambda - \lambda_h, \varphi_h) = 0$$

**Rechnung**

$$\begin{aligned}
 a(u - u_h, u - u_h) &\leq a(u - u_h, e) - (\lambda - \lambda_h, u - (u_h + \delta)) + \\
 &\quad - (\lambda_h, \psi - (u_h + \delta)) \\
 &= a(u - u_h, e) - (\lambda - \lambda_h, e) - (\lambda - \lambda_h, -\delta) + \\
 &\quad - (\lambda_h, \psi - (u_h + \delta)) \\
 &= a(u - u_h, e - e_i) - (\lambda - \lambda_h, e - e_i) - (\lambda - \lambda_h, -\delta) + \\
 &\quad - (\lambda_h, \psi - (u_h + \delta)) \\
 &= (f + \lambda_h, e - e_i) - a(u_h, e - e_i) - (\lambda - \lambda_h, -\delta) + \\
 &\quad - (\lambda_h, \psi - (u_h + \delta))
 \end{aligned}$$

Es bleibt

$$-(\lambda, -\delta) - (\lambda_h, \psi - u_h).$$

Zusammen:

$$\|\nabla e\|^2 \leq (f + \lambda_h, e - e_i) - a(u_h, e - e_i) - (\lambda_h, \psi - u_h) + (\lambda, \delta)$$

□

## 5.4 Sattelpunktsuche

Aufgabe:

$$\min = \frac{1}{2} u^T A u - f^T u$$

unter den Nebenbedingungen

$$B u \leq g.$$

$u, f \in \mathbb{R}^n$ ,  $g \in \mathbb{R}^m$ ,  $A \in M(n, n)$  positiv definit  $B \in M(m, n)$ .

Schreibweise:

$$\begin{aligned} J(u) &= \frac{1}{2}u^T A u - f^T u \\ \Phi(\varphi) &= (B\varphi - g) \end{aligned}$$

**Bemerkung 5.1** Die dem Ausgangsproblem zugeordnete Lagrange-Funktion lautet

$$\mathcal{L}(\varphi, w) := J(\varphi) + w^T \Phi(\varphi), \quad \varphi \in \mathbb{R}^n, w \in \mathbb{R}_+^m.$$

**Wiederholung:**  $(u, \lambda) \in \mathbb{R}^n \times \mathbb{R}_+^m$  heisst Sattelpunkt, falls gilt:

$$\mathcal{L}(u, q) \leq \mathcal{L}(u, \lambda) \leq \mathcal{L}(\varphi, \lambda) \quad \forall \varphi \in \mathbb{R}^n, q \in \mathbb{R}_+^m$$

**Wiederholung:**  $(u, \lambda)$  Sattelpunkt von  $\mathcal{L}$ : Dann ist  $u$  Lösung der Ausgangsaufgabe.

**Algorithmus 5.1** Sattelpunktsuche

1. Start: Wähle  $\lambda_0 \in \mathbb{R}_+^m$ .
2. Für  $k = 0, 1, 2$  definiere  $u_k$  durch

$$\mathcal{L}(u_k, \lambda_k) \leq \mathcal{L}(\varphi, \lambda_k) \quad \forall \varphi \in \mathbb{R}^n.$$

$$\text{D.h. } Au_k = f - B^T \lambda_k.$$

3.  $\lambda_{k+1} = P_{\mathbb{R}_+^m}(\lambda_k + \alpha_k \Phi(u_k))$ ,  $\alpha_k \geq 0$ .

**Lemma 5.5** Sei  $q \in \mathbb{R}_+^m$  gegeben. Sei  $v \in \mathbb{R}^n$  definiert durch

$$\mathcal{L}(v, q) \leq \mathcal{L}(\varphi, q) \quad \forall \varphi \in \mathbb{R}^n.$$

Dann gilt:

$$J'(v)(\varphi - v) + q^T (\Phi(\varphi) - \Phi(v)) \geq 0 \quad \forall \varphi \in \mathbb{R}^n$$

**Beweis:** Da  $q$  fest, betrachte  $\mathcal{L}(v) = \mathcal{L}(v, q)$ .

Also  $\mathcal{L}(v) \leq \mathcal{L}(\varphi) \quad \forall \varphi \in \mathbb{R}^n$

$$\Rightarrow \mathcal{L}'(v)(\varphi - v) \geq 0 \quad \forall \varphi \in \mathbb{R}^n$$

Nun gilt:  $\mathcal{L}'(v) = J'(v) + q^T \Phi'(v)$  und  $q^T \Phi(v) = q^T (Bv - g)$

$$\Rightarrow q^T \Phi'(v) = q^T B$$

Einsetzen liefert:

$$\begin{aligned}
 J'(v)(\varphi - v) + q^T B(\varphi - v) &\geq 0 \quad \forall \varphi \in \mathbb{R}^n \\
 \Leftrightarrow J'(v)(\varphi - v) + q^T \left( \underbrace{(B\varphi - g)}_{\Phi(\varphi)} - \underbrace{(Bv - g)}_{\Phi(v)} \right) &\geq 0
 \end{aligned}$$

□

**Satz 5.3**  $\exists \alpha, \beta > 0$ , so dass mit  $\alpha \leq \alpha_k \leq \beta$  der Algorithmus gegen die Lösung  $u$  konvergiert.

**Beweis:**

1. Algorithmus:  $\mathcal{L}(u_k, \lambda_k) \leq \mathcal{L}(\varphi, \lambda_k) \quad \forall \varphi \in \mathbb{R}^n$

Sattelpunkt:  $\mathcal{L}(u, \lambda) \leq \mathcal{L}(\varphi, \lambda) \quad \forall \varphi \in \mathbb{R}^n$

Benutze Lemma 5.5:

$$\frac{1}{2}(Au_k, \varphi - u_k) - (f, \varphi - u_k) + \left( \lambda_k, \Phi(\varphi) - \Phi(u_k) \right) \geq 0 \quad (5.5)$$

$$\frac{1}{2}(Au, \varphi - u) - (f, \varphi - u) + \left( \lambda, \Phi(\varphi) - \Phi(u) \right) \geq 0 \quad (5.6)$$

Testen mit  $\varphi = u$  in (5.5) und  $\varphi = u_k$  in (5.6) und Addition:

$$\frac{1}{2}(A(u_k - u), u_k - u) + \left( \lambda_k - \lambda, \Phi(u_k) - \Phi(u) \right) \leq 0$$

2. Zeige:  $\lambda = P_{\mathbb{R}_+^m}(\lambda + \bar{\alpha}\Phi(u))$  für  $\bar{\alpha} > 0$

Sattelpunkt:  $\mathcal{L}(u, q) \leq \mathcal{L}(u, \lambda) \quad \forall q \in \mathbb{R}_+^m$

$$\begin{aligned}
 \Leftrightarrow J(u) + q^T \Phi(u) &\leq J(u) + \lambda^T \Phi(u) \\
 \Rightarrow \left( (q - \lambda), \Phi(u) \right) &\leq 0 \\
 \left( (q - \lambda), -\bar{\alpha}\Phi(u) \right) &\geq 0 \quad \bar{\alpha} > 0 \\
 \Leftrightarrow \left( \lambda - (\lambda + \bar{\alpha}\Phi(u)), q - \lambda \right) &\geq 0
 \end{aligned}$$

Benutze Satz 3.16 (Charakterisierung der Projektion über Ungleichung).

$$\Rightarrow \lambda = P_{\mathbb{R}_+^m}(\lambda + \bar{\alpha}\Phi(u))$$

3. Algorithmus:  $\lambda_{k+1} = P_{\mathbb{R}_+^m}(\lambda_k + \alpha_k \Phi(u_k))$

Aus 2.:  $\lambda = P_{\mathbb{R}_+^m}(\lambda + \alpha_k \Phi(u))$

Zusammen:

$$\lambda_{k+1} - \lambda = P_{\mathbb{R}_+^m}(\lambda_k + \alpha_k \Phi(u_k)) - P_{\mathbb{R}_+^m}(\lambda + \alpha_k \Phi(u))$$

Betrachte die Norm und benutze das  $P_{\mathbb{R}_+^m}$  nicht expansiv ist.

$$\begin{aligned} \|\lambda_{k+1} - \lambda\|^2 &\leq \left\| (\lambda_k - \lambda) + \alpha_k (\Phi(u_k) - \Phi(u)) \right\|^2 \\ &\leq \|\lambda_k - \lambda\|^2 + 2\alpha_k (\lambda_k - \lambda, \Phi(u_k) - \Phi(u)) + \\ &\quad + \alpha_k^2 \|\Phi(u_k) - \Phi(u)\|^2 \end{aligned}$$

Benutze 2. und  $L = \|B\|$ .

$$\begin{aligned} \dots &\leq \|\lambda_k - \lambda\|^2 - \alpha_k (A(u_k - u), u_k - u) + \alpha_k^2 L^2 \|u_k - u\|^2 \\ &\leq \|\lambda_k - \lambda\|^2 - \alpha_k \lambda_{\min} (u_k - u, u_k - u) + \alpha_k^2 L^2 \|u_k - u\|^2 \\ &= \|\lambda_k - \lambda\|^2 - (\lambda_{\min} \alpha_k - \alpha_k^2 L^2) \|u_k - u\|^2 \end{aligned}$$

Wähle  $\alpha_k > 0$ , so dass  $\alpha_k (\lambda_{\min} - \alpha_k L^2) \geq \gamma > 0$ .

Dann gilt:

$$\|\lambda_{k+1} - \lambda\|^2 + \gamma \|u - u_k\|^2 \leq \|\lambda_k - \lambda\|^2$$

D.h.  $\|\lambda_{k+1} - \lambda\|$  ist monoton fallend und  $\|\lambda_{k+1} - \lambda\| \geq 0$ .

$$\lim_{k \rightarrow \infty} \|\lambda_{k+1} - \lambda\| = l$$

Also

$$\underbrace{\lim_{k \rightarrow \infty} \|\lambda_{k+1} - \lambda\|^2}_{=l} + \lim_{k \rightarrow \infty} \gamma \|u - u_k\|^2 \leq \underbrace{\lim_{k \rightarrow \infty} \|\lambda_k - \lambda\|^2}_{=l}$$

$$\Rightarrow \lim_{k \rightarrow \infty} \gamma \|u_k - u\|^2 = 0$$

□

## 5.5 Dualitätsargument

Situation: Laplace-Problem auf  $\Omega$  mit Nullrandwerten, FE mit linearen Ansätzen  
Fehlerabschätzungen:

$$\|u - u_h\|_{H^1(\Omega)} \leq ch |u|_{H^2(\Omega)}$$

und suboptimal

$$\|u - u_h\|_{L^2(\Omega)} \leq ch |u|_{H^2(\Omega)}$$

Andererseits Interpolationsresultat:

$$\|u - I_h u\|_{L^2(\Omega)} \leq ch^2 |u|_{H^2(\Omega)}$$

### 5.5.1 A Priori Abschätzung

**Satz 5.4** Sei  $\Omega$  ein konvexes, polygonales Gebiet.

$$\begin{aligned} u : \quad (\nabla u, \nabla \varphi) &= (f, \varphi) \quad \forall \varphi \in V \\ u_h \quad (\nabla u_h, \nabla \varphi) &= (f, \varphi) \quad \forall \varphi \in V_h \end{aligned}$$

Dann gibt es  $c > 0$  unabhängig von  $u$  und  $h$ , so dass gilt:

$$\|u - u_h\|_{L^2(\Omega)} \leq ch^2 |u|_{H^2(\Omega)}$$

**Beweis:** Betrachte das folgende Hilfsproblem (“Duales Problem”)

$$\begin{aligned} -\Delta z &= e \quad \text{auf } \Omega \\ z &= 0 \quad \text{auf } \partial\Omega \end{aligned}$$

**Bemerkung 5.2** (Stabilität des dualen Problems)

Man kann zeigen, falls  $\Omega$  konvex ist, dass gilt:

$$\|z\|_{H^2(\Omega)} \leq C_s \|e\|_{L^2(\Omega)}$$

und  $C_s > 0$  ist unabhängig von  $e$ .

Schreibe das duale Problem in variationeller Formulierung:

$$\begin{aligned} -(\varphi, \Delta z) &= (\varphi, e) \quad \forall \varphi \in V \\ \Leftrightarrow (\nabla \varphi, \nabla z) &= (\varphi, e) \quad \forall \varphi \in V \end{aligned}$$

Wähle spezielle Funktion  $\varphi = e$ .

$$\begin{aligned} (e, e) &= (\nabla e, \nabla z) \\ &= (\nabla e, \nabla z - \nabla I_h z) \\ &\leq \|\nabla e\| \|\nabla z - \nabla I_h z\| \\ &\leq ch |u|_{H^2(\Omega)} ch |z|_{H^2(\Omega)} \\ &\leq ch |u|_{H^2(\Omega)} ch C_s \|e\|_{L^2(\Omega)} \end{aligned}$$

$$\Rightarrow \|e\|_{L^2(\Omega)} \leq Ch^2 |u|_{H^2(\Omega)} \quad \square$$

## 5.5.2 A posteriori Abschätzung

DWR-Methode (Dual Weighted Residual)

Ausgangspunkt: Duales Problem in der Form

$$z \in V : (\nabla\varphi, \nabla z) = J(\varphi).$$

$J$  ein lineares Funktional.

Wähle  $\varphi = e$ :

$$\begin{aligned} J(e) &= (\nabla e, \nabla z) \\ &= (\nabla e, \nabla z - \nabla I_h z) \\ &= (\nabla u - \nabla u_h, \nabla z - \nabla I_h z) \\ &= (f, z - I_h z) - (\nabla u_h, \nabla z - \nabla I_h z) \\ &= (f, z - I_h z) - \sum_T (\nabla u_h, \nabla z - \nabla I_h z)_T \\ &= (f, z - I_h z) - \sum_T \left[ (-\Delta u_h, z - I_h z)_T + \int_{\partial\Omega} \frac{1}{2} [\partial_n u_h] (z - I_h z) d\Gamma \right] \\ &\leq \sum_T \|f + \Delta u_h\|_T \|z - I_h z\|_T + \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T} \|z - I_h z\|_{\partial T} \end{aligned}$$

**Wiederholung:** Spursatz

$$\begin{aligned} \|z - I_h z\|_{\partial T}^2 &\leq \frac{2}{h_T} \|z - I_h z\|_T^2 + h_T \|\nabla(z - I_h z)\|_T^2 \\ &\leq \frac{2}{h_T} ch_T^4 \|z\|_{H^2(T)}^2 + h_T h_T^2 \|z\|_{H^2(T)}^2 \\ &\leq ch_T^3 \|z\|_{H^2(T)}^2 \end{aligned}$$

Also

$$J(e) \leq \sum_T \|f + \Delta u_h\| ch_T^2 \|z\|_{H^2(T)} + \sum_T \left\| \frac{1}{2} [\partial_n u_h] \right\|_{\partial T} ch_T^{3/2} \|z\|_{H^2(T)}$$

Wir haben gezeigt:

**Satz 5.5** *DWR-Abschätzung*

$$J(e) \leq \sum_T ch_T^2 \left( \|f + \Delta u_h\|_T + \frac{\|\partial_n u_h\|_{\partial T}}{2\sqrt{h_T}} \right) \|z\|_{H^2(T)}$$

**Beispiel 5.2** *Für das Fehlerfunktional  $J$*

1.  $J(\varphi) = \delta_{x_0}(\varphi)$  "Dirac". Also

$$J(e) = \delta_{x_0}(e) = e(x_0) = u(x_0) - u_h(x_0)$$

2.  $J(\varphi) = \int_{\Gamma} \varphi \, d\Gamma$  "Mittelwert. Also

$$J(e) = \int_{\Gamma} e \, d\Gamma = \int_{\Gamma} u \, d\Gamma - \int_{\Gamma} u_h \, d\Gamma$$

$$J(\varphi) = (\nabla\varphi, \nabla z) \leftarrow z = ?$$

Satz 5.5 in der Praxis:

- $z$  ist unbekannt.
- Löse  $(\nabla\varphi, \nabla z) = J(\varphi)$  numerisch, z.B. mit FE  $\rightarrow z \approx \tilde{z}$ .
- $\|z\|_{H^2(T)} \approx \|\tilde{z}\|_{H^2(T)}$

Typischerweise berechnet man  $\tilde{z}$  auf  $\mathbb{T}$ .

# 6 Parabolische Probleme

## Modellbeispiel

$$\partial_t u - \Delta u = f \quad \text{auf } \Omega \times I$$

$I$  ist das Zeitintervall  $[0, T]$ ,  $u = u(t, x)$ .

Randbedingungen und Anfangsbedingungen:

$$u(t, x) = 0 \quad t \in I, x \in \partial\Omega$$

$$u(0, x) = u_0(x) \quad x \in \Omega$$

## Numerische Behandlung

### 1. Zeitdiskretisierung

Zerlegung von  $I$  durch Stützstellen  $(t_0, t_1, \dots, t_N = T)$ .

Schreibweise:  $u^n = u(t_n, x)$ ,  $I_n = (t_{n-1}, t_n)$ ,  $K_n = t_n - t_{n-1}$ ,

$a(v, w) = (\nabla v, \nabla w)$

Schwache Formulierung:

$$\int_{I_n} (\partial_t u, \varphi) dt + \int_{I_n} a(u, \varphi) dt = \int_{I_n} (f, \varphi) dt \quad \forall \varphi \in V$$

$$\int_{I_n} a(u, \varphi) dt \approx K_n((1 - \alpha) a(u^{n-1}, \varphi) + \alpha a(u^n, \varphi)) \quad \alpha \in (0, 1]$$

Analog für  $\int_{I_n} (f, \varphi) dt$ .

$$(u^n - u^{n-1}, \varphi) + K_n \alpha a(u^n, \varphi) = \int_{I_n} (f, \varphi) dt - K_n(1 - \alpha) a(u^{n-1}, \varphi)$$

Also insgesamt:

$$(u^n, \varphi) + K_n \alpha a(u^n, \varphi) = \alpha K_n (f^n, \varphi) + (1 - \alpha) K_n (f^{n-1}, \varphi) + (u^{n-1}, \varphi) - K_n(1 - \alpha) a(u^{n-1}, \varphi)$$

### 2. Ortsdiskretisierung

Benutze FE zur Approximation von  $u^n$ .

Numerische Schemata:

1. Implizites Rückwärts-Euler-Verfahren

$$\left( \frac{u_h^n - u_h^{n-1}}{k_n}, \varphi \right) + a(u_h^n, \varphi) = (f^n, \varphi) \quad \forall \varphi \in V_h \subset V$$

$$u_h^0 = u_{0h}$$

Genauigkeit in der Zeit:  $(\Delta t)$  mit  $\Delta t = \max_n k_n$

2. Crank-Nicolson

$$\left( \frac{u_h^n - u_h^{n-1}}{k_n}, \varphi \right) + a \left( \frac{u_h^n + u_h^{n-1}}{2}, \varphi \right) = \left( \frac{f^n + f^{n-1}}{2}, \varphi \right) \quad \forall \varphi \in V_h$$

$$u_h^0 = u_{0h}$$

Genauigkeit in der Zeit:  $O(\Delta t^2)$

3. Implizites Zweischrittverfahren

$$\left( \frac{\frac{3}{2}u_h^n - \frac{3}{2}u_h^{n-1} - \frac{1}{2}u_h^{n-1} + \frac{1}{2}u_h^{n-2}}{k_n}, \varphi \right) + a(u_h^n, \varphi) = (f^n, \varphi) \quad \forall \varphi \in V_h$$

$u_h^0 = u_{0h}$  und  $u_h^1$  gegeben, z.B. durch einen Euler Schritt.

Genauigkeit in der Zeit:  $O(\Delta t^2)$

**Bemerkung 6.1** *Dieses Verfahren sollte dem Crank-Nicolson-Verfahren vorgezogen werden.*

## 6.1 Parabolische Variationsungleichungen

Sei  $K \subset V$  abgeschlossen und konvex.

Aufgabe: Finde  $u(t) \in K$  f.ü. auf  $I$ , so dass:

$$(\partial_t u, \varphi - u) + a(u, \varphi - u) \geq (f, \varphi - u) \quad \forall \varphi \in K$$

mit  $u(0) = u_0 \in K$ .

Numerische Schemata:

1. Euler-Rückwärts-Verfahren

$$\left( \frac{u_h^n - u_h^{n-1}}{k_n}, \varphi - u_h^n \right) + a(u_h^n, \varphi - u_h^n) \geq (f^n, \varphi - u_h^n) \quad \forall \varphi \in K_h = V_h \cap K$$

$$u_h^0 = u_{0h}$$

## 2. Implizites Zweischnittverfahren

$$\left( \frac{\frac{3}{2}u_h^n - \frac{3}{2}u_h^{n-1}}{k_n} - \frac{\frac{1}{2}u_h^{n-1} - \frac{1}{2}u_h^{n-2}}{k_{n-1}}, \varphi - u_h^n \right) + a(u_h^n, \varphi - u_h^n) \geq (f^n, \varphi - u_h^n)$$

$\forall \varphi \in K_h$ ,  $u_h^0 = u_{0h}$  und  $u_h^1$  gegeben, z.B. durch einen Euler Schritt.

# 7 Sattelpunktprobleme

Aufgabe: Variationsprobleme mit Nebenbedingungen

Bezeichnungen:

$X, M$  Hilberträume  
 $a : X \times X \rightarrow \mathbb{R}, b : X \times M \rightarrow \mathbb{R}$  stetige Bilinearformen  
 $X', M'$  Dualräume zu  $X$  und  $M$ .

Paarungen von  $X$  und  $X'$  und  $M$  und  $M'$  werden mit  $\langle \cdot, \cdot \rangle$  bezeichnet.  
Weiterhin gegeben:  $f \in X', g \in M'$

**Problem (PM)** Gesucht wird in  $X$  das Minimum von

$$J(u) = \frac{1}{2} a(u, u) - \langle f, u \rangle$$

unter den Nebenbedingungen

$$b(u, q) = \langle g, q \rangle \quad \text{für } q \in M$$

Umformulierung:

Betrachte die Lagrangefunktion

$$\mathcal{L}(u, \lambda) := J(u) - [b(u, \lambda) - \langle g, \lambda \rangle]$$

Dies führt auf das Sattelpunktproblem:

**Problem (PS)** Gesucht wird  $(u, \lambda) \in X \times M$  mit

$$\begin{aligned} a(u, \varphi) + b(\varphi, \lambda) &= \langle f, \varphi \rangle \quad \forall \varphi \in X \\ b(u, q) &= \langle g, q \rangle \quad \forall q \in M \end{aligned}$$

Für jede Lösung  $(u, \lambda)$  von (PS) kann man die Sattelpunkteigenschaft

$$\mathcal{L}(u, q) \leq \mathcal{L}(u, \lambda) \leq \mathcal{L}(\varphi, \lambda)$$

nachrechnen.

**Beispiel 7.1** Betrachte das Randwertproblem

$$\begin{aligned} -\Delta u &= f \quad \text{auf } \Omega \\ u &= g \quad \text{auf } \partial\Omega \end{aligned}$$

Mit  $X = H^1(\Omega)$  und  $M = L_2(\partial\Omega)$ .

Definiere:

$$\begin{aligned} a(u, \varphi) &= \int_{\Omega} \nabla u \nabla \varphi \, dx \\ \langle f, \varphi \rangle &= \int_{\Omega} f \varphi \, dx = (f, \varphi) \\ b(\varphi, q) &= \int_{\partial\Omega} \varphi q \, d\Gamma \\ \langle g, q \rangle &= \int_{\Gamma} g q \, d\Gamma \end{aligned}$$

Die Nebenbedingung lautet:

$$\int_{\Gamma} u q \, d\Gamma = \int_{\Gamma} g q \, d\Gamma \quad \forall q \in M$$

**Beispiel 7.2** Erneut

$$\begin{aligned} -\Delta u &= f \quad \text{auf } \Omega \subset \mathbb{R}^2 \\ u &= g \quad \text{auf } \partial\Omega \end{aligned}$$

Substitution:  $\sigma = \nabla u$  (Spannung  $\sigma$ ).

Mit  $\Delta u = \operatorname{div} \nabla u$  führt das auf das System

$$\begin{aligned} \sigma &= \nabla u \\ \operatorname{div} \sigma &= -f \end{aligned}$$

und  $u = 0$  auf  $\Gamma$

Herleitung einer variationellen Formulierung:

Ausgangspunkt: "Testen"

$$\begin{aligned} (\sigma, \tau) &= (\nabla u, \tau) \quad \forall \tau \in (L_2(\Omega))^2 \\ (\operatorname{div} \sigma, \varphi) &= (-f, \varphi) \quad \forall \varphi \in H_0^1(\Omega) \end{aligned}$$

Partielle Integration liefert:

$$\begin{aligned} (\sigma, \tau) - (\nabla u, \tau) &= 0 \quad \forall \tau \in (L_2(\Omega))^2 \\ -(\sigma, \nabla \varphi) &= (-f, \varphi) \quad \forall \varphi \in H_0^1(\Omega) \end{aligned}$$

In diesem Beispiel:

$$\begin{aligned} X &= (L_2(\Omega))^2, \quad M = H_0^1(\Omega), \\ a(\sigma, \tau) &= (\sigma, \tau) = \int_{\Omega} (\sigma_1 \tau_1 + \sigma_2 \tau_2) \, dx, \quad b(\tau, \varphi) = -(\tau, \nabla \varphi) \end{aligned}$$

## 7.1 Hilfsmittel aus der Funktionalanalysis

### 7.1.1 Adjungierte Operatoren

Bezeichnungen:

$X, Y$  Banachräume

$X', Y'$  Dualräume

$\langle \cdot, \cdot \rangle$  Paarungen zwischen  $X$  und  $X'$  bzw.  $Y$  und  $Y'$

$L : X \rightarrow Y$  beschränkter linearer Operator

#### Konstruktion eines stetigen linearen Funktionals auf $X$

1. Vorgabe eines  $y^* \in Y'$ .
2.  $x \in X$  wird abgebildet durch einsetzen in  
 $\langle y^*, L \cdot \rangle : X \ni x \mapsto \langle y^*, Lx \rangle \in \mathbb{R}$ .

Wir haben also nach Vorgabe eines  $y^* \in Y'$  durch  $\langle y^*, L \cdot \rangle$  ein Element aus  $X'$  konstruiert. Dieses wird mit  $L'y^* \in X'$  bezeichnet.

$$\langle L'y^*, x \rangle := \langle y^*, Lx \rangle$$

$$L' : Y' \rightarrow X'$$

**Definition 7.1**  $L'$  heisst der zu  $L$  adjungierte Operator.

**Definition 7.2** Sei  $V$  ein abgeschlossener Unterraum von  $X$ . Dann heisst

$$V^0 := \{l \in X' \mid \langle l, v \rangle = 0 \text{ für } v \in V\}$$

die Polare von  $V$ .

**Definition 7.3** Sei  $X$  ein Hilbertraum. Dann heisst

$$V^\perp := \{x \in X \mid (x, v) = 0 \text{ für } v \in V\}$$

orthogonales Komplement.

**Satz 7.1** (closed range theorem)

Mit obigen Voraussetzungen gilt:

1.  $L(X)$  abgeschlossen in  $Y$ .

$\Leftrightarrow$

2.  $L(X) = (\text{Ker } L')^0$

### 7.1.2 Abstrakter Existenzsatz

Bezeichnungen:

- $U, V$  Hilberträume.
- $U', V'$  Dualräume.
- $a : U \times V \rightarrow \mathbb{R}$  eine Bilinearform.

Definiere einen Operator  $L : U \rightarrow V'$ :

$$U \ni u \mapsto \langle Lu, \cdot \rangle := a(u, \cdot) \in V'$$

Die betrachteten Variationsprobleme haben die Struktur:

$$\begin{aligned} a(u, v) &= \langle f, v \rangle \quad \forall v \in V \\ \langle Lu, v \rangle &= \langle f, v \rangle \quad \forall v \in V \end{aligned}$$

Wir können formal schreiben:

$$u = L^{-1}f$$

**Definition 7.4** Seien  $U, V$  normierte Räume. Eine bijektive, lineare Abbildung  $L : U \rightarrow V'$  heisst Isomorphismus, wenn  $L$  und  $L^{-1}$  stetig sind.

**Satz 7.2** Seien  $U, V$  Hilberträume. Eine lineare Abbildung  $L : U \rightarrow V'$  ist ein Isomorphismus, wenn die zugehörige Form  $a : U \times V \rightarrow \mathbb{R}$  folgende Bedingungen erfüllt.

1. Stetigkeit:  $\exists c \geq 0$  mit  $|a(u, v)| \leq c\|u\|_U\|v\|_V$
2. Inf-sup-Bedingung:  $\exists \alpha > 0$  mit  $\sup_{v \in V, v \neq 0} \frac{a(u, v)}{\|v\|} \geq \alpha\|u\|_U \quad \forall u \in U$   
(bzw.  $\inf_{u \in U} \sup_{v \in V} \frac{a(u, v)}{\|u\|_U\|v\|_V} \geq \alpha > 0$ )
3.  $\forall v \in V \exists u \in U$  mit  $a(u, v) \neq 0$ .

**Beweis:**

a) Aus 2. folgt  $L$  ist injektiv

$$Lu_1 = Lu_2 \rightarrow a(u_1, v) = a(u_2, v) \quad \forall v \in V$$

$$\text{Also ist } \sup_{v \in V} \frac{a(u_1 - u_2, v)}{\|v\|_V} = 0 \geq \alpha\|u_1 - u_2\|$$

$$\Rightarrow u_1 = u_2$$

b)  $L^{-1}$  ist stetig auf dem Bild von  $L$ .

Zu  $f \in L(U)$  gibt es wegen a) ein eindeutiges  $u = L^{-1}f$ .

$$\begin{aligned} \alpha \|L^{-1}f\|_U &= \alpha \|u\|_U \\ &\leq \sup_{v \in V} \frac{a(u, v)}{\|v\|_V} \\ &= \sup_{v \in V} \frac{\langle Lu, v \rangle}{\|v\|_V} \\ &= \sup_{v \in V} \frac{\langle f, v \rangle}{\|v\|_V} \\ &= \|f\|_{V'} \end{aligned}$$

$$\Rightarrow \frac{\|L^{-1}f\|_U}{\|f\|_{V'}} \leq \frac{1}{\alpha}$$

$$f \in L(U) \text{ war beliebig, d.h. } \|L^{-1}\| = \sup_{f \in V'} \frac{\|L^{-1}f\|_U}{\|f\|_{V'}} \leq \frac{1}{\alpha}.$$

c)  $L(U)$  ist abgeschlossen.

$L^{-1}$  ist stetig nach b).

$L$  ist stetig.

$\Rightarrow L(U)$  abgeschlossen.

d)  $L$  ist surjektiv.

$$L' : V \rightarrow U'$$

$$V \ni v \mapsto \langle L \cdot, v \rangle = a(\cdot, v)$$

wegen 3. liegt nur  $v = 0 \in V$  im Kern( $L'$ ). Es ist Kern( $L'$ )  $\subset V$

Satz 7.1 sagt:

$$\begin{aligned} L(U) &= (\text{Kern } L')^0 \\ &= \{l \in V' \mid \langle l, v \rangle = 0, v \in \text{Kern } L'\} \\ &= V' \end{aligned}$$

□

### 7.1.3 Abstrakter Konvergenzsatz

Seien  $U_h \subset U$  und  $V_h \subset V$  endlich-dimensionale Räume. Die diskrete Aufgabe lautet:

$$a(u_h, v) = \langle f, v \rangle \quad \forall v \in V_h \tag{7.1}$$

**Satz 7.3** *Es gelten die Bezeichnungen und Bedingungen aus Satz 7.2. Weiterhin seien  $u_h \subset U$ ,  $V_h \subset V$  so gewählt, dass gilt:*

$$2.h \quad \inf_{u_h \in U_h} \sup_{v_h \in V_h} \frac{a(u_h, v_h)}{\|u_h\| \|v_h\|} \geq \alpha$$

und

$$3.h \quad \forall v_h \in V_h \exists u_h \in U_h \text{ mit } a(u_h, v_h) \neq 0.$$

Dann gilt die Fehlerabschätzung:

$$\|u - u_h\| \leq \left(1 - \frac{C}{\alpha}\right) \inf_{w_h \in U_h} \|u - w_h\|$$

**Beweis:** Es gilt die Galerkin-Bedingung:

$$a(u - u_h, v) = 0 \quad \forall v \in V_h$$

Damit gilt für ein beliebiges  $w_h \in U_h$ :

$$a(u - w_h, v) = a(u_h - w_h, v) \quad \forall v \in V_h \tag{7.2}$$

Wir schätzen ab:

$$\begin{aligned} \alpha \|u_h - w_h\| &\leq \sup_{v_h \in V_h} \frac{a(u_h - w_h, v_h)}{\|v_h\|} \\ &\leq \sup_{v_h \in V_h} \frac{a(u - w_h, v_h)}{\|v_h\|} \\ &\leq \sup_{v \in V} \frac{a(u - w_h, v)}{\|v\|} \\ &\leq C \frac{\|u - w_h\| \|v\|}{\|v\|} = C \|u - w_h\| \end{aligned}$$

Zusammen also:

$$\|u_h - w_h\| \leq \frac{C}{\alpha} \|u - w_h\|$$

Benutze nun die Dreiecksungleichung:

$$\begin{aligned} \|u - w_h + w_h - u_h\| &\leq \|u - w_h\| + \|w_h - u_h\| \\ &\leq \left(1 + \frac{C}{\alpha}\right) \|u - w_h\| \end{aligned}$$

□

## 7.2 Die Inf-sup-Bedingung

Aufgabe PS:

Gesucht wird  $(u, \lambda) \in X \times M$  mit

$$\begin{aligned} a(u, \varphi) + b(\varphi, \lambda) &= \langle f, \varphi \rangle \quad \forall \varphi \in X \\ b(u, q) &= \langle g, q \rangle \quad \forall q \in M \end{aligned}$$

Abstrakte Situation PA:

$$L : X \times M \rightarrow X' \times M', (u, \lambda) \mapsto (f, g)$$

**Satz 7.4** Durch das Sattelpunktproblem PS wird mit PA genau dann ein Isomorphismus  $L : X \times M \rightarrow X' \times M'$  erklärt, wenn die beiden folgenden Bedingungen erfüllt sind:

1. Die Bilinearform  $a$  ist  $V$ -elliptisch, d.h.  $\exists \alpha > 0$  und  $V := \{v \in X \mid b(v, q) = 0 \text{ für } q \in M\}$ , so dass gilt:

$$a(v, v) \geq \alpha \|v\|^2 \quad \text{für } v \in V \subset X$$

2.  $\exists \beta > 0$  mit  $\inf_{q \in M} \sup_{v \in X} \frac{b(v, q)}{\|v\| \|q\|} \geq \beta$

**Beweisbemerkung:** Die abstrakte inf-sup-Bedingung aus Satz 7.2 kann hier durch Eigenschaften der Formen  $a$  und  $b$  ausgedrückt werden.

## 7.3 Gemischte Finite-Element-Methoden

Wähle  $X_h \subset X, M_h \subset M$ . Man löst:

$PS_h$ : Gesucht wird  $(u_h, \lambda_h) \in X_h \times M_h$  mit

$$\begin{aligned} a(u_h, \varphi) + b(\varphi, \lambda_h) &= \langle f, \varphi \rangle \quad \forall \varphi \in X_h \\ b(u_h, q) &= \langle g, q \rangle \quad \forall q \in M_h \end{aligned}$$

**Definition 7.5** Eine Familie von FE-Räumen  $X_h, M_h$  erfüllt die Babuska-Brezzi-Bedingung, wenn es von  $h$  unabhängige Zahlen  $\alpha > 0$  und  $\beta > 0$  mit folgenden Eigenschaften gibt:

- 1.<sub>h</sub>  $a$  ist  $V_h$ -elliptisch, d.h. mit  $\alpha > 0$  und  $V_h := \{v_h \in X_h \mid b(v_h, q_h) = 0 \text{ für } q_h \in M_h\}$  gilt:

$$a(v_h, v_h) \geq \alpha \|v_h\|^2 \quad \forall v_h \in V_h$$

- 2.<sub>h</sub> Mit  $\beta > 0$  gilt:

$$\inf_{q_h \in M_h} \sup_{v_h \in X_h} \frac{b(v_h, q_h)}{\|v_h\| \|q_h\|} \geq \beta$$

**Satz 7.5** Die Voraussetzungen von Satz 7.4 seien erfüllt und  $X_h, M_h$  mögen die Babuska-Brezzi-Bedingung erfüllen. Dann gilt die Fehlerabschätzung:

$$\|u - u_h\| + \|\lambda - \lambda_h\| \leq c \left( \inf_{v_h \in X_h} \|u - v_h\| + \inf_{q_h \in M_h} \|\lambda - q_h\| \right)$$

**Beweis 1** Anwendung des abstrakten Konvergenzsatzes 7.3 □

## 7.4 Diskrete Sattelpunktprobleme

Die Diskretisierung von PS durch  $PS_h$  führt auf das folgende System

$$\begin{aligned} Au + B^T \lambda &= f & A \in M(n, n) & \quad f, u \in \mathbb{R}^n \\ Bu &= g & B \in M(m, n) & \quad g, \lambda \in \mathbb{R}^m \end{aligned}$$

wobei  $A$  positiv definit ist.

Bemerkung:  $A$  positiv definit  $\Rightarrow A^{-1}$  existiert.

Somit lässt sich  $u$  schreiben als

$$u = A^{-1}(f - B^T \lambda)$$

Einsetzen in  $Bu = g$  liefert:

$$\begin{aligned} B(A^{-1}f - A^{-1}B^T \lambda) &= g \\ \Leftrightarrow BA^{-1}B^T \lambda &= BA^{-1}f - g \end{aligned}$$

Die implizit gegebene Matrix  $S = BA^{-1}B^T$  heisst Schurkomplement.

Bemerkung:  $A^{-1}$  heisst positiv definit  $\Rightarrow A^{-1} = LL^T$ .

Wegen  $S = BA^{-1}B^T = BLL^TB^T$  ist  $S$  symmetrisch und es gilt:

$$(S\lambda, \lambda) = (BLL^TB^T \lambda, \lambda) = (L^TB^T \lambda, L^TB^T \lambda) \geq 0 \quad \forall \lambda \in \mathbb{R}^m$$

$\Rightarrow S$  positiv definit.

Zur Lösung von  $BA^{-1}B^T \lambda = BA^{-1}f - g$  können also Gradientenverfahren und cg-Verfahren verwendet werden. Das Resultat ist  $\bar{\lambda} \in \mathbb{R}^m$ .

Berechne die Lösung  $u$  durch

$$u = A^{-1}f - A^{-1}B^T \bar{\lambda}.$$

### Simultane cg-Variante

**Algorithmus 7.1** *Uzawa-Algorithmus mit konjugierten Richtungen*

1. *Initialisierung:*  $\lambda_0 \in \mathbb{R}^m$ .

$$Au_1 = f - B^T \lambda_0, \quad d_1 = -q_1 = Bu_1 - g$$

2. Für  $k = 1, 2, \dots$

$$\begin{aligned} p_k &= B^T d_k \\ h_k &= A^{-1} p_k \\ \alpha_k &= \frac{q_k^T q_k}{p_k^T h_k} \\ \lambda_k &= \lambda_{k-1} + \alpha_k d_k \\ u_{k+1} &= u_k - \alpha_k h_k \\ q_{k+1} &= g - B u_{k+1} \\ \beta_k &= \frac{q_{k+1}^T q_{k+1}}{q_k^T q_k} \\ d_{k+1} &= -q_{k+1} + \beta_k d_k \end{aligned}$$

Alternativen:

Algemeiner Ansatz:

$$Gx = b, \quad G \in M(N, N), \quad x, b \in \mathbb{R}^N$$

$$x_{k+1} = x_k + T^{-1}(b - Gx_k).$$

z.B.  $T = \text{diag}(G)$  Jacobi-Verfahren

oder  $T = (0 \setminus G)$  Gauss-Seidel-Verfahren.

Übertragung auf  $S$ :

$$\begin{pmatrix} u_{k+1} \\ \lambda_{k+1} \end{pmatrix} = \begin{pmatrix} u_k \\ \lambda_k \end{pmatrix} + \begin{pmatrix} C & B^T \\ B & 0 \end{pmatrix}^{-1} \begin{pmatrix} f - Au_k - B^T \lambda_k \\ g - Bu_k \end{pmatrix}$$

z.B.  $C = \text{diag}(A)$ .

Allgemein sollte  $C$  eine symmetrische, positiv definite Matrix sein, die

1.  $A$  approximiert.
2. Leicht invertierbar ist.

**Algorithmus 7.2** Über Defektkorrekturen

Für  $k = 1, 2, \dots$

1. Bestimme das Residuum  $\begin{pmatrix} r_k^u \\ r_k^\lambda \end{pmatrix} = \begin{pmatrix} f - Au_k - B^T \lambda_k \\ g - Bu_k \end{pmatrix}$
2. Berechne  $\begin{pmatrix} d_k^u \\ d_k^\lambda \end{pmatrix}$  durch  $\begin{pmatrix} C & B^T \\ B & 0 \end{pmatrix} \begin{pmatrix} d_k^u \\ d_k^\lambda \end{pmatrix} = \begin{pmatrix} r_k^u \\ r_k^\lambda \end{pmatrix}$
3.  $\begin{pmatrix} u_{k+1} \\ \lambda_{k+1} \end{pmatrix} = \begin{pmatrix} u_k \\ \lambda_k \end{pmatrix} + \begin{pmatrix} d_k^u \\ d_k^\lambda \end{pmatrix}$

## 7.5 Laplace-Gleichung als gemischtes Problem

Problem: Ergänzung zu Beispiel 7.2:

$\Delta u = -f$  kann man formal als System schreiben (o.B.d.A. auf  $\Omega \subset \mathbb{R}^2$ ).

$$\begin{aligned} \nabla u &= \sigma & u : \Omega &\rightarrow \mathbb{R} \\ \partial_x \sigma_1 + \partial_y \sigma_2 = \operatorname{div} \sigma &= -f & \sigma &= \begin{pmatrix} \sigma_1 \\ \sigma_2 \end{pmatrix}, \sigma_i : \Omega \rightarrow \mathbb{R} \end{aligned}$$

### 7.5.1 Primal-gemischte variationelle Formulierung

Gesucht ist  $(\sigma, u) \in L_2(\Omega)^2 \times H_0^1(\Omega)$  so, dass gilt:

$$\begin{aligned} (\sigma, \tau) - (\nabla u, \tau) &= 0 & \forall \tau \in L_2(\Omega)^2 \\ -(\sigma, \nabla \varphi) &= -(f, \varphi) & \forall \varphi \in H_0^1(\Omega) \end{aligned}$$

In den abstrakten Rahmen fällt dies mit

$$X = L_2(\Omega)^2, M = H_0^1(\Omega)$$

$$a(\sigma, \tau) = (\sigma, \tau) := (\sigma_1, \tau_1) + (\sigma_2, \tau_2), b(\tau, \varphi) = -(\tau, \nabla \varphi).$$

Nachweis:

1. (Bedingung 1. aus Satz 7.4)

Wegen  $\|\sigma\|^2 = a(\sigma, \sigma)$  ist  $a$  sogar elliptisch auf ganz  $X$ .

2. (Bedingung 2. aus Satz 7.4)

Sei  $v \in H_0^1(\Omega)$  gegeben. Wähle  $\tau = -\nabla v \in L_2(\Omega)^2$ . Es gilt:

$$\begin{aligned} \frac{b(\tau, v)}{\|\tau\|} &= \frac{-(\tau, \nabla v)}{\|\tau\|} \\ &= \frac{(\nabla v, \nabla v)}{\|v\|} \\ &= \|\nabla v\| \\ &\geq \frac{1}{2} \|v\|_1 \end{aligned}$$

$$\text{Also } \sup_{\tau \in L_2(\Omega)^2} \frac{b(\tau, v)}{\|\tau\|} \geq \frac{1}{C} \|v\|_1$$

### Passende Finite Elemente

Triangulierung  $\mathbb{T}_h$  mit Dreiecken. Wähle  $k \geq 1$  und

$$X_h = \{\sigma_h \in L_2(\Omega)^2 \mid \sigma_h|_T \in P_{k-1}\}, M_h = \{\varphi_h \in H_0^1(\Omega) \mid \varphi_h|_T \in P_k\}.$$

Man beachte:  $X_h$  unstetiger,  $M_h$  stetiger Ansatz.

### 7.5.2 Dual-gemischte Formulierung

Vorbereitung:  $H_{div} := \{\tau \in (L_2(\Omega))^2 \mid \text{div } \tau \in L_2(\Omega)\}$ .

Zugehörige Norm:  $\|\tau\|_{div}^2 = \|\tau\|^2 + \|\text{div } \tau\|^2$ .

Für  $v \in H_0^1(\Omega)$  und  $\sigma \in H_{div}$  gilt:

$$(\sigma, \nabla v) = -(\text{div } \sigma, v)$$

Somit lautet die dual-gemischte Formulierung:

Gesucht ist  $(\sigma, u) \in H_{div} \times L_2(\Omega)$ , so dass gilt:

$$\begin{aligned} (\sigma, \tau) + (u, \text{div } \tau) &= 0 & \forall \tau \in H_{div} \\ (\text{div } \sigma, \varphi) &= -(f, \varphi) & \forall \varphi \in L_2(\Omega) \end{aligned}$$

Dies fällt in den abstrakten Rahmen mit

$X := H_{div}$ ,  $M := L_2(\Omega)$ ,  $a(\sigma, \tau) = (\sigma, \tau)$ ,  $b(\tau, v) = (\text{div } \tau, v)$ .

**Bemerkung 7.1**  $V = \{\tau \in H_{div} \mid (\text{div } \tau, v) = 0 \text{ für } v \in L_2(\Omega)\}$ .

Nachweis:

1. 1. aus Satz 7.4:

$$\begin{aligned} a(\tau, \tau) = (\tau, \tau) &= (\tau, \tau) + (\text{div } \tau, \text{div } \tau) \\ &= \|\tau\|^2 + \|\text{div } \tau\|^2 \\ &= \|\tau\|_{div}^2 \quad \forall \tau \in V \end{aligned}$$

2. 2. aus Satz 7.4:

$$\text{zz: } \sup_{\tau \in H_{div}} \frac{(\text{div } \tau, v)}{\|\tau\|_{div}} \geq \beta \|v\|$$

a) Sei  $v \in L_2(\Omega)$  beliebig. Wähle dazu  $w \in C_0^\infty(\Omega)$  mit

$$\|v - w\| \leq \frac{1}{2} \|v\|. \tag{7.3}$$

(Wahl von  $w$  möglich: Dichtheitsargument)

b) Man setze  $\xi := \inf\{x_1 \mid x \in \Omega\}$  und  $\tau_1(x) = \int_{\xi}^{x_1} w(t, x_2) dt$ ,  $\tau_2(x) = 0$ .

Damit ist offenbar  $\text{div } \tau = \partial_{x_1} \tau_1 = w$ .

Nun folgt Argumentation ähnlich wie im Beweis zur Poincaré-Ungleichung:

$$|\tau(x)|^2 = |\tau_1(x)|^2 \leq \int_{\xi}^{x_1} |w(t, x_2)|^2 dt \leq \int_{\xi}^s |w(t, x_2)|^2 dt$$

mit  $s = \max\{x_1 \mid x \in \Omega\}$ .

Integration über  $x_1$ :

$$\int_{\xi}^s |\tau_1(x)|^2 dx_1 \leq c \int_{\xi}^s |w(t, x_2)|^2 dt = c \int_{\xi}^s |w(x_1, x_2)|^2 dx$$

Integration über  $x_2$ :

$$\|\tau\|^2 \leq c \|w\|^2 \tag{7.4}$$

c) Ausgangspunkt ist (7.3).

$$\begin{aligned} (v - w, v - w) &\leq \frac{1}{4} \|v\|^2 \\ \Leftrightarrow 2(v, w) &\geq \frac{3}{4} \|v\|^2 + \|w\|^2 \\ \Rightarrow (v, w) &\geq \frac{3}{8} \|v\|^2 + \frac{1}{2} \|w\|^2 \geq c \|v\|^2 \end{aligned}$$

d) Auswertung von  $\frac{b(\tau, v)}{\|\tau\|_{div}}$ .

$$\begin{aligned} \frac{b(\tau, v)}{\|\tau\|_{div}} &= \frac{(div \tau, v)}{(\|\tau\|^2 + \|div \tau\|^2)^{1/2}} \\ &\geq \frac{(div \tau, v)}{c \|w\|} \\ &= \frac{(w, v)}{c \|w\|} \\ &\geq \frac{c \|v\|^2}{d \|v\|} \\ &= c \|v\| \end{aligned}$$

### Passende Finite Elemente

Das Raviart-Thomas-Element.

Zur Dual-gemischten Formulierung:

$$X_h = \left\{ \tau \in L_2(\Omega)^2 \mid \tau|_T = \begin{pmatrix} a_t \\ b_t \end{pmatrix} + c_T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \text{ für } T \in \mathbb{T}_h, a_T, b_T, c_T \in \mathbb{R}, \right. \\ \left. \tau \cdot n \text{ stetig an den Elementgrenzen} \right\}$$

$$M_h = \{v \in L_2(\Omega) \mid v|_T = d_T \text{ für } T \in \mathbb{T}_h, d_T \in \mathbb{R}\}.$$

**Bemerkung 7.2**  $div \tau|_T = \partial_{x_1} \tau_1 + \partial_{x_2} \tau_2 = c_T + c_T = const$  für  $\tau \in X_h$ .

# 8 Themengebiete

Finanzmathematik

Gletscher

Mechanik

1. Simulation, Diskretisierung, Lösung, Analyse
2. Parameteridentifizierung, Inverse Probleme

**Beispiel 8.1**  $-\mu\Delta u = f$

*Man kennt, z.B. durch Messung,  $u$  an verschiedenen Punkten.*

*Fragestellung  $\mu$ ?*

## Optionshandel

Eine Option ist z.B. das Recht, nicht die Verpflichtung, zu einem Zeitpunkt  $T$  eine Aktie zu kaufen, zum vorher festgelegten Preis  $K$ .

Bezeichnungen:  $S$  ist der Preis der Aktie,  $V$  ist der Wert der Option.

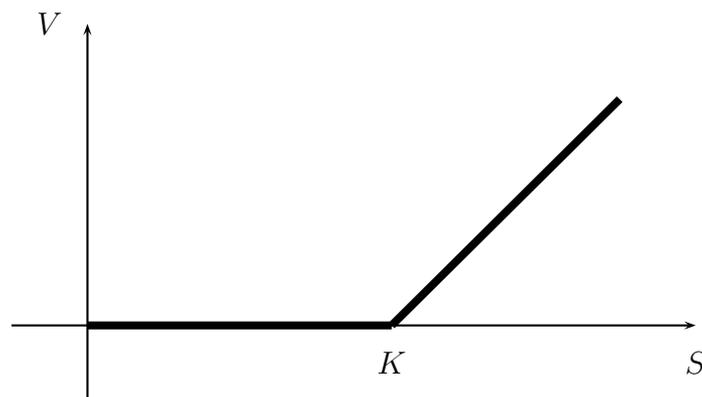


Abbildung 8.1: Pay-Off-Diagramm für  $t = T$ .

## Modellierung des Optionshandels

Black-Scholes-Gleichung:

$$\partial_t V + \frac{1}{2}\sigma^2 S^2 \partial_S^2 V + rS \partial_S V - rV = 0$$

Durch Substitution:

$$\frac{\partial y}{\partial \tau} - \partial_x^2 y = 0$$

+ Randbedingungen, Anfangsbedingungen

**Europäische Optionen:** Handlung nur bei  $T$ .

**Amerikanische Optionen:** Handlung für  $t \leq T$ .

$$V_{Am} \geq V_{Euro}$$

Bei Amerikanischen Optionen gilt in schwacher Form:

$$(\partial_t y, \varphi - y) + (\partial_x y, \partial_x(\varphi - y)) \geq 0$$

+ Randbedingungen, Anfangsbedingungen

### Gletscher

“Verbale” Modellierung:

Eis ist ein langsam fließendes, temperaturabhängiges, inkompressibles Fluid.

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

$u = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}$ ,  $\Delta u = \begin{pmatrix} \Delta u_1 \\ \Delta u_2 \end{pmatrix}$ ,  $u$  ist Vektor der Verschiebungen,  $p$  ist der skalare Druck.

$\Omega$  ist a-priori nicht bekannt.

Über Vereinfachung der Modellierung wird die Höhenfunktion  $h$  gesucht, für die gilt:

$$(\partial_t h, \varphi) + (c(h)\nabla h, \nabla \varphi) = (a, \varphi) \quad \text{auf } I$$

$a$  ist die Akkumulationsfunktion,  $I$  ist unbekannt.

Wähle  $\Omega \supset I$  gross genug:

$$\partial_t h, \varphi - h) + (c(h)\nabla h, \nabla(\varphi - h)) \geq (a, \varphi - h) \quad \text{auf } \Omega$$

$\Omega$  fest gewählt. Nebenbedingung:  $h, \varphi \geq 0$

$$\begin{aligned} -\varepsilon \Delta u + (u \nabla) u + \nabla p &= f \\ \operatorname{div} u &= 0 \end{aligned}$$

Navier-Stokes-Gleichungen beschreiben inkompressible Flüssigkeiten.

Nichtlineare Iteration:

Wähle  $u_0$  als Startlösung

Iteriere  $i = 1, 2, \dots$

$$\begin{aligned} -\varepsilon \Delta u_i + (u_{i-1} \nabla) u_i + \nabla p_i &= f \\ \operatorname{div} u_i &= 0 \end{aligned}$$

# Literaturverzeichnis

- [1] D. Braess. *Finite Elemente*. Springer, 1997.
- [2] R. Glowinski. *Numerical methods for nonlinear variational problems*. Springer Series in Comp. Physics. Springer, 1983.
- [3] C. Grossmann and H.-G. Roos. *Numerik partieller Differentialgleichungen*. Teubner Studienbücher, 1992.
- [4] C. Johnson. *Numerical solution of partial differential equations by the finite element method*. Studentlitteratur, 1987.
- [5] D. Kinderlehrer and Stampacchia G. *An introduction to variational inequalities and their applications*. Academic Press, 1980.
- [6] H.R. Schwarz. *Methode der finiten Elemente*. Teubner, 1984.
- [7] R. Seydel. *Tools for Computational Finance*. Springer, 2002.