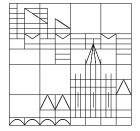
# Skript zum Numerik-Teil

# der Vorlesung

# Theorie und Numerik partieller Differentialgleichungen

Wintersemester 2024/25

### Johannes Schropp



Universität Konstanz

Fachbereich Mathematik und Statistik

Stand: 5. Februar 2025

### Inhaltsverzeichnis

1	Finite Differenzenverfahren für elliptische Differentialgleichungen	į
	a) Das klassische Differenzenverfahren	3
	b) Konsistenz und Stabilität des klassischen Differenzenverfahrens für die Dirichletsche Randwertaufgabe	11
2	Finite Elemente Methoden für elliptische Differentialgleichungen	19
	a) Theoretische Grundlagen zur Finite Elemente Methode	19
	b) Finite Elemente Methoden	25
	c) Stabilität und Konvergenz der Finite Elemente Methode	31
3	Finite Differenzenverfahren für parabolische Differentialgleichungen	38
	a) Das Prinzip der Linienmethode	38
	b) Konsistenz der Differenzenverfahren	45
	c) Stabilität und Konvergenz der Differenzenverfahren	48
4	Hyperbolische Differentialgleichungen	54
	a) Differenzenverfahren für die Wellengleichung	54
	b) Die Courant-Friedrichs-Levy Bedingung	61

# 1. Finite Differenzenverfahren für elliptische Differentialgleichungen

### a) Das klassische Differenzenverfahren

Wir betrachten die Dirichlet'sche Randwertaufgabe

$$-\Delta u(x,y) = g(x,y), \quad (x,y) \in \Omega,$$
  

$$u(x,y) = \gamma(x,y), \quad (x,y) \in \partial \Omega.$$
(1-1)

Hierbei ist  $\Omega \subset \mathbb{R}^2$  ein beschränktes Gebiet mit gegebenen Funktionen  $g:\Omega \longrightarrow \mathbb{R}$ ,  $\gamma:\partial\Omega \longrightarrow \mathbb{R}$ . Eine Funktion  $u\in C^2(\Omega)\cap C^0(\overline{\Omega})$ , welche (1-1) erfüllt, heißt klassische Lösung von (1-1).

**1.1 Satz.** Vorgelegt sei (1-1) mit  $g \in C^0(\overline{\Omega})$ ,  $\gamma \in C^0(\partial \Omega)$  und  $\partial \Omega$  sei stückweise stetig differenzierbar. Dann existiert genau eine klassische Lösung von (1-1).

Für diese Lösung gilt das Maximum-Minimum Prinzip, d.h.

$$g \ge 0 \implies u \ge \min\{\gamma(x) \mid x \in \partial\Omega\},\$$
  
$$g \le 0 \implies u \le \max\{\gamma(x) \mid x \in \partial\Omega\}.$$

Dabei setzen wir hier für ein  $f: A \longrightarrow \mathbb{R}$ :

$$f \ge 0 \iff f(x) \ge 0 \quad \forall x \in A.$$

Seien nun

$$L: C^2(\Omega) \longrightarrow C^0(\Omega), \quad Lu = -\Delta u,$$
  
 $R: C^0(\overline{\Omega}) \longrightarrow C^0(\partial \Omega), \quad Ru = u|_{\partial \Omega},$ 

so lässt sich (1-1) als

$$Lu = q$$
,  $Ru = \gamma$ 

schreiben.

Sei  $\Omega$  ein beschränktes Gebiet mit stückweise glattem Rand. Das Paar (L,R) ist invers monoton, d.h. für die nach Satz 1.1 eindeutige Lösung u von  $Lu=g, Ru=\gamma$  gilt

$$(g, \gamma) \ge 0 \Longrightarrow u \ge 0.$$

Offensichtlich folgt

$$\gamma, g \ge 0 \Longrightarrow u \ge \min\{\gamma(x) \mid x \in \partial\Omega\} \ge 0.$$

© Johannes Schropp 5. Februar 2025

Diese Eigenschaft der Inversmonotonie wollen wir auch in den numerischen Verfahren zu (1-1) wiederfinden.

Nun wenden wir uns der numerischen Behandlung des Problems (1-1) zu.

Sei h > 0. Wir überziehen zunächst  $\mathbb{R}^2$  mit einem äquidistanten Gitter

$$\mathbb{R}^2_h := \{ (ih, jh) \, | \, i, j \in \mathbb{Z} \}.$$

Es kann sein, dass  $\partial\Omega\cap\mathbb{R}^2_h=\varnothing$  gilt, d.h. dass  $\mathbb{R}^2_h$  überhaupt keine Randpunkte von  $\Omega$  enthält. Um dieser Schwierigkeit erst einmal aus dem Weg zu gehen, nehmen wir an

$$\Omega = (0,1)^2$$

und definieren das Gitter

$$\Omega_h = \{(ih, jh) \mid i, j \in \{1, \dots, M-1\}\}$$

für ein  $h = \frac{1}{M} > 0$ .

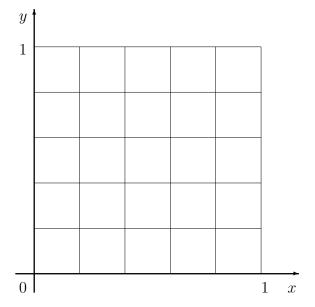


Abbildung 1: Beispiel für ein Gitter über  $(0,1)^2$  mit  $M=5, h=0.2, (M-1)^2=16$ 

Die Methode der Finiten Differenzen beruht nun darauf, Ableitungen von (1-1) durch Differenzenguotienten zu ersetzen.

Sei nun  $v \in C^4([a,b],\mathbb{R})$ , so gilt nach der Taylor-Formel

$$v(s \pm h) = v(s) \pm hv'(s) + \frac{h^2}{2}v''(s) \pm \frac{h^3}{6}v^{(3)}(s) + \frac{h^4}{24}v^{(4)}(\eta_{\pm})$$

für  $s, s \pm h \in [a, b]$  mit  $\eta_- \in [s - h, s], \, \eta_+ \in [s, s + h].$ 

Dann gilt

$$\begin{split} \left| \frac{1}{h^2} (-v(s-h) + 2v(s) - v(s+h)) + v''(s) \right| &= \\ &= \left| h^{-2} \left[ -v(s) + hv'(s) - \frac{h^2}{2} v''(s) + \frac{h^3}{6} v^{(3)}(s) - \frac{h^4}{24} v^{(4)}(\eta_-) + 2v(s) \right. \\ &- v(s) - hv'(s) - \frac{h^2}{2} v''(s) - \frac{h^3}{6} v^{(3)}(s) - \frac{h^4}{24} v^{(4)}(\eta_+) \right] + v''(s) \right| \\ &= \frac{h^2}{24} \left| v^{(4)}(\eta_-) + v^{(4)}(\eta_+) \right| \le Ch^2 = O(h^2). \end{split}$$

Somit ist

$$-v''(s) = \frac{1}{h^2} \left( -v(s-h) + 2v(s) - v(s+h) \right) + O(h^2), \tag{1-2}$$

falls  $v \in C^4([a, b])$ .

Mit der Formel (1-2) für -v'' bekommen wir für  $(x,y)\in\Omega_h$  folgende Approximationen:

$$-u_{xx}(x,y) \sim \frac{1}{h^2} \left( -u(x-h,y) + 2u(x,y) - u(x+h,y) \right),$$
  
$$-u_{yy}(x,y) \sim \frac{1}{h^2} \left( -u(x,y-h) + 2u(x,y) - u(x,y+h) \right),$$

also

$$-\Delta u(x,y) \sim \frac{1}{h^2} \left( -u(x-h,y) - u(x+h,y) - u(x,y-h) - u(x,y+h) + 4u(x,y) \right)$$

oder mit der Schreibweise  $u_{ij} = u(ih, jh)$ 

$$-(\Delta u)_{ij} \sim \frac{1}{h^2} \left( -u_{i-1,j} - u_{i+1,j} - u_{i,j-1} - u_{i,j+1} + 4u_{ij} \right).$$

An den inneren Gitterpunkten

$$\stackrel{\circ}{\Omega}_h = \{(ih, jh) \in \Omega_h \,|\, i, j \in \{2, \dots, M-2\}\}$$

ersetzen wir die Differentialgleichung (1-1) durch

$$h^{-2}(-u_{i-1,j}-u_{i+1,j}-u_{i,j-1}-u_{i,j+1}+4u_{i,j})=g_{ij}$$

 $mit g_{ij} = g(ih, jh).$ 

Bei den randnahen Gitterpunkten  $\Omega_h^R=\Omega\backslash\stackrel{\circ}{\Omega}_h$  wählen wir dieselbe Ersetzung, wobei jetzt

$$u_{i\pm 1,j} = \gamma_{i\pm 1,j}$$
 bzw.  $u_{i,j\pm 1} = \gamma_{i,j\pm 1}$ 

einzusetzen ist, falls  $((i \pm 1)h, jh)$  bzw.  $(ih, (j \pm 1)h) \in \partial \Omega$ .

Insgesamt erhalten wir ein lineares Gleichungssystem

$$Au = r (1-3)$$

der Dimension  $(M-1)^2$  für die gesuchte Approximation u.

Wir können das Differenzenverfahren auch präziser formulieren.

#### **1.2 Definition.** Eine Abbildung

$$w: \Omega_h \longrightarrow \mathbb{R},$$
 $w = (\underbrace{(w(ih, jh))}_{=w_{ij}}, (i, j) \in \{1, \dots, M-1\})$ 

nennen wir eine Gitterfunktion und schreiben kurz  $w \in \mathbb{R}^{\Omega_h}$  oder genauer  $w^h \in \mathbb{R}^{\Omega_h}$ .

Somit gilt für A, u, r aus (1-3)

$$A \in \mathbb{R}^{\Omega_h, \Omega_h}, u, r \in \mathbb{R}^{\Omega_h}.$$

Die explizite Darstellung der Matrix A und des Vektors r in (1-3) hängt von der Nummerierung der Gitterpunkte ab.

Beispielsweise erhalten wir für die natürliche, zeilenweise von links unten nach rechts oben laufende Nummerierung folgende Darstellung:

$$u = (u^1, \dots, u^{M-1}), \quad u^j = (u_{1,j}, u_{2,j}, \dots, u_{M-1,j}) \in \mathbb{R}^{M-1},$$
  
 $g = (g^1, \dots, g^{M-1}), \quad g^j = (g_{1,j}, g_{2,j}, \dots, g_{M-1,j}) \in \mathbb{R}^{M-1}$ 

sowie

$$B = h^{-2} \begin{pmatrix} 4 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 4 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 4 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 4 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 4 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 4 \end{pmatrix} \in \mathbb{R}^{M-1,M-1},$$

$$C = h^{-2}I \in \mathbb{R}^{M-1,M-1}$$

und in Blockschreibweise

$$\begin{pmatrix}
B & -C & 0 & \dots & 0 & 0 & 0 \\
-C & B & -C & \dots & 0 & 0 & 0 & 0 \\
0 & -C & B & \dots & 0 & 0 & 0 & 0 \\
\vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & 0 & \dots & B & -C & 0 \\
0 & 0 & 0 & \dots & -C & B & -C \\
0 & 0 & 0 & \dots & 0 & -C & B
\end{pmatrix}
\begin{pmatrix}
u^1 \\ u^2 \\ u^3 \\ \vdots \\ u^{M-3} \\ u^{M-2} \\ u^{M-1} \end{pmatrix}$$

$$= \underbrace{\begin{pmatrix} g^{1} \\ g^{2} \\ g^{3} \\ \vdots \\ g^{M-3} \\ g^{M-2} \\ g^{M-1} \end{pmatrix}}_{=:r} + h^{-2} \begin{pmatrix} \widetilde{\gamma}^{1} \\ \widetilde{\gamma}^{2} \\ \widetilde{\gamma}^{3} \\ \vdots \\ \widetilde{\gamma}^{M-3} \\ \widetilde{\gamma}^{M-2} \\ \widetilde{\gamma}^{M-1} \end{pmatrix}$$
(1-4)

mit

$$\widetilde{\gamma}^{1} = (\gamma_{1,0} + \gamma_{0,1}, \gamma_{2,0}, \dots, \gamma_{M-2,0}, \gamma_{M-1,0} + \gamma_{M,1}), 
\widetilde{\gamma}^{j} = (\gamma_{0,j}, 0, \dots, 0, \gamma_{M,j}), \quad j = 2, \dots, M-2, 
\widetilde{\gamma}^{M-1} = (\gamma_{0,M-1} + \gamma_{1,M}, \gamma_{2,M}, \dots, \gamma_{M-2,M}, \gamma_{M-1,M} + \gamma_{M,M-1}).$$

Es bleibt nun nachzuweisen, dass das Gleichungssystem (1-3) eine eindeutige Lösung besitzt. Dazu benötigen wir etwas Matrizentheorie.

**1.3 Definition.** Eine Matrix  $A \in \mathbb{R}^{N,N}$  heißt nicht negativ oder monoton, falls  $A_{ij} \geq 0$  für  $i, j \in \{1, \dots, N\}$ . Wir schreiben hierfür einfach  $A \geq 0$ .

Seien  $A, B \in \mathbb{R}^{N,N}$ . Dann ist  $A \leq B \iff B - A \geq 0$ .

- **1.4 Definition.**  $A \in \mathbb{R}^{N,N}$  heißt invers monoton, falls  $A^{-1}$  existiert und  $A^{-1} \geq 0$  gilt.
- **1.5 Definition.** A heißt  $L_0$ -Matrix, falls  $A_{ij} \leq 0$  für  $i, j \in \{1, ..., N\}$  mit  $i \neq j$ .
- **1.6 Definition.** A heißt M-Matrix, falls A eine invers monotone  $L_0$ -Matrix ist.

Bekanntermaßen gilt für  $A \in \mathbb{R}^{N,N}$  die Äquivalenz

$$A \ge 0 \iff (u \le v \Rightarrow Au \le Av, \quad \forall u, v \in \mathbb{R}^N).$$

1.7 Bemerkung. Die Matrix A des klassischen Differenzenverfahrens (vgl. Gleichungssystem (1-4)) ist offensichtlich eine  $L_0$ -Matrix. Um sicherzustellen, dass die Lösung u wohldefiniert ist, müssen wir die Invertierbarkeit von A nachweisen.

Dies geschieht mit

- **1.8 Satz** (*M*-Kriterium). Eine  $L_0$ -Matrix  $A \in \mathbb{R}^{N,N}$  ist genau dann eine *M*-Matrix, wenn ein e > 0,  $e \in \mathbb{R}^N$  existiert mit  $Ae \ge 0$  und der Verbindungseigenschaft: Zu jedem  $i_0 \in \{1, \ldots, N\}$  mit  $(Ae)_{i_0} = 0$  gibt es eine Kette  $i_0, i_1, \ldots, i_r \in \{1, \ldots, N\}$  mit  $(Ae)_{i_r} > 0$  und  $A_{i_{r-1}, i_r} \ne 0$  für alle  $j \in \{1, \ldots, r\}$ .
- **1.9 Bezeichnung.** Ein e > 0 mit  $Ae \ge 0$  und der Verbindungseigenschaft heißt majorisierendes Element für A.

Wir betrachten jetzt wieder das Gleichungssystem  $A^h u = r^h$  (vgl. (1-4)) des klassischen Differenzenverfahrens zu  $-\Delta u = g$  in  $\Omega = (0, 1)^2$ ,  $u = \gamma$  auf  $\partial\Omega$ .

**1.10 Lemma.** Die Matrix  $A^h \in \mathbb{R}^{(M-1)^2,(M-1)^2}$  des klassischen Differenzenverfahrens ist eine M-Matrix.

Beweis: Offensichtlich ist  $A^h$  eine  $L_0$ -Matrix. Ferner gilt für  $\mathbb{I}=(1,\ldots,1)^T\in\mathbb{R}^{(M-1)^2}$ 

$$A^h \mathbb{I} = h^{-2} \left( egin{array}{c} \delta^{(2)} \\ \delta^{(1)} \\ \vdots \\ \delta^{(1)} \\ \delta^{(2)} \end{array} 
ight),$$

wobei

$$\delta^{(2)} = \begin{pmatrix} 2 \\ 1 \\ \vdots \\ 1 \\ 2 \end{pmatrix} \in \mathbb{R}^{M-1}, \quad \delta^{(1)} = \begin{pmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 1 \end{pmatrix} \in \mathbb{R}^{M-1}.$$

Genauer gilt

$$A^{h}\mathbb{I}(x,y) = 0 \text{ für } (x,y) \in \overset{\circ}{\Omega}_{h},$$
  
 $A^{h}\mathbb{I}(x,y) > 0 \text{ für } (x,y) \in \Omega_{h} \setminus \overset{\circ}{\Omega}_{h}.$ 

Wir müssen also zu jedem  $(x_0, y_0) \in \overset{\circ}{\Omega}_h$  eine Kette  $(x_0, y_0), \ldots, (x_r, y_r) \in \Omega_h$  finden mit  $A_{(x_{i-1}, y_{i-1}), (x_i, y_i)} \neq 0$  für  $i = 1, \ldots, r$  und  $(x_r, y_r) \in \Omega_h \setminus \overset{\circ}{\Omega}_h$ . Gemäß Definition von  $A^h$  mittels  $B^h$  und  $C^h$  ist dies aber immer möglich.

**1.11 Bemerkung.** a) Da  $A^h$  eine M-Matrix ist, können wir das lineare Gleichungssystem (1-3) mit dem Gauß-Algorithmus ohne Pivotisierung auflösen. Dabei ist zu berücksichtigen, dass  $A^h$  eine Bandmatrix mit der Bandweite 2(M-1)+1=2M-1 ist. Die Elimination einer Bandmatrix der Dimension  $(M-1)^2$  mit der Bandweite 2M-1 erfordert etwa  $\frac{1}{2}(M-1)^2(2M-1)^2\sim 2M^4$  Multiplikationen. Der Aufwand steigt also mit M sehr stark an.

b) Nach dem M-Kriterium folgt überdies  $(A^h)^{-1} \ge 0$ . Somit erhalten wir, dass die Lösung  $u^h$  von (1-4)  $u^h \ge 0$  erfüllt, falls  $(g, \gamma) \ge 0$ .

Wir behandeln nun das Problem eines krummlinigen Randes. Dazu setzen wir

$$\Omega_h = \Omega \cap \mathbb{R}^2_h$$

und ordnen jedem Punkt  $(x,y) \in \Omega_h$  vier Nachbarpunkte  $N_k = N_k(x,y,h) \in \overline{\Omega}$ , k = 1, 2, 3, 4 zu.

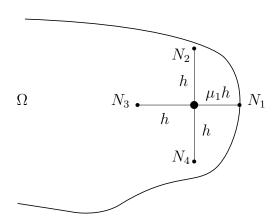


Abbildung 2: Nachbarpunkte

Im Folgenden sei

$$e^{1} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, e^{2} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, e^{3} = \begin{pmatrix} -1 \\ 0 \end{pmatrix}, e^{4} = \begin{pmatrix} 0 \\ -1 \end{pmatrix},$$
$$\mu_{k} = \mu_{k}(x, y, h) = \sup \left\{ \mu \in [0, 1] \mid (x, y) + the^{k} \in \Omega, \forall t \in [0, \mu] \right\},$$
$$N_{k} = N_{k}(x, y, h) = (x, y) + \mu_{k}(x, y, h)he^{k}, \quad k = 1, 2, 3, 4.$$

Die Menge der inneren Gitterpunkte ist

$$\overset{\circ}{\Omega}_{h} = \{ (x, y) \in \Omega_{h} \mid N_{k}(x, y, h) \in \Omega, \quad k = 1, 2, 3, 4 \}.$$

An einem inneren Gitterpunkt (x, y) gilt überdies  $\mu_k(x, y, h) = 1, k = 1, 2, 3, 4$ .

Wir können also die Differentialgleichung  $-\Delta u = g$  in  $\Omega$ ,  $u = \gamma$  auf  $\partial \Omega$  an inneren Gitterpunkten  $(x, y) \in \stackrel{\circ}{\Omega}_h$  wieder durch

$$h^{-2}(-u(x-h,y)-u(x+h,y)-u(x,y-h)-u(x,y+h)+4u(x,y))=q(x,y)$$
 (1-5)

diskretisieren.

Die Menge

$$\Omega_h^R = \Omega_h \setminus \mathring{\Omega}_h = \{(x,y) \in \Omega_h \mid N_k(x,y,h) \in \partial \Omega \text{ für ein } k = 1,2,3,4\}$$

enthält die randnahen Punkte.

Zur Diskretisierung von  $-\Delta u(x,y) = g(x,y), (x,y) \in \Omega_h^R$  benötigen wir zunächst eine Formel für die gewöhnliche Ableitung v'' bei nicht äquidistanten Knoten.

**1.12 Lemma.** Sei  $v \in C^3([-a, a], \mathbb{R})$  für ein a > 0. Dann gilt für alle  $\mu_0, \mu_1 \in (0, 1]$  und  $h \leq a$ 

$$\left| \frac{2}{\mu_0 \mu_1(\mu_0 + \mu_1)h^2} \left\{ \mu_0 v(\mu_1 h) - (\mu_0 + \mu_1)v(0) + \mu_1 v(-\mu_0 h) \right\} - v''(0) \right|$$

$$\leq \frac{2}{3} h \cdot \max\{|v'''(x)| \mid |x| \leq a\}.$$

Mit Hilfe von Lemma 1.12 können wir  $-\Delta u(x,y)=g(x,y),\ (x,y)\in\Omega_h^R$  ersetzen durch

$$g(x,y) = h^{-2} \left\{ -\frac{2}{\mu_1(\mu_1 + \mu_3)} u(N_1) - \frac{2}{\mu_3(\mu_1 + \mu_3)} u(N_3) - \frac{2}{\mu_2(\mu_2 + \mu_4)} u(N_2) - \frac{2}{\mu_4(\mu_2 + \mu_4)} u(N_4) + 2\left(\frac{1}{\mu_1\mu_3} + \frac{1}{\mu_2\mu_4}\right) u(x,y) \right\}.$$
(1-6)

Für jeden Nachbarn  $N_k \in \partial\Omega$ ,  $k \in \{1,2,3,4\}$  ist dabei der Wert  $u(N_k) = \gamma(N_k)$  gemäß  $u = \gamma$  auf  $\partial\Omega$  einzusetzen.

Das hierdurch definierte Verfahren heißt klassisches Differenzenverfahren. Es liefert ein lineares Gleichungssystem

$$Au = r, \quad u \in \mathbb{R}^{\Omega_h}. \tag{1-7}$$

Will man (1-7) mit einem Bandalgorithmus lösen, so sind zur Aufstellung von A die Gitterpunkte durchzunummerieren.

Eine mögliche Nummerierung ist zeilenweise von links unten nach rechts oben, d.h. (x, y) kommt vor  $(\widetilde{x}, \widetilde{y})$ , falls  $y < \widetilde{y}$  oder  $(x < \widetilde{x}, y = \widetilde{y})$ .

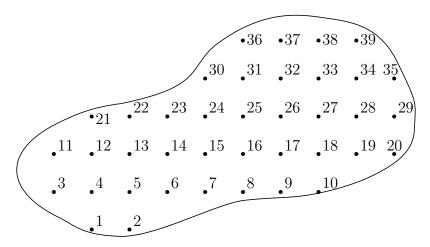


Abbildung 3: Nummerierung der Gitterpunkte

Haben wir statt  $-\Delta u(x,y) = g(x,y)$  in  $\Omega$  die Aufgabe

$$-\Delta u(x,y) = g(x,y,u(x,y)), \quad (x,y) \in \Omega,$$

so ersetzen wir g(x,y) in den Differenzengleichungen (1-5) und (1-6) einfach durch g(x,y,u(x,y)) und erhalten ein nichtlineares Gleichungssystem

$$A^h u = G^h(u), \quad u \in \mathbb{R}^{\Omega_h}$$

mit einem Diagonalfeld

$$(G^h(u))(x,y) = g(x,y,u(x,y)) + \text{Randterme}, \quad (x,y) \in \Omega_h.$$

### b) Konsistenz und Stabilität des klassischen Differenzenverfahrens für die Dirichletsche Randwertaufgabe

Vorgelegt sei

$$-\Delta u(x,y) = g(x,y), \quad (x,y) \in \Omega,$$
  

$$u(x,y) = \gamma(x,y), \quad (x,y) \in \partial\Omega,$$
(1-8)

wobei  $\Omega \subset \mathbb{R}^2$  ein beschränktes Gebiet ist.

Das klassische Differenzenverfahren auf  $\Omega_h = \Omega \cap \mathbb{R}^2_h$  ist für alle  $(x,y) \in \Omega_h$  durch die Formelzeile (1-6) gegeben. Da sich gemäß Lemma 1.12 für  $(x,y) \in \Omega_h^R$  nur der Konsistenzfehler O(h) ergibt, multiplizieren wir diese Gleichungen mit h durch und schreiben das dann entstehende Gleichungssystem in der Form

$$A^h u = r^h, \quad u \in \mathbb{R}^{\Omega_h} \tag{1-9}$$

bzw.

$$T^h u = A^h u - r^h = 0, \quad T^h : \mathbb{R}^{\Omega_h} \longrightarrow \mathbb{R}^{\Omega_h}.$$

Sei h > 0 und sei  $\Omega_h = \Omega \cap \mathbb{R}^2_h$ . Zu  $u \in \mathbb{R}^{\Omega_h}$  sei

$$||u||_{\infty} = \max\{|u(x,y)| \mid (x,y) \in \Omega_h\}.$$

**1.13 Definition.** Sei  $W \subset C^2(\overline{\Omega})$ . Das numerische Modell  $T^h(u) = A^h u - r^h = 0$ ,  $u \in \mathbb{R}^{\Omega_h}$  heißt W-konsistent, falls jede Lösung  $\overline{u} \in W$  der Randwertaufgabe (1-8)

$$||T^h(\overline{u}_h)||_{\infty} \longrightarrow 0$$
 für  $h \longrightarrow 0$ 

erfüllt. Das Verfahren heißt W-konsistent der Ordnung p, falls überdies

$$||T^h(\overline{u}_h)||_{\infty} = O(h^p)$$
 für  $h \longrightarrow 0$ 

gilt. Dabei bezeichnet  $\overline{u}_h$  die Restriktion von  $\overline{u}$  auf  $\Omega_h$ .

**1.14 Definition.** Das Modell  $T^h(u) = 0$  heißt W-konvergent, falls es zu jeder Lösung  $\overline{u} \in W$  der Randwertaufgabe (1-8) so ein  $h_0 > 0$  gibt, dass

$$T^h(u) = 0$$

für  $0 < h \le h_0$  eine Lösung  $u^h \in \mathbb{R}^{\Omega_h}$  besitzt mit

$$\|\overline{u}_h - u^h\|_{\infty} \longrightarrow 0$$
 für  $h \longrightarrow 0$ .

Gilt überdies  $\|\overline{u}_h - u^h\|_{\infty} = O(h^p)$ , so heißt das Modell W-konvergent der Ordnung p.

**1.15 Definition.** Das Modell  $T^h(u) = 0$ ,  $u \in \mathbb{R}^{\Omega_h}$  heißt stabil bezüglich h, falls es von h unabhängige Konstanten  $h_0 > 0$  und C > 0 derart gibt, dass die Stabilitätsungleichung

$$||u - v||_{\infty} \le C||T^h(u) - T^h(v)||_{\infty}$$

für alle  $u, v \in \mathbb{R}^{\Omega_h}$ ,  $0 < h \le h_0$  erfüllt ist.

**1.16 Bemerkung.** Im linearen Fall  $T^h(u) = A^h u - r^h$  ist dies gleichbedeutend zu

$$||u||_{\infty} \le C||A^h u||_{\infty}, \quad \forall u \in \mathbb{R}^{\Omega_h}, \quad 0 < h \le h_0.$$

**1.17 Satz.** Sei das Modell W-konsistent oder W-konsistent der Ordnung p und stabil. Es existiere eine Lösung  $u^h \in \mathbb{R}^{\Omega_h}$  von  $T^h(u) = 0$ . Dann ist das Modell auch W-konvergent bzw. W-konvergent der Ordnung p.

Beweis: Wir setzen  $u=u^h$  und  $v=\overline{u}_h$  in die Stabilitätsungleichung ein und finden

$$||u^{h} - \overline{u}_{h}||_{\infty} \leq C||\underbrace{T^{h}(u^{h})}_{=0} - T(\overline{u}_{h})||_{\infty}$$

$$= C \cdot ||T^{h}(\overline{u}_{h})||_{\infty} \longrightarrow 0 \quad \text{für} \quad h \longrightarrow 0$$

bzw.

$$||u^h - \overline{u}_h||_{\infty} = C \cdot ||T^h(\overline{u}_h)||_{\infty} = O(h^p)$$
 für  $h \longrightarrow 0$ .

**1.18 Satz.** An jeder Lösung  $\overline{u} \in C^4(\overline{\Omega})$  von (1-8) liegt Konsistenz des klassischen Differenzenverfahrens (1-9) der Ordnung 2 vor, d.h.

$$||A^h \overline{u}_h - r^h||_{\infty} = \sup\{|(A^h \overline{u}_h)(x, y) - r^h(x, y)| \mid (x, y) \in \Omega_h\} = O(h^2).$$

 $(d.h. C^4$ -Konsistenz der Ordnung 2)

Beweis: Für  $(x,y) \in \overset{\circ}{\Omega}_h$  gilt mit  $T^h(u) = A^h u - r^h$  nach Formel (1-2)  $|(A^h \overline{u}_h - r^h)(x,y)| = |h^{-2} \left( -\overline{u}(x-h,y) - \overline{u}(x+h,y) - \overline{u}(x,y-h) - \overline{u}(x,y+h) + 4\overline{u}(x,y) \right) - \underbrace{g(x,y)}_{=-\Delta \overline{u}(x,y)}|$   $\leq |h^{-2} (-\overline{u}(x-h,y) + 2\overline{u}(x,y) - \overline{u}(x+h,y)) + \overline{u}_{xx}(x,y)|$   $+ |h^{-2} (-\overline{u}(x,y-h) + 2\overline{u}(x,y) - \overline{u}(x,y+h)) + \overline{u}_{yy}(x,y)|$   $\leq C \cdot h^2 \left( \max\{|\overline{u}_{xxxx}(x+\xi h,y)| \mid |\xi| \leq 1 \right)$ 

mit einer von h und (x, y) unabhängigen Konstanten  $\widetilde{C}$ .

Für  $(x,y) \in \Omega_h^R$  setzen wir

$$v(\xi) = \overline{u}(x+\xi,y),$$
  
 $w(\eta) = \overline{u}(x,y+\eta)$ 

 $+\max\{|\overline{u}_{uuuy}(x,y+\eta h)|\,|\,|\eta|\leq 1\}\} \leq \widetilde{C}h^2$ 

und finden damit die Abschätzung

$$|(A^{h}\overline{u}_{h} - r^{h})(x, y)| = h \left| h^{-2} \left\{ \frac{-2}{\mu_{1}(\mu_{1} + \mu_{3})} \overline{u}(N_{1}) - \frac{2}{\mu_{3}(\mu_{1} + \mu_{3})} \overline{u}(N_{3}) - \frac{2}{\mu_{2}(\mu_{2} + \mu_{4})} \overline{u}(N_{2}) - \frac{2}{\mu_{4}(\mu_{2} + \mu_{4})} \overline{u}(N_{4}) + 2 \left( \frac{1}{\mu_{1}\mu_{3}} + \frac{1}{\mu_{2}\mu_{4}} \right) \overline{u}(x, y) \right\} - \underbrace{g(x, y)}_{=-\Delta \overline{u}(x, y)} \right|.$$

© Johannes Schropp 5. Februar 2025

Unter Berücksichtigung von

$$\overline{u}(N_1) = \overline{u}(x + \mu_1 h, y) = v(\mu_1 h), 
\overline{u}(N_3) = \overline{u}(x - \mu_3 h, y) = v(-\mu_3 h), \quad v''(0) = \overline{u}_{xx}(x, y), 
\overline{u}(N_2) = \overline{u}(x, y + \mu_2 h) = w(\mu_2 h), 
\overline{u}(N_4) = \overline{u}(x, y - \mu_4 h) = w(-\mu_4 h), \quad w''(0) = \overline{u}_{yy}(x, y)$$

erhält man

$$|(A^{h}\overline{u}_{h} - r^{h})(x,y)| \leq h \left| h^{-2} \left\{ \frac{-2}{\mu_{1}(\mu_{1} + \mu_{3})} v(\mu_{1}h) - \frac{2}{\mu_{3}(\mu_{1} + \mu_{3})} v(-\mu_{3}h) + \frac{2}{\mu_{1}\mu_{3}} v(0) \right\} + v''(0) \right|$$

$$+ h \left| h^{-2} \left\{ \frac{-2}{\mu_{2}(\mu_{2} + \mu_{4})} w(\mu_{2}h) - \frac{2}{\mu_{4}(\mu_{2} + \mu_{4})} w(-\mu_{4}h) + \frac{2}{\mu_{2}\mu_{4}} w(0) \right\} + w''(0) \right|$$

$$\leq \frac{2h^{2}}{3} \left( \max\{ |\overline{u}_{xxx}(x + \xi h, y)| - \mu_{3} \leq \xi \leq \mu_{1} \right\}$$

$$\max\{ |\overline{u}_{yyy}(x, y + \eta h)| - \mu_{4} \leq \eta \leq \mu_{2} \} \leq Ch^{2}.$$

**1.19 Satz.** Die Matrix  $A^h$  des klassischen Differenzenverfahrens ist eine M-Matrix. Wählt man ein Rechteck  $(a,b) \times (c,d) \supset \overline{\Omega}$ , so gilt für den positiven Vektor  $e_h = e|_{\Omega_h} \in \mathbb{R}^{\Omega_h}$  mit

$$e(x,y) = (x-a)(b-x) + (y-c)(d-y), \quad (x,y) \in \mathbb{R}^2$$

die Ungleichung

$$A^h e_h \ge \rho e_h$$

 $f\ddot{u}r\ \rho := \frac{16}{(b-a)^2 + (d-c)^2} > 0$ , falls h > 0 hinreichend klein gewählt ist.

A<sup>h</sup> erfüllt die Stabilitätsungleichung

$$||u||_{\infty} \le \frac{1}{\rho} ||A^h u||_{\infty} \quad \forall u \in \mathbb{R}^{\Omega_h}$$

für hinreichend kleine h > 0.

Beweis: Wir wissen, dass  $A^h$  eine  $L_0$ -Matrix ist.

Für  $(x,y) \in \stackrel{\circ}{\Omega}_h$  ist die in (1-6) verwandte Differenzenformel exakt für Polynome bis zum Grad 3, d.h. für e(x,y) = (x-a)(b-x) + (y-c)(d-y) gilt

$$(A^h e_h)(x,y) = h^{-2}(-e(x-h,y) - e(x+h,y) - e(x,y-h) - e(x,y+h) + 4e(x,y))$$

$$= -(\Delta e)(x,y) = -(-2-2) = 4$$

für alle  $(x,y) \in \stackrel{\circ}{\Omega}_h$ .

Für  $(x, y) \in \Omega_h^R$  werden durch die Differenzenformel (1-6) noch Polynome bis zum Grad 2 exakt differenziert (vgl. Lemma 1.12).

Unter Beachtung der Formelzeile (1-6) sowie der Tatsache  $u(N_k) = \gamma(N_k)$ , falls  $N_k \in \partial\Omega$ , folgt also

$$(A^{h}e_{h})(x,y) = h \left[ -\underbrace{\Delta e(x,y)}_{=-4} + h^{-2} \sum_{\substack{k=1 \\ N_{k}(x,y,h) \in \partial \Omega}}^{4} \frac{2}{\mu_{k}(\mu_{k} + \mu_{\widetilde{k}})} e(N_{k}) \right]$$

 $mit \ \widetilde{k} = (k+2) \mod 4.$ 

Nach Voraussetzung  $(x,y) \in \Omega_h^R$  ist die letztere Summe nicht leer und  $\mu_k, \mu_{\widetilde{k}} \in (0,1]$ , d.h.  $0 < \mu_k(\mu_k + \mu_{\widetilde{k}}) \le 2$ ,  $\frac{2}{\mu_k(\mu_k + \mu_{\widetilde{k}})} \ge 1$  und somit

$$(A^h e_h)(x,y) \ge h(4 + h^{-2} \underbrace{\min\{e(x,y) \mid (x,y) \in \overline{\Omega}\}}_{=:e_0}).$$

Da e stetig ist und e(x,y) > 0 für  $(x,y) \in \overline{\Omega}$  gilt, folgt  $e_0 > 0$ .

Wir finden also

$$(A^h e_h)(x,y) \ge h^{-1} e_0$$

für  $(x,y) \in \Omega_h^R$ .

Insgesamt erhalten wir

$$A^h e_h \ge \min\left(4, \frac{e_0}{h}\right) \mathbb{I} > 0, \quad e_h > 0,$$

d.h.  $A^h$  ist eine M-Matrix.

Für  $h \leq \frac{e_0}{4}$  folgt

$$A^h e_h \ge 4 \, \mathbb{I} \ge 4 \frac{e_h}{\|e_h\|_{\infty}} \ge \frac{4}{\left(\frac{b-a}{2}\right)^2 + \left(\frac{d-c}{2}\right)^2} e_h =: \rho e_h$$
 (1-10)

Da  $A^h$  eine M-Matrix ist, ergibt sich

$$\|(A^h)^{-1}\|_{\infty} = \|(A^h)^{-1}\mathbb{I}\|_{\infty} \le \frac{1}{4}\|e_h\|_{\infty} \le \frac{1}{\rho}.$$

Somit folgt die Stabilitätsungleichung

$$||u||_{\infty} = ||(A^h)^{-1}A^h u||_{\infty} \le ||(A^h)^{-1}||_{\infty}||A^h u||_{\infty} \le \frac{1}{\rho}||A^h u||_{\infty}.$$
 (1-11)

Aus Satz 1.18 und Satz 1.19 erhalten wir

**1.20 Korollar.** Das klassische Differenzenverfahren für die Dirichletsche Randwertaufgabe (1-8) in einem beschränkten Gebiet  $\Omega$  ist  $C^4(\overline{\Omega})$ -konvergent der Ordnung 2 bzgl.  $\|\cdot\|_{\infty}$ .

Beweis: Ist  $u^h$  die Lösung des klassischen Differenzenverfahrens, so gilt

$$\|\overline{u}_h - u^h\|_{\infty} \le \frac{1}{\rho} \|A^h \overline{u}_h - \underbrace{A^h u^h}_{=r^h}\|_{\infty} = \frac{1}{\rho} \|A^h \overline{u}_h - r^h\|_{\infty} = O(h^2).$$

**1.21 Bemerkung.** Dass die Lösung  $\overline{u}$  in  $C^4(\overline{\Omega}, \mathbb{R})$  liegt, ist für Gebiete mit "Ecken" im Allgemeinen nicht erfüllt, z.B. im Fall  $-\Delta u = 1$  in  $\Omega = (0, 1)^2$ , u = 0 auf  $\partial\Omega$ . Wäre  $\overline{u} \in C^4(\overline{\Omega}, \mathbb{R})$  Lösung hiervon, so würde gelten

$$-\Delta u(0,0) = 1 = -u_{xx}(0,0) - u_{yy}(0,0) = 0,$$

da u = 0 auf  $\partial \Omega$ .

Wir betrachten jetzt kurz das klassische Differenzenverfahren für die nichtlineare Aufgabe

$$-\Delta u(x,y) = g(x,y,u(x,y)) \text{ in } \Omega,$$
  

$$u(x,y) = \gamma \text{ auf } \partial \Omega.$$
(1-12)

Es lautet

$$g(x, y, u(x, y)) = h^{-2} \left\{ \frac{-2}{\mu_1(\mu_1 + \mu_3)} u(N_1) - \frac{2}{\mu_3(\mu_1 + \mu_3)} u(N_3) - \frac{2}{\mu_2(\mu_2 + \mu_4)} u(N_2) - \frac{2}{\mu_4(\mu_2 + \mu_4)} u(N_4) + 2\left(\frac{1}{\mu_1\mu_3} + \frac{1}{\mu_2\mu_4}\right) u(x, y) \right\}.$$

Dabei beachten wir  $u(N_k) = \gamma(N_k)$  für  $(x, y) \in \Omega_h = \Omega \cap \mathbb{R}_h^2$ , falls  $N_k \in \partial \Omega$ .

Dies liefert ein nichtlineares System der Form

$$T^{h}u = A^{h}u - G^{h}(u) = 0 (1-13)$$

mit einem Diagonalfeld

$$G^h(u)(x,y) = g(x,y,u(x,y)) +$$
Randterme.

Das Gleichungssystem (1-13) lässt sich nun z.B. mit dem Newtonverfahren lösen. Präziser: Man wähle ein  $u^0$  und löse für  $i = 0, 1, 2, \ldots$  das lineare Gleichungssystem

$$DT^h(u^i)d^i = -T^h(u^i)$$

und setze  $u^{i+1} = u^i + d^i$ ,  $i = 0, 1, 2, \dots$ 

Für den Fall  $T^h(u) = A^h u - G^h(u)$  folgt

$$DT^h(u) = A^h - DG^h(u)$$

und mit  $\frac{\partial g}{\partial v}(x,y,v) \le \mu < \rho = \frac{16}{(b-a)^2 + (d-c)^2}$  folgt mit e aus Satz 1.19 und Formel (1-10)

$$DT^{h}(u)e_{h} = A^{h}e_{h} - \underbrace{DG^{h}(u)}_{<\mu I^{h}}e_{h} \ge \rho e_{h} - \mu e_{h} = (\rho - \mu)e_{h} > 0,$$

falls h > 0 hinreichend klein.

Da  $DT^h(u)$  eine  $L_0$ -Matrix ist, ist  $e_h$  ein majorisierendes Element für  $DT^h(u)$ , d.h.  $DT^h(u)$  ist M-Matrix und das Newtonverfahren ist in dieser Situation durchführbar.

Für ein beliebiges Gleichungssystem der Form T(u) = 0 erinnern wir an den lokalen Konvergenzsatz des Newtonverfahrens.

**1.22 Satz** (lokaler Konvergenzsatz, ohne Beweis). Sei  $U \subset \mathbb{R}^N$  offen und sei  $T: U \longrightarrow \mathbb{R}^N$  zweimal stetig differenzierbar. Ferner existiere eine Lösung  $\overline{u} \in U$  von T(u) = 0, und  $DT(\overline{u})$  sei invertierbar. Dann gibt es eine Kugel

$$K_{\rho}(\overline{u}) := \{ u \in \mathbb{R}^N \mid ||u - \overline{u}|| < \rho \} \subset U, \quad \rho > 0$$

so, dass T(u) = 0 keine weitere Lösung in  $K_{\rho}(\overline{u})$  besitzt. Für jeden Startwert  $u^0 \in K_{\rho}(\overline{u})$  existiert die Newtonfolge

$$u^{n+1} = u^n - DT(u^n)^{-1}T(u^n),$$

liegt in  $K_{\rho}(\overline{u})$  und konvergiert gegen  $\overline{u}$ .

Überdies gibt es ein C > 0 mit

$$||u^{n+1} - \overline{u}||_{\infty} \le C||u^n - \overline{u}||_{\infty}^2$$

für alle  $n \geq 0$ ,  $u^0 \in K_{\rho}(\overline{u})$ . Letztere Ungleichung bedeutet lokal quadratische Konvergenz.

**1.23 Bemerkung.** Eine Nullstelle  $\overline{u} \in U$  von T(u) = 0 mit invertierbarem  $DT(\overline{u})$  heißt regulär.

Typische Abbruchkriterien für Newtonverfahren sind

- $i) ||T(u^i)||_{\infty} \leq \text{TOL}, \text{ z.B. TOL} = 10^{-10}, \text{oder/und}$
- *ii*)  $\|u^{i+1} u^i\|_{\infty} \le EPS$ , z.B.  $EPS = 10^{-5}$ .

Die angegebenen Werte beziehen sich auf eine Rechnerarithmetik mit 16 Stellen.

**1.24 Satz** (ohne Beweis). Vorgelegt sei (1-12), und es gelte  $g \in C^1(\overline{\Omega} \times \mathbb{R}, \mathbb{R})$  sowie

$$\frac{\partial g}{\partial u}(x,y,u) \le \mu < \rho := \frac{16}{(b-a)^2 + (d-c)^2},$$

 $\overline{\Omega} \subset (a,b) \times (c,d)$ . Dann ist das klassische Differenzenverfahen (1-13) für die Aufgabe (1-12)  $C^4(\overline{\Omega})$ -konvergent der Ordnung 2.

Vor- und Nachteile der Differenzenverfahren lauten also:

- Differenzenverfahren sind einfach aufzustellen, da lediglich Ableitungen durch Differenzenquotienten ersetzt werden.
- Die bei äquidistantem Gitter entstehenden Gleichungssysteme haben sehr viel Struktur und lassen sich effizient lösen.
- $\ominus$  Die Differenzengleichungen werden kompliziert bei nicht äquidistanten Stützstellen und komplexen Gebieten.
- $\ominus$  Wir brauchen Lösungen in  $C^4(\overline{\Omega})$  um Konvergenz der Ordnung 2 zu sichern.

### 2. Finite Elemente Methoden für elliptische Differentialgleichungen

### a) Theoretische Grundlagen zur Finite Elemente Methode

**2.1 Definition.** Es sei  $\Omega \subset \mathbb{R}^2$  ein Gebiet. Wir definieren

$$L^2(\Omega) := \{u : \Omega \longrightarrow \mathbb{R} \mid |u|^2 \text{ Lebesgue-integrierbar} \}.$$

 $u, v \in L^2(\Omega)$  sind gleich in  $L^2(\Omega)$  genau dann, wenn u(x) = v(x) für fast alle  $x \in \Omega$  gilt.

Versehen mit dem Skalarprodukt

$$\langle u, v \rangle_0 = \int_{\Omega} u(x)v(x)dx$$

und der dadurch erzeugten Norm

$$||u||_{L^2} = \sqrt{\langle u, u \rangle} = \left(\int_{\Omega} |u(x)|^2 dx\right)^{1/2}$$

wird der Raum  $(L^2(\Omega), \langle \cdot, \cdot \rangle_0)$  zum Hilbertraum.

**2.2 Definition.** Zu  $u:\Omega\longrightarrow\mathbb{R}$  heißt

$$\operatorname{supp}(u) = \overline{\{x \in \Omega \mid u(x) \neq 0\}}$$

Träger von u.

#### 2.3 Definition. Der Raum

$$C_0^{\infty}(\Omega) = \{ u \in C^{\infty}(\Omega) \mid \text{supp}(u) \text{ kompakt} \}$$

heißt Raum der Testfunktionen.

Dieser Raum ist nicht leer, denn der sogenannte Friedrichs'sche Glättungskern  $u_{\varepsilon}(\cdot - x_0)$  mit

$$u_{\varepsilon}(x) = \begin{cases} \exp\left(-\frac{\varepsilon^2}{\varepsilon^2 - \|x\|_2^2}\right), & \|x\|_2 < \varepsilon \\ 0, & \text{sonst} \end{cases} \quad (\varepsilon > 0)$$

ist in  $C_0^{\infty}(\Omega)$  für  $x_0 \in \Omega$  und hinreichend kleine  $\varepsilon > 0$ .

Ferner benötigen wir noch den Begriff der schwachen Ableitung.

**2.4 Definition.** Sei  $u \in L^2(\Omega)$ , und sei  $\alpha = (\alpha_1, \alpha_2) \in \mathbb{N}_0^2$  ein Multiindex. u besitzt die schwache Ableitung  $v := \partial^{\alpha} u \in L^2(\Omega)$ , wenn für alle  $\varphi \in C_0^{\infty}(\Omega)$  gilt

$$\int_{\Omega} v \cdot \varphi \, dx = (-1)^{|\alpha|} \int_{\Omega} u \partial^{\alpha} \varphi \, dx.$$

Offensichtlich besitzt  $u \in C^k(\Omega)$  schwache Ableitungen  $\partial^{\alpha} u$  für  $|\alpha| := \alpha_1 + \alpha_2 \le k$ , und diese stimmen mit den klassischen Ableitungen überein.

**2.5 Definition.** Sei  $m \in \mathbb{N}_0$ . Dann heißt der Raum

$$H^{m}(\Omega) = \{ u \in L^{2}(\Omega) \mid \partial^{\alpha} u \in L^{2}(\Omega) \quad \forall \alpha \in \mathbb{N}_{0}^{2} : |\alpha| \leq m \}$$

Sobolev-Raum.

Der Sobolev-Raum  $H^m(\Omega)$  versehen mit dem Skalarprodukt

$$\langle u, v \rangle_{H^m} = \sum_{|\alpha| \le m} \int_{\Omega} \partial^{\alpha} u \cdot \partial^{\alpha} v \, dx = \sum_{|\alpha| \le m} \langle \partial^{\alpha} u, \partial^{\alpha} v \rangle_{0}$$

ist ein Hilbert-Raum.

Es stellt sich nun die Frage, ob  $H^m$ -Funktionen stetig sind. Um diese Frage positiv beantworten zu können, benötigt man Voraussetzungen an den Rand  $\partial\Omega$ .

**2.6 Definition.** Sei  $\Omega \subset \mathbb{R}^n$  beliebig.  $\Omega$  besitzt die gewöhnliche Kegeleigenschaft, wenn es einen endlichen Kegel C so gibt, dass es zu jedem  $x \in \partial \Omega$  einen zu C kongruenten Kegel  $C_x$  gibt, dessen Spitze in x ist und der ganz in  $\Omega$  liegt.

Kongruenz bedeutet dabei, dass C und  $C_x$  durch Kongruenzabbildung ineinander überführbar sind. Kongruenzabbildungen sind Parallelverschiebung, Drehung, Spiegelung und die Vernüpfung dieser Abbildungen.

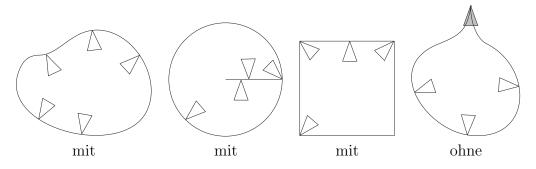


Abbildung 4: Beispiel für Gebiete mit und ohne Kegeleigenschaft

**2.7 Satz** (Sobolevscher Einbettungssatz). Sei  $\Omega \subset \mathbb{R}^d$  ein Gebiet mit gewöhnlicher Kegeleigenschaft. Ist  $m > \frac{d}{2}$ , so gilt

$$H^m(\Omega) \subset C(\overline{\Omega})$$

und die Einbettungsabbildung  $H^m(\Omega) \hookrightarrow C(\overline{\Omega})$  ist stetig.

Also ist für d = 2, 3 die Einbettung stetig, falls  $m \ge 2$ .

Ferner benötigen wir die Hilberträume

$$H_0^m(\Omega) = \{ u \in H^m(\Omega) \mid u = 0 \text{ auf } \partial\Omega \}.$$

Da für Gebiete  $\Omega \subset \mathbb{R}^2$  die Elemente aus  $H^1_0(\Omega)$  nicht unbedingt stetig sein müssen, ist "u=0 auf  $\partial\Omega$ " möglicherweise nicht definiert. Der Ausdruck "u=0 auf  $\partial\Omega$ " muss in einem schwachen Sinne erklärt werden.

**2.8 Definition.** Ein Gebiet  $\Omega \in \mathbb{R}^n$  besitzt die strikte Kegeleigenschaft, falls es eine lokal (d.h. für alle Kompakta) endliche offene Überdeckung  $\{\theta_i\}$  von  $\partial\Omega$  mit zugehörigen Kegeln  $\{K_i\}$  gibt mit

$$\forall x \in \Omega \cap \theta_i : \quad x + K_i \subset \Omega.$$

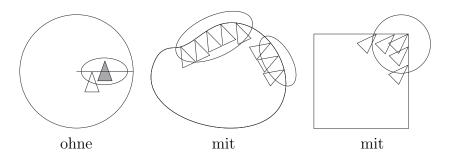


Abbildung 5: Beispiel für Gebiete mit der strikten und ohne strike Kegeleigenschaft

**2.9 Satz.** Sei  $\Omega \subset \mathbb{R}^2$  ein beschränktes Gebiet mit strikter Kegeleigenschaft  $\partial\Omega$ . Dann gibt es eindeutig einen linearen stetigen Operator

$$\Gamma: H^1(\Omega) \longrightarrow L^2(\partial\Omega)$$

 $mit \ \Gamma(u)=u|_{\partial\Omega} \ f\ddot{u}r \ alle \ u\in C^1(\overline{\Omega}).$ 

Man nennt  $\Gamma(u)$  die verallgemeinerten Randwerte von u oder den Spuroperator. Häufig schreibt man "u = 0 auf  $\partial\Omega$ " anstelle von  $\Gamma(u) = 0$ .

Es gilt also

$$H_0^1(\Omega) = \{ u \in H^1(\Omega) \, | \, \Gamma(u) = 0 \}.$$

Es lässt sich auch die Dichtheit von  $C_0^{\infty}(\Omega)$  in  $H_0^m(\Omega)$  zeigen.

### Die Dirichlet'sche Aufgabe mit homogenen Randbedingungen

Nun wenden wir uns unserem Modellproblem zu. Vorgelegt sei die Aufgabe

$$-\Delta u = g \text{ in } \Omega,$$
  

$$u = 0 \text{ auf } \partial \Omega.$$
 (2-1)

 $\Omega \subset \mathbb{R}^2$  sei dabei ein beschränktes Gebiet mit stückweise glattem Rand  $\partial\Omega$  und der strikten Kegeleigenschaft. Es sei  $g \in C(\overline{\Omega})$  und  $u \in C^2(\Omega) \cap C(\overline{\Omega})$  eine klassische Lösung.

Multipliziert man (2-1) mit  $v \in C_0^{\infty}(\Omega)$  und integriert über  $\Omega$ , so ergibt sich unter Benutzung des Divergenzsatzes

$$\int_{\Omega} g \cdot v \, dx = -\int_{\Omega} \Delta u \cdot v \, dx$$

$$= \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\partial \Omega} \frac{\partial u}{\partial n} \cdot v \, dx$$

$$= \int_{\Omega} \nabla u \cdot \nabla v \, dx \text{ für alle } v \in C_0^{\infty}(\Omega). \tag{2-2}$$

Diese Gleichung gilt wegen der Dichtheit von  $C_0^{\infty}(\Omega)$  und der Stetigkeit von  $\langle g, \cdot \rangle - \langle \nabla u, \nabla \cdot \rangle$  in  $H_0^1(\Omega)$  auch für alle  $v \in H_0^1(\Omega)$ .

Für  $V = H_0^1(\Omega)$  setzen wir

$$a: V \times V \longrightarrow \mathbb{R},$$
  
 $a(u,v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx$ 

und

$$b: V \longrightarrow \mathbb{R},$$

$$b(v) = \int_{\Omega} g \cdot v \, dx.$$

Die schwache Formulierung (Variationsgleichung) von (2-1) lautet

$$a(u,v) = b(v), \quad \forall v \in V = H_0^1(\Omega).$$
 (2-3)

Man nennt (2-3) die Variationsgleichung zu (2-1).

**2.10 Definition.**  $u \in V$  heißt schwache Lösung von (2-1), wenn u die Variationsgleichung (2-3) erfüllt.

Ist umgekehrt u eine schwache Lösung von (2-2) mit  $u \in C^2(\Omega) \cap C(\overline{\Omega})$ , so gilt

$$\int_{\Omega} (g + \Delta u)v \, dx = \int_{\Omega} (g \cdot v - \nabla u \cdot \nabla v) \, dx = 0$$

für alle  $v \in C_0^{\infty}(\Omega)$ .

Da $\Omega\subset\mathbb{R}^2$ beschränkt, folgt $g+\Delta u\in L^2(\Omega)$  und somit

$$\int_{\Omega} (\underbrace{g + \Delta u}_{\in L^2(\Omega)}) v \, dx = 0 \quad \forall v \in C_0^{\infty}(\Omega).$$

Das Fundamentallemma der Variationsrechnung (Satz von de la Vallée-Poisson) liefert dann

$$g + \Delta u = 0$$
 fast überall in  $\Omega$ .

Da  $u \in H_0^1(\Omega)$ , gilt  $0 = \Gamma(u) = u|_{\partial\Omega}$ .

**2.11 Definition.** Sei V ein reeller Hilbertraum und  $a: V \times V \to \mathbb{R}$  eine Bilinearform.

i) a heißt stetig auf V, falls ein K > 0 existiert mit

$$|a(u,v)| \le K \cdot ||u|| \cdot ||v||$$

für alle  $u, v \in V$ .

ii) a heißt koerziv (oder elliptisch) auf V, wenn es ein  $\alpha > 0$  gibt mit

$$a(u,u) \geq \alpha \|u\|^2$$

für alle  $u \in V$ .

**2.12 Lemma.** Die Variationsgleichung a(u,v) = b(v),  $\forall v \in V$ , a symmetrisch und koerziv hat die gleichen Lösungen  $u \in V$  wie die Minimierungsaufgabe (Variationsaufgabe)

$$F(v) = \frac{1}{2}a(v,v) - b(v) = \int_{\Omega} \frac{1}{2}|\nabla v|^2 - g \cdot v \, dx \to \min \ f\ddot{u}r \, v \in V.$$

Das Minimierungsproblem heißt Prinzip der minimalen potentiellen Energie oder das Dirichlet'sche Prinzip.

Beweis: Es gilt mit der Symmetrie von a

$$a(w, w) - a(u, u) = a(w + u, w - u) = 2a(u, w - u) + a(w - u, w - u).$$

Es gelte nun  $a(u,v) = b(v) \quad \forall v \in V$ . Dann folgt

$$F(w) = \frac{1}{2}a(w,w) - b(w)$$

$$= \frac{1}{2}a(u,u) - b(u) + a(u,w-u) - b(w-u) + \frac{1}{2}a(w-u,w-u). (2-4)$$

Setze nun v = w - u in (2-4) ein und erhalte mit a koerziv

$$F(u+v) = F(w) = \frac{1}{2}a(u,u) - b(u) + a(u,v) - b(v) + \frac{1}{2}a(v,v)$$
  
  $\geq F(u) \quad \forall v \in V.$ 

Sei jetzt u die Lösung der Minimierungsaufgabe  $F(v) = \frac{1}{2}a(v,v) - b(v) = min.$ 

Annahme: Es existiere ein  $v \in V$  mit  $a(u,v) \neq b(v)$ . Ohne Einschränkung gelte dabei wegen der Linearität a(u,v) < b(v). Mit w := u + tv, t > 0 und (2-4) folgt

$$F(w) = F(u+tv) = \frac{1}{2}a(u+tv, u+tv) - b(u+tv)$$

$$= \frac{1}{2}a(u, u) - b(u) + a(u, tv) - b(tv) + \frac{1}{2}a(tv, tv)$$

$$= F(u) + t(a(u, v) - b(v)) + \frac{t^2}{2}a(v, v).$$

Damit lässt sich t > 0 hinreichend klein so wählen, dass F(w) = F(u + tv) < F(u) im Widerspruch zu u Lösung des Minimierungsproblems.

Es stellt sich nun die Frage nach den Bedingungen an a und b, unter welchen die Existenz und Eindeutigkeit von Lösungen von (2-3) gesichert werden.

**2.13 Satz** (Lax-Milgram). Seien V ein reeller Hilbertraum,  $a: V \times V \to \mathbb{R}$  eine stetige und koerzive Bilinearform sowie  $b: V \to \mathbb{R}$  ein stetiges lineares Funktional. Dann existiert genau ein  $u \in V$  mit

$$a(u, v) = b(v) \quad \forall v \in V.$$

Für diese Lösung gilt

$$||u|| \le \frac{1}{\alpha} ||b||_{V'}.$$

Dabei ist

$$||b||_{V'} = \sup \left\{ \frac{|b(v)|}{||v||} \mid v \in V, v \neq 0 \right\}.$$

### b) Finite Elemente Methoden

#### Die Dirichlet'sche Aufgabe mit inhomogenen Randbedingungen

Wir betrachten das Problem

$$-\Delta u = g \text{ in } \Omega,$$
  

$$u = \gamma \text{ auf } \partial \Omega.$$
 (2-5)

 $\Omega \subset \mathbb{R}^2$  sei wieder ein beschränktes Gebiet mit stückweise glattem Rand  $\partial\Omega$  mit strikter Kegeleigenschaft. Es gelte  $\gamma:\Omega\longrightarrow\mathbb{R}$  glatt. Dann gibt es zwei Möglichkeiten, die Gleichung (2-5) zu lösen.

i) Transformiere die Aufgabe (2-5) auf homogene Randbedingungen.

Sei  $w = u - \gamma$  eine Lösung von

$$-\Delta w = g + \Delta \gamma \text{ in } \Omega,$$
  
$$w = 0 \text{ auf } \partial \Omega.$$

Dann ist  $u = w + \gamma$  eine Lösung von (2-5), denn

$$-\Delta u = -\Delta(w + \gamma) = -\Delta w - \Delta \gamma = g \text{ in } \Omega,$$
  
$$u|_{\partial\Omega} = \underbrace{w|_{\partial\Omega}}_{=0} + \gamma|_{\partial\Omega} = \gamma|_{\partial\Omega}.$$

ii) Suche eine Funktion u mit  $w = u - \gamma \in V = H_0^1(\Omega)$ .

Die schwache Formulierung des Problems hat dann die Form

$$\int_{\Omega} \nabla (w + \gamma) \cdot \nabla v - g \cdot v dx = 0 \quad \forall v \in V.$$

Das allgemeine Ritz'sche Verfahren zur Lösung der Randwertaufgabe (2-5) lautet jetzt: Wähle Ansatzfunktionen  $u_i \in C^1(\overline{\Omega})$ ,  $u_i|_{\partial\Omega} = 0$ , i = 1, ..., m und finde ein  $\widetilde{u} \in \gamma + V_h$ ,  $V_h = \operatorname{span}\{u_1, ..., u_m\}$  mit

$$\int_{\Omega} \nabla(\widetilde{u}) \cdot \nabla v - g \cdot v dx = 0 \quad \forall v \in V_h.$$
 (2-6)

Wir suchen also  $\widetilde{u}$  in der Form

$$\widetilde{u} = \gamma + \sum_{j=1}^{m} c_j u_j, \quad c_j \in \mathbb{R}, \quad j = 1, \dots, m.$$

Es genügt nun, (2-6) auf der Basis  $\{u_1, \ldots, u_m\}$  von V zu fordern, da (2-6) linear in V ist. Damit erhalten wir

$$0 = \int_{\Omega} \nabla \left( \gamma + \sum_{j=1}^{m} c_{j} u_{j} \right) \cdot \nabla u_{i} - g \cdot u_{i} dx$$
$$= \sum_{j=1}^{m} c_{j} \int_{\Omega} \nabla u_{j} \cdot \nabla u_{i} dx + \int_{\Omega} \left( \nabla \gamma \cdot \nabla u_{i} - g u_{i} \right) dx, \quad i = 1, \dots, m.$$

Wir bekommen also für  $c=(c_1,\ldots,c_m)\in\mathbb{R}^m$  ein lineares Gleichungssystem

$$Ac = r, \quad A \in \mathbb{R}^{m,m}, \quad r \in \mathbb{R}^m.$$
 (2-7)

Hier ist

$$A_{ij} = \int_{\Omega} \nabla u_j \cdot \nabla u_i \, dx, \quad 1 \le i, j \le m,$$

d.h. A ist eine symmetrische Matrix, und  $r = (r_1, \ldots, r_m)$  ist definiert durch

$$r_i = -\int_{\Omega} (\nabla \gamma \cdot \nabla u_i - gu_i) \ dx, \quad 1 \le i \le m.$$

Beim klassischen Ritz'schen Verfahren wählt man Ansatzfunktionen  $u_i \in C^1(\overline{\Omega})$ ,  $u_i|_{\partial\Omega} = 0$ , i = 1, ..., m, die im Allgemeinen auf dem gesamten Gebiet nicht verschwinden. A wird dann eine vollbesetzte Matrix.

Die Finite Elemente Verfahren weichen von dem Ritz-Verfahren in zwei wesentlichen Punkten ab:

- 1.) geringere Glattheit der Ansatzfunktionen,
- 2.) Ansatzfunktionen mit lokalem Träger.

**Zu 1.)** Man unterteilt das Gebiet  $\Omega$  in sogenannte finite Elemente durch eine Triangulierung.

Sei es der Einfachheit halber angenommen, dass  $\Omega$  ein polygonales Gebiet ist, d.h. der Rand  $\partial\Omega$  bestehe aus endlich vielen Geradenstücken.

Im allgemeinen Fall wird ein beliebiges beschränktes Gebiet durch  $\Omega_{T_h}$  approximiert.

 $E_h$  ist eine Menge von Dreiecken e. Je zwei Dreiecke  $e, e' \in E_h$  haben entweder eine Seite, einen Eckpunkt oder nichts gemeinsam. h ist die dabei maximal auftretende Kantenlänge.

Man sucht eine Lösung u in  $u_0 + V_{T_h}$ ,  $u_0$  Approximation von  $\gamma$ , wobei

$$V_{T_h} = \{u \in C(\Omega_{T_h}) \mid u \text{ ist ein Polynom mit } \deg(u) \leq r \text{ in } x_1 \text{ und } x_2 \}$$

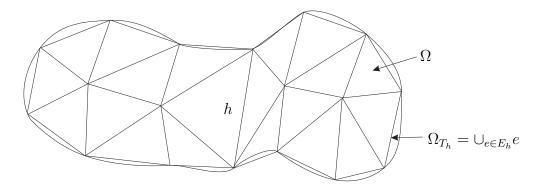


Abbildung 6: Triangulierung eines Gebietes  $\Omega$ 

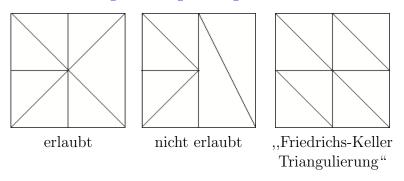


Abbildung 7: Beispiele für Triangulierungen

auf jedem Dreieck  $e \in E_h$  und u = 0 auf  $\partial \Omega_{T_h}$ .

Für r=1 handelt es sich hierbei um "lineare finite Elemente". Für r=2 sprechen wir von "quadratischen finiten Elementen", u.s.w.

Wir beschränken uns hier auf den Fall r = 1, d.h. u ist linear in  $x_1$  und  $x_2$  auf jedem Dreieck  $e \in E_h$ .

**Zu 2.**) Benutze eine Basis von Ansatzfunktionen in  $V_{T_h}$ , welche jeweils im größten Teil des Gebietes  $\Omega$  verschwinden. Dies impliziert, dass die entstehende Matrix dann dünn besetzt ist.

Vorgehensweise (lineare finite Elemente):

Seien  $P_i$ ,  $i=1,\ldots,M$  die durchnummerierten Knoten der Triangulierung  $\Omega_{T_h}$ , wobei die Knoten, die nicht auf  $\partial\Omega$  liegen, gerade die  $P_i$ ,  $i=1,\ldots,M$  mit m< M seien.

Definiere die Funktionen  $u_i$ , i = 1, ..., M auf  $\Omega_{T_h}$  durch

$$u_i(P_j) = \delta_{ij}, \quad 1 \le i, j \le M,$$
 (2-8)  
 $u_i$  ist linear in  $x_1$  und  $x_2$  auf jedem  $e \in E_h, u_i \in C(\Omega_{T_h}), \quad 1 \le i \le M.$ 

Die  $u_i$ , i = 1, ..., M heißen auch Formfunktionen. Es lässt sich zeigen, dass das Interpolationsproblem (2-8) eindeutig lösbar ist.

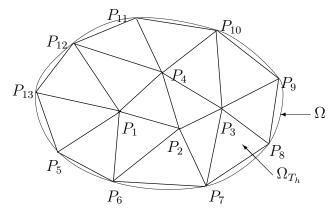


Abbildung 8: Nummerierung der Knoten, M=13, m=4

#### Struktur der Formfunktionen

Die Werte von  $u_i$  auf einer Dreiecksseite sind wegen der Linearität durch die Endwerte in den Knoten bestimmt. u ist daher auch stetig in  $\Omega_{T_h}$  und verschwindet auf allen Dreiecken, bei denen der Knoten  $P_i$  nicht vorkommt.

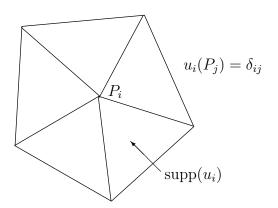


Abbildung 9: Träger  $supp(u_i)$  einer Formfunktion  $u_i$ 

Wir setzen ferner

$$u_0 := \sum_{j=m+1}^{M} \gamma(P_j) u_j,$$

wobei  $P_{m+1},\ldots,P_M\in\partial\Omega$  zu beachten ist, und erhalten

$$u_0(P_i) = \sum_{j=m+1}^{M} \gamma(P_j) \underbrace{u_j(P_i)}_{=\delta_{ij}} = \gamma(P_i), \quad i = m+1, \dots, M.$$

Der Ansatz erhält dann die Form

$$\widetilde{u} = u_0 + \sum_{j=1}^{m} c_j u_j = \sum_{j=1}^{M} c_j u_j$$
 (2-9)

mit  $c_j = \gamma(P_j)$  für  $j = m + 1, \dots, M$ .

Der Koeffizient  $c_i$  gibt dann gerade den Wert von  $\widetilde{u}$  im Knoten  $P_i$ ,  $i=1,\ldots,M$  an, denn

$$\widetilde{u}(P_i) = \sum_{j=1}^{M} c_j \underbrace{u_j(P_i)}_{=\delta_{ij}} = c_i, \quad i = 1, \dots, M.$$

Gesucht sind  $c_1, \ldots, c_m$ . Die Koeffizienten  $c_{m+1}, \ldots, c_M$  sind schon bekannt.

### Detailierte Aufstellung des linearen Gleichungssystems (2-7)

**2.14 Definition.** Die Aufstellung des linearen Gleichungssystems (2-7) heißt Assemblierung.

Unser Gleichungssystem hat die Form Ac = r mit

$$A_{ij} = \int_{\Omega} \nabla u_j \cdot \nabla u_i \, dx, \quad 1 \le i, j \le m,$$

$$r_i = -\int_{\Omega} (\nabla u_0 \cdot \nabla u_i - gu_i) \, dx, \quad 1 \le i \le m.$$

Wir ersetzen  $\Omega$  durch  $\Omega_{T_h}$ , g durch

$$g_{T_h} = \sum_{j=1}^{M} g(P_j) u_j$$

und finden mit  $u_0 = \sum_{j=m+1}^{M} \gamma(P_j) u_j$  dann

$$A_{ij} = \int_{\Omega_{T_h}} \nabla u_i \cdot \nabla u_j \, dx, \quad 1 \le i, j \le m,$$

$$r_i = \int_{\Omega_{T_h}} g_{T_h} u_i - \nabla u_0 \cdot \nabla u_i \, dx$$

$$= \sum_{j=1}^M g(P_j) \int_{\Omega_{T_h}} u_i u_j \, dx$$

$$- \sum_{j=m+1}^M \gamma(P_j) \int_{\Omega_{T_h}} \nabla u_j \cdot \nabla u_i \, dx, \quad i = 1, \dots, m.$$

Da das Integral über  $\Omega_{T_h}$  als die Summe der Integrale über jeweiligen Dreiecken  $e \in E_h$  dargestellt werden kann, d.h.

$$\int_{\Omega_{T_h}} f(x_1, x_2) \, dx = \sum_{e \in E_h} \int_e f(x_1, x_2) \, dx,$$

bleiben damit die folgenden Integrale zu bestimmen

$$\int_{e} \nabla u_i \cdot \nabla u_j \, dx, \quad \int_{e} u_i u_j \, dx, \quad 1 \le i, j \le M. \tag{2-10}$$

Die Integrale in (2-10) lassen sich wie folgt berechnen.

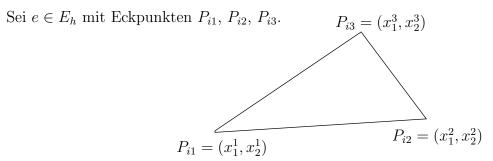


Abbildung 10: Dreieck e mit den Ecken  $P_{i1}$ ,  $P_{i2}$ ,  $P_{i3}$ 

Es sind dann die Integrale  $S^e_{jk} = \int_e \nabla u_{ij} \cdot \nabla u_{ik} \, dx, \, k,j = 1,2,3$  von Null verschieden.

Die symmetrische Matrix  $S^e=(S^e_{jk})_{1\leq j,k\leq 3}\in\mathbb{R}^{3,3}$  ist beim Aufbau der Gesamtmatrix A auf die Untermatrix

$$\begin{pmatrix} A_{i1,i1} & A_{i1,i2} & A_{i1,i3} \\ A_{i2,i1} & A_{i2,i2} & A_{i2,i3} \\ A_{i3,i1} & A_{i3,i2} & A_{i3,i3} \end{pmatrix}$$

aufzuaddieren. Man findet

$$S^e = \frac{1}{F_e} C_e C_e^T$$

mit

$$C_e = \begin{pmatrix} x_2^2 - x_2^3 & x_1^3 - x_1^2 \\ x_2^3 - x_2^1 & x_1^1 - x_1^3 \\ x_2^1 - x_2^2 & x_1^2 - x_1^1 \end{pmatrix} \in \mathbb{R}^{3,2}$$

und  $F_e = (x_1^2 - x_1^1)(x_2^3 - x_2^1) - (x_1^3 - x_1^1)(x_2^2 - x_2^1)$ . Dabei ist  $F_e$  die doppelte Fläche des Dreiecks e mit den Eckpunkten  $P_{i1}$ ,  $P_{i2}$ ,  $P_{i3}$ .

Das zweite Integral in (2-10) ergibt sich als die Matrix  $M^e \in \mathbb{R}^{3,3}$  mit

$$(M^e)_{jk} = \int_e u_{ij} \cdot u_{ik} \, dx, \quad 1 \le j, k \le 3,$$

für welche

$$M^e = \frac{F_e}{24} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} \in \mathbb{R}^{3,3}$$

gilt.

### c) Stabilität und Konvergenz der Finite Elemente Methode

Es sei im Folgenden  $V = H_0^1(\Omega)$  und

$$a: V \times V \longrightarrow \mathbb{R}, \quad a(u, v) = \int_{\Omega} \nabla v \cdot \nabla u \, dx,$$
  
 $b: V \longrightarrow \mathbb{R}, \quad b(v) = \int_{\Omega} g \cdot v \, dx$ 

die schwache Formulierung des homogenen Modellproblems (2-1). Die Variationsaufgabe besteht nun darin, dass ein  $\overline{u} \in V$  gesucht wird mit

$$a(\overline{u}, v) = b(v) \quad \forall v \in V.$$
 (2-11)

Das Galerkin-Verfahren lautet entsprechend: Finde ein  $u_h \in V_h$  mit

$$a(u_h, v) = v(v) \quad \forall v \in V_h \tag{2-12}$$

für einen endlich dimensionalen Teilraum  $V_h \subset V$  von V.

Man nennt (2-12) die Galerkinschen Gleichungen zu (2-11) bzw. zum Modellproblem (2-1).

 $e := \overline{u} - u_h$  bezeichne den dadurch entstehenden Fehler, wobei  $\overline{u}$  und  $u_h$  die Gleichungen (2-11) bzw. (2-12) lösen.

Wir finden die Fehlergleichung

$$a(e, v) = 0 \quad \forall v \in V_h$$

denn aus (2-11) folgt

$$a(\overline{u}, v) = b(v) \quad \forall v \in V_h \subset V$$

und aus (2-12) folgt

$$a(u_h, v) = b(v) \quad \forall v \in V_h.$$

Insgesamt erhalten wir

$$a(e, v) = a(\overline{u} - u_h, v) = b(v) - b(v) = 0 \quad \forall v \in V_h.$$
 (2-13)

Ist a eine symmetrische positiv definite Bilinearform, so besagt (2-13), dass e und v für  $v \in V_h$  orthogonal sind, d.h.  $e \in V_h^{\perp}$ .

(2-13) bedeutet die Orthogonalität des Fehlers auf dem Ansatzraum. Durch das Galerkin-Vefahren (2-12) wird also dasjenige Element  $u_h \in V_h$  charakterisiert, welches bezüglich der Norm

$$\|\cdot\|_a = \sqrt{a(\cdot,\cdot)}$$

den kleinsten Abstand zu  $\overline{u} \in V$  hat.

**2.15 Lemma.** Sei  $V_h \subset V$  ein Unterraum von V, a ein Skalarprodukt auf V und  $||u||_a = \sqrt{a(u,u)}$  die davon erzeugte Norm. Dann gilt für  $u_h \in V_h$ :

$$a(\overline{u} - u_h, v) = 0 \quad \forall v \in V_h,$$

was mit

$$\|\overline{u} - u_h\|_a = \inf\{\|\overline{u} - v\|_a \mid v \in V_h\}.$$

äquivalent ist.

Beweis: Für ein festes  $\overline{u} \in V$  sei  $b(v) = a(\overline{u}, v), v \in V_h$ , d.h. b ist ein lineares Funktional auf  $V_h$ . Somit ist

$$\underbrace{a(\overline{u}, v)}_{=b(v)} = a(u_h, v) \quad v \in V_h$$

eine Variationsgleichung auf  $V_h$ . Diese hat nach Lemma 2.12 die gleiche Lösungen wie

$$F(u_h) = \inf\{F(v) \mid v \in V_h\}$$

mit 
$$F(v) = \frac{1}{2}a(v,v) - b(v) = \frac{1}{2}a(v,v) - a(\overline{u},v).$$

F hat die gleichen Minima wie das Funktional

$$G(v) = (2F(v) + a(\overline{u}, \overline{u}))^{1/2}$$

$$= (a(v, v) - 2a(\overline{u}, v) + a(\overline{u}, \overline{u}))^{1/2}$$

$$= a(\overline{u} - v, \overline{u} - v)^{1/2} = ||\overline{u} - v||_a.$$

Wir setzen jetzt für a das Folgende voraus: Sei  $\|\cdot\|$  eine Norm auf V bezüglich der gelte

- i)  $a: V \times V \longrightarrow \mathbb{R}$  ist stetig bezüglich  $\|\cdot\|$ , d.h. es existiert ein M > 0 mit  $|a(u,v)| < M \cdot \|u\| \cdot \|v\| \quad \forall u,v \in V.$
- ii) a ist V-elliptisch, d.h. es existiert ein  $\alpha > 0$  mit

$$a(u, u) > \alpha ||u||^2 \quad \forall u \in V.$$

Ferner sei das lineare Funktional  $b: V \longrightarrow \mathbb{R}$  stetig.

Dann sichert der Satz 2.13 von Lax und Milgram, dass die Variationsgleichung

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

genau eine Lösung  $u \in V$  besitzt.

In unserem Fall ist  $a: V \times V \longrightarrow \mathbb{R}$  ein Skalarprodukt auf V. Deshalb ergibt sich mit der induzierten Norm  $\|\cdot\| = \|\cdot\|_a$  unter Verwendung der Cauchy-Schwarzschen Ungleichung die Gültigkeit von i) mit M=1 sowie von ii) mit  $\alpha=1$ .

Die V-Elliptizität impliziert die Stabilität der Galerkin-Approximation.

**2.16 Lemma.** Die Lösung  $u_h$  der Galerkin-Gleichung a(u,v) = b(v),  $\forall v \in V_h$  ist stabil im Sinne der Abschätzung

$$||u_h|| \leq \frac{1}{\alpha} ||b||_{V'}$$

 $f\ddot{u}r$  alle h > 0, wobei

$$||b||_{V'} = \sup \left\{ \frac{|b(v)|}{||v||} \mid v \in V, v \neq 0 \right\}.$$

Beweis: Aus  $a(u_h, v) = b(v)$  für alle  $v \in V_h$  folgt

$$\alpha \|u_h\|^2 \stackrel{ii)}{\leq} a(u_h, u_h) = b(u_h)$$
  
 $\leq \frac{|b(u_h)|}{\|u_h\|} \cdot \|u_h\| \leq \|b\|_{V'} \cdot \|u_h\|.$ 

Also folgt mit  $\alpha > 0$ 

$$||u_h|| \leq \frac{||b||_{V'}}{\alpha}.$$

Somit gilt bis auf eine Konstante weiterhin die Approximationsaussage aus Lemma 2.15.

**2.17 Satz** (Cea's Lemma). Unter den Voraussetzungen i) und ii) an die Bilinearform a und der Stetigkeit des linearen Funktionals b gilt für den Fehler der Galerkin-Lösung u<sub>h</sub> die Abschätzung

$$\|\overline{u} - u_h\| \le \frac{M}{\alpha} \cdot \inf\{\|\overline{u} - v\| \mid v \in V_h\}.$$

Beweis: Sei  $v \in V_h$  beliebig. Aus der Fehlergleichung

$$a(\overline{u} - u_h, v) = 0 \quad v \in V_h \tag{Vgl. (2-13)}$$

folgt

$$a(\overline{u} - u_h, u_h - v) = 0,$$

 $da u_h - v \in V_h$ .

Mit der V-Elliptizität erhalten wir dann

$$\alpha \|\overline{u} - u_h\|^2 \leq a(\overline{u} - u_h, \overline{u} - u_h)$$

$$= a(\overline{u} - u_h, \overline{u} - u_h) + \underbrace{a(\overline{u} - u_h, u_h - v)}_{=0}$$

$$= a(\overline{u} - u_h, \overline{u} - v)$$

$$\leq M \cdot \|\overline{u} - u_h\| \cdot \|\overline{u} - v\|$$

für alle  $v \in V_h$ , d.h.

$$\|\overline{u} - u_h\| \le \frac{M}{\alpha} \|\overline{u} - v\| \quad \forall v \in V_h.$$

Also folgt

$$\|\overline{u} - u_h\| \le \frac{M}{\alpha} \cdot \inf\{\|\overline{u} - v\| \mid v \in V_h\}.$$

Es reicht also für eine asymptotische Fehlerdarstellung in h den Restapproximationsfehler von  $V_h$ , d.h. inf $\{\|\overline{u} - v\| \mid v \in V_h\}$  abzuschätzen.

Betrachte nun die konkrete Variationsgleichung

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

mit  $a(u,v)=\int_{\Omega}\nabla u\cdot\nabla v\,dx$ ,  $b(v)=\int_{\Omega}g\cdot v\,dx$ , welche die schwache Formulierung der Dirichlet'schen Randwertaufgabe darstellt. Dabei sei  $V=H^1_0(\Omega)$ .

Zu betrachten bleibt die Wahl der Normen. Offensichtlich kann als die Norm $\|\cdot\|=\|\cdot\|_a$ 

$$||u||_a = \left(\int_{\Omega} |\nabla u|^2 \, dx\right)^{1/2}$$

genommen werden. Alternativ kann aber auch die durch das Skalarprodukt

$$\langle u, v \rangle_1 = \sum_{|\alpha| \le 1} \langle \partial^{\alpha} u, \partial^{\alpha} v \rangle_0$$

induzierte Norm

$$||u||_1 = \left(\int_{\Omega} |u(x)|^2 dx + \int_{\Omega} |\nabla u(x)|^2 dx\right)^{1/2}, \quad u \in H_0^1(\Omega)$$

gewählt werden, falls die Bilinearform a bezüglich  $\|\cdot\|_1$  stetig und V-elliptisch ist. Für die Stetigkeit benutzen wir die Abschätzung

$$||u||_a = \left(\int |\nabla u(x)|^2 dx\right)^{1/2}$$

$$\leq \left( \int_{\Omega} |u(x)|^2 dx + \int_{\Omega} |\nabla u(x)|^2 dx \right)^{1/2} = ||u||_1,$$

und somit folgt mit Cauchy-Schwarz

$$|a(u,v)| \le ||u||_a \cdot ||v||_a \le ||u||_1 \cdot ||v||_1 \quad \forall u, v \in H_0^1(\Omega),$$

d.h. die Bilinearform a ist auch stetig bezüglich  $\|\cdot\|_1$  mit M=1.

Die V-Elliptizität

$$a(u, u) \ge \alpha ||u||_1^2, \quad \forall u \in V \quad (\alpha > 0)$$

gilt nicht im Allgemeinen für  $V=H^1(\Omega)$ . Sie gilt aber für  $V=H^1_0(\Omega)$ . Dazu benötigen wir den folgenden Satz.

**2.18 Satz** (Erste Poincarésche Abschätzung). Sei  $\Omega$  ein beschränktes Gebiet. Dann gibt es ein C>0 mit

$$\left(\int_{\Omega} |u(x)|^2 dx\right)^{1/2} \le C \left(\int_{\Omega} |\nabla u(x)|^2 dx\right)^{1/2} \quad \forall u \in H_0^1(\Omega).$$

Mit dem Satz von Poincaré folgt also

$$||u||_0 = \left(\int_{\Omega} |u(x)|^2 dx\right)^{1/2} \le C||u||_a$$

und damit

$$\|u\|_1^2 = \|u\|_0^2 + \|u\|_a^2 \le (1 + C^2) \|u\|_a^2,$$

d.h.

$$\frac{1}{\sqrt{1+C^2}} \|u\|_1 \le \|u\|_a, \quad u \in H^1_0(\Omega).$$

Dies liefert

$$a(u, u) = ||u||_a^2 \ge \frac{1}{1 + C^2} ||u||_1^2,$$

und somit die V-Elliptizität bezüglich  $\|\cdot\|_1$  mit  $\alpha = \frac{1}{1+C^2}$ .

Damit sind die Normen  $\|\cdot\|_a$  und  $\|\cdot\|_1$  auf  $V=H^1_0(\Omega)$  äquivalent und erzeugen daher denselben Konvergenzbegriff.

Im Lemma 2.17 von Cea erhalten wir für  $\|\cdot\| = \|\cdot\|_1$  die Abschätzung

$$\|\overline{u} - u_h\|_1 \le \frac{1}{\alpha} \cdot \inf\{\|\overline{u} - v\|_1 \,|\, v \in V_h\}$$
  
=  $(1 + C^2) \cdot \inf\{\|\overline{u} - v\|_1 \,|\, v \in V_h\}.$ 

Bis jetzt war unsere Abschätzung des Approximationsfehlers unabhängig von der konkreten Wahl von  $V_h$ . Jetzt möchten wir den Approximationsfehler der Finiten Elemente Methode analysieren. In diesem Fall ist  $V_h = V_{T_h}$  mit

$$V_{T_h} = \{ u \in C(\Omega_{T_h}) \mid u \text{ ist affin linear in } x_1 \text{ und } x_2$$
 auf jedem Dreieck  $e \in E_h \text{ und } u = 0 \text{ auf } \partial \Omega_{T_h} \}.$ 

Es gilt

$$\|\overline{u} - u_h\|_1 \le (1 + C^2) \cdot \inf\{\|\overline{u} - v\|_1 | v \in V_{T_h}\}$$
  
 
$$\le (1 + C^2) \|\overline{u} - w\|_1$$

für jedes  $w \in V_{T_h}$ .

Wir konstruieren nun ein spezielles  $w \in V_{T_h}$  durch Interpolation.

Sei  $\Omega$  polygonal berandet mit strikter Kegeleigenschaft, und sei  $\Omega_{T_h}$  eine Familie von Triangulierungen von  $\Omega$  mit Schrittweite h für  $0 < h \le h_0$ . Betrachte die Abbildung

$$I_h: \{u \in C(\overline{\Omega}) \mid u = 0 \text{ auf } \partial\Omega\} \longrightarrow V_{T_h}$$

$$v \longmapsto I_h(v) = \sum_{i=1}^m v(P_i)u_i$$

mit den Formfunktionen  $u_i$ , i = 1, ..., m zu den Knoten  $P_i$  der Triangulierung  $\Omega_{T_h}$ .

 $I_h(v)$  ist nur für  $v \in H^2(\Omega)$  definiert, da  $H^2(\Omega) \subset C(\overline{\Omega})$  nach dem Sobolevschen Einbettungssatz 2.7.

Abzuschätzen bleibt der Interpolationsfehler  $I_h(v) - v$  in der  $H^1$ -Norm. Für diesen Fehler kann man zeigen

$$||I_h(v) - v||_1 \le \overline{C} \cdot h \cdot ||v||_2 \quad \forall v \in H^2(\Omega)$$

 $(\|\cdot\|_2 \text{ Norm auf } H^2(\Omega))$ , falls  $(T_h)_{0 < h < h_0}$  eine Familie von Triangulierungen von  $\Omega$  ist, so dass für den jeweils maximalen Winkel  $\gamma_{h,\max}$  in einem Dreieck  $e \in T_h$  gilt

$$\gamma_{h,\max} \leq \gamma_{\max} < \pi$$

für  $0 < h \le h_0$ . Letztere Bedingung heißt Maximalwinkelbedingung.

Liegt also die Lösung  $\overline{u}$  der Variationsaufgabe in  $H^2(\Omega) \cap H^1_0(\Omega)$ , so dürfen wir  $v = \overline{u}$  wählen und erhalten für lineare Finite Elemente die Abschätzung

$$\|\overline{u} - u_h\|_1 \le (1 + C^2) \|\overline{u} - I_h(\overline{u})\|_1$$
  
 $\le (1 + C^2) \cdot \overline{C} \cdot h \|\overline{u}\|_2 \le \hat{C}h, \quad 0 < h \le h_0.$ 

Also ist die lineare Finite Elemente Methode konvergent der Ordnung 1 bezüglich  $\|\cdot\|_1$  an Lösungen der Glattheit  $H^2(\Omega) \cap H^1_0(\Omega)$ .

Fehlerabschätzung in der L<sup>2</sup>-Norm

2.19 Lemma (Aubin-Nitsche). Besitzt das duale Problem

$$a(v,z) = b(v) = \int_{\Omega} g \cdot v \, dx \quad \forall v \in V$$

 $f\ddot{u}r\ g\in L^2(\Omega)$  eine Lösung  $z\in H^2(\Omega)$  mit  $||z||_2:=||z||_{H^2}\leq C||g||_0$ , so gilt

$$\|\overline{u} - u_h\|_0 \le \widetilde{C} \cdot h \cdot \|\overline{u} - u_h\|_1.$$

Somit folgt unter den Voraussetzungen von Lemma 2.19

$$\|\overline{u} - u_h\|_0 \le \widetilde{C} \cdot h \cdot \|\overline{u} - u_h\|_1 \le \widetilde{C}\widehat{C}h^2, \quad 0 < h \le h_0,$$

d.h. wir haben Konvergenz der Ordnung 2 bezüglich  $\|\cdot\|_0$ .

- **2.20 Bemerkung.** a) Die Konvergenzaussage für lineare Finite Elemente bleibt auch richtig, falls der Rand  $\partial\Omega$  stückweise stetig differenzierbar ist.
- b) Höhere Konvergenzordnung in der  $H^1$ -Norm bzw. der  $L^2$ -Norm lassen sich mit polynomialen Finiten Elementen erreichen, falls die Lösung der Randwertaufgabe hinreichend oft schwach differenzierbar ist und  $\Omega$  polygonal berandet ist.
- c) Im Fall von höheren polynomialen Finiten Elementen  $(r \geq 2)$ , benötigt man spezielle Randelemente, um höhere Konvergenzordnung bei Gebieten  $\Omega$  mit stückweise glattem Rand  $\partial\Omega$  sicherzustellen.

# 3. Finite Differenzenverfahren für parabolische Differentialgleichungen

# a) Das Prinzip der Linienmethode

Wir betrachten die nichtlineare Anfangswertaufgabe

$$u_{t} = u_{xx} + f(u, u_{x}, x, t) \text{ in } \Omega = (0, 1) \times (0, T),$$

$$u(x, 0) = u_{0}(x) \text{ für } 0 \le x \le 1,$$

$$u(0, t) = \gamma_{0}(t), u(1, t) = \gamma_{1}(t) \text{ für } 0 \le t \le T.$$
(3-1)

Das Prinzip der Linienmethode beruht darauf, zuerst eine Diskretisierung in der Ortsvariablen vorzunehmen. Die Zeitvariable t bleibt noch kontinuerlich.

Um diesen Schritt durchzuführen, betrachten wir zunächst die Randwertaufgabe

$$-w''(x) = f(x, w(x), w'(x)) \text{ in } (0, 1),$$
  

$$w(0) = \gamma_0, w(1) = \gamma_1.$$
(3-2)

Zu  $\Delta x = \frac{1}{M} > 0 \ (M \in \mathbb{N})$  bekommen wir das Gitter  $\Omega_{\Delta x} = \{j\Delta x \mid j = 0, \dots, M\}$ . Wir ersetzen -w''(x) durch

$$-w''(x) \sim \frac{1}{\Delta x^2} (-w(x - \Delta x) + 2w(x) - w(x + \Delta x))$$
 (3-3)

und w'(x) durch

$$w'(x) \sim \frac{1}{2\Delta x} (w(x + \Delta x) - w(x - \Delta x)). \tag{3-4}$$

Präziser gilt

$$\left| \frac{1}{\Delta x^2} (-w(x - \Delta x) + 2w(x) - w(x + \Delta x)) + w''(x) \right| = O(\Delta x^2),$$

falls  $w \in C^4$  (vgl. Kapitel 1) bzw.

$$\left| \frac{1}{2\Delta x} (w(x + \Delta x) - w(x - \Delta x)) - w'(x) \right| = O(\Delta x^2)$$

für  $w \in C^3$ .

Führen wir jetzt den Ersetzungsprozess durch, so erhalten wir

$$w(0) = \gamma_0,$$
  
$$\frac{1}{\Delta x^2} (-w(x - \Delta x) + 2w(x) - w(x + \Delta x)) =$$

$$= f\left(x, w(x), \frac{1}{2\Delta x}(w(x + \Delta x) - w(x - \Delta x))\right), \qquad (3-5)$$

$$x = j\Delta x, j = 1, \dots, M - 1,$$

$$w(1) = \gamma_1.$$

Gesucht ist also ein Vektor  $(w(0), w(\Delta x), \dots, w(1 - \Delta x), w(1))$ , welcher (3-5) löst. Wir eliminieren w(0), w(1) aus dem System (3-5) und schreiben es in kompakter Form

$$\frac{1}{\Delta x^{2}} \begin{pmatrix}
2 & -1 & 0 & \dots & 0 & 0 & 0 \\
-1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\
0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\
\vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & 0 & \dots & 2 & -1 & 0 & 0 \\
0 & 0 & 0 & \dots & 0 & -1 & 2 & -1 & 0 \\
0 & 0 & 0 & \dots & 0 & -1 & 2 & -1 & 0 \\
0 & 0 & 0 & \dots & 0 & -1 & 2 & 0
\end{pmatrix}
\begin{pmatrix}
f\left(\Delta x, w(\Delta x), \frac{1}{2\Delta x}(w(2\Delta x) - \gamma_{0})\right) + \frac{\gamma_{0}}{\Delta x^{2}} \\
f\left(2\Delta x, w(2\Delta x), \frac{1}{2\Delta x}(w(3\Delta x) - w(\Delta x))\right) \\
\vdots \\
f\left(j\Delta x, w(j\Delta x), \frac{1}{2\Delta x}(w((j+1)\Delta x) - w((j-1)\Delta x))\right), j = 2, \dots, M-2 \\
\vdots \\
f\left(1 - 2\Delta x, w(1 - 2\Delta x), \frac{1}{2\Delta x}(w(1 - \Delta x) - w(1 - 3\Delta x))\right) \\
f\left(1 - \Delta x, w(1 - \Delta x), \frac{1}{2\Delta x}(\gamma_{1} - w(1 - 2\Delta x))\right) + \frac{\gamma_{1}}{\Delta x^{2}}
\end{pmatrix}.$$

Das Gleichungssystem (3-6) hat die Form

$$A^{\Delta x}w = G^{\Delta x}(w), \quad w \in \mathbb{R}^{\Omega_{\Delta x}}$$
 (3-7)

mit

$$A^{\Delta x} = \frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 & 0 \end{pmatrix},$$

$$G^{\Delta x}(w) = \begin{pmatrix} f\left(\Delta x, w(\Delta x), \frac{1}{2\Delta x}(w(2\Delta x) - \gamma_0)\right) + \frac{\gamma_0}{\Delta x^2} \\ f\left(j\Delta x, w(j\Delta x), \frac{1}{2\Delta x}(w((j+1)\Delta x) - w((j-1)\Delta x))\right), \\ j = 2, \dots, M - 2 \\ f\left(1 - \Delta x, w(1 - \Delta x), \frac{1}{2\Delta x}(\gamma_1 - w(1 - 2\Delta x))\right) + \frac{\gamma_1}{\Delta x^2} \end{pmatrix}.$$

Insbesondere zeigt dies, dass

$$R^{\Delta x} = -A^{\Delta x}w + G^{\Delta x}(w)$$

eine Diskretisierung für den Operator A ist mit

$$\begin{split} A:D(A)\subset C^2((0,1))\cap C([0,1]) &\;\;\longrightarrow\;\; C((0,1)),\\ w &\;\;\longmapsto\;\; w''(\cdot)+f(\cdot,w(\cdot),w'(\cdot)), \end{split}$$

wobei

$$D(A) := \{ w \in C^2((0,1)) \cap C([0,1]) \mid w(0) = \gamma_0, \ w(1) = \gamma_1 \}.$$

Wir kehren nun zu unserer Anfangsrandwertaufgabe (3-1) zurück. Gemäß des Prinzips der Linienmethode diskretisieren wir für beliebiges  $t \in [0, T]$  zuerst die Raumvariable.

Wir suchen eine Gleichung für die Funktion

$$v(t) = (u(\Delta x, t), u(2\Delta x, t), \dots, u(1 - 2\Delta x, t), u(1 - \Delta x, t))$$
  
=  $(v_1(t), v_2(t), \dots, v_{M-2}(t), v_{M-1}(t)) \in \mathbb{R}^{M-1}, \quad 0 \le t \le T.$ 

Eine einfache Diskretisierung können wir jetzt in der Ortsvariablen durch das klassische Differenzenverfahren vornehmen.

Wir ersetzen den Ausdruck  $u_{xx}(x,t)+f(u(x,t),u_x(x,t),x,t), u(0,t)=\gamma_0(t), u(1,t)=\gamma_1(t)$  für  $t \in [0,T]$  durch

$$-\underbrace{\frac{1}{\Delta x^2}\begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 & 0 \end{pmatrix}\begin{pmatrix} v_1(t) \\ v_2(t) \\ v_3(t) \\ \vdots \\ v_{M-3}(t) \\ v_{M-2}(t) \\ v_{M-1}(t) \end{pmatrix}}_{=x^{\Delta x}} + \underbrace{\begin{pmatrix} f\left(v_1(t), \frac{1}{2\Delta x}(v_2(t) - \gamma_0), \Delta x, t\right) \\ f\left(v_j(t), \frac{1}{2\Delta x}(v_{j+1}(t) - v_{j-1}(t)), j\Delta x, t\right), \\ j = 2, \dots, M - 2 \\ f\left(v_{M-1}(t), \frac{1}{2\Delta x}(\gamma_1 - v_{M-2}(t)), (M-1)\Delta x, t\right) \end{pmatrix}}_{=:H^{\Delta x}(v(t))}$$

$$=: -A^{\Delta x}v(t) + r^{\Delta x}(t) + H^{\Delta x}(v(t)).$$

Präziser wird  $A(t)u(t,\cdot), u(t,\cdot) \in D(A(t))$  durch

$$-A^{\Delta x}v(t) + r^{\Delta x}(t) + H^{\Delta x}(v(t))$$

ersetzt, wobei der zeitabhängige Differentialoperator A(t) wie folgt definiert ist:

$$A(t): D(A(t)) \subset C^2((0,1)) \cap C([0,1]) \longrightarrow C((0,1)),$$
  
$$u \longmapsto u_{xx} + f(u, u_x, \cdot, t)$$

mit dem Definitionsbereich

$$D(A(t)) := \{ u \in C^2((0,1)) \cap C([0,1]) \mid u(0) = \gamma_0(t), \ u(1) = \gamma_1(t) \}.$$

Setzt man  $v'(t) = (u_t(\Delta x, t), \dots, u_t((M-1)\Delta x, t))$ , so ergibt sich unter Beachtung von (3-1) die gewöhnliche Differentialgleichung

$$v'(t) = -A^{\Delta x}v(t) + H^{\Delta x}(v(t)) + r^{\Delta x}(t)$$

$$=: F_{\Delta x}(v(t), t), \quad 0 \le t \le T,$$

$$v(0) = v^{0} = (u_{0}(\Delta x), \dots, u_{0}((M-1)\Delta x))$$
(3-8)

als Ersatz für die Anfangsrandwertaufgabe (3-1).

Der Übergang von (3-1) zu (3-8) wird als Semidiskretisierung bezeichnet. Zu einer vollständigen Diskretisierung von (3-1) kommen wir, wenn wir auf (3-8) eines der handelsüblichen Verfahren zur Lösung von Anfangswertproblemen anwenden. Bei der Durchführung dieser Verfahren sollte man sich jedoch die spezielle Gestalt von  $F_{\Delta x}$  zu Nutze machen, denn das System (3-8) wird für  $\Delta x \longrightarrow 0$  immer größer. Das einfachste Verfahren für (3-8) ist das Euler-Cauchy Verfahren.

Seien  $\Delta t = \frac{T}{N} > 0$ ,  $v^0$  vorgegeben.  $v^j$  approximiere  $v(j\Delta t)$ ,  $j = 0, \ldots, N$ . Dann finden wir mit  $t_i = j\Delta t$ 

$$v^{j+1} = v^{j} + \Delta t \cdot F_{\Delta x}(v^{j}, t_{j}), \quad j = 0, \dots, N - 1,$$

$$v^{0} = (u_{0}(\Delta x), \dots, u_{0}((M - 1)\Delta x))$$
(3-9)

oder explizit mit

$$v^{j} = (u_{1}^{j}, \dots, u_{M-1}^{j}), \quad u_{i}^{j} = u(x_{i}, t_{j}), \quad x_{i} = i\Delta x$$

folgt

$$\begin{pmatrix} u_{1}^{j+1} \\ \vdots \\ u_{i}^{j+1} \\ \vdots \\ u_{M-1}^{j+1} \end{pmatrix} = \begin{pmatrix} u_{1}^{j} \\ \vdots \\ u_{i}^{j} \\ \vdots \\ u_{M-1}^{j} \end{pmatrix} + \Delta t \begin{cases} \frac{1}{\Delta x^{2}} \begin{pmatrix} -2u_{1}^{j} + u_{2}^{j} \\ \vdots \\ u_{i-1}^{j} - 2u_{i}^{j} + u_{i+1}^{j} \\ \vdots \\ u_{M-2}^{j} - 2u_{M-1}^{j} \end{pmatrix}$$

$$+ \frac{1}{\Delta x^{2}} \begin{pmatrix} \gamma_{0}(t_{j}) \\ 0 \\ \vdots \\ 0 \\ \gamma_{1}(t_{j}) \end{pmatrix} + \begin{pmatrix} f\left(u_{1}^{j}, \frac{1}{2\Delta x}(u_{2}^{j} - \gamma_{0}(t_{j})), x_{1}, t_{j}\right) \\ \vdots \\ f\left(u_{M-1}^{j}, \frac{1}{2\Delta x}(u_{i+1}^{j} - u_{i-1}^{j}), x_{i}, t_{j}\right) \\ \vdots \\ f\left(u_{M-1}^{j}, \frac{1}{2\Delta x}(v_{1}^{j}, u_{2}^{j} - v_{M-2}^{j}), x_{M-1}, t_{j}\right) \end{cases}$$

mit den Anfangsbedingungen

$$\begin{pmatrix} u_1^0 \\ \vdots \\ u_i^0 \\ \vdots \\ u_{M-1}^0 \end{pmatrix} = \begin{pmatrix} u_0(x_1) \\ \vdots \\ u_0(x_i) \\ \vdots \\ u_0(x_{M-1}) \end{pmatrix}. \tag{3-11}$$

(3-10), (3-11) heißen auch explizites Differenzenverfahren zur Anfangsrandwertaufgabe (3-1). Man erhält  $u_i^{j+1}$  aus  $u_{i-1}^j$ ,  $u_i^j$ ,  $u_{i+1}^j$ , was man durch folgendes Schema andeutet.

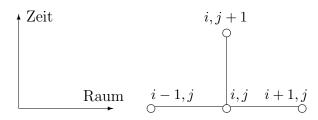


Abbildung 11: "Differenzenstern"

Wir betrachten als nächstes für das Liniensystem (3-8) das von einem Parameter  $\vartheta \in [0,1]$  abhängige Verfahren

$$v_{j+1} = v^{j} + \Delta t \left[ \vartheta F_{\Delta x}(v^{j+1}, t_{j+1}) + (1 - \vartheta) F_{\Delta x}(v^{j}, t_{j}) \right],$$

$$j = 0, \dots, N - 1,$$

$$v^{0} = (u_{0}(x_{1}), \dots, u_{0}(x_{M-1})).$$
(3-12)

Die Spezialfälle des  $\vartheta$ -Verfahrens für das Liniensystem sind:

 $\vartheta = 0$ : Euler-Cauchy Verfahren,

 $\vartheta = \frac{1}{2}$ : Trapezenmethode oder sogenanntes "Crank-Nicholson Verfahren",

 $\vartheta = \overline{1}$ : Implizites Euler-Cauchy Verfahren.

Für  $\vartheta > 0$  erfordert (3-12) die Auflösung eines nichtlinearen Gleichungssystems und das Verfahren heißt implizit. Die Abbildung 12 stellt die zugehörigen Differenzensterne dar.

**3.1 Beispiel** (Chemische Reaktions-Transport-Gleichung). Die dazu gehörige Anfangsrandwertaufgabe lautet

$$u_t = u_{xx} - k \cdot e_0(x)u, \quad 0 \le x \le 1, \ t \ge 0, \quad k \ge 0$$
  
 
$$u(x,0) = 1 - x, \quad 0 \le x \le 1$$
  
 
$$u(0,t) = 1, \ u(1,t) = 0, \quad e_0(x) \ge 0, \ 0 \le x \le 1.$$

$$i, j+1 \qquad \qquad i, j+1 \qquad \qquad i, j+1 \qquad \qquad i, j+1 \qquad \qquad i, j+1 \qquad \qquad i + 1, j + 1 \qquad \qquad$$

Abbildung 12: Differenzensterne für  $\vartheta \in \{0, \frac{1}{2}, 1\}$ 

Wir finden für das Liniensystem

$$v'(t) = F_{\Delta x}(v(t), t), v(0) = v^0$$

mit

$$F_{\Delta x}(v,t) = -\frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix} v + \frac{1}{\Delta x^2} \begin{pmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{pmatrix}$$

$$-k \begin{pmatrix} e_0(x_1)v_1 \\ e_0(x_2)v_2 \\ e_0(x_3)v_3 \\ \vdots \\ e_0(x_{M-3})v_{M-3} \\ e_0(x_{M-2})v_{M-2} \\ e_0(x_{M-1})v_{M-1} \end{pmatrix} = -B_{\Delta x}v + r^{\Delta x},$$

wobei  $x_i = i\Delta x, i = 1, \dots, M-1$ . Dabei ist

$$B_{\Delta x} = \frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix}$$

$$r^{\Delta x} = \frac{k \cdot \operatorname{diag}(e_0(x_1), \dots, e_0(x_{M-1})),}{\Delta x^2}$$

 $B_{\Delta x}$  ist eine  $L_0$ -Matrix mit majorisierendem Element  $\mathbb{I} = (1, 1, \dots, 1, 1)^T$ , d.h. eine M-Matrix, da  $e_0(x) \geq 0$ ,  $0 \leq x \leq 1$  und h > 0.

Diskretisierung der Anfangsrandwertaufgabe mit dem  $\vartheta$ -Verfahren liefert dann die Iteration

$$v^{j+1} = v^{j} + \Delta t \cdot \left[ \vartheta(-B_{\Delta x}v^{j+1} + r) + (1 - \vartheta)(-B_{\Delta x}v^{j} + r) \right],$$
  

$$j = 0, \dots, N - 1$$
  

$$v^{0} = (u_{0}(x_{1}), \dots, u_{0}(x_{M-1})),$$

woraus sich

$$\left(\frac{1}{\Delta t}I + \vartheta B_{\Delta x}\right)v^{j+1} = \left(\frac{1}{\Delta t}I - (1 - \vartheta)B_{\Delta x}\right)v^{j} + r$$

für  $j = 0, \dots, N - 1$  ergibt.

aufgabe sehr wichtig.

Ferner ist  $\frac{1}{\Delta t} + \vartheta B_{\Delta x}$  eine M-Matrix für  $\vartheta \in [0, 1]$ . In jedem Schritt ist also für  $\vartheta > 0$  ein tridiagonales Gleichungssystem zu lösen.

#### Beobachtungen

- a)  $\vartheta = 0$ :
  - Es ergeben sich nur für sehr kleine Zeitschritte gute Näherungen. Bei zu großem  $\Delta t$  treten starke Oszillationen in x-Richtung auf, die sich in t-Richtung verstärken.
- b)  $\vartheta \geq \frac{1}{2}$ :
  Die für  $\vartheta = 0$  beobachteten Oszillationen treten in der numerischen Lösung nicht mehr auf. Der Mehraufwand (d.h. das Lösen eines linearen Gleichungssystems in jedem Zeitschritt) wird aber durch die impliziten Verfahren mehr als wettgemacht. Die impliziten Verfahren sind also für die Anfangsrandwert-
- **3.2 Bemerkung** (Lösbarkeit von (3-12) im nichtlinearen Fall f(u, x, t)). Sei  $f \in C^1(\mathbb{R} \times [0, 1] \times [0, T], \mathbb{R})$ , und es gelte

$$\frac{\partial f}{\partial u}(u, x, t) < \mu$$

für  $u \in \mathbb{R}$ ,  $0 \le x \le 1$  und  $0 \le t \le T$ . Dann lässt sich Folgendes zeigen: Das implizite  $\vartheta$ -Verfahren  $(\vartheta > 0)$ 

$$v^{j+1} = v^j + \Delta t \cdot \left[ \vartheta F_{\Delta x}(v^{j+1}, t_{j+1}) + (1 - \vartheta) F_{\Delta x}(v^j, t_i) \right], j = 0, \dots, N - 1$$

ist auf jedem Zeitlevel  $t_j=j\Delta t$  für alle  $\Delta t$  mit  $\vartheta\mu\Delta t<1$  durchführbar.

# b) Konsistenz der Differenzenverfahren

Wir bringen die Differentialgleichung aus Abschnitt a), d.h. die vollständige Diskretisierung auf die Form

$$T^h(u) = 0, \quad u \in \mathbb{R}^{\Omega_h}, \quad T^h : \mathbb{R}^{\Omega_h} \longrightarrow \mathbb{R}^{\Omega_h}$$

mit der Familie von Gittern  $\{\Omega_h\}_{0 < h \le h_0}$ . Wir betrachten hier die Anfangsrandwertaufgabe

$$u_t = u_{xx} + f(u, x, t) \text{ in } \Omega = (0, 1) \times (0, T),$$
 (3-13)  
 $u(x, 0) = u_0(x) \text{ für } 0 \le x \le 1,$   
 $u(0, t) = \gamma_0(t), u(1, t) = \gamma_1(t) \text{ für } 0 \le t \le T.$ 

Dazu setzen wir  $h=(\Delta x,\Delta t)$  und  $\Delta x=\frac{1}{M},\,\Delta t=\frac{T}{N}$  sowie

$$\Omega_h = \{(x_i, t_j) = (i\Delta x, j\Delta t) \mid i = 1, \dots, M - 1, j = 0, \dots, N\}.$$

Für  $u \in \mathbb{R}^{\Omega_h}$  gilt

$$u = (u_1^0, \dots, u_{M-1}^0, u_1^1, \dots, u_{M-1}^1, \dots, u_1^N, \dots, u_{M-1}^N),$$
  

$$u_i^j = u(x_i, t_j) = u(i\Delta x, j\Delta t), \quad i = 1, \dots, M-1, j = 0, \dots, N.$$

Wir finden dann  $T^h(u) = 0$  mit

$$(T^{h}(u))_{i}^{j} = \begin{cases} u_{i}^{0} - u_{0}(x_{i}), & \begin{cases} j = 0\\ i = 1, \dots, M - 1 \end{cases} \\ \frac{1}{\Delta t} \left( u_{i}^{j} - u_{i}^{j-1} \right) \\ -\vartheta \left[ \frac{1}{\Delta x^{2}} \left( u_{i-1}^{j} - 2u_{i}^{j} + u_{i+1}^{j} \right) + f(u_{i}^{j}, x_{i}, t_{j}) \right] & \begin{cases} j = 1, \dots, M - 1 \end{cases} \\ i = 1, \dots, M - 1 \end{cases}$$

$$(3-14)$$

Hierbei ist

$$u_0^j = \gamma_0(t_j), \quad u_0^{j-1} = \gamma_0(t_{j-1}), \quad u_M^j = \gamma_1(t_j), \quad u_M^{j-1} = \gamma_1(t_{j-1})$$

zu setzen.

Nach Kapitel 1 haben wir die Begriffe Konsistenz, Stabilität und Konvergenz für  $T^h(u) = 0$  definiert. Wir untersuchen die Konsistenz für Anfangswertaufgabe (3-13).

Sei  $\overline{u} \in C^2(\Omega) \cap C(\overline{\Omega})$  eine klassische Lösung von (3-13) und  $\overline{u}|_h = \overline{u}|_{\Omega_h}$  bezeichne die Einschränkung von  $\overline{u}$  auf  $\Omega_h$ . Wir untersuchen den Konsistenzfehler  $||T^h(\overline{u}_h)||_{\infty}$  in der Maximumsnorm.

Sei  $w \in \mathbb{R}^{\Omega_h}$  mit

$$||w||_{\infty} = \max\{|w_i^j| \mid i = 1, \dots, M - 1, \quad j = 0, \dots, N\}.$$

Für j = 0, i = 1, ..., M - 1 gilt

$$(T^h(\overline{u}_h))_i^0 = \underbrace{\overline{u}_i^0}_{=u_0(x_i)} -u_0(x_i) = 0.$$

Für  $j=1,\ldots,N,\,i=1,\ldots,M-1$  erhalten wir

$$\begin{split} (T^{h}(\overline{u}_{h}))_{i}^{j} &= \frac{1}{\Delta t} \left( \overline{u}_{i}^{j} - \overline{u}_{i}^{j-1} \right) \\ &- \vartheta \left[ \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j} - 2 \overline{u}_{i}^{j} + \overline{u}_{i+1}^{j} \right) + f(\overline{u}_{i}^{j}, x_{i}, t_{j}) \right] \\ &- (1 - \vartheta) \left[ \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j-1} - 2 \overline{u}_{i}^{j-1} + \overline{u}_{i+1}^{j-1} \right) + f(\overline{u}_{i}^{j-1}, x_{i}, t_{j-1}) \right] \\ &= \frac{1}{\Delta t} (\overline{u}_{i}^{j} - \overline{u}_{i}^{j-1}) - \vartheta(\overline{u}_{t})_{i}^{j} - (1 - \vartheta)(\overline{u}_{t})_{i}^{j-1} \\ &+ \vartheta \left\{ \underbrace{\left( \overline{u}_{xx} \right)_{i}^{j} + f(\overline{u}_{i}^{j}, x_{i}, t_{j})}_{=(\overline{u}_{t})_{i}^{j}} - \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j} - 2 \overline{u}_{i}^{j} + \overline{u}_{i+1}^{j} \right) - f(\overline{u}_{i}^{j}, x_{i}, t_{j}) \right\} \\ &+ (1 - \vartheta) \left\{ \underbrace{\left( \overline{u}_{xx} \right)_{i}^{j-1} + f(\overline{u}_{i}^{j-1}, x_{i}, t_{j-1})}_{=(\overline{u}_{t})_{i}^{j-1}} - \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j-1} - 2 \overline{u}_{i}^{j-1} + \overline{u}_{i+1}^{j-1} \right) \right. \\ &- f(\overline{u}_{i}^{j-1}, x_{i}, t_{j-1}) \right\} \\ &= \frac{1}{\Delta t} \left( \overline{u}_{i}^{j} - \overline{u}_{i}^{j-1} \right) - \vartheta(\overline{u}_{t})_{i}^{j} - (1 - \vartheta)(\overline{u}_{t})_{i}^{j-1} \\ &+ \vartheta \left\{ (\overline{u}_{xx})_{i}^{j} - \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j} - 2 \overline{u}_{i}^{j} + \overline{u}_{i+1}^{j} \right) \right\} \\ &+ (1 - \vartheta) \left\{ (\overline{u}_{xx})_{i}^{j-1} - \frac{1}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j-1} - 2 \overline{u}_{i}^{j-1} + \overline{u}_{i+1}^{j-1} \right) \right\}. \end{split}$$

Nach Kapitel 1 gilt für ein C > 0

$$\left| \frac{1}{\Delta x^2} \left( \overline{u}_{i-1}^k - 2\overline{u}_i^k + \overline{u}_{i+1}^k \right) - (\overline{u}_{xx})_i^k \right| \le C \cdot \Delta x^2, \quad k = j, j - 1,$$

falls  $\frac{\partial^{\nu} \overline{u}}{\partial x^{\nu}} \in C(\overline{\Omega}), \ \nu = 1, 2, 3, 4.$ 

Zur Abschätzung des übrigen Terms betrachten wir

$$\frac{1}{\Delta t}(w(t+\Delta t)-w(t))-\vartheta w'(t+\Delta t)-(1-\vartheta)w'(t).$$

Mit

$$w(t + \Delta t) = w(t) + \Delta t w'(t) + \frac{\Delta t^2}{2} w''(t) + O(\Delta t^3),$$
  
$$w'(t + \Delta t) = w'(t) + \Delta t w''(t) + O(\Delta t^2)$$

für  $w \in C^3([0,T])$  folgt

$$\begin{split} \frac{1}{\Delta t} (w(t + \Delta t) - w(t)) - \vartheta w'(t + \Delta t) - (1 - \vartheta)w'(t) \\ &= w'(t) + \frac{\Delta t}{2} w''(t) + O(\Delta t^2) - \vartheta(w'(t) + \Delta t w''(t) + O(\Delta t^2)) - (1 - \vartheta)w'(t) \\ &= \Delta t w''(t) \left(\frac{1}{2} - \vartheta\right) + O(\Delta t^2), \end{split}$$

falls  $w \in C^3([0,T])$ .

Somit folgt

$$\left| \frac{1}{\Delta t} \left( \overline{u}_i^j - \overline{u}_i^{j-1} \right) - \vartheta(\overline{u}_t)_i^j - (1 - \vartheta)(\overline{u}_t)_i^{j-1} \right| \leq C \begin{cases} \Delta t, \text{falls } \overline{u}_t, \overline{u}_{tt} \in C(\overline{\Omega}), \\ \Delta t^2, \text{falls } \overline{u}_t, \overline{u}_{tt}, \overline{u}_{ttt} \in C(\overline{\Omega}), \end{cases} \vartheta = \frac{1}{2}.$$

**3.3 Satz.** Das  $\vartheta$ -Verfahren für  $\vartheta \in [0,1]$  zur Anfangsrandwertaufgabe (3-13) ist konsistent der Ordnung 1 in  $\Delta t$  und 2 in  $\Delta x$  bezgl.  $\|\cdot\|_{\infty}$ , d.h.

$$||T^h(\overline{u}_h)||_{\infty} = O(\Delta t + \Delta x^2)$$

an jeder klassischen Lösung von (3-13) mit  $\frac{\partial^{\nu} \overline{u}}{\partial t^{\nu}} \in C(\overline{\Omega}), \ \nu = 1, 2, \ \frac{\partial^{\nu} \overline{u}}{\partial x^{\nu}} \in C(\overline{\Omega}), \ \nu = 0, 1, 2, 3, 4.$ 

Das Crank-Nicholson Verfahren ist sogar  $O(\Delta t^2 + \Delta x^2)$  konsistent, falls zusätzlich  $\frac{\partial^3 \overline{u}}{\partial t^3} \in C(\overline{\Omega})$  gilt.

Unsere Rechnungen zeigen keinen Unterschied im Konsistenzverhalten für das explizite ( $\vartheta=0$ ) und das implizite ( $\vartheta>0$ ) Verfahren. Die im Abschnitt a) beobachteten drastischen Unterschiede müssen daher mit der Stabilität zusammenhängen, die wir als nächstes diskutieren.

**3.4 Bemerkung.** Die Glattheitsbedingungen für  $\overline{u}$  von Satz 3.3 sind oft in den Ecken  $x=0,\ t=0$  bzw.  $x=1,\ t=0$  nicht erfüllt. Beispielsweise erfordert  $\overline{u}\in C(\overline{\Omega})$  für (3-13) die Verträglichkeitsbedingungen (auch Kompatibilitätsbedingungen genannt)  $\gamma_0(0)=u_0(0), \gamma_1(0)=u_0(1)$  und  $u_t, u_{xx}\in C(\overline{\Omega})$  die weiteren Bedingungen

$$\gamma_0'(0) = u_0''(0) + f(\gamma_0(0), 0, 0) \quad (t = 0, x = 0), 
\gamma_1'(0) = u_0''(1) + f(\gamma_1(0), 1, 0) \quad (t = 0, x = 1).$$

Jedoch garantiert die Theorie der parabolischen Differentialgleichungen trotzdem die Glattheit von  $\overline{u}(x,t)$  für  $t>0,\ 0\leq x\leq 1$ , falls die Daten  $f,\ u_0,\ \gamma_0,\ \gamma_1$  hinreichend glatt sind.

# c) Stabilität und Konvergenz der Differenzenverfahren

Wir erinnern an die Stabilitätsdefinition für ein diskretes Modell.

**3.5 Definition.** Es sei  $T^h(u) = 0$ ,  $T^h : \mathbb{R}^{\Omega_h} \longrightarrow \mathbb{R}^{\Omega_h}$  ein diskretes Modell. Sei der Raum  $\mathbb{R}^{\Omega_h}$  mit einer Norm  $\|\cdot\|$  versehen. Gibt es ein von h unabhängiges C > 0 mit

$$||u - v|| \le C||T^h(u) - T^h(v)||, \quad \forall u, v \in \mathbb{R}^{\Omega_h}, \quad 0 < h \le h_0, \quad h = (\Delta x, \Delta t),$$

so heißt das Modell  $T^h$  stabil.

Wir analysieren hier die Stabilität der Wärmeleitungsgleichung

$$u_{t} = u_{xx} \text{ in } (0,1) \times (0,T),$$

$$u(x,0) = u_{0}(x), \quad 0 \le x \le 1,$$

$$u(0,t) = \gamma_{0}, u(1,t) = \gamma_{1} \text{ für } 0 \le t \le T.$$
(3-15)

Mit den Vektoren  $v^j=(u^j_1,\ldots,u^j_{M-1}),\,j=0,\ldots,N,$  sowie den Schrittweiten  $\Delta x=\frac{1}{M}>0,\,\Delta t=\frac{T}{N}>0$  hat das  $\vartheta$ -Verfahren die Form

$$v^{0} = (u_{0}(x_{1}), \dots, u_{0}(x_{M-1})) =: r^{0},$$
  
$$v^{j+1} = v^{j} + \Delta t \left[ \vartheta F_{\Delta x}(v^{j+1}, t_{j+1}) + (1 - \vartheta) F_{\Delta x}(v^{j}, t_{j}) \right]$$

mit

$$F_{\Delta x}(v,t) = -\frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix} v + \frac{1}{\Delta x^2} \begin{pmatrix} \gamma_0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \\ \gamma_1 \end{pmatrix}$$

$$=: -\Gamma v + r^1,$$

d.h.

$$v^{j+1} = v^j + \Delta t \left[ \vartheta(-\Gamma v^{j+1} + r^1) + (1 - \vartheta)(-\Gamma v^j + r^1) \right], \quad v^0 = r^0.$$

Letztere Formelzeile ist mit

$$(I + \Delta t \vartheta \Gamma) v^{j+1} = (I - \Delta t (1 - \vartheta) \Gamma) v^j + \Delta t r^1, \quad v^0 = r^0$$

und somit mit

$$\underbrace{\left(\frac{1}{\Delta t}I + \vartheta\Gamma\right)}_{=:A} v^{j+1} = \underbrace{\left(\frac{1}{\Delta t}I - (1 - \vartheta)\Gamma\right)}_{=:B} v^{j} + r^{1}, j = 0, \dots, N - 1 \quad (3-16)$$

$$v^0 = r^0$$

äquivalent.

In Matrixform lautet die Gleichung (3-16)

$$\begin{pmatrix}
I & 0 & 0 & \dots & 0 & 0 & 0 \\
-B & A & 0 & \dots & 0 & 0 & 0 & 0 \\
0 & -B & A & \dots & 0 & 0 & 0 & 0 \\
\dots & \dots & \dots & \ddots & \dots & \dots & \dots & \dots \\
0 & 0 & 0 & \dots & A & 0 & 0 & 0 \\
0 & 0 & 0 & \dots & -B & A & 0 & 0 \\
0 & 0 & 0 & \dots & 0 & -B & A
\end{pmatrix}
\begin{pmatrix}
v^0 \\
v^1 \\
v^2 \\
\vdots \\
v^{N-2} \\
v^{N-1} \\
v^N
\end{pmatrix}
-
\begin{pmatrix}
r^0 \\
r^1 \\
r^1 \\
\vdots \\
r^1 \\
r^1
\end{pmatrix}
=
\begin{pmatrix}
0 \\
0 \\
0 \\
\vdots \\
0 \\
0 \\
0
\end{pmatrix} (3-17)$$

Wir müssen die Stabilität von Dreiecksblockmatrizen der obigen Form untersuchen.

**3.6 Lemma.** Sei eine Norm  $\|\cdot\|_*$  auf  $\mathbb{R}^m$  gegeben. Die Matrizen  $A(\Delta t), B(\Delta t) \in \mathbb{R}^{m,m}$  mögen von  $\Delta t \in (0,T)$  abhängen.  $A(\Delta t)$  sei invertierbar und erfülle

$$||A(\Delta t)^{-1}||_* \le C_1 \cdot \Delta t, \quad \forall \Delta t \in (0, T]. \tag{3-18}$$

Ferner existiere ein  $C_2 > 0$  mit

$$\|(A^{-1}B)^n(\Delta t)\|_* \le C_2 \quad \text{für} \quad 0 \le n \cdot \Delta t \le T, n \in \mathbb{N}. \tag{3-19}$$

Dann gilt für die (n+1)-blockige Matrix

$$H(\Delta t) = \begin{pmatrix} I & 0 & 0 & \dots & 0 & 0 & 0 \\ -B & A & 0 & \dots & 0 & 0 & 0 \\ 0 & -B & A & \dots & 0 & 0 & 0 \\ \dots & \dots & \dots & \ddots & \dots & \dots \\ 0 & 0 & 0 & \dots & A & 0 & 0 \\ 0 & 0 & 0 & \dots & -B & A & 0 \\ 0 & 0 & 0 & \dots & 0 & -B & A \end{pmatrix}$$
(3-20)

 $mit \ 0 \le n\Delta t \le T \ die \ Stabilit "atsungleichung"$ 

$$||v||_{*,\infty} \le C_2(1+C_1T) \cdot ||H(\Delta t)v||_{*,\infty}, \quad \forall v \in \mathbb{R}^{m(n+1)}, \quad 0 \le n\Delta t \le T, \quad (3-21)$$

wobei

$$||v||_{*,\infty} = ||(v^0, v^1, \dots, v^n)||_{*,\infty} = \max\{||v^i||_* \mid i = 0, \dots, n\}.$$

**3.7 Bemerkung.** Setzt man  $C(\Delta t) = (A^{-1}B)(\Delta t)$ , so lautet (3-19)

$$||C(\Delta t)^n||_* \le C_2$$
 für  $0 \le n\Delta t \le T$ .

Diese Bedingung wird in der Literatur oft zur Definition der Stabilität eines Verfahrens der Form (3-17) herangezogen. Sie erweist sich auch als notwendig für ein konsistentes und konvergentes Verfahren (Lax'scher Äquivalenzsatz).

#### **3.8 Bemerkung.** Hinreichend für (3-19) ist

$$\|(A^{-1}B)(\Delta t)\|_* \le 1 + C_3 \Delta t. \tag{3-22}$$

Dann folgt nämlich

$$\|(A^{-1}B)^n(\Delta t)\|_* \le \|(A^{-1}B)(\Delta t)\|_*^n \le (1 + C_3\Delta t)^n$$
  
 $\le \exp(C_3\Delta t)^n = \exp(C_3n\Delta t)$   
 $\le \exp(C_3T) =: C_2 \text{ für } 0 \le n\Delta t \le T.$ 

Beweis von Lemma 3.6: Die Gleichung

$$H(\Delta t)v = g = (g^0, g^1, \dots, g^n)$$

bedeutet

$$v^0 = g^0, \quad Av^j - Bv^{j-1} = g^j, \quad j = 1, \dots, n$$

oder

$$v^{j} = A^{-1}Bv^{j-1} + A^{-1}g^{j}, \quad j = 1, \dots, n, \quad v^{0} = g^{0}.$$

Wir zeigen nun für  $j \leq n, n\Delta t \leq T$  die Darstellung

$$v^{j} = (A^{-1}B)^{j}g^{0} + \sum_{k=1}^{j} (A^{-1}B)^{j-k}A^{-1}g^{k}.$$

Der Beweis wird durch Induktion über j erbracht:

$$j = 0:$$
  $v^0 = g^0$   
 $j - 1 \to j:$ 

$$v^{j} = A^{-1}Bv^{j-1} + A^{-1}g^{j}$$

$$= (A^{-1}B)\left[ (A^{-1}B)^{j-1}g_{0} + \sum_{k=1}^{j-1} (A^{-1}B)^{j-1-k}A^{-1}g^{k} \right] + A^{-1}g^{j}$$

$$= (A^{-1}B)^{j}g_{0} + \sum_{k=1}^{j-1} (A^{-1}B)^{j-k}A^{-1}g^{j} + A^{-1}g^{j}$$

$$= (A^{-1}B)^{j}g_{0} + \sum_{k=1}^{j} (A^{-1}B)^{j-k}A^{-1}g^{k}.$$

Somit folgt

$$||v^{j}||_{*} \leq ||(A^{-1}B)^{j}||_{*} \cdot ||g^{0}||_{*} + \sum_{k=1}^{j} ||(A^{-1}B)^{j-k}||_{*} \cdot ||A^{-1}||_{*} ||g^{k}||_{*}$$

$$\leq C_{2}||g^{0}||_{*} + \sum_{k=1}^{j} C_{2} \cdot C_{1} \cdot \Delta t \cdot ||g^{k}||_{*}$$

$$\leq C_{2}||g^{0}||_{*} + C_{1}C_{2} \sum_{k=1}^{j} \Delta t \cdot \max\{||g^{k}||_{*} | k = 1, \dots, n\}$$

$$\leq C_{2}(||g^{0}||_{*} + C_{1}T \cdot \max\{||g^{k}||_{*} | k = 1, \dots, n\})$$

$$\leq C_{2}(1 + C_{1}T)||g||_{*,\infty} = C_{2}(1 + C_{1}T)||H(\Delta t)v||_{*,\infty}, \quad j = 0, \dots, n.$$

Dies liefert

$$||v||_{*,\infty} \le C_2(1+C_1T)||H(\Delta t)v||_{*,\infty}.$$

Wir prüfen nun (3-18), (3-19) für das  $\vartheta$ -Verfahren. Wir haben

$$A = \frac{1}{\Delta t}I + \vartheta\Gamma, \quad B = \frac{1}{\Delta t}I - (1 - \vartheta)\Gamma,$$

$$m = M - 1, \quad \|\cdot\|_{*} = \|\cdot\|_{\infty},$$

$$\Gamma = \frac{1}{\Delta x^{2}} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 & 0 \end{pmatrix}.$$

 $A(\Delta t) = \frac{1}{\Delta t}I + \vartheta\Gamma, \ \vartheta \ge 0$  ist eine  $L_0$ -Matrix und

$$A(\Delta t)\mathbb{I} = \frac{1}{\Delta t}\mathbb{I} + \frac{\vartheta}{\Delta x^2} \begin{pmatrix} 1\\0\\\vdots\\0\\1 \end{pmatrix} \ge \frac{1}{\Delta t}\mathbb{I} > 0.$$
 (3-23)

Also ist  $A(\Delta t)$  eine M-Matrix, und es gilt mit (3-23)

$$||A^{-1}(\Delta t)||_{\infty} = ||A^{-1}(\Delta t)\mathbb{I}||_{\infty} \le ||\Delta t\mathbb{I}||_{\infty} = \Delta t.$$

(3-18) gilt also mit  $C_1 = 1$ .

Es gilt  $B(\Delta t) = \frac{1}{\Delta t}I - (1 - \vartheta)\Gamma$ . Somit ist die Bedingung  $B(\Delta t) \geq 0$  mit

$$\frac{1}{\Delta t} - (1 - \vartheta) \frac{2}{\Delta x^2} \ge 0$$

äquivalent, da alle außerdiagonalen Elemente von  $B(\Delta t)$  größer gleich Null sind. Daraus folgt

$$1 \ge 2(1 - \vartheta) \frac{\Delta t}{\Delta x^2},$$

was mit

$$\frac{\Delta t}{\Delta x^2} \le \frac{1}{2(1-\vartheta)}. (3-24)$$

gleichbedeutend ist.

Wir setzen nun (3-24) voraus und schreiben weiter:

$$(A(\Delta t) - B(\Delta t))\mathbb{I} = \vartheta \Gamma \mathbb{I} + (1 - \vartheta)\Gamma \mathbb{I} = \Gamma \mathbb{I} \ge 0 \implies$$
$$A(\Delta t)\mathbb{I} \ge B(\Delta t)\mathbb{I} \implies$$
$$\mathbb{I} \ge (A^{-1}B)(\Delta t)\mathbb{I}, \text{ da } A^{-1}(\Delta t) \ge 0.$$

(3-24) stellt nun  $B(\Delta t) \geq 0$  sicher, und wir haben  $(A^{-1}B)(\Delta t) \geq 0$ . Somit folgt

$$\|(A^{-1}B)(\Delta t)\|_{\infty} = \|(A^{-1}B)(\Delta t)\mathbb{I}\|_{\infty} \le \|\mathbb{I}\|_{\infty} = 1.$$

Also sind die Bedingungen (3-22) und (3-21) mit  $C_3 = 0$  bzw.  $C_2 = 1$  erfüllt. Nach Lemma 3.6 erfüllt  $H(\Delta t)$  die Stabilitätsungleichung (3-21) mit der Stabilitätskonstanten  $C_2(1 + C_1T) = 1 + T$ , die sowohl von  $\Delta x$  als auch von  $\Delta t$  unabhängig ist. Da  $T^h$  bis auf inhomogene Terme mit  $H(\Delta t)$  übereinstimmt, folgt die Stabilitätsungleichung mit  $\|\cdot\| = \|\cdot\|_{*,\infty} = \|\cdot\|_{\infty}$ .

#### 3.1 Satz Unter der Bedingung

$$\frac{\Delta t}{\Delta x^2} \le \frac{1}{2(1-\vartheta)}$$

ist das  $\vartheta$ -Verfahren für die Wärmeleitungsgleichung bezgl.  $\|\cdot\|_{\infty}$  auf  $\mathbb{R}^{\Omega_h}$  stabil. Genauer gilt

$$||u-v||_{\infty} \le (1+T)||T^h(u)-T^h(v)||_{\infty}, \quad \forall u,v \in \mathbb{R}^{\Omega_h}, \quad h = (\Delta x, \Delta t),$$

wobei  $T^h: \mathbb{R}^{\Omega_h} \longrightarrow \mathbb{R}^{\Omega_h}$  wie in (3-14) mit  $f \equiv 0$  definiert ist.

**3.9 Bemerkung.** Die Restriktion  $\frac{\Delta t}{\Delta x^2} \leq \frac{1}{2}$  für das explizite Verfahren und  $\frac{\Delta t}{\Delta x^2} \leq \infty$  für das rein implizite Verfahren stimmen mit unseren Beobachtungen überein. Bei  $\vartheta = \frac{1}{2}$  können wir jedoch  $\frac{\Delta t}{\Delta x^2} \leq 1$  verletzen. Wir werden aber zeigen, dass bzgl. einer anderen Norm das  $\vartheta$ -Verfahren für  $\vartheta \in \left[\frac{1}{2},1\right]$  bedingungslos stabil ist. Diese Norm stellt eine Verallgemeinerung der  $L^2((0,1))$ -Norm auf die Gitterfunktionen dar.

**3.10 Korollar.** Unter der Voraussetzung  $\frac{\Delta t}{\Delta x^2} \leq \frac{1}{2(1-\vartheta)}$  ist das  $\vartheta$ -Verfahren,  $0 \leq \vartheta \leq 1$ , bzgl. der Maximumsnorm konvergent der Ordnung 1 in  $\Delta t$  und 2 in  $\Delta x$  an der Lösung  $\overline{u}$  der Wärmeleitungsgleichung, falls  $\frac{\partial^{\nu}}{\partial t^{\nu}}\overline{u} \in C(\overline{\Omega})$ ,  $\nu = 1, 2, \frac{\partial^{\nu}}{\partial x^{\nu}}\overline{u} \in C(\overline{\Omega})$ ,  $\nu = 0, 1, 2, 3, 4$ ,  $\Omega = (0, 1) \times (0, T)$ , d.h.

$$\max\{|\overline{u}(x_i,t_i)-u^h(x_i,t_i)|\,|\,i=1,\ldots,M-1,\,j=0,\ldots,N\} \le C\cdot(\Delta t + \Delta x^2),$$

wobei  $u^h$  die Lösung von  $T^h(u)=0$  bezeichnet. Für das Crank-Nicholson Verfahren gilt sogar

$$\max\{|\overline{u}(x_i, t_j) - u^h(x_i, t_j)| | i = 1, \dots, M - 1, j = 0, \dots, N\} \le C \cdot (\Delta t^2 + \Delta x^2),$$

falls zusätzlich  $\frac{\partial^3}{\partial t^3}\overline{u} \in C(\overline{\Omega})$  gilt.

Stabilität der Wärmeleitungsgleichung bezüglich  $\|\cdot\|_{2,*}$ 

Setzt man

$$||v||_* = ||v||_2 = \left(\Delta x \sum_{i=1}^{M-1} v_i^2\right)^{1/2}, v \in \mathbb{R}^{M-1},$$

so stimmt  $\|\cdot\|_2$  bis auf den Fehler der Ordnung  $\sqrt{\Delta x}$  mit der euklidischen Norm überein, falls v Riemann-integrierbar ist. Im Grenzwert  $\Delta x \to 0$  finden wir

$$\Delta x \sum_{i=1}^{M-1} v_i^2 \stackrel{\Delta x \to 0}{\longrightarrow} \int_0^1 v^2(x) \, \mathrm{d}x = \|v\|_{L^2(0,1)}^2$$

aufgrund der Konvergenz der Riemannschen Summe gegen das Riemannsche Integral.

Ferner gilt für eine symmetrische Matrix  $A \in \mathbb{R}^{M-1,M-1}$ :

$$||A||_2 = \sup\{||Av||_2 \mid v \in \mathbb{R}^{M-1}, ||v||_2 = 1\}$$
  
=  $\max\{|\lambda| \mid \lambda \in \sigma(A) \subset \mathbb{R}\}.$ 

#### 3.11 Lemma. Es sei

$$C = \frac{1}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix} \in \mathbb{R}^{M-1,M-1}, \quad \Delta x = \frac{1}{M} > 0.$$

C hat die Eigenwerte

$$\lambda_k = \frac{2}{\Delta x^2} \left( 1 - \cos\left(\frac{k\pi}{M}\right) \right), \quad k = 1, \dots, M - 1$$

mit den Eigenvektoren  $v^k = (v_1^k, \dots, v_{M-1}^k) \in \mathbb{R}^{M-1}$ ,

$$v_i^k = \sin\left(\frac{ik\pi}{M}\right), \quad i = 1, ..., M - 1, \quad k = 1, ..., M - 1.$$

Beweis: Man rechne nach!

Wir kehren nun zur Wärmeleitungsgleichung zurück und berechnen  $||A^{-1}(\Delta t)||_2$  und  $||(A^{-1}B)(\Delta t)||_2$  für

$$A(\Delta t) = \frac{1}{\Delta t}I + \vartheta C, \quad B(\Delta t) = \frac{1}{\Delta t}I - (1 - \vartheta)C.$$

Gemäß Lemma ?? hat C die Eigenwerte  $\lambda_k = \frac{2}{\Delta x^2} \left(1 - \cos\left(\frac{k\pi}{M}\right)\right), k = 1, \dots, M - 1$ . Für  $\lambda_k$  gilt die Abschätzung

$$0 < \lambda_k < \frac{4}{\Lambda r^2}, \quad k = 1, \dots, M - 1.$$

Die Matrizen  $A(\Delta t)$ ,  $B(\Delta t)$ ,  $A^{-1}(\Delta t)$  und C sind symmetrisch. Ferner gilt  $(BA)(\Delta t) = (AB)(\Delta t)$  nach Definition von  $A(\Delta t)$  und  $B(\Delta t)$ . Also folgt

$$(A^{-1}B)^{T}(\Delta t) = (BA^{-1})^{T}(\Delta t) = ((A^{-1})^{T}B^{T})(\Delta t) = (A^{-1}B)(\Delta t),$$

d.h.  $(A^{-1}B)(\Delta t)$  ist symmetrisch.

 $A(\Delta t)$  hat die Eigenwerte  $\frac{1}{\Delta t} + \vartheta \lambda_k$ . Somit sind

$$\frac{1}{\frac{1}{\Delta t} + \vartheta \lambda_k} = \frac{\Delta t}{1 + \Delta t \vartheta \lambda_k}, \quad k = 1, \dots, M - 1$$

die Eigenwerte von  $A^{-1}(\Delta t)$ , und wir finden

$$||A^{-1}(\Delta t)||_2 = \max\left\{\left|\frac{\Delta t}{1 + \Delta t \vartheta \lambda_k}\right| \mid k = 1, \dots, M - 1\right\} \le \Delta t,$$

d.h.  $C_1 = 1$  in Lemma 3.6.

Die Eigenwerte von  $(A^{-1}B)(\Delta t) = \left(\frac{1}{\Delta t}I + \vartheta C\right)^{-1} \left(\frac{1}{\Delta t}I - (1-\vartheta)C\right)$  sind

$$\mu_k = \frac{\frac{1}{\Delta t} - (1 - \vartheta)\lambda_k}{\frac{1}{\Delta t} + \vartheta\lambda_k} \le 1, \quad k = 1, \dots, M - 1.$$

Man beachte dazu, dass  $A(\Delta t)$ ,  $B(\Delta t)$ ,  $A^{-1}(\Delta t)$ ,  $(A^{-1}B)(\Delta t)$  dieselbe Basis  $v^1, \ldots, v^{M-1}$  aus Eigenvektoren wie C besitzen. Wir finden also

$$||(A^{-1}B)(\Delta t)||_2 \le 1$$
, falls  $\mu_k \ge -1$ ,  $k = 1, \dots, M-1$ .

Dies ist äquivalent zu

$$\mu_k \ge -1 \Leftrightarrow \frac{1}{\Delta t} - (1 - \vartheta)\lambda_k \ge -\frac{1}{\Delta t} - \vartheta\lambda_k$$
  
 $\Leftrightarrow \frac{2}{\Delta t} \ge (1 - 2\vartheta)\lambda_k, \quad k = 1, \dots, M - 1.$ 

Für  $\vartheta \geq \frac{1}{2}$  ist dies stets erfüllt, während sich für  $\vartheta < \frac{1}{2}$  die Bedingung

$$\lambda_k \le \frac{2}{\Delta t(1-2\vartheta)}, \quad k = 1, \dots, M-1$$

ergibt. Dies ist insbesondere abgesichert, wenn

$$\lambda_k \le \frac{4}{\Delta x^2} \le \frac{2}{\Delta t (1 - 2\vartheta)} \Leftrightarrow \frac{\Delta t}{\Delta x^2} \le \frac{1}{2(1 - 2\vartheta)}.$$

gilt. Das Lemma 3.6 ist also mit  $C_1 = C_2 = 1$  anwendbar.

#### 3.12 Satz. Unter der Bedingung

$$\frac{\Delta t}{\Delta x^2} \le \begin{cases} \infty, & \text{für } \vartheta \ge \frac{1}{2} \\ \frac{1}{2(1-2\vartheta)}, & \text{für } \vartheta < \frac{1}{2} \end{cases}$$

ist das θ-Verfahren für die Wärmeleitungsgleichung bezüglich der Norm

$$||u||_{2,\infty} = \max\left\{\left(\Delta x \sum_{i=1}^{M-1} (u_i^j)^2\right)^{1/2} | j = 0,\dots, N\right\}, \quad u \in \mathbb{R}^{\Omega_h}, \ h = (\Delta x, \Delta t)$$

stabil. An jeder Lösung  $\overline{u}$  der linearen Wärmeleitungsgleichung mit

$$\frac{\partial^{\nu} \overline{u}}{\partial t^{\nu}} \in C([0,1] \times [0,T]), \nu = 1, 2, \quad \frac{\partial^{\nu} \overline{u}}{\partial x^{\nu}} \in C([0,1] \times [0,T]), \nu = 1, 2, 3, 4$$

liegt Konvergenz gemäß

$$\|\overline{u}_h - u^h\|_{2,\infty} = O(\Delta t + \Delta x^2), \quad T^h(u^h) = 0$$

vor. Für das Crank-Nicholson Verfahren gilt sogar

$$\|\overline{u}_h - u^h\|_{2,\infty} = O(\Delta t^2 + \Delta x^2),$$

falls zusätzlich  $\frac{\partial^3 \overline{u}}{\partial t^3} \in C([0,1] \times [0,T]).$ 

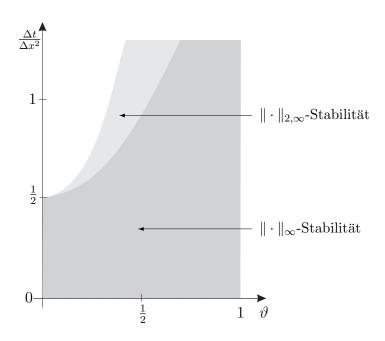


Abbildung 13: Stabilitätsbereiche

# 4. Hyperbolische Differentialgleichungen

# a) Differenzenverfahren für die Wellengleichung

Eines der bekanntesten Beispiele für hyperbolische Differentialgleichungen ist die Wellengleichung. Wir betrachten zunächst die analytische Lösung der Wellengleichung als reine Anfangswertaufgabe

$$u_{tt} = c^2 u_{xx}, \quad x \in \mathbb{R}, \ t \ge 0,$$
 $u(x,0) = u_0(x), \quad u_t(x,0) = u_1(x), \ x \in \mathbb{R}.$ 

$$(4-1)$$

Dies modeliert einen unendlich ausgedehnten Stab oder eine unendliche Saite.

Mit dem Ansatz

$$u(x,t) = f_1(x+ct) + f_2(x-ct), \quad f_1, f_2 \in C^2(\mathbb{R})$$

findet man

$$u_{tt}(x,t) = c^2 f_1''(x+ct) + (-c)^2 f_2''(x-ct)$$
  
=  $c^2 (f_1''(x+ct) + f_2''(x-ct)) = c^2 u_{xx}(x,t),$ 

d.h.  $f_1(x+ct) + f_2(x-ct)$  genügt der Differentialgleichung (4-1).

Man kann nun auch die Anfangsbedingungen abgleichen und erhält

$$u(x,t) = \frac{1}{2}(u_0(x+ct) + u_0(x-ct)) + \frac{1}{2c} \int_{x-ct}^{x+ct} u_1(\xi)d\xi.$$
 (4-2)

(4-2) ist als d'Alembert-Formel bekannt.

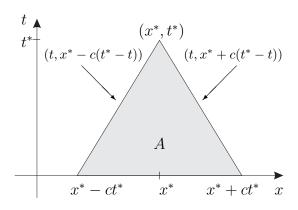


Abbildung 14: Abhängigkeitsbereich von  $(x^*, t^*)$ 

Die auf der Abbildung (14) mit den Pfeilen gekennzeichneten Linien heißen Charakteristiken für (4-1).

Gemäß der Lösungsdarstellung (4-2) hängt der Wert  $u(x^*, t^*)$  von den Anfangsvorgaben im Intervall

$$I_{x^*,t^*} = [x^* - ct^*, x^* + ct^*]$$

ab. Dieses Intervall wird als Abhängigkeitsintervall des Punktes  $(x^*, t^*)$  bezeichnet. Das Dreieck A mit den Ecken  $(x^* - ct^*, 0), (x^* + ct^*, 0), (x^*, t^*)$  wird entsprechend das Abhängigkeitsdreieck von  $(x^*, t^*)$  genannt.

Fixieren wir umgekehrt ein festes  $x^* \in \mathbb{R}$  für t = 0, so beeinflussen  $u_0(x^*)$ ,  $u_1(x^*)$  die Werte der Lösung u in

$$B := \{(x,t) \mid t \ge 0, x^* - ct \le x \le x^* + ct\}.$$

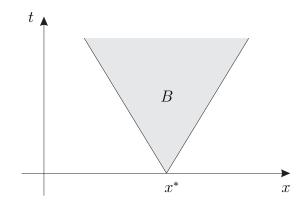


Abbildung 15: Einflussbereich von  $(x^*, t^*)$ 

Wir betrachten nun die Anfangswertaufgabe

$$u_{tt} = c^{2}u_{xx}, \quad 0 \le x \le 1, \ t \ge 0,$$

$$u(x,0) = u_{0}(x), \quad u_{t}(x,0) = u_{1}(x), \quad 0 \le x \le 1,$$

$$u(0,t) = u(1,t) = 0, \quad t \ge 0.$$

$$(4-3)$$

Ferner setzen wir die Verträglichkeitsbedingungen

$$u_0(0) = u_0(1) = 0, \quad u_1(0) = u_1(1) = 0$$

voraus.

Wir möchen nachweisen, dass die d'Alembertsche Lösungsformel (4-2) auch das Problem (4-3) löst. Dazu setzen wir  $u_0, u_1$  2-periodisch folgenderweise auf ganz  $\mathbb{R}$  fort:

$$u_i(1+x) = -u_i(1-x), \quad 0 \le x \le 1,$$
  
 $u_i(x) = u_i(x-2n), \quad \text{falls } 2n \le x \le 2n+2, \ n \in \mathbb{Z}, \ i=0,1.$ 

Hiermit erreicht man

$$u_i(x) = -u_i(-x),$$

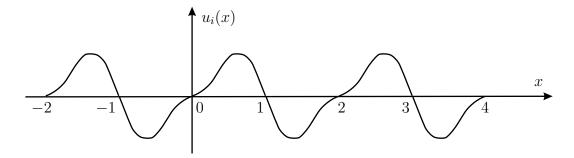


Abbildung 16: Periodische Fortsetzung

$$u_i(1+x) = -u_i(1-x), x \in \mathbb{R}, i = 0, 1,$$

d.h. die fortgesetzten Funktionen sind punktsymmetrisch zu (0,0) und (1,0).

Dann folgt für die Lösung der Wellengleichung (4-1) auf ganz  $\mathbb{R}$  überdies

$$u(0,t) = \frac{1}{2} \left( \underbrace{u_0(ct)}_{=-u_0(-ct)} + u_0(-ct) \right) + \frac{1}{2c} \int_{-ct}^{ct} u_1(\xi) d\xi$$
$$= \frac{1}{2c} \left( \int_{-ct}^{0} -u_1(-\xi) d\xi + \int_{0}^{ct} u_1(\xi) d\xi \right).$$

Mit der Transformation  $\eta=-\xi,\,d\eta=-d\xi$  ergibt sich

$$u(0,t) = \frac{1}{2c} \left( \int_{ct}^{0} -u_1(\eta)(-1)d\eta + \int_{0}^{ct} u_1(\xi)d\xi \right) = 0.$$

Eine analoge Rechnung liefert u(1,t) = 0, d.h. u löst das Problem (4-3).

Allgemeiner betrachten wir nun

$$u_{tt} = c^{2}u_{xx} + f(u_{x}, u, x, t), \quad 0 \le x \le 1, \ 0 \le t \le T,$$

$$u(x, 0) = u_{0}(x), \quad u_{t}(x, 0) = u_{1}(x), \quad 0 \le x \le 1,$$

$$u(0, t) = \gamma_{0}(t), \quad u(1, t) = \gamma_{1}(t), \quad 0 \le t \le T.$$

$$(4-4)$$

Wir führen auf  $[0,1] \times [0,T]$  das Gitter

$$\Omega_h = \{ (i\Delta x, j\Delta t) \mid i = 1, \dots, M - 1, j = 0, \dots, N \},$$
  
$$\Delta x = \frac{1}{M}, \quad \Delta t = \frac{T}{N}, \quad h = (\Delta x, \Delta t)$$

ein. Ohne den Umweg über die Linienmethode stellen wir direkt Differenzengleichungen für die Unbekannten  $u_i^j = u(i\Delta x, j\Delta t), i = 1, ..., M-1, j = 0, ..., N$  auf. Mit  $x_i = i\Delta x, t_j = j\Delta t$  und den symmetrischen Differenzenquotienten für die erste und zweite Ableitung finden wir

$$\frac{1}{\Delta t^2} \left( u_i^{j+1} - 2u_i^j + u_i^{j-1} \right) = \frac{c^2}{\Delta x^2} \left( u_{i+1}^j - 2u_i^j + u_{i-1}^j \right)$$

+ 
$$f\left(\frac{1}{2\Delta x}\left(u_{i+1}^{j}-u_{i-1}^{j}\right),u_{i}^{j},x_{i},t_{j}\right)$$
 (4-5)

für  $i = 1, \dots, M - 1$  und  $j = 1, \dots, N$ .

Für j = 0 setzen wir natürlich

$$u_i^0 = u_0(x_i), \quad i = 1, \dots, M - 1.$$
 (4-6)

Zusätzlich setzen wir

$$u_0^j = \gamma_0(t_j), \quad u_M^j = \gamma_1(t_j), \quad j = 0, \dots, N.$$
 (4-7)

Uns fehlt noch aber eine Gleichung für  $u_i^1$ , z.B.

$$\frac{1}{\Delta t} \left( u_i^1 - u_i^0 \right) = u_1(x_i), \quad i = 1, \dots, M - 1.$$
 (4-8)

Während (4-5), (4-6), (4-7) bekanntermaßen die Konsistenzordnung  $O(\Delta t^2 + \Delta x^2)$  haben, liegt für (4-8) nur  $O(\Delta t)$  vor. Eine  $O(\Delta t^2)$  Approximation gewinnt man aus der Entwicklung

$$\frac{u(x,\Delta t) - u(x,0)}{\Delta t} = u_t(x,0) + \frac{1}{2}\Delta t \cdot u_{tt}(x,0) + O(\Delta t^2).$$
 (4-9)

Für eine glatte Lösung  $\overline{u}$  von (4-4) gilt nämlich

$$\overline{u}_{tt}(x,0) = c^2 \overline{u}_{xx}(x,0) + f(\overline{u}_x, \overline{u}, x, 0)$$
$$= c^2 u_0''(x) + f(u_0', u_0, x, 0).$$

Man kann daher (4-8) durch

$$\frac{1}{\Delta t} \left( u_i^1 - u_i^0 \right) = u_1(x_i) + \frac{\Delta t}{2} \left( c^2 u_0''(x_i) + f(u_0'(x_i), u_0(x_i), x_i, 0) \right) \tag{4-10}$$

für i = 1, ..., M - 1 ersetzen.

Man bezeichnet (4-5)—(4-7), (4-8) oder (4-10) auch als explizites Verfahren.

Wir können aber (4-5) wiederum als Spezialfall eines von einem Parameter abhängigen Verfahrens ansehen. Setze dazu

$$\sigma_i^j = \frac{c^2}{\Delta x^2} \left( u_{i-1}^j - 2u_i^j + u_{i+1}^j \right) + f \left( \frac{1}{2\Delta x} \left( u_{i+1}^j - u_{i-1}^j \right), u_i^j, x_i, t_j \right)$$
(4-11)

und betrachte dann das im Allgemeinen implizite Verfahren

$$\frac{1}{\Delta t^2} \left( u_i^{j+1} - 2u_i^j + u_i^{j-1} \right) = \vartheta \sigma_i^{j+1} + (1 - 2\vartheta) \sigma_i^j + \vartheta \sigma_i^{j-1}, \quad \vartheta \in \left[ 0, \frac{1}{2} \right]. \tag{4-12}$$

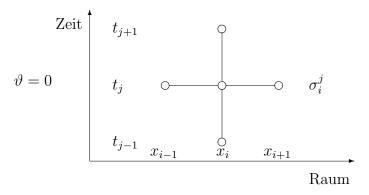


Abbildung 17: Differenzenstern für  $\vartheta=0$ 

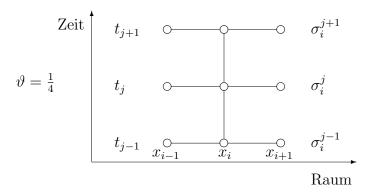


Abbildung 18: Differenzenstern für  $\vartheta=\frac{1}{4}$ 

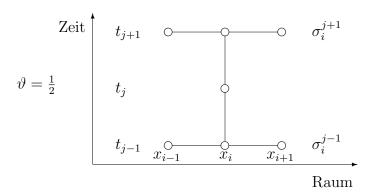


Abbildung 19: Differenzenstern für  $\vartheta = \frac{1}{2}$ 

Die Anfangs- und Randbedingungen werden dabei wieder wie vorhin behandelt.

Auf den Abbildungen 17—19 werden die Differenzensterne für  $\vartheta=0,\frac{1}{4},\frac{1}{2}$  graphisch dargestellt.

Zur Analyse des Konsistenzfehlers schreiben wir das numerische Verfahren wieder

in der Form  $T^h(u) = 0, T^h : \mathbb{R}^{\Omega_h} \longrightarrow \mathbb{R}^{\Omega_h}$  mit

$$(T^{h}(u))_{i}^{j} = \begin{pmatrix} u_{i}^{0} - u_{0}(x_{i}), j = 0, i = 1, \dots, M - 1 \\ \frac{1}{\Delta t} (u_{i}^{1} - u_{i}^{0}) - u_{1}(x_{i}) - \frac{\Delta t}{2} \left[ c^{2} u_{0}''(x_{i}) + f(u_{0}'(x_{i}), u_{0}(x_{i}), x_{i}, 0) \right], \\ j = 1, i = 1, \dots, M - 1 \\ \frac{1}{\Delta t^{2}} \left( u_{i}^{j} - 2u_{i}^{j-1} + u_{i}^{j-2} \right) - \left[ \vartheta \sigma_{i}^{j} + (1 - 2\vartheta) \sigma_{i}^{j-1} + \vartheta \sigma_{i}^{j-2} \right], \\ j = 2, \dots, N, i = 1, \dots, M - 1 \end{pmatrix}$$

und

$$\sigma_i^j = \frac{c^2}{\Delta x^2} \left( u_{i-1}^j - 2u_i^j + u_{i+1}^j \right) + f \left( \frac{1}{2\Delta x} \left( u_{i+1}^j - u_{i-1}^j \right), u_i^j, x_i, t_j \right).$$

**4.1 Lemma.** An einer Lösung  $\overline{u} \in C^4([0,1] \times [0,T])$  von (4-4) ist das parameterabhängige Verfahren  $O(\Delta t^2 + \Delta x^2)$  konsistent, falls f einer Lipschitzbedingung bezüglich der Variablen  $u_x$  gleichmäßig in u, t, x genügt.

Beweis: Sei  $\overline{u}_h$  die Restriktion von  $\overline{u}$  auf  $\Omega_h$ .

Für j = 0 gilt

$$T^h(\overline{u}_h)_i^0 = \overline{u}_i^0 - u_0(x_i) = 0, \quad i = 1, \dots, M - 1.$$

Ist j = 1, so folgt

$$T^{h}(\overline{u}_{h})_{i}^{1} = \frac{1}{\Delta t} \left( \overline{u}_{i}^{1} - \overline{u}_{i}^{0} \right) - u_{1}(x_{i})$$

$$- \frac{\Delta t}{2} \underbrace{\left[ c^{2} u_{0}''(x_{i}) + f(u_{0}'(x_{i}), u_{0}(x_{i}), x_{i}, 0) \right]}_{= c^{2} \overline{u}_{xx}(x_{i}, 0) + f(\overline{u}_{x}(x_{i}, 0), \overline{u}(x_{i}, 0), x_{i}, 0) = \overline{u}_{tt}(x_{i}, 0)}$$

$$= \frac{\overline{u}(x_{i}, \Delta t) - \overline{u}(x_{i}, 0)}{\Delta t} - \overline{u}_{t}(x_{i}, 0) - \frac{\Delta t}{2} \overline{u}_{tt}(x_{i}, 0)$$

$$\stackrel{(4-9)}{=} O(\Delta t^{2}).$$

Für j = 2, ..., N gehen wir wie folgt vor. Es sei

$$\overline{\sigma}_{i}^{j} = \frac{c^{2}}{\Delta x^{2}} \left( \overline{u}_{i-1}^{j} - 2\overline{u}_{i}^{j} + \overline{u}_{i+1}^{j} \right) + f \left( \frac{1}{2\Delta x} \left( \overline{u}_{i+1}^{j} - \overline{u}_{i-1}^{j} \right), \overline{u}_{i}^{j}, x_{i}, t_{j} \right).$$

Dann gilt nach Konstruktion

$$\left| \frac{1}{\Delta x^2} (\overline{u}_{i-1}^j - 2\overline{u}_i^j + \overline{u}_{i+1}^j) - (\overline{u}_{xx})_i^j \right| \leq C_1 \Delta x^2,$$

$$\left| f(\frac{1}{\Delta x} (\overline{u}_{i+1}^j - \overline{u}_{i-1}^j), \overline{u}_i^j, x_i, t_j) - f((\overline{u}_x)_i^j, \overline{u}_i^j, x_i, t_j) \right| \leq$$

$$\leq L \left| \frac{1}{2\Delta x} (\overline{u}_{i+1}^j - \overline{u}_{i-1}^j) - (\overline{u}_x)_i^j \right| \leq C_2 \Delta x^2$$

und somit

$$\overline{\sigma}_i^j = c^2(\overline{u}_{xx})_i^j + f((\overline{u}_x)_i^j, \overline{u}_i^j, x_i, t_j) + O(\Delta x^2)$$
$$= (\overline{u}_{tt})_i^j + O(\Delta x^2).$$

Also folgt für den Konsistenzfehler

$$T^{h}(\overline{u}_{h})_{i}^{j} = \frac{1}{\Delta t^{2}} \left( \overline{u}_{i}^{j} - 2\overline{u}_{i}^{j-1} + \overline{u}_{i}^{j-2} \right)$$

$$- \left[ \vartheta(\overline{u}_{tt})_{i}^{j} + (1 - 2\vartheta)(\overline{u}_{tt})_{i}^{j-1} + \vartheta(\overline{u}_{tt})_{i}^{j-2} \right] + O(\Delta x^{2})$$

$$= (\overline{u}_{tt})_{i}^{j-1} + O(\Delta t^{2}) - \left[ \vartheta(\overline{u}_{tt})_{i}^{j} + (1 - 2\vartheta)(\overline{u}_{tt})_{i}^{j-1} + \vartheta(\overline{u}_{tt})_{i}^{j-2} \right] + O(\Delta x^{2})$$

$$= -\vartheta \left[ \left( \overline{u}_{tt} \right)_{i}^{j} - 2(\overline{u}_{tt})_{i}^{j-1} + (\overline{u}_{tt})_{i}^{j-2} \right] + O(\Delta t^{2} + \Delta x^{2})$$

$$= O(\Delta t^{2} + \Delta x^{2}),$$

wobei hier

$$w(t - \Delta t) - 2w(t) + w(t + \Delta t) = w(t) - \Delta t w'(t) + O(\Delta t^2) - 2w(t) + w(t) + \Delta t w'(t) + O(\Delta t^2) = O(\Delta t^2)$$

für  $w = \overline{u}_{tt}(x,\cdot)$  zu beachten ist.

Im Falle f = 0 in (4-4), d.h. der linearen Wellengleichung, finden wir mit

$$\Gamma = \frac{c^2}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix}, \quad r^j = \frac{c^2}{\Delta x^2} \begin{pmatrix} \gamma_0(t_j) \\ 0 \\ 0 \\ \vdots \\ 0 \\ \gamma_1(t_j) \end{pmatrix}$$

dann die Iteration

$$\frac{1}{\Delta t^2} \left( v^{j+1} - 2v^j + v^{j-1} \right) = \vartheta \left( -\Gamma v^{j+1} + r^{j+1} \right) + (1 - 2\vartheta)(-\Gamma v^j + r^j) + \vartheta(-\Gamma v^{j-1} + r^{j-1}), \quad j = 1, \dots, N - 1.$$

Dies ist äquivalent zu

$$\left(\frac{1}{\Delta t^2}I + \vartheta\Gamma\right)v^{j+1} = \frac{1}{\Delta t^2}\left(2v^j - v^{j-1}\right) + \vartheta r^{j+1} + (1 - 2\vartheta)(-\Gamma v^j + r^j)$$

$$+\vartheta(-\Gamma v^{j-1} + r^{j-1}), \quad j = 1, \dots, N-1$$
 (4-13)

mit den Startdaten

$$v^{0} = (u_{0}(x_{1}), \dots, u_{0}(x_{M-1})),$$

$$\frac{1}{\Delta t}(v^{1} - v^{0}) = (u_{1}(x_{1}), \dots, u_{1}(x_{M-1})) + \frac{\Delta t}{2}c^{2}(u_{0}''(x_{1}), \dots, u_{0}''(x_{M-1})).$$

In jedem Fall erfordert (4-13) auf jedem Zeitlevel die Auflösung eines linearen Gleichungssystems mit der Matrix

$$A = \frac{1}{\Delta t^2} I + \vartheta \Gamma,$$

die offensichtlich eine M-Matrix ist.

# b) Die Courant-Friedrichs-Levy Bedingung

Beim expliziten Differenzenverfahren  $(\vartheta = 0)$  hängt der Wert  $u_i^j$ , also die Näherung für  $\overline{u}(x_i, t_j)$ , von den Startwerten  $u_{i-j}^0, u_{i-j+1}^0, \ldots, u_{i+j}^0$  ab. Man vergleiche hierzu folgenden Differenzenstern:

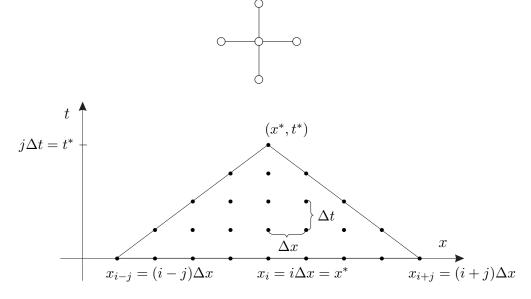


Abbildung 20: Numerisches Abhängigkeitsintervall

Man nennt  $[(i-j)\Delta x, (i+j)\Delta x] = [x_{i-1}, x_{i+j}]$  das numerische Abhängigkeitsintervall des Punktes  $(x^*, t^*) = (x_i, t_j)$ .

Die Wellengleichung selbst liefert das Abhängigkeitsintervall  $[x^* - ct^*, x^* + ct^*]$  für  $(x^*, t^*)$ . Ist dieses Intervall größer als das numerische Abhängigkeitsintervall, so kann

man keine Konvergenz des Verfahrens erwarten, denn Startwerte, die die exakte Lösung beeinflussen, werden für die numerische Lösung gar nicht berücksichtigt. Notwendig ist daher die Courant-Friedrichs-Levy Bedingung (auch CFL-Bedingung genannt).

#### Courant-Friedrich-Levy Bedingung:

Das numerische Abhängigkeitsintervall jedes Punktes  $(x^*, t^*) = (x_i, t_j)$  umfasst das kontinuierliche Abhängigkeitsintervall.

In unserem Fall bedeutet dies  $x_i - ct_j \ge x_{i-j}$  und  $x_{i+j} \ge x_i + ct_j$ .

Dies ist äquivalent zu

$$-ct_j \ge -x_j \text{ und } x_j \ge ct_j,$$

d.h.

$$c \cdot \Delta t \leq \Delta x$$
 bzw.  $\frac{\Delta t}{\Delta x} \leq \frac{1}{c}$ ,

wobei letzteres "CFL-Bedingung" heißt.

Als nächstes wird gezeigt, dass die CFL-Bedingung im Falle der Wellengleichung

$$u_{tt} = c^{2}u_{xx}, \quad 0 \le x \le 1, \ 0 \le t \le T,$$

$$u(x,0) = u_{0}(x), \quad u_{t}(x,0) = u_{1}(x), \quad 0 \le x \le 1,$$

$$u(0,t) = \gamma_{0}(t), \quad u(1,t) = \gamma_{1}(t), \quad 0 \le t \le T$$

$$(4-14)$$

für das explizite Differenzenverfahren ( $\vartheta=0$ ) eine Stabilitätsbedingung bezüglich der Norm  $\|\cdot\|_{2,\infty}$  liefert. Wir gehen analog zum parabolischen Fall vor und setzen  $h=(\Delta x,\Delta t),\,\Delta x=1/M,\,\Delta t=T/N,\,x_i=i\Delta x,\,t_j=j\Delta t$  sowie

$$\Omega_h = \{(i\Delta x, j\Delta t); i = 1, \dots, M - 1, j = 0, \dots, N\}$$

und  $u=(v^0,v^1,\ldots,v^N)\in\mathbb{R}^{\Omega_h},\,v^j=(u^j_1,\ldots,u^j_{M-1}),\,j=0,1,\ldots,N.$  Wir schreiben das explizite diskrete Modell  $(\vartheta=0)$  in der Form  $T^h(u)=0,\,T^h:\mathbb{R}^{\Omega_h}\to\mathbb{R}^{\Omega_h}$  definiert durch

$$(T^h(u))^j = \begin{cases} v^0 - r^0 & \text{für } j = 0, \\ \frac{1}{\Delta t} (v^1 - v^0) - r^1 & \text{für } j = 1, \\ \frac{1}{\Delta t^2} v^j - \left(\frac{2}{\Delta t^2} I - \Gamma\right) v^{j-1} + \frac{1}{\Delta t^2} v^{j-2} - r^j & \text{für } j = 2, \dots, N. \end{cases}$$

 $_{
m mit}$ 

$$r_i^0 = u_0(x_i), \ i = 1, \dots, M - 1,$$

$$r_i^1 = u_1(x_i) + \frac{\Delta t c^2}{2} u_0''(x_i), \ i = 1, \dots, M - 1,$$

$$r^j = \left( -\frac{c^2}{\Delta x^2} \gamma_0(t_{j-1}), 0, \dots, 0, -\frac{c^2}{\Delta x^2} \gamma_1(t_{j-1}) \right), \ j = 2, \dots, N.$$

und

$$\Gamma = \frac{c^2}{\Delta x^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 & 0 & 0 \\ -1 & 2 & -1 & \dots & 0 & 0 & 0 & 0 \\ 0 & -1 & 2 & \dots & 0 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & 2 & -1 & 0 \\ 0 & 0 & 0 & \dots & -1 & 2 & -1 \\ 0 & 0 & 0 & \dots & 0 & -1 & 2 \end{pmatrix}.$$

In Matrixform lautet die Gleichung  $T^h(u) = 0$  dann

$$H_{\Delta x, \Delta t} * \begin{pmatrix} v^0 \\ v^1 \\ v^2 \\ \vdots \\ v^{N-1} \\ v^N \end{pmatrix} - \begin{pmatrix} r^0 \\ r^1 \\ r^2 \\ \vdots \\ r^{N-1} \\ r^N \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{pmatrix}$$

mit

$$H_{\Delta x,\Delta t} = \begin{pmatrix} I & 0 & 0 & \dots & 0 & 0\\ -I/\Delta t & I/\Delta t & 0 & \dots & 0 & 0\\ I/\Delta t^2 & -2I/\Delta t^2 + \Gamma & I/\Delta t^2 & 0 & 0 & 0\\ 0 & \ddots & \ddots & \ddots & 0 & 0\\ 0 & 0 & \ddots & \ddots & \ddots & 0\\ 0 & 0 & 0 & I/\Delta t^2 & -2I/\Delta t^2 + \Gamma & I/\Delta t^2 \end{pmatrix}.$$

$$(4-15)$$

Es ist also die Stabilität von Blockdreiecksmatrizen  $H_{\Delta x,\Delta t}$  der Form (4-15) zu untersuchen. Wir verwenden die Normen

$$||w||_2 = \left(\Delta x \sum_{i=1}^{M-1} w_i^2\right)^{1/2}, \ w \in \mathbb{R}^{M-1},$$

$$||u||_{2,\infty} = \max\{||v^j||_2; \ j = 0, 1, \dots, N\}, \ u = (v^0, v^1, \dots, v^N) \in \mathbb{R}^{\Omega_h}.$$

Die Matrix  $\Gamma$  besitzt die Eigenwerte

$$\lambda_k = \frac{2c^2}{\Delta x^2} (1 - \cos(k\pi \Delta x)), \ k = 1, \dots, M - 1$$

mit den Eigenvektoren  $w^k$ ,

$$w_i^k = \sin(ik\pi\Delta x), \ 1 \le i, k \le M - 1.$$

**4.2 Lemma.** Sei  $H_{\Delta x,\Delta t}$  durch (4-15) gegeben mit einer symmetrischen von  $\Delta x$  abhängigen Matrix  $\Gamma$ , für deren Eigenwerte  $\lambda_k$ ,  $k = 1, \ldots, M-1$  gilt

$$0 < \gamma_1 \le \lambda_k \le \frac{\gamma_2}{\Delta x^2} - \gamma_3, \ k = 1, \dots, M - 1, \ \gamma_2, \gamma_3 > 0.$$
 (4-16)

Dann folgt unter den Schrittweitenbedingungen

$$\Delta t^2 \cdot \max(\gamma_1, \gamma_3) \leq \sigma < 4 \tag{4-17}$$

und

$$\frac{\Delta t}{\Delta x} \le \frac{2}{\sqrt{\gamma_2}} \tag{4-18}$$

die Stabilitätsungleichung

$$||u||_{2,\infty} \le C ||H_{\Delta x,\Delta t}u||_{2,\infty} \quad \forall u = (0, v^1, v^2, \dots, v^N).$$
 (4-19)

**4.3 Bemerkung.** In unserem Fall folgt mit  $\cos(x) \le 1 - \frac{3}{8}x^2$ ,  $0 \le x \le \pi/2$  sogleich

$$\lambda_{k} \geq \lambda_{1} = \frac{2c^{2}}{\Delta x^{2}} (1 - \cos(\pi \Delta x)) \geq \frac{2c^{2}}{\Delta x^{2}} \cdot \frac{3}{8} \pi^{2} \Delta x^{2} = \frac{3}{4} c^{2} \pi^{2},$$

$$\lambda_{k} \leq \lambda_{M-1} = \frac{2c^{2}}{\Delta x^{2}} (1 - \cos(\pi - \pi \Delta x)) = \frac{2c^{2}}{\Delta x^{2}} (1 + \cos(\pi \Delta x))$$

$$\leq \frac{2c^{2}}{\Delta x^{2}} \left(2 - \frac{3}{8} \pi^{2} \Delta x^{2}\right) = \frac{4c^{2}}{\Delta x^{2}} - \frac{3}{4} c^{2} \pi^{2}.$$

(4-16) gilt also mit  $\gamma_2=4c^2, \ \gamma_1=\gamma_3=\frac{3}{4}\pi^2c^2$  und damit geht (4-18) in die CFL-Bedingung über. Ferner ist (4-17) ist äquivalent zu

$$\frac{\Delta t c \pi \sqrt{3}}{2} \leq \sqrt{\sigma} < 2.$$

- **4.4 Bemerkung.** Die Einschränkung  $v^0 = 0$  in Lemma 4.2 ist nicht wesentlich, da der Konsistenzfehler zur Zeit t = 0 verschwindet. Da  $T^h$  bis auf inhomogene Terme mit  $H_{\Delta x, \Delta t}$  übereinstimmt, folgt die Stabilitätsungleichung für  $T^h$  bezgl. der Norm  $\|\cdot\|_{2,\infty}$ .
- **4.5 Korollar.** Die Anfangsrandwertaufgabe (4-14) besitze eine Lösung

$$\overline{u} \in C^4([0,1] \times [0,T]).$$

Dann ist das explizite Verfahren ( $\vartheta = 0$ ) unter der CFL-Bedingung  $\frac{\Delta t}{\Delta x} \leq \frac{1}{c}$  konvergent bezüglich  $\|\cdot\|_{2,\infty}$  der Ordnung 2 in  $\Delta x$  und 2 in  $\Delta t$ .

**4.6 Bemerkung.** Im Falle des allgemeinen  $\vartheta$ -Verfahrens,  $0 \le \vartheta \le \frac{1}{2}$  lässt sich unter der Bedingung

$$\frac{\Delta t}{\Delta x} \leq \begin{cases} \frac{1}{c\sqrt{1-4\vartheta}} & \text{für } 0 \leq \vartheta < 1/4 \\ \infty & \text{für } 1/4 \leq \vartheta \leq 1/2 \end{cases}$$

eine Stabilitätsungleichung für das  $\vartheta$ -Verfahren bezgl.  $\|\cdot\|_{2,\infty}$  nachweisen. Überdies gilt Konvergenz der Ordnung 2 in  $\Delta x$  und  $\Delta t$  bezgl.  $\|\cdot\|_{2,\infty}$  an jeder Lösung  $\overline{u}\in C^4([0,1]\times[0,T])$ .