

Kapitel 5

Finite–Elemente–Methoden (FEM)

5.1 Das Galerkin–Verfahren, Bezeichnungen

Bemerkung 5.1 Grundidee von Finite–Elemente–Methoden, das Ritzsche Verfahren. Sei V ein Hilbert–Raum mit dem Skalarprodukt $a(\cdot, \cdot)$. Wir betrachten das Problem

$$\min_{v \in V} F(v) = \min_{v \in V} \left(\frac{1}{2} a(v, v) - f(v) \right),$$

wobei $f(\cdot) : V \rightarrow \mathbb{R}$ ein beschränktes lineares Funktional ist. Wie bereits bewiesen ist, besitzt das Variationsproblem eine eindeutig bestimmte Lösung $u \in V$, die außerdem die Gleichung

$$a(u, v) = f(v) \quad \forall v \in V \quad (5.1)$$

löst, Satz 4.10 (Rieszscher Darstellungssatz).

Um die Lösung der obigen Probleme mit einem numerischen Verfahren zu approximieren, setzen wir voraus, dass V ein separabler Hilbert–Raum ist, das heißt V besitzt eine abzählbare Basis. Dann gibt es endlich–dimensionale Teilräume $V_1, V_2, \dots \subset V$ mit $\dim V_k = k$, die folgende Eigenschaft besitzen: zu jedem $u \in V$ und $\varepsilon > 0$ gibt es ein $K \in \mathbb{N}$ und ein $u_k \in V_k$ mit

$$\|u - u_k\|_V \leq \varepsilon \quad \forall k \geq K.$$

Es wird dabei nicht verlangt, dass es eine Inklusion der Form $V_k \subset V_{k+1}$ gibt.

Die Ritz–Approximation von (5.1) ist wie folgt definiert. Gesucht ist $u_k \in V_k$ mit

$$a(u_k, v_k) = f(v_k) \quad \forall v_k \in V_k. \quad (5.2)$$

Die wesentliche Idee des Ritzschen Verfahrens besteht also darin, dass man den unendlich–dimensionalen Raum V durch einen endlich–dimensionalen Raum V_k ersetzt. \square

Lemma 5.2 Eigenschaften der Ritzschen Approximation.

1. Der Fehler ist orthogonal zum Raum V_k , das heißt es gilt

$$a(u - u_k, v_k) = 0 \quad \forall v_k \in V_k. \quad (5.3)$$

2. u_k ist die Bestapproximierende von u in V_k bezüglich der von $a(\cdot, \cdot)$ induzierten Norm.



Abbildung 5.1: Walter Ritz (1878 – 1909).

3. Die Folge der Ritz-Approximierenden konvergiert gegen die Lösung von (5.1), das heißt $u_k \rightarrow u$ für $k \rightarrow \infty$.

Beweis: Die Beweise dieser Aussagen sollten bereits in vorangegangenen Vorlesungen gegeben worden sein.

Da endlich-dimensionale Teilräume von Hilbert-Räumen wiederum Hilbert-Räume sind, besitzt nach dem Riesz'schen Darstellungssatz auch die Gleichung der Ritz-Approximation eine eindeutige Lösung, die ebenso ein Minimierungsproblem im Raum V_k löst. Aus der Differenz der Gleichungen (5.1) und (5.2) erhält man die Orthogonalitätsrelation

$$a(u - u_k, v_k) = 0 \quad \forall v_k \in V_k.$$

Das besagt, dass der Fehler $u - u_k$ senkrecht zum Raum V_k ist: $u - u_k \perp V_k$. Demnach ist u_k die orthogonale Projektion von u in den Raum V_k bezüglich des Skalarproduktes von V . Das heißt, u_k ist die Bestapproximierende von u in V_k

$$\|u - u_k\|_V = \inf_{v_k \in V_k} \|u - v_k\|_V.$$

Zum Beweis nutzt man die Orthogonalität (5.3) und die Cauchy-Schwarz-Ungleichung. Sei $w_k \in V_k$ beliebig, dann ist

$$\begin{aligned} \|u - u_k\|_V^2 &= a(u - u_k, u - u_k) = a(u - u_k, u - \underbrace{(u_k - w_k)}_{v_k}) = a(u - u_k, u - v_k) \\ &\leq \|u - u_k\|_V \|u - v_k\|_V. \end{aligned}$$

Da $w_k \in V_k$ beliebig ist, ist auch $v_k \in V_k$ beliebig.

Mit der Bestapproximationseigenschaft erhält man

$$\|u - u_k\|_V = \inf_{v_k \in V_k} \|u - v_k\|_V \leq \varepsilon,$$

woraus schließlich die Konvergenz der Ritz-Approximation $u_k \rightarrow u$ für $k \rightarrow \infty$ folgt. ■

Bemerkung 5.3 Der Fall einer unsymmetrischen Bilinearform. Im nicht-variationellen Fall, also wenn $b(\cdot, \cdot)$ unsymmetrisch ist, kann man: Finde $u \in V$ mit

$$b(u, v) = f(v) \quad \forall v \in V \tag{5.4}$$

auch mit dem Ritz'schen Verfahren approximieren. Die Eigenschaften von $b(\cdot, \cdot)$ seien Beschränktheit

$$|b(u, v)| \leq M \|u\|_V \|v\|_V, \quad M \in \mathbb{R},$$

und Koerzitivität

$$m \|v\|_V^2 \leq b(v, v), \quad m > 0.$$

Das diskrete Problem lautet: Finde $u_k \in V_k$, so dass

$$b(u_k, v_k) = f(v_k) \quad \forall v_k \in V_k. \quad (5.5)$$

Die diskrete Lösung existiert eindeutig nach Satz 4.12 (Lax–Milgram). Sie ist jedoch keine orthogonale Projektion in V_k mehr. Trotzdem kann man die gleiche Fehlerabschätzung wie im variationellen Fall beweisen. \square

Lemma 5.4 Lemma von Cea¹. *Sei die Bilinearform $b(\cdot, \cdot)$ beschränkt und koerzitiv. Dann gilt*

$$\|u - u_k\|_V \leq \frac{M}{m} \inf_{v_k \in V_k} \|u - v_k\|_V. \quad (5.6)$$

Beweis: Aus der Differenz der stetigen Gleichung (5.4) und der diskreten Gleichung (5.5)

$$b(u - u_k, v_k) = 0 \quad \forall v_k \in V_k$$

und

$$m \|v\|_V^2 \leq b(v, v) \quad \text{und} \quad |b(u, v)| \leq M \|u\|_V \|v\|_V$$

folgt sofort

$$\begin{aligned} \|u - u_k\|_V^2 &\leq \frac{1}{m} b(u - u_k, u - u_k) = \frac{1}{m} b(u - u_k, u - v_k) \\ &\leq \frac{M}{m} \|u - u_k\|_V \|u - v_k\|_V. \end{aligned}$$

■

Bemerkung 5.5 Galerkin–Methode. Im unsymmetrischen Fall wird das Ritzsche Verfahren Galerkin–Methode genannt. Das lineare Gleichungssystem wird genauso wie im symmetrischen Fall hergeleitet. Betrachte dazu das Zwei–Punkt–Randwertproblem

$$-\varepsilon u''(x) + b(x)u'(x) + c(x)u(x) = f(x), \quad \text{für } x \in (0, 1), \quad u(0) = u(1) = 0.$$

Die schwache Formulierung lautet: Finde $u \in H_0^1(0, 1)$, so dass für alle $v \in H_0^1(0, 1)$

$$\int_0^1 \left(\varepsilon u'(x)v'(x) + b(x)u'(x)v(x) + c(x)u(x)v(x) \right) dx = \int_0^1 f(x)v(x) dx$$

gilt. Falls (\cdot, \cdot) das Skalarprodukt in $L^2(0, 1)$ bezeichnet, kann die schwache Formulierung übersichtlicher geschrieben werden

$$b(u, v) := \varepsilon(u', v') + (bu', v) + (cu, v) = (f, v).$$

Sei $\{\phi_i\}_{i=1}^k$ eine beliebige Basis von V_k , dann macht man den Ansatz

$$u_k = \sum_{j=1}^k u^j \phi_j$$

mit unbekanntem Koeffizienten $\mathbf{u} = (u^1, \dots, u^k)^T$. Die variationelle Formulierung genau dann erfüllt, wenn sie für alle Basisfunktionen erfüllt ist. Man erhält

$$\sum_{j=1}^k \left[\varepsilon(\phi_j', \phi_i') + (b\phi_j', \phi_i) + (c\phi_j, \phi_i) \right] u^j = (f, \phi_i), \quad i = 1, \dots, k,$$

¹Cea

was äquivalent zu einem Gleichungssystem $\mathbf{A}u = \mathbf{b}$ ist. Die Einträge der Steifigkeitsmatrix sind

$$a_{ij} = \varepsilon(\phi'_j, \phi'_i) + (b\phi'_j, \phi_i) + (c\phi_j, \phi_i).$$

Die Systemmatrix ist nicht symmetrisch.

Die Eigenschaften der Bilinearform wurden im Beispiel 4.9 untersucht. Sind $b, c \in L^\infty(0, 1)$, so ist die Bilinearform beschränkt und die Konstante M ist in der Größenordnung von $\max\{\|b\|_\infty, \|c\|_\infty\}$. Gilt $-b'(x)/2 + c(x) \geq 0$, so ist sie koerzitiv mit $m = \varepsilon$. Falls beide Bedingungen erfüllt sind, dann ist das Lemma von Cea anwendbar und für den Fehler gilt

$$\|u - u_k\|_{H_0^1} \leq C \frac{\max\{\|b\|_\infty, \|c\|_\infty\}}{\varepsilon} \inf_{v_k \in V_k} \|u - v_k\|_{H_0^1}, \quad C \in \mathbb{R}.$$

Im singular gestörten Fall $\varepsilon \ll \|b\|_\infty$ ist der erste Faktor in dieser Fehlerabschätzung sehr groß. \square



Abbildung 5.2: Boris Grigorievich Galerkin (1871 – 1945).

Bemerkung 5.6 Parametrische Finite-Elemente. In der parametrischen Definition von Finite-Elementen geht man von einer Referenzgitterzelle \hat{K} mit einem lokalen Raum $P(\hat{K})$ und Funktionalen $\hat{\Phi}_1, \dots, \hat{\Phi}_N$ sowie einer Klasse regulärer Transformationen $\{F_K : \hat{K} \rightarrow K\}$ aus. Die Bilder $\{K\}$ bilden die Menge der zulässigen Gitterzellen. Die lokalen Räume sind durch

$$P(K) = \left\{ p : p = \hat{p} \circ F_K^{-1}, \hat{p} \in \hat{P}(\hat{K}) \right\}$$

definiert und die lokalen Funktionale durch

$$\Phi_{K,i}(v(\mathbf{x})) = \hat{\Phi}_i(v(F_K(\hat{\mathbf{x}}))),$$

wobei $\hat{\mathbf{x}} = (\hat{x}_1, \dots, \hat{x}_d)^T$ die Koordinaten der Referenzzelle bezeichnet. Es gilt $\mathbf{x} = F_K(\hat{\mathbf{x}})$. Als Funktionale werden Funktionswerte in gewissen Punkten, Ableitungen in gewissen Punkten oder Integrale über den Gitterzellen verwendet. \square

Definition 5.7 Finite-Elemente-Raum. Eine Funktion $v(\mathbf{x})$ definiert auf Ω mit $v|_{\text{int}(K)} \in P(K)$, $\text{int}(K) = K \setminus \partial K$, heißt stetig bezüglich Φ_i , falls

$$\Phi_i(v|_{K_1}) = \Phi_i(v|_{K_2})$$

für alle $K_1, K_2 \in \omega_i$, wobei ω_i die Vereinigung derjenigen Gitterzellen bezeichnet, auf denen Φ_i nicht verschwindet.

Der Raum

$$S = \left\{ v \in L^\infty(\Omega) : v|_{\text{int}(K)} \in P(K) \text{ und } v \text{ ist stetig bezüglich } \Phi_i, i = 1, \dots, N \right\}$$

heißt Finite-Element-Raum. \square

Bemerkung 5.8 Zum parametrischen Konzept. Neben der oben angegebenen Definition von Finite-Element-Räumen, kann man diese Räume auch direkt auf den Zellen des Gitter definieren. Falls die Referenzabbildungen affin sind (lineare Abbildung plus konstante Verschiebung), dann sind beide Definitionen oft äquivalent. Die Referenzabbildungen hängen nur von der Gestalt der Gitterzellen ab, aber nicht vom Finite-Element-Raum.

Das parametrische Konzept besitzt viele Vorteile bei der Implementation von Finite-Element-Methoden, da man alle benötigten Informationen (Basisfunktionen, Funktionale, Quadraturformeln) nur auf der Referenzzelle zu programmieren braucht. \square

Beispiel 5.9 Affines Konzept für P_1 in 1D. Man nimmt als Referenzgitterzelle beispielsweise $\hat{K} = [-1, 1]$. Die Referenzabbildung auf eine Gitterzelle $K = [x_i, x_{i+1}]$ wird so definiert, dass sie affin ist, das heißt es gilt

$$F_K(\hat{x}) = \alpha \hat{x} + \beta = x,$$

den Punkt -1 bildet man auf x_i sowie den Punkt 1 auf x_{i+1} ab. Das heißt

$$\begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \begin{pmatrix} x_i \\ x_{i+1} \end{pmatrix} \implies \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \frac{1}{2} \begin{pmatrix} x_{i+1} - x_i \\ x_{i+1} + x_i \end{pmatrix}.$$

Auf \hat{K} definiert man nun zwei lineare Basisfunktionen

$$\hat{\phi}_1(\hat{x}) = \frac{1}{2}(-\hat{x} + 1), \quad \hat{\phi}_2(\hat{x}) = \frac{1}{2}(\hat{x} + 1).$$

Die erste Funktion ist im Punkt -1 gleich Eins und verschwindet im Punkt 1 , bei der zweiten ist es genau umgekehrt.

Die Rücktransformation $F_K^{-1} : K \rightarrow \hat{K}$ von der Gitterzelle K auf die Referenzzelle \hat{K} besitzt die Gestalt

$$\hat{x} = \frac{x - \beta}{\alpha} = \frac{2x - (x_{i+1} + x_i)}{x_{i+1} - x_i}.$$

Die Basisfunktionen auf K sind definiert durch

$$\phi_i(x) := \hat{\phi}_i(F_K^{-1}(x)), \quad i = 1, 2.$$

Damit erhält man

$$\begin{aligned} \phi_1(x) &= \hat{\phi}_1\left(\frac{2x - (x_{i+1} + x_i)}{x_{i+1} - x_i}\right) = \frac{1}{2} \left(-\frac{2x - (x_{i+1} + x_i)}{x_{i+1} - x_i} + 1 \right) \\ &= -\frac{1}{2} \left(\frac{2x - 2x_{i+1}}{x_{i+1} - x_i} \right) = \frac{x_{i+1} - x}{x_{i+1} - x_i}, \\ \phi_2(x) &= \hat{\phi}_2\left(\frac{2x - (x_{i+1} + x_i)}{x_{i+1} - x_i}\right) = \frac{x - x_i}{x_{i+1} - x_i}. \end{aligned}$$

Das sind gerade die beiden Basisfunktionen, die man mit direkter Definition auf der Zelle K erhält. \square

Bemerkung 5.10 Assemblierung: Berechnung der Matrixeinträge und der rechten Seite. Die Matrixeinträge des Modellproblems besitzen die Form, siehe Bemerkung 5.5,

$$\begin{aligned} a_{ij} &= \int_0^1 \left(\varepsilon \phi_j'(x) \phi_i'(x) + b(x) \phi_j'(x) \phi_i(x) + c(x) \phi_j(x) \phi_i(x) \right) dx \\ &= \sum_{k=0}^{N-1} \int_{x_k}^{x_{k+1}} \left(\varepsilon \phi_j'(x) \phi_i'(x) + b(x) \phi_j'(x) \phi_i(x) + c(x) \phi_j(x) \phi_i(x) \right) dx. \end{aligned}$$

Das heißt, man kann die Integrale auf den einzelnen Gitterzellen berechnen und dann aufsummieren. Beim P_1 -Finite-Element hat man für $i = j$ Integrale auf genau zwei Gitterzellen zu berechnen, für $i = j \pm 1$ auf genau einer Gitterzelle und sonst sind alle Integrale Null. Für die Assemblierung der rechten Seite gilt eine analoge Formel.

Es bestehen die Möglichkeiten, die Integrale direkt auf einer Gitterzelle K zu berechnen oder das Integral auf die Referenzzelle \hat{K} zu transformieren. Der zweite Weg ist für die Implementation von Finite-Element-Methoden günstiger. Nach der Transformation auf \hat{K} kann man eine Quadraturformel anwenden, die für das Referenzelement implementiert ist. Man muss sich jedoch ansehen, wie sich die Terme in den Integralen transformieren.

Die Transformation der Integrale erfolgt natürlich mit der Substitutionsregel unter Verwendung der Referenzabbildung $F_K : [-1, 1] \rightarrow [x_k, x_{k+1}]$

$$\int_{x_k}^{x_{k+1}} f(x) dx = \int_{-1}^1 f(F_K(\hat{x})) F_K'(\hat{x}) d\hat{x}.$$

Es gilt, siehe Beispiel 5.9,

$$F_K(\hat{x}) = \frac{1}{2}((x_{k+1} - x_k)\hat{x} + (x_{k+1} + x_k)) = x \implies F_K'(\hat{x}) = \frac{x_{k+1} - x_k}{2}.$$

Das ist die halbe Länge der Gitterzelle $[x_k, x_{k+1}]$, woraus $F_K'(\hat{x}) > 0$ folgt. Für die Basisfunktion ist die Transformation auf die Referenzzelle durch

$$\phi_i(x) = \hat{\phi}_i(F_K^{-1}(x)) = \hat{\phi}_i(\hat{x})$$

gegeben. Zur Transformation der Ableitung verwendet man die Kettenregel

$$\phi_i'(x) = \frac{d\phi_i(x)}{dx} = \frac{d\hat{\phi}_i(\hat{x})}{d\hat{x}} = \frac{d\hat{\phi}_i(\hat{x})}{d\hat{x}} \frac{d\hat{x}}{dx} = \frac{d\hat{\phi}_i(\hat{x})}{d\hat{x}} \frac{2}{x_{k+1} - x_k}.$$

Die Ableitungen der Basisfunktionen auf der Referenzzelle kann man vorher explizit ausrechnen und dann implementieren. Damit erhält man

$$\begin{aligned} \int_{x_k}^{x_{k+1}} \varepsilon \phi_j'(x) \phi_i'(x) dx &= \frac{2}{x_{k+1} - x_k} \int_{-1}^1 \varepsilon \frac{d\hat{\phi}_j(\hat{x})}{d\hat{x}} \frac{d\hat{\phi}_i(\hat{x})}{d\hat{x}} d\hat{x}, \\ \int_{x_k}^{x_{k+1}} b(x) \phi_j'(x) \phi_i(x) dx &= \int_{-1}^1 b(F_K(\hat{x})) \frac{d\hat{\phi}_j(\hat{x})}{d\hat{x}} \hat{\phi}_i(\hat{x}) d\hat{x}, \\ \int_{x_k}^{x_{k+1}} c(x) \phi_j(x) \phi_i(x) dx &= \frac{x_{k+1} - x_k}{2} \int_{-1}^1 c(F_K(\hat{x})) \hat{\phi}_j(\hat{x}) \hat{\phi}_i(\hat{x}) d\hat{x}, \\ \int_{x_k}^{x_{k+1}} f(x) \phi_i(x) dx &= \frac{x_{k+1} - x_k}{2} \int_{-1}^1 f(F_K(\hat{x})) \hat{\phi}_i(\hat{x}) d\hat{x}. \end{aligned}$$

Nun kann man hinreichend genaue Quadraturformeln auf der Referenzzelle zur Approximation der Integrale verwenden. Eine genaue Quadratur ist vor allem für finite Elemente höherer Ordnung wesentlich, damit die Genauigkeit nicht durch Quadraturfehler beeinträchtigt wird. \square

Beispiel 5.11 Assemblierung: P_1 in 1D. Betrachte den Fall, dass die Parameterfunktionen konstant sind, $b(x) = b$, $c(x) = c$ und $f(x) = f$. Für das P_1 -Finite-Element auf der Referenzzelle $[-1, 1]$ gelten

$$\hat{\phi}_1(\hat{x}) = \frac{1}{2}(-\hat{x} + 1), \quad \hat{\phi}'_1(\hat{x}) = -\frac{1}{2}, \quad \hat{\phi}_2(\hat{x}) = \frac{1}{2}(\hat{x} + 1), \quad \hat{\phi}'_2(\hat{x}) = \frac{1}{2}.$$

Betrachte nun den Matrixeintrag $a_{i,i+1}$, der mit Hilfe der Testfunktion $\phi_i(x)$, die auf $\hat{\phi}_1(\hat{x})$ transformiert wird, und der Ansatzfunktion $\phi_{i+1}(x)$, die auf $\hat{\phi}_2(\hat{x})$ transformiert wird, berechnet wird. Der gemeinsame Träger ist die Gitterzelle $[x_i, x_{i+1}]$. Bezeichne h die Länge dieser Zelle. Dann folgt

$$\begin{aligned} a_{i,i+1} &= \frac{2\varepsilon}{h} \int_{-1}^1 \frac{1}{2} \cdot \left(-\frac{1}{2}\right) d\hat{x} + b \int_{-1}^1 \frac{1}{2} \cdot \frac{1}{2}(-\hat{x} + 1) d\hat{x} \\ &\quad + \frac{ch}{2} \int_{-1}^1 \frac{1}{2}(\hat{x} + 1) \frac{1}{2}(-\hat{x} + 1) d\hat{x} \\ &= -\frac{\varepsilon}{h} + \frac{b}{2} + \frac{ch}{6}. \end{aligned}$$

Für die i -te Komponente der rechten Seite, erhält, man

$$f_i = f \int_0^1 \phi_i(x) dx = hf,$$

da die Fläche unter einer Hütchenfunktion das Maß h besitzt. *andere Einträge als Übungsaufgabe*

Verwendet man zur Approximation der Integrale die Trapezregel, so ergibt sich

$$\frac{ch}{2} \int_{-1}^1 \frac{1}{2}(-\hat{x} + 1) \frac{1}{2}(\hat{x} + 1) d\hat{x} = \frac{ch}{2} 2(0 + 0) = 0.$$

In diesem Fall ist

$$a_{i,i+1} = -\frac{\varepsilon}{h} + \frac{b}{2}.$$

Für $c = 0$ (oder mit Trapezregel) sind das, bis auf den Faktor h , die gleichen Einträge wie beim zentralen Differenzenverfahren, siehe Bemerkung 3.10. Man erhält für $c = 0$

$$-h\varepsilon D^+ D^- u_i + hb_i D^0 u_i = hf_i \quad \iff \quad -\varepsilon D^+ D^- u_i + b_i D^0 u_i = f_i.$$

Dieser Zusammenhang zwischen Finite-Differenzen-Methoden und Finite-Element-Methoden gilt im allgemeinen nicht mehr, wenn die Koeffizientenfunktionen nicht konstant sind. In höheren Dimensionen unterscheiden sich FDM und FEM im allgemeinen auch bei konstanten Koeffizienten. \square

Bemerkung 5.12 Interpolierende. Eine Interpolierende ist eine (vernünftige) Approximation einer Funktion aus einem Sobolev-Raum durch eine Funktion aus dem Finite-Elemente-Raum. Die Frage, wie gut man Funktionen aus einem Sobolev-Raum durch Finite-Elemente-Funktionen approximieren kann ist von zentraler Bedeutung, siehe Eigenschaft 2. der Ritzschen Approximation, Lemma 5.2, und die Fehlerabschätzung (5.6) des Lemmas von Cea.

Die analytische Formulierung auf einer Referenzgitterzelle \hat{K} ist wie folgt. Seien $\hat{K} \subset \mathbb{R}$, zum Beispiel $\hat{K} = [-1, 1]$, $\hat{P}(\hat{K})$ ein Polynomraum der Dimension N , $\hat{\Phi}_1, \dots, \hat{\Phi}_N : C^s(\hat{K}) \rightarrow \mathbb{R}$ stetige lineare Funktionale und $\hat{\phi}_1(\hat{x}), \dots, \hat{\phi}_N(\hat{x}) \in \hat{P}(\hat{K})$ eine lokale Basis. Das heißt, $\{\hat{\phi}_i(\hat{x})\}_{i=1}^N$ ist eine Basis von $\hat{P}(\hat{K})$ und es gilt

$$\hat{\Phi}_i(\hat{\phi}_j) = \delta_{ij}, \quad i, j = 1, \dots, N.$$

Für $\hat{v} \in C^s(\hat{K})$ wird die Interpolierende $(I_{\hat{K}}\hat{v})(\hat{x})$ durch

$$I_{\hat{K}}\hat{v}(\hat{x}) = \sum_{i=1}^N \hat{\Phi}_i(\hat{v}) \hat{\phi}_i(\hat{x})$$

definiert. Der Operator $I_{\hat{K}}$ ist ein stetiger und linearer Operator von $C^s(\hat{K})$ nach $\hat{P}(\hat{K})$. Aus der Linearität folgt, dass $I_{\hat{K}}$ die Identität auf $\hat{P}(\hat{K})$ ist *Übungsaufgabe*

$$(I_{\hat{K}}\hat{p})(\hat{x}) = \hat{p}(\hat{x}) \quad \forall \hat{p} \in \hat{P}(\hat{K}).$$

□

Satz 5.13 Interpolationsabschätzung für eine beliebige Gitterzelle. *Seien eine Referenzgitterzelle \hat{K} , Funktionale $\{\hat{\Phi}_i\}$ und ein Polynomraum $\hat{P}(\hat{K})$ gegeben. Seien weiter $p \in [1, \infty]$ und $(m+1-s)p > 1$. Dann gibt es eine Konstante c unabhängig von $v \in W^{m+1,p}(K)$ mit*

$$\left\| (v - I_K v)^{(k)} \right\|_{L^p(K)} \leq c h_K^{m+1-k} \left\| v^{(m+1)} \right\|_{L^p(K)}, \quad 0 \leq k \leq m+1, \quad (5.7)$$

für alle $v \in W^{m+1,p}(K)$. Man beachte, dass die Potenz von h_K unabhängig von p ist.

Beweis: Der Beweis sollte bereits in der Einführungsvorlesung über Finite-Element-Methoden gegeben worden sein. Eine Darstellung dieser Thematik ist auch im Anhang zu finden. ■

Bemerkung 5.14 CSR-Speicherschema von schwach besetzten Matrizen.

Von schwach besetzten Matrizen speichert man natürlich nur die Einträge, die nicht Null sind und zugehörige Informationen über die Position des Eintrags. Die am weitesten verbreiteste Herangehensweise ist das CSR-Speicherschema (condensed sparse row). Bei diesem Schema werden die Nichtnulleinträge zeilenweise abgespeichert. Innerhalb einer Zeile brauchen sie nicht bezüglich der Spaltenindizes angeordnet zu werden.

Sei eine schwach besetzte Matrix $A \in \mathbb{R}^{m \times n}$ mit nnz Nichtnullelementen zu speichern. Dann braucht man drei Arrays:

- `double-Array entries` der Länge nnz , darin werden die Einträge von A zeilenweise gespeichert,
- `int-Array col_ptr` der Länge nnz , darin stehen die Spaltenindizes der zugehörigen Einträge von `entries`.
- `int-Array row_ptr` der Länge $m+1$, darin wird abgespeichert, an welcher Stelle im Array `entries` die i -te Zeile beginnt, $i = 1, \dots, m$. Der letzte Eintrag von `row_ptr` verweist auf den ersten Speicherplatz nach dem Ende des Arrays `entries`.

□

Beispiel 5.15 Die Matrix

$$A = \begin{pmatrix} 1 & 0 & 0 & 2 & 0 \\ 3 & 4 & 0 & 5 & 0 \\ 6 & 0 & 7 & 8 & 9 \\ 0 & 0 & 10 & 11 & 0 \\ 0 & 0 & 0 & 0 & 12 \end{pmatrix}$$

kann wie folgt gespeichert werden (Numerierung beginnt bei 0):

```

entries  -  1  2  3  4  5  6  7  8  9  10 11 12
col_ptr  -  0  3  0  1  3  0  2  3  4  2  3  4 .
row_ptr  -  0  2  5  9 11 12

```

Eine andere Möglichkeit ist

```

entries  -  2  1  4  5  3  7  9  8  6 11 10 12
col_ptr  -  3  0  1  3  0  2  4  3  0  3  2  4 .
row_ptr  -  0  2  5  9 11 12

```

□

5.2 Stabilisierte Finite-Element-Methoden

Bemerkung 5.16 Zum Lemma von Lax-Milgram für singulär gestörte Probleme. Betrachte das Modellproblem: Finde $u \in V = H_0^1(0,1)$ so dass

$$a(u, v) = f(v) \quad \forall v \in V \quad (5.8)$$

mit

$$a(u, v) := \int_0^1 \left(\varepsilon u'(x)v'(x) + b(x)u'(x)v(x) + c(x)u(x)v(x) \right) dx,$$

$$f(v) := \int_0^1 f(x)v(x) dx.$$

Sei

$$c(x) - \frac{b'(x)}{2} \geq \omega > 0 \quad \text{für alle } x \in [0,1].$$

Mit einer analogen Rechnung wie in Bemerkung 4.9 zeigt man, dass $a(\cdot, \cdot)$ koerzitiv bezüglich der von ε abhängigen Norm

$$\|v\|_\varepsilon^2 := \varepsilon |v|_{1,2}^2 + \|v\|_0^2 = \varepsilon \|v'\|_{L^2(0,1)}^2 + \|v\|_{L^2(0,1)}^2$$

ist. Das heißt, es existiert eine von ε unabhängige Konstante μ , so dass

$$a(v, v) \geq \mu \|v\|_\varepsilon^2 \quad \forall v \in V$$

gilt. Mit partieller Integration (*Übungsaufgabe*) zeigt man, dass es eine von ε unabhängige Konstante β gibt, so dass

$$|a(v, w)| \leq \beta \|v\|_\varepsilon \|w\|_{H^1} \quad \forall (v, w) \in V \times V.$$

Es gibt jedoch keine von ε unabhängige Konstante γ mit

$$|a(v, w)| \leq \gamma \|v\|_\varepsilon \|w\|_\varepsilon \quad \forall (v, w) \in V \times V.$$

Nutzt man die Abschätzungen mit Konstanten die unabhängig von ε sind, erhält man analog zum Beweis des Lemmas von Cea

$$\|u - u_h\|_\varepsilon \leq C \inf_{v_h \in V_h} \|u - v_h\|_{H^1}$$

mit C unabhängig von ε . Ist V^h ein Standard-Finite-Elemente-Raum (stückweise polynomial), dann kann man zeigen, dass in Grenzsichten

$$\inf_{v_h \in V_h} \|u - v_h\|_{H^1} \rightarrow \infty \quad \text{für } \varepsilon \rightarrow 0$$

für festes h gilt. Deswegen hat man keine gleichmäßige Konvergenz $\|u - u_h\|_\varepsilon \rightarrow 0$ für $h \rightarrow 0$. Die $H^1(0,1)$ -Norm ist zur Untersuchung konvektions-dominanter Probleme nicht geeignet. □

5.2.1 Petrov–Galerkin–Methoden und Upwind–Verfahren

Bemerkung 5.17 Petrov²–Galerkin–Methode. Eine Finite–Element–Methode, bei welcher Ansatz– und Testraum unterschiedlich sind, wird Petrov–Galerkin–Methode genannt. Seien S_h der Ansatzraum und T_h der Testraum, mit $\dim(S_h) = \dim(T_h)$, dann lautet eine Petrov–Galerkin–Methode: Finde $u_h \in S_h$, so dass

$$a(u_h, v_h) = f(v_h) \quad \forall v_h \in T_h.$$

□

Beispiel 5.18 Petrov–Galerkin–Methode und Upwind–Verfahren. Betrachte

$$-\varepsilon u''(x) + bu'(x) = 0$$

mit $b \in \mathbb{R} \setminus \{0\}$. Nutze als Ansatzfunktionen stückweise lineare Funktionen

$$\phi_i(x) = \begin{cases} (x - x_{i-1})/h & \text{für } x \in [x_{i-1}, x_i], \\ (x_{i+1} - x)/h & \text{für } x \in [x_i, x_{i+1}], \\ 0 & \text{sonst,} \end{cases} \quad i = 1, \dots, N-1.$$

Definiere die Blasenfunktion

$$\sigma_{i-1/2}(x) = \begin{cases} 4(x - x_{i-1})(x_i - x)/h^2 & \text{für } x \in [x_{i-1}, x_i], \\ 0 & \text{sonst.} \end{cases}$$

Die Testfunktionen werden nun als stückweise quadratische Funktionen definiert

$$\psi_i(x) = \phi_i(x) + \frac{3}{2}\kappa (\sigma_{i-1/2}(x) - \sigma_{i+1/2}(x)), \quad i = 1, \dots, N-1,$$

wobei κ ein zu wählender Upwind–Parameter ist. Direktes Nachrechnen (*Übungsaufgabe*) zeigt, dass man damit das folgende Schema erhält

$$-\varepsilon D^+ D^- u_i + b \left[\left(\frac{1}{2} - \kappa \right) D^+ u_i + \left(\frac{1}{2} + \kappa \right) D^- u_i \right] = 0.$$

Wählt man $\kappa = \text{sgn}(b)/2$, so erhält man das einfache Upwind–Finite–Differenzen–Verfahren, siehe Definition 3.31.

Eine Testfunktion $\psi_i(x)$, definiert in den Knoten $\{0, 0.5, 1\}$ für $\kappa = 1/2$ ist in Abbildung 5.3 dargestellt.

□

Bemerkung 5.19 Man kann auch einige angepasste Upwind–Verfahren mit Hilfe von Petrov–Galerkin–Methoden gewinnen. Das funktioniert auch für nichtkonstante Koeffizientenfunktionen. Bessere Ergebnisse als etwa beim Iljin–Allen–Southwell–Verfahren, Satz 3.54, das heißt lineare Konvergenz, sind aber nicht zu erreichen.

□

5.2.2 Die Stromlinien–Diffusions–Finite–Elemente–Methode

Bemerkung 5.20 Ziele. Die Ziele bestehen darin, ein Verfahren zu konstruieren, welches stabiler als das Galerkin–Verfahren ist und welches für Finite–Elemente beliebiger Ordnung genutzt werden kann. Die Konvergenz dieses Verfahrens soll zudem in einer geeigneten Norm von höherer Ordnung sein.

Es wird das Modellproblem

$$-\varepsilon u''(x) + b(x)u'(x) + c(x)u(x) = f(x) \text{ in } (0, 1), \quad u(0) = u(1) = 0, \quad (5.9)$$

²Petrov

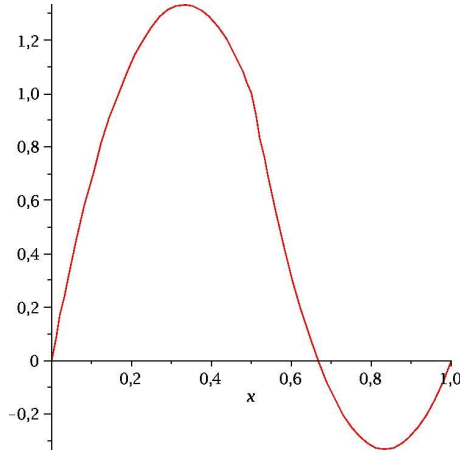


Abbildung 5.3: Stückweise quadratische Testfunktion für $\kappa = 1/2$.

unter der Bedingung

$$c(x) - \frac{b'(x)}{2} \geq \omega > 0 \quad \text{fast überall } x \in (0, 1).$$

betrachtet. □

Bemerkung 5.21 Idee. Eine Idee zur Konstruktion eines stabileren Verfahrens besteht darin, zur Galerkin-Finite-Element-Methode gewichtete Residuen der starken Formulierung der Differentialgleichung (5.9) zu addieren. Dazu wird (5.9) mit $(bv')(x)$ multipliziert, über jedes Teilintervall (x_{i-1}, x_i) , $i = 1, \dots, N$, mit einem Gewicht versehen und integriert, und dann zur Galerkin-Methode addiert. Die Abweichung des Residuums der starken Form der Gleichung von Null wird bestraft. □

Definition 5.22 Stromlinien-Diffusions-Finite-Elemente-Methode, SDFEM, Stromlinien-Upwind-Petrov-Galerkin FEM, SUPG. Die Stromlinien-Diffusions-Finite-Elemente-Methode (SDFEM) oder Stromlinien-Upwind-Petrov-Galerkin (SUPG) FEM ist wie folgt definiert: Finde $u_h \in V_h$, so dass

$$a_h(u_h, v_h) = f_h(v_h) \quad \forall v_h \in V_h \tag{5.10}$$

gilt, mit

$$\begin{aligned} a_h(v, w) &:= \varepsilon(v', w') + (bv' + cv, w) \\ &\quad + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \delta_i \left(-\varepsilon v''(x) + b(x)v'(x) + c(x)v(x) \right) \left(b(x)w'(x) \right) dx, \\ f_h(w) &:= (f, w) + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \delta_i f(x) \left(b(x)w'(x) \right) dx. \end{aligned} \tag{5.11}$$

Dabei sind $\{\delta_i\}_{i=1}^N$ geeignet zu wählende Gewichte, welche SD-Parameter genannt werden. □

Bemerkung 5.23 Zur SDFEM.

- Der Name Stromlinien–Diffusion–FEM wird erst in höheren Dimensionen klar. Dort wird als Testfunktion für das Residuum die Ableitung in Konvektionsrichtung gewählt, $\mathbf{b} \cdot \nabla w$. Das ist die sogenannte Stromlinienrichtung. Man wird beispielsweise in Beispiel 5.24 erkennen, dass diese Diskretisierung zusätzliche Diffusion einführt. In höheren Dimensionen wirkt diese Diffusion nur in Richtung der Stromlinien.
- Der Name SUPG kommt daher, dass man diese Methode als Petrov–Galerkin–Methode mit den Testraum

$$\text{span} \left\{ w(x) + \sum_{i=1}^N \delta_i b(x) w'(x) \right\}$$

betrachten kann.

- Die Methode wurde von Brooks und Hughes in [HB79, BH82] entwickelt.
- Für die Lösung der Galerkin–FEM wird sich das Residuum der starken Formulierung der Gleichung im allgemeinen stark von Null unterscheiden. Bei der SDFEM wird verlangt, dass dieses Residuum (in einem schwachen Sinne) nicht zu groß sein darf. Das Gewicht dieses Residuum in der SDFEM wird durch die SD–Parameter bestimmt.
- Für eine Finite–Elemente–Funktion ist die zweite Ableitung im allgemeinen nur stückweise definiert, und zwar innerhalb der Gitterzellen.
- Die SD–Parameter werden oft in jedem Intervall (x_{i-1}, x_i) konstant gewählt. Das Ziel der Analysis besteht darin, eine möglichst günstige Wahl dieser Parameter aufzuzeigen.

□

Beispiel 5.24 SDFEM für P_1 . Betrachte $V_h = P_1$ auf einem äquidistanten Gitter mit $h_i = h$, $i = 1, \dots, N$. Sind alle Koeffizientenfunktionen konstant, $c = 0$, und wählt man die SD–Parameter auch konstant, so reduziert sich die linke Seite der SDFEM zu

$$\begin{aligned} \varepsilon(u'_h, v'_h) + (bu'_h, v_h) + \sum_{i=1}^N \delta \int_{x_{i-1}}^{x_i} \left(-\varepsilon \cdot 0 + bu'_h(x) \right) (bv'_h(x)) dx \\ = \varepsilon(u'_h, v'_h) + b(u'_h, v_h) + \delta b^2(u'_h, v'_h). \end{aligned}$$

Das entspricht der Galerkin FEM einer Gleichung mit linker Seite

$$-(\varepsilon + \delta b^2) u''(x) + bu'(x).$$

Aus Beispiel 5.11 ist bekannt, dass die Galerkin FEM wiederum einem zentralen Differenzenverfahren entspricht. Die rechte Seite ist

$$(f, v_h) + \sum_{i=1}^N \delta \int_{x_{i-1}}^{x_i} f b v'_h(x) dx = (f, v_h) + \underbrace{\delta f b \sum_{i=1}^N \int_{x_{i-1}}^{x_i} v'_h(x) dx}_{=0} = (f, v_h).$$

Die Summe verschwindet, da sich jede Testfunktion $v'_h(x)$ als Linearkombination der Basisfunktionen $\{\phi_i(x)\}$ von P_1 schreiben lässt und das Integral über die Ableitung jeder Basisfunktion gleich Null ist. Oder man verwendet partielle Integration um dies zu sehen.

Insgesamt entspricht die SDFEM unter den obigen Voraussetzungen dem angepassten Finite–Differenzen–Upwind–Verfahren (3.7)

$$-\varepsilon \left(1 + \delta \frac{b^2}{\varepsilon} \right) D^+ D^- u_i + b D^0 u_i = f_i,$$

das heißt $\sigma(q) = 1 + \delta b^2/\varepsilon$, $q = bh/(2\varepsilon)$. Wählt man den SD-Parameter als

$$\delta(q) = \frac{h}{2b} \left(\coth(q) - \frac{1}{q} \right),$$

so ist

$$\sigma(q) = 1 + \frac{hb^2}{2b\varepsilon} \left(\coth(q) - \frac{1}{q} \right) = 1 + q \left(\coth(q) - \frac{1}{q} \right) = q \coth(q).$$

Damit erhält man das Iljin-Allen-Southwell-Verfahren.

Mit $\delta = h/(2b)$ erhält man das einfache Upwind-Verfahren.

Achtung: diese einfachen Zusammenhänge gelten in höheren Dimensionen nicht mehr! \square

Definition 5.25 Konsistente Finite-Element-Methode. Sei $u(x)$ eine hinreichend glatte Lösung von (5.9). Eine Finite-Element-Methode: Finde $u_h \in V_h$, so dass

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

wird konsistent genannt, wenn gilt

$$a_h(u, v_h) = f_h(v_h) \quad \forall v_h \in V_h. \quad (5.12)$$

\square

Bemerkung 5.26 Zum Begriff Konsistenz. Konsistenz einer Finite-Element-Methode ist nicht das gleiche wie Konsistenz einer Finiten-Differenzen-Methode, siehe Definition 3.6. Für Finite-Element-Methoden bedeutet Konsistenz, dass eine hinreichend glatte Lösung auch die diskrete Gleichung erfüllt. \square

Lemma 5.27 Galerkin-Orthogonalität. Eine konsistente Finite-Element-Methode besitzt die Eigenschaft der Galerkin-Orthogonalität

$$a_h(u - u_h, v_h) = 0 \quad \forall v_h \in V_h. \quad (5.13)$$

Man sagt auch, dass der Fehler „senkrecht“ auf dem Finite-Element-Raum steht.

Beweis: Die Aussage folgt sofort aus der Gültigkeit von (5.10) und (5.12) durch Subtraktion dieser beiden Gleichungen. \blacksquare

Lemma 5.28 Konsistenz der SDFEM. Die SDFEM (5.10) – (5.11) ist konsistent.

Beweis: Für eine hinreichend glatte Lösung $u(x)$ von (5.9) ist das Residuum der starken Form der Gleichung gleich Null. Damit verschwinden die SDFEM-Terme in (5.11). Durch partielle Integration erhält man aus den übrigen Termen, dass

$$\int_0^1 \left(-\varepsilon u''(x) + b(x)u'(x) + c(x)u(x) - f(x) \right) v_h(x) dx = 0 \quad \forall v_h \in V_h$$

gilt. Für eine hinreichend glatte Lösung verschwindet der Ausdruck in der Klammer und diese Aussage ist wahr. Damit ist die SDFEM konsistent. \blacksquare

Bei der Analysis stabilisierter FEM ist es wichtig, dass man geeignete Normen verwendet.

Definition 5.29 Stromlinien–Diffusions–Norm, SD–Norm. Auf V_h wird die Stromlinien–Diffusions–Norm

$$\|v_h\|_{SD} := \left(\varepsilon |v_h|_1^2 + \omega \|v_h\|_0^2 + \sum_{i=1}^N \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i}^2 \right)^{1/2}$$

definiert. Hierbei ist $I_i := (x_{i-1}, x_i)$ und $\|\cdot\|_{0,I_i}$ ist die Norm in $L^2(I_i)$. \square

Satz 5.30 Koerzitivität der SD–Bilinearform. Sei

$$0 < \delta_i \leq \frac{1}{2} \min \left\{ \frac{h_i^2}{\varepsilon c_{\text{inv}}^2}, \frac{\omega}{\|c\|_{L^\infty(I_i)}^2} \right\}, \quad (5.14)$$

wobei c_{inv} die Konstante der inversen Ungleichung

$$\|v_h''\|_{0,I_i} \leq c_{\text{inv}} h_i^{-1} \|v_h'\|_{0,I_i} \quad (5.15)$$

ist. Dann ist die SD–Bilinearform (5.11) koerzitiv bezüglich der SD–Norm, das heißt es gilt

$$a_h(v_h, v_h) \geq \frac{1}{2} \|v_h\|_{SD}^2 \quad \forall v_h \in V_h.$$

Beweis: Mit partieller Integration folgt, siehe Beispiel 4.9,

$$(b v_h' + c v_h, v_h) = \left(\left(-\frac{b'}{2} + c \right) v_h, v_h \right) \quad \forall v_h \in V_h.$$

Mit der Definition von ω ergibt sich

$$\begin{aligned} a_h(v_h, v_h) &= \varepsilon |v_h|_1^2 + \underbrace{\int_0^1 \left(c(x) - \frac{b'(x)}{2} \right) v_h^2(x) dx}_{\geq \omega > 0} + \sum_{i=1}^N \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i}^2 \\ &\quad + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \delta_i (-\varepsilon v_h''(x) + c(x) v_h(x)) (b(x) v_h'(x)) dx \\ &\geq \|v_h\|_{SD}^2 + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \delta_i (-\varepsilon v_h''(x) + c(x) v_h(x)) (b(x) v_h'(x)) dx. \end{aligned}$$

Nun wird der zweite Term nach oben abgeschätzt, womit man insgesamt eine Abschätzung nach unten erhält, wenn man die Abschätzung des zweiten Terms vom ersten Term sub-

trahiert. In der Abschätzung wird die Definition des SD-Parameters verwendet. Es ist

$$\begin{aligned}
& \left| \int_{x_{i-1}}^{x_i} \delta_i (-\varepsilon v_h''(x) + c(x)v_h(x)) (b(x)v_h'(x)) dx \right| \\
& \leq \int_{x_{i-1}}^{x_i} \left(\delta_i^{1/2} \varepsilon |v_h''(x)| \right) \left(\delta_i^{1/2} |b(x)v_h'(x)| \right) dx \\
& \quad + \int_{x_{i-1}}^{x_i} \left(\delta_i^{1/2} |c(x)| |v_h(x)| \right) \left(\delta_i^{1/2} |b(x)v_h'(x)| \right) dx \\
& \stackrel{\text{CSU}}{\leq} \left(\delta_i^{1/2} \varepsilon \|v_h''\|_{0,I_i} + \delta_i^{1/2} \|c\|_{L^\infty(I_i)} \|v_h\|_{0,I_i} \right) \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i} \\
& \stackrel{(5.15)}{\leq} \left(\delta_i^{1/2} \frac{\varepsilon c_{\text{inv}}}{h_i} \|v_h'\|_{0,I_i} + \delta_i^{1/2} \|c\|_{L^\infty(I_i)} \|v_h\|_{0,I_i} \right) \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i} \\
& \stackrel{(5.14)}{\leq} \left(\frac{h_i}{\sqrt{2\varepsilon c_{\text{inv}}}} \frac{\varepsilon c_{\text{inv}}}{h_i} \|v_h'\|_{0,I_i} + \frac{\sqrt{\omega}}{\sqrt{2}\|c\|_{L^\infty(I_i)}} \|c\|_{L^\infty(I_i)} \|v_h\|_{0,I_i} \right) \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i} \\
& = \left(\sqrt{\frac{\varepsilon}{2}} \|v_h'\|_{0,I_i} + \sqrt{\frac{\omega}{2}} \|v_h\|_{0,I_i} \right) \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i} \\
& \stackrel{\text{Young Ugl.}}{\leq} \frac{\varepsilon}{2} \|v_h'\|_{0,I_i}^2 + \frac{1}{4} \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i}^2 + \frac{\omega}{2} \|v_h\|_{0,I_i}^2 + \frac{1}{4} \left\| \sqrt{\delta_i} b v_h' \right\|_{0,I_i}^2 \\
& = \frac{1}{2} \|v_h\|_{SD,I_i}^2.
\end{aligned}$$

Summation über alle Gitterzellen und Einsetzen in die erste Abschätzung ergibt die Aussage des Satzes ■

Folgerung 5.31 Koerzitivität der SD-Bilinearform für lineare finite Elemente. Für stückweise lineare finite Elemente ist die SD-Bilinearform (5.11) koerzitiv bezüglich der SD-Norm mit der Parameterwahl

$$0 < \delta_i \leq \frac{\omega}{\|c\|_{L^\infty(I_i)}^2}. \quad (5.16)$$

Beweis: Der Beweis ist wie für Satz 5.30, wobei man ausnutzt, dass für stückweise lineare finite Elemente $v_h''(x) = 0$ in I_i , $i = 1, \dots, N$, ist und die entsprechenden Terme im Beweis entfallen. ■

Bemerkung 5.32 Zur Koerzitivität der SD-Bilinearform.

- Der Beweis von Satz 5.30 ist typisch für die Untersuchung stabilisierter Finite-Element-Methoden. Man versucht die störenden Terme irgendwie mit der verwendeten Norm abzuschätzen. Das geht im allgemeinen nur, wenn man eine geeignete Norm verwendet. Insbesondere muss die Stabilisierung in dieser Norm irgendwie auftauchen.
- Satz 5.30 gibt eine hinreichende obere Schranke für die SD-Parameter.
- Aus Satz 5.30 folgt die Stabilität der SDFEM bezüglich der SD-Norm. Stabilität bedeutet, dass eine geeignete Norm der Lösung durch geeignete Normen der Daten abgeschätzt werden kann. Es gilt

$$\begin{aligned}
\|u_h\|_{SD}^2 & \leq 2a_h(u_h, u_h) = f_h(u_h) \\
& = (f, u_h) + \sum_{i=1}^N \int_{x_{i-1}}^{x_i} \delta_i f(x) (b(x)u_h'(x)) dx \\
& \stackrel{\text{CSU}}{\leq} \frac{1}{\sqrt{\omega}} \|f\|_0 \sqrt{\omega} \|u_h\|_0 + \sum_{i=1}^N \sqrt{\delta_i} \|f\|_{0,I_i} \left\| \sqrt{\delta_i} b u_h' \right\|_{0,I_i} \\
& \stackrel{\text{Young}}{\leq} C \|f\|_0^2 + \frac{1}{2} \left(\omega \|u_h\|_0^2 + \sum_{i=1}^N \left\| \sqrt{\delta_i} b u_h' \right\|_{0,I_i}^2 \right).
\end{aligned}$$

Daraus ergibt sich die Stabilität. Die Konstante hängt von ω und von der oberen Schranke für δ_i ab.

- Alle $v_h \in V_h$ erfüllen

$$\|v_h\|_{SD} \geq \min\{1, \omega\} \|v_h\|_\varepsilon.$$

Damit folgt, dass die SDFEM auch bezüglich der Norm $\|\cdot\|_\varepsilon$ stabil ist. Bezüglich $\|\cdot\|_\varepsilon$ ist auch die Galerkin-FEM stabil, jedoch nicht bezüglich $\|\cdot\|_{SD}$. Damit ist die Stabilitätsaussage von Satz 5.30 stärker als die Stabilitätsaussage für die Galerkin-FEM. □

Beispiel 5.33 Fortsetzung: SDFEM für P_1 . Im Beispiel 5.24 wurde gezeigt, dass man unter gewissen Bedingungen mit

$$\delta(q) = \frac{h}{2b} \left(\coth(q) - \frac{1}{q} \right), \quad q = \frac{bh}{2\varepsilon},$$

das Iljin-Allen-Southwell-Verfahren, also ein gleichmäßig konvergentes Verfahren, erhält. Man braucht aber auch Parameter im Falle von nichtkonstanten Koeffizientenfunktionen, Finite-Elementen höherer Ordnung und für Probleme in höheren Dimensionen. Dabei kann man versuchen, den Spezialfall zu verallgemeinern. Mit Taylor-Entwicklung, Übungsaufgabe, sieht man, dass

$$\begin{aligned} \coth q - \frac{1}{q} &= \frac{q}{3} + \mathcal{O}(q^3) \quad \text{für } q \rightarrow 0, \\ \coth q - \frac{1}{q} &= 1 + \mathcal{O}\left(\frac{1}{q}\right) \quad \text{für } q \rightarrow \infty. \end{aligned}$$

Ist ε konstant und geht $h \rightarrow 0$, so folgt $q \rightarrow 0$ und es ist $\delta(q) \approx hq/(6b)$. Für festes h und $\varepsilon \rightarrow 0$ folgt $q \rightarrow \infty$ und es ist $\delta(q) \approx h/(2b)$. Damit ist die folgende Wahl des SD-Parameters motiviert

$$\delta(q) = \begin{cases} \frac{h^2}{12\varepsilon} & \text{für } 0 < q \ll 1, \\ \frac{h}{2b} & \text{für } q \gg 1. \end{cases}$$

Falls das Gitter sehr grob im Vergleich zu ε ist, also $q \gg 1$, dann geht die SDFEM in 1D in das einfache Upwind-Verfahren über. □

Satz 5.34 Konvergenz des SDFEM. Gelte für die Lösung von (5.9) $u \in H^{k+1}(0, 1)$, für die Koeffizientenfunktionen $b \in W^{1,\infty}(0, 1)$, $c \in L^\infty(0, 1)$ und betrachte die SDFEM mit P_k -Finite-Elementen. Die SD-Parameter seien wie folgt gegeben

$$\delta_i = \begin{cases} C_0 \frac{h_i^2}{\varepsilon} & \text{für } h_i < \varepsilon, \\ C_0 h_i & \text{für } \varepsilon \leq h_i, \end{cases} \quad (5.17)$$

wobei die Konstante $C_0 > 0$ klein genug ist, um (5.14) für $k \geq 2$ beziehungsweise (5.16) für $k = 1$ zu erfüllen. Dann erfüllt die Lösung $u_h \in P_k$ die Fehlerabschätzung

$$\|u - u_h\|_{SD} \leq C \left(\varepsilon^{1/2} h^k + h^{k+1/2} \right) |u|_{k+1}$$

mit einer von ε unabhängigen Konstanten C und $h = \max_{i=1,\dots,N} h_i$.



Abbildung 5.4: Joseph-Louis Lagrange (1736 – 1813).

Beweis: Sei $u_h^I \in V_h$ die Lagrange-Interpolierende von $u(x)$. Mit Dreiecksungleichung erhält man

$$\| \|u - u_h\| \|_{SD} \leq \| \|u - u_h^I\| \|_{SD} + \| \|u_h^I - u_h\| \|_{SD}.$$

Der erste Term auf der rechten Seite ist der Interpolationsfehler. Mit Hilfe der Interpolationsfehlerabschätzung (5.7), die man für jeden Term der SD-Norm anwendet, erhält man

$$\begin{aligned} \| \|u - u_h^I\| \|_{SD} &\leq \left(C\varepsilon h^{2k} |u|_{k+1}^2 + C\omega h^{2(k+1)} |u|_{k+1}^2 + C \sum_{i=1}^N \delta_i \|b\|_{\infty, I_i}^2 h_i^{2k} |u|_{k+1, I}^2 \right)^{1/2} \\ &\leq C \left(\varepsilon h^{2k} + h^{2(k+1)} + h^{2k+1} \right)^{1/2} |u|_{k+1} \\ &\leq C \left(\varepsilon^{1/2} h^k + h^{k+1/2} \right) |u|_{k+1}. \end{aligned}$$

Hierbei wurde $\delta_i \leq h_i \leq h$ ausgenutzt.

Betrachte nun den zweiten Term auf der rechten Seite. Die Koerzitivität, Satz 5.30 und die Galerkin-Orthogonalität (5.13) ergeben

$$\frac{1}{2} \| \|u_h^I - u_h\| \|_{SD}^2 \leq a_h(u_h^I - u_h, u_h^I - u_h) = a_h(u_h^I - u, u_h^I - u_h).$$

Nun wird die Dreiecksungleichung auf $a_h(u_h^I - u, u_h^I - u_h)$ angewandt und dann jeder Term einzeln abgeschätzt. Wesentlich dabei ist die Interpolationsabschätzung (5.7). Sei $w_h = u_h^I - u_h$. Für den Diffusionsterm gilt

$$\begin{aligned} \left| \varepsilon \left((u_h^I - u)', w_h' \right) \right| &\stackrel{\text{CSU}}{\leq} \varepsilon \| \|u_h^I - u\| \|_0 \|w_h'\|_0 = \varepsilon^{1/2} \| \|u_h^I - u\| \|_0 \varepsilon^{1/2} \|w_h'\|_0 \\ &\stackrel{(5.7)}{\leq} C\varepsilon^{1/2} h^k |u|_{k+1} \varepsilon^{1/2} \|w_h'\|_0 \leq C\varepsilon^{1/2} h^k |u|_{k+1} \| \|w_h\| \|_{SD}. \end{aligned}$$

Für den reaktiven Term erhält man auf ähnliche Art und Weise

$$\begin{aligned} \left| \left(c(u_h^I - u), w_h \right) \right| &\stackrel{\text{CSU}}{\leq} \|c\|_{\infty} \| \|u_h^I - u\| \|_0 \|w_h\|_0 = \omega^{-1/2} \|c\|_{\infty} \| \|u_h^I - u\| \|_0 \omega^{1/2} \|w_h\|_0 \\ &\stackrel{(5.7)}{\leq} Ch^{k+1} |u|_{k+1} \| \|w_h\| \|_{SD}. \end{aligned}$$

Als nächstes werden die Terme betrachtet, die man bei der SDFEM-Stabilisierung erhält. Wegen $\varepsilon \delta_i \leq C_0 h_i^2$ folgt

$$\begin{aligned}
& \left| \sum_{i=1}^N \left(-\varepsilon (u_h^I - u)'' , \delta_i b w_h' \right)_{I_i} \right| \\
& \stackrel{\text{CSU}}{\leq} \sum_{i=1}^N \varepsilon^{1/2} \left\| (u_h^I - u)'' \right\|_{0,I_i} \varepsilon^{1/2} \delta_i^{1/2} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \\
& \leq C_0^{1/2} \sum_{i=1}^N h_i \varepsilon^{1/2} \left\| (u_h^I - u)'' \right\|_{0,I_i} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \\
& \stackrel{\text{CSU}}{\leq} C_0^{1/2} \varepsilon^{1/2} h \left(\sum_{i=1}^N \left\| (u_h^I - u)'' \right\|_{0,I_i}^2 \right)^{1/2} \left(\sum_{i=1}^N \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i}^2 \right)^{1/2} \\
& \stackrel{(5.7)}{\leq} C \varepsilon^{1/2} h \left(\sum_{i=1}^N h_i^{2(k-1)} |u|_{k+1,I_i}^2 \right)^{1/2} \left(\sum_{i=1}^N \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i}^2 \right)^{1/2} \\
& \leq C \varepsilon^{1/2} h^k |u|_{k+1} \|w_h\|_{SD}.
\end{aligned}$$

Für die anderen Terme erhält man unter Nutzung von $\delta_i \leq C_0 h_i$

$$\begin{aligned}
& \left| \sum_{i=1}^N \left(b(u_h^I - u)' + c(u_h^I - u), \delta_i b w_h' \right) \right| \\
& \stackrel{\text{CSU}}{\leq} \sum_{i=1}^N \|b\|_\infty \left\| (u_h^I - u)' \right\|_{0,I_i} \delta_i^{1/2} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \varepsilon^{1/2} \|w_h'\|_0 \\
& \quad + \sum_{i=1}^N \|c\|_\infty \left\| (u_h^I - u) \right\|_{0,I_i} \delta_i^{1/2} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \\
& \leq C \left(\sum_{i=1}^N h_i^{1/2} \left\| (u_h^I - u)' \right\|_{0,I_i} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \right. \\
& \quad \left. + \sum_{i=1}^N h_i^{1/2} \left\| (u_h^I - u) \right\|_{0,I_i} \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i} \right) \\
& \leq C h^{1/2} \left[\left(\sum_{i=1}^N \left\| (u_h^I - u)' \right\|_{0,I_i}^2 \right)^{1/2} + \left(\sum_{i=1}^N \left\| (u_h^I - u) \right\|_{0,I_i}^2 \right)^{1/2} \right] \\
& \quad \times \left(\sum_{i=1}^N \left\| \sqrt{\delta_i} b w_h' \right\|_{0,I_i}^2 \right)^{1/2} \\
& \stackrel{(5.7)}{\leq} C \left(h^{k+1/2} + h^{k+3/2} \right) |u|_{k+1} \|w_h\|_{SD}.
\end{aligned}$$

Für eine optimale Abschätzung des konvektiven Terms muss man diesen erst partiell integrieren

$$\begin{aligned}
\left(b(u_h^I - u)', w_h \right) &= \left((u_h^I - u)', b w_h \right) = - \left(u_h^I - u, (b w_h)' \right) \\
&= - \left(u_h^I - u, b' w_h \right) - \left(u_h^I - u, b w_h' \right).
\end{aligned}$$

Nun schätzt man die letzten beiden Terme einzeln ab. Mit den gleichen Techniken wie bei den bisherigen Abschätzungen erhält man

$$\begin{aligned}
\left| \left(u_h^I - u, b' w_h \right) \right| &\leq \omega^{-1/2} \|b'\|_\infty \left(\sum_{i=1}^N \left\| u_h^I - u \right\|_{0,I_i}^2 \right)^{1/2} \omega^{1/2} \|w_h\|_0 \\
&\leq C h^{k+1} |u|_{k+1} \|w_h\|_{SD}.
\end{aligned}$$

Bei der Abschätzung des anderen Terms muss man unterscheiden, ob im Intervall I_i gilt $\varepsilon \leq h_i$ oder $\varepsilon > h_i$. Man erhält

$$\begin{aligned}
& \left| \left(u_h^I - u, bw_h' \right) \right| \\
& \stackrel{\text{CSU}}{\leq} \sum_{\varepsilon \leq h_i} \delta_i^{-1/2} \|u_h^I - u\|_{0,I_i} \left\| \sqrt{\delta_i} bw_h' \right\|_{0,I_i} + \sum_{\varepsilon > h_i} \|b\|_\infty \|u_h^I - u\|_{0,I_i} \|w_h'\|_{0,I_i} \\
& \stackrel{(5.7)}{\leq} C \left(\sum_{\varepsilon \leq h_i} \delta_i^{-1/2} h_i^{k+1} |u|_{k+1,I_i} \left\| \sqrt{\delta_i} bw_h' \right\|_{0,I_i} + \sum_{\varepsilon > h_i} h_i^{k+1} |u|_{k+1,I_i} \|w_h'\|_{0,I_i} \right) \\
& \stackrel{C_0 h_i \leq \delta_i, \varepsilon > h_i}{\leq} C \left(\sum_{\varepsilon \leq h_i} C_0^{-1/2} h_i^{-1/2} h_i^{k+1} |u|_{k+1,I_i} \left\| \sqrt{\delta_i} bw_h' \right\|_{0,I_i} \right. \\
& \quad \left. + \sum_{\varepsilon > h_i} h_i^{k+1/2} |u|_{k+1,I_i} \varepsilon^{1/2} \|w_h'\|_{0,I_i} \right) \\
& \stackrel{\text{CSU}}{\leq} Ch^{k+1/2} |u|_{k+1} \left[\left(\sum_{i=1}^N \left\| \sqrt{\delta_i} bw_h' \right\|_{0,I_i}^2 \right)^{1/2} + \varepsilon |w_h|_1 \right] \\
& \leq Ch^{k+1/2} |u|_{k+1} \|w_h\|_{SD}.
\end{aligned}$$

Fasst man nun alle Abschätzungen zusammen, so erhält man die Aussage des Satzes. \blacksquare

Bemerkung 5.35 Zur Konvergenzabschätzung.

- Wesentlich für die Abschätzung mit einer von ε unabhängigen Konstanten C ist, dass der Term

$$\left(\sum_{i=1}^N \left\| \sqrt{\delta_i} bw_h' \right\|_{0,I_i}^2 \right)^{1/2}$$

Bestandteil der Norm ist, in der man den Fehler abschätzt. Eine solche Abschätzung gilt für die von ε abhängigen Norm $\|\cdot\|_\varepsilon$ nicht.

- Auf der anderen Seite ist der Wert einer von ε unabhängigen Konstanten fraglich, da im allgemeinen $|u|_{k+1}$ von ε abhängen wird.
- In numerischen Simulationen sieht man oft eine Konvergenz der Ordnung h^{k+1} . Dabei spielen zum Beispiel strukturierte Gitter eine Rolle. In [Zho97] wurden Beispiele in zwei Dimensionen konstruiert, welche zeigen, dass die Abschätzung aus Satz 5.34 scharf ist.

\square

Beispiel 5.36 SDFEM. Das Standardbeispiel

$$-\varepsilon u'' + u' = 1 \quad \text{auf } (0, 1), \quad u(0) = u(1) = 0,$$

passt nicht in die Theorie der SDFEM, da $c(x) - \frac{b'(x)}{2} = 0$ ist. Trotzdem kann man auch auf dieses Beispiel die SDFEM–Stabilisierung anwenden. Eine Konvergenzanalyse ist in der Norm

$$\left(\varepsilon |v_h|_1^2 + \sum_{i=1}^N \left\| \sqrt{\delta_i} bv_h' \right\|_{0,I_i}^2 \right)^{1/2}$$

möglich. Man hatte also keine Kontrolle über den Fehler in der $L^2(0, 1)$ –Norm.

Ein grundlegendes Problem der Anwendung der SDFEM ist die freie Konstante C_0 in der Parameterdefinition (5.17). Am Standardbeispiel sieht man sehr gut, dass man für verschiedene Konstanten stark unterschiedliche Ergebnisse erhält, siehe

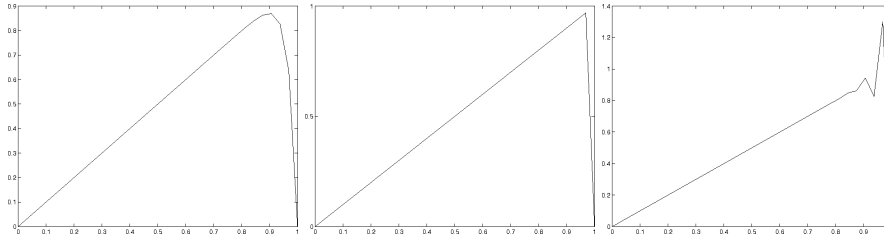


Abbildung 5.5: Mit der SDFEM berechnete Ergebnisse für das Standardbeispiel, $C_0 = 1$, $C_0 = 0.5$, $C_0 = 0.25$ von links nach rechts, $h = 1/32$, P_1 Finites-Element.

Abbildung 5.5. Ist C_0 zu groß, dann ist die Grenzschicht verschmiert, für ein geeignetes C_0 findet man eine Lösung die (fast) knotenexakt ist, und ist C_0 zu klein, dann entstehen an der Grenzschicht unphysikalische Oszillationen.

Für allgemeine Probleme ist es schwierig, C_0 geeignet zu wählen. In höheren Dimensionen wird man im allgemeinen auch kein C_0 mehr finden, so dass man eine (fast) knotenexakte Lösung erhält. Dafür besitzen mit der SDFEM berechnete Lösungen in höheren Dimensionen im allgemeinen unphysikalische Oszillationen an Grenzschichten. \square

Bemerkung 5.37 Andere Wahl des SDFEM-Parameters. Man nimmt auch statt (5.17) den Parameter aus Beispiel 5.33

$$\delta_i = \frac{h_i}{2 \|b\|_{L^\infty(I_i)}} \left(\coth(\text{Pe}_i) - \frac{1}{\text{Pe}_i} \right), \quad \text{Pe}_i = \frac{\|b\|_{L^\infty(I_i)} h_i}{2\varepsilon},$$

wobei Pe_K die lokale Péclet-Zahl ist. In dieser Definition hat man zwar keinen freien Parameter mehr, aber man stellt fest, dass bei Gleichungen in höheren Dimensionen unphysikalische Oszillationen in den berechneten Lösungen auftreten. \square

5.3 Lokale-Projektions-Stabilisierungs-Verfahren

Bemerkung 5.38 Motivation. Die bisher vorgestellten Stabilisierungsverfahren versuchen, mit einer möglichst geschickten Wahl von zusätzlicher Diffusion vernünftige Lösungen von konvektions-dominanten Problemen zu berechnen. Das funktioniert in einer Dimension ganz gut, in höheren Dimensionen aber bei weitem schlechter. Dort gibt es dann Stabilisierungsverfahren, welche noch komplizierte Ausdrücke zur Stabilisierung vorschlagen, als sie bisher präsentiert wurden. Oft sind diese Ausdrücke sogar nichtlinear.

Lokale-Projektions-Stabilisierungs-Verfahren (LPS-Verfahren) verfolgen eine andere Strategie. In diesen Verfahren wird eine einfache zusätzliche Diffusion verwendet, diese jedoch nur auf bestimmte Anteile (Skalen) der berechneten Lösung direkt angewandt. Dazu zerlegt man die Lösung in sogenannte große und kleine Skalen, wobei die zusätzliche Diffusion direkt nur auf die kleinen Skalen angewandt wird. Es gibt unterschiedliche grundsätzliche Herangehensweisen, wie man eine solche Zerlegung der Lösung definieren kann, zum Beispiel kann man große Skalen durch räumliche Mittelwerte definieren. Im Rahmen von Finite-Elemente-Methoden bieten sich jedoch Projektionen in geeignete Funktionenräume an, da man bei diesen Methoden ohnehin Funktionenräume im Rahmen der variationellen Formulierung zur Verfügung hat. Die kleinen Skalen sind dann als Differenz von allen Skalen und den großen Skalen definiert.

In diesem Abschnitt wird eine LPS-Methode vorgestellt. Die Darstellung stützt sich auf [Tob09]. \square

Bemerkung 5.39 Problemstellung. Es wird die schwache Formulierung (5.8) mit hinreichend glatten Koeffizientenfunktionen und mit der üblichen Bedingung für die Koerzivität

$$-\frac{1}{2}b'(x) + c(x) \geq \omega > 0 \quad \forall x \in [0, 1],$$

betrachtet.

Für die Finite-Elemente-Methode wird das Intervall mit einem Gitter $0 = x_0, \dots, x_N = 1$ zerlegt. \square

Bemerkung 5.40 Klassen von LPS-Methoden. Man kann zwei Klassen von LPS-Methoden unterscheiden. Zur Definition der großen Skalen mittels Projektion wird ein zusätzlicher Finite-Elemente-Raum benötigt. Diesen kann man auf dem gleichen Gitter wie oben wählen. In diesem Fall wird es nötig sein, mit Finite-Elementen höherer Ordnung für die Lösung zu arbeiten. Alternativ kann man den Raum für die Projektion auch auf einem gröberen Gitter definieren. Dann erhält man eine Zweigitter-Methode. Zweigitter-Methoden sind oft recht aufwendig zu implementieren. In diesem Abschnitt wird ein Verfahren aus der ersten Klasse betrachtet, eine sogenannte Eingitter-LPS-Methode. \square

Bemerkung 5.41 Eine Eingitter-LPS-Methode. Für ein $r \in \mathbb{N}$ seien der Lösungs- und der Projektionsraum der Eingitter-LPS-Methode wie folgt definiert

$$\begin{aligned} V_h &:= \{v_h \in H_0^1(0, 1) : v_h|_K \in P_r^+(K) \forall K \in \mathcal{T}_h\}, \\ D_h &:= \{q_h \in L^2(0, 1) : q_h|_K \in P_{r-1}(K) \forall K \in \mathcal{T}_h\}. \end{aligned}$$

Hierbei ist

$$P_r^+(K) = P_r(K) + \Theta_K P_{r-1}(K)$$

der angereicherte Raum zu $P_r(K)$, wobei Θ_K das Polynom kleinsten Grades bezeichnet, welches auf dem Rand von K verschwindet. In einer Dimension ist dies ein quadratisches Polynom, so dass $P_r^+(K) = P_{r+1}(K)$ gilt. In höheren Dimensionen $d \geq 2$ gilt $P_r^+(K) \subsetneq P_{r+d}(K)$.

Sei

$$\pi_h : L^2(0, 1) \rightarrow D_h, \quad r \mapsto \pi_h r : (r - \pi_h r, q_h) = 0 \quad \forall q_h \in D_h,$$

die $L^2(0, 1)$ -Projektion in D_h , $id : L^2(0, 1) \rightarrow L^2(0, 1)$ die Identität und $\kappa_h := id - \pi_h$ der Fluktuationsoperator. Der Stabilisierungsterm der LPS-Methode wird wie folgt definiert

$$S_{\text{LP}}(u_h, v_h) := \sum_{K \in \mathcal{T}_h} \tau_K (\kappa_h(bu'_h), \kappa_h(bv'_h))_K$$

mit den Stabilisierungsparametern $\{\tau_K\}$. Im Gegensatz zur SUPG-Methode ist diese Stabilisierung symmetrisch, das heißt $S_{\text{LP}}(u_h, v_h) = S_{\text{LP}}(v_h, u_h)$. Das wird in der Literatur oft als Vorteil der Methode angeführt.

Für die numerische Analysis benötigt man wieder eine geeignete Norm, in welcher die Stabilisierung auftaucht

$$\|v\|_{\text{LP}} := \left(\varepsilon |v|_1^2 + \omega \|v\|_0^2 + \sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h(bv')\|_{0,K}^2 \right)^{1/2}.$$

Da D_h ein unstetiger Finite-Elemente-Raum ist, ist die L^2 -Projektion gitterzellenweise definiert: $\pi_{h,K} : L^2(K) \rightarrow D_h(K)$

$$(\pi_{h,K} r - r, q_h)_K = 0 \quad \forall q_h \in L^2(K).$$

Diese Lokalität ist wichtig für die Effizienz der Methode. Auf dem Teilraum $D_h(K) = P_{r-1}(K) \subset L^2(K)$ ist dies die identische Abbildung und demzufolge gilt $\kappa_h r_h = r_h - \pi_{h,K} r_h = 0$ für alle $r_h \in P_{r-1}(K)$. Aus dem Bramble–Hilbert–Lemma folgt damit

$$\|\kappa_h q\|_{0,K} \leq Ch_K^l |q|_{l,K} \quad \forall q \in H^l(K), \quad 0 \leq l \leq r. \quad (5.18)$$

Die Eingitter–LPS–Methode lautet nun wie folgt: Finde $u_h \in V_h$, so dass für alle $v_h \in V_h$ gilt

$$a(u_h, v_h) + S_{\text{LP}}(u_h, v_h) = (f, v_h). \quad (5.19)$$

Der Schlüssel zur Fehlerabschätzung wird der Nachweis der Existenz eines Interpolationsoperators mit einer zusätzlichen Eigenschaft sein. \square

Lemma 5.42 Existenz eines speziellen Interpolationsoperators. *Es gibt einen Interpolationsoperator $j_h : H_0^1(0, 1) \rightarrow V_h$, so dass*

$$(j_h w - w, q_h) = 0 \quad \forall q_h \in D_h, \quad w \in H_0^1(0, 1), \quad (5.20)$$

und

$$|j_h w - w|_{m,K} \leq Ch_K^{l+1-m} \|w\|_{l+1,K} \quad \forall w \in H^{l+1}(K), \quad K \in \mathcal{T}_h, \quad (5.21)$$

für $l = 0, \dots, r+1$, $m \leq l+1$, gelten.

Beweis: Siehe [Tob09]. Das Besondere an diesem Interpolationsoperator ist nicht die Interpolationsgenauigkeit, sondern die Orthogonalität zu D_h . Der Interpolationsoperator wird nur in der Analysis der Methode verwendet und nicht in der Implementation. Er kann in einer Dimension sogar explizit charakterisieren werden. \blacksquare

Satz 5.43 Konvergenz der LPS–Methode. *Sei u die Lösung von (5.8) und sei u_h die Lösung von (5.19) mit $\tau_K \sim h_K$. Dann gilt die Fehlerabschätzung*

$$\|u - u_h\|_{\text{LP}} \leq C \left(\sum_{K \in \mathcal{T}_h} (\varepsilon + h_K) h_K^{2r} \|u\|_{r+1,K}^2 \right)^{1/2},$$

unter der Voraussetzung, dass $u \in H_0^1(0, 1) \cap H^{r+1}(0, 1)$.

Beweis: Mit partieller Integration zeigt man zunächst

$$a(v_h, v_h) \geq \varepsilon |v_h|_1^2 + \omega \|v_h\|_0^2 \quad \forall v_h \in V_h,$$

vergleiche Beispiel 4.9. Daraus folgt sofort

$$a(v_h, v_h) + S_{\text{LP}}(v_h, v_h) \geq \|v_h\|_{\text{LP}}^2 \quad \forall v_h \in V_h.$$

Setze nun $w_h = j_h u - u_h \in V_h$, dann folgt damit

$$\begin{aligned} & \|w_h\|_{\text{LP}}^2 \\ & \leq a(w_h, w_h) + S_{\text{LP}}(w_h, w_h) \\ & = a(j_h u - u, w_h) + a(u - u_h, w_h) + S_{\text{LP}}(j_h u - u, w_h) + S_{\text{LP}}(u - u_h, w_h) \\ & = a(j_h u - u, w_h) + (f, w_h) + S_{\text{LP}}(j_h u - u, w_h) + S_{\text{LP}}(u, w_h) \\ & \quad - a(u_h, w_h) - S_{\text{LP}}(u_h, w_h) \\ & = a(j_h u - u, w_h) + (f, w_h) + S_{\text{LP}}(j_h u - u, w_h) + S_{\text{LP}}(u, w_h) - (f, w_h) \\ & = a(j_h u - u, w_h) + S_{\text{LP}}(j_h u - u, w_h) + S_{\text{LP}}(u, w_h), \end{aligned}$$

wobei die Definition (5.19) der LPS-Methode verwendet wurde. Nun wird jeder Term einzeln abgeschätzt. Für den Interpolationsfehler wird (5.21) verwendet. Beim Diffusionsterm erhält man mit der Cauchy–Schwarz–Ungleichung für Integrale und für Summen

$$\begin{aligned}
|\varepsilon((j_h u - u)', w_h')| &= \left| \sum_{K \in \mathcal{T}_h} (\varepsilon(j_h u - u)', w_h')_K \right| \leq \sum_{K \in \mathcal{T}_h} \varepsilon |j_h u - u|_{1,K} |w_h|_{1,K} \\
&\leq C \sum_{K \in \mathcal{T}_h} \varepsilon h_K^r \|u\|_{r+1,K} |w_h|_{1,K} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} \varepsilon h_K^{2r} \|u\|_{r+1,K}^2 \right)^{1/2} \left(\sum_{K \in \mathcal{T}_h} \varepsilon |w_h|_{1,K}^2 \right)^{1/2} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} \varepsilon h_K^{2r} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}}.
\end{aligned}$$

Der konvektive Term wird partiell integriert

$$(b(j_h u - u)', w_h) = - (b(j_h u - u), w_h') - (b'(j_h u - u), w_h)$$

und der zweite Term wird zusammen mit dem reaktiven Term in der gleichen Art und Weise wie der Diffusionsterm abgeschätzt

$$|((c - b')(j_h u - u), w_h)| \leq C \left(\sum_{K \in \mathcal{T}_h} h_K^{2(r+1)} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}}.$$

Für den nächsten Term erhält man mit Cauchy–Schwarz–Ungleichungen, (5.18) und (5.21)

$$\begin{aligned}
|S_{\text{LP}}(j_h u - u, w_h)| &\leq \left(\sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h (b(j_h u - u)')\|_{0,K}^2 \right)^{1/2} \left(\sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h (b w_h')\|_{0,K}^2 \right)^{1/2} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} \tau_K \|b\|_{L^\infty(K)} \|(j_h u - u)'\|_{0,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} \tau_K \|b\|_{L^\infty(K)} h_K^{2r} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} h_K^{2r+1} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}}.
\end{aligned}$$

Für den letzten Term ergibt sich auf ähnliche Art und Weise

$$\begin{aligned}
S_{\text{LP}}(u, w_h) &\leq \left(\sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h (b u')\|_{0,K}^2 \right)^{1/2} \left(\sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h (b w_h')\|_{0,K}^2 \right)^{1/2} \\
&\leq \left(\sum_{K \in \mathcal{T}_h} \tau_K h_K^{2r} \|b u'\|_{r,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}} \\
&\leq C \left(\sum_{K \in \mathcal{T}_h} h_K^{2r+1} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}}.
\end{aligned}$$

Jetzt muss noch der erste Teil des konvektiven Terms abgeschätzt werden. Es gilt wegen der Orthogonalität des speziellen Interpolationsoperators

$$\begin{aligned}
(b(j_h u - u), w_h') &= (j_h u - u, b w_h') = (j_h u - u, \kappa_h (b w_h')) + (j_h u - u, b w_h' - \kappa_h (b w_h')) \\
&= (j_h u - u, \kappa_h (b w_h')) + (j_h u - u, \pi_h (b w_h')) = (j_h u - u, \kappa_h (b w_h')),
\end{aligned}$$

da $\pi_h(bw'_h) \in D_h$. Damit folgt unter Verwendung von (5.21)

$$\begin{aligned}
& |(j_h u - u, \kappa_h(bw'_h))| \\
& \leq \left(\sum_{K \in \mathcal{T}_h} \frac{1}{\tau_K} \|j_h u - u\|_{0,K}^2 \right)^{1/2} \left(\sum_{K \in \mathcal{T}_h} \tau_K \|\kappa_h(bw'_h)\|_{0,K}^2 \right)^{1/2} \\
& \leq C \left(\sum_{K \in \mathcal{T}_h} \frac{h_K^{2(r+1)}}{\tau_K} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}} \\
& \leq C \left(\sum_{K \in \mathcal{T}_h} h_K^{2r+1} \|u\|_{r+1,K}^2 \right)^{1/2} \|w_h\|_{\text{LP}}.
\end{aligned}$$

Der Beweis wird mit Hilfe der Dreiecksungleichung abgeschlossen

$$\|u - u_h\|_{\text{LP}} \leq \|u - j_h u\|_{\text{LP}} + \|j_h u - u_h\|_{\text{LP}},$$

wobei der erste Term mit (5.21) abgeschätzt wird und zur Abschätzung des zweiten Terms die Abschätzungen des Beweises verwendet werden. ■

Bemerkung 5.44 Zu LPS-Methoden.

- Durch die Blasenfunktionen und den Projektionsraum hat man bei einer Ein-Level-LPS-Methode wesentlich mehr Freiheitsgrade auf dem gleichen Gitter, als bei Methoden, welche auf Standard-Finite-Elemente-Räumen basieren, zum Beispiel der SUPG-Methode.
- Es stellt sich bei numerischen Studien heraus, dass die Wahl des Stabilisierungsparameters schwierig ist. Sowohl zu kleine als auch zu große Parameter führen zu schlechten Ergebnissen, das heißt, es sind große unphysikalische Oszillationen vorhanden. Falls man bei einer Wahl ein schlechtes Ergebnis hat, weiß man nicht, ob man die Parameter verkleinern oder vergrößern soll. In [Kno10] wurde eine LPS-Methode vorgestellt, bei der die Projektionsgebiete überlappend sind. Bei dieser Methode tritt diese Schwierigkeit nicht mehr auf. Zu kleine Parameter führen immer noch zu großen Oszillationen, aber zu große Parameter nur noch zu Verschmierungen.

Kapitel 6

Finite–Volumen–Methoden

Bemerkung 6.1 Grundlegende Idee. Finite–Volumen–Methoden (FVM) basieren auf Integralbilanzen über sogenannten Kontrollvolumen. Dazu wird das Intervall zunächst in kleine Gebiete, eben diese Kontrollvolumen, zerlegt und die Differentialgleichung wird über jedem Kontrollvolumen integriert. Im Anschluss wird partielle Integration (Gaußscher Satz in mehreren Dimensionen) angewandt, um die Integrale über den Kontrollvolumen, die Ableitungen enthalten, in Integrale auf dem Rand der Kontrollvolumen zu überführen. In einer Dimension, sind das Punktwerte in den Randpunkten. Dann verwendet man geeignete Approximationen für die Randintegrale, womit man ein Differenzenverfahren erhält.

Die Integralbilanzen können oft als Erhaltungsgesetze für physikalische Größen interpretiert werden. Deshalb werden Finite–Volumen–Methoden vor allem bei solchen Problemen mit Erfolg verwendet, bei denen die Erhaltung von Größen sehr wichtig ist, da diese Verfahren die Erhaltungseigenschaft bei der Approximation bewahren. Ein Beispiel sind inkompressible Strömungen, bei denen die Masse des Fluids in einem festen Strömungsgebiet konstant ist. \square

Beispiel 6.2 Standard–Finite–Volumen–Methode. Betrachte

$$-\varepsilon u''(x) + (b(x)u(x))' + c(x)u(x) = f(x) \text{ für } x \in (0, 1), \quad u(0) = u(1) = 0,$$

mit $b(x) \geq \beta > 0$ und $c(x) \geq 0$. Das Intervall wird mit Hilfe eines Gitters mit den Gitterpunkten $0 = x_0, \dots, x_N = 1$ zerlegt. Der Einfachheit halber sei das Gitter äquidistant mit Gitterweite h .

Finite–Volumen–Methoden benötigen ein Zweit–Gitter (secondary grid). Dieses wird in einer Dimension mit Hilfe der Mittelpunkte der Teilintervalle definiert. Setze

$$x_{i+1/2} := \frac{x_i + x_{i-1}}{2}, \quad i = 0, \dots, N - 1.$$

Die Kontrollvolumen werden mit Hilfe des Zweit–Gitters definiert

$$(0, x_{1/2}), (x_{1/2}, x_{3/2}), \dots, (x_{N-1/2}, 1).$$

Integration der Gleichung über ein Kontrollvolumen ergibt

$$\begin{aligned} & \int_{x_{i-1/2}}^{x_{i+1/2}} \left(-\varepsilon u''(x) + (b(x)u(x))' + c(x)u(x) \right) dx \\ &= -\varepsilon u'(x) \Big|_{x_{i+1/2}}^{x_{i-1/2}} + (bu)(x_{i+1/2}) - (bu)(x_{i-1/2}) + \int_{x_{i-1/2}}^{x_{i+1/2}} c(x)u(x) dx \\ &= \int_{x_{i-1/2}}^{x_{i+1/2}} f(x) dx. \end{aligned} \tag{6.1}$$

Die Terme werden nun durch Werte auf dem Originalgitter approximiert: die Ableitungen im ersten Term auf der linken Seite durch Differenzenquotienten, die Funktionswerte durch Mittelwerte und die Integrale durch Quadraturlformeln. Mögliche Varianten sind

$$\begin{aligned} u'(x_{i+1/2}) &\approx \frac{u_{i+1}^N - u_i^N}{h}, & u'(x_{i-1/2}) &\approx \frac{u_i^N - u_{i-1}^N}{h}, \\ g(x_{i\pm 1/2}) &\approx \frac{g(x_i) + g(x_{i\pm 1})}{2} & \int_{x_{i-1/2}}^{x_{i+1/2}} g(x) dx &\approx g(x_i)h. \end{aligned}$$

Dafür braucht man vernünftige Approximationen für $u'(0)$ und $u'(1)$.

Für konstantes $b(x)$ erhält man mit diesen Approximationen

$$-\varepsilon \left(\frac{u_{i+1}^N - u_i^N}{h} - \frac{u_i^N - u_{i-1}^N}{h} \right) + b \left(\frac{u_i^N + u_{i+1}^N}{2} - \frac{u_i^N + u_{i-1}^N}{2} \right) + c_i h u_i^N = f_i h$$

mit $c_i = c(x_i)$, $f_i = f(x_i)$. Das ist äquivalent zum zentralen Differenzschema

$$-\varepsilon \frac{u_{i+1}^N - 2u_i^N + u_{i-1}^N}{h^2} + b \frac{u_{i+1}^N - u_{i-1}^N}{2h} + c_i u_i^N = f_i.$$

□

Bemerkung 6.3 Cell-centered Finite-Volumen-Methoden. Finite-Volumen-Methoden, welche ein Zweit-Grid nutzen um die Kontrollvolumen zu definieren, werden cell-centered Finite-Volumen-Methoden genannt. Man kann auch das Originalgitter zur Definition der Kontrollvolumen verwenden. Diese Methoden heißen dann cell-vertex Finite-Volumen-Methoden. Die letzteren Methoden sind aber nicht besonders populär, da sie instabil sind. □

Bemerkung 6.4 Finite-Volumen-Methoden für singular gestörte Probleme. Um für singular gestörte Probleme eine stabile Finite-Volumen-Methode zu erhalten, muss man den Konvektionsterm $(bu)(x_{i\pm 1/2})$ in (6.1) durch einen Upwind-Term approximieren, zum Beispiel durch

$$(bu)(x_{i+1/2}) \approx b(x_{i+1/2}) (\lambda_i u_{i+1}^N + (1 - \lambda_i) u_i^N),$$

mit $\lambda_i \in [0, 1/2]$. Für $\lambda_i = 1/2$ erhält man das zentrale Differenzschema und für $\lambda_i = 0$ das einfache Upwind-Verfahren aus Definition 3.31. Mit Werten zwischen 0 und 1/2 kann man die Größe des Upwindings variieren.

Seien $b(x)$ und $\lambda_i = \lambda$ konstant. Dann erhält man mit der Upwind-Approximation

$$\begin{aligned} &(bu)(x_{i+1/2}) - (bu)(x_{i-1/2}) \\ &\approx b \left((\lambda u_{i+1}^N + (1 - \lambda) u_i^N) - (\lambda u_i^N + (1 - \lambda) u_{i-1}^N) \right) \\ &= b \left(\lambda u_{i+1}^N + (1 - 2\lambda) u_i^N - (1 - \lambda) u_{i-1}^N \right) \\ &= b \left(\frac{u_{i+1}^N - u_{i-1}^N}{2} + \left(\lambda - \frac{1}{2} \right) u_{i+1}^N + (1 - 2\lambda) u_i^N - \left(\frac{1}{2} - \lambda \right) u_{i-1}^N \right) \\ &= b \left(\frac{u_{i+1}^N - u_{i-1}^N}{2} \right) - \frac{bh(1 - 2\lambda)}{2} \left(\frac{u_{i+1}^N - 2u_i^N + u_{i-1}^N}{h} \right). \end{aligned}$$

Nun kann das stabilisierte Finite-Volumen-Verfahren als angepasstes Upwind-Verfahren (3.7) mit

$$\sigma(q) = 1 + q(1 - 2\lambda), \quad q = \frac{bh}{2\varepsilon},$$

interpretiert werden. Damit übertragen sich auch alle Eigenschaften von angepassten Upwind-Verfahren auf diese stabilisierte Finite-Volumen-Methode.

Insbesondere ist es auch möglich, dass Iljin-Allen-Southwell-Verfahren aus Definition 3.53 mit Hilfe einer Finiten-Volumen-Methode zu generieren, siehe [RST08]. \square

Bemerkung 6.5 Finite-Volumen-Methoden in höheren Dimensionen. Anders als in einer Dimension, sind Finite-Volumen-Methoden in höheren Dimensionen grundsätzlich von Finite-Differenzen-Methoden und Finite-Element-Methoden verschieden ! \square

Kapitel 7

Zusammenfassung und Ausblick

Bemerkung 7.1 Verfahren. Zur Diskretisierung von partiellen Differentialgleichungen gibt es im wesentlichen drei Verfahren:

- Finite-Differenzen-Methoden:
 - approximieren die Ableitungen der starken Form der Gleichung mit Hilfe von Differenzenquotienten,
 - einfach zu verstehen und zu implementieren,
 - Taylor-Entwicklung wesentlich in der Analysis,
- Finite-Element-Methoden:
 - basieren auf der schwachen (variationellen) Formulierung der zu Grunde liegenden Gleichung in Sobolev-Räumen,
 - approximieren den unendlich-dimensionalen Sobolev-Raum durch einen endlich-dimensionalen Raum,
 - Analysis basiert auf Konzepten der Funktionalanalysis,
- Finite-Volumen-Methoden:
 - basieren auf Integration der zu Grunde liegenden Gleichung,
 - sichern die Erhaltung von Größen in Kontrollvolumen.

□

Bemerkung 7.2 Dimension.

- In einer Dimension lassen sich die Verfahren oft ineinander überführen.
- In höheren Dimensionen sind die drei Herangehensweisen grundsätzlich verschieden. Alle Verfahren besitzen Vor- und Nachteile, zum Beispiel:
 - in komplizierten Gebieten sind Finite-Elemente und Finite-Volumen flexibler als Finite-Differenzen,
 - die Implementierung von Finite-Elementen ist wesentlich aufwändiger als die von Finite-Differenzen und Finite-Volumen,
 - Finite-Volumen-Methoden sind dort erfolgreich, wo man Erhaltungssätze erfüllen muss,
 - für Finite-Element-Methoden ist die Theorie am weitesten entwickelt.
- Ein neues Problem in höheren Dimensionen ist, dass komplizierte Gebiete auftreten können. Das hat sowohl Auswirkungen in der Analysis (Regularität der Lösung) als auch in der Praxis (Gittergenerierung).
- Die Gitterzellen in d Dimensionen sind d -dimensional. Diese Gitterzellen müssen geeignet angeordnet werden, damit ein vernünftiges Gitter entsteht. Das ist insbesondere bei komplizierten Gebieten nicht trivial. Gittergenerie-

rung, insbesondere in drei Dimensionen, ist ein wichtiges Forschungsgebiet.

□

Bemerkung 7.3 Singulär gestörte Probleme. Standard-Diskretisierungen berechnen nutzlose Lösungen für singulär gestörte Probleme, schon bei konstanten Koeffizienten. Man benötigt geeignete Stabilisierungen.

- In einer Dimension findet man Verfahren, um sehr gute Lösungen zu erhalten, zum Beispiel das Iljin–Allen–Southwell–Verfahren.
- In höheren Dimensionen ist die Entwicklung geeigneter stabiler Verfahren ein aktueller Forschungsgegenstand. Die bisher entwickelten Verfahren führen oft zu nicht zufriedenstellenden Lösungen (Grenzschichten zu stark verschmiert, unphysikalische Oszillationen).

□