

Eine neue dPG-Finite-Elemente-Methode niedriger Ordnung für Stokes Gleichungen

Masterarbeit

zur Erlangung des akademischen Grades
Master of Science (M. Sc.)

eingereicht von:	Sophie Louise Puttkammer
geboren am:	29.04.1991
geboren in:	Berlin
Betreuer:	Prof. Dr. Carsten Carstensen

Berlin, den 29. September 2015

Inhaltsverzeichnis

1	Einleitung	1
2	Vorbereitung	7
2.1	Triangulierungen	7
2.2	Vektor und Matrix Notation	9
2.3	Funktionenräume	10
2.4	Diskrete Funktionenräume	17
2.5	Weitere Hilfsmittel	19
2.6	Modellierung	23
3	Analysis der dPG-Methode	29
3.1	Allgemeine Informationen zur dPG-Methode	29
3.2	Problemformulierung und kontinuierliche Bedingungen	32
3.3	Diskrete Inf-Sup-Bedingung	41
3.4	Fortin-Interpolator	51
4	Implementierung	61
4.1	Löser und Fehlerschätzer	61
4.2	Basiswahl	64
4.3	Berechnung der Bestandteile des linearen Gleichungssystems	68
4.3.1	Koeffizientenmatrix	68
4.3.2	Normmatrix	79
4.3.3	Vektor zum Funktional F und zur linearen Nebenbedingung Λ . .	84
4.4	Exakter Fehler	90
4.5	Skalierung	93
4.6	Realisierung	94
5	Numerische Experimente	99
5.1	Das "colliding flow" Beispiel	99
5.2	Ein Poisson-Problem Beispiel	102
5.3	Ein Standardbeispiel im L -Gebiet	103
5.4	Ein Beispiel im Schlitzgebiet	106
5.5	Das "backward facing step" Beispiel	108
6	Zusammenfassung	110
	Literaturverzeichnis	111

1 Einleitung

Seien es Bauteile von Windrädern oder Flugzeugen, Wärmeverteilungen in Reinigungsanlagen, elektromagnetische Effekte in Mikrochips, in zahlreichen Anwendungen müssen partielle Differentialgleichungen gelöst werden, um Ergebnisse zu optimieren oder vorherzusagen. Das liegt daran, dass die mathematischen Modelle vieler physikalischer Prozesse genau solche Gleichungen sind. So gibt es die Maxwell Gleichungen für den Elektromagnetismus, die Navier-Lamé Gleichungen in der linearen Elastizität oder die Wärmeleitungsgleichungen, um die Temperaturverteilung in einem Medium zu charakterisieren. Diese Gleichungen haben gemein, dass sie nur in Spezialfällen analytisch gelöst werden können. Daher wurden eine Vielzahl numerischer Methoden entwickelt, um diese Probleme anzugehen.

In dieser Arbeit wird eine neue diskontinuierliche Petrov-Galerkin (dPG) Finite-Elemente-Methode zur Lösung von Stokes Gleichungen vorgestellt.

In der ersten Hälfte des 19. Jahrhunderts formulierten Claude Navier und Georg Stokes unabhängig voneinander ein System partieller Differentialgleichungen zur Beschreibung des Strömungsverhaltens Newtonscher Flüssigkeiten. Ein Spezialfall und grundlegendes Modellproblem der Strömungsmechanik sind die Stokes Gleichungen für inkompressible Newtonsche Flüssigkeiten. In Abschnitt 2.6 werden einige der zahlreichen Formulierungen dieses Systems erläutert und hergeleitet. Das Modell in dieser Arbeit basiert auf der Pseudospannungs-Geschwindigkeits Formulierung. Gegeben ein Gebiet Ω , eine äußere Kraftdichte $f \in L^2(\Omega; \mathbb{R}^n)$ und Randdaten $g \in L^2(\partial\Omega; \mathbb{R}^n)$, werden dabei ein Geschwindigkeitsfeld $u : \Omega \rightarrow \mathbb{R}^n$ und eine Pseudospannung $\sigma : \Omega \rightarrow \mathbb{R}^{n \times n}$ gesucht, für die gilt

$$\begin{aligned} \operatorname{div} \sigma + f &= 0 & \text{in } \Omega, \\ \operatorname{dev} \sigma - D u &= 0 & \text{in } \Omega, \\ u &= g & \text{entlang } \partial\Omega. \end{aligned}$$

Diese Gleichungen spielen beispielsweise eine Rolle, wenn Strömungen an einer rückwärts gewandten Stufe betrachtet werden, wie im "backward facing step" Beispiel. Für diese Anwendung kann im Allgemeinen analytisch keine Lösung bestimmt werden. Gegeben sind ein leicht verzerrtes L -Gebiet, $\Omega = ((-2, 8 \times (-1, 1)) \setminus ((-2, 0) \times (-1, 0)))$, sowie die konstante rechte Seite $f \equiv 0$ und die Dirichlet Randdaten für $(x_1, x_2)^\top \in \partial\Omega$

$$g((x_1, x_2)^\top) = \begin{cases} (0, 0)^\top & \text{für } -2 < x_1 < 8, \\ 1/10(-x_2(x_2 - 1), 0)^\top & \text{für } x_1 = -2, \\ 1/80(-(x_2 - 1)(x_2 + 1), 0)^\top & \text{für } x_1 = 8. \end{cases}$$

Mit dem im Zuge dieser Arbeit implementierten Programm kann für dieses Problem das in Abbildung 1.1 dargestellte Geschwindigkeitsfeld ermittelt werden. In Kapitel 5

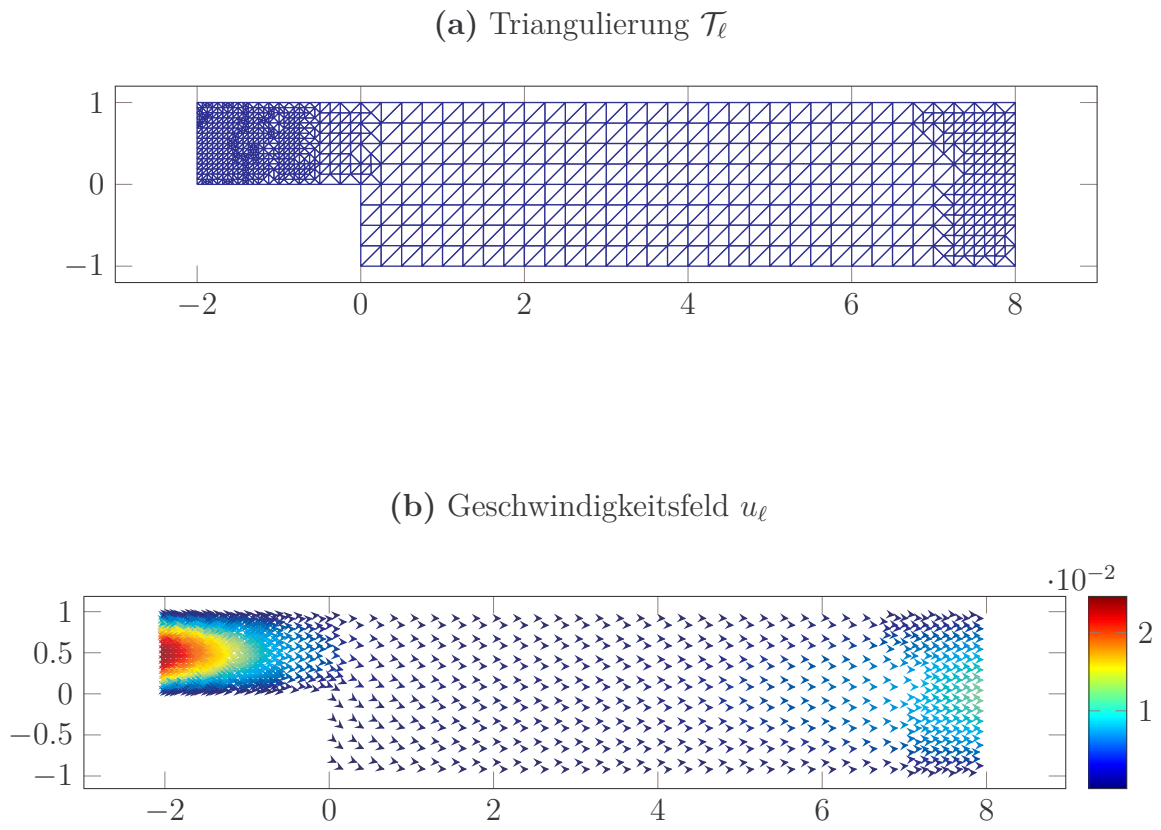


Abbildung 1.1: Triangulierung mit 1551 Elementen (15511 Freiheitsgraden) und entsprechendes Geschwindigkeitsfeld im "backward facing Step" Beispiel bei adaptiver Verfeinerung mit $\theta = 0.3$

werden dieses und weitere mit der neuen dPG-Methode untersuchte Beispiele genauer diskutiert.

Dem Namen nach ist die dPG-Methode eine Verbindung aus einer diskontinuierlichen Galerkin-Methode, bei welcher gebrochene Test- und Ansatzfunktionen möglich sind, und einer Petrov-Galerkin-Methode, bei der sich Test- und Ansatzraum unterscheiden dürfen. Der Begriff diskontinuierliche Petrov-Galerkin Methode wurde zuerst in [BMS02] verwendet. Dort wurde bereits eine ultra-schwache Formulierung betrachtet, es wurden also alle Ableitungen auf die Testfunktionen verschoben. Außerdem wurden die numerischen Flüsse der diskontinuierlichen Galerkin-Methode durch unabhängige, nur auf Kanten definierte Variablen ersetzt. Die Idee, die optimalen Testräume, welche die Methode ursprünglich motivierten, zu bestimmen, wurde von Leszek Demkowicz und Jay Gopalakrishnan aufgegriffen. In [DG10, DG11a, DG11b, DGN12] untersuchten sie diese Klasse von Methoden und ihre Anwendung beispielsweise für Wellentransportprobleme und Konvektions-Diffusions-Gleichungen. In [RBTD14] betrachteten Demkowicz, Roberts und Bui-Than die Anwendungsmöglichkeiten der dPG-Methode für Stokes Probleme. Sie verwenden dabei stückweise Polynome vom Grad $k \geq 1$ zur Approximation der vorkommenden L^2 -Funktionen, also des Spannungstensors und des Geschwindigkeitsfeldes, und für die Normalenspur der Spannung in $H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$. Die Spuren der Geschwindigkeit in $H^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ werden mit stückweisen Polynomen des Grades $k + 1$ diskretisiert. Der Testraum wird bei dieser Methode durch sukzessives Anreichern und Vergleichen der Ergebnisse ermittelt. Ziel ist es, die minimale Anreicherung zu finden, bei der die Ergebnisse ungefähr so gut sind wie bei größeren Anreicherungen. In ihren Experimenten wird meist der Anreicherungsgrad 1 verwendet, d.h. stückweise Polynome und Raviart-Thomas Funktionen vom Grad $k + 2$.

Die in dieser Arbeit vorgestellte neue dPG-Methode verringert den Rechenaufwand enorm, denn sie kommt mit dem diskreten Ansatzraum

$$X_h := P_0(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_0(\mathcal{T}; \mathbb{R}^n) \times S_0^1(\mathcal{E}; \mathbb{R}^n) \times P_0(\mathcal{E}; \mathbb{R}^n)$$

und dem diskreten Testraum

$$Y_h := RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^n)$$

aus. Dass hier sämtliche theoretischen Betrachtungen für eine beliebige Dimension $n \in \mathbb{N}$ durchgeführt werden, erlaubt dabei eine problemlose Verallgemeinerung der vorgestellten zweidimensionalen Implementierung.

Im Allgemeinen kann eine dPG-Methode als Least-Squares Finite-Elemente-Methode, die das Residuum der ultra-schwachen Formulierung in einer Nicht-Standardnorm minimiert, oder als Finite-Elemente-Methode mit Nicht-Standardtesträumen aufgefasst werden. Den Betrachtungen liegt dabei eine wohlgestellte Variationsformulierung einer partiellen Differentialgleichung zu Grunde. Seien dazu zwei Banachräume X und Y , eine stetige Bilinearform $b : X \times Y \rightarrow \mathbb{R}$ und ein Funktional $F \in Y^*$ gegeben. Gesucht wird dann

eine Lösung $x \in X$, welche für jedes $y \in Y$ die Gleichung

$$b(x, y) = F(y)$$

erfüllt. Wie bereits erwähnt ist eine Besonderheit der dPG-Methoden, dass schon auf diesem kontinuierlichen Level eine Triangulierung zu Grunde gelegt wird, also gebrochene Testfunktionen betrachtet werden. In der praktischen dPG-Methode wird für diskrete Unterräume $X_h \subseteq X$ und $Y_h \subseteq Y$ ein $x_h \in X_h$ gesucht, für das gilt

$$x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}.$$

Aus [CDG14] und [CH15] geht hervor, dass die folgenden Bedingungen genügen, um gewisse Fehlerabschätzungen zu garantieren. Im Hilbertraum Fall ist zu zeigen, dass die Bilinearform b beschränkt ist und einen trivialen Kern besitzt und dass die folgende kontinuierliche Inf-Sup-Bedingung erfüllt ist

$$0 < \beta := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y}.$$

Für die konforme Diskretisierung ist dann lediglich die diskrete Inf-Sup-Bedingung

$$0 < \beta_h := \inf_{x_h \in X_h \setminus \{0\}} \sup_{y_h \in Y_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y}$$

zu überprüfen. Dieser letzte Punkt ist natürlich leichter zu garantieren, je größer der diskrete Testraum $Y_h \subseteq Y$ gewählt wird. Sind die genannten Bedingungen gleichzeitig erfüllt, so sind optimale Konvergenz und diskrete Stabilität im Sinn der folgenden a-priori-Abschätzung gewährleistet

$$\|x - x_h\|_X \leq C \min_{\xi_h \in X_h} \|x - \xi_h\|_X,$$

wobei $C \leq \|b\|^2 / \beta \beta_h$ eine generische nur von Ω abhängige Konstante ist. Außerdem gilt die a-posteriori-Abschätzung

$$\begin{aligned} \beta \|x - \xi_h\|_X &\leq \|b\| / \beta_h \|F - b(\xi_h, \bullet)\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \\ &\leq 2\|b\|^2 / \beta_h \|x - \xi_h\|_X \end{aligned}$$

für jedes $\xi_h \in X_h$, die Methode ist also auch bei inexaktem Lösen stabil. In dieser Abschätzung kommt neben dem Fehlerschätzer $\|F - b(\xi_h, \bullet)\|_{Y_h^*}$, der für adaptives Verfeinern verwendet wird, auch einen Datenapproximationsfehler $\|F \circ (1 - \Pi)\|_{Y^*}$ vor. Dabei wird mit $\Pi : Y \rightarrow Y_h$ ein Fortin-Interpolator bezeichnet, also ein beschränkter, linearer Operator, der die Eigenschaft $b(\bullet, (1 - \Pi)y) = 0$ in X^* für alle $y \in Y$ erfüllt. Die Existenz dieses Operators ist im Grunde zu der diskreten Inf-Sup-Bedingung äquivalent. In den meisten Arbeiten wird dieser Operator daher direkt konstruiert und auf den Beweis der

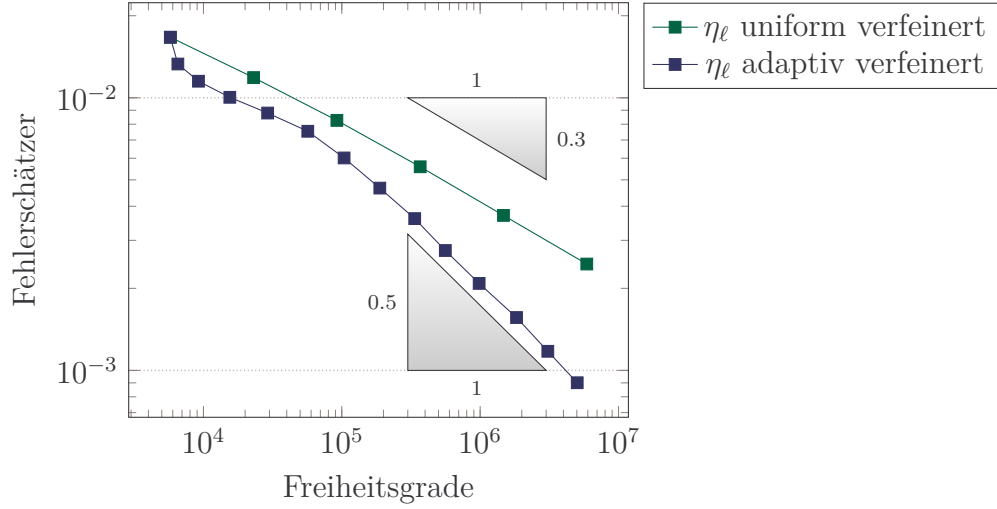


Abbildung 1.2: Konvergenzplot des Fehlerschätzers für das "backward facing step" Beispiel bei uniformem Verfeinern und adaptivem Verfeinern mit Bulkparameter $\theta = 0.3$

letzten Bedingung verzichtet. Die hier gewählte Herangehensweise, die diskrete Inf-Sup-Bedingung direkt nachzuweisen, hat den Vorteil, dass ein deutlich kleinerer Testraum gewählt werden kann. Dafür ist der Datenapproximationsfehler schwerer zu bestimmen. Im zweidimensionalen Fall kann allerdings ein solcher Operator Π angegeben werden. Damit schließt diese Arbeit direkt an die aktuelle Forschung an. Die numerischen Ergebnisse zeigen, dass bereits mit Hilfe von $\|F - b(\xi_h, \cdot)\|_{Y_h^*}$ optimale Konvergenzraten des adaptiven Verfahrens erzielt werden, wie z.B. in Abbildung 1.2 zu erkennen.

In diesem allgemeinen Rahmen wird die neue dPG-Methode hergeleitet. Dabei wird das Residuum der ultra-schwachen Formulierung in der Dualnorm des gebrochenen Testraumes minimiert. In der ultra-schwachen Form der Pseudostress-Geschwindigkeits Formulierung der Stokes Gleichungen werden die Spuren der Geschwindigkeit und Normalspuren der Pseudospannung auf dem Rand der einzelnen Elemente der Triangulierung durch neue Variablen substituiert. Dies führt zu folgender Bilinearform $b : X \times Y \rightarrow \mathbb{R}$, für alle

$$\begin{aligned} x &= (\boldsymbol{\sigma}, u, s, t) \in X := L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n), \\ &\quad \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n) \\ y &= (\boldsymbol{\tau}, v) \in Y := H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times H^1(\mathcal{T}; \mathbb{R}^n) \end{aligned}$$

sei

$$\begin{aligned} b(x, y) &:= -\langle t, \gamma_0 v \rangle_{\partial\mathcal{T}} - \langle \gamma_\nu \boldsymbol{\tau}, s \rangle_{\partial\mathcal{T}} \\ &\quad + \int_{\Omega} \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx + \int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \, dx. \end{aligned}$$

Für diese Bilinearform die zuvor erwähnten Bedingungen nachzuweisen, wird einen Großteil dieser Arbeit ausmachen. Der Beweis der kontinuierlichen Inf-Sup-Bedingung wird durch das Splitting Lemma und die Wohlgestelltheit der dualen gemischten Formulierung ermöglicht. Die diskrete Inf-Sup-Bedingung wird nachgewiesen, in dem zu jedem $x_h \in X_h$ ein $\tilde{y}_h \in Y_h$ so konstruiert wird, dass die Ungleichungen $\|\tilde{y}_h\|_Y \leq C_{\tilde{y}_h} b(x_h, \tilde{y}_h)$ und $\|x_h\|_X \leq C_{x_h} b(x_h, \tilde{y}_h)$ für zwei direkt angegebene generische Konstanten $C_{\tilde{y}_h}$, C_{x_h} gelten. Dabei erweisen sich die Spurterme als besonders kompliziert, weshalb Fortsetzungen $w_c \in P_1(\mathcal{T}; \mathbb{R}^n)$ von $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^n)$ bzw. $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ von $t \in P_0(\mathcal{E}; \mathbb{R}^n)$ verwendet werden.

Diese Arbeit ist wie folgt aufgebaut. Der nächste Abschnitt beschäftigt sich mit der Einführung bzw. Wiederholung grundlegender Notationen und Definitionen. Des Weiteren werden die kontinuierlichen und diskreten Funktionsräume vorgestellt und einige bekannte bzw. häufiger benötigte Lemmata zitiert bzw. nachgewiesen. Zuletzt erfolgt eine kurze Einführung zu den Stokes Gleichungen und einigen numerischen Verfahren zu deren Lösung.

Im dritten Kapitel wird die neue dPG-Methode analysiert. Dazu werden zunächst die allgemeinen Grundlagen zu dieser Klasse numerischer Methoden erläutert. Es werden die verschiedenen Problemformulierungen eingeführt und die Bedingungen für die bereits erwähnten a-priori- und a-posteriori-Fehlerabschätzungen vorgestellt. Diese Bedingungen werden anschließend nachgewiesen – zunächst für das auf kontinuierlichem Level formulierte Problem und danach für die gewählte Diskretisierung. Als Letztes wird in diesem Abschnitt noch der Fortin-Interpolator dieser Methode untersucht, der Gegenstand aktueller Forschung ist.

Der vierte Abschnitt widmet sich der Implementierung. Dazu werden zunächst die grundlegenden Berechnungen durchgeführt und anschließend die im Zuge dieser Arbeit entstandenen Programme kurz vorgestellt.

Vor der Zusammenfassung werden im fünften Kapitel einige numerischen Experimente ausgewertet, die mit Hilfe der entstandenen Programme durchgeführt wurden. Diese Benchmark Probleme für die Stokes Gleichungen erlauben, die Ergebnisse mit denen anderer Methoden zu vergleichen.

2 Vorbereitung

In diesem Kapitel werden grundlegende Notationen eingeführt bzw. wiederholt und gewisse nützliche Eigenschaften der verwendeten Räume und Operatoren festgehalten. In dieser Arbeit wird die Schreibweise $a \lesssim b$ verwendet, wenn eine generische Konstante C existiert, so dass $a \leq Cb$ gilt. Des Weiteren wird $a \lesssim b \lesssim a$ durch $a \approx b$ abgekürzt.

2.1 Triangulierungen

In dieser Arbeit sei $\Omega \subseteq \mathbb{R}^n$ stets ein beschränktes, polygonal berandetes Lipschitz-Gebiet mit Dirichlet Rand Γ . Dieses Gebiet wird in abgeschlossene n -Simplexe unterteilt, so dass sich eine reguläre Triangulierung ergibt.

Definition 2.1 (ℓ -Simplex). Sei $\ell \in \mathbb{N}$ und $\ell \leq n$. Eine Menge $T \subset \mathbb{R}^n$ heißt ℓ -Simplex, falls T die konvexe Hülle von $\ell + 1$ paarweise linear unabhängigen Punkten im \mathbb{R}^n ist.

Bei einem n -Simplex T bezeichne $\mathcal{N}(T) := \{P_1, \dots, P_{n+1}\}$ die Menge aller Eckpunkte bzw. Knoten von T , $\mathcal{E}(T) = \{\text{conv}(\mathcal{P}) : \mathcal{P} \subset \mathcal{N}(T), |\mathcal{P}| = n\}$ die Menge aller Seiten von T und $\text{int}(T)$ alle inneren Punkte.

Definition 2.2 (Reguläre Triangulierung). Eine *reguläre Triangulierung* von Ω ist eine endlich Menge \mathcal{T} aus abgeschlossenen n -Simplexen mit folgenden Eigenschaften

- (i) $\text{int}(T) \neq \emptyset$ für alle $T \in \mathcal{T}$,
- (ii) $\bar{\Omega} = \bigcup_{T \in \mathcal{T}} T$,
- (iii) $\text{int}(T) \cap \text{int}(\hat{T}) = \emptyset$ für alle $T, \hat{T} \in \mathcal{T}$ mit $T \neq \hat{T}$,
- (iv) wenn für $T, \hat{T} \in \mathcal{T}$ mit $T \neq \hat{T}$ gilt $T \cap \hat{T} \neq \emptyset$, so ist der Schnitt von T und \hat{T} die konvexe Hülle der $\ell \leq n$ gemeinsamen Knoten und $T \cap \hat{T} = \text{conv}(\mathcal{N}(T) \cap \mathcal{N}(\hat{T}))$ ist ein $\ell - 1$ Simplex, d.h. es gibt in der Triangulierung keine hängenden Knoten.

Für eine reguläre Triangulierung \mathcal{T} von Ω sei $\mathcal{N} := \bigcup_{T \in \mathcal{T}} \mathcal{N}(T)$ die Menge aller Knoten der Triangulierung und $\mathcal{N}(\Omega)$ bzw. $\mathcal{N}(\partial\Omega)$ die Mengen der inneren Knoten bzw. Randknoten. Mit \mathcal{E} , $\mathcal{E}(\Omega)$ bzw. $\mathcal{E}(\partial\Omega)$ werden die Menge aller Seiten, der inneren Seiten bzw. der Randseiten der Triangulierung beschrieben.

Definition 2.3 (Skelett). Das *Skelett* einer regulären Triangulierung \mathcal{T} ist durch

$$\partial\mathcal{T} := \bigcup_{T \in \mathcal{T}} \bigcup_{E \in \mathcal{E}(T)} E$$

definiert.

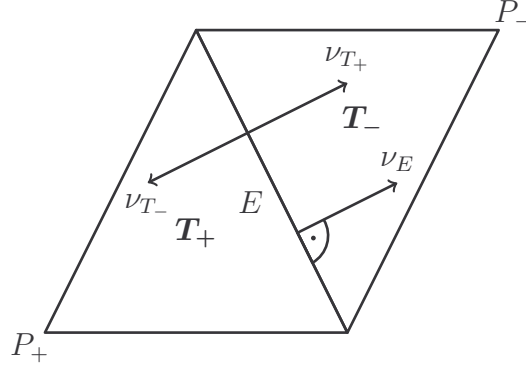


Abbildung 2.1: Seitenpatch ω_E für $E \in \mathcal{E}(\Omega)$ für $\Omega \subseteq \mathbb{R}^2$

Auf jedem n -Simplex $T \in \mathcal{T}$ bezeichne ν_T das äußere Einheitsnormalenvektorfeld entlang des Randes ∂T . Außerdem ist jede Seite $E \in \mathcal{E}$ mit einer Orientierung der Einheitsnormale ν_E , wie in Abbildung 2.1 dargestellt, ausgestattet. Entlang einer Außenseite, $E \in \mathcal{E}(\partial\Omega)$, zeige ν_E stets nach außen. Entlang einer inneren Seite, $E = \partial T_+ \cap \partial T_- \in \mathcal{E}(\Omega)$, werde eine der beiden möglichen Orientierungen der Normale ν_E fixiert. Die Nachbarelemente werden stets so benannt, dass ν_E von T_+ nach T_- zeigt, also $\nu_E = \nu_{T_+}$. Damit wird die folgende Vorzeichenfunktion definiert

$$\text{sgn}(T, E) := \nu_E \cdot \nu_T \in \{\pm 1\}, \text{ für alle } T \in \mathcal{T}, E \in \mathcal{E}(T). \quad (2.1)$$

Die Konvention der Elementbenennung wird auch bei der Definition von Sprüngen verwendet, welche bei nur stückweise stetigen Funktionen auf der Triangulierung \mathcal{T} interessant werden.

Definition 2.4 (Sprung). Der *Sprung* einer gegebenen Lebesgue integrierbaren Funktion $v \in L^2(\mathcal{T}; \mathbb{M})$ (siehe Abschnitt 2.3) mit $\mathbb{M} \subset \mathbb{R}^{m \times n}$, $n, m \in \mathbb{N}_{>0}$ entlang einer inneren Seite $E \in \mathcal{E}(\Omega)$ ist definiert als $[v]_E := (v|_{T_+} - v|_{T_-})|_E \in L^2(E; \mathbb{M})$ und entlang einer Randseite $E \in \mathcal{E}(\partial\Omega)$ durch $[v]_E := v|_E \in L^2(E; \mathbb{M})$.

Die folgenden Größen der Triangulierung bzw. der Elemente werden in dieser Arbeit ebenfalls eine Rolle spielen.

Es bezeichne $P_0(\mathcal{T}; \mathbb{M})$ mit $\mathbb{M} \subset \mathbb{R}^{m \times n}$, $n, m \in \mathbb{N}_{>0}$, die stückweise konstanten Funktionen, die auf jedem Element $T \in \mathcal{T}$ einen konstanten Wert in \mathbb{M} haben, also in $P_0(T; \mathbb{M})$ liegen (siehe auch Abschnitt 2.4).

Für jedes Element $T \in \mathcal{T}$ sei

$$h_T := \text{diam}(T) = \max_{x, y \in T} |x - y| \in P_0(T)$$

der Durchmesser und

$$h_{\max} := \max_{T \in \mathcal{T}} h_T$$

bezeichne die maximale Netzweite der Triangulierung \mathcal{T} . Des Weiteren sei die Funktion $h_{\mathcal{T}} \in P_0(\mathcal{T}; \mathbb{R})$ für jedes $T \in \mathcal{T}$ definiert durch $h_{\mathcal{T}}|_T := h_T$.

Der Schwerpunkt eines n -Simplexes $T \in \mathcal{T}$ wird mit

$$\text{mid}(T) := \oint_T x \, dx = |T|^{-1} \int_T x \, dx = \frac{1}{n+1} \sum_{z \in \mathcal{N}(T)} z$$

bezeichnet. Die Funktion $\text{mid}(\mathcal{T}) \in P_0(\mathcal{T}; \mathbb{R}^n)$ habe den Wert $\text{mid}(T)$ auf jedem $T \in \mathcal{T}$. Wenn außerdem \bullet die Identitätsabbildung in \mathbb{R}^n und \perp die Orthogonalität in $L^2(T)$ bezeichnen, so gilt

$$\bullet - \text{mid}(T) \perp P_0(\mathcal{T}) \quad \text{und} \quad \|\bullet - \text{mid}(T)\|_{L^\infty(\Omega)} \leq h_{\max} n/(n+1). \quad (2.2)$$

Die erste Tatsache ergibt sich sofort aus der Integraldefinition des Schwerpunkts, die zweite ergibt sich auf $T = \text{conv}\{P_1, \dots, P_{n+1}\} \in \mathcal{T}$, da das Maximum von $|x - \text{mid}(T)|$ für $x \in T$ in einem Eckpunkt angenommen wird. O.B.d.A werde das Maximum bei P_1 erreicht, dann gilt für alle $x \in T$

$$\begin{aligned} |x - \text{mid}(T)| &\leq \left| P_1 - \frac{1}{n+1} \sum_{j=1}^{n+1} P_j \right| = \frac{1}{n+1} \left| nP_1 - \sum_{j=2}^{n+1} P_j \right| \\ &\leq \frac{1}{n+1} \sum_{j=2}^{n+1} |P_1 - P_j| \leq h_T \frac{n}{n+1}. \end{aligned}$$

2.2 Vektor und Matrix Notation

Dieser Paragraph stellt Details der in dieser Arbeit verwendeten Matrix und Vektor Notation klar.

Für zwei Vektoren $a, b \in \mathbb{R}^n$ wird das Skalarprodukt wie folgt notiert,

$$a \cdot b = a^\top b = \sum_{j=1}^n a_j b_j \in \mathbb{R},$$

wohingegen $A : B$ das Skalarprodukt von zwei $n \times n$ Matrizen $A, B \in \mathbb{R}^{n \times n}$ symbolisiert,

$$A : B = \sum_{j,k=1}^n A_{jk} B_{jk} \in \mathbb{R}.$$

Für einen Vektor $a \in \mathbb{R}^n$ und eine Matrix $B \in \mathbb{R}^{n \times m}$ sei $a \cdot B = a^\top B \in \mathbb{R}^{1 \times m}$.

Das dyadische Produkt zweier Vektoren $a, b \in \mathbb{R}^n$ ist definiert als

$$a \otimes b := ab^\top \in \mathbb{R}^{n \times n}.$$

Es gilt $|a \otimes b| = |a||b|$.

Mit \bullet wird jeweils die Identitätsabbildung bezeichnet. Die Notation $|\bullet|$ ist kontextabhängig, sie steht für die durch \cdot bzw. $:$ induzierte Norm auf \mathbb{R}^n bzw. $\mathbb{R}^{n \times n}$, die Kardinalität einer endlichen Menge oder das n - bzw. $(n-1)$ -dimensionale Lebesgue-Maß einer Teilmenge von \mathbb{R}^n .

In einem normierten Vektorraum $(X, \|\bullet\|_X)$ wird die Sphäre wie folgt definiert

$$S(X) := \{x \in X : \|x\|_X = 1\}.$$

Der Deviator einer Matrix $A \in \mathbb{R}^{n \times n}$ ist ein linearer Operator definiert durch $\text{dev } A := A - 1/n \text{ tr}(A) \text{ I}_{n \times n}$, wobei der lineare Operator $\text{tr } A = A_{11} + \dots + A_{nn} = \sum_{j=1}^n A_{jj}$ die Spur der Matrix bezeichnet. Der Raum $\mathbb{R}_{\text{dev}}^{n \times n}$ enthält alle deviatorischen, also spurfreien, $n \times n$ Matrizen, d.h.

$$\mathbb{R}_{\text{dev}}^{n \times n} := \{A \in \mathbb{R}^{n \times n} : \text{dev } A = A\} = \{A \in \mathbb{R}^{n \times n} : \text{tr } A = 0\}. \quad (2.3)$$

Lemma 2.5 (Eigenschaften des Deviators). *Es gilt*

- (i) $\text{tr dev } A = 0$ für alle $A \in \mathbb{R}^{n \times n}$ und damit $\text{dev dev } A = \text{dev } A$,
- (ii) $\text{dev } A : \text{dev } B = \text{dev } A : B = A : \text{dev } B$ für alle $A, B \in \mathbb{R}^{n \times n}$,
- (iii) $\|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 = 1/n \|\text{tr } \boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\text{dev } \boldsymbol{\tau}\|_{L^2(\Omega)}^2$ für alle $\boldsymbol{\tau} \in L^2(\Omega; \mathbb{R}^{n \times n})$, insbesondere $\|\text{dev } \boldsymbol{\tau}\|_{L^2(\Omega)} \leq \|\boldsymbol{\tau}\|_{L^2(\Omega)}$.

Beweis. Punkt (i) gilt, da $\text{tr}(\text{I}_{n \times n}) = n$. Punkt (ii) ergibt sich, weil $\text{I}_{n \times n} : A = \text{tr}(A)$ für alle $A \in \mathbb{R}^{n \times n}$ und damit $\text{dev } A : \text{dev } B = A : B - 1/n \text{ tr}(A) \text{ tr}(B) = \text{dev } A : B = A : \text{dev } B$. Punkt (iii) folgt aus der Orthogonalität von deviatorischem und isochorischem Anteil, denn

$$\text{dev } \boldsymbol{\tau} : \text{tr}(\boldsymbol{\tau}) \text{ I}_{n \times n} = (\boldsymbol{\tau} - 1/n \text{ tr}(\boldsymbol{\tau}) \text{ I}_{n \times n}) : \text{tr}(\boldsymbol{\tau}) \text{ I}_{n \times n} = \text{tr}(\boldsymbol{\tau})^2 - \text{tr}(\boldsymbol{\tau})^2 = 0$$

und $\boldsymbol{\tau} = \text{dev } \boldsymbol{\tau} + 1/n \text{ tr}(\boldsymbol{\tau}) \text{ I}_{n \times n}$. Dies ist lediglich in die Definition der L_2 -Norm aus Abschnitt 2.3 einzusetzen. \square

2.3 Funktionenräume

In diesem Paragraphen werden die benötigten Funktionenräume im Überblick dargestellt, um die Notationen zu klären.

Der Lebesgue-Raum $L^2(\Omega; \mathbb{R}^n)$ wird mit der von dem L^2 -Skalarprodukt $(u, v)_{L^2(\Omega)} := \int_{\Omega} u \cdot v \, dx$ induzierten Norm $\|\bullet\|_{L^2(\Omega)}$ ausgestattet, wobei $u, v \in L^2(\Omega; \mathbb{R}^n)$, d.h. messbar

mit $\|v\|_{L^2(\Omega)}, \|u\|_{L^2(\Omega)} < \infty$. Bei dem matrixwertigen Äquivalent $L^2(\Omega, \mathbb{R}^{n \times n})$ lautet das entsprechende L^2 -Skalarprodukt $(\mathbf{u}, \mathbf{v})_{L^2(\Omega)} := \int_{\Omega} \mathbf{u} : \mathbf{v} \, dx$ für $\mathbf{u}, \mathbf{v} \in L^2(\Omega, \mathbb{R}^{n \times n})$. In dieser Arbeit spielen auch folgende Quotientenräume eine Rolle,

$$Q/\mathbb{R} := \left\{ \boldsymbol{\tau} \in Q : \int_{\Omega} \text{tr } \boldsymbol{\tau} \, dx = 0 \right\}, \quad (2.4)$$

wobei ein Unterraum $Q \subseteq L^2(\Omega; \mathbb{R}^{n \times n})$ betrachtet wird. Es wird also zusätzlich eine lineare Nebenbedingung verlangt.

Für $k \in \mathbb{N}$ sei $C^k(\Omega)$ der Raum der k -fach stetig differenzierbaren Funktionen auf Ω , $C^\infty(\Omega) := \bigcap_{k \in \mathbb{N}} C^k(\Omega)$ und $C_0^\infty := \{f \in C^\infty(\Omega) : \text{supp}(f) \subset \Omega \text{ kompakt}\}$ der Raum der Testfunktionen.

Eine Zusammenfassung der Differentialoperatoren scheint ebenfalls angebracht. Für $\partial/\partial x_k$ für $1 \leq k \leq n$ wird die Kurzschreibweise ∂_k verwendet. Für $u \in C^1(\Omega)$ sei $\nabla u := (\partial_1 u, \dots, \partial_n u)^\top \in \mathbb{R}^n$ der Gradient und für $u \in C^1(\Omega; \mathbb{R}^n)$ sei $D u \in \mathbb{R}^{n \times n}$ mit $(D u)_{jk} := \partial_k u_j$ für $1 \leq j, k \leq n$ die Jakobi-Matrix. Die Divergenz sei $\text{div } u := \sum_{j=1}^n \partial_j u_j = \text{tr}(D u) \in \mathbb{R}$ für $u \in C^1(\Omega; \mathbb{R}^n)$. Für matrix-wertige Funktionen $\mathbf{u} \in C^1(\Omega; \mathbb{R}^{n \times n})$ wird die Divergenz zeilenweise angewendet, d.h. $\text{div } \mathbf{u} := (\text{div } u_1, \dots, \text{div } u_n)^\top \in \mathbb{R}^n$, wobei u_j für $1 \leq j \leq n$ die j -te Zeile der Matrix $\mathbf{u} \in C^1(\Omega; \mathbb{R}^{n \times n})$ beschreibt.

Um den Sobolevraum $H^1(\Omega; \mathbb{R}^n)$ zu definieren wird der Begriff der schwachen Ableitung benötigt [Wer11, S. 210]. Für $1 \leq j \leq n$ sei die j -te schwache Ableitung von $u \in L^2(\Omega)$ eine Funktion $\tilde{u} \in L^2(\Omega)$, so dass für alle Testfunktionen $\varphi \in C_0^\infty(\Omega)$ gilt $\int_{\Omega} u \, \partial_j \varphi \, dx = - \int_{\Omega} \tilde{u} \varphi \, dx$. Diese schwache Ableitung \tilde{u} wird mit $\partial_j u$ bezeichnet. Im mehrdimensionalen Fall $u \in L^2(\Omega; \mathbb{R}^n)$ ist die Existenz der j -ten schwachen Ableitung von u_k für alle $1 \leq j, k \leq n$ äquivalent zu der Gleichheit $\int_{\Omega} u \cdot \text{div } \boldsymbol{\varphi} \, dx = - \int_{\Omega} D u : \boldsymbol{\varphi} \, dx$ für alle Testfunktionen $\boldsymbol{\varphi} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$, $D u$ bezeichnet dann die schwache Ableitung. Es sei

$$H^1(\Omega; \mathbb{R}^n) := \left\{ u \in L^2(\Omega; \mathbb{R}^n) : \text{im schwachen Sinne existiert } D u \in L^2(\Omega; \mathbb{R}^{n \times n}) \right\}.$$

Dieser Sobolevraum wird mit der H^1 -Norm

$$\|u\|_{H^1(\Omega)}^2 := \|u\|_{L^2(\Omega)}^2 + \|D u\|_{L^2(\Omega)}^2$$

versehen. Wie die Friedrichs Ungleichung, Lemma 2.16, zeigt, ist die Energienorm

$$\|u\| := \|D u\|_{L^2(\Omega)}$$

dazu äquivalent. Dirichlet Randdaten werden in diesen Raum mit Hilfe des Spuoperators $\gamma_0 : H^1(\Omega; \mathbb{R}^n) \rightarrow L^2(\partial\Omega; \mathbb{R}^n)$ mit $\gamma_0 u = u|_{\partial\Omega}$ (siehe Definition 2.7) einbezogen

$$H_0^1(\Omega; \mathbb{R}^n) := \left\{ u \in H^1(\Omega; \mathbb{R}^n) : \gamma_0(u) = 0 \text{ fast überall auf } \Gamma \right\}.$$

Es ist bekannt, dass $H_0^1(\Omega; \mathbb{R}^n)$ mit $\|\cdot\|_{H^1(\Omega)}$ bzw. $\|\cdot\|$ ein Banachraum ist [Alt06, S.68].

Um die Räume $H(\operatorname{div}, \Omega)$ und $H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ zu definieren, wird der Begriff der schwachen Divergenz benötigt. Die schwache Divergenz der Funktion $u \in L^2(\Omega; \mathbb{R}^n)$ sei $\hat{u} : \Omega \rightarrow \mathbb{R}$ bezeichnet als $\operatorname{div} u$, wenn $\int_{\Omega} \varphi \hat{u} \, dx = - \int_{\Omega} \nabla \varphi \cdot u \, dx$ für alle Testfunktionen $\varphi \in C_0^\infty(\Omega)$. Definiere damit

$$\begin{aligned} H(\operatorname{div}, \Omega) &:= \left\{ u \in L^2(\Omega; \mathbb{R}^n) : \text{im schwachen Sinne existiert } \operatorname{div} u \in L^2(\Omega) \right\}, \\ H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) &:= \left\{ \mathbf{u} \in L^2(\Omega; \mathbb{R}^{n \times n}) : u_j \in H(\operatorname{div}, \Omega) \, \forall 1 \leq j \leq n \right\}, \end{aligned}$$

wobei $u_j \in L^2(\Omega; \mathbb{R}^n)$ wieder die j -te Zeile der Matrix $\mathbf{u} \in L^2(\Omega; \mathbb{R}^{n \times n})$ sei. Auch im schwachen Sinne wird also der Divergenzoperator auf Matrizen zeilenweise angewendet. Diese Räume werden ausgestattet mit der Norm

$$\|u\|_{H(\operatorname{div})}^2 := \|u\|_{L^2(\Omega)}^2 + \|\operatorname{div} u\|_{L^2(\Omega)}^2$$

für $u \in H(\operatorname{div}, \Omega)$ bzw. $u \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ zu Banachräumen [Bra13, S.139].

Die in dieser Arbeit verwendeten Testräume erfüllen die Regularitätsbedingungen nur stückweise auf den einzelnen Elementen $T \in \mathcal{T}$ der regulären Triangulierung \mathcal{T} . Es seien

$$\begin{aligned} H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) &:= \left\{ \boldsymbol{\tau} \in L^2(\Omega; \mathbb{R}^{n \times n}) : \forall T \in \mathcal{T} \, \boldsymbol{\tau}|_T \in H(\operatorname{div}, T; \mathbb{R}^{n \times n}) \right\}, \\ H^1(\mathcal{T}; \mathbb{R}^n) &:= \left\{ v \in L^2(\Omega; \mathbb{R}^n) : \forall T \in \mathcal{T} \, v|_T \in H^1(T; \mathbb{R}^n) \right\}. \end{aligned}$$

Für $\boldsymbol{\tau} \in H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})$ bzw. $v \in H^1(\mathcal{T}; \mathbb{R}^n)$ werden die stückweise Divergenz mit $\operatorname{div}_{\text{NC}}$ bzw. die schwachen Ableitungen mit \mathbf{D}_{NC} bezeichnet, d.h.

$$\begin{aligned} \operatorname{div}_{\text{NC}} \boldsymbol{\tau}|_T &:= \operatorname{div}(\boldsymbol{\tau}|_T) \text{ für alle } T \in \mathcal{T}, \\ \mathbf{D}_{\text{NC}} v|_T &:= \mathbf{D}(v|_T) \text{ für alle } T \in \mathcal{T}. \end{aligned}$$

Damit ergeben sich die Normen für $H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})$ und $H^1(\mathcal{T}; \mathbb{R}^n)$, als

$$\begin{aligned} \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T})}^2 &:= \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})}^2 := \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\operatorname{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)}^2, \\ \|v\|_{H^1(\mathcal{T})}^2 &:= \|v\|_{H^1(\mathcal{T}; \mathbb{R}^n)}^2 := \|v\|_{L^2(\Omega)}^2 + \|\mathbf{D}_{\text{NC}} v\|_{L^2(\Omega)}^2 := \|v\|_{L^2(\Omega)}^2 + \|v\|_{\text{NC}}^2. \end{aligned}$$

Zusätzlich werden noch Spurräume benötigt. In [Alt06, S.265] und [GR86, S.27 f.] werden die Spursätze bewiesen, die die Definitionen der Spurooperatoren ermöglichen.

Satz 2.6. *Es sei $U \subseteq \mathbb{R}^n$ ein offenes, beschränktes Lipschitz-Gebiet. Dann existiert*

genau eine stetige, lineare Abbildung $\gamma_0 : H^1(U) \rightarrow L^2(\partial U)$, so dass

$$\gamma_0 w = w|_{\partial U}, \quad \text{für alle } w \in H^1(U) \cap C^0(\bar{U}).$$

Definition 2.7. Der Bildraum dieses Spuoperators wird mit $H^{1/2}(\partial U) := \gamma_0(H^1(U))$, sein Dualraum mit $H^{-1/2}(\partial U) = (H^{1/2}(\partial U))^*$ bezeichnet.

Satz 2.8. Es sei $U \subseteq \mathbb{R}^n$ ein offenes, beschränktes Lipschitz-Gebiet. Dann existiert genau eine stetige, lineare Abbildung $\gamma_\nu : H(\text{div}, U) \rightarrow H^{-1/2}(\partial U)$, so dass

$$\gamma_\nu q = (q|_{\partial U}) \cdot \nu, \quad \text{für alle } q \in H(\text{div}, U).$$

Für alle $q \in H(\text{div}, U)$ und $w \in H^1(U)$ gilt

$$\langle \gamma_\nu q, \gamma_0 w \rangle_{\partial U} = \int_U q \nabla w \, dx + \int_U \text{div}(q) w \, dx.$$

Eine komponentenweise Anwendung dieser Abbildung erlaubt die Definition der mehrdimensionalen Spuoperatoren

$$\begin{aligned} \gamma_0 : H^1(U; \mathbb{R}^n) &\rightarrow H^{1/2}(\partial U; \mathbb{R}^n), \\ \gamma_\nu : H(\text{div}, U; \mathbb{R}^{n \times n}) &\rightarrow H^{-1/2}(\partial U; \mathbb{R}^n). \end{aligned}$$

Sei \mathcal{T} eine regulären Triangulierung von Ω , so können Spuabbildung auf das Skelett $\partial \mathcal{T}$ definiert werden. Die Funktion $\gamma_0^\mathcal{T} : H^1(\mathcal{T}; \mathbb{R}^n) \rightarrow \prod_{T \in \mathcal{T}} H^{1/2}(\partial T; \mathbb{R}^n)$ bildet $w \in H^1(\mathcal{T}; \mathbb{R}^n)$ auf $\gamma_0^\mathcal{T} w := (s_T)_{T \in \mathcal{T}}$ mit $s_T := \gamma_0(w|_T)$ für alle $T \in \mathcal{T}$ ab. Die Funktion $\gamma_\nu^\mathcal{T} : H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n}) \rightarrow \prod_{T \in \mathcal{T}} H^{-1/2}(\partial T; \mathbb{R}^n)$ bildet $\mathbf{q} \in H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n})$ auf $\gamma_\nu^\mathcal{T} \mathbf{q} := (t_T)_{T \in \mathcal{T}}$ mit $t_T := \gamma_\nu(\mathbf{q}|_T)$ für alle $T \in \mathcal{T}$ ab. Die assoziierten Spurräume werden wie folgt definiert.

$$H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n) := \gamma_0^\mathcal{T}(H_0^1(\Omega; \mathbb{R}^n)) \tag{2.5}$$

$$H^{-1/2}(\partial \mathcal{T}; \mathbb{R}^n) := \gamma_\nu^\mathcal{T}(H(\text{div}, \Omega; \mathbb{R}^{n \times n})) \tag{2.6}$$

Diese Räume werden mit den folgenden Normen ausgestattet

$$\begin{aligned} \|s\|_{H_0^{1/2}(\partial \mathcal{T})} &:= \|s\|_{H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)} := \inf_{\substack{w \in H_0^1(\Omega; \mathbb{R}^n) \\ \gamma_0^\mathcal{T} w = s}} \|w\|_{H^1(\Omega)}, \\ \|t\|_{H^{-1/2}(\partial \mathcal{T})} &:= \|t\|_{H^{-1/2}(\partial \mathcal{T}; \mathbb{R}^n)} := \inf_{\substack{\mathbf{q} \in H(\text{div}, \Omega; \mathbb{R}^{n \times n}) \\ \gamma_\nu^\mathcal{T} \mathbf{q} = t}} \|\mathbf{q}\|_{H(\text{div}, \Omega)}. \end{aligned}$$

Definition 2.9. Die Dualitätsklammer auf dem Skelett sei für alle $t = (t_T)_{T \in \mathcal{T}} \in \prod_{T \in \mathcal{T}} H^{-1/2}(T; \mathbb{R}^n)$ und $s = (s_T)_{T \in \mathcal{T}} \in \prod_{T \in \mathcal{T}} H^{1/2}(\partial T; \mathbb{R}^n)$ definiert durch

$$\langle t, s \rangle_{\partial \mathcal{T}} := \sum_{T \in \mathcal{T}} \langle t_T, s_T \rangle_{\partial T}.$$

Bemerkung 2.10 Sei $v \in H_0^1(\mathcal{T}; \mathbb{R}^n)$ mit $\gamma_0^\mathcal{T} v \in H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ so gilt bereits $v \in H_0^1(\Omega; \mathbb{R}^n)$. Für den Nachweis sei $w \in H_0^1(\Omega; \mathbb{R}^n)$ mit $\gamma_0^\mathcal{T} w = \gamma_0^\mathcal{T} v$ und $\varphi \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$ eine beliebige Testfunktion, dann gilt die Formel aus Satz 2.8 für jedes $T \in \mathcal{T}$ und damit

$$\begin{aligned} \int_{\Omega} D_{\text{NC}} v : \varphi \, dx &= \langle \gamma_\nu^\mathcal{T} \varphi, \gamma_0^\mathcal{T} v \rangle_{\partial \mathcal{T}} - \int_{\Omega} \operatorname{div} \varphi \cdot v \, dx \\ &= \langle \gamma_\nu^\mathcal{T} \varphi, \gamma_0^\mathcal{T} w \rangle_{\partial \mathcal{T}} - \int_{\Omega} \operatorname{div} \varphi \cdot v \, dx \\ &= \int_{\Omega} D w : \varphi \, dx + \int_{\Omega} \operatorname{div} \varphi \cdot w \, dx - \int_{\Omega} \operatorname{div} \varphi \cdot v \, dx \\ &= \int_{\partial \Omega} w \varphi \cdot \nu \, dx - \int_{\Omega} \operatorname{div} \varphi \cdot v \, dx = - \int_{\Omega} \operatorname{div} \varphi \cdot v \, dx. \end{aligned}$$

Somit gilt nach Definition der schwachen Ableitung, dass $D_{\text{NC}} v = D v$. Analog folgt für $\tau \in H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})$ mit $\gamma_\nu^\mathcal{T} \tau \in H^{-1/2}(\partial \mathcal{T}; \mathbb{R}^n)$, dass bereits $\tau \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ gilt. Sei dazu $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ mit $\gamma_\nu^\mathcal{T} \mathbf{q} = \gamma_\nu^\mathcal{T} \tau$ und $\varphi \in C_0^\infty(\Omega; \mathbb{R}^n)$ eine beliebige Testfunktion, so gilt

$$\begin{aligned} \int_{\Omega} \operatorname{div}_{\text{NC}} \tau \cdot \varphi \, dx &= \langle \gamma_\nu^\mathcal{T} \tau, \gamma_0^\mathcal{T} \varphi \rangle_{\partial \mathcal{T}} - \int_{\Omega} D \varphi : \tau \, dx \\ &= \langle \gamma_\nu^\mathcal{T} \mathbf{q}, \gamma_0^\mathcal{T} \varphi \rangle_{\partial \mathcal{T}} - \int_{\Omega} D \varphi : \tau \, dx \\ &= \int_{\Omega} D \varphi : \mathbf{q} \, dx + \int_{\Omega} \operatorname{div} \mathbf{q} \cdot \varphi \, dx - \int_{\Omega} D \varphi : \tau \, dx \\ &= \int_{\partial \Omega} \varphi \mathbf{q} \cdot \nu \, dx - \int_{\Omega} D \varphi : \tau \, dx = - \int_{\Omega} D \varphi : \tau \, dx. \end{aligned}$$

Also gilt nach Definition der schwachen Ableitung, dass $\operatorname{div}_{\text{NC}} \tau = \operatorname{div} \tau$.

Bemerkung 2.11 Mit Hilfe der vorangehenden Bemerkung 2.10 wird umgehend deutlich, dass die Normen $\|\cdot\|_{H_0^{1/2}(\partial \mathcal{T})}$ bzw. $\|\cdot\|_{H^{-1/2}(\partial \mathcal{T})}$ Minimumsnormen sind.

Sei $s \in H_0^{1/2}(\Omega; \mathbb{R}^n)$, so ist die Funktion $w \in H_0^1(\Omega; \mathbb{R}^n)$ mit $\gamma_0^\mathcal{T} w = s$ und $\|s\|_{H_0^{1/2}(\partial \mathcal{T})} = \|w\|_{H^1(\Omega)}$ nach den Euler-Lagrange Gleichungen für jedes $T \in \mathcal{T}$ Lösung der folgenden Differentialgleichung in $H^1(T; \mathbb{R}^n)$

$$\begin{aligned} w &= s \quad \text{entlang } \partial T, \\ -\operatorname{div}(D w) + w &= 0 \quad \text{in } T. \end{aligned}$$

Für $t \in H^{-1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ ist die Funktion $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ mit $\gamma_\nu^\mathcal{T} \mathbf{q} = t$ und $\|t\|_{H^{-1/2}(\partial \mathcal{T})} =$

$\|\mathbf{q}\|_{H(\text{div}, \Omega)}$ auf jedem $T \in \mathcal{T}$ Lösung der folgenden Differentialgleichung in $H(\text{div}, T; \mathbb{R}^{n \times n})$

$$\begin{aligned} \mathbf{q} \cdot \nu &= t \quad \text{entlang } \partial T, \\ -\text{D}(\text{div } \mathbf{q}) + \mathbf{q} &= 0 \quad \text{in } T. \end{aligned}$$

Somit ist es stets möglich eine Fortsetzung $w \in H_0^1(\Omega; \mathbb{R}^n)$ bzw. $\mathbf{q} \in H(\text{div}, \Omega; \mathbb{R}^{n \times n})$ mit derselben Norm zu finden.

Elementweise sind $H^{1/2}(\partial T)$ und $H^{-1/2}(\partial T)$ nach Definition 2.7 dual. Dies ist nicht auf die Räume $H^{1/2}(\partial \mathcal{T})$ und $H^{-1/2}(\partial \mathcal{T})$ übertragbar. Es können allerdings die folgenden nützlichen Dualitätsaussagen wie in [CDG15b, Lemma 3.1] bzw. [CDG15a, Theorem 2.3] bewiesen werden.

Lemma 2.12 (Dualitätslemma I). *Für alle $s \in H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ gilt*

$$\|s\|_{H^{1/2}(\partial \mathcal{T})} = \sup_{\substack{\boldsymbol{\tau} \in H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R} \\ \boldsymbol{\tau} \neq 0}} \frac{\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, s \rangle_{\partial \mathcal{T}}}{\|\boldsymbol{\tau}\|_{H(\text{div}, \mathcal{T})}}.$$

Beweis. Sei zunächst $w \in H_0^1(\Omega; \mathbb{R}^n)$ die Fortsetzung von s mit derselben Norm aus Bemerkung 2.11. Setze $\mathbf{q} := \text{D } w$, dann gilt nach obiger Bemerkung $\text{div } \mathbf{q}|_T = \text{div}(\text{D } w|_T) = w|_T$ auf jedem Element $T \in \mathcal{T}$ und

$$\int_\Omega \text{tr } \mathbf{q} \, dx = \int_\Omega \text{div } w \, dx = \int_{\partial \Omega} w \cdot \nu \, ds = 0,$$

also $\mathbf{q} \in H(\text{div}, \mathcal{T}; \mathbb{R}^n)/\mathbb{R}$. Weiter ergibt sich nach der Wahl von \mathbf{q}

$$\|\mathbf{q}\|_{H(\text{div}, \mathcal{T})}^2 = \|\mathbf{q}\|_{L^2(\Omega)}^2 + \|\text{div}_{\text{NC}} \mathbf{q}\|_{L^2(\Omega)}^2 = \|w\|_{H^1(\Omega)}^2.$$

Nun bedeutet Satz 2.8

$$\langle \gamma_\nu^\mathcal{T} \mathbf{q}, s \rangle_{\partial \mathcal{T}} = \langle \gamma_\nu^\mathcal{T} \mathbf{q}, \gamma_0^\mathcal{T} w \rangle_{\partial \mathcal{T}} = \int_\Omega \mathbf{q} : \text{D } w \, dx + \int_\Omega \text{div}_{\text{NC}} \mathbf{q} \cdot w \, dx = \|w\|_{H^1(\Omega)}^2.$$

Dies ist die eine Richtung der behaupteten Gleichheit, denn

$$\|s\|_{H^{1/2}(\partial \mathcal{T})} = \|w\|_{H^1(\Omega)} = \frac{\langle \gamma_\nu^\mathcal{T} \mathbf{q}, s \rangle_{\partial \mathcal{T}}}{\|\mathbf{q}\|_{H(\text{div}, \mathcal{T})}} \leq \sup_{\substack{\boldsymbol{\tau} \in H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R} \\ \boldsymbol{\tau} \neq 0}} \frac{\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, s \rangle_{\partial \mathcal{T}}}{\|\boldsymbol{\tau}\|_{H(\text{div}, \mathcal{T})}}.$$

Für die andere Richtung betrachte außer der minimalen Fortsetzung $w \in H_0^1(\Omega; \mathbb{R}^n)$ ein beliebiges $\mathbf{q} \in S(H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R})$. Dann zeigen partielle Integration und Cauchy-Schwarz-Ungleichung auf jedem Element $T \in \mathcal{T}$

$$\langle \gamma_0 w, \gamma_\nu \mathbf{q}|_T \rangle_{\partial T} = \int_T \mathbf{q} : \text{D } v \, dx + \int_T v \cdot \text{div } \mathbf{q} \, dx \leq \|v\|_{H^1(T)} \|\mathbf{q}\|_{H(\text{div}, T)} \leq \|v\|_{H^1(T)},$$

da \mathbf{q} normiert ist. In der Summe ergibt sich damit

$$\left\langle \gamma_0^\mathcal{T} w, \gamma_\nu^\mathcal{T} \mathbf{q} \right\rangle_{\partial\mathcal{T}} \leq \sum_{T \in \mathcal{T}} \|w\|_{H^1(T)} = \|w\|_{H^1(\Omega)}.$$

Daraus folgt die Behauptung, da $\mathbf{q} \in S(H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R})$ beliebig war. \square

Lemma 2.13 (Dualitätslemma II). Für alle $t \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ gilt

$$\|t\|_{H^{-1/2}(\partial\mathcal{T})} = \sup_{v \in H^1(\mathcal{T}; \mathbb{R}^n) \setminus \{0\}} \frac{\left\langle t, \gamma_0^\mathcal{T} v \right\rangle_{\partial\mathcal{T}}}{\|v\|_{H^1(\mathcal{T})}}.$$

Beweis. Sei zunächst $w \in H^1(\mathcal{T}; \mathbb{R}^n)$ auf jedem Element $T \in \mathcal{T}$ die Lösung des folgenden Neumann Problems

$$\operatorname{div} \operatorname{D} w - w = 0 \text{ in } T \quad \text{und} \quad \operatorname{D} w \cdot \nu_T = t_T \text{ auf } \partial T.$$

Setze $\mathbf{q} := \operatorname{D}_{\text{NC}} w \in H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})$, dann gilt $\gamma_\nu^\mathcal{T} \mathbf{q} = t \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ also nach Bemerkung 2.10 $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$. Außerdem ist $\operatorname{div} \mathbf{q} = w$ und daher ergibt sich

$$\|\mathbf{q}\|_{H(\operatorname{div})}^2 = \|\mathbf{q}\|_{L^2(\Omega)}^2 + \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)}^2 = \|\operatorname{D}_{\text{NC}} w\|_{L^2(\Omega)}^2 + \|w\|_{L^2(\Omega)}^2 = \|w\|_{H^1(\mathcal{T})}^2.$$

Nun bedeutet Satz 2.8

$$\left\langle t, \gamma_0^\mathcal{T} w \right\rangle_{\partial\mathcal{T}} = \left\langle \gamma_\nu^\mathcal{T} \mathbf{q}, \gamma_0^\mathcal{T} w \right\rangle_{\partial\mathcal{T}} = \int_{\Omega} \mathbf{q} : \operatorname{D}_{\text{NC}} w \, dx + \int_{\Omega} \operatorname{div} \mathbf{q} \cdot w \, dx = \|\mathbf{q}\|_{H(\operatorname{div})}^2.$$

Damit ergibt sich die erste Richtung

$$\begin{aligned} \|t\|_{H^{-1/2}(\partial\mathcal{T})} &= \inf_{\substack{\boldsymbol{\tau} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) \\ \gamma_\nu^\mathcal{T} \boldsymbol{\tau} = t}} \|\boldsymbol{\tau}\|_{H(\operatorname{div})} \leq \|\mathbf{q}\|_{H(\operatorname{div})} = \frac{\left\langle t, \gamma_0^\mathcal{T} w \right\rangle_{\partial\mathcal{T}}}{\|w\|_{H^1(\mathcal{T})}} \\ &\leq \sup_{v \in H^1(\mathcal{T}; \mathbb{R}^n) \setminus \{0\}} \frac{\left\langle t, \gamma_0^\mathcal{T} v \right\rangle_{\partial\mathcal{T}}}{\|v\|_{H^1(\mathcal{T})}}. \end{aligned}$$

Der Beweis der anderen Ungleichung funktioniert analog zu Lemma 2.12. Diesmal wird allerdings zu dem $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ mit $\gamma_\nu^\mathcal{T} \mathbf{q} = t$ und minimaler Norm aus Bemerkung 2.11 ein beliebiges $w \in S(H_0^1(\mathcal{T}; \mathbb{R}^n))$ gewählt. Damit ergibt sich mit elementweise partieller Integration und Cauchy-Schwarz Ungleichung diesmal

$$\left\langle \gamma_0^\mathcal{T} w, \gamma_\nu^\mathcal{T} \mathbf{q} \right\rangle_{\partial\mathcal{T}} = \sum_{T \in \mathcal{T}} \left\langle \gamma_0 w, \gamma_\nu \mathbf{q} \right\rangle_{\partial T} \leq \sum_{T \in \mathcal{T}} \|w\|_{H^1(T)} \|\mathbf{q}\|_{H(\operatorname{div}, T)} \leq \|\mathbf{q}\|_{H(\operatorname{div}, \Omega)}.$$

Da dies für jedes $w \in S(H_0^1(\mathcal{T}; \mathbb{R}^n))$ gilt, folgt der zweite Teil der Behauptung. \square

Mit den Definitionen aus diesem Paragraphen können nun die in der dPG-Formulierung des Stokes Problems auftretenden Produkträume X und Y eingeführt werden. Der

Ansatzraum X , in dem die Lösung der Differentialgleichung gesucht wird, lautet

$$X := L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n) \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n), \quad (2.7)$$

und sei mit folgender Norm ausgestattet

$$\|x\|_X^2 := \|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 + \|s\|_{H_0^{1/2}(\partial\mathcal{T})}^2 + \|t\|_{H^{-1/2}(\partial\mathcal{T})}^2 \quad \text{für jedes } x = (\boldsymbol{\sigma}, u, s, t) \in X.$$

Der Testraum Y , mit dessen Elementen die Variationsformulierung der Bilinearform ausgewertet wird, liest sich

$$Y := H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times H^1(\mathcal{T}; \mathbb{R}^n) \quad (2.8)$$

und sei mit folgender Norm ausgestattet

$$\|y\|_Y^2 := \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T})}^2 + \|v\|_{H^1(\mathcal{T})}^2 \quad \text{für jedes } y = (\boldsymbol{\tau}, v) \in Y.$$

Damit liegt den Räumen X und Y bereits auf kontinuierlichem Level eine Triangulierung zu Grunde.

Bemerkung 2.14 Sowohl der Ansatzraum als auch der Testraum sind Hilberträume. Die Norm $\|\cdot\|_Y$ erfüllt die Parallelogrammgleichung

$$\|y + \hat{y}\|_Y^2 + \|y - \hat{y}\|_Y^2 = 2(\|y\|_Y^2 + \|\hat{y}\|_Y^2).$$

Damit lässt sich mittels

$$\begin{aligned} \langle y, \hat{y} \rangle_Y &:= \frac{1}{4} (\|y + \hat{y}\|_Y^2 - \|y - \hat{y}\|_Y^2) \\ &= \int_{\Omega} \boldsymbol{\tau} : \hat{\boldsymbol{\tau}} \, dx + \int_{\Omega} \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot \operatorname{div}_{\text{NC}} \hat{\boldsymbol{\tau}} \, dx + \int_{\Omega} v \cdot \hat{v} \, dx + \int_{\Omega} \operatorname{D}_{\text{NC}} v : \operatorname{D}_{\text{NC}} \hat{v} \, dx \end{aligned}$$

ein Skalarprodukt für alle $y = (\boldsymbol{\tau}, v)$, $\hat{y} = (\hat{\boldsymbol{\tau}}, \hat{v}) \in Y$ definieren [Wer11, S.205].

2.4 Diskrete Funktionenräume

In diesem Paragraphen werden die in der Diskretisierung verwendeten Räume vorgestellt. Die endlich-dimensionalen Unterräume des Ansatzraumes $X_h \subset X$ und Testraumes $Y_h \subset Y$ sind Räume stückweise polynomialer Funktionen.

Zunächst seien für $\mathbb{M} \subseteq \mathbb{R}^m$ oder $\mathbb{M} \subseteq \mathbb{R}^{m \times n}$ und $k, m, n \in \mathbb{N}_{>0}$,

$$\begin{aligned} P_k(T) &:= \{p_k \in L^\infty(T) : p_k \text{ ist auf } T \text{ ein Polynom vom Grad } \leq k\}, \\ P_k(\mathcal{T}) &:= \{p_k \in L^\infty(\Omega) : \forall T \in \mathcal{T}, p_k|_T \in P_k(T)\}, \\ P_k(T; \mathbb{M}) &:= \{\mathbf{q}_k \in L^\infty(T; \mathbb{M}) : \text{jede Komponente von } \mathbf{q}_k \text{ gehört zu } P_k(T)\}, \\ P_k(\mathcal{T}, \mathbb{M}) &:= \{\mathbf{q}_k \in L^\infty(\Omega; \mathbb{M}) : \forall T \in \mathcal{T}, \mathbf{q}_k|_T \in P_k(T; \mathbb{M})\}. \end{aligned}$$

Auf dem Skelett seien analog definiert

$$\begin{aligned} P_k(E) &:= \{p_k \in L^\infty(E) : p_k|_E \text{ ist ein Polynom vom Grad } \leq k\}, \\ P_k(\mathcal{E}) &:= \left\{p_k \in L^\infty\left(\bigcup \mathcal{E}\right) : \forall E \in \mathcal{E}, p_k|_E \in P_k(E)\right\}, \\ P_k(E; \mathbb{M}) &:= \{\mathbf{q}_k \in L^\infty(E; \mathbb{M}) : \text{jede Komponente von } \mathbf{q}_k \text{ gehört zu } P_k(E)\}, \\ P_k(\mathcal{E}; \mathbb{M}) &:= \left\{\mathbf{q}_k \in L^\infty\left(\bigcup \mathcal{E}; \mathbb{M}\right) : \forall E \in \mathcal{E}, \mathbf{q}_k|_E \in P_k(E; \mathbb{M})\right\}. \end{aligned}$$

Die stetigen und stückweise P_k -Finite-Elemente-Funktionen auf \mathcal{T} werden mit

$$S^k(\mathcal{T}; \mathbb{M}) := P_k(\mathcal{T}) \cap C(\bar{\Omega}) \text{ und } S_0^k(\mathcal{T}; \mathbb{M}) := S^k(\mathcal{T}; \mathbb{M}) \cap C_0(\Omega)$$

und auf dem Skelett mit

$$S_0^k(\mathcal{E}; \mathbb{M}) := P_k(\mathcal{E}; \mathbb{M}) \cap C_0\left(\bigcup \mathcal{E}; \mathbb{M}\right)$$

bezeichnet.

Die (stückweise) Raviart-Thomas Funktionen niedrigster Ordnung sind definiert als

$$\begin{aligned} RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) &:= \left\{ \mathbf{q}_{\text{RT}} \in L^\infty(\Omega; \mathbb{R}^{n \times n}) : \exists A \in P_0(\mathcal{T}; \mathbb{R}^{n \times n}), \exists b \in P_0(\mathcal{T}; \mathbb{R}^n), \right. \\ &\quad \left. \mathbf{q}_{\text{RT}} = A + b \otimes (\cdot - \text{mid}(\mathcal{T})) \right\}, \\ RT_0(\mathcal{T}; \mathbb{R}^{n \times n}) &:= RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) \cap H(\text{div}, \Omega; \mathbb{R}^{n \times n}). \end{aligned}$$

Es kann jede Raviart-Thomas Funktion $\mathbf{q}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n})$ auf jedem Simplex $T \in \mathcal{T}$ geschrieben werden als $\mathbf{q}_{\text{RT}}|_T = A + 1/n \text{ div } \mathbf{q}_{\text{RT}} \otimes (\cdot - \text{mid}(T))$ für ein $A \in \mathbb{R}^{n \times n}$. Wenn also Π_0 die P_0 -Projektion von $f \in L^2(\Omega; \mathbb{M})$ auf $P_0(\mathcal{T})$,

$$\Pi_0 f|_T := |T|^{-1} \int_T f \, dx = \int_T f \, dx,$$

bezeichnet, dann gilt nach (2.2)

$$(1 - \Pi_0) \mathbf{q}_{\text{RT}} = 1/n \text{ div } \mathbf{q}_{\text{RT}} \otimes (\cdot - \text{mid}(\mathcal{T})) \perp P_0(\mathcal{T}). \quad (2.9)$$

Außerdem gilt für alle $f, g \in L^2(\Omega)$

$$\int_\Omega f \Pi_0 g \, dx = \int_\Omega \Pi_0 f \Pi_0 g \, dx, \quad (2.10)$$

da auf jedem Element $\int_T f \Pi_0 g \, dx = \Pi_0 g \int_T f \, dx = \Pi_0 g |T| \int_T f \, dx = \int_T \Pi_0 f \Pi_0 g \, dx$.

Die Raviart-Thomas Funktionen ermöglichen außerdem eine Einbettung von $P_0(\mathcal{E}; \mathbb{R}^n)$ in $H^{-1/2}(\partial T; \mathbb{R}^n)$.

Bemerkung 2.15 Es gilt für jedes $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ und für alle $E \in \mathcal{E}$, $\mathbf{q}_{\text{RT}}|_E \nu_E \in$

$P_0(E; \mathbb{R}^n)$. Umgekehrt, existiert zu jedem $t_0 \in P_0(\mathcal{E}; \mathbb{R}^n)$ genau ein $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$, so dass $\mathbf{q}_{\text{RT}}|_E \nu_E = t_0|_E$. Beides folgt direkt aus [BC05, Lemma 4.1].

Nun ist $P_0(\mathcal{E}; \mathbb{R}^n)$ ein Unterraum von $H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ vermöge der Einbettung

$$P_0(\mathcal{E}; \mathbb{R}^n) \hookrightarrow H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n), \quad t_0 \mapsto t = (t_T)_{T \in \mathcal{T}} \text{ mit } t_T = \mathbf{q}_{\text{RT}} \nu_T|_{\partial T},$$

wobei $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ die laut Bemerkung 2.15 eindeutige Funktion mit $\mathbf{q}_{\text{RT}}|_E \nu_E = t_0|_E$ auf allen Seiten $E \in \mathcal{E}$ ist. Offensichtlich gilt auf jedem Element $T \in \mathcal{T}$ und jeder Seite $E \in \mathcal{E}(T)$ die Gleichheit $(\mathbf{q}_{\text{RT}} \nu_T)|_E = \text{sgn}(T, E) t_0$.

Damit lassen sich nun für eine gegebene reguläre Triangulierung \mathcal{T} von Ω die endlich dimensionalen Unterräume der konformen Diskretisierung angeben

$$X_h := P_0(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_0(\mathcal{T}; \mathbb{R}^n) \times S_0^1(\mathcal{E}; \mathbb{R}^n) \times P_0(\mathcal{E}; \mathbb{R}^n), \quad (2.11)$$

$$Y_h := RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^n). \quad (2.12)$$

2.5 Weitere Hilfsmittel

In diesem Paragraphen werden einige nützliche Tatsachen vorgestellt.

Die folgenden Lemmata werden für verschiedene Normabschätzung in dieser Arbeit benötigt.

Lemma 2.16 (Friedrichs Ungleichung). *Für jedes beschränkte, polygonal berandete Lipschitzgebiet Ω existiert eine nur von Ω abhängige Konstante C_F , so dass für jedes $v \in H_0^1(\Omega; \mathbb{R}^n)$ gilt, dass*

$$\|v\|_{L^2(\Omega)} \leq C_F \|D v\|_{L^2(\Omega)}.$$

Die nur von Ω abhängige Friedrichskonstante kann auf verschiedenen Wegen abgeschätzt werden. Laut [Bra13, S.29] gilt, falls Ω in einem n -dimensionalen Würfel der Kantenlänge q enthalten ist, so gilt $C_F \leq q$. In [Car09b, Theorem 0.49] ist bewiesen, dass $C_F \leq \text{width}(\Omega)/\pi$, wobei $\text{width}(\Omega) := L := \beta - \alpha$ die kürzeste Länge ist, so dass Ω zwischen den parallelen Hyperebenen $\{x \cdot \nu = \alpha\}$ und $\{x \cdot \nu = \beta\}$ mit Abstand L für einen Einheitsvektor $\nu \in \mathbb{R}^n$ liegt.

Lemma 2.17 (Poincaré Ungleichung). *Für jedes beschränkte, konvexe Gebiet Ω und jedes $v \in H^1(\Omega; \mathbb{R}^n)$ mit $\bar{v} := f_\Omega v \, dx := |\Omega|^{-1} \int_\Omega v \, dx$ existiert eine nur von Ω abhängige Konstante C_P , so dass*

$$\|v - \bar{v}\|_{L^2(\Omega)} \leq C_P \|D v\|_{L^2(\Omega)}.$$

Auch diese Konstante kann genauer angegeben werden. Für beliebige Dimensionen ist in [PW60] bewiesen, dass $C_P \leq \text{diam}(\Omega)/\pi$, wobei $\text{diam}(\Omega) := \sup_{x,y \in \Omega} |x - y|$. Für den zweidimensionalen Fall ist in [LS10] eine optimale Konstante hergeleitet. Sei \mathcal{T} eine reguläre Triangulierung des beschränkten, konvexen Gebietes $\Omega \subseteq \mathbb{R}^2$. So gilt auf jeden Dreieck $T \in \mathcal{T}$ mit $0 < C_P(T) \leq h_T/j_{1,1}$, wobei $j_{1,1} > \pi$ die erste Wurzel der ersten Besselfunktion ist, dass

$$\left\| v - \oint_T v \, dx \right\|_{L^2(T)} \leq C_P(T) \|D v\|_{L^2(T)}.$$

So ergibt sich auf ganz Ω folgende Abschätzung für die Poincaré Konstante $0 < C_P \leq h_{\max}/j_{1,1}$.

Das folgende häufig verwendete Lemma ist aus [CGS14] entnommen.

Lemma 2.18 (Tr-Div-Dev). *Sei $\Sigma_0 \subseteq H(\text{div}, \Omega; \mathbb{R}^{n \times n})$ ein abgeschlossener Unterraum, der den konstanten Tensor $I_{n \times n}$ nicht enthält, dann existiert eine nur von Ω abhängige Konstante $C_{td} < \infty$, so dass für alle $\boldsymbol{\tau} \in \Sigma_0$ gilt*

$$\|\text{tr } \boldsymbol{\tau}\|_{L^2(\Omega)} \leq C_{td} (\|\text{dev } \boldsymbol{\tau}\|_{L^2(\Omega)} + \|\text{div } \boldsymbol{\tau}\|_{L^2(\Omega)}).$$

Bemerkung 2.19 In dieser Arbeit wird Lemma 2.18 lediglich auf den Raum $\Sigma_0 := \left\{ \boldsymbol{\tau} \in H(\text{div}, \Omega; \mathbb{R}^{n \times n}) : \int_{\Omega} \text{tr } \boldsymbol{\tau} \, dx = 0 \right\}$ angewendet. Es ist wie in [Hel14, S.27] leicht zu zeigen, dass dieses Σ_0 die in Lemma 2.18 geforderten Bedingungen erfüllt. Σ_0 enthält $I_{n \times n}$ offensichtlich nicht. Um außerdem die Abgeschlossenheit in $H(\text{div}, \Omega; \mathbb{R}^{n \times n})$ zu prüfen, sei $(\boldsymbol{\tau}^j)_{j \in \mathbb{N}} \in \Sigma_0$ eine konvergente Folge mit Grenzwert $\boldsymbol{\tau} \in H(\text{div}, \Omega; \mathbb{R}^{n \times n})$. Dann gilt wegen der Linearität der Spur, $\boldsymbol{\tau}^j \in \Sigma_0$ und Cauchy-Schwarz für $L^2(\Omega)$ für alle $j \in \mathbb{N}$

$$\begin{aligned} \left| \int_{\Omega} \text{tr } \boldsymbol{\tau} \, dx \right| &= \left| \int_{\Omega} \text{tr } (\boldsymbol{\tau} - \boldsymbol{\tau}^j) \, dx \right| \leq \int_{\Omega} |\text{tr } (\boldsymbol{\tau} - \boldsymbol{\tau}^j)| \, dx \\ &= \left\| \text{tr } (\boldsymbol{\tau} - \boldsymbol{\tau}^j) \right\|_{L^1(\Omega)} = \left\| \sum_{k=1}^n (\tau_{kk} - \tau_{kk}^j) \right\|_{L^1(\Omega)} \\ &\leq \sum_{k=1}^n \left\| (\tau_{kk} - \tau_{kk}^j) \right\|_{L^1(\Omega)} \leq |\Omega|^{1/2} \sum_{k=1}^n \left\| (\tau_{kk} - \tau_{kk}^j) \right\|_{L^2(\Omega)}. \end{aligned}$$

Die Cauchy-Schwarz-Ungleichung in \mathbb{R}^n erlaubt weiter

$$\begin{aligned} \sum_{k=1}^n \left\| (\tau_{kk} - \tau_{kk}^j) \right\|_{L^2(\Omega)} &= (1, \dots, 1) \cdot \left(\left\| (\tau_{11} - \tau_{11}^j) \right\|_{L^2(\Omega)}, \dots, \left\| (\tau_{nn} - \tau_{nn}^j) \right\|_{L^2(\Omega)} \right)^{\top} \\ &\leq \sqrt{n} \sqrt{\sum_{k=1}^n \left\| (\tau_{kk} - \tau_{kk}^j) \right\|_{L^2(\Omega)}^2}. \end{aligned}$$

Außerdem gilt nach der Definition des L_2 -Skalarprodukts für Funktionen nach $\mathbb{R}^{n \times n}$,

$$\sum_{k=1}^n \left\| \left(\tau_{kk} - \tau_{kk}^j \right) \right\|_{L^2(\Omega)}^2 \leq \sum_{k,l=1}^n \left\| \left(\tau_{kl} - \tau_{kl}^j \right) \right\|_{L^2(\Omega)}^2 = \left\| \boldsymbol{\tau} - \boldsymbol{\tau}^j \right\|_{L^2(\Omega)}^2 \rightarrow 0, \quad \text{für } j \rightarrow \infty.$$

Daher gilt $\int_{\Omega} \operatorname{tr} \boldsymbol{\tau} \, dx = 0$, also $\boldsymbol{\tau} \in \Sigma_0$.

In den Beweisen der folgenden Kapitel ist auch die im Folgende zu einer Raviart-Thomas Funktion \mathbf{q}_{RT} definierte Hilfsfunktion $\tilde{\mathbf{q}}_{\text{RT}}$ nützlich. Einige ihrer Eigenschaften werden hier aufgeführt.

Bemerkung 2.20 Sei $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ und definiere dazu

$$\tilde{\mathbf{q}}_{\text{RT}} := \mathbf{q}_{\text{RT}} - 1/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \mathbf{I}_{n \times n} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}).$$

Dann gilt offensichtlich $\int_{\Omega} \operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}} \, dx = 0$, d.h. Lemma 2.18 ist auf $\tilde{\mathbf{q}}_{\text{RT}}$ anwendbar. Außerdem gilt

- (i) $\|\operatorname{dev} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} = \|\operatorname{dev} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)}$ und $\|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} = \|\operatorname{div} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)}$, da $\operatorname{dev} \mathbf{I}_{n \times n} = \operatorname{div} \mathbf{I}_{n \times n} = 0$,
- (ii) Für $\boldsymbol{\sigma}_0 \in L^2(\Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ gilt

$$\|\boldsymbol{\sigma}_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \leq \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)},$$

denn $A : \mathbf{I}_{n \times n} = \operatorname{tr} A$ für alle $A \in \mathbb{R}^{n \times n}$ und $\int_{\Omega} \operatorname{tr} \boldsymbol{\sigma}_0 \, dx = 0$ und damit

$$\begin{aligned} \|\boldsymbol{\sigma}_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)}^2 &= \left\| \boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}} + 1/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \mathbf{I}_{n \times n} \right\|_{L^2(\Omega)}^2 \\ &= \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \left\| 1/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \mathbf{I}_{n \times n} \right\|_{L^2(\Omega)}^2 \\ &\quad + 2/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \left(\int_{\Omega} (\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}) : \mathbf{I}_{n \times n} \, dx \right) \\ &= \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + 1/n^2 \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right)^2 \left(\int_{\Omega} \mathbf{I}_{n \times n} : \mathbf{I}_{n \times n} \, dx \right) \\ &\quad + 2/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \left(\int_{\Omega} \operatorname{tr} \boldsymbol{\sigma}_0 - \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \\ &= \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 - 1/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \\ &= \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 - |\Omega|/n \left(\int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right)^2 \\ &\leq \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned}$$

- (iii) Wenn $f \in L^2(\Omega; \mathbb{R})$ und $\int_{\Omega} f \, dx = 0$, dann ist $\int_{\Omega} f \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx = \int_{\Omega} f \operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}} \, dx$, da $\operatorname{tr} \mathbf{q}_{\text{RT}} - \operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}} = f_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx$ konstant ist.

Desweiteren ist folgende Zerlegung sehr hilfreich.

Lemma 2.21. *Für alle deviatorischen Funktionen $u \in P_0(\mathcal{T}; \mathbb{R}^{n \times n})$ existieren $z \in H_0^1(\Omega; \mathbb{R}^n)$ mit $\operatorname{div} z = 0$ und $b \in L^2(\Omega; \mathbb{R}^{n \times n})$, so dass*

$$u = \mathbf{D} z + \operatorname{dev} b$$

eine orthogonal Zerlegung in $L^2(\Omega; \mathbb{R}^{n \times n})$ ist.

Beweis. Die Lösbarkeit des Stokes Problems in n Dimensionen (siehe beispielsweise [Bra13, S.153 ff.]) garantiert, dass es $\zeta \in L_0^2(\Omega) := \{v \in L^2(\Omega) : \int_{\Omega} v \, dx = 0\}$ und $z \in H_0^1(\Omega; \mathbb{R}^n)$ gibt, so dass für alle Testfunktionen $\varphi \in H_0^1(\Omega; \mathbb{R}^n)$ gilt

$$\int_{\Omega} u : \mathbf{D} \varphi \, dx = \int_{\Omega} \mathbf{D} z : \mathbf{D} \varphi \, dx - \int_{\Omega} \zeta \cdot \operatorname{div} \varphi \, dx \quad (2.13)$$

$$\operatorname{div} z = 0. \quad (2.14)$$

Sei

$$b := u - \mathbf{D} z \in L^2(\Omega; \mathbb{R}^{n \times n}),$$

so gilt wegen Lemma 2.5 Punkt (i) und (2.14)

$$\operatorname{tr} b = \operatorname{tr}(u - \mathbf{D} z) = 0 - \operatorname{div} z = 0.$$

Da u deviatorisch und z divergenzfrei sind, ist also auch b deviatorisch und $b = \operatorname{dev} b$. Weil außerdem $z \in H_0^1(\Omega; \mathbb{R}^n)$ liegt, gilt nach (2.13)–(2.14)

$$\begin{aligned} \int_{\Omega} b : \mathbf{D} z \, dx &= \int_{\Omega} u : \mathbf{D} z \, dx - \int_{\Omega} \mathbf{D} z : \mathbf{D} z \, dx \\ &= \int_{\Omega} \mathbf{D} z : \mathbf{D} z \, dx - \int_{\Omega} \zeta \cdot \operatorname{div} z \, dx - \int_{\Omega} \mathbf{D} z : \mathbf{D} z \, dx = 0. \end{aligned}$$

Insbesondere ergibt sich

$$\|\mathbf{D} z\|_{L^2(\Omega)}^2 + \|b\|_{L^2(\Omega)}^2 = \|u\|_{L^2(\Omega)}^2. \quad \square$$

Dieser Paragraph schließt mit einer nützlichen Abschätzung, die später Schreibarbeit erspart.

Bemerkung 2.22 Seien $a, b \in \mathbb{R}^+$, wobei $\mathbb{R}^+ := \{x \in \mathbb{R} : 0 < x\}$, und $C_1, C_2, K_1, K_2 \in \mathbb{R}^+$ Konstanten. In den Beweisen werden häufiger Abschätzungen der folgenden Form benötigt

$$C_1 a + C_2 b \leq C \sqrt{K_1 a^2 + K_2 b^2}.$$

Folgender Trick führt zu einer kleinen Konstante C

$$\begin{aligned} C_1 a + C_2 b &= \sqrt{(C_1 a + C_2 b)^2} \leq \min_{\lambda \in \mathbb{R}^+} \sqrt{(1 + \lambda) C_1^2 a^2 + (1 + 1/\lambda) C_2^2 b^2} \\ &\leq \min_{\lambda \in \mathbb{R}^+} \max \left\{ \sqrt{(1 + \lambda) C_1^2 / K_1}, \sqrt{(1 + 1/\lambda) C_2^2 / K_2} \right\} \sqrt{K_1 a^2 + K_2 b^2}. \end{aligned}$$

C wird offenbar minimal, wenn $(1 + \lambda) C_1^2 / K_1 = (1 + 1/\lambda) C_2^2 / K_2$, also $\lambda = (C_2^2 K_1) / (C_1^2 K_2)$ und $C = \sqrt{C_1^2 / K_1 + C_2^2 / K_2}$.

2.6 Modellierung

In diesem Abschnitt sollen die Stokes bzw. Navier-Stokes Gleichungen genauer betrachtet werden, die es numerisch mit der dPG-Methode zu lösen gilt. Die Ausführungen richten sich nach [CTVW10], [CW07] und [Bra13, III.6].

Mit Hilfe der Navier-Stokes Gleichungen wird die Bewegung einer inkompressiblen, zähen Flüssigkeit in einem n -dimensionalen Körper Ω (mit $n = 2$ oder $n = 3$) beschrieben.

Sei $f : \Omega \rightarrow \mathbb{R}^n$ die gegebene äußere Kraftdichte und $\eta/\rho > 0$ die kinematische Viskosität, der Quotient aus dynamischer Viskosität η und Dichte ρ . Die Zähigkeit kann durch geeignete Skalierung später als $\eta/\rho = 1$ angesetzt werden, da die Dichte der Flüssigkeit als praktisch konstant betrachtet wird. Zu bestimmen sind das Geschwindigkeitsfeld $u : \Omega \rightarrow \mathbb{R}^n$, die Druckverteilung $p : \Omega \rightarrow \mathbb{R}$ und der Spannungstensor $\sigma_{\text{ph}} : \Omega \rightarrow \mathbb{R}^{n \times n}$. Die Standardgleichungen für inkompressible Newtonsche Flüssigkeiten, die diese Größen verbinden sind, das Materialgesetz

$$\sigma_{\text{ph}} + p \mathbf{I}_{n \times n} - 2 \frac{\eta}{\rho} \epsilon(u) = 0, \quad (2.15)$$

mit dem symmetrischen Gradienten $\epsilon(u) := (\mathbf{D} u + (\mathbf{D} u)^\top) / 2$, die Impulserhaltung

$$\frac{\partial u}{\partial t} + u \cdot \mathbf{D} u - \operatorname{div} \sigma_{\text{ph}} = f \quad (2.16)$$

und das Gesetz über die Erhaltung der Masse, das in einem System ohne Quellen und Senken besagt, dass

$$\operatorname{div} u = 0. \quad (2.17)$$

Zu diesem System gehören noch eine Anfangsbedingung $u|_{t=0} = u_0$ und Randbedingungen. In den meisten Anwendungen inkompressibler Newtonscher Flüsse handelt es sich um Dirichlet Randbedingungen an die Geschwindigkeit. Davon soll auch hier ausgegangen

werden. Es ist also

$$u = g \quad \text{entlang } \partial\Omega \quad (2.18)$$

gefordert. Aus dem Gaußschen Integralsatz ergibt sich eine Verträglichkeitsbedingung an die Randwerte $g : \Omega \rightarrow \mathbb{R}^n$, denn

$$\int_{\partial\Omega} g \cdot \nu \, ds = \int_{\partial\Omega} u \cdot \nu \, ds = \int_{\Omega} \operatorname{div} u \, dx = 0. \quad (2.19)$$

Da der Druck durch (2.15)–(2.17) lediglich bis auf eine additive Konstante bestimmt wird, ist weiterhin eine Normierung der Form

$$\int_{\Omega} p \, dx = 0$$

üblich. Die Gleichungen (2.15)–(2.17) werden als Spannungs-Geschwindigkeits-Druck Formulierung der inkompressiblen Navier-Stokes Gleichungen bezeichnet.

Aus diesem System kann die Spannung σ_{ph} eliminiert werden, in dem (2.15) in (2.16) eingesetzt wird. Da bei inkompressiblen Flüssigkeiten die Divergenz von $(\operatorname{D} u)^\top$ verschwindet, lautet das entstehende System

$$\begin{aligned} \frac{\partial u}{\partial t} + u \cdot \operatorname{D} u + \nabla p - \frac{\eta}{\rho} \Delta(u) &= f \quad \text{in } \Omega, \\ \operatorname{div} u &= 0 \quad \text{in } \Omega \end{aligned} \quad (2.20)$$

und wird als Geschwindigkeits-Druck Formulierung der Navier-Stokes Gleichung bezeichnet. Hierbei sei auch der Laplace-Operator zeilenweise definiert $\Delta u := (\Delta u_1, \dots, \Delta u_n)^\top = \sum_{k=1}^n (\partial_{kk} u_1, \dots, \partial_{kk} u_n)^\top = \operatorname{div} \operatorname{D} u$. Ohne den nichtlinearen Term $u \cdot \operatorname{D} u$ werden diese Gleichungen als Stokes Gleichungen bezeichnet. Das klassische stationäre Stokes Problem entsteht, wenn die Zähigkeit zu 1 skaliert wird. Es werden also ein Geschwindigkeitsfeld $u : \Omega \rightarrow \mathbb{R}^n$ und eine Druckverteilung $p : \Omega \rightarrow \mathbb{R}$ gesucht, die folgenden Gleichungen erfüllen

$$\begin{aligned} \Delta u + \nabla p &= -f \quad \text{in } \Omega, \\ \operatorname{div} u &= 0 \quad \text{in } \Omega, \\ u &= g \quad \text{entlang } \partial\Omega. \end{aligned} \quad (2.21)$$

Dabei genüge g der Verträglichkeitsbedingung (2.19). Als klassische Lösung des Stokes Problem werden $u \in C^2(\Omega, \mathbb{R}^n) \cap C(\bar{\Omega}; \mathbb{R}^n)$ und $p \in C^0(\Omega)$ bezeichnet, die das System (2.21) erfüllen, sollten diese existieren.

Unter Verwendung der homogenen Randbedingungen $g \equiv 0$ mit den Räumen $\tilde{X} := H_0^1(\Omega; \mathbb{R}^n)$, $\tilde{Y} := L_0^2 := \{g \in L^2(\Omega) : \int_{\Omega} g \, dx = 0\}$ ergibt sich das folgende Sattelpunkt-

problem. Zu gegebenem $f \in L^2(\Omega; \mathbb{R}^n)$ sind $u \in \tilde{X}$ und $p \in \tilde{Y}$ gesucht mit

$$\begin{aligned} \int_{\Omega} \mathbf{D} u : \mathbf{D} v \, dx + \int_{\Omega} \operatorname{div} v \cdot p \, dx &= \int_{\Omega} f \cdot v \, dx \quad \text{für alle } v \in \tilde{X}, \\ \int_{\Omega} \operatorname{div} u \cdot q \, dx &= 0 \quad \text{für alle } q \in \tilde{Y}. \end{aligned}$$

Die eindeutige Lösbarkeit dieser gemischten Formulierung wird mit Hilfe des Brezzi-Splitting Theorems bewiesen, in dem das Ladyshenskaya Lemma und Eigenschaften von \mathbf{D} ausgenutzt werden, um die Inf-Sup-Bedingung zu beweisen [Bra13, S.154ff]. Zu dieser variationellen Formulierung sind zahlreiche Diskretisierungen bekannt. Beispielsweise das Taylor-Hood-Element mit $\tilde{X}_h := C_0(\Omega; \mathbb{R}^2) \cap P_2(\mathcal{T}; \mathbb{R}^2)$ und $\tilde{Y}_h := C(\bar{\Omega}) \cap L_0^2(\Omega) \cap P_1(\mathcal{T})$ [Bra13, S.163] oder das MINI-Element welches in \tilde{X}_h bubble Funktionen, die sich als Produkt der Baryzentrischen Koordinaten ergeben, einsetzt [Bra13, S.164].

Die Beschäftigung mit nicht Newtonschen Flüssigkeiten lenkt mehr Interesse auf die Spannungs-Geschwindigkeits-Druck Formulierung, denn in dieser Situation ist (2.15) nicht linear, die Spannung kann damit nicht einfach eliminiert werden. Ein anderer Vorteil dieses Systems ist, dass die Spannung direkt berechnet wird. Dies verhindert, dass durch die numerische Ableitung der berechneten Geschwindigkeit, die Approximationsordnung von σ_{ph} vermindert wird. Die Spannungs-Geschwindigkeits-Druck Formulierung weist allerdings auch Nachteile auf. Zum einen ist die Zahl der Unbekannten deutlich höher, zum anderen ist die Symmetrieanforderung an die Spannung kompliziert. Daher wird die Pseudospannungs-Geschwindigkeits Formulierung betrachtet, die akkurate Approximation der physikalischen Größen Spannung und Wirbelstärke ermöglicht und keine Zusatzbedingungen an die Approximationsräume stellt. Dazu wird die Pseudospannung

$$\sigma := \frac{\eta}{\rho} \mathbf{D} u - p \mathbf{I}_{n \times n}$$

definiert. Auf Grund der Divergenzfreiheit von u , (2.17), gilt

$$p = -\frac{1}{n} \operatorname{tr} \sigma.$$

Der Druck ist also leicht zu bestimmen. Für den spurfreien Anteil $\operatorname{dev} \sigma = \sigma - 1/n \operatorname{tr} \sigma \mathbf{I}_{n \times n}$, gilt also

$$\operatorname{dev} \sigma = \frac{\eta}{\rho} \mathbf{D} u.$$

Da zusätzlich nach (2.15) und wegen $\operatorname{div} ((\mathbf{D} u)^\top) = 0$ gilt

$$\operatorname{div} \sigma_{\text{ph}} = \operatorname{div} \sigma,$$

ergeben sich die folgenden Gleichungen als Pseudo-Spannung-Geschwindigkeits Formulie-

rung

$$\frac{\partial u}{\partial t} + u \cdot \mathbf{D} u - \operatorname{div} \boldsymbol{\sigma} = f \quad \text{in } \Omega, \quad (2.22)$$

$$\operatorname{dev} \boldsymbol{\sigma} - \frac{\eta}{\rho} \mathbf{D} u = 0 \quad \text{in } \Omega, \quad (2.23)$$

$$u = g \quad \text{entlang } \partial\Omega.$$

wobei g wie zuvor (2.19) erfülle. Die Normierung $\int_{\Omega} p \, dx = 0$, wird zu der linearen Nebenbedingung

$$\int_{\Omega} \operatorname{tr} \boldsymbol{\sigma} \, dx = 0,$$

die in $H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ gefordert wird. Aus der Pseudospannung $\boldsymbol{\sigma}$ lassen sich neben dem Druck, p , auch ohne Verringerung der Approximationsgüte die Spannung, $\boldsymbol{\sigma}_{\text{ph}}$, und die Wirbelstärke, $\omega = \nabla \times u = 1/2 (\mathbf{D} u - (\mathbf{D} u)^{\top})$, wie folgt berechnen

$$\begin{aligned} \boldsymbol{\sigma}_{\text{ph}} &= -p \mathbf{I}_{n \times n} + \eta/\rho (\mathbf{D} u + (\mathbf{D} u)^{\top}) = \boldsymbol{\sigma} + \eta/\rho (\mathbf{D} u)^{\top} = \boldsymbol{\sigma} + (\operatorname{dev} \boldsymbol{\sigma})^{\top}, \\ \omega &= \rho/(2\eta) (\operatorname{dev} \boldsymbol{\sigma} - (\operatorname{dev} \boldsymbol{\sigma})^{\top}). \end{aligned}$$

In dieser Arbeit wird das stationäre Stokes Problem in der Pseudospannungs-Geschwindigkeits Formulierung betrachtet, bei dem die als konstant angenommene Viskosität η/ρ durch Skalierung von $\boldsymbol{\sigma}$ und f mit dieser Größe nicht mehr auftritt.

Problem 1 (Pseudospannungs-Geschwindigkeits Formulierung)

Zu gegebenem $f \in L^2(\Omega; \mathbb{R}^n)$ und $g \in H^{1/2}(\partial\Omega; \mathbb{R}^n)$ sind $u \in H^1(\Omega; \mathbb{R}^n)$ und $\boldsymbol{\sigma} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ gesucht mit

$$\begin{aligned} \operatorname{div} \boldsymbol{\sigma} + f &= 0 \quad \text{in } \Omega, \\ \operatorname{dev} \boldsymbol{\sigma} - \mathbf{D} u &= 0 \quad \text{in } \Omega, \\ u &= g \quad \text{entlang } \partial\Omega, \end{aligned}$$

wobei g entlang jeder der M Zusammenhangskomponenten Γ_m , $1 \leq m \leq M$, des Randes $\Gamma = \bigcup_{m=1}^M \Gamma_m$ die Kompatibilitätsbedingung

$$\int_{\Gamma_m} \nu \cdot g \, ds = 0$$

erfüllt.

Eine gemischte Formulierung zu diesem Problem verwendet $\tilde{X} := H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})/\mathbb{R}$ und $\tilde{Y} := L^2(\Omega)$ und sucht $\boldsymbol{\sigma} \in \tilde{X}$ und $u \in \tilde{Y}$ mit (2.18) zu gegebenem $f \in L^2(\Omega)$, so

dass

$$\begin{aligned} \int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot \operatorname{div} \boldsymbol{\tau} \, dx &= \int_{\partial\Omega} g \cdot (\boldsymbol{\tau} \boldsymbol{\nu}) \, ds \quad \text{für alle } \boldsymbol{\tau} \in \tilde{X}, \\ \int_{\Omega} \operatorname{div} \boldsymbol{\sigma} \cdot v \, dx &= - \int_{\Omega} f \cdot v \, dx \quad \text{für alle } v \in \tilde{Y}. \end{aligned}$$

Dieses Variationsproblem hat nach [CTVW10, Theorem 2.3] eine eindeutige Lösung.

Bemerkung 2.23 Sei $Z := H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega)$ für jedes $(\boldsymbol{\rho}, w) \in Z$ versehen mit der Norm $\|(\boldsymbol{\rho}, w)\|_Z^2 := \|\boldsymbol{\rho}\|_{H(\operatorname{div}, \Omega)}^2 + \|w\|_{L^2(\Omega)}^2$. In [CTVW10, Theorem 2.3] wird mit Hilfe des Brezzi-Splitting Theorems [Bra13, Satz 4.3] bewiesen, dass die lineare Abbildung $L : Z \rightarrow Z^*$, die durch

$$\langle L(z), \hat{z} \rangle = \tilde{b}(z, \hat{z}) \quad \text{für } \hat{z} \in Z$$

definiert ist, ein Isomorphismus ist. Also erfüllt die Bilinearform $\tilde{b} : Z \times Z \rightarrow \mathbb{R}$

$$\tilde{b}((\boldsymbol{\tau}, v), (\boldsymbol{\rho}, w)) := \int_{\Omega} \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\rho} \, dx + \int_{\Omega} v \cdot \operatorname{div} \boldsymbol{\rho} \, dx + \int_{\Omega} \operatorname{div} \boldsymbol{\tau} \cdot w \, dx,$$

die folgende Inf-Sup-Bedingung

$$\gamma := \inf_{z \in Z \setminus \{0\}} \sup_{\hat{z} \in Z \setminus \{0\}} \frac{\tilde{b}(z, \hat{z})}{\|z\|_Z \|\hat{z}\|_Z} > 0.$$

Dies wird später für den Beweis einer ähnlichen Inf-Sup-Bedingung in der ultra-schwachen Formulierung verwendet, siehe Theorem 3.8.

In der dPG-Version der Pseudospannungs-Geschwindigkeits Formulierung wird schon auf dem kontinuierlichen Level eine reguläre Triangulierung \mathcal{T} zu Grunde gelegt. Wie auch bei der gemischten Formulierung werden die Gleichungen aus Problem 1 mit jeweils einer Testfunktion multipliziert und über Ω integriert. Anschließend werden beide Gleichungen partiell integriert, da die Testfunktionen im vorliegenden Setting gebrochen sind, entstehen dabei Randintegrale auf dem Skelett. Es ergeben sich folgende Gleichungen

$$\int_{\Omega} \boldsymbol{\sigma} : \operatorname{D}_{\text{NC}} v \, dx - \langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\sigma}, \gamma_0^{\mathcal{T}} v \rangle_{\partial\mathcal{T}} = \int_{\Omega} f \cdot v \, dx \quad \text{für alle } v \in H^1(\mathcal{T}; \mathbb{R}^n), \quad (2.24)$$

$$\int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} + \int_{\Omega} u \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \, dx - \langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\tau}, \gamma_0^{\mathcal{T}} u \rangle = 0 \quad \text{für alle } \boldsymbol{\tau} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}. \quad (2.25)$$

Es werden neue Variablen für die Spuren der Geschwindigkeit u auf dem Skelett $\gamma_0^{\mathcal{T}} u = s$ und die Normalenspur des Pseudostress $\boldsymbol{\sigma}$ auf dem Skelett $\gamma_{\nu}^{\mathcal{T}} \boldsymbol{\sigma} = t$ eingeführt. Dadurch ist es möglich weder an u noch an $\boldsymbol{\sigma}$ Differenzierbarkeitsbedingungen zu stellen. Auch die Geschwindigkeitsrandbedingungen werden auf s übertragen, also wird $s|_{\Gamma} = g$ gefordert. Die Komponenten einer Lösung zu (2.24)–(2.25) werden also in folgenden Räumen gesucht, $\boldsymbol{\sigma} \in L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$, $u \in L^2(\Omega; \mathbb{R}^n)$, $s \in H^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ und $t \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$.

Die dPG Methode löst dieses System, in dem $x \in X$ wie in (2.7) gesucht wird, so dass für alle $y \in Y$ wie in (2.8)

$$b(x, y) = F(y). \quad (2.26)$$

Dabei ist die Bilinearform b , wie in Definition 3.2 definiert, die Summe der linken Seiten von (2.24) und (2.25). Das Funktional F aus Definition 3.2 ist leicht modifiziert, es enthält neben den Summen der rechten Seiten noch den Term $\langle \gamma_\nu^T \boldsymbol{\tau}, \hat{s} \rangle_{\partial \mathcal{T}}$ wobei $\hat{s} \in H^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ im Wesentlichen eine Funktion ist, die die Randdaten $\hat{s}|_\Gamma = g$ realisiert. Dies ermöglicht ein diskretes Problem mit homogenen Randbedingungen, also $s \in H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$, zu lösen. Aus einer Lösung u und $\boldsymbol{\sigma}$ von Problem 1 kann die Lösung dieses modifizierten Problems, leicht gewonnen werden. Es gilt nach Satz 3.3, dass $x = (u, \boldsymbol{\sigma}, \gamma_0^T u - \hat{s}, \gamma_\nu^T \boldsymbol{\sigma})$ (2.26) erfüllt.

3 Analysis der dPG-Methode

Zunächst einmal werden in Abschnitt 3.1 einige allgemeine Informationen zu der Klasse der dPG-Methoden gegeben. In den folgenden Paragraphen wird dann die neue Methode detaillierter vorgestellt und analysiert.

3.1 Allgemeine Informationen zur dPG-Methode

Dieser Paragraph gibt einen Überblick über die allgemeine Theorie der dPG-Methoden inklusive der a-posteriori- Fehlerabschätzung aus [CDG14].

Die diskontinuierliche Petrov-Galerkin-Methode (dPG-Methode) kann als Finite-Elemente-Methode mit Nicht-Standardtesträumen oder als Least-Squares Finite-Elemente-Methode, die das Residuum der ultra-schwachen Formulierung in einer Nicht-Standardnorm minimiert, aufgefasst werden [DG10],[DG11b]. Diese Dualität wird auch in den verschiedenen Problemformulierungen deutlich. Dazu werden zunächst die Banachräume $(X, \|\cdot\|_X)$ und $(Y, \|\cdot\|_Y)$ und zwei abgeschlossene, endlich dimensionale Unterräume $X_h \subseteq X$ und $Y_h \subseteq Y$ betrachtet. Außerdem sei $b: X \times Y \rightarrow \mathbb{R}$ eine reelle, beschränkte Bilinearform. In diesem Setting können folgende Formulierungen betrachtet werden.

Problem 2 (Variationsformulierung des exakten Problems)

Zu gegebenem $F \in Y^*$ ist $x \in X$ gesucht mit $b(x, \cdot) = F$ in Y^* .

Dieses Problem ist wegen $\|F - b(x, \cdot)\|_{Y^*} = 0 \leq \|F - b(\xi, \cdot)\|_{Y^*}$ für alle $\xi \in X$ äquivalent zu folgender Minimierungsaufgabe.

Problem 3 (Exaktes Problem als Minimierung des Residuums)

Zu gegebenem $F \in Y^*$ ist $x \in X$ gesucht mit $x \in \operatorname{argmin}_{\xi \in X} \|F - b(\xi, \cdot)\|_{Y^*}$.

Diese Zielfunktion wird unter Umständen auch quadriert betrachtet, weshalb die dPG-Methode als Least-Square-Methode verstanden werden kann. In der idealisierten dPG-Methode wird das Problem im endlich dimensionalen Unterraum $X_h \subseteq X$ gelöst.

Problem 4 (Idealisierte dPG-Methode als Minimierung des Residuums)

Zu gegebenem $F \in Y^*$ ist $x_h \in X_h$ gesucht mit $x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \cdot)\|_{Y^*}$.

Ist Y insbesondere ein Hilbertraum mit Skalarprodukt $\langle \cdot, \cdot \rangle_Y$, so erlaubt der sogenannte *trial-to-test*-Operator $T: X \rightarrow Y$ mit $\langle Tx, y \rangle_Y = b(x, y)$ für alle $y \in Y$ eine äquivalente Formulierung, [RBTD14]. Wenn $R: Y \rightarrow Y^*$ der Riesz-Isomorphismus ist [Alt06, S.163] und $B: X \rightarrow Y^*$ die Abbildung $Bx(y) := b(x, y)$, so gilt $T = R^{-1} \circ B$.

Problem 5 (Idealisierte dPG-Methode als Variationsformulierung)

Zu gegebenem $F \in Y^*$ ist $x_h \in X_h$ gesucht mit $b(x_h, \bullet) = F$ in $(T(X_h))^*$.

Im Allgemeinen ist Y ein unendlich dimensionaler Vektorraum und der trial-to-test-Operator T bzw. $\|F - b(\xi_h, \bullet)\|_{Y^*}$ sind in der idealisierten dPG-Methode kaum zu berechnen. Darum wird in der praktischen dPG-Methode statt Y ein ebenfalls endlich dimensionaler Unterraum $Y_h \subseteq Y$ betrachtet.

Problem 6 (Praktische dPG-Methode als Minimierung des Residuums)

Zu gegebenem $F \in Y^*$ ist $x_h \in X_h$ gesucht mit $x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}$.

Dazu existiert ebenfalls eine äquivalente Variationsformulierung. Dafür wird der approximierte trial-to-test-Operator $T_h : X \rightarrow Y_h$ mit $\langle Tx, y_h \rangle_Y = b(x, y_h)$ für alle $y_h \in Y_h$ definiert.

Problem 7 (Praktische dPG-Methode als Variationsformulierung)

Zu gegebenem $F \in Y^*$ ist $x_h \in X_h$ gesucht mit $b(x_h, \bullet) = F$ in $(T_h(X_h))^*$.

Für die praktische dPG Methode wurde in [CDG14] eine a-posteriori-Fehlerkontrolle vorgestellt. Dabei wird vorausgesetzt, dass $(X, \|\cdot\|_X)$ ein reflexiver reeller Banachraum ist und $(Y, \langle \cdot, \cdot \rangle_Y)$ ein reeller Hilbertraum. Dies ist für die in dieser Arbeit verwendeten Räume aus (2.7) bzw. (2.8) gewährleistet. Um die Stabilität, die quasi-optimale Konvergenz sowie die Zuverlässigkeit der dPG-Methode und die Effizienz des Residuumschätzers zu folgern, werden außerdem die folgenden Bedingung gefordert, deren Nachweis sich der Rest des Kapitels widmet.

Bedingung (B1) „Beschränktheits- und Eindeutigkeitsbedingung“

Es ist $b : X \times Y \rightarrow \mathbb{R}$ eine beschränkte Bilinearform mit $\|b\| := \sup_{x \in S(X)} \sup_{y \in S(Y)} b(x, y) < \infty$ und trivialem Kern $N := \{y \in Y : b(\bullet, y) = 0 \text{ in } X^*\} = \{0\}$.

Bedingung (B2) „Kontinuierliche Inf-Sup-Bedingung“

Es gilt $0 < \beta := \inf_{x \in S(X)} \sup_{y \in S(Y)} b(x, y)$.

Bedingung (B3) „Existenz eines Fortin-Operators“

Die abgeschlossenen Unterräume $X_h \subseteq X$, $Y_h \subseteq Y$ erlauben die Definition einer beschränkten, linearen Abbildung $\Pi : Y \rightarrow Y_h$ mit Operatornorm $\|\Pi\|$ und der Eigenschaft $b(\xi_h, (1 - \Pi)y) = 0$ für alle $\xi_h \in X_h$, $y \in Y$.

Aus (B1) und (B2) folgt nach [Bra13, Satz 3.6], dass Problem 2 bzw. Problem 3 eine eindeutige Lösung $x \in X$ besitzen. Außerdem ergibt sich, dass die Zielfunktion der Minimierungsaufgabe, $\|F - b(\xi_h, \bullet)\|_{Y^*}$, in Problem 4 bzw. Problem 5 ein zuverlässiger und effizienter Fehlerschätzer ist.

Lemma 3.1 (Zuverlässigkeit und Effizienz der idealisierten dPG-Methode).

Unter den Voraussetzungen (B1) und (B2) gilt für jedes $\xi_h \in X_h$

$$\beta \|x - \xi_h\|_X \leq \|F - b(\xi_h, \cdot)\|_{Y^*} \leq \|b\| \|x - \xi_h\|_X.$$

Beweis. Es gilt

$$\beta = \inf_{\zeta \in X \setminus \{0\}} \frac{\|b(\zeta, \cdot)\|_{Y^*}}{\|\zeta\|_X} \leq \frac{\|b(x - \xi_h, \cdot)\|_{Y^*}}{\|x - \xi_h\|_X} \leq \frac{\|b\| \|x - \xi_h\|_X}{\|x - \xi_h\|_X}.$$

Dies ist die Behauptung, da $\|b(x - \xi_h, \cdot)\|_{Y^*} = \|F - b(\xi_h, \cdot)\|_{Y^*}$. \square

Für die praktische dPG Methode wurde in [GQ14, Theorem 2.1] unter der Zusatzbedingung, dass X ein Hilbertraum ist, bereits gezeigt, dass aus den Bedingungen (B1), (B2) und (B3) die diskrete Stabilität und quasi-optimale Konvergenz in dem Sinn folgen, dass

$$\|x - x_h\|_X \leq \frac{\|\Pi\| \|b\|}{\beta} \min_{\xi_h \in X_h} \|x - \xi_h\|_X$$

gilt. Wie bereits erwähnt, wurde die a-posteriori-Fehlerabschätzung für die praktische dPG-Methode unter den genannten Bedingungen in [CDG14] vervollständigt. Die Einbeziehung eines Datenapproximationsfehlers neben dem Residuumsfehler ergibt einen zuverlässigen und effizienten Schätzer für die praktische dPG-Methode. Genauer gilt für alle $\xi_h \in X_h$

$$\begin{aligned} \beta \|x - \xi_h\|_X &\leq \sqrt{1 + \|\Pi\|^2} \|F - b(\xi_h, \cdot)\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \\ &\leq \|b\| \left(\sqrt{1 + \|\Pi\|^2} + \|1 - \Pi\| \right) \|x - \xi_h\|_X \end{aligned}$$

Da diese Abschätzungen für alle Elemente aus X_h und nicht nur für die diskrete Lösung x_h gelten, ist Stabilität beim inexakten Lösen gewährleistet.

In [CH15, Lemma 2.5] wird bewiesen, dass die letzte Bedingung durch die folgende diskrete Inf-Sup-Bedingung ersetzt werden kann.

Bedingung (B4) „Diskrete Inf-Sup-Bedingung“

Es gilt $0 < \beta_h := \inf_{x_h \in S(X_h)} \sup_{y_h \in S(Y_h)} b(x_h, y_h)$ für die abgeschlossenen Unterräume $X_h \subseteq X$, $Y_h \subseteq Y$.

Zusätzlich wird lediglich gefordert, dass auch X wie im vorliegenden Fall ein Hilbertraum ist. Ist (B2) erfüllt, so sind (B3) und (B4) sogar äquivalent und es gilt $\beta \leq \beta_h \|\Pi\|$ bzw. $\|\Pi\| \leq \|b\| / \beta_h$. Nach [CH15, Lemma 2.5] gilt dann die folgende Fehlerabschätzung für alle $\xi_h \in X_h$

$$\begin{aligned} \beta \|x - \xi_h\|_X &\leq \|b\| / \beta_h \|F - b(\xi_h, \cdot)\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \\ &\leq 2\|b\|^2 / \beta_h \|x - \xi_h\|_X. \end{aligned} \quad (3.1)$$

In dieser Arbeit wird der Ansatz verfolgt, die Bedingungen (B1), (B2) und (B4) für die vorgestellte neue Methode zu überprüfen. In dem gegebenen Setting folgt aus (B4) die eindeutige Lösbarkeit von Problem 6 bzw. Problem 7. In [Hel14, Kapitel 3] sind auch Bedingungen und Zuverlässig- bzw. Effizienzabschätzungen für den Fall, dass Y kein Hilbertraum ist, zu finden.

3.2 Problemformulierung und kontinuierliche Bedingungen

In diesem Paragraphen werden die Bilinearform b und das Funktional F definiert. Anschließend werden die Annahmen (B1) und (B2) aus dem vorangegangenen Abschnitt überprüft.

Gegeben seien dazu der Ansatzraum X und der Testraum Y

$$\begin{aligned} X &:= L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n) \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n), \\ Y &:= H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times H^1(\mathcal{T}; \mathbb{R}^n) \end{aligned}$$

ausgestattet mit den Normen

$$\begin{aligned} \|(\boldsymbol{\sigma}, u, s, t)\|_X^2 &:= \|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 + \|s\|_{H_0^{1/2}(\partial\mathcal{T})}^2 + \|t\|_{H^{-1/2}(\partial\mathcal{T})}^2, \\ \|(\boldsymbol{\tau}, v)\|_Y^2 &:= \|\boldsymbol{\tau}\|_{H(\operatorname{div}, \mathcal{T})}^2 + \|v\|_{H^1(\mathcal{T})}^2. \end{aligned}$$

X und Y sind nach Bemerkung 2.14 insbesondere Hilberträume.

Zunächst werden aus den gegebenen Stokes Gleichungen die Bilinearform b und das Funktional F zur Formulierung der dPG-Methode hergeleitet, dabei wird die bereits erwähnte Pseudospannungs-Geschwindigkeits Formulierung verwendet.

Problem 8 (Pseudospannungs-Geschwindigkeits Formulierung)

Zu gegebenem $f \in L^2(\Omega; \mathbb{R}^n)$ und $g \in H^{1/2}(\partial\Omega; \mathbb{R}^n)$ sind $u \in H^1(\Omega; \mathbb{R}^n)$ und $\boldsymbol{\sigma} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$ gesucht mit

$$\operatorname{div} \boldsymbol{\sigma} + f = 0 \quad \text{in } \Omega, \quad (3.2)$$

$$\operatorname{dev} \boldsymbol{\sigma} - D u = 0 \quad \text{in } \Omega, \quad (3.3)$$

$$u = g \quad \text{entlang } \Gamma, \quad (3.4)$$

wobei g entlang jeder der M Zusammenhangskomponenten Γ_m , $1 \leq m \leq M$, des Randes $\Gamma = \bigcup_{m=1}^M \Gamma_m$ die Kompatibilitätsbedingung

$$\int_{\Gamma_m} \boldsymbol{\nu} \cdot g \, ds = 0 \quad (3.5)$$

erfüllt.

Definition 3.2. Die Bilinearform $b : X \times Y \rightarrow \mathbb{R}$ sei für alle $x = (\boldsymbol{\sigma}, u, s, t) \in X$ und $y = (\boldsymbol{\tau}, v) \in Y$ definiert durch

$$b(x, y) := \int_{\Omega} \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx + \int_{\Omega} \text{dev } \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx - \langle t, \gamma_0^{\mathcal{T}} v \rangle_{\partial \mathcal{T}} - \langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\tau}, s \rangle_{\partial \mathcal{T}}. \quad (3.6)$$

Zu den gegebenen Daten $f \in L^2(\Omega; \mathbb{R}^n)$ und $\hat{s} \in H^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ wird das Funktional $F \in Y^*$ für alle $y = (\boldsymbol{\tau}, v) \in Y$ definiert durch

$$F(y) := \int_{\Omega} f \cdot v \, dx + \langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\tau}, \hat{s} \rangle_{\partial \mathcal{T}}. \quad (3.7)$$

Mit (2.7)–(2.8) und (3.6)–(3.7) lässt sich die ultra-schwache Formulierung analog Problem 2 aufstellen.

Problem 9 (Ultra-schwache Formulierung)

Zu gegebenem $f \in L^2(\Omega; \mathbb{R}^n)$ und $\hat{s} \in H^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ ist $x = (\boldsymbol{\sigma}, u, t, s) \in X$ gesucht mit

$$b(x, y) = F(y) \quad \text{für alle } y = (\boldsymbol{\tau}, v) \in Y. \quad (3.8)$$

Im Folgenden wird der Zusammenhang zwischen der ultra-schwachen und der Pseudo-spannungs-Geschwindigkeits Formulierung untersucht und gezeigt, dass für die theoretischen Betrachtungen homogene Nullranddaten genügen, wenn $\hat{s} = g$ entlang Γ gilt.

Satz 3.3. Es sei $(u, \boldsymbol{\sigma}) \in H^1(\Omega; \mathbb{R}^n) \times H(\text{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$ eine Lösung von Problem 8. Außerdem sei $\hat{u} \in H^1(\Omega; \mathbb{R}^n)$ mit $\hat{u} = g$ auf Γ und $\hat{s} := \gamma_0^{\mathcal{T}} \hat{u} \in H^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$. Setze nun $s := \gamma_0^{\mathcal{T}} u - \hat{s} \in H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$ und $t := \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\sigma} \in H^{-1/2}(\partial \mathcal{T}; \mathbb{R}^n)$.

Dann löst $(\boldsymbol{\sigma}, u, s, t) \in X$ Problem 9 zu den gewünschten Daten $f \in L^2(\Omega; \mathbb{R}^n)$ und $\hat{s} \in H^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$.

Beweis. Zunächst einmal gilt nun $(\boldsymbol{\sigma}, u, s, t) \in X$, da $u - \hat{u} = 0$ auf Γ und somit $s \in H_0^{1/2}(\partial \mathcal{T}; \mathbb{R}^n)$. Es seien nun $(\boldsymbol{\tau}, v) \in Y$ und $T \in \mathcal{T}$ beliebig. Mit Hilfe der partiellen Integration aus Satz 2.8 und Gleichung (3.2) aus Problem 8 ergibt sich

$$\langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\sigma}, \gamma_0^{\mathcal{T}} v \rangle_{\partial T} = \int_T \boldsymbol{\sigma} : \mathbf{D} v \, dx + \int_T \text{div } \boldsymbol{\sigma} \cdot v \, dx = \int_T \boldsymbol{\sigma} : \mathbf{D} v \, dx - \int_T f \cdot v \, dx$$

und damit in Summe über alle Simplices $T \in \mathcal{T}$, da $t = \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\sigma}$,

$$\langle t, \gamma_0^{\mathcal{T}} v \rangle_{\partial \mathcal{T}} = \int_{\Omega} \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx - \int_{\Omega} f \cdot v \, dx. \quad (3.9)$$

Ebenso ergibt sich aus Satz 2.8 und Gleichung (3.3) aus Problem 8

$$\langle \gamma_\nu \boldsymbol{\tau}, \gamma_0 u \rangle_{\partial T} = \int_T \boldsymbol{\tau} : \mathbf{D} u \, dx + \int_T \operatorname{div} \boldsymbol{\tau} \cdot u \, dx = \int_T \boldsymbol{\tau} : \operatorname{dev} \boldsymbol{\sigma} \, dx + \int_T \operatorname{div} \boldsymbol{\tau} \cdot u \, dx$$

und nach Summation über alle Elemente $T \in \mathcal{T}$, da $\gamma_0^\mathcal{T} u = \hat{s} + s$,

$$\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, \hat{s} + s \rangle_{\partial \mathcal{T}} = \int_\Omega \boldsymbol{\tau} : \operatorname{dev} \boldsymbol{\sigma} \, dx + \int_\Omega \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot u \, dx. \quad (3.10)$$

Die Addition von Gleichung (3.9) und Gleichung (3.10) ergibt

$$\begin{aligned} \int_\Omega f \cdot v \, dx + \langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, \hat{s} \rangle_{\partial \mathcal{T}} &= - \langle t, \gamma_0^\mathcal{T} v \rangle_{\partial \mathcal{T}} + \int_\Omega \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx \\ &\quad - \langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, s \rangle_{\partial \mathcal{T}} + \int_\Omega \boldsymbol{\tau} : \operatorname{dev} \boldsymbol{\sigma} \, dx + \int_\Omega \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot u \, dx, \end{aligned}$$

also ist nach Definition 3.2 Gleichung (3.8) erfüllt und die Behauptung bewiesen. \square

Zunächst wird nun der erste Teil von (B1), die Beschränktheit von b , bewiesen.

Satz 3.4 (Beschränktheit von b). *Die in (3.6) definierten Bilinearform b ist beschränkt, denn es gilt*

$$|b(x, y)| \leq M \|x\|_X \|y\|_Y \quad \text{für alle } x \in X, y \in Y,$$

wobei $M = \|b\| \leq \sqrt{3}$.

Beweis. Es seien $x = (\boldsymbol{\sigma}, u, t, s) \in X$ und $y = (\boldsymbol{\tau}, v) \in Y$ gegeben. Nach Bemerkung 2.11 existieren $w \in H_0^1(\Omega; \mathbb{R}^n)$ mit $s = \gamma_0^\mathcal{T} w$ und $\|w\|_{H^1(\Omega)} = \|s\|_{H^{1/2}(\partial \mathcal{T})}$ sowie $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ mit $t = \gamma_\nu^\mathcal{T} \mathbf{q}$ und $\|\mathbf{q}\|_{H(\operatorname{div})} = \|t\|_{H^{-1/2}(\partial \mathcal{T})}$. Die partielle Integrationsformel (2.8) bedeutet nach Summation über alle $T \in \mathcal{T}$

$$\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, s \rangle_{\partial \mathcal{T}} = \int_\Omega \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot w \, dx + \int_\Omega \boldsymbol{\tau} : \mathbf{D} w \, dx$$

und

$$\langle t, \gamma_0^\mathcal{T} v \rangle_{\partial \mathcal{T}} = \int_\Omega \operatorname{div} \mathbf{q} \cdot v \, dx + \int_\Omega \mathbf{q} : \mathbf{D}_{\text{NC}} v \, dx.$$

Damit kann b wie folgt geschrieben werden

$$\begin{aligned} b(x, y) &= - \int_\Omega \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot w \, dx - \int_\Omega \boldsymbol{\tau} : \mathbf{D} w \, dx - \int_\Omega \operatorname{div} \mathbf{q} \cdot v \, dx - \int_\Omega \mathbf{q} : \mathbf{D}_{\text{NC}} v \, dx \\ &\quad + \int_\Omega \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx + \int_\Omega \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_\Omega u \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \, dx. \end{aligned}$$

Die Anwendung von Dreiecks- und Cauchy-Schwarz-Ungleichung sowie Lemma 2.5 Punkt (iii) ergibt

$$|b(x, y)| \leq \|\boldsymbol{\tau}\|_{L^2(\Omega)} \left(\|\mathbf{D} w\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right) + \|\operatorname{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)} \left(\|w\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)} \right)$$

$$\|v\|_{L^2(\Omega)} \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)} + \|\mathbf{D}_{\text{NC}} v\|_{L^2(\Omega)} \left(\|\mathbf{q}\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right).$$

Dies kann als \mathbb{R}^4 -Skalarprodukt der Vektoren

$$\left(\|\mathbf{D} w\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)}, \|w\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)}, \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)}, \|\mathbf{q}\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right)^\top.$$

und

$$\left(\|\boldsymbol{\tau}\|_{L^2(\Omega)}, \|\operatorname{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)}, \|v\|_{L^2(\Omega)}, \|\mathbf{D}_{\text{NC}} v\|_{L^2(\Omega)} \right)^\top.$$

betrachtet werden. Dafür ergibt die Cauchy-Schwarz-Ungleichung folgende Abschätzung

$$\begin{aligned} |b(x, y)| &\leq \|y\|_Y \left(\left(\|\mathbf{D} w\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right)^2 + \left(\|w\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)} \right)^2 + \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)}^2 \right. \\ &\quad \left. + \left(\|\mathbf{q}\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right)^2 \right)^{1/2}. \end{aligned}$$

Der zweite Faktor kann wie folgt gegen $\|x\|_X$ abgeschätzt werden, in dem die Forderung an die Norm von w und \mathbf{q} genutzt werden

$$\begin{aligned} &\left(\|\mathbf{D} w\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right)^2 + \left(\|w\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)} \right)^2 + \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)}^2 + \left(\|\mathbf{q}\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}\|_{L^2(\Omega)} \right)^2 \\ &\leq \min_{\lambda, \mu, \kappa \in \mathbb{R}^+} (1 + \lambda) \|\mathbf{D} w\|_{L^2(\Omega)}^2 + (1 + 1/\lambda) \|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + (1 + \mu) \|w\|_{L^2(\Omega)}^2 + \|\operatorname{div} \mathbf{q}\|_{L^2(\Omega)}^2 \\ &\quad + (1 + 1/\mu) \|u\|_{L^2(\Omega)}^2 + (1 + \kappa) \|\mathbf{q}\|_{L^2(\Omega)}^2 + (1 + 1/\kappa) \|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 \\ &\leq \max \left\{ \min_{\lambda, \mu, \kappa \in \mathbb{R}^+} \{2 + 1/\lambda + 1/\kappa, 1 + 1/\mu, 1 + \mu, 1 + \lambda, 1 + \kappa, 1\} \right\} \\ &\quad \left(\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 + \|w\|_{H^1(\Omega)}^2 + \|\mathbf{q}\|_{H(\operatorname{div})}^2 \right) \\ &\leq 3 \left(\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 + \|s\|_{H^{1/2}(\partial\mathcal{T})}^2 + \|t\|_{H^{-1/2}(\partial\mathcal{T})}^2 \right) \\ &= 3 \|x\|_X^2. \end{aligned}$$

Die Konstante 3 ergibt sich dabei mit $\lambda = \kappa = 2$ und $\mu = 1$. Es gilt also

$$|b(x, y)| \leq \sqrt{3} \|x\|_X \|y\|_Y. \quad \square$$

Des Weiteren ist der Kern der Bilinearform b , $N := \{y \in Y : b(\cdot, y) = 0 \in X^*\}$, zu untersuchen, um (B1) zu zeigen.

Lemma 3.5. *Für Bilinearform b aus Definition 3.2 gilt*

$$N = \{y \in Y : b(\cdot, y) = 0 \in X^*\} = \{0\}.$$

Beweis. Sei $y = (\boldsymbol{\tau}, v) \in Y$, sodass $b(x, y) = 0$ für jedes $x = (\boldsymbol{\sigma}, u, s, t) \in X$.

Die Wahl $x = (0, u, 0, 0) \in X$ führt für jedes $u \in C_0^\infty(\Omega; \mathbb{R}^n) \subseteq L^2(\Omega; \mathbb{R}^n)$ zu der Gleichung $0 = \int_\Omega \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \cdot u \, dx$. Nach dem Fundamentallemma der Variationsrechnung

gilt daher $\operatorname{div}_{\text{NC}} \boldsymbol{\tau} \equiv 0$, d.h. insbesondere $\boldsymbol{\tau} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$ mit $\operatorname{div} \boldsymbol{\tau} = 0$. Also gilt nach (2.6) $\gamma_\nu^\mathcal{T} \boldsymbol{\tau} \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ und damit $x = (0, 0, 0, \gamma_\nu^\mathcal{T} \boldsymbol{\tau}) \in X$. Daraus ergibt sich

$$0 = \langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, \gamma_0^\mathcal{T} v \rangle_{\partial\mathcal{T}} = \int_\Omega v \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau} \, dx + \int_\Omega \mathbf{D}_{\text{NC}} v : \boldsymbol{\tau} \, dx = \int_\Omega \mathbf{D}_{\text{NC}} v : \boldsymbol{\tau} \, dx.$$

Mit $x = (\boldsymbol{\tau}, 0, 0, 0) \in X$, da $H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \subseteq L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$, folgt

$$0 = \int_\Omega \boldsymbol{\tau} : \mathbf{D}_{\text{NC}} v \, dx + \int_\Omega \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\tau} \, dx = \int_\Omega \operatorname{dev} \boldsymbol{\tau} : \boldsymbol{\tau} \, dx = \|\operatorname{dev} \boldsymbol{\tau}\|_{L^2(\Omega)}^2,$$

also $\operatorname{dev} \boldsymbol{\tau} = 0$. Nach dem Tr-Div-Dev Lemma, Lemma 2.18, für $\boldsymbol{\tau} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$, gilt außerdem

$$\|\operatorname{tr} \boldsymbol{\tau}\|_{L^2(\Omega)} \leq C_{\text{td}} (\|\operatorname{dev} \boldsymbol{\tau}\|_{L^2(\Omega)} + \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(\Omega)}) = 0,$$

also auch $\operatorname{tr} \boldsymbol{\tau} = 0$ und somit $\boldsymbol{\tau} \equiv 0$.

Damit ergibt sich für jedes $x = (\boldsymbol{\sigma}, 0, 0, 0) \in X$,

$$0 = \int_\Omega \boldsymbol{\sigma} : \mathbf{D}_{\text{NC}} v \, dx. \quad (3.11)$$

Sei $w \in H_0^1(\Omega; \mathbb{R}^n)$ die eindeutig bestimmt schwache Lösung des Laplace Problems

$$\operatorname{div} \mathbf{D} w = \Delta w = v \text{ in } \Omega \quad \text{and} \quad w = 0 \text{ auf } \Gamma.$$

Dann gilt $\mathbf{q} := \mathbf{D} w \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ mit $\operatorname{div} \mathbf{q} = v$ und $\gamma_\nu^\mathcal{T} \mathbf{q} \in H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$. Nach dem Gaußschen Integralsatz ist insbesondere

$$\int_\Omega \operatorname{tr} \mathbf{q} \, dx = \int_\Omega \operatorname{div} w \, dx = \int_{\partial\Omega} w \cdot \nu \, dx = 0,$$

also $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \subseteq L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$. Mit $x = (0, 0, 0, \gamma_\nu^\mathcal{T} \mathbf{q}) \in X$ ergibt sich aus Satz 2.8 und (3.11)

$$0 = \langle \gamma_\nu^\mathcal{T} \mathbf{q}, \gamma_0^\mathcal{T} v \rangle_{\partial\mathcal{T}} = \int_\Omega v \cdot \operatorname{div} \mathbf{q} \, dx + \int_\Omega \mathbf{q} : \mathbf{D}_{\text{NC}} v \, dx = \|v\|_{L^2(\Omega)},$$

also $v = 0$.

Damit ist $N = \{0\}$ bewiesen. □

Es gilt also (B1) für die gegebene Bilinearform und es bleibt (B2) zu beweisen. Dazu wird folgendes Lemma, [CDG15b, Lemma 2.1], verwendet. Dieses Lemma ist in einer speziell an das vorliegende Setting angepassten Form auch in [CDG15a, Theorem 3.1] zu finden.

Lemma 3.6 (Splitting Lemma). *Es seien X und Y Banachräume mit Unterräumen $X_1, X_2 \subseteq X$ und $Y_1 \subseteq Y$, so dass X (nicht notwendigerweise direkt) zerlegbar ist in $X = X_1 + X_2$. Des Weiteren gelte*

$$b(x_2, y_1) = 0 \quad \text{für alle } x_2 \in X_2, y_1 \in Y_1$$

sowie die Inf-Sup-Bedingungen

$$0 < \beta_1 := \inf_{x_1 \in S(X_1)} \sup_{y_1 \in S(Y_1)} b(x_1, y_1) \quad \text{und} \quad 0 < \beta_2 := \inf_{x_2 \in S(X_2)} \sup_{y \in S(Y)} b(x_2, y).$$

In diesem Fall gilt (B2) mit

$$\beta \geq \frac{\beta_1 \beta_2}{\beta_1 + \beta_2 + \|b\|}.$$

Dieses Lemma wird mit der Zerlegung $X_1 := L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times \{0\} \times \{0\}$ und $X_2 := \{0\} \times \{0\} \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n) \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ und dem Unterraum $Y_1 := H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times H_0^1(\Omega; \mathbb{R}^n)$ verwendet um (B2) zu zeigen. Dafür wird jedoch eine Hilfsbetrachtung benötigt.

Lemma 3.7. *Es gilt die folgende Inf-Sup-Bedingung*

$$0 < \hat{\beta} := \inf_{(\sigma, u) \in S(\hat{X})} \sup_{(\tau, v) \in S(Y_1)} \left(\int_{\Omega} \sigma : D v \, dx + \int_{\Omega} \operatorname{dev} \sigma : \tau \, dx + \int_{\Omega} u \cdot \operatorname{div} \tau \, dx \right) \quad (3.12)$$

mit $\hat{X} = L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n)$ ausgestattet mit der Norm $\|(\sigma, u)\|_{\hat{X}}^2 := \|\sigma\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2$ und $Y_1 = H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times H_0^1(\Omega; \mathbb{R}^n)$ mit der Norm von Y .

Beweis. Es sei $0 \neq (\sigma, u) \in \hat{X}$ vorgegeben. Dazu ist $(\tau, v) \in Y$ zu finden, so dass die Bedingung (3.12) erfüllt ist. Zu diesem Zweck sei $\tilde{F} \in Z^*$ definiert durch

$$\tilde{F}(\rho, w) := \int_{\Omega} \sigma : \rho + u \cdot w \, dx \quad \text{für alle } (\rho, w) \in Z := H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n),$$

also mit Z wie in Bemerkung 2.23. Die Cauchy-Schwarz-Ungleichung zeigt

$$\|\tilde{F}\|_{Z^*} \leq \|\sigma\|_{L^2(\Omega)} + \|u\|_{L^2(\Omega)} \leq \sqrt{2} \|(\sigma, u)\|_{\hat{X}}.$$

Sei für alle $(\tau, v), (\rho, w) \in Z$

$$\tilde{b}((\tau, v), (\rho, w)) := \int_{\Omega} \operatorname{dev} \tau : \rho \, dx + \int_{\Omega} v \cdot \operatorname{div} \rho \, dx + \int_{\Omega} \operatorname{div} \tau \cdot w \, dx,$$

ebenfalls wie in Bemerkung 2.23. Nach ebendieser Bemerkung existiert $(\tau, -v) \in Z$ mit $\tilde{b}((\tau, -v), \cdot) = \tilde{F}$ in Z^* . Es gilt also für jedes $(\rho, w) \in Z$

$$0 = \int_{\Omega} (\sigma - \operatorname{dev} \tau) : \rho \, dx + \int_{\Omega} \operatorname{div} \rho \cdot v \, dx + \int_{\Omega} (u - \operatorname{div} \tau) \cdot w \, dx. \quad (3.13)$$

Das Fundamentallemma der Variationsrechnung liefert $\operatorname{div} \boldsymbol{\tau} = u$. Nach der folgenden Argumentation gilt sogar $v \in H_0^1(\Omega; \mathbb{R}^n)$ mit $D v = \boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}$.

Um dies zu sehen, ist zunächst für alle $\boldsymbol{\varphi} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$ zu zeigen, dass

$$\int_{\Omega} \operatorname{div} \boldsymbol{\varphi} \cdot v \, dx = - \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\varphi} \, dx. \quad (3.14)$$

Betrachte dazu statt $\boldsymbol{\varphi} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$ wie in Bemerkung 2.20

$$\tilde{\boldsymbol{\varphi}} := \boldsymbol{\varphi} - 1/n \left(\int_{\Omega} \operatorname{tr} \boldsymbol{\varphi} \, dx \right) \mathbf{I}_{n \times n} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R}.$$

Es gilt $\operatorname{div} \boldsymbol{\varphi} = \operatorname{div} \tilde{\boldsymbol{\varphi}}$. Außerdem gilt nach Lemma 2.5 und da $\boldsymbol{\sigma} \in L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$

$$\begin{aligned} \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \tilde{\boldsymbol{\varphi}} \, dx &= \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\varphi} \, dx - 1/n \int_{\Omega} \operatorname{tr} \boldsymbol{\varphi} \, dx \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \mathbf{I}_{n \times n} \, dx \\ &= \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\varphi} \, dx - 1/n \int_{\Omega} \operatorname{tr} \boldsymbol{\varphi} \, dx \int_{\Omega} \operatorname{tr} \boldsymbol{\sigma} - 0 \, dx \\ &= \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\varphi} \, dx. \end{aligned}$$

Da $C_0^\infty(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \subseteq H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R}$, gilt nach (3.13) auch (3.14), d.h. es liegt $v \in H^1(\Omega; \mathbb{R}^n)$ mit schwacher Ableitung $D v = \boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}$. Für alle $\boldsymbol{\varphi} \in C_0^\infty(\Omega; \mathbb{R}^{n \times n})$ ergibt sich mit partieller Integration und Anwendung von (3.13), nach analoger Argumentation,

$$\begin{aligned} \int_{\partial\Omega} v \cdot \boldsymbol{\varphi} \boldsymbol{\nu} \, ds &= \int_{\Omega} D v : \boldsymbol{\varphi} \, dx + \int_{\Omega} v \cdot \operatorname{div} \boldsymbol{\varphi} \, dx \\ &= \int_{\Omega} (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) : \boldsymbol{\varphi} \, dx + \int_{\Omega} \operatorname{div} \boldsymbol{\varphi} \cdot v \, dx = 0, \end{aligned}$$

also hat v homogene Randdaten. Insbesondere gilt also $(\boldsymbol{\tau}, v) \in Y_1$.

Wird $(\boldsymbol{\tau}, v) \in Y_1$ in die zu untersuchende Bilinearform aus (3.12) eingesetzt, so ergibt sich nach Lemma 2.5 Punkt (ii)

$$\begin{aligned} \int_{\Omega} \boldsymbol{\sigma} : D v \, dx + \int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot \operatorname{div} \boldsymbol{\tau} \, dx \\ = \int_{\Omega} \boldsymbol{\sigma} : (\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}) \, dx + \int_{\Omega} \boldsymbol{\sigma} : \operatorname{dev} \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot u \, dx \\ = \|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 = \|(\boldsymbol{\sigma}, u)\|_{\hat{X}}^2. \end{aligned}$$

Die Norm $\|(\boldsymbol{\tau}, v)\|_Z$ lässt sich mit der Inf-Sup-Bedingung aus Bemerkung 2.23 abschätzen

$$\gamma \|(\boldsymbol{\tau}, v)\|_Z = \gamma \|(\boldsymbol{\tau}, -v)\|_Z \leq \left\| \tilde{b}((\boldsymbol{\tau}, -v), \cdot) \right\|_{Z^*} = \|F\|_{Z^*} \leq \sqrt{2} \|(\boldsymbol{\sigma}, u)\|_{\hat{X}}.$$

Desweiteren gilt mit der Dreiecksungleichung

$$\|D v\|_{L^2(\Omega)}^2 = \|\boldsymbol{\sigma} - \operatorname{dev} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 \leq 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + 2\|\operatorname{dev} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 \leq 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + 2\|\boldsymbol{\tau}\|_{L^2(\Omega)}^2.$$

Diese letzten beiden Ungleichungen erlauben die Abschätzung

$$\begin{aligned} \|(\boldsymbol{\tau}, v)\|_Y^2 &= \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|v\|_{L^2(\Omega)}^2 + \|D v\|_{L^2(\Omega)}^2 \\ &\leq 3\|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\operatorname{div} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|v\|_{L^2(\Omega)}^2 + 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 \\ &\leq 3\|(\boldsymbol{\tau}, v)\|_Z^2 + 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 \\ &\leq 6/\gamma^2 \left(\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 + \|u\|_{L^2(\Omega)}^2 \right) + 2\|\boldsymbol{\sigma}\|_{L^2(\Omega)}^2 \\ &\leq (6/\gamma^2 + 2) \|(\boldsymbol{\sigma}, u)\|_{\hat{X}}^2. \end{aligned}$$

Nun ergibt sich

$$\begin{aligned} &\sup_{(\boldsymbol{\rho}, w) \in Y_1 \setminus \{0\}} \frac{\int_{\Omega} \boldsymbol{\sigma} : D w \, dx + \int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\rho} \, dx + \int_{\Omega} u \cdot \operatorname{div} \boldsymbol{\rho} \, dx}{\|(\boldsymbol{\rho}, w)\|_Y} \\ &\geq \frac{\int_{\Omega} \boldsymbol{\sigma} : D v \, dx + \int_{\Omega} \operatorname{dev} \boldsymbol{\sigma} : \boldsymbol{\tau} \, dx + \int_{\Omega} u \cdot \operatorname{div} \boldsymbol{\tau} \, dx}{\|(\boldsymbol{\tau}, v)\|_Y} \\ &\geq \frac{\|(\boldsymbol{\sigma}, u)\|_{\hat{X}}^2}{\sqrt{6/\gamma^2 + 2} \|(\boldsymbol{\sigma}, u)\|_{\hat{X}}} \\ &= (6/\gamma^2 + 2)^{-1/2} \|(\boldsymbol{\sigma}, u)\|_{\hat{X}}. \end{aligned}$$

Da $(\boldsymbol{\sigma}, u) \in \hat{X}$ beliebig war, beweist dies die behauptete Inf-Sup-Bedingung mit

$$\hat{\beta} \geq (6/\gamma^2 + 2)^{-1/2} > 0. \quad \square$$

Mit dieser Vorarbeit ist es möglich die kontinuierliche Inf-Sup-Bedingung für b zu beweisen.

Theorem 3.8 (Kontinuierliche Inf-Sup-Bedingung für b). *Es gilt für die Bilinearform b definiert in Gleichung (3.6)*

$$0 < \beta = \inf_{x \in S(X)} \sup_{y \in S(Y)} b(x, y).$$

Beweis. Wie bereits angedeutet wird für diesen Beweis das Splitting Lemma, Lemma 3.6, verwendet. Dabei seien

$$\begin{aligned} X_1 &:= L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times \{0\} \times \{0\} \subseteq X, \\ X_2 &:= \{0\} \times \{0\} \times H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n) \times H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n) \subseteq X, \\ Y_1 &:= Y_1 = H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times H_0^1(\Omega; \mathbb{R}^n) \subseteq Y. \end{aligned}$$

Für diese Räume sind nun die Voraussetzungen von Lemma 3.6 zu prüfen. Es gilt $X = X_1 + X_2$. Für jedes $y_1 := (\boldsymbol{\tau}, v) \in Y_1$ und alle $\mathbf{q} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ und $w \in H_0^1(\Omega; \mathbb{R}^n)$, also jedes $x_2 := (0, 0, \gamma_0^\mathcal{T} w, \gamma_\nu^\mathcal{T} \mathbf{q}) \in X_2$, zeigt sich außerdem mit partieller Integration

$$\begin{aligned} b(x_2, y_1) &= -\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}, \gamma_0^\mathcal{T} w \rangle_{\partial\mathcal{T}} - \langle \gamma_\nu^\mathcal{T} \mathbf{q}, \gamma_0^\mathcal{T} v \rangle_{\partial\mathcal{T}} \\ &= -\int_{\Omega} \mathbf{q} : \mathbf{D} v \, dx - \int_{\Omega} \operatorname{div} \mathbf{q} \cdot v \, dx - \int_{\Omega} \boldsymbol{\tau} : \mathbf{D} w \, dx - \int_{\Omega} \operatorname{div} \boldsymbol{\tau} \cdot w \, dx \\ &= -\int_{\partial\Omega} v \cdot \mathbf{q} \nu \, ds - \int_{\partial\Omega} w \cdot \boldsymbol{\tau} \nu \, ds \\ &= 0. \end{aligned}$$

Die Bedingung

$$0 < \beta_2 := \inf_{x_2 \in S(X_2)} \sup_{y \in S(Y)} b(x_2, y)$$

wird mit Hilfe von Lemma 2.12 und Lemma 2.13 gezeigt. Für jedes $x_2 = (0, 0, s, t) \in X_2$ gilt

$$\begin{aligned} \|x_2\|_X &= \sqrt{\|s\|_{H^{1/2}(\partial\mathcal{T})}^2 + \|t\|_{H^{-1/2}(\partial\mathcal{T})}^2} \\ &\leq \|s\|_{H^{1/2}(\partial\mathcal{T})} + \|t\|_{H^{-1/2}(\partial\mathcal{T})} \\ &= \sup_{\mathbf{q} \in S(H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R})} \langle \gamma_\nu^\mathcal{T} \mathbf{q}, s \rangle_{\partial\mathcal{T}} + \sup_{w \in S(H^1(\mathcal{T}; \mathbb{R}^n))} \langle t, \gamma_0^\mathcal{T} w \rangle_{\partial\mathcal{T}} \\ &\leq \sup_{\substack{\mathbf{q} \in S(H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R}) \\ w \in S(H^1(\mathcal{T}; \mathbb{R}^n))}} -b(x_2, (\mathbf{q}, w)) \\ &= \sup_{\substack{\mathbf{q} \in S(H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})/\mathbb{R}) \\ w \in S(H^1(\mathcal{T}; \mathbb{R}^n))}} \frac{-b(x_2, (\mathbf{q}, w)) \|(\mathbf{q}, w)\|_Y}{\|(\mathbf{q}, w)\|_Y} \\ &\leq \sqrt{2} \sup_{y \in S(Y)} b(x_2, y), \end{aligned}$$

also gilt die Inf-Sup-Bedingung mit $\beta_2 \geq \sqrt{2}$. Die Bedingung

$$0 < \beta_1 := \inf_{x_1 \in S(X_1)} \sup_{y_1 \in S(Y_1)} b(x_1, y_1)$$

gilt nach Lemma 3.7 mit $\beta_1 = \hat{\beta} \geq (6/\gamma^2 + 2)^{-1/2}$, wobei γ die Inf-Sup-Konstante aus Bemerkung 2.23 ist. Damit lässt sich die kontinuierliche Inf-Sup-Konstante wie in Lemma 3.6 durch

$$\beta \geq \frac{\beta_1 \beta_2}{\beta_1 + \beta_2 + \|b\|}$$

abschätzen. □

Der folgende Abschnitt widmet sich nach einigen Vorbereitungen dem Nachweis von (B4).

3.3 Diskrete Inf-Sup-Bedingung

Die vorgeschlagenen diskreten Räume lauten

$$\begin{aligned} X_h &:= P_0(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_0(\mathcal{T}; \mathbb{R}^n) \times S_0^1(\mathcal{E}; \mathbb{R}^n) \times P_0(\mathcal{E}; \mathbb{R}^n), \\ Y_h &:= RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_1(\mathcal{T}; \mathbb{R}^n). \end{aligned}$$

Damit lautet das diskrete zu lösende Problem analog zu Problem 6 wie folgt.

Problem 10 (Praktische dPG-Methode im gegebenen Setting)

Zu gegebenem $f \in L^2(\Omega; \mathbb{R}^n)$ und $\hat{s} \in H^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ ist $x_h \in X_h$ gesucht mit

$$x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*},$$

wobei b und F wie in Definition 3.2 definiert sind.

Die diskreten Spuren in $S_0^1(\mathcal{E}; \mathbb{R}^n)$ bzw. $P_0(\mathcal{E}; \mathbb{R}^n)$ können eindeutig zu globalen Funktionen in $S_0^1(\mathcal{T}; \mathbb{R}^n)$ bzw. $RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ fortgesetzt werden. Mit Hilfe dieser Fortsetzungen wird in Bemerkung 3.11 eine äquivalente Bilinearform definiert, die den Beweis der diskreten Inf-Sup-Bedingung vereinfacht. Der Vollständigkeit halber wird in den folgenden Lemmata die Normäquivalenz zwischen den diskreten Spuren und ihren Fortsetzungen bewiesen.

Lemma 3.9. *Die Spur $s_1 := \gamma_0 w_c \in S_0^1(\mathcal{E}; \mathbb{R}^n)$ von einer Funktion $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^n)$ gehört zu $H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ und genügt folgender Normäquivalenz*

$$\|s_1\|_{H_0^{1/2}(\partial\mathcal{T})} \approx \|w_c\|_{H_0^1(\Omega)}.$$

Beweis. Es gilt, die nicht triviale Richtung

$$\|w_c\|_{H_0^1(\Omega)} \lesssim \|s_1\|_{H_0^{1/2}(\partial\mathcal{T})} = \inf_{\substack{v \in H_0^1(\Omega, \mathbb{R}^n) \\ \gamma_0 v = s_1}} \|v\|_{H_0^1(\Omega)},$$

zu beweisen. Für alle $v \in H_0^1(\Omega; \mathbb{R}^n)$ mit Spur $\gamma_0 v = s_1$ und für die elementweise konstante Ableitung $D w_c$ gilt auf jedem n -Simplex $T \in \mathcal{T}$

$$D w_c|_T = |T|^{-1} \int_T D w_c \, dx = |T|^{-1} \int_{\partial T} \nu_T s_1 \, ds = |T|^{-1} \int_T D v \, dx.$$

Daher gilt $D w_c = \Pi_0 D v$ und für die Norm

$$\|D w_c\|_{L^2(\Omega)} = \|\Pi_0 D v\|_{L^2(\Omega)} \leq \|D v\|_{L^2(\Omega)} \leq \|v\|_{H_0^1(\Omega)}.$$

Zusammen mit der Friedrichs Ungleichung aus Lemma 2.16, $\|w_c\|_{L^2(\Omega)} \leq C_F \|D w_c\|_{L^2(\Omega)}$, ergibt sich

$$\|w_c\|_{H_0^1(\Omega)} \leq \sqrt{1 + C_F^2} \|v\|_{H_0^1(\Omega)}.$$

Da $v \in H_0^1(\Omega; \mathbb{R}^n)$ eine beliebige Fortsetzung ist, beweist dies das Lemma. Genauer gilt mit $C_s := \sqrt{1 + C_F^2}$

$$\|s_1\|_{H_0^{1/2}(\partial\mathcal{T})} \leq \|w_c\|_{H_0^1(\Omega)} \leq C_s \|s_1\|_{H_0^{1/2}(\partial\mathcal{T})}. \quad (3.15)$$

□

Lemma 3.10. *Die Normalenspur $t_0 = \gamma_\nu \mathbf{q}_{RT} \in P_0(\mathcal{E}; \mathbb{R}^n)$ der Raviart-Thomas Funktion $\mathbf{q}_{RT} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ gehört zu $H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ und genügt folgender Normäquivalenz*

$$\|t_0\|_{H^{-1/2}(\partial\mathcal{T})} \approx \|\mathbf{q}_{RT}\|_{H(\text{div}, \mathcal{T})}.$$

Beweis. Es gilt, die nicht triviale Richtung

$$\|\mathbf{q}_{RT}\|_{H(\text{div}, \mathcal{T})} \lesssim \|t_0\|_{H^{-1/2}(\partial\mathcal{T})} = \inf_{\substack{\boldsymbol{\rho} \in H(\text{div}, \mathcal{T}; \mathbb{R}^{n \times n}) \\ \gamma_\nu \boldsymbol{\rho} = t_0}} \|\boldsymbol{\rho}\|_{H(\text{div}, \mathcal{T})},$$

zu beweisen. Für alle $\boldsymbol{\rho} \in H(\text{div}, \mathcal{T}; \mathbb{R}^n)$ mit Normalenspur $\gamma_\nu \boldsymbol{\rho} = t_0$ und die elementweise konstante Divergenz $\text{div } \mathbf{q}_{RT}$ gilt auf jedem n-Simplex $T \in \mathcal{T}$

$$\text{div } \mathbf{q}_{RT}|_T = |T|^{-1} \int_T \text{div } \mathbf{q}_{RT} \, dx = |T|^{-1} \int_{\partial T} t \, ds = |T|^{-1} \int_T \text{div } \boldsymbol{\rho} \, dx.$$

Also ist $\text{div } \mathbf{q}_{RT} = \Pi_0 \text{div } \boldsymbol{\rho}$ und für die Norm gilt

$$\|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)} = \|\Pi_0 \text{div } \boldsymbol{\rho}\|_{L^2(\Omega)} \leq \|\text{div } \boldsymbol{\rho}\|_{L^2(\Omega)} \leq \|\boldsymbol{\rho}\|_{H(\text{div}, \mathcal{T})}. \quad (3.16)$$

Außerdem gilt für jedes $x \in T \in \mathcal{T}$

$$\begin{aligned} \mathbf{q}_{RT}|_T(x) &= \Pi_0 \mathbf{q}_{RT} + 1/n \, \text{div } \mathbf{q}_{RT} \otimes (x - \text{mid}(T)) \\ &= D \left(\Pi_0 \mathbf{q}_{RT} x + 1/(2n) \, \text{div } \mathbf{q}_{RT} |x - \text{mid}(T)|^2 \right). \end{aligned}$$

Daher existiert für alle $T \in \mathcal{T}$ ein $g_T \in H^1(T; \mathbb{R}^n)$, so dass $\mathbf{q}_{RT} = D g_T$ und per Normierung $\int_T g_T \, dx = 0$. Mit Hilfe zweier partieller Integrationen und der Poincaré Ungleichung aus Lemma 2.17 mit Konstante $C_P \leq \text{diam}(\Omega)/\pi$, zeigt sich

$$\begin{aligned} \|D g_T\|_{L^2(T)}^2 &= \int_T D g_T : \mathbf{q}_{RT} \, dx \\ &= - \int_T g_T \cdot \text{div } \mathbf{q}_{RT} \, dx + \int_{\partial T} g_T \cdot t_0 \, ds \end{aligned}$$

$$\begin{aligned}
&= - \int_T g_T \cdot \Pi_0 \operatorname{div} \boldsymbol{\rho} \, dx + \int_{\partial T} g_T \cdot \gamma_\nu \boldsymbol{\rho} \, ds \\
&= \int_T g_T \cdot (1 - \Pi_0) \operatorname{div} \boldsymbol{\rho} \, dx + \int_T \boldsymbol{\rho} : \mathbf{D} g_T \, dx \\
&\leq C_P \|\mathbf{D} g_T\|_{L^2(T)} \|(1 - \Pi_0) \operatorname{div} \boldsymbol{\rho}\|_{L^2(T)} + \|\mathbf{D} g_T\|_{L^2(T)} \|\boldsymbol{\rho}\|_{L^2(T)}.
\end{aligned}$$

Daher gilt

$$\|\mathbf{q}_{\text{RT}}\|_{L^2(T)} = \|\mathbf{D} g_T\|_{L^2(T)} \leq C_P \|(1 - \Pi_0) \operatorname{div} \boldsymbol{\rho}\|_{L^2(T)} + \|\boldsymbol{\rho}\|_{L^2(T)}.$$

Das Aufsummieren über alle Dreiecke und (3.16) zeigen

$$\begin{aligned}
\|\mathbf{q}_{\text{RT}}\|_{H(\operatorname{div}, \mathcal{T})}^2 &= \|\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\
&\leq \|\operatorname{div} \boldsymbol{\rho}\|_{L^2(\Omega)}^2 + \left(C_P \|\operatorname{div} \boldsymbol{\rho}\|_{L^2(\Omega)} + \|\boldsymbol{\rho}\|_{L^2(\Omega)} \right)^2 \\
&\leq \min_{\lambda \in \mathbb{R}^+} \left\{ (1 + \lambda) C_P^2 \|\operatorname{div} \boldsymbol{\rho}\|_{L^2(\Omega)}^2 + (1 + 1/\lambda) \|\boldsymbol{\rho}\|_{L^2(\Omega)}^2 \right\} \\
&\leq \min_{\lambda \in \mathbb{R}^+} \max \left\{ 1 + (1 + \lambda) C_P^2, 1 + 1/\lambda \right\} \|\boldsymbol{\rho}\|_{H(\operatorname{div}, \mathcal{T})}^2 \\
&\leq C_t^2 \|\boldsymbol{\rho}\|_{H(\operatorname{div}, \mathcal{T})}^2,
\end{aligned}$$

wobei die Wahl $\lambda = 1/2 \left(\sqrt{C_P^2 + 4/C_P} - 1 \right)$, so dass $1 + (1 + \lambda) C_P^2 = (1 + 1/\lambda)$, die Konstante

$$C_t = \sqrt{1 + \left(C_P + \sqrt{C_P^2 + 4} \right) C_P/2}$$

ergibt. Da $\boldsymbol{\rho} \in H(\operatorname{div}, \mathcal{T}; \mathbb{R}^{n \times n})$ eine beliebige Fortsetzung war beweist dies das Lemma und insbesondere gilt

$$\|t_0\|_{H^{-1/2}(\partial \mathcal{T})} \leq \|\mathbf{q}_{\text{RT}}\|_{H(\operatorname{div}, \mathcal{T})} \leq C_t \|t_0\|_{H^{-1/2}(\partial \mathcal{T})}. \quad (3.17)$$

□

Nun folgt wie angekündigt die Definition einer alternativen Bilinearform.

Bemerkung 3.11 (Alternative Bilinearform) Die Fortsetzung der Spurterme $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^n)$ bzw. $t_0 \in P_0(\mathcal{E}; \mathbb{R}^n)$ durch $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^n)$ bzw. $\mathbf{q}_{\text{RT}} \in RT_0(\mathcal{T}; \mathbb{R}^{n \times n})$ erlaubt eine weitere partielle Integration, die zu folgender alternativer Bilinearform führt

$$\begin{aligned}
\hat{b}(\hat{x}_h, y_h) &= \int_\Omega \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}} v_1 \, dx + \int_\Omega \operatorname{dev} \boldsymbol{\sigma}_0 : \boldsymbol{\tau}_{\text{RT}} \, dx + \int_\Omega u_0 \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau}_{\text{RT}} \, dx \\
&\quad - \int_\Omega (v_1 \cdot \operatorname{div} \mathbf{q}_{\text{RT}} + \mathbf{q}_{\text{RT}} : \mathbf{D}_{\text{NC}} v_1) \, dx - \int_\Omega (\boldsymbol{\tau}_{\text{RT}} : \mathbf{D} w_c + w_c \cdot \operatorname{div}_{\text{NC}} \boldsymbol{\tau}_{\text{RT}}) \, dx.
\end{aligned}$$

Hier ist der Testraum unverändert, d.h. $y_h \in Y_h$ wie in (2.12), und die Ansatzfunktion $\hat{x}_h = (\boldsymbol{\sigma}_0, u_0, w_c, \mathbf{q}_{\text{RT}}) \in \hat{X}_h$ stammt aus dem folgenden Raum

$$\hat{X}_h := P_0(\mathcal{T}; \mathbb{R}^{n \times n}) / \mathbb{R} \times P_0(\mathcal{T}; \mathbb{R}^n) \times S_0^1(\mathcal{T}, \mathbb{R}^n) \times RT_0(\mathcal{T}; \mathbb{R}^{n \times n}). \quad (3.18)$$

Dieser Raum sei mit der Norm des kontinuierlichen Äquivalents $\hat{X} := L^2(\Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \times L^2(\Omega; \mathbb{R}^n) \times H_0^1(\Omega; \mathbb{R}^n) \times H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n})$ versehen,

$$\|x_h\|_{\hat{X}}^2 = \|\sigma_0\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2 + \|w_c\|_{H_0^1(\Omega)}^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\operatorname{div}, \Omega)}^2. \quad (3.19)$$

Die diskrete Inf-Sup-Bedingung wird nun mit Hilfe des Tr-dev-div Lemmas, Lemma 2.18, und der Zerlegung aus Lemma 2.21, aus Paragraph 2.5, zunächst für die alternative Bilinearform aus Bemerkung 3.11 bewiesen. Dabei kann die Inf-Sup-Konstante direkt berechnet werden.

Theorem 3.12. *In den diskreten Räumen aus (2.11)–(2.12), erfüllt die Bilinearform b (3.6) die Inf-Sup-Bedingung*

$$1 \lesssim \beta_h := \inf_{x_h \in \hat{X}_h \setminus \{0\}} \sup_{y_h \in Y_h} \frac{b(x_h, y_h)}{\|y_h\|_Y \|x_h\|_X}.$$

Die Konstante β_h wird im Beweis berechnet.

Beweis. Zunächst wird wie angekündigt mit der alternativen Bilinearform aus Bemerkung 3.11 gearbeitet.

Schritt 1. Diskrete Testfunktionen. Zu einem gegebenen $\hat{x}_h = (\sigma_0, u_0, w_c, \mathbf{q}_{\text{RT}}) \in \hat{X}_h \setminus \{0\}$ wie in (3.18) wird folgendes $\tilde{y}_h = (\tilde{\tau}_{\text{RT}}, \tilde{v}_1)$ gewählt

$$\begin{aligned} \tilde{\tau}_{\text{RT}} &:= \operatorname{dev} \sigma_0 - \operatorname{D} w_c + 1/n (u_0 - \Pi_0 w_c) \otimes (\cdot - \operatorname{mid}(T)), \\ \tilde{v}_1 &:= -\operatorname{div} \mathbf{q}_{\text{RT}} + (\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}) (\cdot - \operatorname{mid}(T)). \end{aligned}$$

Dies ist eine zulässige Testfunktion, denn offensichtlich liegen $\tilde{\tau}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{n \times n})$ und $\tilde{v}_1 \in P_1(\mathcal{T}; \mathbb{R}^n)$ und die globale Nebenbedingung $\int_{\Omega} \operatorname{tr} \tilde{\tau}_{\text{RT}} \, dx = 0$ ist erfüllt. Letzteres gilt, da $\operatorname{dev} \sigma_0$ per Definition spurfrei ist, $u_0 - \Pi_0 w_c$ elementweise konstant und nach (2.2) $\int_T \cdot - \operatorname{mid}(T) \, dx = 0$ sowie $\int_{\Omega} \operatorname{div} w_c \, dx = \int_{\partial\Omega} w_c \cdot \nu \, ds = 0$. Außerdem gilt, $\operatorname{div}_{\text{NC}} \tilde{\tau}_{\text{RT}} = u_0 - \Pi_0 w_c$, $\operatorname{D}_{\text{NC}} \tilde{v}_1 = \sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}$ und $\Pi_0 \tilde{v}_1 = -\operatorname{div} \mathbf{q}_{\text{RT}}$.

Schritt 2. Es gilt

$$\begin{aligned} \hat{b}(\hat{x}_h, \tilde{y}_h) &= \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\operatorname{dev} \sigma_0 - \operatorname{D} w_c\|_{L^2(\Omega)}^2 \\ &\quad + \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned}$$

Beweis von Schritt 2. Die Bilinearform \hat{b} mit $\hat{x}_h = (\sigma_0, u_0, w_c, \mathbf{q}_{\text{RT}})$ und $y_h = (\tau_{\text{RT}}, v_1)$ lautet

$$\begin{aligned} \hat{b}(\hat{x}_h, y_h) &= \int_{\Omega} (u_0 - w_c) \cdot \operatorname{div}_{\text{NC}} \tau_{\text{RT}} \, dx + \int_{\Omega} (\operatorname{dev} \sigma_0 - \operatorname{D} w_c) : \tau_{\text{RT}} \, dx \\ &\quad + \int_{\Omega} (\sigma_0 - \mathbf{q}_{\text{RT}}) : \operatorname{D}_{\text{NC}} v_1 \, dx - \int_{\Omega} v_1 \cdot \operatorname{div} \mathbf{q}_{\text{RT}} \, dx. \end{aligned}$$

Das Einsetzen der in *Schritt 1.* zu $\hat{x}_h \in \hat{X}_h$ gewählten Testfunktion $\tilde{y}_h \in Y_h$ ergibt nach obigen Überlegungen zu $\text{div}_{\text{NC}} \tilde{\tau}_{\text{RT}}$, $\text{D}_{\text{NC}} \tilde{v}_1$ und $\Pi_0 \tilde{v}_1$

$$\begin{aligned} \hat{b}(\hat{x}_h, \tilde{y}_h) &= \int_{\Omega} (u_0 - w_c) \cdot (u_0 - \Pi_0 w_c) \, dx + \int_{\Omega} (\sigma_0 - \mathbf{q}_{\text{RT}}) : (\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}) \, dx \\ &\quad + \int_{\Omega} (\text{dev } \sigma_0 - \text{D } w_c) : (\text{dev } \sigma_0 - \text{D } w_c + (u_0 - \Pi_0 w_c) \otimes (\cdot - \text{mid}(T)) / n) \, dx \\ &\quad - \int_{\Omega} \left(-\text{div } \mathbf{q}_{\text{RT}} + (\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}) (\cdot - \text{mid}(T)) \right) \cdot \text{div } \mathbf{q}_{\text{RT}} \, dx. \end{aligned}$$

Daraus folgt mit (2.10) und (2.2), genauer weil $(\cdot - \text{mid}(T))$ senkrecht zu den elementweise konstanten Funktionen $P_0(\mathcal{T})$, wie $\text{dev } \sigma_0$, $\text{D } w_c$ und $\text{div } \mathbf{q}_{\text{RT}}$ steht, sofort die behauptete Darstellung von \hat{b} . \square

Um die diskrete Inf-Sup-Bedingung zu beweisen werden nun $\|\hat{x}_h\|_{\hat{X}}$ und $\|\tilde{y}_h\|_Y$ gegen die Bilinearform $\hat{b}(\hat{x}_h, \tilde{y}_h)$ in der soeben entwickelten Form abgeschätzt.

Schritt 3. Es gilt $\|\tilde{y}_h\|_Y^2 \leq \left(1 + h_{\max}^2 n^2 / (n+1)^2\right) \hat{b}(\hat{x}_h, \tilde{y}_h)$.

Beweis von Schritt 3. Die Gesetzmäßigkeiten aus (2.2) sowie die Gleichungen für $\text{div}_{\text{NC}} \tilde{\tau}_{\text{RT}}$ und $\text{D}_{\text{NC}} \tilde{v}_1$ ergeben

$$\begin{aligned} \|\tilde{y}_h\|_Y^2 &= \|\tilde{\tau}_{\text{RT}}\|_{H(\text{div}, \mathcal{T})}^2 + \|\tilde{v}_1\|_{H^1(\mathcal{T})}^2 \\ &= \left\| \text{dev } \sigma_0 - \text{D } w_c + 1/n (u_0 - \Pi_0 w_c) \otimes (\cdot - \text{mid}(T)) \right\|_{L^2(\Omega)}^2 + \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 \\ &\quad + \left\| -\text{div } \mathbf{q}_{\text{RT}} + (\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}) (\cdot - \text{mid}(T)) \right\|_{L^2(\Omega)}^2 + \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &\leq \|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)}^2 + \left(1/n^2 \|\cdot - \text{mid}(T)\|_{L^\infty(\Omega)}^2 + 1\right) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 \\ &\quad + \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \left(\|\cdot - \text{mid}(T)\|_{L^\infty(\Omega)}^2 + 1\right) \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &\leq \|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)}^2 + \left(h_{\max}^2 / (n+1)^2 + 1\right) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 \\ &\quad + \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \left(h_{\max}^2 n^2 / (n+1)^2 + 1\right) \|\sigma_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2. \end{aligned}$$

Zusammen mit der Gleichheit aus *Schritt 2.* folgt die Behauptung. \square

Der zentrale Schritt, um eine ähnliche Abschätzung für $\|\hat{x}_h\|_{\hat{X}}$ herzuleiten, ist die ein wenig technische Abschätzung von $\|w_c\| := \|\text{D } w_c\|_{L^2(\Omega)}$ gegen $\hat{b}(\hat{x}_h, \tilde{y}_h)$. Dabei sind die Zerlegung aus Lemma 2.21, und das Tr-Div-Dev Lemma, Lemma 2.18, entscheidend.

Schritt 4. Es gilt $\|w_c\|^2 \leq C_{w_c} \hat{b}(\hat{x}_h, \tilde{y}_h)$.

Beweis von Schritt 4. Es wird die Zerlegung

$$\|w_c\|^2 \leq \min_{\lambda \in \mathbb{R}^+} \left\{ (1 + \lambda) \|w_c + z\|^2 + (1 + 1/\lambda) \|z\|^2 \right\} \quad (3.20)$$

betrachtet. Die Kontrolle beider Terme auf der rechten Seite durch $B^2 := \hat{b}(\hat{x}_h, \tilde{y}_h)$ gelingt, wenn $z \in H_0^1(\Omega; \mathbb{R}^n)$ divergenzfrei und $b \in L^2(\Omega; \mathbb{R}^{n \times n})$ so gewählt werden, dass

$$\operatorname{dev}(\boldsymbol{\sigma}_0 - \operatorname{D} w_c) = \operatorname{D} z + \operatorname{dev} b \quad (3.21)$$

eine orthogonale Zerlegung ist. Eine solche Wahl ist nach Lemma 2.21 möglich. Die Normgleichheit aus Lemma 2.5 Punkt (iii) ermöglicht mit dieser orthogonalen Zerlegung folgende Abschätzung

$$\begin{aligned} \|\operatorname{dev} \boldsymbol{\sigma}_0 - \operatorname{D} w_c\|_{L^2(\Omega)}^2 &= 1/n \|\operatorname{div} w_c\|_{L^2(\Omega)}^2 + \|\operatorname{dev}(\operatorname{D} w_c - \operatorname{dev} \boldsymbol{\sigma}_0)\|_{L^2(\Omega)}^2 \\ &= 1/n \|\operatorname{div} w_c\|_{L^2(\Omega)}^2 + \|z\|^2 + \|\operatorname{dev} b\|_{L^2(\Omega)}^2. \end{aligned}$$

Aus *Schritt 2.* folgt die Beschränkung der Summe und damit aller Summanden nach oben durch $\|\operatorname{dev} \boldsymbol{\sigma}_0 - \operatorname{D} w_c\|_{L^2(\Omega)}^2 \leq B^2$. Der eine Summand in Gleichung (3.20) ist also bereits kontrolliert. Lemma 2.5 Punkt (iii) und die Divergenzfreiheit von $z \in H_0^1(\Omega; \mathbb{R}^n)$ ermöglichen folgende Zerlegung des zweiten Summanden

$$\begin{aligned} A^2 &:= \|w_c + z\|^2 = \|\operatorname{dev} \operatorname{D}(w_c + z)\|_{L^2(\Omega)}^2 + 1/n \|\operatorname{div} w_c\|_{L^2(\Omega)}^2 \\ &\leq \|\operatorname{dev} \operatorname{D}(w_c + z)\|_{L^2(\Omega)}^2 + \|\operatorname{dev} \boldsymbol{\sigma}_0 - \operatorname{D} w_c\|_{L^2(\Omega)}^2 \\ &\leq \|\operatorname{dev} \operatorname{D}(w_c + z)\|_{L^2(\Omega)}^2 + B^2 \end{aligned} \quad (3.22)$$

Die Rechenregel für den deviatorischen Anteil aus Lemma 2.5 sowie Gleichung (3.21) erlauben

$$\begin{aligned} \|\operatorname{dev} \operatorname{D}(w_c + z)\|_{L^2(\Omega)}^2 &= \int_{\Omega} (\operatorname{dev} \operatorname{D}(w_c + z)) : \operatorname{D}(w_c + z) \, dx \\ &= \int_{\Omega} (\operatorname{dev} \boldsymbol{\sigma}_0 - \operatorname{dev} b) : \operatorname{D}(w_c + z) \, dx. \end{aligned}$$

Die Cauchy-Schwarz-Ungleichung, die Definition und Rechenregeln für den Deviator und die Divergenzfreiheit von z ermöglichen weiter

$$\begin{aligned} \|\operatorname{dev} \operatorname{D}(w_c + z)\|_{L^2(\Omega)}^2 &\leq \|\operatorname{dev} b\|_{L^2(\Omega)} \|w_c + z\| + \int_{\Omega} \operatorname{dev}(\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}) : \operatorname{D}(w_c + z) \, dx \\ &\quad + \int_{\Omega} \mathbf{q}_{\text{RT}} : (\operatorname{D}(w_c + z) - \operatorname{div} w_c / n \, \mathbf{I}_{n \times n}) \, dx \\ &\leq (\|\operatorname{dev} b\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}) \|w_c + z\| \\ &\quad + \int_{\Omega} \mathbf{q}_{\text{RT}} : \operatorname{D}(w_c + z) \, dx - \frac{1}{n} \int_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \operatorname{div} w_c \, dx. \end{aligned} \quad (3.23)$$

Eine stückweise partielle Integration ergibt

$$\int_{\Omega} \mathbf{q}_{\text{RT}} : \operatorname{D}(w_c + z) \, dx = - \int_{\Omega} (w_c + z) \cdot \operatorname{div} \mathbf{q}_{\text{RT}} \, dx,$$

da die normalen Komponenten der Raviart-Thomas Funktionen stückweise konstant

und stetig sind, während die Sprünge der stetigen Funktionen $z \in H_0^1(\Omega; \mathbb{R}^n)$ und $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^n)$ verschwinden und somit keine Randintegrale auftreten. Da außerdem $\int_{\Omega} \operatorname{div} w_c \, dx = 0$ gilt, kann nach Bemerkung 2.20 (iii) \mathbf{q}_{RT} im letzten Term von (3.23) durch

$$\tilde{\mathbf{q}}_{\text{RT}} := \mathbf{q}_{\text{RT}} - 1/n \left(\oint_{\Omega} \operatorname{tr} \mathbf{q}_{\text{RT}} \, dx \right) \mathbf{I}_{n \times n} \in H(\operatorname{div}, \Omega; \mathbb{R}^{n \times n}) / \mathbb{R} \quad (3.24)$$

ersetzt werden. Durch Einsetzen dieser Erkenntnisse in die Formel für $\|\operatorname{dev} \mathbf{D}(w_c + z)\|_{L^2(\Omega)}^2$, ergibt sich unter Verwendung der Cauchy-Schwarz und der Friedrichsungleichung mit Konstante $C_F \leq \operatorname{width}(\Omega)/\pi$ aus Lemma 2.16

$$\begin{aligned} \|\operatorname{dev} \mathbf{D}(w_c + z)\|_{L^2(\Omega)}^2 &\leq 1/n \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \|\operatorname{div} w_c\|_{L^2(\Omega)} + \|w_c + z\| \\ &\quad \times \left(\|\operatorname{dev} b\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_F \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \right) \\ &\leq 1/n \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \|\operatorname{div} w_c\|_{L^2(\Omega)} + \|w_c + z\| \\ &\quad \times \left(\|\operatorname{dev} \boldsymbol{\sigma}_0 - \mathbf{D} w_c\|_{L^2(\Omega)} + \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_F \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \right) \\ &\leq 1/\sqrt{n} \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} B + A \left(\sqrt{1 + C_F^2} B + \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \right). \end{aligned}$$

Der Faktor $1/\sqrt{n} = \sqrt{n}/n$ vor $\|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} B$ ergibt sich, da $1/n \|\operatorname{div} w_c\|_{L^2(\Omega)}^2 \leq B^2$. Die Konstante $\sqrt{1 + C_F^2}$ ergibt sich aus Bemerkung 2.22 mit $a = \|\operatorname{dev} \boldsymbol{\sigma}_0 - \mathbf{D} w_c\|_{L^2(\Omega)}$, $b = \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}$, $C_1 = K_1 = K_2 = 1$ und $C_2 = C_F$. Die Eigenschaft (2.9) für Raviart-Thomas Funktionen und (2.2) erlauben weiter die Abschätzung

$$\begin{aligned} \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 &= \|(1 - \Pi_0) \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &= 1/n^2 \|\operatorname{div} \mathbf{q}_{\text{RT}}(\cdot - \operatorname{mid}(T))\|_{L^2(\Omega)}^2 + \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &\leq h_{\max}^2/(n+1)^2 \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\boldsymbol{\sigma}_0 - \Pi_0 \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &\leq \max \{1, h_{\max}^2/(n+1)^2\} B^2 \end{aligned} \quad (3.25)$$

Das in (2.20) definierte $\tilde{\mathbf{q}}_{\text{RT}}$ erfüllt die Voraussetzungen des Tr-Div-Dev Lemmas, Lemma 2.18. Zusammen mit der Dreiecksungleichung und Bemerkung 2.20 (i) ergibt sich

$$\begin{aligned} \frac{1}{C_{\text{td}}} \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} &\leq \|\operatorname{dev} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{div} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &= \|\operatorname{dev} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \|\operatorname{dev} (\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}})\|_{L^2(\Omega)} + \|\operatorname{dev} \boldsymbol{\sigma}_0\|_{L^2(\Omega)} + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \max \{1, h_{\max}/(n+1)\} B + \|\operatorname{div} \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\operatorname{dev} \boldsymbol{\sigma}_0\|_{L^2(\Omega)} \\ &\leq \left(1 + \max \{1, h_{\max}/(n+1)\}\right) B + \|\operatorname{dev} \boldsymbol{\sigma}_0\|_{L^2(\Omega)}. \end{aligned}$$

Eine erneute Anwendung der Zerlegung (3.21) erlaubt gemeinsam mit der Dreiecksun-

gleichung und der Definition des Deviators, die Abschätzung des letzten Terms

$$\begin{aligned}
 \|\operatorname{dev} \boldsymbol{\sigma}_0\|_{L^2(\Omega)} &\leq \|D(z + w_c) - 1/n \operatorname{div}(w_c) \mathbf{I}_{n \times n} + \operatorname{dev} b\|_{L^2(\Omega)} \\
 &\leq \|z + w_c\| + \sqrt{1/n} \|\operatorname{div} w_c\|_{L^2(\Omega)} + \|\operatorname{dev} b\|_{L^2(\Omega)} \\
 &\leq A + \sqrt{2} \|\operatorname{dev} \boldsymbol{\sigma}_0 - D w_c\|_{L^2(\Omega)} \\
 &\leq A + \sqrt{2} B.
 \end{aligned}$$

Der Faktor $\sqrt{2}$ ergibt sich wie in Bemerkung 2.22 mit $a = \|\operatorname{div} w_c\|_{L^2(\Omega)}$, $b = \|\operatorname{dev} b\|_{L^2(\Omega)}$, $C_1 = \sqrt{1/n}$, $K_1 = 1/n$ und $C_2 = K_2 = 1$. Insgesamt hat sich damit aus der Ungleichung (3.22) Folgendes ergeben

$$\begin{aligned}
 A^2 &\leq B^2 + A \left(B \sqrt{1 + C_F^2} + \|\boldsymbol{\sigma}_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \right) + B/\sqrt{n} \|\operatorname{tr} \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\
 &\leq B^2 + AB \left(\sqrt{1 + C_F^2} + \max \{1, h_{\max}/(n+1)\} \right) \\
 &\quad + BC_{\text{tdd}}/\sqrt{n} \left(B \left(1 + \max \{1, h_{\max}/(n+1)\} \right) + \|\operatorname{dev} \boldsymbol{\sigma}_0\|_{L^2(\Omega)} \right) \\
 &\leq B^2 \left(1 + C_{\text{tdd}}/\sqrt{n} \left(1 + \max \{1, h_{\max}/(n+1)\} \right) \right) + BC_{\text{tdd}}/\sqrt{n} (A + B\sqrt{2}) \\
 &\quad + AB \left(\sqrt{1 + C_F^2} + \max \{1, h_{\max}/(n+1)\} \right) \\
 &\leq B^2 \left(1 + C_{\text{tdd}}/\sqrt{n} \left(1 + \sqrt{2} + \max \{1, h_{\max}/(n+1)\} \right) \right) \\
 &\quad + AB \left(\sqrt{1 + C_F^2} + C_{\text{tdd}}/\sqrt{n} + \max \{1, h_{\max}/(n+1)\} \right).
 \end{aligned}$$

Mit $C_1 := 1 + C_{\text{tdd}}/\sqrt{n} \left(1 + \sqrt{2} + \max \{1, h_{\max}/(n+1)\} \right)$ und $C_2 := \sqrt{1 + C_F^2} + C_{\text{tdd}}/\sqrt{n} + \max \{1, h_{\max}/(n+1)\}$ liest sich dies

$$A^2 \leq C_1 B^2 + C_2 AB.$$

Aus dieser Ungleichung für A^2 ergibt sich

$$0 \leq (A - C_2 B)^2 = A^2 - 2C_2 AB + C_2^2 B^2 \leq C_1 B^2 - C_2 AB + C_2^2 B^2$$

und damit $A^2 \leq C_1 B^2 + C_2 AB \leq 2C_1 B^2 + C_2^2 B^2 \leq (2C_1 + C_2^2) B^2$. Also garantiert die Wahl $C_{w_c} = \min_{\lambda \in \mathbb{R}^+} \left\{ (1 + \lambda)(2C_1 + C_2^2) + (1 + 1/\lambda) \right\}$

$$\|w_c\|^2 \leq C_{w_c} \hat{b}(\hat{x}_h, \tilde{y}_h).$$

$(1 + \lambda)(2C_1 + C_2^2) + (1 + 1/\lambda)$ wird durch $\lambda = 1/\sqrt{2C_1 + C_2^2}$ minimiert, dies ergibt $C_{w_c} = \left(1 + \sqrt{2C_1 + C_2^2} \right)^2$. \square

Schritt 5. Es gilt $\|\hat{x}_h\|_{\hat{X}}^2 \lesssim \hat{b}(\hat{x}_h, \tilde{y}_h)$.

Beweis von Schritt 5. Hier werden die einzelnen Teile von $\|\hat{x}_h\|_{\hat{X}}^2$ gegen die Terme aus $\hat{b}(\hat{x}_h, y_h)$ in der Form von *Schritt 2.* abgeschätzt. Zunächst erlauben die Friedrichsungleichung 2.16 mit Konstant $C_F \leq \text{width}(\Omega)/\pi$ und die Dreiecksungleichung folgende Abschätzung

$$\begin{aligned} \|\hat{x}_h\|_{\hat{X}}^2 &= \|\sigma_0\|_{L^2(\Omega)}^2 + \|u_0\|_{L^2(\Omega)}^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2 + \|w_c\|_{H^1(\Omega)}^2 \\ &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + \left(\|u_0 - \Pi_0 w_c\|_{L^2(\Omega)} + \|\Pi_0 w_c\|_{L^2(\Omega)} \right)^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2 + (1 + C_F^2) \|w_c\|^2 \\ &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + \left(\|u_0 - \Pi_0 w_c\|_{L^2(\Omega)} + C_F \|w_c\| \right)^2 + (1 + C_F^2) \|w_c\|^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2 \\ &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + \left(1 + C_F^2 \right) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + (2 + 2C_F^2) \|w_c\|^2 + \|\mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2. \end{aligned}$$

Für die Raviart-Thomas Funktion \mathbf{q}_{RT} folgt aus (2.9) und (2.2) die Ungleichung

$$\begin{aligned} \|\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 &= \|(1 - \Pi_0)\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \|\Pi_0\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 \\ &\leq 1/n^2 \|\bullet - \text{mid}(T)\|_{L^\infty(\Omega)}^2 \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + \left(\|\sigma_0 - \Pi_0\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|\sigma_0\|_{L^2(\Omega)} \right)^2 \\ &\leq h_{\max}^2/(n+1)^2 \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + 2\|\sigma_0 - \Pi_0\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}^2 + 2\|\sigma_0\|_{L^2(\Omega)}^2. \end{aligned}$$

Das Einfügen von $\tilde{\mathbf{q}}_{\text{RT}}$ aus (2.20), welches die Anwendung des Tr-Dev-Div Lemmas und der Bemerkung 2.20 erlaubt, sowie die Eigenschaften des Deviators aus Lemma 2.5 ermöglichen $\|\sigma_0\|_{L^2(\Omega)}$ abzuschätzen

$$\begin{aligned} \|\sigma_0\|_{L^2(\Omega)} &\leq \|\sigma_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + \|\tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \|\sigma_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + \|\text{dev } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + \|\text{tr } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \|\sigma_0 - \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + C_{\text{tdd}} \|\text{div } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} + (1 + C_{\text{tdd}}) \|\text{dev } \tilde{\mathbf{q}}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \|\sigma_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_{\text{tdd}} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + (1 + C_{\text{tdd}}) \|\text{dev } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\leq \|\sigma_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_{\text{tdd}} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\quad + (1 + C_{\text{tdd}}) \left(\|\text{dev } (\sigma_0 - \mathbf{q}_{\text{RT}})\|_{L^2(\Omega)} + \|\text{dev } \sigma_0\|_{L^2(\Omega)} \right) \\ &\leq (2 + C_{\text{tdd}}) \|\sigma_0 - \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + C_{\text{tdd}} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + (1 + C_{\text{tdd}}) \|\text{dev } \sigma_0\|_{L^2(\Omega)}. \end{aligned}$$

Weitere Dreiecksungleichungen und eine erneute Anwendung von (2.9) und (2.2) führen zu folgender Darstellung

$$\begin{aligned} \|\sigma_0\|_{L^2(\Omega)} &\leq (2 + C_{\text{tdd}}) \left(\|\sigma_0 - \Pi_0\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + \|(1 - \Pi_0)\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \right) + C_{\text{tdd}} \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} \\ &\quad + (1 + C_{\text{tdd}}) \left(\|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)} + \|w_c\| \right) \\ &\leq (2 + C_{\text{tdd}}) \|\sigma_0 - \Pi_0\mathbf{q}_{\text{RT}}\|_{L^2(\Omega)} + (1 + C_{\text{tdd}}) \left(\|\text{dev } \sigma_0 - \text{D } w_c\|_{L^2(\Omega)} + \|w_c\| \right) \\ &\quad + ((2 + C_{\text{tdd}})h_{\max}/(n+1) + C_{\text{tdd}}) \|\text{div } \mathbf{q}_{\text{RT}}\|_{L^2(\Omega)}. \end{aligned}$$

Es sind somit alle Terme von $\|\hat{x}_h\|_{\hat{X}}$ durch Terme aus *Schritt 2.* und $\|w_c\|^2$ kontrolliert.

Insgesamt ergibt sich folgende Abschätzung zur Berechnung der Konstanten.

$$\begin{aligned}
\|\hat{x}_h\|_{\hat{X}}^2 &\leq \|\sigma_0\|_{L^2(\Omega)}^2 + (1 + C_F^2) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + (2 + 2C_F^2) \|w_c\|^2 + \|\mathbf{q}_{RT}\|_{H(\text{div}, \Omega)}^2 \\
&\leq (1 + C_F^2) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + (2 + 2C_F^2) \|w_c\|^2 + 2\|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 \\
&\quad + (1 + h_{\max}^2/(n+1)^2) \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + 3\|\sigma_0\|_{L^2(\Omega)}^2 \\
&\leq (1 + C_F^2) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + (2 + 2C_F^2) \|w_c\|^2 + 2\|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 \\
&\quad + (1 + h_{\max}^2/(n+1)^2) \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + 3 \left[(2 + C_{\text{tdd}}) \|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)} \right. \\
&\quad \left. + ((2 + C_{\text{tdd}}) h_{\max}/(n+1) + C_{\text{tdd}}) \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)} \right. \\
&\quad \left. + (1 + C_{\text{tdd}}) (\|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)} + \|w_c\|) \right]^2.
\end{aligned}$$

Der letzte quadrierte Term kann als Skalarprodukt zweier Vektoren im \mathbb{R}^4 betrachtet werden. Das Anwenden der Cauchy-Schwarz Ungleichung ergibt dann folgende verkürzte Darstellung

$$\begin{aligned}
&3 \left[(2 + C_{\text{tdd}}) \|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)} + ((2 + C_{\text{tdd}}) h_{\max}/(n+1) + C_{\text{tdd}}) \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)} \right. \\
&\quad \left. + (1 + C_{\text{tdd}}) (\|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)} + \|w_c\|) \right]^2 \\
&\leq 3 \left((2 + C_{\text{tdd}})^2 + ((2 + C_{\text{tdd}}) h_{\max}/(n+1) + C_{\text{tdd}})^2 + 2(1 + C_{\text{tdd}})^2 \right) \\
&\quad \left(\|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)}^2 + \|w_c\|^2 \right) \\
&= \left(\frac{3(C_{\text{tdd}} + 2)^2}{(n+1)^2} h_{\max}^2 + \frac{6(C_{\text{tdd}}^2 + 2C_{\text{tdd}})}{n+1} h_{\max} + 12(C_{\text{tdd}} + 1)^2 + 6 \right) \\
&\quad \left(\|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)}^2 + \|w_c\|^2 \right) \\
&=: C_3 \left(\|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)}^2 + \|w_c\|^2 \right).
\end{aligned}$$

Eingesetzt in die ursprüngliche Abschätzung ergibt dies

$$\begin{aligned}
\|\hat{x}_h\|_{\hat{X}}^2 &\leq (1 + C_F^2) \|u_0 - \Pi_0 w_c\|_{L^2(\Omega)}^2 + (2 + 2C_F^2 + C_3) \|w_c\|^2 + C_3 \|\text{dev } \sigma_0 - D w_c\|_{L^2(\Omega)}^2 \\
&\quad + (2 + C_3) \|\sigma_0 - \Pi_0 \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + (1 + h_{\max}^2/(n+1)^2 + C_3) \|\text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 \\
&\leq \max \left\{ 1 + C_F^2, C_3 + \max \left\{ 2, 1 + h_{\max}^2/(n+1)^2 \right\} \right\} \hat{b}(\hat{x}_h, \tilde{y}_h) \\
&\quad + (C_3 + 2C_F^2 + 2) \|w_c\|^2.
\end{aligned}$$

Mit den Konstanten

$$\begin{aligned}
C_3 &:= h_{\max}^2 3(C_{\text{tdd}} + 2)^2/(n+1)^2 + h_{\max} 6(C_{\text{tdd}}^2 + 2C_{\text{tdd}})/(n+1) + 12(C_{\text{tdd}} + 1)^2 + 6, \\
C_{x_h} &:= \max \left\{ 1 + C_F^2, C_3 + \max \left\{ 2, 1 + h_{\max}^2/(n+1)^2 \right\} \right\} + (C_3 + 2C_F^2 + 2) C_{w_c}
\end{aligned}$$

gilt also

$$\|\hat{x}_h\|_{\hat{X}}^2 \leq C_{x_h} \hat{b}(\hat{x}_h, \tilde{y}_h). \quad \square$$

Schritt 6. Bestimmung der Inf-Sup-Konstante. Für jedes $\hat{x}_h \in \hat{X}_h \setminus \{0\}$ erfüllt das $\tilde{y}_h \in Y_h$ aus Schritt 1.

$$0 < \frac{1}{\sqrt{C_{x_h} (1 + h_{\max}^2 n^2 / (n+1)^2)}} \leq \frac{\hat{b}(\hat{x}_h, \tilde{y}_h)}{\|\hat{x}_h\|_{\hat{X}} \|\tilde{y}_h\|_Y} \leq \sup_{y_h \in Y_h} \frac{\hat{b}(\hat{x}_h, y_h)}{\|\hat{x}_h\|_{\hat{X}} \|y_h\|_Y},$$

daher gilt $\hat{\beta}_h \geq \left(C_{x_h} + (1 + h_{\max}^2 n^2 / (n+1)^2) \right)^{-1/2}$. In (3.15) und (3.17) wird klar, dass das gesuchte $\beta_h \geq \hat{\beta}_h$, denn

$$\begin{aligned} b(x_h, y_h) &= \hat{b}(\hat{x}_h, y_h) \geq \hat{\beta}_h \inf_{\hat{x}_h \in \hat{X} \setminus \{0\}} \sup_{y_h \in Y_h} \|\hat{x}_h\|_{\hat{X}_h} \|y_h\|_Y \\ &\geq \beta_h \inf_{x_h \in X_h \setminus \{0\}} \sup_{y_h \in Y_h} \|x_h\|_X \|y_h\|_Y. \end{aligned} \quad \square$$

3.4 Fortin-Interpolator

In [GQ14],[DG11a] und [CDG14] wird der Fortin-Operator explizit konstruiert, der Ansatz stattdessen die diskrete Inf-Sup Bedingung zu beweisen wurde erst in [CGHW14] und [CH15] eingeführt. Die a-posteriori-Abschätzung aus (3.1) für alle $\xi_h \in X_h$

$$\begin{aligned} \beta \|x - \xi_h\|_X &\leq \|b\| / \beta_h \|F - b(\xi_h, \cdot)\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \\ &\leq 2\|b\|^2 / \beta_h \|x - \xi_h\|_X. \end{aligned}$$

besteht aus einem Datenapproximationsfehler $\|F \circ (1 - \Pi)\|_{Y^*}$ und dem residualen Fehlerschätzer $\|F - b(\xi_h, \cdot)\|_{Y_h^*}$. In diesem Abschnitt sollen für den zweidimensionalen, in dieser Arbeit implementierten Fall der Datenapproximationsfehler und seine Konvergenzordnung genauer untersucht werden. Es ist für $n = 2$ möglich auch für die vorgestellte Diskretisierung niedriger Ordnung eine beschränkte, lineare Abbildung $\Pi: Y \rightarrow Y_h$ mit $b(x_h, (1 - \Pi)y) = 0$ für alle $x_h \in X_h$ und alle $y \in Y$ direkt und einfach anzugeben. Diese Abbildung wird Fortin-Interpolator genannt.

Dafür werden zwei diskrete Helmholtz Zerlegungen verwendet, in denen die folgenden Operatoren und Räume vorkommen. Es werden zunächst der Curl-Operator, der gedrehte Gradienten, und der curl-Operator, dessen Spur, definiert. Für eine Funktion $\beta \in C^1(\Omega; \mathbb{R}^2)$, seien $\text{Curl } \beta := \begin{pmatrix} -\partial\beta_1/\partial x_2 & \partial\beta_1/\partial x_1 \\ -\partial\beta_2/\partial x_2 & \partial\beta_2/\partial x_1 \end{pmatrix} \in \mathbb{R}^{2 \times 2}$ und $\text{curl } \beta := \text{tr}(\text{Curl } \beta) = \partial\beta_2/\partial x_1 - \partial\beta_1/\partial x_2 \in \mathbb{R}$. Mit Curl_{NC} wird die elementweise Anwendung des Curl-Operators bezeichnet.

Damit wird der folgende Teilraum von $S^1(\mathcal{T}; \mathbb{R}^2)$ definiert

$$X_{\text{curl}} := \left\{ v_C \in S^1(\mathcal{T}; \mathbb{R}^2) : \int_{\Omega} v_C \, dx = 0 \text{ und } \int_{\Omega} \text{curl } v_C \, dx = 0 \right\}. \quad (3.26)$$

Des Weiteren werden die Crouzeix-Raviart Funktionen mit und ohne homogene Dirichlet Randbedingungen in den Mittelpunkten der Außenkanten

$$\begin{aligned} \text{CR}^1(\mathcal{T}; \mathbb{R}^2) &:= \{v \in P_1(\mathcal{T}, \mathbb{R}^2) : v \text{ ist stetig in } \text{mid}(E) \text{ für alle } E \in \mathcal{E}(\Omega)\}, \\ \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2) &:= \{v \in \text{CR}^1(\mathcal{T}; \mathbb{R}^2) : v(\text{mid}(E)) = 0 \text{ für alle } E \in \mathcal{E}(\partial\Omega)\} \end{aligned}$$

und die diskret divergenzfreien Crouzeix-Raviart Funktionen

$$Z_{\text{CR}} := \{v \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2) : \text{div}_{\text{NC}} v = 0\} \quad (3.27)$$

benötigt. Die Crouzeix-Raviart Funktionen $w_{\text{CR}} \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ haben insbesondere die Eigenschaft, dass $\int_E [w_{\text{CR}}]_E \, ds = 1/|E| \int_E [w_{\text{CR}}]_E \, ds = 0$ für alle $E \in \mathcal{E}$ gilt. Um eine ähnliche Eigenschaft zu garantieren wird folgender Interpolant definiert. Für $v \in H^1(\mathcal{T}; \mathbb{R}^2)$ sei der lokale nicht konforme Interpolant $I_{\text{NC}}^{\text{pw}} v \in P_1(\mathcal{T}; \mathbb{R}^2)$ für alle $T \in \mathcal{T}$ durch $(I_{\text{NC}}^{\text{pw}} v)|_T(\text{mid}(E)) = \int_E v|_T \, dx$ definiert. Es ist schnell nachzurechnen, dass für alle $T \in \mathcal{T}$ und alle $E \in \mathcal{E}(T)$ gilt $\int_E I_{\text{NC}}^{\text{pw}} v|_T = \int_E v|_T \, ds$ und damit auch $\Pi_0 \text{D}_{\text{NC}} v = \text{D}_{\text{NC}} I_{\text{NC}}^{\text{pw}} v$.

Damit können nun die benötigten zweidimensionalen diskreten Helmholtz Zerlegungen zitiert werden. Mit den Räumen (3.26) und (3.27) können stückweise konstante deviatorische Matrizen charakterisiert werden, wie in [CPR13] gezeigt.

Lemma 3.13 (Diskrete Helmholtz Zerlegung für deviatorische Matrizen). *Die Zerlegung*

$$P_0(\mathcal{T}; \mathbb{R}_{\text{dev}}^{2 \times 2}) = \text{D}_{\text{NC}} Z_{\text{CR}} \oplus \text{dev Curl } X_{\text{curl}}$$

ist orthogonal in $L^2(\Omega; \mathbb{R}_{\text{dev}}^{2 \times 2})$, wenn Ω einfach zusammenhängend ist.

Des Weiteren gilt nach [AF89] bzw. [Car09b, Theorem 3.32] die folgende Zerlegung der stückweise konstanten Funktionen.

Lemma 3.14 (Diskrete Helmholtz Zerlegung). *Ist Ω einfach zusammenhängend, so ist die Zerlegung*

$$P_0(\mathcal{T}; \mathbb{R}^{2 \times 2}) = \text{D } S_0^1(\mathcal{T}; \mathbb{R}^2) \oplus \text{Curl}_{\text{NC}}(\text{CR}^1(\mathcal{T})/\mathbb{R})^2$$

orthogonal in $L^2(\Omega; \mathbb{R}^{2 \times 2})$.

Dabei bedeutet $\text{CR}^1(\mathcal{T})/\mathbb{R}$ lediglich bis auf Konstante und garantiert die Eindeutigkeit der Zerlegung. Damit diese Helmholtz Zerlegungen verwendet werden können, wird Ω

als einfach zusammenhängend angenommen (im Zweifel sind alle Betrachtungen auf den einzelnen Zusammenhangskomponenten durchzuführen).

Mit dieser Vorarbeit ist es möglich den Fortin-Operator Π und den diskreten Kern N_h anzugeben.

Lemma 3.15. *Zu $y = (\boldsymbol{\tau}, v) \in Y$ erfülle $\alpha_{CR} \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ für alle $w_{CR} \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ die folgende Gleichung*

$$\int_{\Omega} \text{D}_{NC} \alpha_{CR} : \text{D}_{NC} w_{CR} \, dx = \int_{\Omega} w_{CR} \cdot (1 - \Pi_0) \text{div}_{NC} \boldsymbol{\tau} \, dx + \int_{\Omega} \Pi_0 \boldsymbol{\tau} : \text{D}_{NC} w_{CR} \, dx. \quad (3.28)$$

Außerdem seien $z_{CR} \in Z_{CR}$ und $\beta_c \in X_{\text{curl}}$ die Funktionen aus der nach Lemma 3.13 existierenden Helmholtz Zerlegung

$$\text{dev}(\Pi_0 \boldsymbol{\tau} - \text{D}_{NC} \alpha_{CR}) = \text{D}_{NC} z_{CR} + \text{dev} \text{Curl} \beta_c. \quad (3.29)$$

Dann besitzt die folgende Abbildung

$$\Pi: H(\text{div}, \mathcal{T}; \mathbb{R}^{2 \times 2}) / \mathbb{R} \times H^1(\mathcal{T}; \mathbb{R}^2) \rightarrow RT_0^{pw}(\mathcal{T}) / \mathbb{R} \times P_1(\mathcal{T})$$

mit

$$\begin{aligned} \Pi(\boldsymbol{\tau}) &:= \boldsymbol{\tau}_{RT} := \text{D}_{NC} \alpha_{CR} + \text{Curl} \beta_c + (\Pi_0 \text{div}_{NC} \boldsymbol{\tau} / 2) \otimes (\cdot - \text{mid}(\mathcal{T})), \\ \Pi(v) &:= v_1 := I_{NC}^{pw} v + z_{CR}, \end{aligned}$$

die gewünschten Eigenschaften, ist also beschränkt und erfüllt $b(x_h, (1 - \Pi)y) = 0$ für alle $x_h \in X_h$ und alle $y \in Y$. Außerdem gilt für den diskreten Kern von b

$$\begin{aligned} N_h &:= \{y_h \in Y_h : b(x_h, y_h) = 0 \text{ für alle } x_h \in X_h\} \\ &= \left\{ (\text{Curl}_{NC} \beta_{CR}, v_{CR}) \in \text{Curl}_{NC} \left(\text{CR}_0^1(\mathcal{T}) / \mathbb{R} \right)^2 \times Z_{CR} : \right. \\ &\quad \left. - \text{D}_{NC} v_{CR} = \text{dev} \text{Curl}_{NC} \beta_{CR} \right\}. \end{aligned}$$

Beweis. Schritt 1. Zunächst wird Π untersucht. Offensichtlich gilt $v_1 \in P_1(\mathcal{T}; \mathbb{R}^2)$ und $\boldsymbol{\tau}_{RT} \in RT_0^{pw}(\mathcal{T}; \mathbb{R}^{2 \times 2})$. Nach dem Gaußschem Integralsatz, den Eigenschaften der Sprünge von $\text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ und den Nebenbedingungen in X_{curl} ist

$$\begin{aligned} \int_{\Omega} \text{tr} \boldsymbol{\tau}_{RT} \, dx &= \int_{\Omega} \text{div}_{NC} \alpha_{CR} \, dx + \int_{\Omega} \text{tr} \left((\Pi_0 \text{div}_{NC} \boldsymbol{\tau} / 2) \otimes (\cdot - \text{mid}(\mathcal{T})) \right) \, dx \\ &\quad + \int_{\Omega} \text{curl} \beta_c \, dx \\ &= \sum_{E \in \mathcal{E}} \int_E [\alpha_{CR}]_E \nu_E \, ds + 0 + 0 = 0 \end{aligned}$$

erfüllt. Also ist $\boldsymbol{\tau}_{RT} \in RT_0^{pw}(\mathcal{T}; \mathbb{R}^{2 \times 2}) / \mathbb{R}$.

Schritt 2. Offenbar ist die Forderung $b(x_h, (1 - \Pi)y) = 0$ für alle $x_h \in X_h$ und alle $y \in Y$ äquivalent dazu, dass $y = (\boldsymbol{\tau}, v)$ und $\Pi y = (\boldsymbol{\tau}_{\text{RT}}, v_1)$ für alle $x_h = (\boldsymbol{\sigma}_0, u_0, s_1, t_0) \in X_h$ folgendes Gleichungssystem erfüllen

$$0 = b((\boldsymbol{\sigma}_0, 0, 0, 0), y - y_h) = \int_{\Omega} \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}}(v - v_1) \, dx + \int_{\Omega} \text{dev } \boldsymbol{\sigma}_0 : (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx, \quad (3.30)$$

$$0 = b((0, u_0, 0, 0), y - y_h) = \int_{\Omega} u_0 \cdot \text{div}_{\text{NC}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx, \quad (3.31)$$

$$0 = b((0, 0, s_1, 0), y - y_h) = - \left\langle \gamma_{\nu}^{\mathcal{T}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}), s_1 \right\rangle_{\partial \mathcal{T}}, \quad (3.32)$$

$$0 = b((0, 0, 0, t_0), y - y_h) = - \left\langle t_0, \gamma_0^{\mathcal{T}}(v - v_1) \right\rangle_{\partial \mathcal{T}}. \quad (3.33)$$

Da $\mathbf{D}_{\text{NC}} \alpha_{\text{CR}}, \text{Curl } \beta_c \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})$ gilt, ist $\text{div}_{\text{NC}} \boldsymbol{\tau}_{\text{RT}} = \Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau}$. Daraus ergibt sich (3.31), denn nach (2.10) gilt für alle $u_0 \in P_0(\mathcal{T})$

$$\int_{\Omega} u_0 \cdot \text{div}_{\text{NC}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx = \int_{\Omega} u_0 \cdot \Pi_0 \text{div}_{\text{NC}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx = 0.$$

Die Gleichung (3.33) lässt sich aus den Eigenschaften von $I_{\text{NC}}^{\text{pw}}$ und den Crouzeix-Raviart Funktionen folgern. Es gilt für alle $t_0 \in P_0(\mathcal{E}; \mathbb{R}^2)$

$$\begin{aligned} - \left\langle t_0, \gamma_0^{\mathcal{T}}(v - v_1) \right\rangle_{\partial \mathcal{T}} &= \sum_{E \in \mathcal{E}} \int_E t_0 [v - v_1]_E \, ds \\ &= \sum_{E \in \mathcal{E}} t_0|_E \left(\int_E [v - I_{\text{NC}}^{\text{pw}} v]_E \, ds - \int_E [z_{\text{CR}}]_E \, ds \right) = 0. \end{aligned}$$

Die Wahl von α_{CR} wie in (3.28) garantiert (3.32). Betrachte zu jedem $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^2)$ die Fortsetzung $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$ mit $\gamma_0^{\mathcal{T}} w_c = s_1$, dann gilt nach Satz 2.8

$$- \left\langle \gamma_{\nu}^{\mathcal{T}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}), s_1 \right\rangle_{\partial \mathcal{T}} = \int_{\Omega} \mathbf{D} w_c : \Pi_0(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx + \int_{\Omega} w_c \cdot \text{div}_{\text{NC}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx.$$

Die Orthogonalität in Lemma 3.14, Gleichung (3.28) und die Tatsache $\text{div}_{\text{NC}} \boldsymbol{\tau}_{\text{RT}} = \Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau}$ zeigen, dass für alle $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2) \subseteq \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ gilt

$$\begin{aligned} &\int_{\Omega} \mathbf{D} w_c : \Pi_0(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx + \int_{\Omega} w_c \cdot \text{div}_{\text{NC}}(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx \\ &= \int_{\Omega} \mathbf{D} w_c : \Pi_0 \boldsymbol{\tau} \, dx - \int_{\Omega} \mathbf{D} w_c : \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} \, dx - \int_{\Omega} \mathbf{D} w_c : \text{Curl } \beta_c \, dx \\ &\quad + \int_{\Omega} w_c \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx \\ &= \int_{\Omega} \mathbf{D} w_c : \Pi_0 \boldsymbol{\tau} \, dx - \int_{\Omega} w_c \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx - \int_{\Omega} \Pi_0 \boldsymbol{\tau} : \mathbf{D} w_c \, dx \\ &\quad + \int_{\Omega} w_c \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx \\ &= 0. \end{aligned}$$

Die linke Seite der Gleichung (3.30) lässt sich mit Lemma 2.5 und (2.10) wie folgt umschreiben

$$\begin{aligned} & \int_{\Omega} \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}}(v - v_1) \, dx + \int_{\Omega} \text{dev } \boldsymbol{\sigma}_0 : (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx \\ &= \int_{\Omega} \boldsymbol{\sigma}_0 : \Pi_0 \mathbf{D}_{\text{NC}}(v - v_1) \, dx + \int_{\Omega} \boldsymbol{\sigma}_0 : \text{dev } \Pi_0(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx. \end{aligned}$$

Die Gleichheit $\Pi_0 \mathbf{D}_{\text{NC}} v = \mathbf{D}_{\text{NC}} I_{\text{NC}}^{\text{pw}} v$ ergibt zusammen mit (3.29)

$$\begin{aligned} & \int_{\Omega} \boldsymbol{\sigma}_0 : \Pi_0 \mathbf{D}_{\text{NC}}(v - v_1) \, dx + \int_{\Omega} \boldsymbol{\sigma}_0 : \text{dev } \Pi_0(\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}) \, dx \\ &= - \int_{\Omega} \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}} z_{\text{CR}} \, dx + \int_{\Omega} \boldsymbol{\sigma}_0 : \text{dev}(\Pi_0 \boldsymbol{\tau} - \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} - \text{Curl } \beta_c) \, dx \\ &= - \int_{\Omega} \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}} z_{\text{CR}} \, dx + \int_{\Omega} \boldsymbol{\sigma}_0 : \mathbf{D}_{\text{NC}} z_{\text{CR}} \, dx = 0. \end{aligned}$$

Damit gilt $b(x_h, (1 - \Pi)y) = 0$ für alle $x_h \in X_h$ und alle $y \in Y$. Es lässt sich direkt verifizieren, dass Π idempotent ist, also $\Pi(\Pi(y)) = \Pi(y)$ gilt.

Schritt 3. Die Beschränktheit der Abbildung lässt sich ebenfalls nachrechnen. Sei $y = (\boldsymbol{\tau}, v) \in Y$ beliebig, dann gilt für $\Pi(y) \in Y_h$ wie in Lemma 3.15

$$\begin{aligned} \|\Pi(y)\|_Y^2 &= \|\Pi(\boldsymbol{\tau})\|_{H(\text{div}, \Omega)}^2 + \|\Pi(v)\|_{H^1(\Omega)}^2 \\ &= \|\mathbf{D}_{\text{NC}} \alpha_{\text{CR}}\|_{L^2(\Omega)}^2 + \|\text{Curl } \beta_c\|_{L^2(\Omega)}^2 + \|\Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau} / 2 \otimes (\cdot - \text{mid}(\mathcal{T}))\|_{L^2(\Omega)}^2 \\ &\quad + \|\Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|I_{\text{NC}}^{\text{pw}} v + z_{\text{CR}}\|_{L^2(\Omega)}^2 + \|\mathbf{D}_{\text{NC}}(I_{\text{NC}}^{\text{pw}} v + z_{\text{CR}})\|_{L^2(\Omega)}^2. \end{aligned}$$

Gleichung (3.28), die Cauchy-Schwarz und die diskrete Friedrichs Ungleichung mit Konstante C_{dF} aus [BS08, 10.6.14] ermöglichen die Abschätzung des ersten Terms

$$\begin{aligned} \int_{\Omega} \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} : \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} \, dx &= \int_{\Omega} \alpha_{\text{CR}} \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx + \int_{\Omega} \Pi_0 \boldsymbol{\tau} : \mathbf{D}_{\text{NC}} \alpha_{\text{CR}} \, dx \\ &\leq \|\alpha_{\text{CR}}\|_{L^2(\Omega)} \|(1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)} + \|\Pi_0 \boldsymbol{\tau}\|_{L^2(\Omega)} \|\alpha_{\text{CR}}\|_{\text{NC}} \\ &\leq \|\alpha_{\text{CR}}\|_{\text{NC}} \left(C_{\text{dF}} \|\text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)} + \|\boldsymbol{\tau}\|_{L^2(\Omega)} \right). \end{aligned}$$

Nach Cauchy-Schwarz und der Normabschätzung in (2.2) gilt

$$\begin{aligned} \|\Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau} / 2 \otimes (\cdot - \text{mid}(\mathcal{T}))\|_{L^2(\Omega)} &\leq \|\Pi_0 \text{div}_{\text{NC}} \boldsymbol{\tau} / 2\|_{L^2(\Omega)} \|\cdot - \text{mid}(\mathcal{T})\|_{L^\infty(\Omega)} \\ &\leq h_{\text{max}} / 3 \|\text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)}. \end{aligned}$$

Außerdem gilt wegen der bereits erwähnten diskreten Friedrichs Ungleichung und $\Pi_0 \mathbf{D}_{\text{NC}} v = \mathbf{D}_{\text{NC}} I_{\text{NC}}^{\text{pw}} v$, dass

$$\begin{aligned} \|I_{\text{NC}}^{\text{pw}} v\|_{L^2(\Omega)}^2 + \|\mathbf{D}_{\text{NC}} I_{\text{NC}}^{\text{pw}} v\|_{L^2(\Omega)}^2 &\leq (1 + C_{\text{dF}}^2) \|\mathbf{D}_{\text{NC}} I_{\text{NC}}^{\text{pw}} v\|_{L^2(\Omega)}^2 = (1 + C_{\text{dF}}^2) \|\Pi_0 \mathbf{D}_{\text{NC}} v\|_{L^2(\Omega)}^2 \\ &\leq (1 + C_{\text{dF}}^2) \|\mathbf{D}_{\text{NC}} v\|_{L^2(\Omega)}^2. \end{aligned}$$

Auf Grund der diskreten Friedrichs Ungleichung genügt es zuletzt die folgende Summe zu betrachten,

$$\|\mathbf{D}_{\text{NC}} z_{\text{CR}}\|_{L^2(\Omega)}^2 + \|\text{Curl } \beta_c\|_{L^2(\Omega)}^2$$

Für $\beta_c \in X_{\text{curl}}$ gilt $\text{div Curl } \beta_c \equiv 0$ und $\int_{\Omega} \text{tr Curl } \beta_c \, dx = 0$, damit liegt $\text{Curl } \beta_c \in H(\text{div}, \Omega; \mathbb{R}^{2 \times 2})/\mathbb{R}$ und das Tr-Div-Dev Lemma, Lemma 2.18, ist anwendbar. Es gilt zusammen mit Lemma 2.5

$$\begin{aligned} \|\text{Curl } \beta_c\|_{L^2(\Omega)}^2 &= \|\text{dev Curl } \beta_c\|_{L^2(\Omega)}^2 + 1/2 \|\text{tr Curl } \beta_c\|_{L^2(\Omega)}^2 \\ &\leq (1 + C_{\text{td}}^2) \|\text{dev Curl } \beta_c\|_{L^2(\Omega)}^2 + C_{\text{td}}^2 \|\text{div Curl } \beta_c\|_{L^2(\Omega)}^2 \\ &= (1 + C_{\text{td}}^2) \|\text{dev Curl } \beta_c\|_{L^2(\Omega)}^2. \end{aligned}$$

Mit der Abkürzung $C := (1 + C_{\text{td}}^2)$, mit (3.29) und der Orthogonalität in eben dieser Helmholtz Zerlegung ergibt folgende Abschätzung

$$\begin{aligned} \|\mathbf{D}_{\text{NC}} z_{\text{CR}}\|_{L^2(\Omega)}^2 + \|\text{Curl } \beta_c\|_{L^2(\Omega)}^2 &\leq \|\mathbf{D}_{\text{NC}} z_{\text{CR}}\|_{L^2(\Omega)}^2 + C \|\text{dev Curl } \beta_c\|_{L^2(\Omega)}^2 \\ &\leq \max\{1, C\} \|\mathbf{D}_{\text{NC}} z_{\text{CR}} + \text{dev Curl } \beta_c\|_{L^2(\Omega)}^2 \\ &\leq C \|\text{dev}(\Pi_0 \boldsymbol{\tau} - \mathbf{D}_{\text{NC}} \alpha_{\text{CR}})\|_{L^2(\Omega)}^2 \\ &\leq 2C \left(\|\Pi_0 \boldsymbol{\tau}\|_{L^2(\Omega)}^2 + \|\mathbf{D}_{\text{NC}} \alpha_{\text{CR}}\|_{L^2(\Omega)}^2 \right) \\ &\leq C \left(6 \|\boldsymbol{\tau}\|_{L^2(\Omega)}^2 + 4C_{\text{dF}}^2 \|\text{div}_{\text{NC}} \boldsymbol{\tau}\|_{L^2(\Omega)}^2 \right). \end{aligned}$$

Damit sind alle Terme durch Teile von $\|y\|_Y$ kontrolliert und die angegebene Abbildung ist beschränkt.

Schritt 4. Es bleibt N_h zu bestimmen. Für alle $\tilde{y}_h = (\tilde{\boldsymbol{\tau}}_{\text{RT}}, \tilde{v}_1) \in N_h$, ergibt sich aus

$$0 = b((0, 0, 0, t_0), (\tilde{\boldsymbol{\tau}}_{\text{RT}}, \tilde{v}_1)) = - \left\langle t_0, \gamma_0^T \tilde{v}_1 \right\rangle_{\partial \mathcal{T}} = - \sum_{E \in \mathcal{E}} t_0|_E \int_E [\tilde{v}_1]_E \, ds$$

für alle $t_0 \in P_0(\mathcal{E}; \mathbb{R}^2)$, die Bedingung $v_1 \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$. Aus

$$0 = b((0, u_0, 0, 0), (\tilde{\boldsymbol{\tau}}_{\text{RT}}, \tilde{v}_1)) = \int_{\Omega} u_0 \cdot \text{div}_{\text{NC}} \tilde{\boldsymbol{\tau}}_{\text{RT}} \, dx$$

für alle $u_0 \in P_0(\mathcal{T}; \mathbb{R}^2)$ folgt $\text{div}_{\text{NC}} \tilde{\boldsymbol{\tau}}_{\text{RT}} = 0$, also $\Pi_0 \tilde{\boldsymbol{\tau}}_{\text{RT}} = \tilde{\boldsymbol{\tau}}_{\text{RT}}$. Des Weiteren kann jedes $s_1 \in S_0^1(\mathcal{E}; \mathbb{R}^2)$ wieder durch ein $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$ mit $\gamma_0^T w_c = s_1$ fortgesetzt werden und es ergibt sich nach (2.8)

$$0 = b((0, 0, s_1, 0), (\tilde{\boldsymbol{\tau}}_{\text{RT}}, \tilde{v}_1)) = - \left\langle \gamma_{\nu}^T (\boldsymbol{\tau} - \boldsymbol{\tau}_{\text{RT}}), s_1 \right\rangle_{\partial \mathcal{T}} = \int_{\Omega} \mathbf{D} w_c : \tilde{\boldsymbol{\tau}}_{\text{RT}} \, dx$$

für alle $w_c \in S_0^1(\mathcal{T}; \mathbb{R}^2)$. Nach der Helmholtz Zerlegung Lemma 3.14 gilt $\tilde{\boldsymbol{\tau}}_{\text{RT}} = \text{Curl}_{\text{NC}} \beta_{\text{CR}}$

für ein $\beta_{\text{CR}} \in \left(\text{CR}^1(\mathcal{T})/\mathbb{R}\right)^2$. Nun folgt aus

$$0 = b((\sigma_0, 0, 0, 0), (\tilde{\tau}_{\text{RT}}, \tilde{v}_1)) = \int_{\Omega} \sigma_0 : D_{\text{NC}} \tilde{v}_1 \, dx + \int_{\Omega} \text{dev } \sigma_0 : \tilde{\tau}_{\text{RT}} \, dx$$

für alle $\sigma_0 \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ noch $\text{dev } \tilde{\tau}_{\text{RT}} = -D_{\text{NC}} \tilde{v}_1$, da $\tilde{v}_1 \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$ divergenzfrei ist. Insgesamt ergibt sich

$$N_h = \left\{ (\text{Curl}_{\text{NC}} \beta_{\text{CR}}, v_{\text{CR}}) \in \text{Curl}_{\text{NC}} \left(\text{CR}^1(\mathcal{T})/\mathbb{R}\right)^2 \times Z_{\text{CR}} : \right. \\ \left. -D_{\text{NC}} v_{\text{CR}} = \text{dev } \text{Curl}_{\text{NC}} \beta_{\text{CR}} \right\}.$$

Die Dimension dieses Raumes lässt sich bestimmen, in dem eine Zerlegung $Y_h = N_h \oplus M_h$ betrachtet wird. Dabei sei wie in [Wer11, V.3.4]

$$M_h = N_h^\perp = \left\{ m_h \in Y_h : \langle m_h, n_h \rangle_Y = 0 \text{ für alle } n_h \in N_h \right\}$$

als das orthogonale Komplement zu N_h im Hilbertraum Y gewählt. Sei wieder T_h der trial-to-test Operator aus Abschnitt 3.1, dann gilt $T_h(X_h) = M_h$. Denn für alle $x_h \in X_h$ und $n_h \in N_h$ gilt $\langle T_h(x_h), n_h \rangle_Y = b(x_h, n_h) = 0$, also $T_h(X_h) \subseteq M_h$. Umgekehrt definiert $m_h \in M_h$ eine Abbildung $\langle \cdot, m_h \rangle_Y \in M_h^*$ und eine Lösung $x_h \in X_h$ mit $b(x_h, \cdot) = \langle m_h, \cdot \rangle \in M_h^*$. Für alle $\hat{y}_h = \hat{m}_h + \hat{n}_h \in Y_h = N_h \oplus M_h$ gilt $\langle m_h, y_h \rangle_Y = \langle m_h, \hat{m}_h \rangle_Y = b(x_h, \hat{m}_h) = b(x_h, \hat{y}_h)$, also ist $m_h = T_h(x_h)$ und $M_h \subseteq T_h(X_h)$. Nun gilt wegen der eindeutigen Lösbarkeit

$$\begin{aligned} \dim(N_h) &= \dim(Y_h) - \dim(M_h) = \dim(Y_h) - \dim(X_h) \\ &= 12|\mathcal{T}| - 1 - (6|\mathcal{T}| - 1 + 2|\mathcal{N}(\Omega)| + 2|\mathcal{E}|). \end{aligned}$$

Nach den Euler Formeln [Car09b, Remark 1.21] gilt $2|\mathcal{T}| + 1 = |\mathcal{N}| + |\mathcal{E}(\Omega)|$. Da $|\mathcal{E}(\partial\Omega)| = |\mathcal{N}(\partial\Omega)|$ ergibt sich

$$\begin{aligned} \dim(N_h) &= 6|\mathcal{T}| - 2|\mathcal{N}(\Omega)| - 2|\mathcal{E}| = 6|\mathcal{T}| - 2(|\mathcal{N}| - |\mathcal{N}(\partial\Omega)| + |\mathcal{E}|) \\ &= 6|\mathcal{T}| - 2(|\mathcal{N}| + |\mathcal{E}(\Omega)|) = 2(|\mathcal{T}| - 1). \end{aligned} \quad \square$$

Vor der genaueren Untersuchung des Datenapproximationsfehlers, sollen noch einige Überlegungen den Kern N_h betreffend festgehalten werden. Die Untersuchung des Kerns N_h bzw. des Bildes $M_h := T_h(X_h)$ des trial-to-test Operators ist lohnend, da die Kenntnis von M_h erlaubt, die dPG-Methode schlicht als gemischte Methode zu betrachten.

Bemerkung 3.16 Es gilt, wie soeben gezeigt,

$$N_h = \left\{ (\text{Curl}_{\text{NC}} \beta_{\text{CR}}, v_{\text{CR}}) \in \text{Curl}_{\text{NC}} \left(\text{CR}^1(\mathcal{T})/\mathbb{R}\right)^2 \times Z_{\text{CR}} : \right. \\ \left. -D_{\text{NC}} v_{\text{CR}} = \text{dev } \text{Curl}_{\text{NC}} \beta_{\text{CR}} \right\}.$$

Nach der Helmholtz Zerlegung, Lemma 3.13, ist $\operatorname{dev} \operatorname{Curl}_{\text{NC}} \beta_{\text{CR}} \perp \operatorname{dev} \operatorname{Curl} X_{\text{curl}}$. Sei $b_{\text{CR}} \in (\operatorname{CR}^1(\mathcal{T})/\mathbb{R})^2$ gegeben, so müsste $\operatorname{dev} \operatorname{Curl}_{\text{NC}} b_{\text{CR}}$ auf $\operatorname{dev} \operatorname{Curl} X_{\text{curl}}$ L^2 -projiziert werden. Der Fehler dieser Projektion ergibt ein geeignetes β_{CR} . Die genauere Betrachtung von $\beta_{\text{CR}} \in (\operatorname{CR}^1(\mathcal{T})/\mathbb{R})^2$ und $b_c \in X_{\text{curl}}$ verdeutlicht einen Zusammenhang zur linearen Elastizität. Nach Definition des symmetrischen Ableitung $\varepsilon(\cdot) := \operatorname{sym} D(\cdot) = (D(\cdot) + D(\cdot)^\top)/2$ gilt

$$0 = \int_{\Omega} \operatorname{dev} \operatorname{Curl}_{\text{NC}} \beta_{\text{CR}} : \operatorname{dev} \operatorname{Curl} b_c \, dx = \int_{\Omega} \varepsilon_{\text{NC}}(\beta_{\text{CR}}) : \varepsilon(b_c) \, dx.$$

Gleichungen dieser Art werden in der Elastizität häufiger untersucht.

Nach der anderen Helmholtz Zerlegung, Lemma 3.14, ist $\operatorname{Curl}_{\text{NC}} \beta_{\text{CR}} \in \operatorname{Curl}(\operatorname{CR}^1(\mathcal{T})/\mathbb{R})^2$ bezüglich der L^2 -Norm senkrecht zu $D S_0^1(\mathcal{T}; \mathbb{R}^2)$. D.h. da $\operatorname{Curl}_{\text{NC}} \beta_{\text{CR}}$ divergenzfrei ist gilt

$$D S_0^1(\mathcal{T}; \mathbb{R}^2) \times \{0\} \subset M_h.$$

Im Folgenden soll nun der Datenapproximationsfehler untersucht werden. Handelt es sich um ein Problem mit homogenen Randdaten so lässt sich der Datenapproximationsfehler,

$$\|F \circ (1 - \Pi)\|_{Y^*} = \sup_{(v, \tau) \in S(Y)} \int_{\Omega} f \cdot (v - \Pi v) \, dx,$$

mit Hilfe des folgenden Lemmas, [CG14, Theorem 4], abschätzen.

Lemma 3.17. *Sei $j_{1,1}$ die erste Wurzel der ersten Besselfunktion, dann gilt mit der Konstante $\kappa := \sqrt{1/48 + j_{1,1}^{-2}} = 0.298234942888$ für alle $v \in H^1(\mathcal{T}; \mathbb{R}^2)$*

$$\|v - I_{\text{NC}}^{\text{pw}} v\|_{L^2(\Omega)} \leq \kappa \|h_{\mathcal{T}} (v - I_{\text{NC}}^{\text{pw}} v)\|_{\text{NC}}.$$

Nach der Definition des Fortin-Operators Π in Lemma 3.15, der Cauchy-Schwarz und der diskreten Friedrichs Ungleichung aus [BS08, 10.6.14] mit Konstante C_{dF} gilt für jedes $y = (\tau, v) \in Y$,

$$\begin{aligned} F((1 - \Pi)(y)) &= \int_{\Omega} f \cdot (v - \Pi v) \, dx \\ &= \int_{\Omega} f \cdot (v - I_{\text{NC}}^{\text{pw}} v) \, dx - \int_{\Omega} f z_{\text{CR}} \, dx \\ &\leq \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \|h_{\mathcal{T}}^{-1} (v - I_{\text{NC}}^{\text{pw}} v)\|_{L^2(\Omega)} + C_{\text{dF}} \|f\|_{L^2(\Omega)} \|z_{\text{CR}}\|_{\text{NC}}. \end{aligned}$$

Der erste Term kann mit Lemma 3.17 abgeschätzt werden

$$\|h_{\mathcal{T}}^{-1} (v - I_{\text{NC}}^{\text{pw}} v)\|_{L^2(\Omega)} \leq \kappa \|v - I_{\text{NC}}^{\text{pw}} v\|_{\text{NC}} = \kappa \|(1 - \Pi_0) D_{\text{NC}} v\|_{L^2(\Omega)} \leq \kappa \|v\|_{\text{NC}}.$$

Zur Betrachtung des zweiten Summanden sei an die Definition von z_{CR} erinnert. Es gilt nach (3.29), der Orthogonalität in der Helmholtz Zerlegung, Lemma 3.13, und Lem-

ma 2.5

$$\begin{aligned}
 \|z_{\text{CR}}\|_{\text{NC}}^2 &= \int_{\Omega} \mathbf{D}_{\text{NC}} z_{\text{CR}} : \mathbf{D}_{\text{NC}} z_{\text{CR}} \, dx \\
 &= \int_{\Omega} \mathbf{D}_{\text{NC}} z_{\text{CR}} : (\text{dev}(\Pi_0 \boldsymbol{\tau} - \mathbf{D}_{\text{NC}} \alpha_{\text{CR}}) + \text{dev} \text{Curl} \beta_c) \, dx \\
 &= \int_{\Omega} \mathbf{D}_{\text{NC}} z_{\text{CR}} : (\Pi_0 \boldsymbol{\tau} - \mathbf{D}_{\text{NC}} \alpha_{\text{CR}}) \, dx.
 \end{aligned}$$

Da insbesondere $z_{\text{CR}} \in \text{CR}_0^1(\mathcal{T}; \mathbb{R}^2)$, ergibt die Gleichung (3.28)

$$\begin{aligned}
 \|z_{\text{CR}}\|_{\text{NC}}^2 &= \int_{\Omega} \mathbf{D}_{\text{NC}} z_{\text{CR}} : \Pi_0 \boldsymbol{\tau} \, dx - \int_{\Omega} z_{\text{CR}} \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx - \int_{\Omega} \Pi_0 \boldsymbol{\tau} : \mathbf{D}_{\text{NC}} z_{\text{CR}} \, dx \\
 &= - \int_{\Omega} z_{\text{CR}} \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx.
 \end{aligned}$$

Dies kann nun mit der Cauchy-Schwarz-Ungleichung und der Poincaré Ungleichung, Lemma 2.17, abgeschätzt werden

$$\begin{aligned}
 \|z_{\text{CR}}\|_{\text{NC}}^2 &= - \int_{\Omega} (1 - \Pi_0) z_{\text{CR}} \cdot (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \, dx \\
 &\leq \| (1 - \Pi_0) z_{\text{CR}} \|_{L^2(\Omega)} \| (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \|_{L^2(\Omega)} \\
 &\leq \| (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \|_{L^2(\Omega)} \sum_{T \in \mathcal{T}} \left\| z_{\text{CR}} - \oint_T z_{\text{CR}} \, dx \right\|_{L^2(T)} \\
 &\leq \| (1 - \Pi_0) \text{div}_{\text{NC}} \boldsymbol{\tau} \|_{L^2(\Omega)} \sum_{T \in \mathcal{T}} h_T / j_{1,1} \| \mathbf{D} z_{\text{CR}} \|_{L^2(T)} \\
 &\leq h_{\max} / j_{1,1} \| z_{\text{CR}} \|_{\text{NC}} \| \text{div}_{\text{NC}} \boldsymbol{\tau} \|_{L^2(\Omega)},
 \end{aligned}$$

wobei $j_{1,1}$ wie in Lemma 3.17 die erste Wurzel der ersten Bessel-Funktion ist. Insgesamt ergibt sich somit

$$\begin{aligned}
 \|F \circ (1 - \Pi)\|_{Y^*} &= \sup_{(v, \boldsymbol{\tau}) \in S(Y)} \int_{\Omega} f \cdot (v - \Pi v) \, dx \\
 &\leq \sup_{(v, \boldsymbol{\tau}) \in S(Y)} \left(h_{\max} C_{\text{dF}} / j_{1,1} \|f\|_{L^2(\Omega)} \| \text{div}_{\text{NC}} \boldsymbol{\tau} \|_{L^2(\Omega)} \right. \\
 &\quad \left. + \kappa \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \|v\|_{\text{NC}} \right) \\
 &\leq h_{\max} C_{\text{dF}} / j_{1,1} \|f\|_{L^2(\Omega)} + \kappa \|h_{\mathcal{T}} f\|_{L^2(\Omega)}.
 \end{aligned}$$

Dies ist kein Term höherer Ordnung. Für das Poisson-Modell Problem und für die Navier-Lamé Gleichungen in [CH15, 6.2] wurde gezeigt, dass eine leichte Vergrößerung des Testraums, welche keinen Einfluss auf den Beweis der diskreten und kontinuierlichen Inf-Sup Bedingung hat, die höhere Ordnung liefert. Dabei werden beispielsweise bubble-Funktionen bzw. RT_1 -Funktionen verwendet. Ein ähnliches Vorgehen würde sich auch für die Stokes Gleichungen anbieten.

Zusammenfassend ist festzustellen, dass der in dieser Arbeit zur Verringerung des Re-

chenaufwands minimal gehaltende Testraum, im Allgemeinen nicht zu einen Datenapproximationsfehler von höherer Ordnung führt. Da die a-posteriori-Abschätzung aus [CDG14] allerdings für jedes Element des Ansatzraumes gilt, ist es möglich die Lösung des Problems mit kleinerem Testraum mit dem Datenfehler einer anderen Methode mit vergrößertem Testraum zu kombinieren. Die diskrete Lösung $x_h \in X_h$ zu Problem 10 kann in (3.1) mit erweitertem Testraum $\hat{Y}_h \supset Y_h$ eingesetzt werden. In diesem Fall wäre als Fehlerschätzer $\|F - b(x_h, \cdot)\|_{\hat{Y}_h^*}$ zu verwenden und der Datenapproximationsfehler $\|F \circ (1 - \hat{\Pi})\|_{Y^*}$ mit dem Fortin-Operator $\hat{\Pi}: Y \rightarrow \hat{Y}_h$ wäre von höherer Ordnung. Es würde sich also beispielsweise anbieten $\|F - b(x_h, \cdot)\|_{\hat{Y}_h^*}$ bei der adaptiven Verfeinerung zu betrachten.

Diese Fortin-Operatoren und Vergrößerungen des Testraumes sind aktuell Forschungsgegenstand. Es wird untersucht in wie weit die numerischen Lösungen und die Fehlerschätzer davon beeinflusst werden. Interessant wäre in diesem Kontext auch die Betrachtung inhomogener Randdaten, also die Einbeziehung des Terms $\langle (1 - \Pi)\gamma_\nu^\mathcal{T} \boldsymbol{\tau}, \hat{s} \rangle_{\partial\mathcal{T}}$.

4 Implementierung

In diesem Kapitel werden die Rechnungen vorgestellt, die der Implementierung der Methode in Dimension $n = 2$ zu Grunde liegen, und wesentliche Informationen zu den als Erweiterung des AFEM-Software-Paketes [Car09a] implementierten Programmen gegeben. Der entstandene Code befindet sich vollständig auf der dieser Arbeit beiliegenden CD.

Zunächst wird in Abschnitt 4.1 der allgemeine Aufbau des zu lösenden Gleichungssystems und eines impliziten Fehlerschätzers beschrieben. Anschließend werden in Abschnitt 4.2 die verwendeten Basisfunktionen vorgestellt, mit denen in Abschnitt 4.3 die in (4.1) benötigten Matrizen explizit berechnet werden. Diese Rechnungen werden durch die entsprechenden Codezeilen ergänzt. Der Paragraph 4.4 widmet sich dem Problem des exakten Fehlers. Danach beschäftigt sich Abschnitt 4.5 mit der Gebietsskalierung, um verschiedene Gewichtungen der Fehlerterme zu vermeiden. Zuletzt werden in Abschnitt 4.6 die einzelnen implementierten Programme kurz vorgestellt.

4.1 Löser und Fehlerschätzer

Zur Lösung des Problems 10 ist zunächst eine Basis für die endlich dimensionalen Unterräume X_h und Y_h zu wählen. Für die theoretischen Untersuchungen seien $X_h = \text{span} \{\xi_1, \dots, \xi_N\}$ und $Y_h = \text{span} \{\zeta_1, \dots, \zeta_M\}$. Bezüglich dieser Basen sei die Steifigkeitsmatrix $\mathbb{B} \in \mathbb{R}^{M \times N}$ definiert durch

$$\mathbb{B}_{m,n} := b(\xi_n, \zeta_m) \quad \text{für } n = 1, \dots, N \text{ und } m = 1, \dots, M.$$

Die positiv definite und symmetrische Normmatrix $\mathbb{M} \in \mathbb{R}^{M \times M}$ sei so definiert, dass für jeden Koeffizientenvektor $y \in \mathbb{R}^M$ zu $y_h = \sum_{m=1}^M y_m \zeta_m \in Y_h$ gilt

$$\|y_h\|_Y^2 = y \cdot \mathbb{M}y.$$

Da Y_h in der gewählten Diskretisierung ein Hilbertraum ist, ist folgende Definition von \mathbb{M} möglich

$$\mathbb{M}_{m,k} := \langle \zeta_m, \zeta_k \rangle_Y \quad \text{für } m, k = 1, \dots, M.$$

Außerdem sei der Vektor der rechten Seite $\mathbb{F} \in \mathbb{R}^M$ als

$$\mathbb{F}_m := F(\zeta_m) \quad \text{für } m = 1, \dots, M$$

festgelegt.

In der in Abschnitt 4.2 gewählten Basis für X_h wird die lineare Nebenbedingung, $\int_{\Omega} \text{tr } \sigma_0 \, dx = 0$ für alle $x_h = (\sigma_0, u_0, s_1, t_0) \in X_h$ nicht direkt realisiert. Sei daher

$\Lambda : X_h \rightarrow \mathbb{R}$ das lineare Funktional, das diese Nebenbedingung, durch $\Lambda(x_h) = 0$ beschreibt und $\mathbb{L} \in \mathbb{R}^N$ der zugehörige Vektor

$$\mathbb{L}_n := \Lambda(\xi_n) \quad \text{für } n = 1, \dots, N.$$

Mit Hilfe dieser Definitionen kann das gegebene Minimierungsproblem als Lösung eines linearen Gleichungssystems erkannt werden.

Theorem 4.1. *Es sei $x_h \in X_h$ und $x \in \mathbb{R}^N$ der zugehörige Koeffizientenvektor mit $x_h = \sum_{n=1}^N x_n \xi_n$, dann gilt*

$$x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}$$

genau dann, wenn für ein $z \in \mathbb{R}$ gilt

$$\begin{pmatrix} \mathbb{B}^\top \mathbb{M}^{-1} \mathbb{B} & \mathbb{L} \\ \mathbb{L}^\top & 0 \end{pmatrix} \begin{pmatrix} x \\ z \end{pmatrix} = \begin{pmatrix} \mathbb{B}^\top \mathbb{M}^{-1} \mathbb{F} \\ 0 \end{pmatrix}. \quad (4.1)$$

Beweis. Dieser Beweis ist ähnlich zu dem in [Hel14, Satz 7.10]. Da die Norm nicht negativ ist, kann statt $x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}$ auch

$$x_h \in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}^2$$

untersucht werden. Es sei $\xi \in \mathbb{R}^N$ der Koeffizientenvektor zu einem beliebigen $\xi_h \in X_h$ und $y \in \mathbb{R}^M$ der Koeffizientenvektor zu einem beliebigen $y_h \in Y_h$. Dann gilt nach den obigen Definitionen $b(\xi_h, y_h) = y \cdot \mathbb{B}\xi$, $\|y_h\|_Y^2 = y \cdot \mathbb{M}y$ und $F(y_h) = y \cdot \mathbb{F}$. So ergibt sich nach Definition der Dualnorm

$$\|F - b(\xi_h, \bullet)\|_{Y_h^*}^2 = \sup_{\substack{y_h \in Y_h, \\ \|y_h\|_Y = 1}} (F(y_h) - b(\xi_h, y_h))^2 = \sup_{\substack{y \in \mathbb{R}^M, \\ y \cdot \mathbb{M}y = 1}} (y \cdot (\mathbb{F} - \mathbb{M}\xi))^2.$$

Aus Gründen der Übersichtlichkeit sei $t := (\mathbb{F} - \mathbb{M}\xi)$. Es handelt sich also um die Extremalaufgabe $\sup_y (y \cdot t)^2$ zu finden unter der Nebenbedingung $y \in \varphi^{-1}(0)$, wobei $\varphi(y) := y \cdot \mathbb{M}y - 1$. Daher kann die Multiplikatorregel von Lagrange angewendet werden [Kön04, S.124]. Der Gradient der Zielfunktion $f(y) := (y \cdot t)^2$ ist $\nabla f(y) = 2t$ und für die Nebenbedingungsfunktion gilt $\nabla \varphi(y) = 2\mathbb{M}y$. Nach der Multiplikatorregel ist $y_0 \in \mathbb{R}^M$ eine Extremalstelle, wenn $\nabla f(y_0) \in \operatorname{span} \{\nabla \varphi(y_0)\}$, d.h. für ein $\lambda \in \mathbb{R}$ gilt

$$t = \lambda \mathbb{M}y_0. \quad (4.2)$$

Diese Bedingung ist bereits hinreichend für eine Maximalstelle y_0 , da die stetige Zielfunktion f ihr Maximum auf der kompakten Menge $\varphi^{-1}(0)$ annimmt. Außerdem erfüllt y_0 die Nebenbedingung, also

$$y_0 \cdot \mathbb{M}y_0 = 1. \quad (4.3)$$

Auf der einen Seite ergibt (4.3) eingesetzt in (4.2) die Gleichung $y_0 \cdot t = \lambda$, d.h.

$$\sup_{\substack{y \in \mathbb{R}^M, \\ y \cdot \mathbb{M}y = 1}} f(y) = f(y_0) = (y_0 \cdot t)^2 = \lambda^2. \quad (4.4)$$

Andererseits ist die Normmatrix \mathbb{M} positiv definit und symmetrisch also invertierbar. Daher ergibt sich aus (4.2) $\mathbb{M}^{-1}t = \lambda y_0$ und damit $t \cdot \mathbb{M}^{-1}t = \lambda t \cdot y_0 = \lambda y_0 \cdot t = \lambda^2$ nach Gleichung (4.4)

$$\max_{\substack{y \in \mathbb{R}^M, \\ y \cdot \mathbb{M}y = 1}} f(y) = (y_0 \cdot t)^2 = \lambda^2 = t \cdot \mathbb{M}^{-1}t.$$

Die gesuchte Norm kann daher wie folgt geschrieben werden

$$\|F - b(\xi_h, \bullet)\|_{Y_h^*}^2 = (\mathbb{F} - \mathbb{B}\xi) \cdot \mathbb{M}^{-1}(\mathbb{F} - \mathbb{B}\xi). \quad (4.5)$$

Durch Ausmultiplizieren ergibt sich nun

$$\|F - b(\xi_h, \bullet)\|_{Y_h^*}^2 = \xi \cdot (\mathbb{B}^\top \mathbb{M} \mathbb{B}) \xi - 2(\mathbb{F} \cdot \mathbb{M}^{-1} \mathbb{B}) \xi + \mathbb{F} \cdot \mathbb{M}^{-1} \mathbb{F}.$$

Die Definition von \mathbb{L} bedeutet, dass $\mathbb{L} \cdot \xi = \Lambda(\xi)$, daher gilt $\xi_h \in X_h$ sobald für den Koeffizientenvektor $\mathbb{L} \cdot \xi = 0$. Also ergibt sich folgende Äquivalenz

$$\begin{aligned} x_h &\in \operatorname{argmin}_{\xi_h \in X_h} \|F - b(\xi_h, \bullet)\|_{Y_h^*}^2 \\ &\Leftrightarrow x \in \operatorname{argmin}_{\substack{\xi \in \mathbb{R}^N, \\ \mathbb{L} \cdot \xi = 0}} \xi \cdot (\mathbb{B}^\top \mathbb{M} \mathbb{B}) \xi - 2(\mathbb{F} \cdot \mathbb{M}^{-1} \mathbb{B}) \xi + \mathbb{F} \cdot \mathbb{M}^{-1} \mathbb{F} \\ &\Leftrightarrow x \in \operatorname{argmin}_{\substack{\xi \in \mathbb{R}^N, \\ \mathbb{L} \cdot \xi = 0}} 1/2 \xi \cdot (\mathbb{B}^\top \mathbb{M} \mathbb{B}) \xi - (\mathbb{F} \cdot \mathbb{M}^{-1} \mathbb{B}) \xi. \end{aligned}$$

Eine erneute Anwendung der Multiplikatorregel von Lagrange mit Parameter $z \in \mathbb{R}$ zur Berechnung einer Minimalstelle x unter Nebenbedingung $\mathbb{L} \cdot x = 0$ ergibt

$$(\mathbb{B}^\top \mathbb{M} \mathbb{B}) x - \mathbb{B}^\top \mathbb{M}^{-1} \mathbb{F} = -z \mathbb{L}.$$

Dies und die Nebenbedingung $\mathbb{L} \cdot x = 0$ ergeben das behauptete Gleichungssystem. \square

In Abschnitt 3.1 wurde die folgende Fehlerabschätzung für $\xi_h \in X_h$ und die exakte Lösung $x \in X$ bei einer dPG-Methode vorgestellt

$$\|x - \xi_h\|_X \lesssim \|F - b(\xi_h, \bullet)\|_{Y_h^*} + \|F \circ (1 - \Pi)\|_{Y^*} \lesssim \|x - \xi_h\|_X.$$

Im Allgemeinen ist dabei der Datenapproximationsfehler $\|F \circ (1 - \Pi)\|_{Y^*}$ von höherer Ordnung. In Abschnitt 3.4 konnte dies leider nicht nachgewiesen werden, doch numerische Beispiele zeigen, dass

$$\eta := \|F - b(x_h, \bullet)\|_{Y_h^*}$$

ein angemessener Fehlerschätzer ist, wenn $x_h \in X_h$ die diskrete Lösung bezeichnet.

Bemerkung 4.2 (Globaler und lokaler Fehlerschätzer) Sei $x \in \mathbb{R}^N$ der Koeffizientenvektor der diskreten Lösung $x_h \in X_h$. Aus dem Beweis von Theorem 4.1 ist in Gleichung (4.5) direkt eine Formel zur Berechnung von η abzulesen. Ein geeigneter globaler Fehlerschätzer kann mittels

$$\eta^2 = \|F - b(x_h, \bullet)\|_{Y_h^*}^2 = (\mathbb{F} - \mathbb{B}x) \cdot \mathbb{M}^{-1} (\mathbb{F} - \mathbb{B}x)$$

berechnet werden. Bei den Basisfunktionen aus Abschnitt 4.2 existiert für jede der $1 \leq m \leq M$ gebrochenen Testfunktionen genau ein $T \in \mathcal{T}$, so dass $\text{supp } \zeta_m \subseteq T$. Damit gilt für

$$Y_h(T) := \{\zeta_m : 1 \leq m \leq M \text{ und } \text{supp } \zeta_m \subseteq T\},$$

dass $Y_h = \bigcup_{T \in \mathcal{T}} Y_h(T)$, und dass die einzelnen $Y_h(T)$ disjunkt sind. Daher zerfällt die Normmatrix \mathbb{M} in lokale Blöcke, wie in Abschnitt 4.3.2 deutlich wird,

$$\mathbb{M}(T) := \begin{bmatrix} \mathbf{T}(T) & 0 \\ 0 & \mathbf{V}(T) \end{bmatrix} \in \mathbb{R}^{12 \times 12},$$

wobei $\mathbb{M}(T)$, $\mathbf{T}(T)$ und $\mathbf{V}(T)$ jeweils alle Beiträge von Basisfunktionen aus $Y_h(T)$ enthalten. Damit kann für die gewählte Basis ein lokaler Fehlerschätzer für $T \in \mathcal{T}$ durch

$$\eta^2(T) := \|F - b(x_h, \bullet)\|_{Y_h(T)^*}^2 = (\mathbb{F} - \mathbb{B}x) \cdot \mathbb{M}^{-1}(T) (\mathbb{F} - \mathbb{B}x)$$

definiert werden. Aufgrund der Blockstruktur der inversen Normmatrix (siehe Abschnitt 4.3.2) und der erwähnten Eigenschaften von $Y_h(T)$ ergibt sich

$$\eta^2 = \sum_{T \in \mathcal{T}} \eta^2(T).$$

In dem Programm wird dieser lokale Fehlerschätzer zur Adaptiven Verfeinerung verwendet.

Um die Matrizen \mathbb{B} , \mathbb{M} , \mathbb{F} und \mathbb{L} zur Lösung des linearen Gleichungssystems (4.1) in Abschnitt 4.3 berechnen zu können, wird zunächst im folgenden Paragraphen die gewählte Basis vorgestellt.

4.2 Basiswahl

Bei der Basiswahl für die diskreten Räume X_h und Y_h zu der Triangulierung \mathcal{T} werden die linearen Nebenbedingungen zunächst nicht berücksichtigt. Im Ansatzraum wird die Nebenbedingung, wie in Abschnitt 4.1 erläutert, durch einen Lagrange-Multiplikator mittels des Vektors \mathbb{L} garantiert und im Testraum wird sie laut dem folgenden Lemma 4.3 nicht benötigt.

Lemma 4.3. Sei $\hat{Y}_h = RT_0^{pw}(\mathcal{T}; \mathbb{R}^{2 \times 2}) \times P_1(\mathcal{T}; \mathbb{R}^2)$ und Y_h wie in (2.12), dann gilt für jedes $x_h = (\sigma_0, u_0, s_1, t_0) \in X_h$ wie in (2.11)

$$\sup_{\hat{y}_h \in \hat{Y}_h \setminus \{0\}} \frac{b(x_h, \hat{y}_h)}{\|\hat{y}_h\|_Y} = \sup_{y_h \in Y_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|y_h\|_Y}.$$

Beweis. Dieser Beweis verläuft analog zu dem von [Hel14, Lemma 5.8]. Zunächst gilt wegen $Y_h \subseteq \hat{Y}_h$, dass

$$\sup_{y_h \in Y_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|y_h\|_Y} \leq \sup_{\hat{y}_h \in \hat{Y}_h \setminus \{0\}} \frac{b(x_h, \hat{y}_h)}{\|\hat{y}_h\|_Y}.$$

Um die andere Abschätzung zu zeigen, sei $\hat{y}_h = (\hat{\tau}_{RT}, v_1) \in \hat{Y}_h \setminus \{0\}$ beliebig. Wähle $y_h = (\tau_{RT}, v_1) \in Y_h$ mit $\tau_{RT} = \hat{\tau}_{RT} - cI_{2 \times 2}$ wobei $c = \int_{\Omega} \text{tr } \hat{\tau}_{RT} \, dx / 2$. Es gilt

$$\langle y_h, (cI_{2 \times 2}, 0) \rangle_Y = \int_{\Omega} \tau_{RT} : cI_{2 \times 2} \, dx + \int_{\Omega} \text{div}_{NC} \tau_{RT} \cdot \text{div}_{NC} I_{2 \times 2} \, dx = c \int_{\Omega} \text{tr } \tau_{RT} \, dx = 0$$

und damit nach Satz des Pythagoras

$$\|y_h\|_Y^2 \leq \|y_h\|_Y^2 + \|(cI_{2 \times 2}, 0)\|_Y^2 = \|\hat{y}_h\|_Y^2.$$

Außerdem gilt wegen den Randbedingungen an s_1 , der Nebenbedingung an σ_0 und Lemma 2.5

$$\begin{aligned} b(x_h, (cI_{2 \times 2}, 0)) &= -\langle \gamma_{\nu}^T cI_{2 \times 2}, s_1 \rangle_{\partial \mathcal{T}} + \int_{\Omega} \text{dev } \sigma_0 : cI_{2 \times 2} \, dx + \int_{\Omega} u \cdot \text{div}_{NC} cI_{2 \times 2} \, dx \\ &= -\int_{\partial \Omega} \gamma_{\nu} cI_{2 \times 2} \cdot s_1 \, ds + c \int_{\Omega} \text{tr } \sigma_0 \, dx = 0 \end{aligned}$$

und damit $b(x_h, y_h) = b(x_h, \hat{y}_h)$. Da \hat{y}_h beliebig war, ergibt sich

$$\sup_{\hat{y}_h \in \hat{Y}_h \setminus \{0\}} \frac{b(x_h, \hat{y}_h)}{\|\hat{y}_h\|_Y} \leq \sup_{y_h \in Y_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|y_h\|_Y},$$

und somit die Behauptung. □

Im Folgenden werden Basen für die neu definierten Räume

$$\begin{aligned} X_h &:= P_0(\mathcal{T}; \mathbb{R}^{2 \times 2}) \times P_0(\mathcal{T}; \mathbb{R}^2) \times S_0^1(\mathcal{E}; \mathbb{R}^2) \times P_0(\mathcal{E}; \mathbb{R}^2), \\ Y_h &:= RT_0^{pw}(\mathcal{T}; \mathbb{R}^{2 \times 2}) \times P_1(\mathcal{T}; \mathbb{R}^2) \end{aligned}$$

vorgestellt.

Dazu werden die Dreiecke T_1, \dots, T_J mit $J := |\mathcal{T}|$ der regulären Triangulierung \mathcal{T} des polygonal berandeten Lipschitz Gebietes $\Omega \subseteq \mathbb{R}^2$, die Menge der Knoten $\mathcal{N} = \{z_1, \dots, z_{|\mathcal{N}|}\}$ und der Kanten $\mathcal{E} = \{E_1, \dots, E_{|\mathcal{E}|}\}$ betrachtet. Die lokale Nummerierung

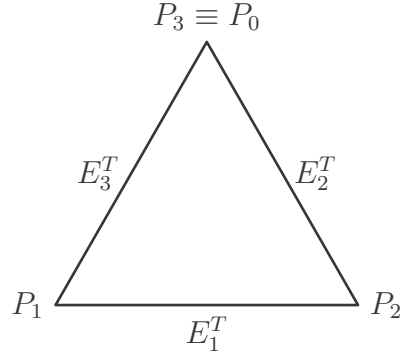


Abbildung 4.1: Lokale Nummerierung der Seiten und Eckpunkte eines Dreiecks $T \in \mathcal{T}$

der Seiten und Eckpunkte in einem Dreieck $T \in \mathcal{T}$ ist dabei stets wie in Abbildung 4.1 dargestellt.

Daneben wird als Anzahl der Freiheitsgrade mit $\text{nrDoF} = \dim(X_h)$ die Dimension des diskreten Ansatzraums bezeichnet.

Des Weiteren sei im Folgenden $\chi(T) \in L^2(\Omega)$ die charakteristische Funktion für ein Dreieck $T \in \mathcal{T}$, also gilt $\chi(T) \equiv 1$ auf T und $\chi(T) \equiv 0$ sonst, und $\chi(E) \in L^2(\partial\mathcal{T})$ die charakteristische Funktion für eine Kante $E \in \mathcal{E}$, also gilt $\chi(E) \equiv 1$ auf E und $\chi(E) \equiv 0$ sonst.

Außerdem sei die Standardbasis des \mathbb{R}^2 durch $e_1 := (1, 0)^\top$ und $e_2 := (0, 1)^\top$ gegeben, wohingegen

$$\mathbf{e}_1 := \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{e}_2 := \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}, \quad \mathbf{e}_3 := \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \mathbf{e}_4 := \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}$$

die Basis des $\mathbb{R}^{2 \times 2}$ sei.

Ebenfalls benötigt werden φ_k , die konforme nodale Basisfunktion erster Ordnung zum Knoten $z_k \in \mathcal{N}$ für $k = 1, \dots, |\mathcal{N}|$ respektive der Triangulierung \mathcal{T} , d.h. $\varphi_k \in P_1(\mathcal{T})$ mit $\varphi_k(z_j) = \delta_{j,k}$ für $z_j \in \mathcal{N}$. Für jedes $T \in \mathcal{T}$ und $\gamma = 1, 2, 3$ sei die lokale nodale Basisfunktion $\varphi(T, \gamma) := \chi(T)\varphi_k$, wobei γ die lokale Nummer des Knoten z_k in T sei. Schließlich werden die lokalen kantenbasierten Raviart-Thomas Basisfunktionen für jedes $T \in \mathcal{T}$ und die lokale Kantenummerierung $\alpha = 1, 2, 3$ mit

$$\psi(T, \alpha) := \chi(T) (x - P_{\alpha-1}) \Big| E_\alpha^T \Big| / |2T|$$

bezeichnet, wobei nach der lokalen Nummerierung aus Abbildung 4.1 $P_{\alpha-1}$ der Eckpunkt gegenüber der Seite E_α^T in T ist.

Aus Gründen der Übersichtlichkeit werden im Folgenden mit σ^b Basisfunktionen zur Lösungskomponente σ_0 bezeichnet, mit u^b die zu u_0 , mit s^b die zu s_1 , mit t^b die zu t_0 und genauso mit τ^b die Basisfunktionen zur Komponente τ_{RT} und mit v^b die zu v_1 . In der Basis der Produkträume X_h und Y_h werden für jeden Teilraum die Basisfunktionen,

die Träger auf einem Element $T \in \mathcal{T}$ haben, hintereinander sortiert. Dies ermöglicht eine Blockstruktur der Matrix.

Spannungsansatzfunktionen. Für die Elemente $j = 1, \dots, J$ und die Komponenten $\lambda = 1, \dots, 4$ sei durch

$$\sigma_{4(j-1)+\lambda}^b := \sigma^B(T_j, \lambda) := \chi(T_j) \mathbf{e}_\lambda$$

die Basis von $P_0(\mathcal{T}; \mathbb{R}^{2 \times 2}) = \text{span} \{ \sigma_1^b, \dots, \sigma_{4J}^b \}$ definiert.

Verschiebungsansatzfunktionen. Für die Elemente $j = 1, \dots, J$ und die Komponenten $\kappa = 1, 2$ sei durch

$$u_{2(j-1)+\kappa}^b := u^B(T_j, \kappa) := \chi(T_j) e_\kappa$$

die Basis von $P_0(\mathcal{T}; \mathbb{R}^2) = \text{span} \{ u_1^b, \dots, u_{2J}^b \}$ beschrieben.

Oberflächenverschiebungsfunktionen. Für die inneren Knoten $k = 1, \dots, |\mathcal{N}(\Omega)|$ und die Komponenten $\kappa = 1, 2$ sei mit den nodalen Basisfunktionen durch

$$s_{2(k-1)+\kappa}^b := s^B(z_k, \kappa) := \gamma_0^\mathcal{T} \varphi_k e_\kappa$$

die Basis von $S_0^1(\mathcal{E}; \mathbb{R}^2) = \text{span} \{ s_1^b, \dots, s_{2|\mathcal{N}(\Omega)|}^b \}$ gegeben.

Die Nullrandbedingungen werden im Programm wie üblich realisiert, in dem zwar die Beiträge für alle Knoten $z_k \in \mathcal{N}$ berechnet werden, das entstehende Gleichungssystem (4.1) allerdings nur in den freien Variablen, d.h. nur für die inneren Knoten, gelöst wird.

Traktionansatzfunktionen. Für die Kanten $\ell = 1, \dots, |\mathcal{E}|$ und die Komponenten $\kappa = 1, 2$ sei durch

$$t_{2(\ell-1)+\kappa}^b := t^B(E_\ell, \kappa) := \chi(E_\ell) e_\kappa$$

die Basis von $P_0(\mathcal{E}; \mathbb{R}^2) = \text{span} \{ t_1^b, \dots, t_{2|\mathcal{E}|}^b \}$ definiert.

Spannungstestfunktionen. Für die Elemente $j = 1, \dots, J$, die lokalen Kanten $\alpha = 1, 2, 3$ und die Komponenten $\kappa = 1, 2$ sei mit den lokalen Raviart-Thomas Basisfunktionen durch

$$\tau_{6(j-1)+2(\alpha-1)+\kappa}^b := \tau^B(T_j, \alpha, \kappa) := e_\kappa \otimes \psi(T_j, \alpha)$$

die Basis von $RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2}) = \text{span} \{ \tau_1^b, \dots, \tau_{6J}^b \}$ gegeben.

Verschiebungstestfunktionen Für die Elemente $j = 1, \dots, J$, die Knoten $\gamma = 1, 2, 3$ und die Komponenten $\kappa = 1, 2$ definiere durch

$$v_{6(j-1)+2(\gamma-1)+\kappa}^b := v^B(T_j, \gamma, \kappa) := \varphi(T_j, \gamma) e_\kappa$$

die Basis von $P_1(\mathcal{T}; \mathbb{R}^2) = \text{span} \{ v_1^b, \dots, v_{6J}^b \}$ mit den lokalen nodalen Basisfunktionen.

Damit sind die Koeffizientenvektoren $x \in \mathbb{R}^N$ mit $N = 6J + 2|\mathcal{N}(\Omega)| + 2|\mathcal{E}|$ und $y \in \mathbb{R}^M$ mit $M = 12J$ so organisiert, dass

$$\begin{aligned} x_h &= \left(\sum_{n=1}^{4J} x_n \boldsymbol{\sigma}_n^b, \sum_{n=1}^{2J} x_{4J+n} u_n^b, \sum_{n=1}^{2|\mathcal{N}(\Omega)|} x_{6J+n} s_n^b, \sum_{n=1}^{2|\mathcal{E}|} x_{6J+2|\mathcal{N}(\Omega)|+n} t_n^b \right), \\ y_h &= \left(\sum_{m=1}^{6J} x_m \boldsymbol{\tau}_m^b, \sum_{m=1}^{6J} x_{6J+m} v_m^b \right). \end{aligned}$$

Daraus ergibt sich die in den folgenden Paragraphen erklärte Assemblierung der Koeffizientenmatrix \mathbb{B} sowie der Normmatrix \mathbb{M} und ihrer Inversen \mathbb{M}^{-1} sowie der Vektoren der rechten Seite \mathbb{F} und der linearen Nebenbedingung \mathbb{L} .

4.3 Berechnung der Bestandteile des linearen Gleichungssystems

4.3.1 Koeffizientenmatrix

Wie in Abschnitt 4.1 ist die Koeffizientenmatrix $\mathbb{B} \in \mathbb{R}^{M \times N}$ definiert durch

$$\mathbb{B}_{m,n} := b(\xi_n, \zeta_m) \quad \text{für } n = 1, \dots, N \text{ und } m = 1, \dots, M.$$

Damit gilt für den Koeffizientenvektor $x \in \mathbb{R}^N$ von $x_h = \sum_{n=1}^N x_n \xi_n \in X_h$ und den $y \in \mathbb{R}^M$ von $y_h = \sum_{m=1}^M y_m \zeta_m \in Y_h$, dass $b(x_h, y_h) = y \cdot \mathbb{B}x$.

Für die Bilinearform aus Gleichung (3.6)

$$\begin{aligned} b(x_h, y_h) &:= -\langle t_0, \gamma_0 v_1 \rangle_{\partial\mathcal{T}} - \langle \gamma_\nu \boldsymbol{\tau}_{\text{RT}}, s_1 \rangle_{\partial\mathcal{T}} \\ &\quad + \int_{\Omega} \boldsymbol{\sigma}_0 : \text{D}_{\text{NC}} v_1 \, dx + \int_{\Omega} \text{dev } \boldsymbol{\sigma}_0 : \boldsymbol{\tau}_{\text{RT}} \, dx + \int_{\Omega} u_0 \cdot \text{div}_{\text{NC}} \boldsymbol{\tau}_{\text{RT}} \, dx. \end{aligned}$$

ergibt sich folgende 2×4 -Blockstruktur

$$\mathbb{B} = \begin{bmatrix} \mathbf{A} & \mathbf{B} & \mathbf{C} & \mathbf{0} \\ \mathbf{D} & \mathbf{0} & \mathbf{0} & \mathbf{E} \end{bmatrix} \quad (4.6)$$

mit $\mathbf{A} \in \mathbb{R}^{6J \times 4J}$, $\mathbf{B} \in \mathbb{R}^{6J \times 2J}$, $\mathbf{C} \in \mathbb{R}^{6J \times 2|\mathcal{N}(\Omega)|}$, $\mathbf{D} \in \mathbb{R}^{6J \times 4J}$ und $\mathbf{E} \in \mathbb{R}^{6J \times 2|\mathcal{E}|}$.

Die folgenden Matrizen vereinfachen die Notation dieser Blöcke. Für $T \in \mathcal{T}$ mit den lokalen Seiten E_1^T, E_2^T, E_3^T definiere

$$L(T) := \text{diag} \left(|E_1^T|, |E_1^T|, |E_2^T|, |E_2^T|, |E_3^T|, |E_3^T| \right) \in \mathbb{R}^{6 \times 6} \quad (4.7)$$

und für $\alpha = 1, 2, 3$ sei $\text{sgn}_\alpha^T := \text{sgn}(T, E_\alpha^T) = \nu_{E_\alpha^T} \cdot \nu_T \in \{\pm 1\}$ wie in (2.1) und

$$S(T) := \text{diag}(\text{sgn}_1^T, \text{sgn}_1^T, \text{sgn}_2^T, \text{sgn}_2^T, \text{sgn}_3^T, \text{sgn}_3^T) \in \mathbb{R}^{6 \times 6}. \quad (4.8)$$

Hilfreich sind auch die folgenden konstanten Matrizen

$$H_1 := \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{pmatrix}, \quad H_2 := \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad \text{und} \quad H_3 := \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \end{pmatrix}. \quad (4.9)$$

Im Folgenden wird die genaue Berechnung der in (4.6) erwähnten Blöcke erläutert. Die gebrochenen Testfunktionen haben jeweils nur ein Element als Träger, daher müssen die vorkommenden Integrale über das Gebiet Ω bzw. das Skelett $\partial\mathcal{T}$ jeweils nur auf diesem einen Dreieck bzw. seinem Rand ausgewertet werden.

Block A enthält die Anteile der Bilinearform für $\sigma_0 \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})$ und $\tau_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})$, die Einträge ergeben sich also aus folgendem Integral der Basisfunktionen

$$A_{mn} = \int_{\Omega} \tau_m^b : \text{dev} \sigma_n^b \, dx \quad \text{für } m = 1, \dots, 6J \text{ und } n = 1, \dots, 4J.$$

Dieses Integral verschwindet, es sei denn τ_m^b und σ_n^b haben als Träger dasselbe Dreieck $T \in \mathcal{T}$, d.h. es gilt $\tau_m^b = \tau^B(T, \alpha, \kappa)$ und $\sigma_n^b = \sigma^B(T, \lambda)$ für ein $T \in \mathcal{T}$, $\alpha = 1, 2, 3$, $\kappa = 1, 2$ und $\lambda = 1, \dots, 4$.

Daher ist es angebracht die lokalen Blöcke, die Integrale über Basisfunktionen in $Y_h(T)$ bzw. $X_h(T)$ für ein $T \in \mathcal{T}$, zu betrachten.

Lemma 4.4. Für jedes Dreieck $T \in \mathcal{T}$ sei der lokale Block $\mathbf{A}(T)$ definiert via

$$\mathbf{A}(T; 2(\alpha - 1) + \kappa, \lambda) := \int_T \tau^B(T, \alpha, \kappa) : \text{dev} \sigma^B(T, \lambda) \, dx$$

für $\kappa = 1, 2$, $\alpha = 1, 2, 3$ und $\lambda = 1, 2, 3, 4$. Es gilt

$$\mathbf{A}(T; 2(\alpha - 1) + \kappa, \lambda) = \begin{cases} 0 & \text{für } \lambda = 1, \\ \mathbf{e}_\lambda : |E_\alpha^T| / 2 (e_k \otimes (\text{mid}(T) - P_{\alpha-1})) & \text{für } \lambda = 2, 3, 4, \end{cases}$$

d.h. der Block $\mathbf{A}(T)$ hat folgende Form

$$\mathbf{A}(T) = \frac{1}{2} L(T) \left(H_1 K(T) \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix} + H_2 K(T) \begin{pmatrix} 0 & 0 & 1 & -1 \\ 0 & -1 & 0 & 0 \end{pmatrix} \right) \in \mathbb{R}^{6 \times 4},$$

vermöge der Matrix $L(T)$ aus (4.7) und der konstanten Matrizen aus (4.9) und mit der lokalen Matrix $K(T) \in \mathbb{R}^{3 \times 2}$, die für die lokalen Seiten E_α^T mit $\alpha = 1, 2, 3$ und nach Abbildung 4.1 gegenüberliegenden Eckpunkt $P_{\alpha-1}$ und $\kappa = 1, 2$, die Einträge

$$K(T; \alpha, \kappa) := (\text{mid}(T) - P_{\alpha-1}) \cdot e_\kappa$$

hat.

Beweis. Auf dem Element $T \in \mathcal{T}$ verschwindet der Deviator von $\sigma^B(T, \lambda)$ für $\lambda = 1$, wohingegen die Funktionen $\sigma^B(T, \lambda)$ für $\lambda = 2, 3, 4$ deviatorisch sind. Außerdem können die stückweise linearen lokalen Raviart-Thomas Funktionen mit der Mittelpunktsformel exakt integriert werden. Also gilt

$$\int_\Omega \psi(T, \alpha) = \int_T (x - P_{\alpha-1}) |E_\alpha^T| / (2|T|) \, dx = (\text{mid}(T) - P_{\alpha-1}) |E_\alpha^T| / 2.$$

Daher ergibt sich für $\lambda = 2, 3, 4$, $\alpha = 1, 2, 3$ und $\kappa = 1, 2$

$$\begin{aligned} \int_T \tau(T, \alpha, \kappa) : \text{dev } \sigma(T, \lambda) \, dx &= \mathbf{e}_\lambda : |E_\alpha^T| / 2 \left(e_\kappa \otimes (\text{mid}(T) - P_{\alpha-1}) \right) \\ &= \begin{cases} |E_\alpha^T| / 2 (\text{mid}(T) - P_{\alpha-1}) \cdot e_1 & \text{für } (\lambda = 2, \kappa = 1), (\lambda = 3, \kappa = 2), \\ -|E_\alpha^T| / 2 (\text{mid}(T) - P_{\alpha-1}) \cdot e_1 & \text{für } (\lambda = 4, \kappa = 2), \\ |E_\alpha^T| / 2 (\text{mid}(T) - P_{\alpha-1}) \cdot e_2 & \text{für } (\lambda = 3, \kappa = 1), (\lambda = 4, \kappa = 1), \\ -|E_\alpha^T| / 2 (\text{mid}(T) - P_{\alpha-1}) \cdot e_2 & \text{für } (\lambda = 2, \kappa = 2). \end{cases} \end{aligned}$$

$\mathbf{A}(T)$ hat also vermöge der Matrizen $K(T)$ und $L(T)$ die folgende Form

$$\mathbf{A}(T) = \frac{1}{2} L(T) \begin{pmatrix} 0 & K(T; 1, 1) & K(T; 1, 2) & K(T; 1, 2) \\ 0 & -K(T; 1, 2) & K(T; 1, 1) & -K(T; 1, 1) \\ 0 & K(T; 2, 1) & K(T; 2, 2) & K(T; 2, 2) \\ 0 & -K(T; 2, 2) & K(T; 2, 1) & -K(T; 2, 1) \\ 0 & K(T; 3, 1) & K(T; 3, 2) & K(T; 3, 2) \\ 0 & -K(T; 3, 2) & K(T; 3, 1) & -K(T; 3, 1) \end{pmatrix} \in \mathbb{R}^{6 \times 4}.$$

Diese ist identisch zu dem behaupteten Matrixprodukt. \square

In dem Code in `computeBlocksStokesDPG` bzw. `computeBlocksStokesDPGFass` bzw. `computeBlocksStokesDPGIntegrate` ist die Berechnung der lokalen Blöcke wie folgt implementiert

```
% initialisation of matrix blocks
A4e = zeros(6,4,nrElems); % local matrices for block A
% constant auxillary matrices
H1 = sparse([1,3,5],[1,2,3],[1,1,1],6,3);
H2 = sparse([2,4,6],[1,2,3],[1,1,1],6,3);
```

```

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    coords = coord4e(:, :, elem); % coordinates of three nodes
    mid = mid4e(elem, :); % midpoint of T
    slength = sl4e(elem, :); % length of sides
    L = [slength, slength]'; L=diag(L(:));
    K = [mid; mid; mid] - [coords(3, :); coords(1, :); coords(2, :)];

    % block A int_T(dev(sigma):tau)
    A4e(:, :, elem) = 1/2 * L * (H1 * K * [0, 1, 0, 0; 0, 0, 1, 1] + H2 * K * [0, 0, 1, -1; 0, -1, 0, 0]);
end

```

Der Matrixblock \mathbf{A} wird aus den lokalen Matrizen für alle $\alpha = 1, 2, 3$, $\kappa = 1, 2$, $\lambda = 1, 2, 3, 4$ und $j = 1, \dots, J$ wie folgt assembliert

$$\mathbf{A}_{6(j-1)+2(\alpha-1)+\kappa, 4(j-1)+\lambda} = \mathbf{A}(T_j; 2(\alpha-1) + \kappa, \lambda).$$

In diesem Fall, entsteht also folgende Blockdiagonalform

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}(T_1) & 0 & \dots & 0 \\ 0 & \mathbf{A}(T_2) & \dots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & \mathbf{A}(T_J) \end{pmatrix} \in \mathbb{R}^{6J \times 4J}.$$

Die globale Assemblierung geschieht in dem Programm mittels folgender Zeilen

```

%% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems, 4, 1), 4, 6, nrElems), [2, 1, 3]);
tempJ = repmat(1:4*nrElems, 6, 1);
A = sparse(tempI(:), tempJ(:), A4e(:)); % block A

```

Block B enthält die Beiträge für $u_0 \in P_0(\mathcal{T}; \mathbb{R}^2)$ und $\boldsymbol{\tau}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})$. Die Einträge ergeben sich also aus folgendem Integral

$$\mathbf{B}_{mn} = \int_{\Omega} \text{div}_{\text{NC}}(\boldsymbol{\tau}_m^b) \cdot u_n^b \, dx \quad \text{für } m = 1, \dots, 6J \text{ und } n = 1, \dots, 2J.$$

Dieses Integral verschwindet, außer $\boldsymbol{\tau}_m^b$ und u_n^b haben dasselbe Dreieck $T \in \mathcal{T}$ als Träger, d.h. $\boldsymbol{\tau}_m^b = \boldsymbol{\tau}^B(T, \alpha, \kappa)$ und $u_n^b = u^B(T, \varkappa)$ für ein $T \in \mathcal{T}$, $\alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$, und in diesem Fall ist es lediglich ein Integral über dieses Element $T \in \mathcal{T}$

Daher ist es sinnvoll, die im folgenden Lemma berechneten lokalen Blöcke zu definieren.

Lemma 4.5. Für jedes Dreieck $T \in \mathcal{T}$ sei der lokale Block $\mathbf{B}(T)$ definiert via

$$\mathbf{B}(T; 2(\alpha-1) + \kappa, \varkappa) := \int_T \text{div} \boldsymbol{\tau}^B(T; \alpha, \kappa) \cdot u^B(T, \varkappa) \, dx$$

für $\alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$. Es gilt

$$\mathbf{B}(T; 2(\alpha-1) + \kappa, \varkappa) = |E_{\alpha}^T| \delta_{\kappa, \varkappa},$$

d.h. der Block $\mathbf{B}(T)$ lässt sich wie folgt berechnen

$$\mathbf{B}(T) = L(T) \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{6 \times 2},$$

vermöge der Matrix $L(T)$ aus (4.7).

Beweis. Für die lokalen Raviart-Thomas Basisfunktionen $\psi(T, \alpha)$ gilt

$$\int_{\Omega} \operatorname{div} \psi(T, \alpha) \, dx = \int_T |E_{\alpha}^T| / |T| \, dx = |E_{\alpha}^T|.$$

Daher ergibt sich für $\kappa, \varkappa = 1, 2$ und $\alpha = 1, 2, 3$

$$\int_T \operatorname{div} \boldsymbol{\tau}^B(T, \alpha, \kappa) \cdot u^B(T, \varkappa) \, dx = |E_{\alpha}^T| \delta_{\kappa, \varkappa}$$

und $\mathbf{B}(T)$ hat die behauptete Form. □

Die entsprechenden Codezeilen für die Berechnung der lokalen Blöcke lauten

```
% initialisation of matrix blocks
B4e = zeros(6,2,nrElems);           % local matrices for block B

%% Parfor-loop to fill local matrices
for elem = 1:nrElems                %elem as current triangle T
    slength = sl4e(elem,:);          % length of sides
    L       = [slength,slength]';    L=diag(L(:));

    % block B int_T(u.div(tau))
    B4e(:, :, elem)=L*[1,0;0,1;1,0;0,1;1,0;0,1];
end
```

In der Koeffizientenmatrix \mathbb{B} wird der Block \mathbf{B} aus den lokalen Matrizen $\mathbf{B}(T)$, für alle $j = 1, \dots, J$, $\alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$ durch

$$\mathbf{B}_{6(j-1)+2(\alpha-1)+\kappa, 2(j-1)+\varkappa} = \mathbf{B}(T_j; 2(\alpha-1) + \kappa, \varkappa)$$

bestimmt. Auch hier ergibt sich die bekannte Blockdiagonalform

$$\mathbf{B} = \begin{pmatrix} \mathbf{B}(T_1) & 0 & \dots & 0 \\ 0 & \mathbf{B}(T_2) & \dots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & \mathbf{B}(T_J) \end{pmatrix} \in \mathbb{R}^{6J \times 2J}.$$

In dem Programm geschieht dieses Assemblieren in

```

%% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems,2,1),2,6,nrElems),[2,1,3]);
tempJ = repmat(1:2*nrElems,6,1);
B=sparse(tempI(:),tempJ(:),B4e(:));      % block B

```

Block C enthält die Beiträge für $s_1 \in S_0^2(\partial\mathcal{T}; \mathbb{R}^2)$ und $\boldsymbol{\tau}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})$. Die Einträge sind also durch folgendes Integral bestimmt

$$\mathbf{C}_{mn} = - \left\langle \gamma_\nu \boldsymbol{\tau}_m^b, s_n^b \right\rangle_{\partial\mathcal{T}} \quad \text{für } m = 1, \dots, 6J \text{ und } n = 1, \dots, 2|\mathcal{N}(\Omega)|.$$

Auch dieses verschwindet, es sei denn die Träger von $\boldsymbol{\tau}_m^b$ und s_n^b sind Teil desselben Dreiecks $T \in \mathcal{T}$, d.h. $\boldsymbol{\tau}_m^b = \boldsymbol{\tau}^B(T, \alpha, \kappa)$ und $s_n^b = s^B(z, \varkappa)$ mit $z \in \mathcal{N}(T)$ für ein $T \in \mathcal{T}$, $\alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$ und in diesem Fall wird lediglich über den Rand dieses $T \in \mathcal{T}$ integriert.

Die im folgenden Lemma bestimmten lokalen Blöcke sind daher nützlich zur Berechnung von \mathbf{C} .

Lemma 4.6. *Für jedes Element $T \in \mathcal{T}$, wird der lokale Block $\mathbf{C}(T)$ definiert durch*

$$\mathbf{C}(T; 2(\alpha - 1) + \kappa, 2(\gamma - 1) + \varkappa) := - \int_{\partial T} \gamma_\nu \left(\boldsymbol{\tau}^B(T, \alpha, \kappa) \right) \cdot s^B(P_\gamma, \varkappa) \, ds,$$

für $\kappa, \varkappa = 1, 2$, die lokale Nummerierung der Kanten durch $\alpha = 1, 2, 3$ und die lokale Nummerierung der Knoten durch $\gamma = 1, 2, 3$. Es ergeben sich folgende Einträge

$$\mathbf{C}(T; 2(\alpha - 1) + \kappa, 2(\gamma - 1) + \varkappa) = \begin{cases} -\delta_{\kappa, \varkappa} |E_\alpha^T| / 2 & \text{für } P_\gamma \in \mathcal{N}(E_\alpha^T), \\ 0 & \text{sonst.} \end{cases}$$

und damit folgende Form des lokalen Blocks

$$\mathbf{C}(T) = -\frac{1}{2} L(T) H_3 \in \mathbb{R}^{6 \times 6},$$

wobei $L(T)$ wie in (4.7) und H_3 wie in in (4.9) seien.

Beweis. Die lokalen Raviart-Thomas Funktionen $\psi(T, \alpha)$ verschwinden entlang von Kanten $F \notin \mathcal{E}(T)$. Außerdem ist $(x - P_{\alpha-1})$ entlang der anliegenden Kanten, also für $x \in F \in \mathcal{E}(T) \setminus \{E_\alpha^T\}$, tangential zu ∂T und entlang der gegenüberliegenden Kante, also für $x \in E_\alpha^T$, ist $(x - P_{\alpha-1}) \cdot \nu_T$ genau die Höhe des Dreiecks, d.h. die Konstante $2|T| / |E_\alpha^T|$. Daher gilt

$$(\psi(T, \alpha) \cdot \nu_T)|_F = |E_\alpha^T| / (2|T|) ((x - P_{\alpha-1}) \cdot \nu_T|_F) = \begin{cases} 1 & \text{für } F = E_\alpha^T, \\ 0 & \text{für } F \neq E_\alpha^T. \end{cases}$$

Dies ist auch in [BC05, Lemma 4.1] bewiesen. Zusätzlich gilt für die Nodalen Basisfunktion auf dem Rand $\gamma_0^T \varphi_k$, dass

$$\int_E \gamma_0 \varphi_k \, ds = \begin{cases} |E|/2 & \text{für } z_k \in \mathcal{N}(E), \\ 0 & \text{sonst.} \end{cases}$$

Damit ergibt sich für die lokale Nummerierung der Seiten mit $\alpha = 1, 2, 3$, $\kappa, \varkappa = 1, 2$, und die lokale Knotenummerierung $\gamma = 1, 2, 3$,

$$\int_{\partial T} \left(\tau^B(T, \alpha, \kappa) \nu_T \right) \cdot s^B(P_\gamma, \varkappa) \, dx = \begin{cases} \delta_{\kappa, \varkappa} |E_\alpha^T|/2 & \text{für } P_\gamma \in \mathcal{N}(E_\alpha^T), \\ 0 & \text{sonst.} \end{cases}$$

In Abbildung 4.1 wird deutlich, dass sich daraus die behauptete Form ergibt

$$\mathbf{C}(T) = -\frac{1}{2}L(T) \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \end{pmatrix} = -\frac{1}{2}L(T) H_3 \in \mathbb{R}^{6 \times 6}.$$

□

In dem Programmcode werden die folgenden lokalen Rechnungen durchgeführt

```
% initialisation of matrix blocks
C4e = zeros(6,6,nrElems); % local matrices for block C
IndJC4e=zeros(6,6,nrElems); % columns for block C
% constant auxillary matrices
H3 = sparse([1:6,1:6],[1:6,3:6,1:2],ones(12,1),6,6);

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    nodes = n4e(elem,:); % nodes for element
    slength = sl4e(elem,:); % length of sides
    L = [slength,slength]'; L=diag(L(:));

    % block C int_dT(gamma_nu(tau).s)
    C4e(:, :, elem)=-1/2*L*H3;
    % index transformation (local to global)
    temp = [2*(nodes-1)+1;2*(nodes-1)+2];
    IndJC4e(:, :, elem)=repmat(temp(:)',6,1);
end
```

Aus diesen lokalen Matrizen wird der Matrixblock \mathbf{C} , für alle $j = 1, \dots, J$, $\alpha = 1, 2, 3$, $\kappa, \varkappa = 1, 2$, und die globale Nummerierung $k \in \{1, \dots, |\mathcal{N}(\Omega)|\}$ der Knoten P_γ für $\gamma = 1, 2, 3$ wie folgt assembliert

$$\mathbf{C}_{6(j-1)+2(\alpha-1)+\kappa, 2(k-1)+\varkappa} = \mathbf{C}(T_j; 2(\alpha-1) + \kappa, 2(\gamma-1) + \varkappa).$$

In diesem Fall, werden die lokalen Blöcke also aufgeteilt, indem die einzelnen Spalten der jeweiligen globalen Knotennummer zugeordnet und so umsortiert werden. Dies geschieht in folgenden Zeilen

```
%% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems,6,1),6,6,nrElems),[2,1,3]);
C=sparse(tempI(:),IndJC4e(:),C4e(:)); % block C
```

Block D enthält die Anteile der Bilinearform für $\sigma_0 \in P_0(\mathcal{T}; \mathbb{R}^{2 \times 2})$ und $v_1 \in P_1(\mathcal{T}; \mathbb{R}^2)$, die Einträge ergeben sich also durch

$$\mathbf{D}_{mn} = \int_{\Omega} D_{\text{NC}} v_m^b : \sigma_n^b \, dx \quad \text{für } m = 1, \dots, 6J \text{ und } n = 1, \dots, 4J.$$

Sie verschwinden, es sei denn v_m^b und σ_n^b haben dasselbe $T \in \mathcal{T}$ als Träger, d.h. $v_m^b = v^B(T, \gamma, \kappa)$ und $\sigma_n^b = \sigma^B(T, \lambda)$ für ein $T \in \mathcal{T}$, $\gamma = 1, 2, 3$, $\kappa = 1, 2$ und $\lambda = 1, \dots, 4$. In diesem Fall ist lediglich über dieses Dreieck $T \in \mathcal{T}$ zu integrieren.

Es erweist sich daher als hilfreich, die lokalen Blöcke zu untersuchen.

Lemma 4.7. *Für jedes Dreieck $T \in \mathcal{T}$ wird der lokale Block $\mathbf{D}(T)$ definiert als*

$$\mathbf{D}(T; 2(\gamma - 1) + \kappa, \lambda) := \int_T D v^B(T, \gamma, \kappa) : \sigma^B(T, \lambda) \, dx$$

für $\kappa = 1, 2$, $\gamma = 1, 2, 3$, und $\lambda = 1, 2, 3, 4$. Es ergibt sich

$$\mathbf{D}(T; 2(\gamma - 1) + \kappa, \lambda) = |T| \left((e_\kappa \otimes \nabla \varphi_{P_\gamma}) : \mathbf{e}_\lambda \right)$$

und damit folgende Form des Blockes

$$\mathbf{D}(T) = |T| \left(H_1 G(T) \begin{pmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{pmatrix} + H_2 G(T) \begin{pmatrix} 0 & 0 & 1 & -1 \\ 1 & -1 & 0 & 0 \end{pmatrix} \right) \in \mathbb{R}^{6 \times 4},$$

dabei sei die Matrix $G(T)$ für die lokale Nummerierung P_1, P_2, P_3 der Eckpunkte von T

$$G(T) := \begin{pmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{3 \times 2}$$

und H_1, H_2 seien die konstanten Matrizen aus (4.9).

Beweis. Wie in [ACF99, S.122] sind die Ableitungen der Nodalen Basisfunktionen zu P_1, P_2, P_3 auf T durch die folgenden Konstanten beschrieben

$$\begin{pmatrix} \nabla \varphi_{P_1} \\ \nabla \varphi_{P_2} \\ \nabla \varphi_{P_3} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}.$$

Daher ergibt sich für $\kappa = 1, 2$, $\lambda = 1, \dots, 4$, und $\gamma = 1, 2, 3$

$$\begin{aligned} \int_T D v^B(T, \gamma, \kappa) : \sigma^B(T, \lambda) \, dx \\ = |T| \left((e_\kappa \otimes \nabla \varphi_{P_\gamma}) : e_\lambda \right) \\ = \begin{cases} |T| \nabla \varphi_{P_\gamma} \cdot e_1 & \text{für } (\lambda = 1, \kappa = 1), (\lambda = 2, \kappa = 1), (\lambda = 3, \kappa = 2), \\ -|T| \nabla \varphi_{P_\gamma} \cdot e_1 & \text{für } (\lambda = 4, \kappa = 2), \\ |T| \nabla \varphi_{P_\gamma} \cdot e_2 & \text{für } (\lambda = 1, \kappa = 2), (\lambda = 3, \kappa = 1), (\lambda = 4, \kappa = 1), \\ -|T| \nabla \varphi_{P_\gamma} \cdot e_2 & \text{für } (\lambda = 2, \kappa = 2). \end{cases} \end{aligned}$$

Vermöge der Matrix $G(T)$ ergibt sich

$$\mathbf{D}(T) = |T| \begin{pmatrix} G(T; 1, 1) & G(T; 1, 1) & G(T; 1, 2) & G(T; 1, 2) \\ G(T; 1, 2) & -G(T; 1, 2) & G(T; 1, 1) & -G(T; 1, 1) \\ G(T; 2, 1) & G(T; 2, 1) & G(T; 2, 2) & G(T; 2, 2) \\ G(T; 2, 2) & -G(T; 2, 2) & G(T; 2, 1) & -G(T; 2, 1) \\ G(T; 3, 1) & G(T; 3, 1) & G(T; 3, 2) & G(T; 3, 2) \\ G(T; 3, 2) & -G(T; 3, 2) & G(T; 3, 1) & -G(T; 3, 1) \end{pmatrix} \in \mathbb{R}^{6 \times 4},$$

somit hat $\mathbf{D}(T)$ die behauptete Form. □

In dem Programm werden die lokalen Blöcke wie folgt bestimmt

```
% initialisation of matrix blocks
D4e = zeros(6,4,nrElems); % local matrices for block D
% constant auxillary matrices
H1 = sparse([1,3,5],[1,2,3],[1,1,1],6,3);
H2 = sparse([2,4,6],[1,2,3],[1,1,1],6,3);

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    coords = coord4e(:, :, elem); % coordinates of three nodes
    area = area4e(elem); % area of T
    G = [1, 1, 1; coords'] \ [0, 0; 1, 0; 0, 1];

    % block D int_T(sigma:Dv)
    D4e(:, :, elem) = area * (H1 * G * [1, 1, 0, 0; 0, 0, 1, 1] + H2 * G * [0, 0, 1, -1; 1, -1, 0, 0]);
end
```

Aus diesen lokalen Matrixblöcken wird der Block \mathbf{D} für alle $j = 1, \dots, J$, $\gamma = 1, 2, 3$, $\kappa = 1, 2$ und $\lambda = 1, 2, 3, 4$ wie folgt assembliert

$$\mathbf{D}_{6(j-1)+2(\gamma-1)+\kappa, 4(j-1)+\lambda} = \mathbf{D}(T_j; 2(\gamma-1) + \kappa, \lambda).$$

Es ergibt sich also wieder eine Diagonalstruktur

$$\mathbf{D} = \begin{pmatrix} \mathbf{D}(T_1) & 0 & \dots & 0 \\ 0 & \mathbf{D}(T_2) & \dots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & \mathbf{D}(T_J) \end{pmatrix} \in \mathbb{R}^{6J \times 4J}.$$

Dies wird realisiert via

```
% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems,4,1),4,6,nrElems),[2,1,3]);
tempJ = repmat(1:4*nrElems,6,1);
D=sparse(tempI(:),tempJ(:),D4e(:));      % block D
```

Block E enthält die Anteile der Bilinearform für $t_0 \in P_0(\mathcal{T}; \mathbb{R}^2)$ und $v_1 \in P_1(\mathcal{T}; \mathbb{R}^2)$ und daher sind die Einträge via

$$\mathbf{E}_{mn} = - \left\langle t_n^b, \gamma_0 v_m^b \right\rangle_{\partial \mathcal{T}} \quad \text{für } m = 1, \dots, 6J \text{ und } n = 1, \dots, 2|\mathcal{E}|$$

definiert. Sie verschwinden, es sei denn die Träger der beiden Funktionen v_m^b und t_n^b liegen im selben Dreieck $T \in \mathcal{T}$, d.h. $v_m^b = v^B(T, \gamma, \kappa)$ und $t_n^b = t^B(E_\ell, \varkappa)$ mit $E_\ell \in \mathcal{E}(T)$ für ein $T \in \mathcal{T}$, $\gamma = 1, 2, 3$, $\kappa, \varkappa = 1, 2$ und $\ell = 1, \dots, |\mathcal{E}|$ und in diesem Fall ist jeweils nur über den Rand dieses $T \in \mathcal{T}$ zu integrieren.

Im Folgenden wird die Berechnung der zu betrachtenden lokalen Blöcke genauer analysiert.

Lemma 4.8. Für jedes Element $T \in \mathcal{T}$ wird der lokale Block $\mathbf{E}(T)$ definiert via

$$\mathbf{E}(T; 2(\gamma - 1) + \kappa, 2(\alpha - 1) + \varkappa) := - \int_{\partial T} \gamma_0 \left(v^B(T, \gamma, \kappa) \right) \cdot t^B(E_\alpha^T, \varkappa) \, ds$$

für $\kappa, \varkappa = 1, 2$, $\gamma = 1, 2, 3$ und die lokale Nummerierung $\alpha = 1, 2, 3$ der Kanten von T . Damit ergibt sich

$$\mathbf{E}(T; 2(\gamma - 1) + \kappa, 2(\alpha - 1) + \varkappa) = \begin{cases} -\delta_{\kappa, \varkappa} \operatorname{sgn}_\alpha^T |E_\alpha^T| / 2 & \text{für } P_\gamma \in \mathcal{N}(E_\alpha^T), \\ 0 & \text{sonst.} \end{cases}$$

und die Blockform

$$\mathbf{E}(T) = -\frac{1}{2} H_3^\top L(T) S(T) \in \mathbb{R}^{6 \times 6},$$

vermöge $S(T)$, $L(T)$ und H_3 wie in (4.7), (4.8) und (4.9), respektive.

Beweis. Es sei erinnert, dass das Integral der Nodalen Basisfunktion φ_k entlang einer Kante den Wert $\int_E \gamma_0^\top \varphi_k \, ds = |E|/2$ hat, wenn $z_k \in \mathcal{N}(E)$, und sonst verschwindet. Damit ergibt sich für die lokale Nummerierung $\gamma = 1, 2, 3$ der Eckpunkte, $\alpha = 1, 2, 3$

der Seiten und $\kappa, \varkappa = 1, 2$ vermöge der Einbettung $P_0(\mathcal{E}; \mathbb{R}^2) \hookrightarrow H^{-1/2}(\partial\mathcal{T}; \mathbb{R}^2)$ aus Paragraph 2.4 folgende Form

$$\int_{\partial T} v^B(T, \gamma, \kappa) t^B(E_\alpha^T, \varkappa) dx = \begin{cases} \delta_{\kappa, \varkappa} \operatorname{sgn}_\alpha^T |E_\alpha^T| / 2 & \text{für } P_\gamma \in \mathcal{N}(E_\alpha^T), \\ 0 & \text{sonst.} \end{cases}$$

Die Nummerierungskonvention aus Abbildung 4.1 zeigt, dass

$$\mathbf{E}(T) = -\frac{1}{2} \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 1 \\ 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \end{pmatrix} L(T) S(T) = -\frac{1}{2} H_3^\top L(T) S(T) \in \mathbb{R}^{6 \times 6},$$

also die Behauptung. □

In dem Programm werden folgende lokale Rechnungen durchgeführt

```
% initialisation of matrix blocks
E4e = zeros(6,6,nrElems); % local matrices for block E
IndJE4e=zeros(6,6,nrElems); % columns for block E
% constant auxillary matrices
H3 = sparse([1:6,1:6],[1:6,3:6,1:2],ones(12,1),6,6);

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    sides = s4e(elem,:); % sides of T
    slength = sl4e(elem,:); % length of sides
    sign = diag(ns4e(:, :, elem)*normal4e(:, :, elem)'); % vector with sgn(sides,elem)
    L = [slength,slength]'; L=diag(L(:));
    S = [sign,sign]'; S=diag(S(:));

    % block E int_dT(t.gamma_0(v))
    E4e(:, :, elem)=-1/2*H3'*S*L;
    % index transformation (local to global)
    temp = [2*(sides-1)+1;2*(sides-1)+2];
    IndJE4e(:, :, elem)=repmat(temp(:)',6,1);
end
```

Die lokalen Blöcke werden in \mathbf{E} für alle $j = 1, \dots, J$, $\gamma = 1, 2, 3$, $\kappa, \varkappa = 1, 2$, und den globalen Nummern $\ell \in \{1, \dots, |\mathcal{E}|\}$ der lokalen Seiten E_α^T für $\alpha = 1, 2, 3$ assembliert via

$$\mathbf{E}_{6(j-1)+2(\gamma-1)+\kappa, 2(\ell-1)+\varkappa} = \mathbf{E}(T_j; 2(\gamma-1) + \kappa, 2(\alpha-1) + \varkappa).$$

Auch hier werden die lokalen Blöcke also aufgespalten und die Spalten den jeweiligen globalen Kantennummern zugeordnet. Dies geschieht durch

```
% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems,6,1),6,6,nrElems),[2,1,3]);
E=sparse(tempI(:),IndJE4e(:),E4e(:)); % block E
```

Damit kann die Koeffizientenmatrix \mathbb{B} nun vollständig assembliert werden. Die lokalen Träger der Testfunktionen garantieren dabei, dass keine Einträge aufsummiert werden, die lokalen Beiträge müssen nur wie beschrieben geeignet in die globalen Matrixblöcke einsortiert werden. Es ergibt sich, dann

```
% coefficient matrix BB=sparse(dimY,dimX)
BB = [A,B,C,sparse(6*nrElems,2*nrSides); ...
      D,sparse(6*nrElems,2*nrElems+2*nrNodes),E];
```

Im nächsten Paragraphen wird nun der Aufbau der Normmatrix für den Testraum Y_h beschrieben.

4.3.2 Normmatrix

Im Hilbertraum $Y_h = \text{span}\{\zeta_1, \dots, \zeta_M\}$ kann die Normmatrix $\mathbb{M} \in \mathbb{R}^{M \times M}$, die für den Koeffizientenvektor $y \in \mathbb{R}^M$ von $y_h = \sum_{m=1}^M y_m \zeta_m \in Y_h$ erfüllt, dass $\|y_h\|_Y^2 = y \cdot \mathbb{M} y$, insbesondere wie folgt definiert werden

$$\mathbb{M}_{m,k} := \langle \zeta_m, \zeta_k \rangle_Y \quad \text{für } m, k = 1, \dots, M.$$

Für das Skalarprodukt aus Bemerkung 2.14 bzw. die Norm von $y_h = (\tau_{\text{RT}}, v_1) \in Y_h$

$$\|y_h\|_Y^2 = \|\tau_{\text{RT}}\|_{H(\text{div}, \mathcal{T})}^2 + \|v_1\|_{H_0^1(\mathcal{T})}^2$$

ergibt sich in der gegebenen Basis folgende Blockstruktur

$$\mathbb{M} = \begin{bmatrix} \mathbf{T} & \mathbf{0} \\ \mathbf{0} & \mathbf{V} \end{bmatrix}$$

mit $\mathbf{T}, \mathbf{V} \in \mathbb{R}^{6J \times 6J}$.

Die in (4.7), (4.8) und (4.9) definierten Matrizen sind erneut nützlich, um die Blöcke \mathbf{V} und \mathbf{T} genauer zu beschreiben. In (4.1) wird lediglich die Inverse der Normmatrix verwendet, daher sind die entsprechenden Codezeilen aus `computeBlocksStokesDPG` bzw. `computeBlocksStokesDPGFass` bzw. `computeBlocksStokesDPGIntegrate` erst in diesem Zusammenhang zu finden.

Block T beschreibt den Norm Anteil von $\tau_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})$. Seine Einträge ergeben sich als

$$\mathbf{T}_{mk} = \int_{\Omega} \tau_m^b : \tau_k^b \, dx + \int_{\Omega} \text{div}_{\text{NC}} \tau_m^b \cdot \text{div}_{\text{NC}} \tau_k^b \, dx \quad \text{für } m, k = 1, \dots, 6J.$$

Diese Integrale verschwinden außer τ_m^b und τ_k^b haben dasselbe Dreieck $T \in \mathcal{T}$ als Träger, d.h. $\tau_m^b = \tau^B(T, \alpha, \kappa)$ und $\tau_k^b = \tau^B(T, \varkappa, \gamma)$, für ein $T \in \mathcal{T}$, $\alpha, \gamma = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$ und in diesem Fall muss nur über dieses $T \in \mathcal{T}$ integriert werden.

Das folgende Lemma beschreibt daher die Berechnung der lokalen Blöcke.

Lemma 4.9. Für jedes Element $T \in \mathcal{T}$ sei der lokale Beitrag $\mathbf{T}(T)$ wie folgt definiert

$$\begin{aligned} \mathbf{T}(T; 2(\alpha - 1) + \kappa, 2(\gamma - 1) + \varkappa) &:= \int_T \boldsymbol{\tau}^B(T, \alpha, \kappa) : \boldsymbol{\tau}^B(T, \gamma, \varkappa) \, dx \\ &\quad + \int_T \operatorname{div} \boldsymbol{\tau}^B(T, \alpha, \kappa) \cdot \operatorname{div} \boldsymbol{\tau}^B(T, \gamma, \varkappa) \, dx \end{aligned}$$

für $\alpha, \gamma = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$. Dann gilt

$$\begin{aligned} \mathbf{T}(T; 2(\alpha - 1) + \kappa, 2(\gamma - 1) + \varkappa) &= \delta_{\kappa, \varkappa} \frac{|E_\alpha^T| |E_\gamma^T|}{4|T|} \left((\operatorname{mid}(T) - P_{\alpha-1}) \cdot (\operatorname{mid}(T) - P_{\gamma-1}) \right. \\ &\quad \left. + s(T)^2/36 + 4 \right) \end{aligned}$$

mit

$$s(T)^2 := |E_1^T|^2 + |E_2^T|^2 + |E_3^T|^2 \in \mathbb{R},$$

wobei die lokalen Seiten E_1^T, E_2^T, E_3^T und die lokalen Eckpunkte $P_1, P_2, P_3 \equiv P_0$ wie in Abbildung 4.1 angeordnet sind, d.h. so dass $P_{\alpha-1}$ gegenüber von E_α^T liegt. Damit ergibt sich folgende Blockform

$$\mathbf{T}(T) = H_1 \tilde{\mathbf{T}}(T) H_1^\top + H_2 \tilde{\mathbf{T}}(T) H_2^\top \in \mathbb{R}^{6 \times 6},$$

mit

$$\tilde{\mathbf{T}}(T) := \frac{1}{|T|} \left(1 + \frac{s(T)^2}{144} \right) \tilde{\mathbf{L}}(T) \otimes \tilde{\mathbf{L}}(T) + \frac{1}{4|T|} \operatorname{diag}(\tilde{\mathbf{L}}(T)) \mathbf{K}(T)^\top \mathbf{K}(T) \operatorname{diag}(\tilde{\mathbf{L}}(T)) \in \mathbb{R}^{3 \times 3}$$

sowie

$$\tilde{\mathbf{L}}(T) := \left(|E_1^T|, |E_2^T|, |E_3^T| \right)^\top \in \mathbb{R}^3,$$

und $\mathbf{K}(T) \in \mathbb{R}^{3 \times 2}$ wie in Lemma 4.4

$$\mathbf{K}(T; \alpha, \kappa) = (\operatorname{mid}(T) - P_{\alpha-1}) \cdot e_\kappa \quad \text{für } \alpha = 1, 2, 3 \text{ und } \kappa = 1, 2.$$

Beweis. Für die lokalen Raviart-Thomas Basisfunktionen $\psi(T, \alpha)$ und $\psi(T, \gamma)$ auf einem Element $T \in \mathcal{T}$ gilt, nach Gleichung (2.2)

$$\begin{aligned} &\int_T \psi(T, \alpha) \cdot \psi(T, \gamma) \\ &= \frac{|E_\alpha^T| |E_\gamma^T|}{4|T|^2} \int_T \left((x - \operatorname{mid}(T) + \operatorname{mid}(T) - P_{\alpha-1}) \cdot (x - \operatorname{mid}(T) + \operatorname{mid}(T) - P_{\gamma-1}) \right) dx \\ &= \frac{|E_\alpha^T| |E_\gamma^T|}{4|T|^2} \left(\|\bullet - \operatorname{mid}(T)\|_{L^2(T)}^2 + |T| (\operatorname{mid}(T) - P_{\alpha-1}) \cdot (\operatorname{mid}(T) - P_{\gamma-1}) \right). \end{aligned}$$

Da $\operatorname{div} \psi(T, \alpha) = |E_\alpha^T|/|T|$, ergibt sich

$$\int_T \operatorname{div} \psi(T, \alpha) \operatorname{div} \psi(T, \gamma) = \frac{|E_\alpha^T| |E_\gamma^T|}{|T|}.$$

Mit $s(T)^2 := |E_1^T|^2 + |E_2^T|^2 + |E_3^T|^2$ gilt nach [Car09b, Ex. 3.13] $\|\bullet - \operatorname{mid}(T)\|_{L^2(T)}^2 = s(T)^2|T|/36$. Also ergibt sich, für $\alpha, \gamma = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$

$$\begin{aligned} & \int_T \left(\boldsymbol{\tau}^B(T, \alpha, \kappa) : \boldsymbol{\tau}^B(T, \gamma, \varkappa) + \operatorname{div} \boldsymbol{\tau}^B(T, \alpha, \kappa) \cdot \operatorname{div} \boldsymbol{\tau}^B(T, \gamma, \varkappa) \right) dx \\ &= \delta_{\kappa, \varkappa} \frac{|E_\alpha^T| |E_\gamma^T|}{4|T|} \left(s(T)^2/36 + (\operatorname{mid}(T) - P_{\alpha-1}) \cdot (\operatorname{mid}(T) - P_{\gamma-1}) + 4 \right). \end{aligned}$$

Daher führt

$$\begin{aligned} \tilde{\mathbf{T}}(T) &= \frac{1}{4|T|} \left(4 + \frac{s(T)^2}{36} \right) \begin{pmatrix} |E_\alpha^T| \\ |E_\beta^T| \\ |E_\gamma^T| \end{pmatrix} \begin{pmatrix} |E_\alpha^T| & |E_\beta^T| & |E_\gamma^T| \end{pmatrix} + \frac{1}{4|T|} \begin{pmatrix} |E_\alpha^T| (\operatorname{mid}(T) - P_{\alpha-1})^\top \\ |E_\beta^T| (\operatorname{mid}(T) - P_{\beta-1})^\top \\ |E_\gamma^T| (\operatorname{mid}(T) - P_{\gamma-1})^\top \end{pmatrix} \\ &\quad \times \begin{pmatrix} |E_\alpha^T| (\operatorname{mid}(T) - P_{\alpha-1}) & |E_\beta^T| (\operatorname{mid}(T) - P_{\beta-1}) & |E_\gamma^T| (\operatorname{mid}(T) - P_{\gamma-1}) \end{pmatrix}, \end{aligned}$$

zu folgender Blockform von $\mathbf{T}(T)$

$$\mathbf{T}(T) = \begin{pmatrix} \tilde{\mathbf{T}}(T; 1, 1) & 0 & \tilde{\mathbf{T}}(T; 1, 2) & 0 & \tilde{\mathbf{T}}(T; 1, 3) & 0 \\ 0 & \tilde{\mathbf{T}}(T; 1, 1) & 0 & \tilde{\mathbf{T}}(T; 1, 2) & 0 & \tilde{\mathbf{T}}(T; 1, 3) \\ \tilde{\mathbf{T}}(T; 2, 1) & 0 & \tilde{\mathbf{T}}(T; 2, 2) & 0 & \tilde{\mathbf{T}}(T; 2, 3) & 0 \\ 0 & \tilde{\mathbf{T}}(T; 2, 1) & 0 & \tilde{\mathbf{T}}(T; 2, 2) & 0 & \tilde{\mathbf{T}}(T; 2, 3) \\ \tilde{\mathbf{T}}(T; 3, 1) & 0 & \tilde{\mathbf{T}}(T; 3, 2) & 0 & \tilde{\mathbf{T}}(T; 3, 3) & 0 \\ 0 & \tilde{\mathbf{T}}(T; 3, 1) & 0 & \tilde{\mathbf{T}}(T; 3, 2) & 0 & \tilde{\mathbf{T}}(T; 3, 3) \end{pmatrix} \in \mathbb{R}^{6 \times 6}.$$

Dies ist das behauptete Matrixprodukt. \square

Aus diesen lokalen Blöcken wird Block \mathbf{T} für alle $j = 1, \dots, J$, $\alpha, \gamma = 1, 2, 3$, und $\kappa, \varkappa = 1, 2$ via

$$\mathbf{T}_{6(j-1)+2(\alpha-1)+\kappa, 6(j-1)+2(\gamma-1)+\varkappa} = \mathbf{T}(T_j; 2(\alpha-1) + \kappa, 2(\gamma-1) + \varkappa).$$

assembliert. Damit ergibt sich auch hier eine Blockdiagonalstruktur

$$\mathbf{T} = \begin{pmatrix} \mathbf{T}(T_1) & 0 & \dots & 0 \\ 0 & \mathbf{T}(T_2) & \dots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & \mathbf{T}(T_J) \end{pmatrix} \in \mathbb{R}^{6J \times 6J}.$$

Block V beschreibt den Norm Anteil von $v_1 \in P_1(\mathcal{T}; \mathbb{R}^2)$, damit ergeben sich die Einträge durch

$$\mathbf{V}_{mk} = \int_{\Omega} v_m^b \cdot v_k^b \, dx + \int_{\Omega} \mathbf{D}_{\text{NC}} v_m^b : \mathbf{D}_{\text{NC}} v_k^b \, dx \quad \text{für } m, k = 1, \dots, 6J.$$

Diese Integrale verschwinden es sei denn v_m^b und v_k^b haben dasselbe Dreieck $T \in \mathcal{T}$ als Träger, d.h. $v_m^b = v^B(T, \gamma, \kappa)$ und $v_k^b = v^B(T, \varkappa, \alpha)$, für ein $T \in \mathcal{T}$, $\gamma, \alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$ und in diesem Fall wird lediglich über dieses $T \in \mathcal{T}$ integriert.

Lemma 4.10. Für jedes Dreieck $T \in \mathcal{T}$ sei der lokale Block $\mathbf{V}(T)$ definiert durch

$$\begin{aligned} \mathbf{V}(T; 2(\gamma - 1) + \kappa, 2(\alpha - 1) + \varkappa) &:= \int_T v^B(T, \gamma, \kappa) \cdot v^B(T, \alpha, \varkappa) \, dx \\ &+ \int_T \mathbf{D} v^B(T, \gamma, \kappa) : \mathbf{D} v^B(T, \alpha, \varkappa) \, dx \end{aligned}$$

für $\gamma, \alpha = 1, 2, 3$ und $\kappa, \varkappa = 1, 2$. Es gilt

$$\mathbf{V}(T; 2(\gamma - 1) + \kappa, 2(\alpha - 1) + \varkappa) = |T| \delta_{\kappa, \varkappa} \left(\nabla \varphi_{\alpha} \cdot \nabla \varphi_{\gamma} + 1/12 (1 + \delta_{\gamma, \alpha}) \right)$$

und es ergibt sich damit folgende Struktur

$$\mathbf{V}(T) = \left(H_1 \tilde{\mathbf{V}}(T) H_1^{\top} + H_2 \tilde{\mathbf{V}}(T) H_2^{\top} \right) \in \mathbb{R}^{6 \times 6},$$

wobei

$$\tilde{\mathbf{V}}(T) := |T| \mathbf{G}(T) \mathbf{G}(T)^{\top} + \frac{|T|}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} \in \mathbb{R}^{3 \times 3}$$

und $\mathbf{G}(T)$ sei wie in Lemma 4.7 definiert durch

$$\mathbf{G}(T) := \begin{pmatrix} 1 & 1 & 1 \\ P_1 & P_2 & P_3 \end{pmatrix}^{-1} \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{3 \times 2},$$

für die lokale Nummerierung P_1, P_2, P_3 der Eckpunkte von T .

Beweis. Es gilt

$$\begin{aligned} \int_T (\varphi(T, \gamma) \cdot \varphi(T, \alpha) + \nabla \varphi(T, \gamma) : \nabla \varphi(T, \alpha)) \, dx &= |T| / 12 (1 + \delta_{\gamma, \alpha}) + |T| \nabla \varphi_{\alpha} \cdot \nabla \varphi_{\gamma} \\ &= \tilde{B}(T; \gamma, \alpha) + \tilde{M}(T; \gamma, \alpha), \end{aligned}$$

wobei $\tilde{B}(T) := |T| \mathbf{G}(T) \otimes \mathbf{G}(T)$ die lokale Steifigkeitsmatrix und $\tilde{M}(T) := \frac{|T|}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}$ die lokale Massematrix für die nodalen Basisfunktionen im Poisson Modell Problem sind.

Dies wird in [ACF99, S.122/127] hergeleitet. Da $\tilde{V}(T) = \tilde{B}(T; \gamma, \alpha) + \tilde{M}(T; \gamma, \alpha)$, hat $\mathbf{V}(T)$ auf Grund der Basisanordnung die folgende Form

$$\mathbf{V}(T) = \begin{pmatrix} \tilde{\mathbf{V}}(T; 1, 1) & 0 & \tilde{\mathbf{V}}(T; 1, 2) & 0 & \tilde{\mathbf{V}}(T; 1, 3) & 0 \\ 0 & \tilde{\mathbf{V}}(T; 1, 1) & 0 & \tilde{\mathbf{V}}(T; 1, 2) & 0 & \tilde{\mathbf{V}}(T; 1, 3) \\ \tilde{\mathbf{V}}(T; 2, 1) & 0 & \tilde{\mathbf{V}}(T; 2, 2) & 0 & \tilde{\mathbf{V}}(T; 2, 3) & 0 \\ 0 & \tilde{\mathbf{V}}(T; 2, 1) & 0 & \tilde{\mathbf{V}}(T; 2, 2) & 0 & \tilde{\mathbf{V}}(T; 2, 3) \\ \tilde{\mathbf{V}}(T; 3, 1) & 0 & \tilde{\mathbf{V}}(T; 3, 2) & 0 & \tilde{\mathbf{V}}(T; 3, 3) & 0 \\ 0 & \tilde{\mathbf{V}}(T; 3, 1) & 0 & \tilde{\mathbf{V}}(T; 3, 2) & 0 & \tilde{\mathbf{V}}(T; 3, 3) \end{pmatrix} \in \mathbb{R}^{6 \times 6},$$

wird also durch das behauptete Matrixprodukt beschrieben. \square

Aus diesen lokalen Blöcken wird der Matrixblock \mathbf{V} für alle $j = 1, \dots, J$, $\gamma, \alpha = 1, 2, 3$, und $\kappa, \varkappa = 1, 2$ via

$$\mathbf{V}_{6(j-1)+2(\gamma-1)+\kappa, 6(j-1)+2(\alpha-1)+\varkappa} = \mathbf{V}(T_j; 2(\gamma-1) + \kappa, 2(\alpha-1) + \varkappa)$$

assembliert. Erneut ergibt sich also eine Blockdiagonalform

$$\mathbf{V} = \begin{pmatrix} \mathbf{V}(T_1) & 0 & \dots & 0 \\ 0 & \mathbf{V}(T_2) & \dots & 0 \\ \vdots & 0 & \ddots & \vdots \\ 0 & \dots & 0 & \mathbf{V}(T_J) \end{pmatrix} \in \mathbb{R}^{6J \times 6J}.$$

Um Blockdiagonalmatrizen zu invertieren müssen nur die einzelnen Blöcke invertiert werden. D.h. die Blockdiagonalform von \mathbf{T} und \mathbf{V} erleichtert die Berechnung von \mathbb{M}^{-1} , es gilt nämlich

$$\begin{aligned} \mathbb{M}^{-1} &= \begin{pmatrix} \mathbf{T}^{-1} & 0 \\ 0 & \mathbf{V}^{-1} \end{pmatrix} \text{ mit} \\ \mathbf{T}^{-1} &= \text{diag}(\mathbf{T}(T_1)^{-1}, \mathbf{T}(T_2)^{-1}, \dots, \mathbf{T}(T_J)^{-1}) \in \mathbb{R}^{6J} \\ \text{und } \mathbf{V}^{-1} &= \text{diag}(\mathbf{V}(T_1)^{-1}, \mathbf{V}(T_2)^{-1}, \dots, \mathbf{V}(T_J)^{-1}) \in \mathbb{R}^{6J}. \end{aligned}$$

Da für $A, B, C \in \mathbb{R}^{3 \times 3}$ und die Hilfsmatrizen aus (4.9) gilt, wenn $AB = C$, dann

$$(H_1 A H_1^\top + H_2 A H_2^\top) (H_1 B H_1^\top + H_2 B H_2^\top) = (H_1 C H_1^\top + H_2 C H_2^\top),$$

genügt es statt der $12J \times 12J$ -Matrix \mathbb{M} für alle $j = 1, \dots, J$ die 3×3 -Blöcke $\tilde{\mathbf{T}}(T_j)$ und $\tilde{\mathbf{V}}(T_j)$ numerisch zu invertieren. Die Inverse Matrix wird dann mit Hilfe von $\mathbf{T}(T_j)^{-1} = H_1 \tilde{\mathbf{T}}(T_j)^{-1} H_1^\top + H_2 \tilde{\mathbf{T}}(T_j)^{-1} H_2^\top$ und $\mathbf{V}(T_j)^{-1} = H_1 \tilde{\mathbf{V}}(T_j)^{-1} H_1^\top + H_2 \tilde{\mathbf{V}}(T_j)^{-1} H_2^\top$ bestimmt.

Die lokalen Blöcke werden in dem Programm wie folgt assembliert

```
% initialisation of matrix blocks
Tinv4e = zeros(6,6,nrElems); % local matrices for block T^{-1}
Vinv4e = zeros(6,6,nrElems); % local matrices for block V^{-1}
% constant auxillary matrices
H1 = sparse([1,3,5],[1,2,3],[1,1,1],6,3);
H2 = sparse([2,4,6],[1,2,3],[1,1,1],6,3);

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    coords = coord4e(:, :, elem); % coordinates of three nodes
    mid = mid4e(elem, :); % midpoint of T
    area = area4e(elem); % area of T
    slength = sl4e(elem, :)' ; % length of sides
    G = [1, 1, 1; coords'] \ [0, 0; 1, 0; 0, 1];
    K = [mid; mid; mid] - [coords(3, :); coords(1, :); coords(2, :)];
    ssq = sum(slength.^2); % s^2

    % block T ||tau||_H(div)
    temp = 1/(4*area)*(diag(slength)*K*K'*diag(slength))...
        +1/area*(1+ssq/144)*(slength*slength');
    % inverse of T
    temp = temp \ eye(3);
    Tinv4e(:, :, elem) = H1*temp*H1' + H2*temp*H2';

    % block V ||v||_H^1
    temp = area*G*G' + area/12*(ones(3,3)+eye(3));
    % inverse of V
    temp = temp \ eye(3);
    Vinv4e(:, :, elem) = H1*temp*H1' + H2*temp*H2';
end
```

Daraus ergeben sich die globalen Blöcke via

```
%% Constructing global matrix blocks
% with tempI as vector of row indices and tempJ vector of column indices
tempI = permute(reshape(repmat(1:6*nrElems,6,1),6,6,nrElems),[2,1,3]);
tempJ = repmat(1:6*nrElems,6,1);
Tinv=sparse(tempI(:),tempJ(:),Tinv4e(:)); % block T^{-1}
Vinv=sparse(tempI(:),tempJ(:),Vinv4e(:)); % block V^{-1}
```

Mit diesen Bestandteilen wird in dem Programm \mathbb{M}^{-1} wie folgt assembliert

```
% norm matrix MMinv=sparse(dimY,dimY)
MMinv = [Tinv,sparse(6*nrElems,6*nrElems);...
        sparse(6*nrElems,6*nrElems),Vinv];
```

Im Folgenden werden nun noch die Vektoren \mathbb{F} der rechten Seite und \mathbb{L} der linearen Nebenbedingung bestimmt.

4.3.3 Vektor zum Funktional F und zur linearen Nebenbedingung Λ

Zunächst wird das Funktional F aus Definition 3.2 betrachtet. Zu gegebenen Daten $f \in L^2(\Omega; \mathbb{R}^n)$ und $\hat{s} \in H^{1/2}(\partial\mathcal{T}; \mathbb{R}^n)$ mit $\hat{s} = g$ auf Γ ist für alle $y = (\boldsymbol{\tau}, v) \in Y$

$$F(y) = \int_{\Omega} f \cdot v \, dx + \left\langle \gamma_{\nu}^{\mathcal{T}} \boldsymbol{\tau}, \hat{s} \right\rangle_{\partial\mathcal{T}}.$$

Für $Y_h = \text{span} \{ \zeta_1, \dots, \zeta_M \}$ ist

$$\mathbb{F} := (F(\zeta_m))_{1 \leq m \leq M} \in \mathbb{R}^M.$$

Die Basisanordnung führt zu folgender Blockstruktur

$$\mathbb{F} = \begin{pmatrix} \mathbf{F} \\ \mathbf{G} \end{pmatrix} \in \mathbb{R}^{12J}$$

mit den Blöcken $\mathbf{F}, \mathbf{G} \in \mathbb{R}^{6J}$.

Block F enthält die Beiträge, die sich durch die Testfunktionen $\boldsymbol{\tau}_{\text{RT}} \in RT_0^{\text{pw}}(\mathcal{T}; \mathbb{R}^{2 \times 2})$ ergeben. Daher sei

$$\mathbf{F}_m = \left\langle \gamma_\nu^\mathcal{T} \boldsymbol{\tau}_m^b, \hat{s} \right\rangle_{\partial\mathcal{T}} \quad \text{für } m = 1, \dots, 6J.$$

Da die Testfunktionen als Träger immer nur ein Element $T \in \mathcal{T}$ haben, ist dieses Integral stets nur auf den Kanten eines Dreiecks auszuwerten. Vor der Berechnung der sich ergebenden lokalen Blöcke, gilt es zunächst eine Funktion \hat{s} zu wählen, die den Bedingungen aus Definition 3.2 genügt. Für die reguläre Triangulierung \mathcal{T} von Ω definiere \hat{s} so, dass $\hat{s}|_E = g|_E$ für alle Randkanten $E \in \mathcal{E}(\partial\Omega)$, also $\hat{s}(z) = g(z)$ für alle Randknoten $z \in \mathcal{N}(\partial\Omega)$, des Weiteren sei $\hat{s}(z) = 0$ für alle Innenknoten $z \in \mathcal{N}(\Omega)$ und $\hat{s}|_E \in P_1(E)$ für alle inneren Kanten $E \in \mathcal{E}(\Omega)$. In diesem Setting berechnet sich Block \mathbf{F} wie folgt.

Lemma 4.11. *Für jedes Dreieck $T \in \mathcal{T}$ sei der lokale Block $\mathbf{F}(T)$ definiert als*

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) := \int_{\partial T} \hat{s} \cdot \boldsymbol{\tau}^B(T, \alpha, \kappa) \nu_T \, ds$$

Die Einträge genügen folgender Fallunterscheidung

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) = \begin{cases} \int_{E_\alpha^T} g \cdot e_\kappa \, ds & \text{für } E_\alpha^T \in \mathcal{E}(\partial\Omega) \\ (\hat{s}(P_\alpha) + \hat{s}(P_{\alpha+1})) \cdot e_\kappa |E_\alpha^T|/2 & \text{für } E_\alpha^T \in \mathcal{E}(\Omega). \end{cases}$$

Damit ist das Aussehen dieses lokalen Blocks abhängig von der Art der Kanten.

Beweis. Für die lokalen Raviart-Thomas Funktionen gilt wie bereits in Lemma 4.6 erklärt,

$$(\psi(T, \alpha) \cdot \nu_T)|_F = |E_\alpha^T| / (2|T|) ((x - P_{\alpha-1}) \cdot \nu_T|_F) = \begin{cases} 1 & \text{für } F = E_\alpha^T, \\ 0 & \text{für } F \neq E_\alpha^T, \end{cases}$$

d.h. es gilt für $\alpha = 1, 2, 3$ und $\kappa = 1, 2$

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) = \int_{E_\alpha^T} \hat{s} \cdot e_\kappa \, ds.$$

Es ist also lediglich \hat{s} zu integrieren. Für die Randkanten, $E_\alpha^T \in \mathcal{E}(\partial\Omega)$, gilt nach Definition von \hat{s}

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) = \int_{E_\alpha^T} g \cdot e_\kappa \, ds.$$

Für die inneren Kanten, $E_\alpha^T \in \mathcal{E}(\Omega)$ ist $\hat{s}|_{E_\alpha^T} \in P_1(E_\alpha^T)$. Da nach der Nummerierung aus Abbildung 4.1 $E_\alpha^T = \text{conv}\{P_\alpha, P_{\alpha+1}\}$ gilt, folgt also

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) = (\hat{s}(P_\alpha) + \hat{s}(P_{\alpha+1})) \cdot e_\kappa |E_\alpha^T|/2.$$

Für die Auswertung entlang innerer Kanten ist eine Unterscheidung nach Eckpunkten zu treffen

$$\mathbf{F}(T; 2(\alpha - 1) + \kappa) = \begin{cases} 0 & \text{für } P_\alpha, P_{\alpha+1} \in \mathcal{N}(\Omega), \\ g(P_{\alpha+1}) \cdot e_\kappa |E_\alpha^T|/2 & \text{für } P_\alpha \in \mathcal{N}(\Omega), P_{\alpha+1} \in \mathcal{N}(\partial\Omega), \\ g(P_\alpha) \cdot e_\kappa |E_\alpha^T|/2 & \text{für } P_\alpha \in \mathcal{N}(\partial\Omega), P_{\alpha+1} \in \mathcal{N}(\Omega), \\ (g(P_\alpha) + g(P_{\alpha+1})) \cdot e_\kappa |E_\alpha^T|/2 & \text{für } P_\alpha, P_{\alpha+1} \in \mathcal{N}(\partial\Omega). \end{cases}$$

□

Aus den lokalen Vektoren wird der Matrixblock \mathbf{F} , für alle $j = 1, \dots, J$, $\alpha = 1, 2, 3$ und $\kappa = 1, 2$ wie folgt assembliert

$$\mathbf{F}_{6(j-1)+2(\alpha-1)+\kappa} = \mathbf{F}(T_j; 2(\alpha - 1) + \kappa).$$

Also werden die einzelnen Blöcke einfach aneinander gehängt

$$\mathbf{F} = \begin{pmatrix} \mathbf{F}(T_1) \\ \mathbf{F}(T_2) \\ \vdots \\ \mathbf{F}(T_J) \end{pmatrix} \in \mathbb{R}^{6J}.$$

In dem Code liest sich das wie folgt

```
% Constructing global matrix blocks
F=sparse(1:6*nrElems, ones(6*nrElems,1), F4e(:), 6*nrElems, 1); % block F
```

Um das numerische Integrieren über die Randkanten zu vermeiden, kann g auf allen Kanten linear approximiert werden, d.h. $\int_E g \cdot e_\kappa \approx (g(z_k) + g(z_\ell)) \cdot e_\kappa |E|/2$ für $E = \text{conv}\{z_k, z_\ell\} \in \mathcal{E}$. In diesem Fall ist \mathbf{F} leicht zu bestimmen. Dazu definiere den Vektor $sh \in \mathbb{R}^{|\mathcal{N}|\times 2}$ durch

$$sh(k) = \begin{cases} g(z_k)^\top & \text{für } z_k \in \mathcal{N}(\partial\Omega), \\ 0 & \text{für } z_k \in \mathcal{N}(\Omega). \end{cases}$$

Für ein Dreieck $T \in \mathcal{T}$, $\alpha = 1, 2, 3$ die lokalen Knotennummern und $\kappa = 1, 2$ seien die lokalen Werte von \hat{s} in den Knoten durch

$$sh(T; 2(\alpha - 1) + \kappa) := sh(k) \cdot e_\kappa,$$

wobei k die globale Nummer des Knoten P_α sei. Des Weiteren sei gegeben

$$Sh(T) := sh(T) + \begin{pmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{pmatrix} sh(T) \in \mathbb{R}^6.$$

Dann ergibt sich vermöge $L(T)$ aus (4.7) folgende approximierte Version von Block \mathbf{F}

$$\mathbf{F}(T) = 1/2 L(T) Sh(T).$$

Diese lokalen Blöcke werden in `computeBlocksStokesDPG` wie folgt assembliert

```
% initialisation of matrix blocks
F4e = zeros(6,nrElems); % local matrices for block F

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    slength = sl4e(elem,:); % length of sides
    sh = hats4e(:, :, elem)'; sh = sh(:); % values of hat(s) in nodes
    L = [slength, slength]'; L=diag(L(:));
    Sh = sh + (diag(ones(4,1),2)+diag(ones(2,1),-4))*sh;

    %% Blocks RHS (right hand side)
    % block F int_dT(gamma_nu tau.hat s)
    F4e(:, elem) = 0.5*L*Sh;
end
```

Diese lineare Approximation erweist sich jedoch nicht immer als praktikabel. Bessere Ergebnisse sind zu erzielen, wenn die Integrale über die Randkanten beispielsweise mit der Fassregel, $\int_E g \, ds \approx |E|/6 (g(z_k) + g(z_\ell) + 4g((z_k + z_\ell)/2))$ für $E = \text{conv}\{z_k, z_\ell\} \in \mathcal{E}(\partial\Omega)$, approximiert werden. In `computeBlocksStokesDPGFass` werden daher zunächst die Nummern der Seiten auf dem Dirichlet Rand bestimmt und anschließend die lokalen Blöcke wie folgt assembliert

```
%sides on Dirichlet boundary
DbSides=zeros(size(n4sDb,1),1);
for i=1:size(n4sDb,1)
    DbSides(i)=s4n(n4sDb(i,1), n4sDb(i,2));
end
% initialisation of matrix blocks
F4e = zeros(6,nrElems); % local matrices for block F

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    slength = sl4e(elem,:); % length of sides
    sh = hats4e(:, :, elem)'; sh = sh(:); % values of hat(s) in nodes
    L = [slength, slength]'; L=diag(L(:));
    Sh = sh + (diag(ones(4,1),2)+diag(ones(2,1),-4))*sh;

    % block F int_dT(gamma_nu tau.hat s)
    F4e(:, elem) = 0.5*L*Sh;

    %treatment of boundary sides
    Db=ismember(sides,DbSides); % boolean, if side is on boundary
    if max(Db)~=0
```

```

        Db= repmat(Db,2,1); Db=Db(:);
        %integration of along Dirichlet boundary sides applying Fass
        %rule
        smid = mids4e(:, :, elem); % midpoints of sides
        temp = hats4e(:, :, elem)+4.*u4Db(c.*smid)+hats4e([2,3,1], :, elem);
        temp = temp'; temp = temp(:);

        F4e(:, elem)=(ones(6,1)-Db).*F4e(:, elem)+Db.*diag(L)./6.*temp;
    end
end

```

Dabei geht der Skalierungsparameter c ein wie in Abschnitt 4.5 erläutert.

In `computeBlocksStokesDPGIntegrate` wird lediglich vor der Schleife mit Hilfe der AFEM-Funktion `integrate` ein Vektor mit den Integralen von g entlang der Randseiten erzeugt und die `if`-Anweisung ein wenig verändert

```

% sides on Dirichlet boundary
DbSides=zeros(size(n4sDb,1),1);
for i=1:size(n4sDb,1)
    DbSides(i)=s4n(n4sDb(i,1), n4sDb(i,2));
end

% numerical integration of g along Dirichlet boundary
int4DbSides = integrate(c4n,n4sDb,@(n4p,Gpts4p,Gpts4ref)(u4Db(c.*Gpts4p)),12);
int4sides = zeros(nrSides,2);
int4sides(DbSides,:)=int4DbSides;
% slicing for parfor loop
ints4e = permute(reshape(int4sides(s4e',:))',2,3,nrElems),[2,1,3]);
% integral of g along boundary sides

% initialisation of matrix blocks
F4e = zeros(6,nrElems) % local matrices for block F

%% Parfor-loop to fill local matrices
for elem = 1:nrElems %elem as current triangle T
    slength = sl4e(elem,:); % length of sides
    sh = hats4e(:, :, elem)'; sh = sh(:); % values of hat(s) in nodes
    L = [slength,slength]'; L=diag(L(:));
    Sh = sh+ (diag(ones(4,1),2)+diag(ones(2,1),-4))*sh;

    % block F int_dT(gamma_nu tau.hat s)
    F4e(:, elem)= 0.5*L*Sh;

    %treatment of boundary sides
    Db=ismember(sides,DbSides); % boolean, if side is on boundary
    if max(Db)~=0
        Db = repmat(Db,2,1); Db=Db(:);
        %integration of along Dirichlet boundary sides insert
        %numerical integrals of g
        sint = ints4e(:, :, elem); % integral along sides
        sint = sint'; sint = sint(:);
        F4e(:, elem) = (ones(6,1)-Db).*F4e(:, elem)+Db.*sint;
    end
end

```

Block G enthält die Beiträge, die sich durch die Testfunktionen $v_1 \in P_1^{\text{pw}}(\mathcal{T}; \mathbb{R}^2)$ ergeben. Seine Einträge lauten

$$\mathbf{G}_m = \int_{\mathcal{T}} f \cdot v_m^b \, dx \text{ für } m = 1, \dots, 6J.$$

Auch hier wird jeweils nur über das zu v_m^b assoziierte $T \in \mathcal{T}$ integriert.

Lemma 4.12. *Der lokale Block $\mathbf{G}(T)$ sei definiert als*

$$\mathbf{G}(T, 2(\gamma - 1) + \kappa) := \int_T f \cdot v^B(T, \gamma, \kappa) \, dx$$

für $\gamma = 1, 2, 3$ und $\kappa = 1, 2$, ergibt sich

$$\mathbf{G}(T, 2(\gamma - 1) + \kappa) = \int_T f \cdot \varphi(T; \gamma) \, e_\kappa \, dx.$$

Beweis. Dies ist lediglich die Definition von $v^B(T, \gamma, \kappa)$. □

Für alle $j = 1, \dots, J$, $\gamma = 1, 2, 3$, und $\kappa = 1, 2$ wird der Block \mathbf{G} aus diesen lokalen Blöcken via

$$\mathbf{G}_{6(j-1)+2(\gamma-1)+\kappa} = \mathbf{G}(T_j; 2(\gamma - 1) + \kappa).$$

assembliert. Also werden die einzelnen Blöcke einfach aneinander gehängt

$$\mathbf{G} = \begin{pmatrix} \mathbf{G}(T_1) \\ \mathbf{G}(T_2) \\ \vdots \\ \mathbf{G}(T_J) \end{pmatrix} \in \mathbb{R}^{6J}.$$

Diese Integrale werden wie üblich nicht durch numerische Integration sondern mit der Mittelpunkts Regel, $\int_T g \, dx \approx |T|g(\text{mid}(T))$ für $g \in L^2(\Omega)$, approximiert. Es gilt $\varphi_k(\text{mid}(T)) = 1/3$, falls $z_k \in \mathcal{N}(T)$, und Null sonst. Daraus ergibt sich die folgende Form des approximierten Blocks

$$\mathbf{G}(T) = \frac{|T|}{3} \begin{pmatrix} f(\text{mid}(T)) \\ f(\text{mid}(T)) \\ f(\text{mid}(T)) \end{pmatrix} \in \mathbb{R}^{6 \times 1}.$$

Dieser Block kann global in folgenden Zeilen bestimmt werden

```
nrElems = size(n4e,1);           % number of elements
mid4e    = computeMid4e(c4n,n4e); % midpoints of elements
area4e   = computeArea4e(c4n,n4e); % area of elements

%% Constructing global matrix blocks
temp = repmat(c^2.*[area4e,area4e].*f(c*mid4e)/3,1,3)';
G=sparse(1:6*nrElems,ones(6*nrElems,1),temp(:),6*nrElems,1); % block G
```

Dabei geht der Skalierungsparameter c ein, der in Abschnitt 4.5 genauer erläutert wird. Der gesamte Vektor \mathbb{F} ergibt sich in allen drei Fällen als

```
% vector of right hand side FF=sparse(dimY,1)
FF = [F;G];
```

Es bleibt lediglich die lineare Nebenbedingung zu untersuchen. Das lineare Funktional $\Lambda : X_h \rightarrow \mathbb{R}$, das diese Nebenbedingung $\int_\Omega \text{tr } \boldsymbol{\sigma}_0 \, dx = 0$ für alle $x_h = (\boldsymbol{\sigma}_0, u_0, s_1, t_0) \in X_h$

durch $\Lambda(x_h) = 0$ beschreibt, lautet

$$\Lambda(x_h) = \int_{\Omega} \operatorname{tr} \boldsymbol{\sigma}_0 \, dx \quad \text{für alle } x_h = (\boldsymbol{\sigma}_0, u_0, s_1, t_0) \in X_h.$$

Damit kann der zugehörige Vektor $\mathbb{L} \in \mathbb{R}^N$ mit

$$\mathbb{L}_n := \Lambda(\xi_n) \quad \text{für } n = 1, \dots, N,$$

wie im folgenden Lemma berechnet werden.

Lemma 4.13. *Der Vektor \mathbb{L} hat folgende Einträge*

$$\mathbb{L}_n = \begin{cases} \int_{T_j} \operatorname{tr} \left(\boldsymbol{\sigma}_{4(j-1)+\lambda}^b \right) \, dx & \text{für } 1 \leq n \leq 4J, n = 4(j-1) + \lambda \text{ mit } \lambda = 1, \dots, 4, \\ 0 & \text{sonst.} \end{cases}$$

Außerdem gilt für alle $T \in \mathcal{T}$ und $\lambda = 1, \dots, 4$

$$\int_T \operatorname{tr} \left(\boldsymbol{\sigma}^B(T, \lambda) \right) = \begin{cases} 2|T| & \text{für } \lambda = 1, \\ 0 & \text{sonst.} \end{cases}$$

Also

$$\mathbb{L} = \begin{pmatrix} 2|T_1| & 0 & 0 & 0 & 2|T_2| & 0 & 0 & 0 & 2|T_3| & \dots & 2|T_J| & 0 & 0 & 0 & 0 & \dots & 0 \end{pmatrix}.$$

Beweis. Diese Darstellung geht direkt aus der Definition hervor, da $\operatorname{tr}(\mathbf{e}_1) = 2$ und $\operatorname{tr}(\mathbf{e}_\lambda) = 0$ für $\lambda = 2, 3, 4$ gilt. \square

Die entsprechenden Codezeilen zur Assemblierung lauten in `computeBlocksStokesDPG` bzw. `computeBlocksStokesDPGFass` bzw. `computeBlocksStokesDPGIntegrate`

```
nrNodes = size(c4n,1);           % number of nodes
nrElems = size(n4e,1);           % number of elements
nrSides = size(n4s,1);           % number of sides
area4e   = computeArea4e(c4n,n4e); % area of elements
dimX      = 6*nrElems+2*nrNodes+2*nrSides; % N, dimension of X_h

% vector for linear side condition LL=sparse(dimX,1)
LL       = sparse(1:4:4*nrElems, ones(nrElems,1), 2*area4e, dimX, 1);
```

4.4 Exakter Fehler

Für den Fall, dass eine exakte Lösung $\boldsymbol{\sigma} \in H(\operatorname{div}, \Omega; \mathbb{R}^{2 \times 2}) / \mathbb{R}$, $u \in H^1(\Omega; \mathbb{R}^2)$ zu Problem 8 bekannt ist, soll der exakte Fehler die Abweichung von $(\boldsymbol{\sigma}, u, \gamma_0^\mathcal{T} u, \gamma_\nu^\mathcal{T} \boldsymbol{\sigma})$ von der assoziierten numerischen Lösung (vgl. Satz 3.3) in der Norm von X angeben.

Hier ist zu beachten, dass die wie in Abschnitt 4.1 bestimmte Lösung $(\boldsymbol{\sigma}_0, u_0, s_1, t_0) \in X_h$, eine Lösung zu Problem 10 ist. Die Lösungskomponente s_1 besitzt also Nullrandbedingungen. In Satz 3.3 ist für den kontinuierlichen Fall zu erkennen, dass $\gamma_0^T u = s + \hat{s}$ gilt, wobei $\hat{s} \in H^{1/2}(\partial\mathcal{T}; \mathbb{R}^2)$ die in der Definition der rechten Seite verwendete Funktion mit $\hat{s} = g$ auf Γ ist und $s \in H_0^{1/2}(\partial\mathcal{T}; \mathbb{R}^2)$ die Lösungskomponente. Um wie gewünscht die nodale P_1 -Approximation \tilde{s}_1 von $\gamma_0^T u$ zu erhalten, werden lediglich in den Randknoten die Werte von g statt den Nullwerten von s_1 gesetzt, d.h.

$$\tilde{s}_1(z_k) = \begin{cases} s_1(z_k) & \text{für } z_k \in \mathcal{N}(\Omega), \\ g(z_k) & \text{für } z_k \in \mathcal{N}(\partial\Omega). \end{cases}$$

Also gilt, wenn $x \in \mathbb{R}^N$ der berechnete Koeffizientenvektor ist,

$$\tilde{s}_1 = \sum_{n=1}^{2|\mathcal{N}(\Omega)|} x_{6J+n} s_n^b + \sum_{\substack{k \in \mathcal{N}(\partial\Omega), \\ \kappa=1, 2}} (g(z_k) \cdot e_\kappa) \gamma_0^T \varphi_k e_\kappa.$$

Des Weiteren werden zur Berechnung des exakten Fehlers die Fortsetzungen von \tilde{s}_1 , also $w_c \in S_1(\Omega; \mathbb{R}^2)$ mit $\gamma_0^T w_c = \tilde{s}_1$, und von t_0 , also $\mathbf{q}_{\text{RT}} \in RT_0(\Omega; \mathbb{R}^{2 \times 2})$ mit $\gamma_\nu^T \mathbf{q}_{\text{RT}} = t_0$, betrachtet, schließlich sind die Normen der Spurräume über solche Fortsetzungen definiert und in Lemma 3.9 bzw. Lemma 3.10 wird die jeweilige Normäquivalenz bewiesen. Dazu sei für $k = 1, \dots, |\mathcal{N}(\Omega)|$ und $\kappa = 1, 2$

$$w_{2(k-1)+\kappa}^b := \varphi_k e_\kappa$$

die Fortsetzung von $s_{2(k-1)+\kappa}^b$, wobei φ_k die nodale Basisfunktion zum Knoten k sei. Die Basis für \mathbf{q}_{RT} wird mit Hilfe der globalen kantenbasierten Raviart-Thomas Funktionen bestimmt, diese sind für alle $E_\ell \in \mathcal{E}$ wie folgt definiert

$$\psi_\ell = \begin{cases} \text{sgn}(T, \alpha) \psi(T, \alpha) & \text{falls } E_\ell \in \mathcal{E}(T) \text{ und } E_\alpha^T = E_\ell, \\ 0 & \text{sonst,} \end{cases}$$

wobei wie in Abschnitt 4.2 $\psi(T, \alpha)$ die lokale Raviart-Thomas Funktion und $\text{sgn}(T, \alpha) = \nu_{E_\alpha^T} \cdot \nu_T \in \{\pm 1\}$. Es sei für $\ell = 1, \dots, |\mathcal{E}|$ und $\kappa = 1, 2$

$$q_{2(\ell-1)+\kappa}^b := \psi_\ell \otimes e_\kappa$$

die Fortsetzung für $t_{2(\ell-1)+\kappa}^b$.

Sind also die exakte Spannung $\boldsymbol{\sigma} \in H(\text{div}, \Omega; \mathbb{R}^{2 \times 2})/\mathbb{R}$ und die exakte Verschiebung $u \in H^1(\Omega; \mathbb{R}^2)$ bekannt, wird der exakte Fehler wie folgt bestimmt

$$E^2 := \|\boldsymbol{\sigma} - \boldsymbol{\sigma}_0\|_{L^2(\Omega)}^2 + \|u - u_0\|_{L^2(\Omega)}^2 + \|u - w_c\|_{H^1(\Omega)}^2 + \|\boldsymbol{\sigma} - \mathbf{q}_{\text{RT}}\|_{H(\text{div}, \Omega)}^2,$$

wobei aus dem Koeffizientenvektor $x \in \mathbb{R}^N$, der Lösung zu (4.1), die Komponenten wie folgt berechnet werden

$$\begin{aligned}\sigma_0 &:= \sum_{n=1}^{4J} x_n \sigma_n^b, \quad u_0 := \sum_{n=1}^{2J} x_{4J+n} u_n^b, \quad \mathbf{q}_{\text{RT}} := \sum_{n=1}^{2|\mathcal{E}|} x_{6J+2|\mathcal{N}(\Omega)|+n} q_n^b, \\ w_c &:= \sum_{n=1}^{2|\mathcal{N}(\Omega)|} x_{6J+n} w_n^b + \sum_{\substack{k \in \mathcal{N}(\partial\Omega), \\ \kappa=1,2}} (g(z_k) \cdot e_\kappa) \varphi_k e_\kappa.\end{aligned}$$

In der matlab-Methode `SolveStokesDPG` werden der Koeffizientenvektor `x` der Lösung zu (4.1) und der Koeffizientenvektor `xh` zu der Approximation von $(\sigma, u, \gamma_0^\top u, \gamma_\nu^\top \sigma)$ bestimmt

```
nrNodes = size(c4n,1); % number of nodes
nrElems = size(n4e,1); % number of elements
nrSides = size(e4s,1); % number of sides
dimX = 6*nrElems+2*nrNodes+2*nrSides; % dimension of X_h
% (surpressing boundary conditions)
% computation of coefficient matrix, norm matrix, vektor for F and linear
% side condition
if Int==0
    [B,Minv, F, D] = computeBlocksStokesDPG(f,u4Db,c4n,n4e,n4sDb,c);
elseif Int ==1
    [B,Minv, F, D] = computeBlocksStokesDPGFass(f,u4Db,c4n,n4e,n4sDb,c);
else
    [B,Minv, F, D] = computeBlocksStokesDPGIntegrate(f,u4Db,c4n,n4e,n4sDb,c);
end
%% Linear system Ax=b corresponding to discrete problem
x = zeros(dimX+1,1);

% compute free components and components fixed by boundary condition for s
DbNodes = unique(n4sDb); % Dirichlet boundary nodes
fixed = [2*(DbNodes'-1)+1;2*(DbNodes'-1)+2];
fixed = 6*nrElems+fixed(:);
free = setdiff(1:dimX,fixed);
% define matrix A
A=[B(:,free)'*Minv*B(:,free),D(free);D(free)',0];
% define vector b
b=[B(:,free)'*Minv*F;0];
% compute coefficient vector in free components
x([free,dimX+1]) = A\b;
% remove component from Lagrange multiplicator for linear side condition
x=x(1:dimX);

%% Solution corresponding to original problem
% computing boundary data
DbCoords = c4n(DbNodes,:); % coordinates for DbNodes
values4DbNodes = u4Db(c.*DbCoords)'; % values in DbNodes

% computing coefficient vector corresponding to original problem
xh = x;
%boundary values
xh(fixed) = values4DbNodes(:);
```

Der Methode `computeErrorStokesDPG` zur Fehlerberechnung wird nur der Koeffizientenvektor `xh` übergeben.

4.5 Skalierung

In dem `matlab`-Code ist eine Skalierung des Gebietes Ω mit dem Faktor $c := \max_{x \in \Omega} |x|$ möglich, damit wird garantiert, dass $\hat{\Omega} := 1/c \Omega \subseteq [-1, 1]^2$. Denn in die Norm des Testraums geht die Gebietsgröße in beiden Komponenten mit unterschiedlicher Gewichtung ein, dies kann dazu führen, dass bei der Minimierung eine der Komponenten bevorzugt wird.

Diese Skalierung wird vermöge der Transformation

$$\Phi = c \text{id} : \hat{\Omega} \rightarrow \Omega$$

vorgenommen. Dazu wird mit $\hat{\Gamma} := \Phi^{-1}(\Gamma)$ der Rand von $\hat{\Omega}$ bezeichnet und es werden die Variablen

$$\hat{u} := u \circ \Phi : \hat{\Omega} \rightarrow \mathbb{R}^2 \quad \text{und} \quad \hat{\sigma} := \sigma \circ \Phi : \hat{\Omega} \rightarrow \mathbb{R}^{2 \times 2}$$

sowie die Daten

$$\hat{f} := f \circ \Phi : \hat{\Omega} \rightarrow \mathbb{R}^2 \quad \text{und} \quad \hat{g} := g \circ \Phi : \hat{\Gamma} \rightarrow \mathbb{R}^2$$

betrachtet. Mit Hilfe der Kettenregel können die mit $D_{\hat{x}}$ und $\text{div}_{\hat{x}}$ bezeichneten Ableitungen in den neuen Koordinaten $\hat{x} := \Phi^{-1}(x)$ berechnet werden

$$\begin{aligned} D_{\hat{x}} \hat{u} &= D_{\hat{x}} (u \circ \Phi) = D u \circ D_{\hat{x}} \Phi = D u \circ (c \text{id}) = c D u, \\ \text{div}_{\hat{x}} \hat{\sigma} &= c \text{div} \sigma. \end{aligned}$$

Hier bezeichnen D und div die Differentialoperatoren bezüglich der ursprünglichen Koordinaten x . Betrachte nun das ursprüngliche System aus Problem 8, zu gegebenem $f \in L^2(\Omega; \mathbb{R}^2)$ und $g \in H^{-1/2}(\partial\Omega; \mathbb{R}^2)$ finde $u \in H^1(\Omega; \mathbb{R}^2)$ und $\sigma \in H(\text{div}, \Omega; \mathbb{R}^{2 \times 2}) / \mathbb{R}$ mit

$$\begin{aligned} \text{div} \sigma &= -f && \text{in } \Omega, \\ \text{dev} \sigma &= D u && \text{in } \Omega, \\ u &= g && \text{entlang } \Gamma. \end{aligned}$$

Falls σ und u also Problem 8 lösen, so gilt für $\hat{\sigma}$ und \hat{u} für jedes $\hat{x} \in \hat{\Omega}$

$$\begin{aligned} \text{div}_{\hat{x}} \hat{\sigma}(\hat{x}) &= c \text{div} \sigma(x) = -c f(x) = -c \hat{f}(\hat{x}), \\ \text{dev} \hat{\sigma}(\hat{x}) &= \text{dev} \sigma(x) = D(u)(x) = 1/c D(\hat{u})(\hat{x}). \end{aligned}$$

Also lösen $\hat{\sigma}$ und \hat{u} das folgende System auf $\hat{\Omega}$. Finde zu $\hat{f} \in L^2(\hat{\Omega}; \mathbb{R}^2)$ und $\hat{g} \in$

$H^{-1/2}(\partial\hat{\Omega}; \mathbb{R}^2)$ Funktionen $\hat{u} \in H^1(\hat{\Omega}; \mathbb{R}^2)$ und $\hat{\sigma} \in H(\operatorname{div}, \hat{\Omega}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ mit

$$\begin{aligned} \operatorname{div}_{\hat{x}} \hat{\sigma} &= -c\hat{f} && \text{in } \hat{\Omega}, \\ \operatorname{dev} \hat{\sigma} &= 1/c D_{\hat{x}} \hat{u} && \text{in } \hat{\Omega}, \\ \hat{u} &= \hat{g} && \text{entlang } \hat{\Gamma}. \end{aligned}$$

Werden die ersten beiden Gleichungen mit c multipliziert und die folgenden Substitutionen vorgenommen $\tilde{f} = c^2 f$ und $\tilde{\sigma} = c\hat{\sigma}$ so entsteht eine schönere Darstellung des Problems. Zu $\tilde{f} \in L^2(\hat{\Omega}; \mathbb{R}^2)$ und $\hat{g} \in H^{-1/2}(\partial\hat{\Omega}; \mathbb{R}^2)$ finde $\tilde{\sigma} \in H(\operatorname{div}, \hat{\Omega}; \mathbb{R}^{2 \times 2})/\mathbb{R}$ und $\hat{u} \in H^1(\hat{\Omega}; \mathbb{R}^2)$, so dass

$$\begin{aligned} \operatorname{div}_{\hat{x}} \tilde{\sigma} &= -\tilde{f} && \text{in } \hat{\Omega}, \\ \operatorname{dev} \tilde{\sigma} &= D_{\hat{x}} \hat{u} && \text{in } \hat{\Omega}, \\ \hat{u} &= \hat{g} && \text{entlang } \hat{\Gamma}. \end{aligned}$$

Die Resubstitution $\tilde{\sigma} = 1/c \hat{\sigma}$ ergibt die Lösung auf $\hat{\Omega}$ und $\sigma = \hat{\sigma} (1/c \cdot)$ bzw. $u = \hat{u} (1/c \cdot)$ die auf Ω . Da die Spannung σ in dem Lösungstupel $(\sigma_0, u_0, s_1, t_0) \in X_h$ durch die Größen σ_0 und t_0 beschrieben wird, müssen diese beiden Größen nach dem Lösen auf $\hat{\Omega}$ mit $1/c$ auf Ω zurückskaliert werden.

In der Methode `solveStokesDPG` findet die Skalierung statt. Vor der Berechnung der in (4.1) benötigten Matrizen werden die Koordinaten der Knoten skaliert

```
c4n = 1/c.*c4n; % scaling
```

Während der Berechnung in `computeBlocksStokesDPG`, `computeBlocksStokesDPGFass` bzw. `computeBlocksStokesDPGInetgrate` ist lediglich zu beachten, dass die Auswertungen der Funktionen g und f in den nicht skalierten Punkten erfolgt und die Auswertung von f wie beschrieben zu behandeln ist

```
DbNodes      = unique(n4sDb);           % Dirichlet boundary nodes
DbCoords      = c4n(DbNodes,:);         % Coordinates for DbNodes
values4DbNodes = u4Db(c.*DbCoords);     % Values in DbNodes

temp = repmat(c^2.*[area4e,area4e].*f(c*mid4e)/3,1,3)';
G=sparse(1:6*nrElems,ones(6*nrElems,1),temp(:),6*nrElems,1); % block G
```

Anschließend wird die zum Originalproblem korrespondierende Lösungskomponente σ in `solveStokesDPG` wie beschrieben resubstituiert

```
% rescale
xh(1:6*nrElems,6*nrElems+2*nrNodes:end)= ...
1/c.* xh(1:6*nrElems,6*nrElems+2*nrNodes:end);
```

4.6 Realisierung

Als Erweiterung des AFEM-Software-Paketes [Car09a] wurden die im Folgenden erläuterten `matlab`-Methoden für die Dimension $n = 2$ implementiert. In dem Programm

StokesDPG läuft im Wesentlichen der folgende Standard AFEM-Algorithmus ab.

Algorithm 1: AFEM

input : Reguläre Anfangstriangulierung \mathcal{T}_0 und Bulk-Parameter $0 < \theta \leq 1$
for $\ell = 0, 1, 2, \dots$ **do**
 solve Berechnung der Lösung x_ℓ zu dem diskreten Problem auf \mathcal{T}_ℓ
 estimate Berechnung der lokalen Beiträge $\eta_\ell^2(T)$ für alle $T \in \mathcal{T}_\ell$ und des globalen Fehlerschätzers $\eta_\ell^2 = \sum_{T \in \mathcal{T}} \eta_\ell^2(T)$
 mark Auswahl einer Menge $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell$ von minimaler Kardinalität $|\mathcal{M}_\ell|$ mit $\theta \sum_{T \in \mathcal{T}_\ell} \eta_\ell^2(T) \leq \sum_{T \in \mathcal{M}_\ell} \eta_\ell^2(T)$
 refine Erzeugen einer regulären Triangulierung $\mathcal{T}_{\ell+1}$ minimaler Kardinalität, die eine Verfeinerung von \mathcal{T}_ℓ ist mit $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell \setminus \mathcal{T}_{\ell+1}$
output : Folge von Triangulierungen \mathcal{T}_ℓ , diskreten Lösungen x_ℓ und Fehlerschätzern η_ℓ

Für den Schritt **SOLVE** wurde die Methode **solveDPGStokes** zum Aufstellen und Lösen des linearen Gleichungssystems (4.1) implementiert. Die einzelnen Matrixblöcke werden dabei, wie in Abschnitt 4.3 erläutert, in dem Unterprogramm **computeBlocksStokesDPG** berechnet. Der Schritt **ESTIMATE** erfolgt in **estimateDPGStokes**. Dabei wird der Fehlerschätzer aus Bemerkung 4.2 verwendet. Die Schritte **MARK** und **REFINE** werden mit Hilfe der bereits existierenden AFEM-Methoden **markBulk** zum Dörfler Markieren und **refineRGB** zur Rot-Grün-Blau-Verfeinerung bzw. im uniformen Fall, $\theta = 1$, mit **refineUniformRed** realisiert.

Ist zu einem Problem die exakte Lösung bekannt, so steht zur Berechnung der Abweichung dieser Lösung von der numerischen die Methode **computeErrorStokesDPG** zur Verfügung, wie in Abschnitt 4.4 erläutert. Außerdem sind die Programme **visualisation** und **visualisationExSol** zur Ausgabe von Plots der einzelnen Variablen bzw. der exakten Lösungskomponenten, sofern bekannt, implementiert.

Die verwendeten Datenstrukturen entsprechen denen im AFEM-Software-Paket. Die Triangulierung wird daher mit folgenden Matrizen beschrieben

$$\begin{array}{ll}
 \mathbf{n4e} \in \mathbb{R}^{|\mathcal{T}| \times 3} & \left| \begin{array}{l} \text{Knoten der Elemente,} \\ \text{Koordinaten der Knoten,} \end{array} \right. \\
 \mathbf{c4n} \in \mathbb{R}^{|\mathcal{N}| \times 2} & \\
 \mathbf{n4sDb} \in \mathbb{R}^{|\mathcal{E}(\Gamma)| \times 2} & \left| \begin{array}{l} \text{Knoten der Seiten, die zum Dirichlet Rand gehören,} \\ \text{Knoten der Seiten, die zum Neumann Rand gehören} \end{array} \right. \\
 \mathbf{n4sNb} = \emptyset & \left| \begin{array}{l} \text{(die letzte Matrix wird als leer angenommen, da Probleme} \\ \text{mit reinem Dirichlet Rand betrachtet werden).} \end{array} \right.
 \end{array}$$

Im Folgenden werden die einzelnen Programme mit Eingabeparametern und Ausgaben kurz erläutert.

`StokesDPG(problem,minNrDoF,theta,scale,int,plotdata,savedata,foldername)`

Über **problem** wird im Hauptprogramm der Datensatz zu einem Standard Stokes Problem ausgewählt. Zur Verfügung stehen **'bf'** ("backward facing step" Beispiel), **'cavity'** ("lid driven cavity" Beispiel), **'colliding'** ("colliding flow" Beispiel), **'Lshape'** (Standardbeispiel im L-Gebiet), **'LshapeHighOsc'** (Beispiel im L-Gebiet mit stark oszillierender rechter Seite), **'slit'** (Standardbeispiel im Schlitzgebiet) und **'noPress'** (ein Poisson-Problem Beispiel). Einige dieser Probleme werden detaillierter in Kapitel 5 vorgestellt. Die zugehörigen Datenfiles wurden aus [Bri12] bzw. [Bri14] entnommen und angepasst. Die minimale Zahl der Freiheitsgrade, **minNrDoF**, legt fest, bis zu welchem Level ℓ verfeinert wird. Der Bulkparameter **theta** wird im Dörfler Marking verwendet. Wird der Parameter **scale=1** gesetzt, wird eine problemgerechte Skalierung, wie in Abschnitt 4.5 erläutert, durchgeführt.

Die Eingabe **int** bestimmt mit welcher Genauigkeit die Funktion $\hat{s} \in H^{1/2}(\partial\mathcal{T}; \mathbb{R}^2)$ aus Definition 3.2 über den Rand von Ω integriert wird. Bei **int=0** wird mit linearen Funktionen approximiert, bei **int=1** wird mit der Fassregel integriert und bei **int=2** wird die AFEM-Funktion **integrate** mit **deg=12** verwendet.

Der Parameter **plotdata=1** führt zur Ausgabe von Plots der im letzten Schritt berechneten Lösung und sofern bekannt der exakten Lösungskomponenten. Bei der Belegung **plotdata=2** werden die entsprechenden Plots in jedem Verfeinerungsschritt erzeugt. Außerdem wird in diesen beiden Fällen ein Konvergenzplot des Fehlerschätzers und soweit bekannt des exakten Fehlers und seiner Komponenten erzeugt.

Mit **savedata=1** werden die Fehlerdaten (Fehlerschätzer und falls bekannt exakter Fehler und seine Komponenten) in **Error4lvl.dat** nach jedem Verfeinerungsschritt gespeichert. Bei **savedata=2** werden zusätzlich in jedem Schritt die Koeffizientenvektoren der Lösung (**x.dat** und **xh.dat**) sowie die Triangulierungsdaten (**c4n.dat** und **n4e.dat**) gespeichert. Durch **foldername** wird der Speicherort spezialisiert. Die Daten werden in **tmp/StokesDPG/problem-foldername** oder falls **foldername** leer ist in **tmp/StokesDPG/problem-date** gespeichert.

```
[x,xh,nrDoF,B,Minv,F]=solveStokesDPG(f,u4Db,c4n,n4e,n4sDb,c,int)
```

In diesem Fall werden in **f** das Funktionshändel zu $f \in L^2(\Omega; \mathbb{R}^2)$ und in **u4Db** das zu $g \in L^2(\Gamma; \mathbb{R}^2)$, den Randbedingungen für $u|_\Gamma$ aus Problem 8, übergeben. Wie bereits erläutert wird mit **c4n**, **n4e** und **n4sDb** die Triangulierung beschrieben. Der Parameter **c** ist, soweit eine Skalierung stattfindet, der Skalierungsparameter. Mit diesen Daten werden abhängig vom Parameter **int**, also dem gewünschten Genauigkeitsgrad der Randintegrale, in **computeBlocksStokesDPG**, **computeBlocksStokesDPGFass** bzw. **computeBlocksStokesDPGIntegrate** die Bestandteile des linearen Gleichungssystems (4.1) bestimmt. Dieses wird dann in den freien Komponenten (unter der Beachtung der Nullrandbedingungen für die Komponente s) im unter Umständen skalierten Gebiet gelöst. Die Koeffizientenmatrix **B**, die inverse Normmatrix **Minv** und **F**, der Vektor zum Funktional F , werden zurückgegeben. Außerdem wird die Lösung in Form des Koeffizientenvektors **x** zurückgegeben. Anschließend wird die Lösung des zu Problem 10 korrespondierenden Originalproblems 8, wie in Satz 3.3 erläutert, bestimmt. Der entsprechende Koeffizientenvektor wird in **xh** übergeben. Die Anzahl der Freiheitsgrade wird als **nrDoF** zurückgegeben.

```
[B,Minv,F,L]=computeBlocksStokesDPG(f,u4Db,c4n,n4e,n4sDb,c)
```

Mit denselben Eingabedaten wie in `solvesStokesDPG` werden zunächst u.a. mit Hilfe bereits existierender AFEM-Programme gewisse Triangulierungs- und Problemdaten berechnet und die benötigten lokalen Matrizen initialisiert. Anschließend werden in einer `parfor`-Schleife über die $j = 1, \dots, J$ Elemente T_j , die lokalen Beiträge zu den in Abschnitt 4.3 beschriebenen Matrixblöcken berechnet. Dabei werden die Randintegrale durch lineare Approximation von `u4Db` bestimmt. Die globalen Blöcke werden dann mit Hilfe des `sparse`-Befehls assembliert. Anschließend werden die Koeffizientenmatrix `B`, die inverse Normmatrix `Minv`, `F`, der Vektor zum Funktional F , und der Vektor der linearen Nebenbedingung `L` aus diesen globalen Blöcken zusammengefügt und zurückgegeben.

```
[B,Minv,F,L]=computeBlocksStokesDPGFass(f,u4Db,c4n,n4e, n4sDb,c)
```

Der Unterschied zu `computeBlocksStokesDPG` besteht lediglich darin, dass die Seitennummern der Dirichlet Seiten bestimmt werden und die Funktion `u4Db` entlang dieser mit Hilfe der Fassregel integriert wird. Dadurch kann der Vektor zum Funktional F besser approximiert werden.

```
[B,Minv,F,L]=computeBlocksStokesDPGIntegrate(f,u4Db,c4n,n4e,n4sDb,c)
```

Im Unterschied zu `computeBlocksStokesDPG` werden hier die Seitennummern der Dirichlet Seiten und mit Hilfe der AFEM-Funktion `integrate` die Integrale von `u4Db` (exakt bis Polynomgrad 12) entlang dieser Seiten bestimmt. Diese Ergebnisse fließen in die Berechnung des Vektors zum Funktional F ein.

```
[eta,eta4e]=estimateStokesDPG(n4e,B,Minv,F,x)
```

Mit Hilfe von `n4e`, Knoten der Elemente, `B`, der Koeffizientenmatrix, `Minv`, der inversen Normmatrix, und `F`, des Vektors zum Funktional F , wird der in Bemerkung 4.2 beschriebene residuale Fehlerschätzer für den Koeffizientenvektor `x` zu der Lösung von Problem 10 berechnet. Nachdem die inverse Normmatrix und der residual Vektor entsprechend vorbereitet wurden, wird in einer `parfor`-Schleife der lokale Fehlerschätzer bestimmt. In `eta` wird der globale und in `eta4e` wird der Vektor aller lokalen Fehlerschätzer zurückgegeben.

```
[Eu,Esigma,Es,Et,E]=computeErrorStokesDPG(c4n,n4e,xh,f,uExact1,uExact2,...  
sigmaExact1,sigmaExact2,gradExact1,gradExact2)
```

Hier müssen die Triangulierungsdaten `c4n` und `n4e`, der Koeffizientenvektor zum Originalproblem `xh`, das Funktionshändel `f` zur rechten Seite $f \in L^2(\Omega; \mathbb{R}^2)$ und die Funktionshändel zu den bekannten Komponenten der Lösung übergeben werden. Dabei seien `uExact1` die erste Komponente der Verschiebungsfunktion $u \in H^1(\Omega; \mathbb{R}^2)$, und `uExact2` die zweite. Mit `sigmaExact1` bzw. `sigmaExact2` werden die ersten beiden bzw. letzten

beiden Komponenten (zeilenweise) des Pseudospannungstensors $\sigma \in H(\text{div}, \Omega; \mathbb{R}^{2 \times 2})$ bezeichnet. Ebenso zeilenweise werden die Komponenten des Gradienten des Geschwindigkeitsfeldes $D u \in L^2(\Omega; \mathbb{R}^{2 \times 2})$ mit `gradExact1` bzw. `gradExact2` benannt. Nachdem die Koeffizienten der einzelnen Lösungskomponenten aus `xh` ausgelesen wurden, werden die Komponenten des Fehlers mit Hilfe der `integrate`-Funktion und bereits existierender Fehler-Funktionen aus dem AFEM-Software-Paket bestimmt. Die P_0 -Fehler $\|\sigma - \sigma_0\|_{L^2(\Omega)}^2$, `Esigma`, und $\|u - u_0\|_{L^2(\Omega)}^2$, `Eu`, werden in `error4eP0L2` mit `integrate` bestimmt. Der P_1 -Fehler $\|u - w_c\|_{H^1(\Omega)}^2 = \|u - w_c\|_{L^2(\Omega)}^2 + \|u - w_c\|^2$, `Es`, wird mit Hilfe der AFEM-Funktionen `error4eP1L2` und `error4eP1Energy` berechnet. Zuletzt wird der RT -Fehler $\|\sigma - \mathbf{q}_{RT}\|_{H(\text{div}, \Omega)}^2 = \|\sigma - \mathbf{q}_{RT}\|_{L^2(\Omega)}^2 + \|\text{div } \sigma - \text{div } \mathbf{q}_{RT}\|_{L^2(\Omega)}^2$, `Et`, mit Hilfe einer modifizierten Methode aus [Bri12] bestimmt. Diese Fehlerkomponenten und der gesamte Fehler, `E`, werden zurückgegeben.

```
visualisation(c4n,n4e,xh,eta4e,nrDoF)
```

Mit Hilfe der Triangulierungsdaten, `c4n` und `n4e`, des Koeffizientenvektors zum Originalproblem `xh`, des Fehlerschätzers pro Element `eta4n` und der Zahl der Freiheitsgrade `nrDoF` werden in mehreren Figuren Plots erzeugt. Zunächst werden die Koeffizienten der einzelnen Lösungskomponenten aus `xh` ausgelesen. Dann wird in einer Figur ein `quiver`-plot der Komponente u_0 erzeugt, in einer anderen werden beide Komponenten von u_0 einzeln als `P04e`-Plot dargestellt. Auch die vier Komponenten von σ_0 werden in je einem `P04e`-Plot dargestellt. Ein weiterer `P04e`-Plot enthält den Fehlerschätzer pro Element `eta4n`. Außerdem gibt es noch einen `P1`-Plot der Fortsetzungen der beiden Komponenten von s_1 . Eine weitere Figur enthält einen Plot der Triangulierung.

```
visualisationExSol(c4n,n4e,xh,eta4e,nrDoF,uExact1,uExact2,sigmaExact1,...
                  sigmaExact2)
```

Zusätzlich zu den Eingabedaten aus `visualisation` müssen hier noch die Funktionsfelder zur exakten Lösung, `uExact1` bzw. `uExact2` für die Komponenten der Verschiebungsfunktion $u \in H^1(\Omega; \mathbb{R}^2)$ sowie `sigmaExact1` bzw. `sigmaExact2` für die Zeilen des Pseudospannungstensors $\sigma \in H(\text{div}, \Omega; \mathbb{R}^{2 \times 2})$, übergeben werden. Dann werden dieselben Figuren wie in `visualisation` erzeugt, nur dass diese neben den Plots der berechneten Lösung auch die der entsprechenden Approximation der exakten Lösung zum Vergleich enthalten.

```
convergenceExSol(nrDoF4lv1,eta4lv1,E4lv1,Eu,Esigma,Es,Et)
```

Diese Funktion erzeugt einen Konvergenzplot des Fehlerschätzers, `eta4lv1`, des exakten Fehlers, `E4lv1` sowie der Fehleranteile `Eu`, `Esigma`, `Es`, `Et`. Wie üblich werden diese Größen gegen die Anzahl der Freiheitsgrade, `nrDoF4lv1`, geplottet.

5 Numerische Experimente

In diesem Kapitel werden einige numerische Experimente vorgestellt, die mit dem Programm aus Abschnitt 4 durchgeführt wurden. Es handelt sich um Benchmark Probleme, deren exakte Lösung z.T. bekannt ist und zur Validierung der Methode herangezogen werden kann. Zunächst werden zwei Beispiele auf dem Einheitsquadrat $(-1, 1)^2$ vorgestellt, danach werden Beispiele auf nicht konvexen Gebieten untersucht, dabei kann zu dem letzten keine analytische Lösung bestimmt werden.

Die Experimente, bei denen eine exakte Lösung bekannt ist, zeigen, dass der verwendete Fehlerschätzer den Fehler zwar um einen gewissen Faktor unterschätzt, aber die richtige Konvergenzrate aufweist. Damit wird die Annahme unterstützt, dass der Term $\|F - b(x_h, \bullet)\|_{Y_h^*}$ ein angemessener Schätzer ist. Denn im Allgemeinen verschwindet der Datenapproximationsfehler, $\|F \circ (1 - \Pi)\|_{Y^*}$, nicht, obwohl z.T. konstante rechte Seite f angesetzt werden. In das Funktional F aus Definition 3.2 fließen schließlich auch die nicht konstanten Dirichlet Randbedingungen g durch den Beitrag $\langle \gamma_\nu^T \tau, \hat{s} \rangle_{\partial\mathcal{T}}$ ein.

5.1 Das "colliding flow" Beispiel

Das erste Benchmark Problem verwendet $f \equiv 0$ als rechte Seite auf dem Einheitsquadrat $\Omega = (-1, 1)^2$. Durch die exakte Lösung für das Geschwindigkeitsfeld

$$u((x_1, x_2)^\top) = (20x_1x_2^4 - 4x_1^5, 20x_1^4x_2 - 4x_2^5)^\top \quad \text{für alle } (x_1, x_2)^\top \in \Omega,$$

sind die Randdaten $g(x) = u(x)$ für alle $x \in \Gamma$ vorgegeben. Die Druckverteilung lautet

$$p((x_1, x_2)^\top) = 120x_1^2x_2^2 - 20x_1^4 - 20x_2^4 - 16/3 \quad \text{für alle } (x_1, x_2)^\top \in \Omega.$$

Wie in Abschnitt 2.6 erläutert, kann aus der diskreten Lösung $x_h = (\sigma_0, u_0, s_1, t_0) \in X_h$ eine P_0 -Projektion der diskreten Druckverteilung mittels

$$p_h = -1/2 \operatorname{tr}(\sigma_0) \in P_0(\mathcal{T})$$

bestimmt werden. In Abbildung 5.1 sind diese Druckverteilung und das Geschwindigkeitsfeld, u_0 , auf einer adaptiv mit Bulkparameter $\theta = 0.3$ bestimmten Triangulierung \mathcal{T}_ℓ mit 436 Dreiecken dargestellt. Die stärkste Verfeinerung ist in den Ecken des Gebietes zu erkennen, da dort die Funktion g die größte Steigung hat. In Abbildung 5.2 sind der Fehlerschätzer und die einzelnen Fehleranteile bei uniformer und adaptiver Verfeinerung mit $\theta = 0.3$ dargestellt. Wie auf einem konvexen Gebiet zu erwarten, hat die Verfeinerungsstrategie keinen Einfluss auf die Konvergenzrate. In beiden Fällen konvergiert der Fehlerschätzer von Beginn an mit der optimalen Rate von 0.5. Eine Betrachtung der einzelnen Fehleranteile zeigt, dass die Approximation der Pseudospannung zu Beginn etwas schlechter ist. Dies führt bei der uniformen Methode zu einem leichten vorasymptotischen

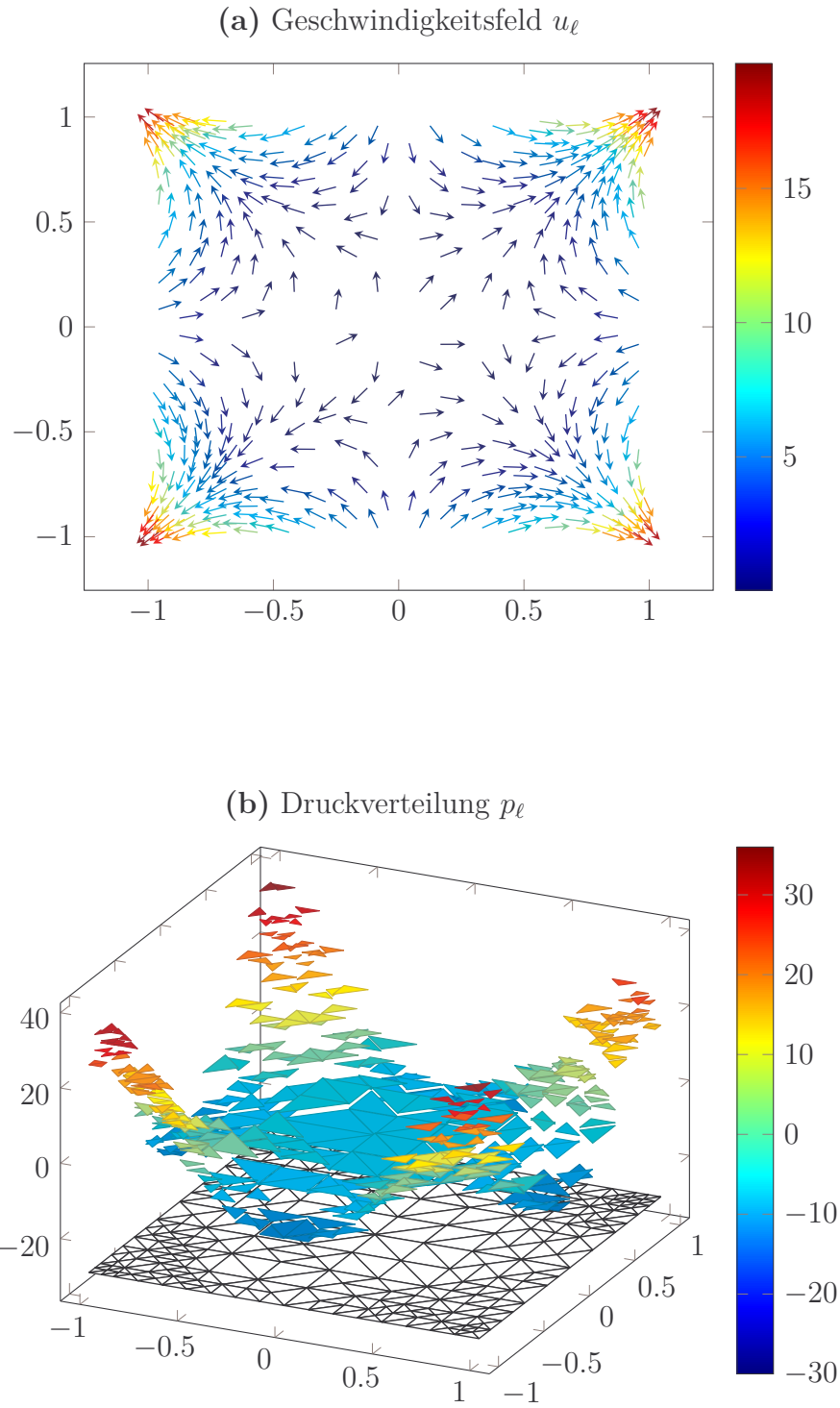
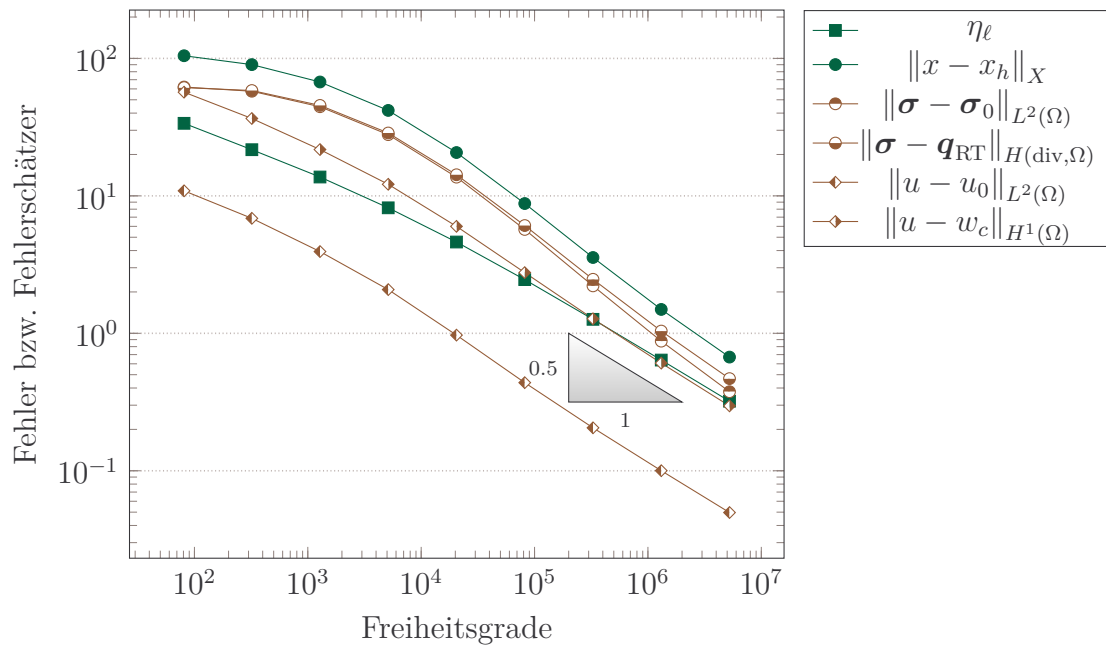


Abbildung 5.1: Plot des Geschwindigkeitsfeldes und der Druckverteilung für das "colliding flow" Beispiel auf einer Triangulierung mit 436 Elementen (4361 Freiheitsgrade), die bei adaptiver Verfeinerung mit $\theta = 0.3$ entstanden ist

(a) Konvergenzplots bei uniformer Verfeinerung



(b) Konvergenzplots bei adaptiver Verfeinerung mit $\theta = 0.3$

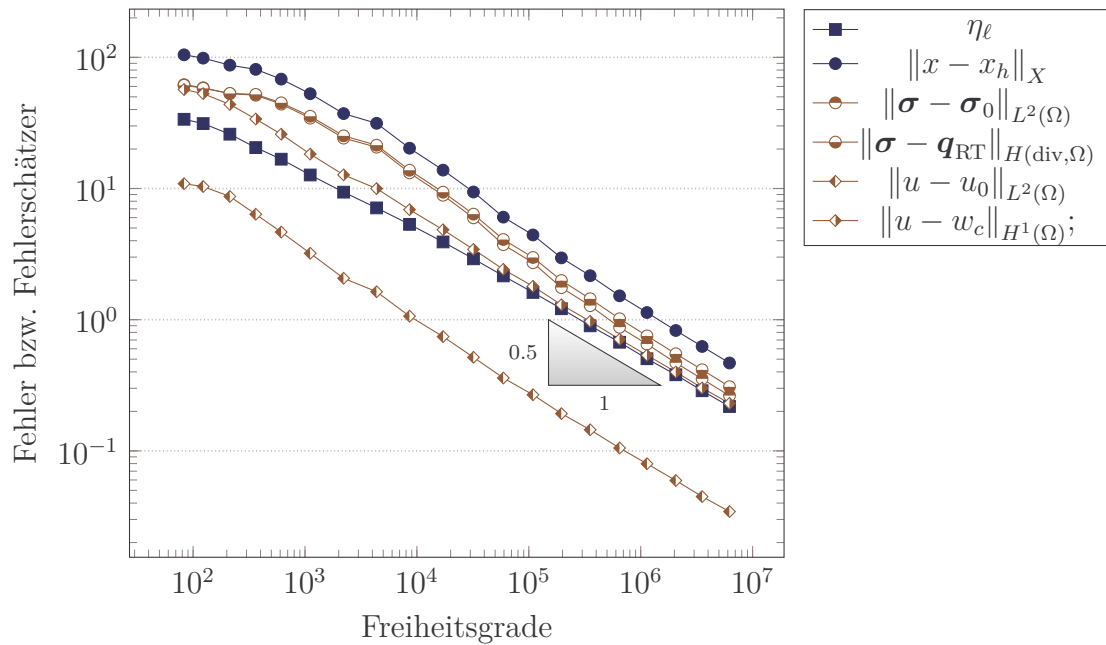


Abbildung 5.2: Konvergenzplots des Fehlerschätzers und der einzelnen Fehleranteile für das "colliding flow" Beispiel bei uniformer und adaptiver Verfeinerung

Effekt beim exakten Fehler. Eine Betrachtung der Komponenten von σ_0 im Vergleich zur exakten Lösung veranschaulicht, dass zu Beginn die Extremalstellen in den Ecken des Gebietes nur schlecht approximiert werden können. Durch die adaptive Verfeinerung wird der Effekt daher deutlich schwächer. Eine genauere Approximation der Randdaten g der Geschwindigkeit beispielsweise durch Verwendung der Fassregel beeinflusst dieses Verhalten nicht. Die Fehler in den Lösungskomponenten, die das Geschwindigkeitsfeld approximieren, konvergieren auch von Beginn an wie erwartet.

5.2 Ein Poisson-Problem Beispiel

Auch hier wird als Gebiet das Einheitsquadrat $\Omega = (-1, 1)^2$ betrachtet. Die rechte Seite ist definiert als

$$f\left((x_1, x_2)^\top\right) = (2x_2, -2x_1) \quad \text{für alle } (x_1, x_2)^\top \in \Omega.$$

Damit ergibt sich als exakte Lösung zu den implizierten Randdaten für das Geschwindigkeitsfeld

$$u\left((x_1, x_2)^\top\right) = \left(-x_1^2 x_2, x_1 x_2^2\right)^\top \quad \text{für alle } (x_1, x_2)^\top \in \Omega.$$

Da hier der Druck $p \equiv 0$ ist, kann dieses Problem als Poisson Problem in zwei Komponenten betrachtet werden. Abbildung 5.3 zeigt, dass sowohl im uniformen als auch im adaptiven Fall die Konvergenzrate des Fehlers und des Fehlerschätzers mit 0.5 optimal sind. Es sind allerdings bei der adaptiven Verfeinerung leichte Abweichungen von der optimalen Kurve zu erkennen. Diese können durch eine exaktere Integration der Randdaten verringert, aber nicht komplett aufgehoben werden. Eine genauere Betrachtung des lokalen Fehlers zeigt, dass dieser im ersten Level bei sechs der acht Elemente beinahe identisch ist, diese Dreiecke aber nicht alle gleichermaßen verfeinert werden. Eine uniforme Verfeinerung vor Beginn des adaptiven Verfeinerns genügt, diesen Effekt zu verhindern.

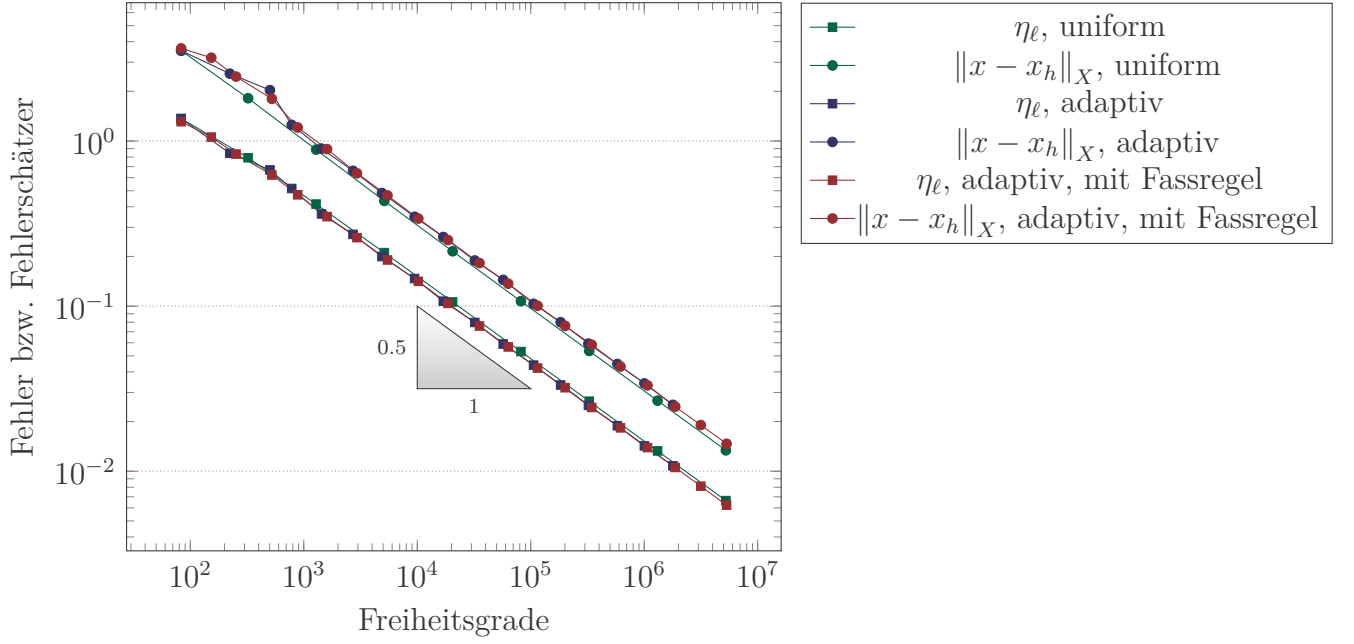


Abbildung 5.3: Konvergenzplot des Fehlerschätzers und exakten Fehlers in dem Poisson-Problem Beispiel bei uniformer Verfeinerung sowie adaptiver Verfeinerung mit $\theta = 0.3$ und verschiedener Behandlung der Randdaten

5.3 Ein Standardbeispiel im L -Gebiet

Dieses Benchmark Problem mit konstanter rechter Seite $f \equiv 0$ ist auf dem L -Gebiet, $\Omega = (-1, 1)^2 \setminus ([0, 1] \times [-1, 0])$, zu lösen. Für die passenden Randdaten $g(x) = u(x)$ für alle $x \in \Gamma$ ist in [Ver99, S.324] eine exakte Lösung in Polarkoordinaten angegeben. Für alle $(r, \varphi) \in [0, \infty) \times [0, 3\pi/2]$ ist mit den Konstanten $\omega = 3\pi/2$ und $\alpha = 856399/1572864$ und der Hilfsfunktion

$$w(\varphi) = \frac{\sin((1 + \alpha)\varphi) \cos(\alpha\omega)}{1 + \alpha} - \cos((1 + \alpha)\varphi) + \frac{\sin((\alpha - 1)\varphi) \cos(\alpha\omega)}{1 - \alpha} + \cos((\alpha - 1)\varphi)$$

das exakte Geschwindigkeitsfeld gegeben durch

$$u(r, \varphi) = r^\alpha \left((1 + \alpha) \sin(\varphi) w(\varphi) + \cos(\varphi) w'(\varphi), -(1 + \alpha) \cos(\varphi) w(\varphi) + \sin(\varphi) w'(\varphi) \right)^\top.$$

Als Druckverteilung ergibt sich

$$p = -r^{\alpha-1} \left((1 + \alpha)^2 w'(\varphi) + w'''(\varphi) \right) / (1 - \alpha).$$

In Abbildung 5.4 ist ein mit adaptiver Verfeinerung mit $\theta = 0.3$ erzeugtes Gitter mit 371 Dreiecken zu sehen. Wie zu erwarten, tritt die stärkste Verfeinerung an der einspringenden Ecke auf. Außerdem ist eine Verfeinerung der Ränder zu erkennen, für die

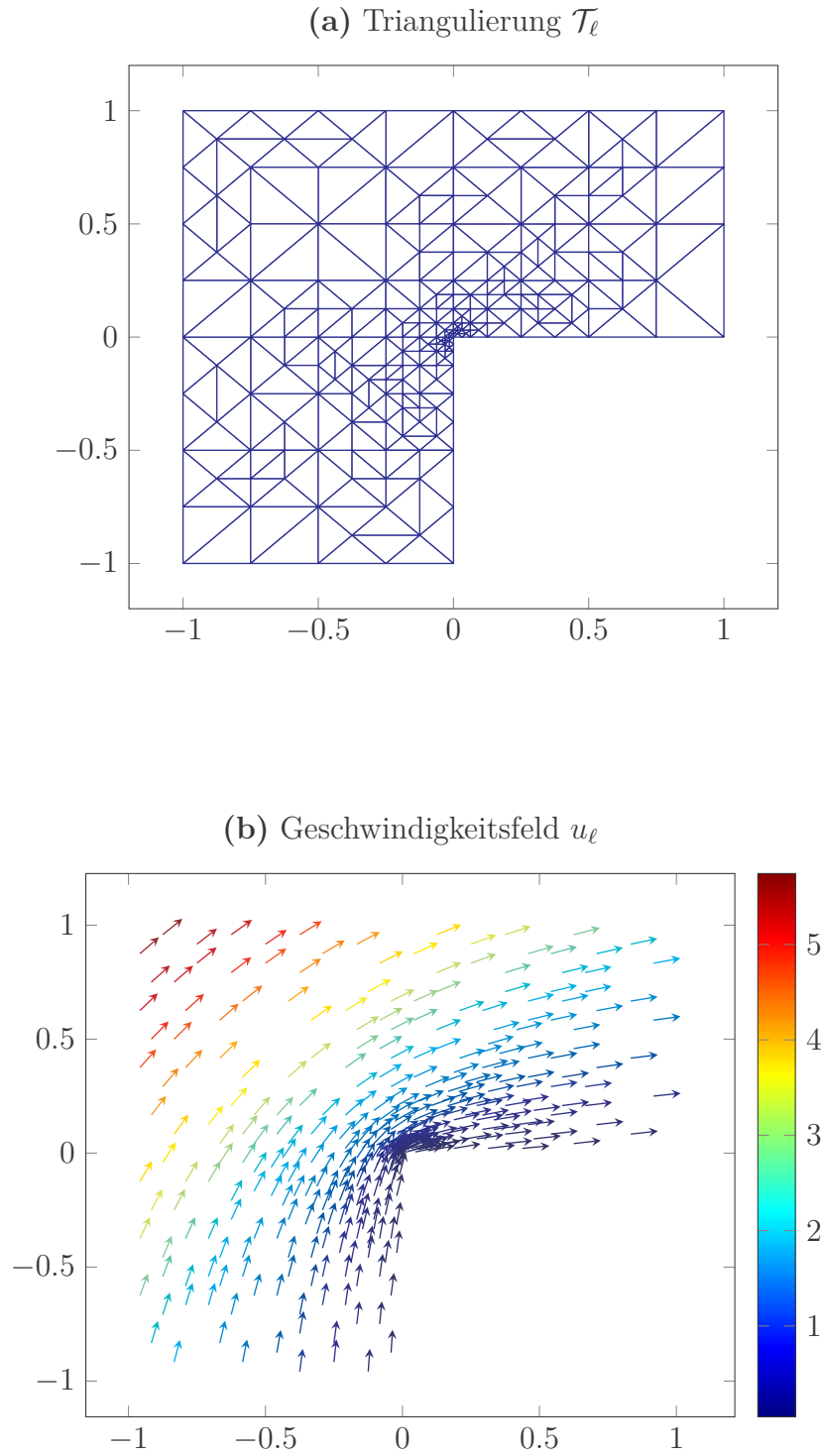
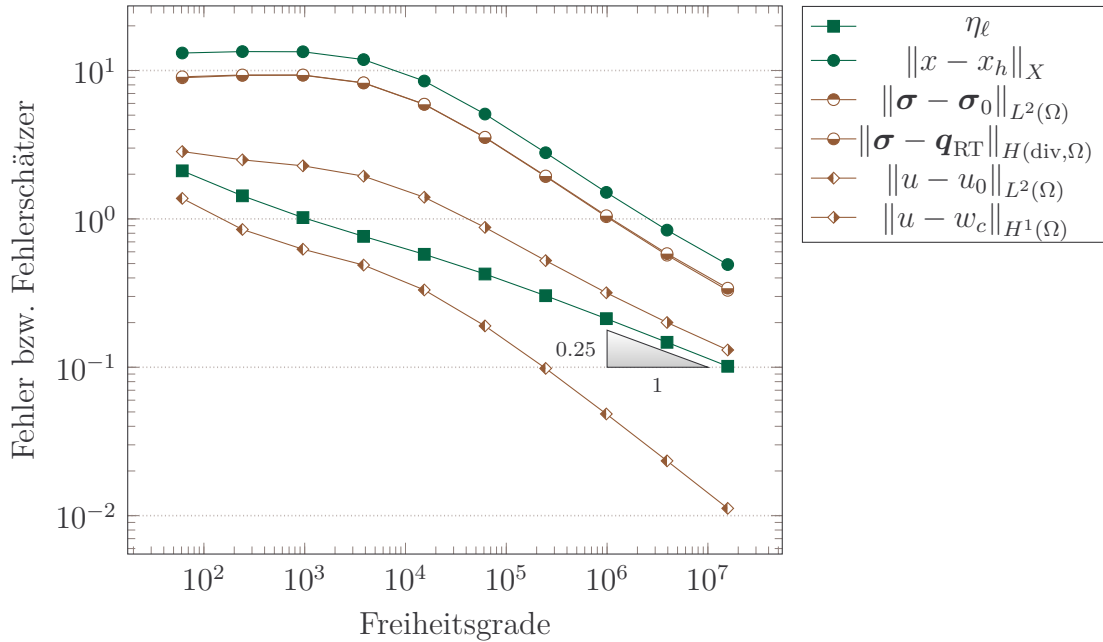
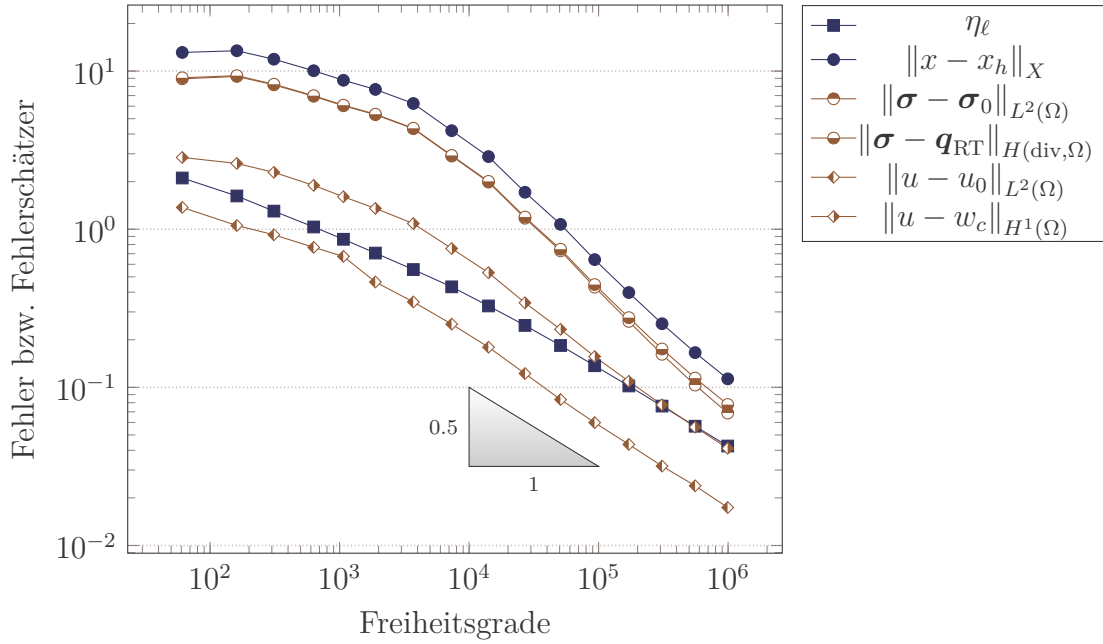


Abbildung 5.4: Triangulierung mit 371 Elementen (3711 Freiheitsgraden) und zugehöriges Geschwindigkeitsfeld in dem Standardbeispiel im L -Gebiet bei adaptiver Verfeinerung mit $\theta = 0.3$

inhomogene Dirichlet Daten gegeben sind. Die Konvergenzplots in Abbildung 5.5 bestätigen, dass die adaptive Verfeinerungsstrategie der uniformen überlegen ist. Im adaptiven Fall wird die optimale Konvergenzrate von 0.5 erreicht. Auf dem nicht konvexen L -Gebiet mit Singularität in der einspringenden Ecke liegt die zu erwartende Konvergenzrate für uniforme Verfeinerung bei 0.25. Diesen Wert erreicht die Methode. Beim exakten Fehler ist ähnlich wie im "colliding flow" Beispiel aus Abschnitt 5.1 ein vorasymptotischer Bereich zu erkennen, der bei der uniformen Methode stärker ausgeprägt ist. Ursache ist wieder die offenbar schlechte Approximation der Pseudospannung. Diese Größe mit Singularität in der einspringenden Ecke ist auf groben Gittern natürlich nur schlecht mit stückweise konstanten bzw. linearen Polynomen zu approximieren. Eine Verbesserung der Randdatenapproximation für die Geschwindigkeit beeinflusst dieses Ergebnis daher nicht. Nach dem vorasymptotischen Bereich scheinen die exakten Fehler zu gut zu konvergieren. Allerdings werden diese durch numerische Integration bestimmt und daher möglicherweise nahe der Singularität unterschätzt. Bei der Betrachtung des gesamten Graphs stimmt auch die Konvergenzrate des Fehlers mit den Erwartungen überein. Insgesamt werden ähnliche Ergebnisse erzielt wie in [Bri12] mit der konformen Least-Squares Methode.

(a) Konvergenzplots bei uniformer Verfeinerung



(b) Konvergenzplots bei adaptiver Verfeinerung mit $\theta = 0.3$

 Abbildung 5.5: Konvergenzplot des Fehlerschätzers und der einzelnen Fehleranteile in dem Standardbeispiel im L -Gebiet bei uniformer und adaptiver Verfeinerung

5.4 Ein Beispiel im Schlitzgebiet

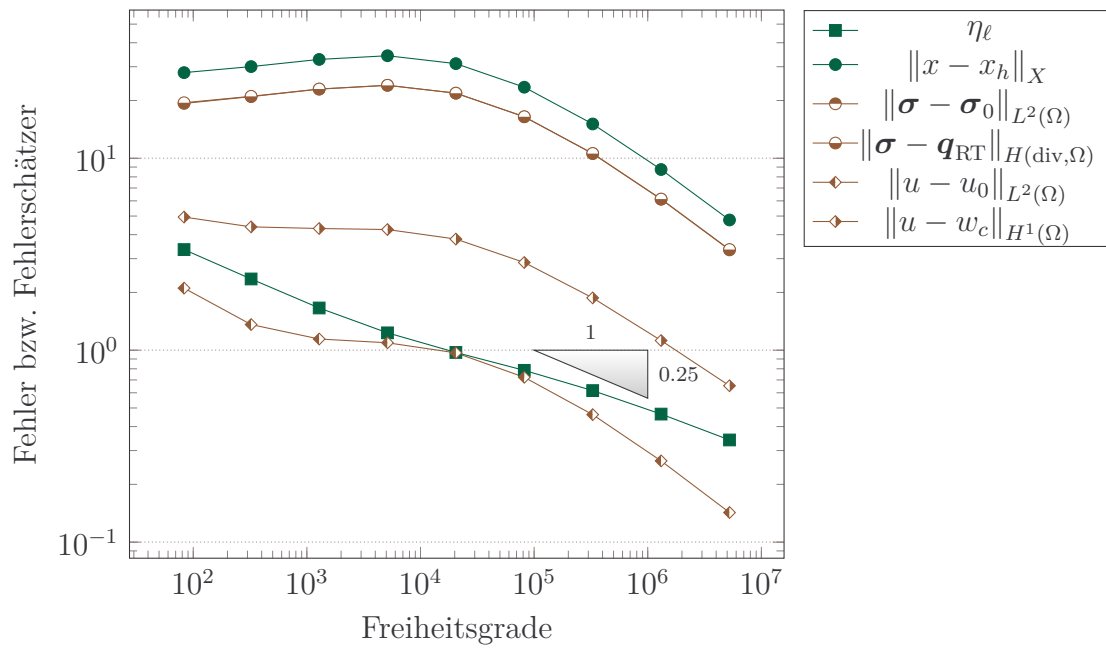
Sei das Schlitzgebiet $\Omega = (-1, 1)^2 \setminus ([0, 1] \times 0)$ vorgegeben. Zu konstanter rechter Seite $f \equiv 0$ und entsprechenden Randdaten kann die exakte Lösung in Polarkoordinaten angegeben werden. Für alle $(r, \varphi) \in [0, \infty) \times [0, 2\pi)$ gilt

$$u(r, \varphi) = \frac{3\sqrt{r}}{2} (\cos(\varphi/2) - \cos(3\varphi/2), 3\sin(\varphi/2) - \sin(3\varphi/2)),$$

$$p(r, \varphi) = \frac{6}{\sqrt{r}} \cos(\varphi/2).$$

In den Konvergenzplots Abbildung 5.6 ist zu erkennen, dass der Fehlerschätzer wieder mit den erwarteten Raten von 0.25 bzw. 0.5 konvergiert. Im Unterschied dazu weist der exakte Fehler einen stark ausgeprägten vorasymptotischen Bereich auf, der auch durch adaptives Verfeinern nicht völlig verschwindet. Danach scheint der exakte Fehler wieder beinah zu gut zu konvergieren. Am Ursprung liegt hier wieder ein Singularität vor, die durch Polynome niedrigen Grades nur schlecht zu approximieren ist. Vergleichbare Ergebnisse werden auch mit der uniformen Least-Squares Methode (siehe [Bri12]) erzielt.

(a) Konvergenzplots bei uniformer Verfeinerung



(b) Konvergenzplots bei adaptiver Verfeinerung mit $\theta = 0.3$

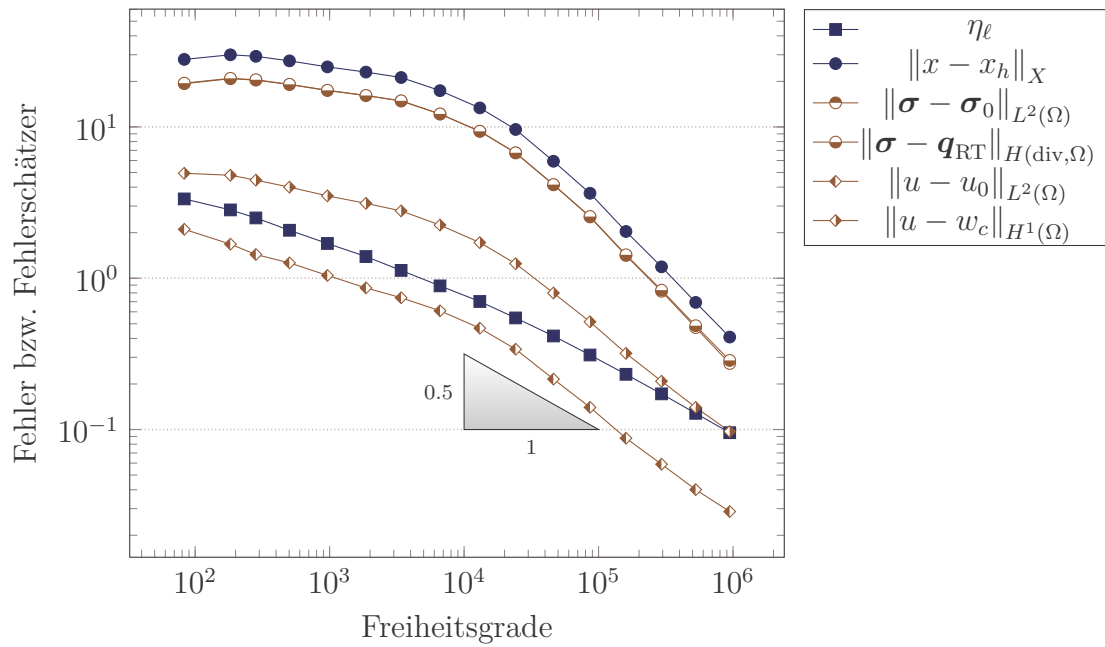


Abbildung 5.6: Konvergenzplot des Fehlerschätzers und der einzelnen Fehleranteile für das Beispiel im Schlitzgebiet bei uniformer und adaptiver Verfeinerung

5.5 Das "backward facing step" Beispiel

Dieses Benchmarkproblem ist anwendungsorientierter als die vorangegangenen Beispiele. Es beschreibt den Fluss einer Newtonschen Flüssigkeit in einer Röhre mit Verengung. Bei diesem Standardbeispiel ist die exakte Lösung nicht bekannt. Im Modell wird die Röhre durch ein verzerrtes L -Gebiet, $\Omega = ((-2, 8) \times (-1, 1)) \setminus ((-2, 0) \times (-1, 0))$ beschrieben, Ein- bzw. Ausfluss werden durch folgende Dirichlet Randdaten für $(x_1, x_2)^\top \in \Gamma$ gegeben

$$g((x_1, x_2)^\top) = \begin{cases} 1/10(-x_2(x_2 - 1), 0)^\top & \text{für } x_1 = -2, \\ 1/80(-(x_2 - 1)(x_2 + 1), 0)^\top & \text{für } x_1 = 8, \\ (0, 0)^\top & \text{sonst,} \end{cases}$$

und die rechte Seite ist konstant $f \equiv 0$. Bei der adaptiven Verfeinerung mit $\theta = 0.3$ entsteht das in Abbildung 5.7 abgebildete Gitter mit 1551 Elementen. Wie zu erwarten, wird an Ein- bzw. Ausfluss besonders stark verfeinert, stärker noch an der kürzeren Seite, wo die Funktion g stärker ansteigt bzw. abfällt. Die einspringende Ecke wird erst relativ spät verfeinert. In Abbildung 5.8 sind die Fehlerschätzer bei adaptiver Verfeinerung mit variierendem Bulkparameter θ dargestellt. Es wird deutlich, dass gilt, je kleiner der Bulkparameter desto besser wird die Konvergenzrate. Auf der anderen Seite wird auch die Anzahl der Rechenschritte größer, die benötigt wird, um eine bestimmte Anzahl von Freiheitsgraden zu erreichen. Die optimalen Konvergenzrate von 0.5 ist auf jeden Fall bereits bei $\theta = 0.3$ gegeben. Der in den anderen Beispielen verwendete Parameter zur adaptiven Verfeinerung erscheint also geeignet. Dieses Gebiet ist leicht verzerrt, eine Skalierung wie in Abschnitt 4.5 ist somit möglich. Der Skalierungsfaktor $c = 10$ ist allerdings so klein, dass er keine sichtbaren Auswirkungen auf die Ergebnisse hat. Ebenso verhält es sich, wenn die Randintegrale mit Hilfe der Fassregel bzw. mit der Methode `integrate` exakter berechnet werden.

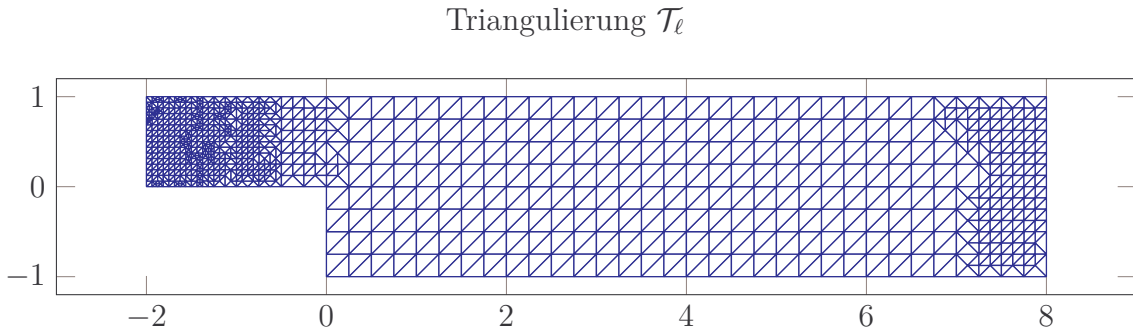


Abbildung 5.7: Triangulierung mit 1551 Elementen (15511 Freiheitsgraden) im "backward facing step" Beispiel bei adaptiver Verfeinerung mit $\theta = 0.3$

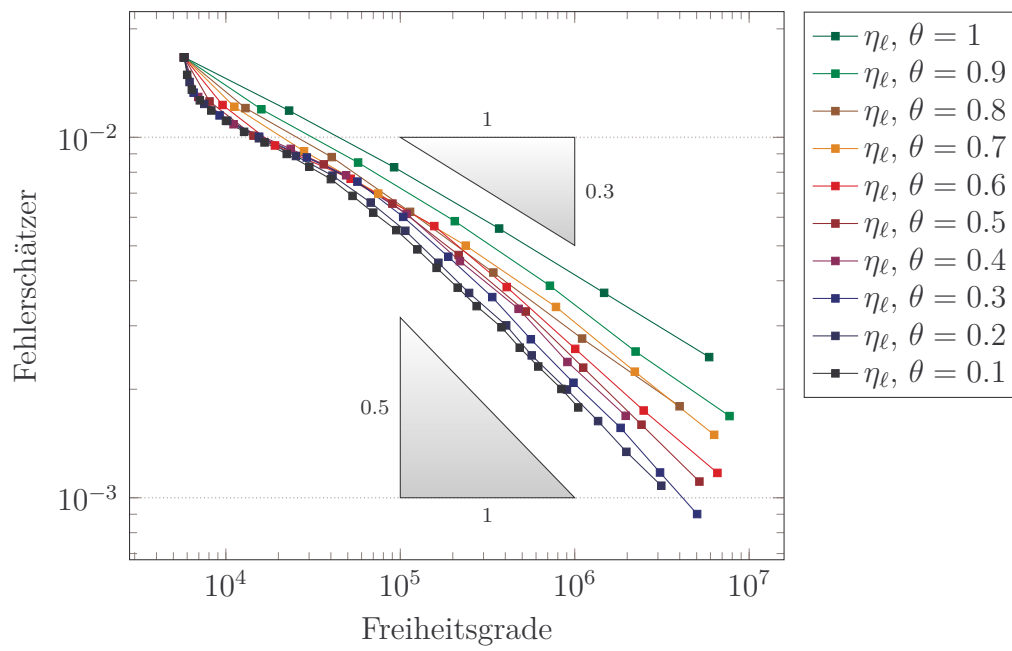


Abbildung 5.8: Konvergenzplot des Fehlerschätzers für das "backward facing step" Beispiel bei verschiedenen Bulkparametern θ

6 Zusammenfassung

In dieser Arbeit wurde eine neue dPG-Methode zum Lösen der Stokes Gleichungen untersucht, die zur Modellierung von Problemen in der Strömungsmechanik verwendet werden. Ausgangspunkt der Untersuchung sind bekannte Ergebnisse zur a-priori- und a-posteriori-Stabilität dieser Methoden, die in Abschnitt 3.1 komprimiert dargestellt werden. Sind Test- und Ansatzraum Hilberträume und ist die zugehörige Bilinearform beschränkt mit trivialen Kern, folgt diese Stabilität bereits aus der kontinuierlichen und diskreten Inf-Sup-Bedingung. Dieser Ansatz ermöglicht die Wahl deutlich kleinerer Räume als bisher in dPG-Methoden für die Stokes Gleichungen vorgeschlagen. Die alternative Herangehensweise basiert auf der Bestimmung eines Fortin-Interpolators. Da dieser Operator in der a-posteriori-Abschätzung im Datenapproximationsfehler auftritt, wird er für den implementierten zweidimensionalen Fall bestimmt und untersucht. Es ist leider nicht möglich zu beweisen, dass dieser Term von höherer Ordnung ist. Die numerischen Experimente bestätigen allerdings, dass der residuale Fehlerschätzer zur adaptiven Verfeinerung genügt.

Die verwendete Bilinearform ergibt sich aus der ultra-schwachen Formulierung der Stokes Gleichungen. In dieser werden die durch die gebrochenen Testfunktionen entstehenden Spurterme auf dem Skelett durch neue Variablen ersetzt. Die Behandlung beliebiger Dirichlet Randdaten wird durch die Einführung eines Spurterms im Funktional der rechten Seite realisiert. Zum Beweis der kontinuierlichen Inf-Sup-Bedingung werden das Splitting Lemma, zwei Dualitätslemmata und die eindeutige Lösbarkeit der gemischten Formulierung verwendet. Die diskrete Inf-Sup-Bedingung wird gezeigt, in dem zu gegebenen Ansatzfunktionen konkrete Testfunktionen gewählt werden. Dabei kommt eine Zerlegung stückweise konstanter deviatorischer Funktionen zum Einsatz, die auf der Lösbarkeit des Stokes Problems beruht. Außerdem erweist sich das Tr-Div-Dev Lemma als äußerst nützlich. Die Beweise werden für beliebige Raumdimensionen durchgeführt und verfolgen die Konstanten explizit.

Das Verfahren wurde für den zweidimensionalen Fall implementiert. Eine ausführliche Dokumentation ist in Kapitel 4 zu finden. Die numerische Experimente zeigen die Stabilität der Methode und optimale Konvergenzraten des verwendeten Fehlerschätzers. In den akademischen Beispielen mit bekannter exakter Lösung wird deutlich, dass der Fehlerschätzer asymptotisch die richtige Konvergenzrate aufweist. Daneben wurden Anwendungsbeispiele untersucht, die analytisch nicht lösbar sind. In allen Versuchen auf nicht konvexen Gebieten ist die adaptive Verfeinerungsstrategie der uniformen deutlich überlegen, wie die Konvergenzgraphen des Schätzers und des tatsächlichen Fehlers zeigen.

Im direkten Anschluss an diese Arbeit wären die Implementierung für höhere Raumdimensionen und Betrachtung höherer Polynomgrade lohnend. Die aktuelle Forschung beschäftigt sich mit dem Nachweis der optimalen Konvergenz für die adaptiven dPG-Methoden. Auf diesem Gebiet liegen noch keine analytischen Ergebnisse vor, daher ist dies äußerst interessant. In diesem Zusammenhang ist sicher auch eine genauere Untersuchung des Einflusses des Datenapproximationsfehlers aufschlussreich.

Literatur

- [ACF99] J. Alpert, C. Carstensen, and S.A. Funken. Remarks around 50 lines of Matlab: short finite element implementation. *Numer. Algorithms*, 20(2–3):117–137, 1999.
- [AF89] D. N. Arnold and R. S. Falk. A uniformly accurate finite element method for the Reissner-Mindlin plate. *SIAM Journal on Numerical Analysis*, 26(6):1276–1290, 1989.
- [Alt06] H. W. Alt. *Lineare Funktionalanalysis: Eine anwendungsorientierte Einführung*, 5. Auflg, 2006.
- [BC05] C. Bahriawati and C. Carstensen. Three Matlab implementations of the lowest-order Raviart-Thomas MFEM with a posteriori error control. *CMAM*, 5(4):333–361, 2005.
- [BMS02] C L Bottasso, S Micheletti, and R Sacco. The discontinuous Petrov–Galerkin method for elliptic problems. *Computer Methods in Applied Mechanics and Engineering*, 191(31):3391–3409, 2002.
- [Bra13] D. Braess. *Finite Elemente: Theorie, schnelle Löser und Anwendungen in der Elastizitätstheorie*. Springer-Verlag, 2013.
- [Bri12] P. Bringmann. *Least-Squares Finite-Elemente-Methoden für die Stokes Gleichungen*. Bachelorarbeit, Humboldt-Universität zu Berlin, 2012.
- [Bri14] P. Bringmann. *Optimal Mesh-Refinement for Incompressible Fluid Dynamics*. Masterarbeit, Humboldt-Universität zu Berlin, 2014.
- [BS08] S.C. Brenner and R. Scott. *The Mathematical Theory of Finite Element Methods*. Texts in Applied Mathematics. Springer, 2008.
- [Car09a] C. Carstensen. AFEM Software-Paket. Humboldt Universität zu Berlin, 2009.
- [Car09b] C. Carstensen. Yonsei Lectures at the WCU Department Computational Science and Engineering on Finite Element Method. Lokal verfügbar, 2009.
- [CDG14] C. Carstensen, L. Demkowicz, and J. Gopalakrishnan. A posteriori error control for DPG methods. *SIAM Journal on Numerical Analysis*, 52(3):1335–1353, 2014.
- [CDG15a] C Carstensen, L Demkowicz, and J Gopalakrishnan. Breaking spaces and forms for the DPG method and applications including Maxwell equations. *arXiv preprint arXiv:1507.05428*, 2015.
- [CDG15b] C. Carstensen, L. Demkowicz, and J. Gopalakrishnan. The paradigm of broken test-functions in DPG discretisations of elliptic second-order PDEs. in Vorbereitung, 2015.

- [CG14] C. Carstensen and D. Gallistl. Guaranteed lower eigenvalue bounds for the biharmonic equation. *Numer. Math.*, 126(1):33–51, 2014.
- [CGHW14] C. Carstensen, D. Gallistl, F. Hellwig, and L. Weggler. Low-order dPG-FEM for an elliptic PDE. *Comput. Math. Appl.*, 68(11):1503–1512, 2014.
- [CGS14] C. Carstensen, D. Gallistl, and M. Schedensack. L2 Best-Approximation of the Elastic Stress in the Arnold-Winther FEM. *Preprint 2014-15, Humboldt-Universität zu Berlin, Institut für Mathematik*, 2014.
- [CH15] C. Carstensen and F. Hellwig. Low-order DPG-FEMs for linear elasticity. *SIAM Journal on Numerical Analysis*, eingereicht, 2015.
- [CPR13] C. Carstensen, D. Peterseim, and H. Rabus. Optimal adaptive nonconforming FEM for the Stokes problem. *Numer. Math.*, 123(2):291–308, 2013.
- [CTVW10] Z. Cai, C. Tong, P. S Vassilevski, and C. Wang. Mixed finite element methods for incompressible flow: Stationary Stokes equations. *Numerical Methods for Partial Differential Equations*, 26(4):957–978, 2010.
- [CW07] Z. Cai and Y. Wang. A multigrid method for the pseudostress formulation of Stokes problems. *SIAM Journal on Scientific Computing*, 29(5):2078–2095, 2007.
- [DG10] L. Demkowicz and J. Gopalakrishnan. A class of discontinuous Petrov–Galerkin methods. Part I: The transport equation. *Computer Methods in Applied Mechanics and Engineering*, 199(23):1558–1572, 2010.
- [DG11a] L. Demkowicz and J. Gopalakrishnan. Analysis of the DPG method for the Poisson equation. *SIAM Journal on Numerical Analysis*, 49(5):1788–1809, 2011.
- [DG11b] L. Demkowicz and J. Gopalakrishnan. A class of discontinuous Petrov–Galerkin methods. Part II. Optimal test functions. *Numerical Methods for Partial Differential Equations*, 27(1):70–105, 2011.
- [DGN12] L. Demkowicz, J. Gopalakrishnan, and A. H. Niemi. A class of discontinuous Petrov–Galerkin methods. Part III: adaptivity. *Applied numerical mathematics*, 62(4):396–427, 2012.
- [GQ14] J. Gopalakrishnan and W. Qiu. An analysis of the practical DPG method. *Mathematics of Computation*, 83(286):537–552, 2014.
- [GR86] Vivette Girault and Pierre-Arnaud Raviart. *Finite element methods for Navier-Stokes equations: theory and algorithms*. Springer, 1986.
- [Hel14] F. Hellwig. *Drei dPG-Methoden niedriger Ordnung für Lineare Elastizität*. Masterarbeit, Humboldt-Universität zu Berlin, 2014.
- [Kön04] K. Königsberger. *Analysis 2*. Springer-Verlag Berlin Heidelberg New York, 4. auflage edition, 2004.

- [LS10] R. S. Laugesen and B. A. Siudeja. Minimizing Neumann fundamental tones of triangles: an optimal Poincaré inequality. *Journal of Differential Equations*, 249(1):118–135, 2010.
- [PW60] L. E. Payne and H. F. Weinberger. An optimal Poincaré inequality for convex domains. *Archive for Rational Mechanics and Analysis*, 5(1):286–292, 1960.
- [RBTD14] N. V. Roberts, T. Bui-Thanh, and L. Demkowicz. The DPG method for the Stokes problem. *Computers & Mathematics with Applications*, 67(4):966 – 995, 2014. High-order Finite Element Approximation for Partial Differential Equations.
- [Ver89] R. Verfürth. A posteriori error estimators for the Stokes equations. *Numerische Mathematik*, 55(3):309–325, 1989.
- [Wer11] D. Werner. *Funktionalanalysis*. Springer Berlin Heidelberg, 2011.

Danksagung

Es gibt ohne Frage viele Menschen, die mir bei meinem Studium und dieser Arbeit geholfen haben.

Besonders möchte ich mich an dieser Stelle bei Prof. Carstensen für die umfangreiche Betreuung bedanken. Die vielen Gespräche und Diskussionen haben mir sehr geholfen. Außerdem möchte ich den Mitgliedern der Arbeitsgruppe vom Prof. Carstensen danken. Sie alle hatten immer ein offenes Ohr für Fragen und haben Probleme mit mir diskutiert. Friedrike Hellwig konnte mir bei vielen Problemen mit den dPG-Methoden helfen, Philipp Bringmann hat sich meinen `matlab`-Fragen angenommen und Karoline Köhler hat mir zugehört, wann immer ich nicht weiter kam.

Meiner Familie kann ich gar nicht genug danken. Es ist schön zu wissen, dass ich immer kommen kann und ihr versuchen werdet mir zu helfen. Papa, ich freue mich, wenn du auch zukünftig meine Arbeiten liest und mir sagst an welcher Stelle nur noch Eingeweihte verstehen, was ich meine. Ein großes Danke geht an Giordana Tornow, Linda Nguyen und Maren Strobl. Mit euch hat sogar das Studieren Spaß gemacht. Lucian, bei dir möchte ich mich auch bedanken, nicht nur für das Finden von gefühlt tausenden Tippfehlern, aber das weißt du.

Selbstständigkeitserklärung

Ich erkläre, dass ich die vorliegende Arbeit selbstständig und nur unter Verwendung der angegebenen Quellen und Hilfsmittel angefertigt habe und ich zum ersten Mal eine Masterarbeit in diesem Studiengang einreiche.

Berlin, den 29. September 2015

Unterschrift