

Chapter 1

Introduction

1.1 Strong Form and Variational Form (06.10.2017)

The finite element method is a scheme for the numerical solution of partial differential equations. In this chapter, we introduce the basic concepts for elliptic problems in the frame of the Riesz theorem. To that end, we consider the most standard example, namely the Poisson equation with mixed Dirichlet-Neumann boundary conditions. We aim to solve

v1:
04.10.2017

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \Gamma_D, \\ \partial u / \partial n &= \phi && \text{on } \Gamma_N, \end{aligned} \tag{1.1}$$

which is said to be the **strong form** of the boundary value problem. Here, Ω denotes a domain in \mathbb{R}^d , $d = 2, 3$. The boundary $\Gamma := \partial\Omega$ is split into the Dirichlet boundary Γ_D and the Neumann boundary Γ_N , respectively. To be more precise, we assume that Γ_D and Γ_N are (relatively) open subsets of Γ with $\Gamma_D \cap \Gamma_N = \emptyset$ and $\Gamma = \bar{\Gamma}_D \cup \bar{\Gamma}_N$. The source term $f : \Omega \rightarrow \mathbb{R}$ as well as the Neumann data $\phi : \Gamma_N \rightarrow \mathbb{R}$ are given, and $u : \Omega \rightarrow \mathbb{R}$ is the unknown solution. Moreover,

$$\Delta u(x) := \sum_{j=1}^d \frac{\partial^2 u}{\partial x_j^2}(x) \tag{1.2}$$

denotes the Laplace operator, which is defined in the classical sense for a function $u \in C^2(\bar{\Omega})$, where $C^k(\bar{\Omega}) := \{w|_{\Omega} \mid w \in C^k(\mathbb{R}^d)\}$. If $u \in C^2(\bar{\Omega})$ solves (1.1), u is said to be a **strong solution** of the mixed boundary value problem.

Throughout the lecture, we shall assume that Ω is a **Lipschitz domain** in \mathbb{R}^d , i.e.,

- Ω is a bounded, open, and connected subset of \mathbb{R}^d ,
- Ω is locally on one side of Γ ,
- Γ can locally be parametrized by Lipschitz continuous functions.

An important consequence of this assumption is the validity of the **integration by parts formula**

$$\int_{\Omega} \frac{\partial u}{\partial x_j} v \, dx + \int_{\Omega} u \frac{\partial v}{\partial x_j} \, dx = \int_{\Gamma} u v n_j \, ds \quad \text{for all } u, v \in C^1(\bar{\Omega}), \tag{1.3}$$

where n_j denotes the j -th component of the outer normal vector of Ω on Γ and where ds denotes the surface measure on Γ . For a precise definition and details, we refer, e.g., to [McL].

Let $u \in C^2(\overline{\Omega})$ be a strong solution of (1.1) and $v \in C_D^1(\overline{\Omega}) := \{w \in C^1(\overline{\Omega}) \mid w|_{\Gamma_D} = 0\}$. Multiplication of $-\Delta u = f$ by v , integration over Ω , and integration by parts yield that

$$\int_{\Omega} f v \, dx = - \int_{\Omega} (\Delta u) v \, dx = - \sum_{j=1}^d \int_{\Omega} \frac{\partial^2 u}{\partial x_j^2} v \, dx = \sum_{j=1}^d \left[\int_{\Omega} \frac{\partial u}{\partial x_j} \frac{\partial v}{\partial x_j} \, dx - \int_{\Gamma} \frac{\partial u}{\partial x_j} v n_j \, ds \right].$$

With $x \cdot y = \sum_{j=1}^d x_j y_j$ the usual scalar product in \mathbb{R}^d , we obtain the **first Green formula**

$$\int_{\Omega} f v \, dx = \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\Gamma} \frac{\partial u}{\partial n} v \, ds, \quad (1.4)$$

where we have used $\nabla u \cdot n = \partial u / \partial n$. Together with $v|_{\Gamma_D} = 0$ and $\Gamma_N = \Gamma \setminus \overline{\Gamma_D}$, we may plug-in the Neumann data to see that

$$\int_{\Omega} f v \, dx = \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\Gamma_N} \frac{\partial u}{\partial n} v \, ds = \int_{\Omega} \nabla u \cdot \nabla v \, dx - \int_{\Gamma_N} \phi v \, ds.$$

Altogether we thus have proven the following proposition:

Proposition 1.1. *Let $u \in C^2(\overline{\Omega})$ solve the strong form (1.1). Then, it holds that*

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx + \int_{\Gamma_N} \phi v \, ds \quad \text{for all } v \in C_D^1(\overline{\Omega}), \quad (1.5)$$

*which is the **variational form** of the boundary value problem (1.1). ■*

This proposition gives a necessary condition for a function u to solve the strong form (1.1). We stress that any strong solution belongs to $C_D^1(\overline{\Omega})$ and that the variational form (1.5) can be understood for $u \in C_D^1(\overline{\Omega})$. This leads to a symmetric variational formulation: Find $u \in C_D^1(\overline{\Omega})$ such that (1.5) holds.

Exercise 1. Prove the following well-known integral formulae:

- For $f \in C^1(\Omega)^d$, let $\operatorname{div} f := \sum_{j=1}^d \frac{\partial f_j}{\partial x_j}$ denote the divergence operators. Then, there holds the **Gauss divergence theorem**

$$\int_{\Omega} \operatorname{div} f \, dx = \int_{\Gamma} f \cdot n \, ds \quad \text{for all } f \in C^1(\overline{\Omega})^d. \quad (1.6)$$

- Besides the first Green formula, there holds the **second Green formula**

$$\int_{\Omega} (-\Delta u) v \, dx + \int_{\Gamma} \frac{\partial u}{\partial n} v \, ds = \int_{\Omega} u (-\Delta v) \, dx + \int_{\Gamma} u \frac{\partial v}{\partial n} \, ds \quad \text{for all } u, v \in C^2(\overline{\Omega}). \quad (1.7)$$

Both are easily obtained from the integration by parts formula. □

1.2 Solvability of Variational Form (06.10.2017)

To look for solutions of the weak form (1.5), we will employ the following Riesz theorem.

Theorem 1.2 (Riesz). *For a Hilbert space H (over \mathbb{R}), the mapping*

$$I_H : H \rightarrow H^*, \quad I_H(u) := (u ; \cdot)_H \quad (1.8)$$

is linear, isometric, and bijective, i.e., for any $F \in H^$ there is a unique $u \in H$ such that*

$$(u ; v)_H = F(v) \quad \text{for all } v \in H. \quad (1.9)$$

Moreover, it holds that $\|u\|_H = \|F\|_{H^}$. ■*

First, we observe that the left-hand side

$$(u ; v) := \int_{\Omega} \nabla u \cdot \nabla v \, dx$$

of the variational form (1.5) defines a scalar product on $C_D^1(\overline{\Omega})$, provided the Dirichlet boundary Γ_D is nontrivial: Clearly, $(u ; v)$ is a symmetric bilinear form on $C_D^1(\overline{\Omega})$. It thus only remains to prove definiteness. Note that $0 = (u ; u) = \|\nabla u\|_{L^2(\Omega)}^2$ implies $\nabla u = 0$, whence u is constant in Ω . Together with $u|_{\Gamma_D} = 0$, this proves $u = 0$. Moreover, the right-hand side

$$F(v) := \int_{\Omega} f v \, dx + \int_{\Gamma_N} \phi v \, ds$$

defines a linear functional on $C_D^1(\overline{\Omega})$ which is continuous with respect to the induced norm $\|v\| := (v ; v)^{1/2}$. We prove this claim only in the special situation $\Gamma = \Gamma_D$ and postpone the abstract proof to a subsequent section.

Lemma 1.3 (Friedrichs' inequality). *Suppose that $\Omega = [a, b] \times [c, d] \subset \mathbb{R}^2$ and $\Gamma_D = \partial\Omega$. Then, it holds that $\|v\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|\nabla v\|_{L^2(\Omega)}$ for all $v \in C_D^1(\overline{\Omega})$.*

Proof. For $x = (x_1, x_2) \in \Omega$, it holds that $v(x_1, c) = 0$. Therefore, the fundamental theorem of calculus yields that

$$v(x) = \int_c^{x_2} \partial_2 v(x_1, t) \, dt.$$

The Hölder inequality yields that

$$|v(x)| \leq |d - c|^{1/2} \left(\int_c^{x_2} |\partial_2 v(x_1, t)|^2 \, dt \right)^{1/2}.$$

Integration over Ω gives

$$\begin{aligned} \|v\|_{L^2(\Omega)}^2 &= \int_{\Omega} |v(x)|^2 dx \leq |d-c| \int_{\Omega} \int_c^{x_2} |\partial_2 v(x_1, t)|^2 dt dx \\ &= |d-c| \int_c^d \int_a^b \int_c^{x_2} |\partial_2 v(x_1, t)|^2 dt dx_1 dx_2 \\ &\leq |d-c| \int_c^d \|\partial_2 v\|_{L^2(\Omega)}^2 dx_2 \\ &= |d-c|^2 \|\partial_2 v\|_{L^2(\Omega)}^2. \end{aligned}$$

This results in $\|v\|_{L^2(\Omega)} \leq |d-c| \|\partial_2 v\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|\nabla v\|_{L^2(\Omega)}$. ■

According to the Hölder and the Friedrichs inequality, we obtain that

$$|F(v)| \leq \|f\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \leq \text{diam}(\Omega) \|f\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} = \text{diam}(\Omega) \|f\|_{L^2(\Omega)} \|v\|.$$

Therefore, the linear functional F is continuous with respect to $\|\cdot\| := \|\nabla(\cdot)\|_{L^2(\Omega)}$ with operator norm $\|F\|_* \leq \text{diam}(\Omega) \|f\|_{L^2(\Omega)}$. If $C_D^1(\overline{\Omega})$ associated with the norm $\|\cdot\|$ were a Hilbert space, the Riesz theorem *would* therefore imply the unique solvability of the variational form (1.5). However, $C_D^1(\overline{\Omega})$ is *not* complete and therefore the Riesz theorem does *not* apply.

The remedy is to consider the (unique) completion of $C_D^1(\overline{\Omega})$ with respect to $\|\cdot\|$. This leads to a so-called **Sobolev space** $H_D^1(\Omega)$, which is —by definition— complete and hence a Hilbert space. Density arguments then lead to an extended variational form: Find $u \in H_D^1(\Omega)$ such that

$$\int_{\Omega} \nabla u \cdot \nabla v dx = \int_{\Omega} f v dx + \int_{\Gamma_N} \phi v ds \quad \text{for all } v \in H_D^1(\Omega), \quad (1.10)$$

which is the **weak form** of the boundary value problem (1.1). Now, the Riesz theorem applies and proves the unique existence of a **weak solution** $u \in H_D^1(\Omega)$ of (1.10). Later on, we are going to show that

- each strong solution $u \in C^2(\overline{\Omega})$ of (1.1) belongs to $H_D^1(\Omega)$ and is also the unique weak solution of (1.10).
- provided the weak solution $u \in H_D^1(\Omega)$ is smooth, i.e., $u \in C^2(\overline{\Omega})$, the weak solution also solves the strong form (1.1).

In this sense, the strong form (1.1) and the weak form (1.10) are equivalent.

1.3 Finite Element Method (10.10.2017)

The finite element method for (1.10) essentially consists of replacing the (infinite dimensional) Sobolev space $H_D^1(\Omega)$ by a finite dimensional subspace $X_h \subset H_D^1(\Omega)$: Find $u_h \in X_h$ such that

$$\int_{\Omega} \nabla u_h \cdot \nabla v_h dx = \int_{\Omega} f v_h dx + \int_{\Gamma_N} \phi v_h ds \quad \text{for all } v_h \in X_h. \quad (1.11)$$

This problem is equivalent to the solution of a system of linear equations $\mathbf{Ax} = \mathbf{b}$, where the system matrix \mathbf{A} is symmetric and positive definite. Of course, the question of convergence depends on the choice of X_h . Thus, there remain some topics for mathematical discussions later on.

The finite element method is a special **Galerkin scheme**. In this section, we collect the most simple properties of Galerkin schemes. Throughout, H is a (real) Hilbert space, and $\langle \cdot ; \cdot \rangle$ is an equivalent scalar product on H , i.e., there are constants $\alpha, \beta > 0$ such that

$$\alpha \|v\|_H \leq \|v\| \leq \beta \|v\|_H \quad \text{for all } v \in H, \quad (1.12)$$

where $\|v\| := \langle v ; v \rangle^{1/2}$ denotes the induced norm. We stress that $\langle \cdot ; \cdot \rangle$ and $\| \cdot \|$ are often called **energy scalar product** and **energy norm**, respectively (see also Exercise 5).

Remark. In the following, we state all results with respect to the norm $\| \cdot \|_H$, which involves the constants $\alpha, \beta > 0$. Analogously, one may state the results with respect to the energy norm $\| \cdot \| = \| \cdot \|_H$, which corresponds to $\alpha = \beta = 1$. \square

For given $F \in H^*$, the Riesz theorem proves the existence and uniqueness of a solution $u \in H$ of

$$\langle u ; v \rangle = F(v) \quad \text{for all } v \in H, \quad (1.13)$$

for what we use the short-hand notation

$$\langle u ; \cdot \rangle = F \in H^* \quad (1.14)$$

to implicitly indicate that this equation holds (pointwise) for all $v \in H$. Now, the Galerkin method simply consists in replacing the continuous space H by some finite dimensional subspace: Let X_h be a finite-dimensional (and hence closed) subspace of H . Since the Riesz theorem applies to the Hilbert space X_h as well, there is a unique **Galerkin solution** $u_h := \mathbb{G}_h u \in X_h$ such that

$$\langle \mathbb{G}_h u ; \cdot \rangle = F \in X_h^*. \quad (1.15)$$

For $u \in H$ and the corresponding functional $\langle u ; \cdot \rangle \in H^*$, this defines the **Galerkin projection**

$$\mathbb{G}_h : H \rightarrow X_h \quad \text{where } \mathbb{G}_h u \in X_h \text{ solves } \langle \mathbb{G}_h u ; \cdot \rangle = \langle u ; \cdot \rangle \in X_h^*. \quad (1.16)$$

Note that $\mathbb{G}_h u \in X_h$ is characterized by the **Galerkin orthogonality**

$$\langle u - \mathbb{G}_h u ; v_h \rangle = 0 \quad \text{for all } v_h \in X_h. \quad (1.17)$$

Before we proceed with the theoretical analysis of Galerkin schemes, we treat an implementational issue. The following theorem is the fundamental observation: Usually, only the scalar product $\langle \cdot ; \cdot \rangle$ and the right-hand side $F \in H^*$ are known, while the exact solution $u \in H$ of (1.13) is unknown. Then, the Galerkin solution $\mathbb{G}_h u \in X_h$ can be computed by solving a linear system of equations — without knowledge of u .

Theorem 1.4. *Let $\{\phi_1, \dots, \phi_N\}$ be a basis of X_h . We define the Galerkin matrix $A \in \mathbb{R}^{N \times N}$ and the vector $b \in \mathbb{R}^N$ by*

$$A_{jk} := \langle \phi_k ; \phi_j \rangle \quad \text{and} \quad b_j := F(\phi_j). \quad (1.18)$$

Then, A is symmetric and positive definite and, in particular, a regular matrix. Moreover, there holds $\mathbb{G}_h u = \sum_{j=1}^N x_j \phi_j$, where the vector $x \in \mathbb{R}^N$ solves $Ax = b$.

Proof. 1. step. Symmetry of A clearly follows from the symmetry of $\langle \cdot ; \cdot \rangle$.

2. step. For any $x \in \mathbb{R}^N$ and $v_h := \sum_{j=1}^N x_j \phi_j$, it holds that

$$\|v_h\|^2 = \langle v_h ; v_h \rangle = \sum_{j,k=1}^N x_j x_k \langle \phi_j ; \phi_k \rangle = x \cdot Ax.$$

This proves $Ax \cdot x > 0$ for all $x \neq 0$. By definition, A is positive definite and hence regular.

3. step. Determine Galerkin solution: Let $x \in \mathbb{R}^n$ be the unique solution of the linear Galerkin system $Ax = b$. We use the basis representation $\mathbb{G}_h u = \sum_{j=1}^N y_j \phi_j$ of the Galerkin solution with some coefficient vector $y \in \mathbb{R}^n$. By use of the linearity of $\langle \cdot ; \cdot \rangle$, equation (1.15) becomes

$$b_k = F(\phi_k) = \langle \mathbb{G}_h u ; \phi_k \rangle = \sum_{j=1}^N y_j \langle \phi_j ; \phi_k \rangle = (Ay)_k \quad \text{for all } k = 1, \dots, N.$$

Therefore, the coefficient vector $y \in \mathbb{R}^N$ satisfies $Ay = b$. This proves $x = y$, i.e., we obtain $\mathbb{G}_h u$ by solving $Ax = b$. ■

Remark. We just remark that Theorem 1.4 can be applied for *any* orthogonal-type projection, e.g., the L^2 -orthogonal projection onto a discrete space. □

We now proceed with the abstract analysis of Galerkin schemes. The following two lemmata provide elementary properties of the Galerkin projection. The first lemma proves stability of the method with respect to changes of the right-hand side F .

v2:
06.10.2017

Lemma 1.5. *The Galerkin projection \mathbb{G}_h is a linear and continuous projection onto X_h with*

$$\|\mathbb{G}_h u\|_H \leq \frac{\beta}{\alpha} \|u\|_H \quad \text{for all } u \in H, \tag{1.19}$$

where $\alpha, \beta > 0$ are the norm equivalence constants from (1.12). Moreover, \mathbb{G}_h is the orthogonal projection onto X_h with respect to the energy scalar product $\langle \cdot ; \cdot \rangle$.

Proof. For $u_h \in X_h$, the Galerkin orthogonality (1.17) implies $\mathbb{G}_h u_h = u_h$. Therefore \mathbb{G}_h is a projection onto X_h . Also the linearity of \mathbb{G}_h follows from the Galerkin orthogonality (1.17). To see the continuity of \mathbb{G}_h , it remains to estimate the operator norm: For $u \in H$ holds

$$\|\mathbb{G}_h u\|^2 = \langle \mathbb{G}_h u ; \mathbb{G}_h u \rangle = \langle u ; \mathbb{G}_h u \rangle \leq \|u\| \|\mathbb{G}_h u\|,$$

whence $\|\mathbb{G}_h u\| \leq \|u\|$ and

$$\alpha \|\mathbb{G}_h u\|_H \leq \|\mathbb{G}_h u\| \leq \|u\| \leq \beta \|u\|_H,$$

where we have used the norm equivalence (1.12) on H as well as the Cauchy inequality for the scalar product $\langle \cdot ; \cdot \rangle$. This proves that $\|\mathbb{G}_h u\|_H \leq (\alpha/\beta) \|u\|_H$ and thus continuity of \mathbb{G}_h . Finally,

we remark that the *unique* orthogonal projection with respect to $\langle \cdot ; \cdot \rangle$, is characterized by the orthogonality relation (1.17). ■

The following Céa lemma states that the **Galerkin error** $\|u - \mathbb{G}_h u\|_H$ is quasi-optimal, i.e., it behaves like the best approximation error up to multiplicative constants, which depend only on the continuous setting but not on X_h .

Lemma 1.6 (Céa). *The Galerkin error is quasi-optimal, i.e.,*

$$\|u - \mathbb{G}_h u\|_H \leq \frac{\beta}{\alpha} \min_{v_h \in X_h} \|u - v_h\|_H \quad \text{for all } u \in H, \quad (1.20)$$

where $\alpha, \beta > 0$ are the norm equivalence constants from (1.12). With respect to the energy norm, it holds that

$$\|u - \mathbb{G}_h u\| = \min_{v_h \in X_h} \|u - v_h\| \quad \text{for all } u \in H, \quad (1.21)$$

i.e., the Galerkin solution $\mathbb{G}_h u$ is the best approximation of u with respect to the energy norm.

Proof. For arbitrary $v_h \in X_h$, the Galerkin orthogonality (1.17) proves that

$$\|u - \mathbb{G}_h u\|^2 = \langle u - \mathbb{G}_h u ; u - v_h \rangle \leq \|u - \mathbb{G}_h u\| \|u - v_h\|,$$

which yields (1.21) with an infimum on the right-hand side. Of course, the minimum in (1.21) is attained for $v_h = \mathbb{G}_h u$. With the same arguments as in the proof of the last lemma, we even see that

$$\alpha \|u - \mathbb{G}_h u\|_H \leq \|u - \mathbb{G}_h u\| \leq \|u - v_h\| \leq \beta \|u - v_h\|_H,$$

which implies (1.20) with an infimum on the right-hand side. This minimum is attained for $v_h = \Pi_h u$ with $\Pi_h : X \rightarrow X_h$ being the orthogonal projection onto X_h with respect to $\|\cdot\|_H$. ■

Exercise 2. Let X be a normed space and $X_h \subseteq X$ be a finite dimensional subspace of X . Then, for any $x \in X$, there exists some (not necessarily unique) $x_h \in X_h$ such that

$$\|x - x_h\|_X = \min_{v_h \in X_h} \|x - v_h\|_X,$$

i.e., best approximation errors on finite dimensional spaces as in (1.20) are always attained. Prove that the set of minimizers is convex, closed and bounded (and hence even compact). □

A major advantage of Galerkin methods is that one can prove convergence for any exact solution $u \in H$ if one knows that smooth functions can be approximated well. In the following, think of the subscript $h > 0$ as a mesh-size parameter with corresponding finite dimensional spaces X_h :

Proposition 1.7. *For all $h > 0$, let X_h be a finite-dimensional subspace of H . We assume that there is a dense subspace D of H with approximation property, namely*

$$\lim_{h \rightarrow 0} \min_{v_h \in X_h} \|v - v_h\|_H = 0 \quad \text{for all } v \in D. \quad (1.22)$$

Then, for any $u \in H$, it holds that

$$\lim_{h \rightarrow 0} \|u - \mathbb{G}_h u\|_H = 0, \quad (1.23)$$

i.e., the sequence of Galerkin solutions converges to the exact solution u .

Proof. For $v \in D$, the quasi-optimality (1.20) yields that

$$\|u - \mathbb{G}_h u\|_H \leq \frac{\beta}{\alpha} \min_{v_h \in X_h} \|u - v_h\|_H \leq \frac{\beta}{\alpha} (\|u - v\|_H + \min_{v_h \in X_h} \|v - v_h\|_H).$$

We have to show that

$$\exists C > 0 \forall \varepsilon > 0 \exists h_0 > 0 \forall h \in (0, h_0) \quad \|u - \mathbb{G}_h u\|_H \leq C \varepsilon.$$

For $\varepsilon > 0$, let $v \in D$ with $\|u - v\|_H \leq \varepsilon$. Choose $h_0 > 0$ according to the approximation assumption (1.23) so that $\min_{v_h \in X_h} \|v - v_h\|_H \leq \varepsilon$ for all $h \in (0, h_0)$. We thus finally obtain $\|u - \mathbb{G}_h u\|_H \leq 2\beta\varepsilon/\alpha$, which concludes the proof. \blacksquare

Although the result of the preceding lemma seems to be very attractive, we stress, however, that the convergence of a Galerkin scheme can be arbitrarily slow. We argue in the abstract setting: If H is a separable Hilbert space, e.g., H is a Sobolev space, there is a countable orthonormal basis $\{\phi_j \mid j \in \mathbb{N}\}$. Any $u \in H$ can be written as $u = \sum_{j=1}^{\infty} x_j \phi_j$ with coefficients $(x_n) \in \ell_2$. If we define $X_j := \text{span}\{\phi_1, \dots, \phi_j\}$, it holds that

$$\min_{v_h \in X_h} \|u - v_h\|_H^2 = \sum_{j=k+1}^{\infty} x_j^2.$$

Finally, the decay of the right-hand side can be very slow. One may think of, e.g., $x_j^2 = j^{-(1+\varepsilon)}$ for any $\varepsilon > 0$, so that the series converges but is — in the beginning — almost the divergent harmonic series.

The following exercise shows that the approximation property (1.22) in particular implies that the Hilbert space H has to be separable.

Exercise 3. Suppose that X is a normed space with finite dimensional subspaces $X_\ell \subseteq X_{\ell+1} \subseteq X$ for all $\ell \in \mathbb{N}$. Suppose that $\mathcal{D} \subseteq X$ is a dense subspace such that, for all $x \in X$,

$$\lim_{\ell \rightarrow \infty} \min_{x_\ell \in X_\ell} \|x - x_\ell\|_X = 0. \quad (1.24)$$

Then, X is separable, i.e., there is a countable and dense subset $M \subseteq X$. \square

Exercise 4. Let $X = \ell_\infty$ and $X_\ell := \{(x_n) \in \ell_\infty \mid x_j = 0 \text{ for all } j \geq \ell\}$. Prove that (1.24) fails to hold for any dense subspace \mathcal{D} . Note that this also follows if one proves that ℓ_∞ is not separable. \square

Remark. All foregoing results of this section hold (in a slightly modified form) in case that $\langle \cdot ; \cdot \rangle$ only is a continuous and elliptic bilinear form on the Hilbert space H , i.e., in all proofs, one can avoid to use the symmetry of $\langle \cdot ; \cdot \rangle$. \square

The following exercise explains why $\|\cdot\|$ is called energy norm. In many situations, the function $J(\cdot)$ has the interpretation of a physical energy.

Exercise 5. Let $\langle\langle \cdot ; \cdot \rangle\rangle$ be a scalar product on the Hilbert space H such that the norm $\|\cdot\|$ is equivalent to $\|\cdot\|_H$. Let $F \in H^*$ and $u \in H$. Then, the following assertions are equivalent:

- $\langle\langle u ; \cdot \rangle\rangle = F \in H^*$;
- $J(u) = \min_{v \in H} J(v)$, where $J(v) := \frac{1}{2} \langle\langle v ; v \rangle\rangle - F(v)$.

In particular, the variational formulation is equivalent to energy minimization, and this result also covers the discrete setting. Derive a formula for the energy error $J(\mathbb{G}_h u) - J(u)$, where $\mathbb{G}_h : H \rightarrow X_h$ denotes the Galerkin projection. \square

Finally, we comment on an extension of the concept of Galerkin schemes to some nonlinear problems. We note that this framework does, in particular, cover the frame of the Lax–Milgram lemma.

Exercise 6 (Main Theorem on Strongly Monotone Operators (Zarantonello '60)). Let H be a Hilbert space and $A : H \rightarrow H^*$ be a Lipschitz continuous and strongly monotone operator, i.e.,

$$\|Au - Av\|_{H^*} \leq L\|u - v\|_H \quad \text{and} \quad \langle Au - Av ; u - v \rangle_{H^* \times H} \geq M\|u - v\|_H^2 \quad \text{for all } u, v \in H$$

with constants $L, M > 0$ that only depend on A . Then, A is bijective. **Hint:** Injectivity of A follows from the monotonicity of A . To prove surjectivity, we apply a fixed point argument: Let $I_H : H \rightarrow H^*$, $I_H(u) := (u ; \cdot)_H$ denote the Riesz mapping. For given $F \in H^*$ and a certain choice of $C > 0$, the mapping $\Phi(u) := u - CI_H^{-1}(Au - F)$ is a contraction on H . Therefore, the Banach contraction theorem applies and provides a unique $u \in H$ with $u = \Phi(u)$. \square

Exercise 7 (Lemma of Lax–Milgram). Use Exercise 6 to derive the Lemma of Lax–Milgram: Let H be a Hilbert space and $a(\cdot, \cdot)$ be a continuous and elliptic bilinear form on H , i.e.,

$$a(u, v) \leq L\|u\|_H\|v\|_H \quad \text{and} \quad a(u, u) \geq M\|u\|_H^2 \quad \text{for all } u, v \in H,$$

where the constants $L, M > 0$ depend only on $a(\cdot, \cdot)$. Then, given a right-hand side $F \in H^*$, there is a unique $u \in H$ with $a(u, \cdot) = F \in H^*$. \square

Exercise 8. Define the Galerkin method in the context of monotone operators: Under the assumptions of Exercise 6, we aim to approximate the solution $u \in H$ of $Au = F \in H^*$. How does the Galerkin method look like in this setting? Prove that the Galerkin operator $\mathbb{G}_h : H \rightarrow X_h$ onto a finite dimensional subspace $X_h \subset H$ is a well-defined (in general nonlinear) and Lipschitz-continuous projection, i.e., $\mathbb{G}_h^2 = \mathbb{G}_h$ with

$$\|\mathbb{G}_h u - \mathbb{G}_h v\|_H \leq C\|u - v\|_H \quad \text{for all } u, v \in H.$$

Céa lemma

$$\|u - \mathbb{G}_h u\|_H \leq C \min_{v_h \in X_h} \|u - v_h\|_H \quad \text{for all } u \in H.$$

Show that the constants $C > 0$ depend only on A . □

Exercise 9. We stick with the setting of monotone operators from Exercise 7 and 8: How can one compute the Galerkin approximation $u_h = \mathbb{G}_h u \in X_h$ of a solution $u \in H$ of $Au = F \in H^*$? For $N = \dim X_h$, provide a (nonlinear) system of equations in \mathbb{R}^N which characterizes the unique solution $u_h = \mathbb{G}_h u \in X_h$. What happens if the operator A is linear? □

Chapter 2

Sobolev Spaces and Poisson Problem

2.1 Sobolev Spaces on Domains (10.10.2017)

This section briefly recalls the definition of Sobolev spaces $H^m(\Omega)$, for integer order $m \in \mathbb{N}_0$, on domains $\Omega \subseteq \mathbb{R}^d$. While this section requires Ω only to be open and connected, the following sections will implicitly assume that Ω is a bounded Lipschitz domain.

Definition. A function $u \in L^1_{loc}(\Omega) := \{w : \Omega \rightarrow \mathbb{R} \text{ measurable} \mid \forall K \subset \Omega \text{ compact } w \in L^1(K)\}$ has a **weak partial derivative** $\partial_j u \in L^1_{loc}(\Omega)$, if the pair $(u, \partial_j u)$ satisfies the integration by parts formula with smooth test functions that vanish on the boundary, i.e., it holds that

$$\int_{\Omega} u(\partial_j v) dx = - \int_{\Omega} (\partial_j u)v dx \quad \text{for all } v \in \mathcal{D}(\Omega) := C_c^\infty(\Omega). \quad (2.1)$$

Note that $\partial_j u$ is (so far) only a symbol, whereas $\partial_j v := \partial v / \partial x_j$ is the classical j -th derivative of $v \in \mathcal{D}(\Omega)$. We say that $u \in L^1_{loc}(\Omega)$ is **weakly differentiable with weak gradient** $\nabla u \in L^1_{loc}(\Omega)$, if all weak derivatives $\partial_j u$, for $j = 1, \dots, d$, exist. \square

From the main theorem of calculus, we infer that the weak derivative is unique, if it exists. Moreover, the weak derivative and the classical derivative coincide, if the classical derivative exists.

Theorem 2.1 (Fundamental Theorem of Calculus of Variations). *Let $f \in L^1_{loc}(\Omega)$ satisfy $\int_{\Omega} f v dx = 0$ for all $v \in \mathcal{D}(\Omega)$. Then, it holds that $f = 0$ almost everywhere in Ω . \blacksquare*

Remark. Note that $C(\Omega) \subset L^1_{loc}(\Omega)$. For $f \in C(\Omega)$, the fundamental theorem of calculus of variations can be proven by elementary calculus: Note that for any $x \in \mathbb{R}^d$ and any radius $\varepsilon > 0$, there is a function $\psi \in \mathcal{D}(\mathbb{R}^d)$ such that $\{y \in \mathbb{R}^d \mid \psi(y) > 0\} = U(x, \varepsilon) := \{y \in \mathbb{R}^d \mid |x - y| < \varepsilon\}$; see the following Exercise 10. Provided $f \in C(\Omega)$ with $f(x) \neq 0$ for some $x \in \Omega$, we may assume $f(x) > 0$. By continuity, there is a small radius $\varepsilon > 0$ such that $U(x, \varepsilon) \subset \Omega$ and that $f(y) > 0$ for all $y \in U(x, \varepsilon)$. With the associated function $\psi \in \mathcal{D}(\Omega)$, we thus see that $\int_{\Omega} f \psi dx > 0$. Note that this argument provides the (logically equivalent) contraposition of the fundamental theorem of calculus of variations in the case of a continuous function f . \square

Exercise 10. (i) Show that the following definition provides $\phi \in C^\infty(\mathbb{R})$ with $\text{supp}(\phi) = [-1, 1]$:

$$\phi(t) := \begin{cases} \exp(-1/(1-t^2)), & \text{for } |t| < 1, \\ 0 & \text{else.} \end{cases}$$

(ii) For $\varepsilon > 0$ and $x \in \mathbb{R}^d$, define the function $\psi_{x,\varepsilon}(y) := \phi(|x-y|^2/\varepsilon)$. Show that $\psi_{x,\varepsilon} \in C^\infty(\mathbb{R}^d)$ with $\text{supp}(\psi_{x,\varepsilon}) = \{y \in \mathbb{R}^d \mid |x-y| \leq \varepsilon\}$ and $\psi_{x,\varepsilon}(y) > 0$ for all $y \in \{y \in \mathbb{R}^d \mid |x-y| < \varepsilon\}$. \square

Corollary 2.2. (i) *The weak derivative $\partial_j u$ is unique, if it exists: If $\partial_j u, \widetilde{\partial_j u} \in L^1_{loc}(\Omega)$ satisfy (2.1), it holds that $\partial_j u = \widetilde{\partial_j u}$ almost everywhere in Ω .*

(ii) *A function $u \in C^1(\Omega)$ is weakly differentiable, and the weak derivative coincides with the classical derivative.*

Proof. (i) It holds that $\int_\Omega (\partial_j u - \widetilde{\partial_j u})v \, dx = 0$ for all $v \in \mathcal{D}(\Omega)$ and thus $\partial_j u - \widetilde{\partial_j u} = 0$ almost everywhere in Ω . (ii) follows from (i) and the integration by parts formula. \blacksquare

A deeper result is the following, which is somehow, nevertheless, quite natural and expected.

Theorem 2.3. *If $u \in L^1_{loc}(\Omega)$ is weakly differentiable with $\nabla u = 0$, then the function u is constant, i.e., there is a constant $c \in \mathbb{R}$ such that $u = c$ almost everywhere in Ω .* \blacksquare

Definition. For $m = 0$, we define $H^0(\Omega) := L^2(\Omega)$ as the classical Lebesgue space of square integrable functions. For $m = 1$, the **Sobolev space** $H^1(\Omega)$ is defined by

$$H^1(\Omega) := \{u \in L^2(\Omega) \mid u \text{ weakly differentiable, } \nabla u \in L^2(\Omega)\} \quad (2.2)$$

and associated with the graph norm

$$\|u\|_{H^1(\Omega)} := (\|u\|_{L^2(\Omega)}^2 + \|\nabla u\|_{L^2(\Omega)}^2)^{1/2}. \quad (2.3)$$

Higher-order Sobolev spaces of integer order $m \in \mathbb{N}$ may be defined inductively by

$$H^m(\Omega) := \{u \in L^2(\Omega) \mid u \text{ weakly differentiable, } \nabla u \in H^{m-1}(\Omega)\}, \quad (2.4)$$

with associated norm

$$\|u\|_{H^m(\Omega)} := (\|u\|_{L^2(\Omega)}^2 + \|\nabla u\|_{H^{m-1}(\Omega)}^2)^{1/2}. \quad (2.5)$$

Remark. Clearly, $C^1(\overline{\Omega}) \subseteq H^1(\Omega)$ and we note below that $C^1(\overline{\Omega})$ is even dense in $H^1(\overline{\Omega})$. \square

Theorem 2.4. *For all $m \in \mathbb{N}_0$, the Sobolev space $H^m(\Omega)$ is a Hilbert space.*

Proof. The proof uses the (hopefully) well-known fact that $H^0(\Omega) = L^2(\Omega)$ is a Hilbert space. We shall proceed by induction on m . However, we explicitly consider the case $m = 1$ first: Obviously, the H^1 -norm is induced by the scalar product

$$(u ; v)_{H^1(\Omega)} := (u ; v)_{L^2(\Omega)} + (\nabla u ; \nabla v)_{L^2(\Omega)} \quad \text{for all } u, v \in H^1(\Omega),$$

i.e., $\|u\|_{H^1(\Omega)}^2 = (u ; u)_{H^1(\Omega)}$. Therefore, it only remains to prove the completeness of $H^1(\Omega)$. Let (u_n) be a Cauchy sequence in $H^1(\Omega)$. Note that, by definition of the H^1 -norm, (u_n) as well as (∇u_n) are Cauchy sequences in $L^2(\Omega)$. Since $L^2(\Omega)$ is complete, there are unique $u \in L^2(\Omega)$ and $g \in L^2(\Omega)^d$ such that

$$\lim_{n \rightarrow \infty} \|u - u_n\|_{L^2(\Omega)} = 0 = \lim_{n \rightarrow \infty} \|g - \nabla u_n\|_{L^2(\Omega)}.$$

By definition of $H^1(\Omega)$, it thus only remains to prove that u is weakly differentiable with gradient $\nabla u = g$. Let $v \in \mathcal{D}(\Omega)$ be an arbitrary test function. From the weak differentiability of each u_n and L^2 -convergence, we obtain that

$$(u ; \partial_j v)_{L^2(\Omega)} = \lim_{n \rightarrow \infty} (u_n ; \partial_j v)_{L^2(\Omega)} = - \lim_{n \rightarrow \infty} (\partial_j u_n ; v)_{L^2(\Omega)} = -(g_j ; v)_{L^2(\Omega)}.$$

Therefore, g_j is the j -th weak derivative of u and consequently $g = \nabla u$. This concludes the case $m = 1$. The induction step for $H^m(\Omega)$ is left to the reader, but obviously follows from the same arguments, where we replace $g \in L^2(\Omega)^d$ by $g \in H^{m-1}(\Omega)^d$. ■

2.2 Main Theorems on Sobolev Spaces (10.10.2017)

From now on, it will be important and thus assumed that $\Omega \subset \mathbb{R}^d$ is a bounded Lipschitz domain. By definition of the Sobolev spaces $H^m(\Omega)$, there holds $H^m(\Omega) \subset H^{m-1}(\Omega)$ with $\|u\|_{H^{m-1}(\Omega)} \leq \|u\|_{H^m(\Omega)}$. In other words, the identity operator $id : H^m(\Omega) \rightarrow H^{m-1}(\Omega)$ is well-defined and continuous. The following Rellich theorem states that it is also compact. This is a pretty strong result. The impact of which will become clear in our proofs of the Poincaré inequality and the Friedrichs inequality.

Theorem 2.5 (Rellich Compactness Theorem). *For any integer order $m \in \mathbb{N}$, the embedding $H^m(\Omega) \subseteq H^{m-1}(\Omega)$ is compact.* ■

We recall that an operator $A \in L(X; Y)$ between normed spaces X and Y is compact, if each bounded sequence (x_n) in X , i.e., $\sup_{n \in \mathbb{N}} \|x_n\|_X < \infty$, allows for a subsequence (x_{n_k}) such that the image (Ax_{n_k}) is convergent in Y , i.e., there is an element $y \in Y$ such that $\lim_{k \rightarrow \infty} \|Ax_{n_k} - y\|_Y = 0$. The next lemma states that compact operators turn weakly convergent sequences to strongly convergent sequences.

Lemma 2.6. *Suppose that $A \in L(X; Y)$ is a compact operator between normed spaces X and Y and that (x_n) is a weakly convergent sequence, i.e., $x_n \rightharpoonup x \in X$. Then, the image (Ax_n) is strongly convergent to Ax in Y , i.e., $Ax_n \rightarrow Ax \in Y$.*

Proof by contradiction. Using the adjoint operator $A^* \in L(Y^*; X^*)$, one sees that $Ax_n \rightharpoonup Ax \in Y$. Assume that (Ax_n) does not strongly converge to Ax . Then, there is a subsequence (Ax_{n_k}) with $\inf_{k \in \mathbb{N}} \|Ax_{n_k} - Ax\|_Y \geq \varepsilon$ for some $\varepsilon > 0$. Recall that weakly convergent sequence are always bounded. Compactness thus provides a further subsequence $(Ax_{n_{k_\ell}})$ of (Ax_{n_k}) with $Ax_{n_{k_\ell}} \rightarrow y \in Y$. In particular, $Ax_{n_{k_\ell}} \rightharpoonup y \in Y$ and therefore $y = Ax$. This contradicts the choice of the subsequence (Ax_{n_k}) . \blacksquare

Exercise 11. Let X be a reflexive Banach space and Y be a normed space. Suppose that $A \in L(X, Y)$ is completely continuous, i.e., for all (x_n) in X , weak convergence $x_n \rightharpoonup x$ in X implies strong convergence $Ax_n \rightarrow Ax$ in Y . Prove that A is compact, i.e., for X being reflexive, the operator A is compact if and only if it is completely continuous. \square

Before the statement and the proof of the Poincaré inequality, we need a further technical lemma. The result is rather standard in the analysis of variational problems.

Lemma 2.7. A continuous and convex functional $f : X \rightarrow \mathbb{R}$ on a normed space X is weakly lower semicontinuous, i.e., for each weakly convergent sequence (x_n) in X with $x_n \rightharpoonup x \in X$, it holds that

$$f(x) \leq \liminf_{n \in \mathbb{N}} f(x_n). \tag{2.6}$$

Proof. 1. step. We prove that the epigraph $G := \{(x, \alpha) \in X \times \mathbb{R} \mid f(x) \leq \alpha\}$ is convex: For $(x, \alpha), (y, \beta) \in G$ and $0 \leq \theta \leq 1$, the convexity of f proves that

$$f(\theta x + (1 - \theta)y) \leq \theta f(x) + (1 - \theta)f(y) \leq \theta\alpha + (1 - \theta)\beta,$$

whence $\theta(x, \alpha) + (1 - \theta)(y, \beta) \in G$, i.e., $G \subseteq X \times \mathbb{R}$ is convex.

2. step. We use the continuity of f to prove that G is also closed: Let (x_n, α_n) be a convergent sequence in G , i.e., it holds that $x_n \rightarrow x \in X$ and $\alpha_n \rightarrow \alpha \in \mathbb{R}$. We prove that $(x, \alpha) \in G$, which follows from

$$f(x) = \lim_{n \rightarrow \infty} f(x_n) \leq \lim_{n \rightarrow \infty} \alpha_n = \alpha.$$

3. step. The following step in the proof is known as *Mazur's lemma*: We prove that the closed and convex set G is also weakly closed in $X \times \mathbb{R} =: Y$, i.e., closed with respect to the weak topology on Y . We argue by contradiction and assume that G is not weakly closed. Then, there is an element $y \in \overline{G}^\sigma \setminus G$, where \overline{G}^σ denotes the weak closure of G . According to the Hahn-Banach separation theorem, there is a functional $\phi \in Y^*$ and a scalar $\lambda \in \mathbb{R}$ such that $\phi(y) < \lambda \leq \inf \phi(G)$. Therefore $U := \phi^{-1}(-\infty, \lambda)$ is weakly open with $y \in U$ and $U \cap G = \emptyset$. This contradicts topologically that y is in the weak closure of G . Hence, $G = \overline{G}^\sigma$ is weakly closed, and we may proceed with the proof of (2.6).

4. step. We show the weak lower semicontinuity of f : Suppose that $x_n \rightharpoonup x \in X$. Considering a subsequence, we may assume without loss of generality that $\alpha := \liminf_n f(x_n) = \lim_n f(x_n)$. For $\alpha = \infty$, (2.6) is trivial. We thus may assume $\alpha < \infty$. Let $\beta > \alpha$ and define $\alpha_n := \max\{\beta, f(x_n)\} \rightarrow$

β . Clearly, $(x_n, \alpha_n) \in G$. Moreover, this sequence is weakly convergent $(x_n, \alpha_n) \rightharpoonup (x, \beta)$. We deduce $(x, \beta) \in G$. Thus, $f(x) \leq \beta$ for all $\beta > \alpha$ and therefore finally $f(x) \leq \alpha = \lim_{n \rightarrow \infty} f(x_n)$. ■

A first consequence of the preceding abstract results is that one can easily construct equivalent norms on the Sobolev space $H^1(\Omega)$.

v3:
11.10.2017

Proposition 2.8. *Let $|\cdot|_{H^1}$ be a continuous seminorm on $H^1(\Omega)$ which is definite on the constant functions, i.e., $|c|_{H^1} = 0$ implies $c = 0$ for all $c \in \mathbb{R}$. Then, there are constants $C_1, C_2 > 0$ such that*

$$|v|_{H^1} \leq C_1 \|v\|_{H^1(\Omega)} \quad \text{as well as} \quad C_2^{-1} \|v\|_{L^2(\Omega)} \leq \|v\| := \|\nabla v\|_{L^2(\Omega)} + |v|_{H^1} \quad \text{for all } v \in H^1(\Omega).$$

In particular, $\| \cdot \|$ defines an equivalent norm on $H^1(\Omega)$, i.e.,

$$(1 + C_1)^{-1} \|v\| \leq \|v\|_{H^1(\Omega)} \leq (1 + C_2) \|v\| \quad \text{for all } v \in H^1(\Omega).$$

Proof. 1. step. Existence of C_1 : We argue by contradiction and assume that there is no constant $C_1 > 0$ such that $|v|_{H^1} \leq C_1 \|v\|_{H^1(\Omega)}$ for all $v \in H^1(\Omega)$. Consequently, there exists a sequence (v_n) in $H^1(\Omega)$ with $|v_n|_{H^1} > n \|v_n\|_{H^1(\Omega)}$. We may define $w_n := v_n/|v_n|_{H^1} \in H^1(\Omega)$ and obtain that $\|w_n\|_{H^1(\Omega)} < 1/n$. Therefore, (w_n) converges to zero in $H^1(\Omega)$ and thus $\lim_{n \rightarrow \infty} |w_n|_{H^1} = |0|_{H^1} = 0$ according to the continuity of $|\cdot|_{H^1}$. However, this contradicts $|w_n|_{H^1} = 1$ which holds by definition of w_n . This concludes the existence of C_1 . In particular, we hence observe $\|v\| \leq (1 + C_1) \|v\|_{H^1(\Omega)}$.

2. step. Existence of C_2 : We assume that there is no constant $C_2 > 0$ such that $\|v\|_{L^2(\Omega)} \leq C_2 \|v\|$ for all $v \in H^1(\Omega)$. Therefore, there exists a sequence (v_n) in $H^1(\Omega)$ such that

$$\frac{1}{n} \|v_n\|_{L^2(\Omega)} > \|v_n\| = \|\nabla v_n\|_{L^2(\Omega)} + |v_n|_{H^1}$$

The definition of $w_n := v_n/\|v_n\|_{L^2(\Omega)}$ leads to a sequence (w_n) in $H^1(\Omega)$ such that

$$\|w_n\|_{L^2(\Omega)} = 1, \quad \|\nabla w_n\|_{L^2(\Omega)} \leq 1/n, \quad |w_n|_{H^1} \leq 1/n.$$

Therefore, (w_n) is a bounded sequence in the Hilbert space $H^1(\Omega)$. A Hilbert space is reflexive. By virtue of the Eberlein-Šmulian theorem, each bounded sequence thus has a weakly convergent subsequence. Therefore, we may assume that $w_n \rightharpoonup w \in H^1(\Omega)$. An application of Lemma 2.7 proves that

$$\|\nabla w\|_{L^2(\Omega)} \leq \liminf_{n \rightarrow \infty} \|\nabla w_n\|_{L^2(\Omega)} = 0,$$

whence the weak limit w is constant. Another application of Lemma 2.7 proves that

$$|w|_{H^1} \leq \liminf_{n \rightarrow \infty} |w_n|_{H^1} = 0$$

since a seminorm is always convex. Therefore, $w = 0$. On the other hand, the Rellich theorem states the strong convergence $w_n \rightarrow w \in L^2(\Omega)$ and thus $\|w\|_{L^2(\Omega)} = \lim_{n \rightarrow \infty} \|w_n\|_{L^2(\Omega)} = 1$. This contradiction concludes the existence of C_2 . In particular, we hence observe $\|v\|_{H^1(\Omega)} \leq \|v\|_{L^2(\Omega)} + \|\nabla v\|_{L^2(\Omega)} \leq (C_2 + 1) \|v\|$. ■

Corollary 2.9 (Poincaré Inequality). *It holds that*

$$\|v\|_{L^2(\Omega)} \leq \tilde{C}_P \left(\|\nabla v\|_{L^2(\Omega)} + \left| \int_{\Omega} v \, dx \right| \right) \quad \text{for all } v \in H^1(\Omega), \quad (2.7)$$

where the constant $\tilde{C}_P > 0$ depends only on Ω . Moreover, $\|v\| := \|\nabla v\|_{L^2(\Omega)} + \left| \int_{\Omega} v \, dx \right|$ defines even an equivalent norm on $H^1(\Omega)$.

Proof. According to Proposition 2.8, it only remains to show that

$$|v|_{H^1} := \left| \int_{\Omega} v \, dx \right| \quad \text{for } v \in H^1(\Omega)$$

defines a continuous seminorm on $H^1(\Omega)$ which is definite on the constant functions. The equality $|c|_{H^1} = |\Omega||c|$ for $c \in \mathbb{R}$ verifies the definiteness. Lipschitz continuity follows from

$$\left| |v|_{H^1} - |w|_{H^1} \right| \leq \left| \int_{\Omega} v - w \, dx \right| \leq |\Omega|^{1/2} \|v - w\|_{L^2(\Omega)} \leq |\Omega|^{1/2} \|v - w\|_{H^1(\Omega)}$$

and from the boundedness of Ω . ■

Corollary 2.10 (Poincaré Inequality). *There is a constant $C_P > 0$, which depends only on the shape of Ω but not on its diameter, such that*

$$\|v\|_{L^2(\Omega)} \leq C_P \operatorname{diam}(\Omega) \|\nabla v\|_{L^2(\Omega)} \quad \text{for all } v \in H_*^1(\Omega) := \{w \in H^1(\Omega) \mid \int_{\Omega} w \, dx = 0\}, \quad (2.8)$$

where $\operatorname{diam}(\Omega) := \sup \{|x - y| \mid x, y \in \Omega\}$ denotes the diameter of Ω .

Proof. The proof is a so-called **scaling argument**: We define $\lambda := \operatorname{diam}(\Omega)$ and $\tilde{\Omega} := \lambda^{-1}\Omega$. Note that the scaled domain $\tilde{\Omega}$ satisfies $\operatorname{diam}(\tilde{\Omega}) = 1$ and depends only on the shape of Ω . We consider the affine bijection $\Phi : \Omega \rightarrow \tilde{\Omega}$, $\Phi(x) := \lambda^{-1}x$. Recall the transformation theorem, which holds for arbitrary diffeomorphisms $\Phi : \Omega \rightarrow \tilde{\Omega}$ and states that

$$\int_{\tilde{\Omega}} \tilde{f} \, dy = \int_{\Omega} \tilde{f}(\Phi(x)) |\det D\Phi(x)| \, dx \quad \text{for all } \tilde{f} \in L^1(\tilde{\Omega}).$$

Note that $\det D\Phi(x) = \lambda^{-d}$ since $\Phi = \lambda^{-1}\mathbf{I}$ in our case. For $v \in H^1(\Omega)$, we define $\tilde{v} := v \circ \Phi^{-1} \in H^1(\tilde{\Omega})$. Then,

$$\|\tilde{v}\|_{L^2(\tilde{\Omega})}^2 = \int_{\tilde{\Omega}} |\tilde{v}|^2 \, dy = \lambda^{-d} \int_{\Omega} |v|^2 \, dx = \lambda^{-d} \|v\|_{L^2(\Omega)}^2.$$

According to the chain rule, it holds that $\nabla \tilde{v} = \lambda (\nabla v) \circ \Phi^{-1}$ and consequently that

$$\|\nabla \tilde{v}\|_{L^2(\tilde{\Omega})}^2 = \lambda^{2-d} \|\nabla v\|_{L^2(\Omega)}^2.$$

With $\tilde{C}_P > 0$ the Poincaré constant from (2.7) for $\tilde{\Omega}$, we thus infer

$$\|v\|_{L^2(\Omega)}^2 = \lambda^d \|\tilde{v}\|_{L^2(\tilde{\Omega})}^2 \leq \lambda^d \tilde{C}_P^2 \|\nabla \tilde{v}\|_{L^2(\tilde{\Omega})}^2 = \lambda^2 \tilde{C}_P^2 \|\nabla v\|_{L^2(\Omega)}^2.$$

Note that \tilde{C}_P depends only on $\tilde{\Omega}$ and thus only on the shape of Ω . This concludes the proof. ■

Remark. We stress that $Iv := \int_{\Omega} v \, dx$ defines a linear and continuous functional on $H^1(\Omega)$. In particular, $H_*^1(\Omega) = \ker(I)$ is a closed subspace of $H^1(\Omega)$ and hence a Hilbert space. According to the Poincaré inequality, it holds that $\|\nabla v\|_{L^2(\Omega)} \leq \|v\|_{H^1(\Omega)} \leq (1 + \tilde{C}_P^2)^{1/2} \|\nabla v\|_{L^2(\Omega)}$ for all $v \in H_*^1(\Omega)$. In particular, $\|\nabla v\|_{L^2(\Omega)}$ defines an equivalent Hilbert norm on $H_*^1(\Omega)$ with associated scalar product $(\nabla u ; \nabla v)_{L^2(\Omega)}$. □

Theorem 2.11 (Meyers-Serrin). For each integer order $m \in \mathbb{N}$, $C^\infty(\overline{\Omega})$ and, in particular, $C^\infty(\Omega) \cap H^m(\Omega)$ are dense subspaces of $H^m(\Omega)$. ■

Theorem 2.12 (Trace Operator). There is a unique operator $\gamma \in L(H^1(\Omega); L^2(\Gamma))$ such that $\gamma v = v|_{\Gamma}$ for all $v \in C^1(\overline{\Omega})$, i.e., γ extends the classical trace defined as restriction $v|_{\Gamma}$ on the boundary for smooth functions v . ■

As a first corollary to Theorem 2.12, we can prove that the integration by parts formula also holds for Sobolev functions $u, v \in H^1(\Omega)$.

Corollary 2.13 (Integration by Parts). For all $u, v \in H^1(\Omega)$, it holds that

$$\int_{\Omega} u \frac{\partial v}{\partial x_j} \, dx + \int_{\Omega} \frac{\partial u}{\partial x_j} v \, dx = \int_{\Gamma} \gamma u \gamma v n_j \, ds. \quad (2.9)$$

Proof. The formula (2.9) holds for $u, v \in C^1(\overline{\Omega})$. All three terms define continuous bilinear forms on $H^1(\Omega) \times H^1(\Omega)$. Therefore (2.9) follows, for arbitrary $u, v \in H^1(\Omega)$ from the density of $C^1(\overline{\Omega})$ in $H^1(\Omega)$: Given $u, v \in H^1(\Omega)$, there are sequences (u_n) and (v_n) in $C^1(\overline{\Omega})$ which converge to u resp. v in $H^1(\Omega)$. Therefore, if $a(\cdot, \cdot) : H^1(\Omega) \times H^1(\Omega) \rightarrow \mathbb{R}$ is continuous, then it holds that $\lim_{n \rightarrow \infty} a(u_n, v_n) = a(u, v)$. This concludes the proof. ■

The analytical treatment of the Dirichlet problem makes use of the so-called Friedrichs inequality, whereas the analytical treatment of the Neumann problem uses the previously proven Poincaré inequality.

Corollary 2.14 (Friedrichs Inequality). Assume that the Dirichlet boundary $\Gamma_D \subseteq \Gamma$ has positive surface measure $|\Gamma_D| > 0$. Then, it holds that

$$\|v\|_{L^2(\Omega)} \leq \tilde{C}_F (\|\nabla v\|_{L^2(\Omega)} + \|\gamma v\|_{L^2(\Gamma_D)}) \quad \text{for all } v \in H^1(\Omega) \quad (2.10)$$

with a constant $\tilde{C}_F > 0$, which depends only on Ω and Γ_D . Moreover, the right-hand side $\|v\| := \|\nabla v\|_{L^2(\Omega)} + \|\gamma v\|_{L^2(\Gamma_D)}$ even defines an equivalent norm on $H^1(\Omega)$.

Proof. We again apply Proposition 2.8. It only remains to show that

$$\|v\|_{H^1} := \|\gamma v\|_{L^2(\Gamma_D)} \quad \text{for } v \in H^1(\Omega)$$

defines a continuous seminorm on $H^1(\Omega)$ which is definite on the constant functions. The definiteness is again easily obtained from $|c|_{H^1} = |\Gamma_D|^{1/2}|c|$ for $c \in \mathbb{R}$. Lipschitz continuity follows from

$$\left| \|v\|_{H^1} - \|w\|_{H^1} \right| \leq \|\gamma v - \gamma w\|_{L^2(\Gamma_D)} = \|\gamma(v-w)\|_{L^2(\Gamma_D)} \leq C \|v-w\|_{H^1(\Omega)}$$

according to the continuity of the trace operator $\gamma \in L(H^1(\Omega); L^2(\Gamma))$. ■

Definition. We define $H_0^1(\Omega) := \overline{\mathcal{D}(\Omega)}^{\|\cdot\|_{H^1}}$ and $H_D^1(\Omega) := \overline{C_D^1(\overline{\Omega})}^{\|\cdot\|_{H^1}}$, where the subscript D indicates the Dirichlet boundary Γ_D . By definition, $H_0^1(\Omega)$ as well as $H_D^1(\Omega)$ are closed subspaces of $H^1(\Omega)$ and thus Hilbert spaces. In particular, it holds that $H_0^1(\Omega) \subseteq H_D^1(\Omega)$. □

The same scaling argument as for the Poincaré inequality proves the following variant of the Friedrichs inequality, where we note that continuity of the trace operator γ proves that $\gamma v = 0$, for $v \in H_0^1(\Omega)$, as well as $(\gamma v)|_{\Gamma_D} = 0$, for $v \in H_D^1(\Omega)$.

Corollary 2.15 (Friedrichs Inequality). *It holds that*

$$\|v\|_{L^2(\Omega)} \leq C_F \text{diam}(\Omega) \|\nabla v\|_{L^2(\Omega)} \quad \text{for all } v \in H_D^1(\Omega) \tag{2.11}$$

with a constant $C_F > 0$ that depends only on the shape of Ω and Γ_D . ■

We finally note the relation between $H_D^1(\Gamma)$ and the trace operator, cf. the Theorem of Meyers-Serrin.

Theorem 2.16. *There holds $H_0^1(\Omega) = \ker(\gamma)$ with $\gamma \in L(H^1(\Omega); L^2(\Gamma))$ the trace operator. Moreover, $H_D^1(\Omega) = \{v \in H^1(\Omega) \mid (\gamma v)|_{\Gamma_D} = 0\}$.* ■

Exercise 12. Usually, one defines the range of the trace operator as $H^{1/2}(\Gamma) := \text{range}(\gamma) \subseteq L^2(\Gamma)$. This space is associated with the norm $\|v\|_{H^{1/2}(\Gamma)} := \inf \{ \|\widehat{v}\|_{H^1(\Omega)} \mid \widehat{v} \in H^1(\Omega) \text{ with } \gamma \widehat{v} = v \}$. Prove that $H^{1/2}(\Gamma)$ associated with this norm is a Hilbert space with continuous inclusion $H^{1/2}(\Gamma) \subseteq L^2(\Gamma)$. **Hint:** Recall the definition and the standard results on quotient spaces and the associated quotient norm! □

For $X = H^1(\Omega)$ and $Y = L^2(\Omega)$, the following exercise shows that the L^2 -scalar products $(f; \cdot)_{L^2(\Omega)}$ for $f \in L^2(\Omega)$ give (up to density) all linear and continuous functionals on $H^1(\Omega)$, i.e., the embedding $L^2(\Omega) \rightarrow H^1(\Omega)^*$, $f \mapsto (f; \cdot)_{L^2(\Omega)}$ is well-defined, linear, continuous, and injective with dense image.

Exercise 13. Let X and Y be Hilbert spaces with continuous embedding $X \subseteq Y$. Show that the mapping $I : Y^* \rightarrow X^*$, $Iy^* := y^*|_X$ is well-defined, linear, and continuous. Prove that $I(Y^*) \subseteq X^*$ is a dense subspace. Moreover, if $X \subseteq Y$ is dense with respect to $\|\cdot\|_Y$, then the embedding I is even injective. □

2.3 Weak Form of Laplace Problem (18.10.2017)

2.3.1 Dirichlet Problem

In this section, we generalize the variational form derived in the introductory section to our Hilbert space setting. We start with the homogeneous Dirichlet problem

$$\begin{aligned} -\Delta u &= f & \text{in } \Omega, \\ u &= 0 & \text{on } \Gamma. \end{aligned} \tag{2.12}$$

Recall that this formulation is called the **strong form** of the boundary value problem. The following proposition provides the — in some sense — equivalent and always uniquely solvable weak form of the boundary value problem.

Proposition 2.17. (i) *Provided that $u \in C^2(\overline{\Omega})$ solves (2.12) for a given source term $f \in C(\overline{\Omega})$, it holds that $u \in H_0^1(\Omega)$ as well as*

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} \quad \text{for all } v \in H_0^1(\Omega). \tag{2.13}$$

(ii) *Given $f \in L^2(\Omega)$, the **weak form** (2.13) has a unique solution $u \in H_0^1(\Omega)$. It holds that*

$$\|u\|_{H^1(\Omega)} \leq C \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{(f ; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} \leq C \|f\|_{L^2(\Omega)}, \tag{2.14}$$

where the constant $C > 0$ depends only on Ω .

(iii) *Provided that $f \in C(\overline{\Omega})$ and that the weak solution $u \in H_0^1(\Omega)$ of (2.13) additionally satisfies $u \in C^2(\overline{\Omega})$, then u even solves the strong form (2.12).*

Proof. (i) We have already seen before that a strong solution $u \in C^2(\overline{\Omega})$ solves the variational form (2.13) for test functions $v \in C_0^1(\overline{\Omega}) := \{w \in C^1(\overline{\Omega}) \mid w|_{\Gamma} = 0\}$ replacing $H_0^1(\Omega)$; see Proposition 1.1. If we keep u fixed, the left-hand side as well as the right-hand side of (2.13) define continuous and linear functionals on $H^1(\Omega)$. Note that the closure of $C_0^1(\overline{\Omega})$ with respect to the H^1 -norm leads to the Hilbert space $H_0^1(\Omega)$. Therefore, standard density arguments prove (2.13).

(ii) According to the Friedrichs inequality, it holds that

$$\|\nabla v\|_{L^2(\Omega)}^2 \leq \|v\|_{H^1(\Omega)}^2 \leq (1 + \tilde{C}_F^2) \|\nabla v\|_{L^2(\Omega)}^2 \quad \text{for all } v \in H_0^1(\Omega).$$

Therefore, the left-hand side of (2.13) defines an equivalent scalar product on $H_0^1(\Omega)$. The Riesz theorem thus provides a unique weak solution $u \in H_0^1(\Omega)$ of (2.13). Plugging-in $u = v \in H_0^1(\Omega)$, the weak form yields that

$$(1 + \tilde{C}_F^2)^{-1} \|u\|_{H^1(\Omega)}^2 \leq \|\nabla u\|_{L^2(\Omega)}^2 = (f ; u)_{L^2(\Omega)} \leq \sup_{v \in H_0^1(\Omega) \setminus \{0\}} \frac{(f ; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} \|u\|_{H^1(\Omega)}$$

which results in the first estimate of (2.14). The second estimate follows from the Cauchy inequality

$$(f ; v)_{L^2(\Omega)} \leq \|f\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \leq \|f\|_{L^2(\Omega)} \|v\|_{H^1(\Omega)}.$$

(iii) Since the weak solution u is smooth, we may use integration by parts to see that

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (-\Delta u ; v)_{L^2(\Omega)} \quad \text{for all } v \in H_0^1(\Omega).$$

The difference with the weak form (2.13) thus yields that

$$0 = (f + \Delta u ; v)_{L^2(\Omega)} \quad \text{for all } v \in H_0^1(\Omega).$$

Note that $F := f + \Delta u \in C(\overline{\Omega})$. With $\mathcal{D}(\Omega) \subseteq H_0^1(\Omega)$, Theorem 2.1 proves $F = 0$; see also the remark right after Theorem 2.1. Consequently, it holds that $-\Delta u = f$ in Ω . The Dirichlet boundary conditions (in the strong form) follow from $0 = \gamma u = u|_\Gamma$. Altogether, u solves (2.12) ■

2.3.2 Mixed Boundary Value Problem

Second, we consider the mixed boundary value problem

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \Gamma_D, \\ \partial u / \partial n &= \phi && \text{on } \Gamma_N, \end{aligned} \tag{2.15}$$

with $\Gamma = \overline{\Gamma}_D \cup \overline{\Gamma}_N$, $\Gamma_D \cap \Gamma_N = \emptyset$, and $|\Gamma_D| > 0$. The limit case $|\Gamma_D| = 0$ corresponds to the Neumann problem which is treated in Section 2.3.3. Recall the trace norm $\|\cdot\|_{H^{1/2}(\Gamma)}$ from Exercise 12. Then, the main proposition reads as follows:

Proposition 2.18. (i) *Suppose that Γ_N is smooth, i.e., the outer normal vector depends continuously on $x \in \Gamma_N$. Provided that $u \in C^2(\overline{\Omega})$ solves the **strong form** (2.15) for a given source term $f \in C(\overline{\Omega})$ and Neumann data $\phi \in C(\overline{\Gamma}_N)$, it holds that $u \in H_D^1(\Omega)$ as well as*

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} + (\phi ; \gamma v)_{L^2(\Gamma_N)} \quad \text{for all } v \in H_D^1(\Omega). \tag{2.16}$$

(ii) *Given $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma_N)$, the **weak form** (2.16) has a unique solution $u \in H_D^1(\Omega)$. It holds that*

$$\begin{aligned} \|u\|_{H^1(\Omega)} &\leq C_1 \left(\sup_{v \in H_D^1(\Omega) \setminus \{0\}} \frac{(f ; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma) \setminus \{0\}} \frac{(\phi ; w)_{L^2(\Gamma_N)}}{\|w\|_{H^{1/2}(\Gamma)}} \right) \\ &\leq C_2 (\|f\|_{L^2(\Omega)} + \|\phi\|_{L^2(\Gamma_N)}) \end{aligned} \tag{2.17}$$

where the constants $C_1, C_2 > 0$ depend only on Ω and Γ_D .

(iii) *Provided that $f \in C(\overline{\Omega})$ and $\phi \in C(\overline{\Gamma}_N)$ and that the weak solution $u \in H_D^1(\Omega)$ of (2.16) additionally satisfies $u \in C^2(\overline{\Omega})$, then u even solves the strong form (2.15).*

Proof is done in the exercises. ■

v4:
13.10.2017

2.3.3 Neumann Problem

Finally, we consider the Neumann problem

$$\begin{aligned} -\Delta u &= f & \text{in } \Omega, \\ \partial u / \partial n &= \phi & \text{on } \Gamma. \end{aligned} \tag{2.18}$$

Note that the solution u of (2.18) cannot be unique: If $u \in C^2(\overline{\Omega})$ solves the **strong form** (2.18), also $u + c$ solves (2.18), for all $c \in \mathbb{R}$. To fix the additive constant, we seek a solution which additionally satisfies, e.g., that

$$\int_{\Omega} u \, dx = 0. \tag{2.19}$$

Moreover, the Gauss divergence theorem shows

$$-\int_{\Omega} f \, dx = \int_{\Omega} \Delta u \, dx = \int_{\Omega} \operatorname{div}(\nabla u) \, dx = \int_{\Gamma} \frac{\partial u}{\partial n} \, ds = \int_{\Gamma} \phi \, ds.$$

Therefore, the data f and ϕ have to satisfy the compatibility condition

$$\int_{\Omega} f \, dx + \int_{\Gamma} \phi \, ds = 0 \tag{2.20}$$

to allow for the existence of (strong) solutions. Recall the trace norm $\|\cdot\|_{H^{1/2}(\Gamma)}$ from Exercise 12.

Proposition 2.19. (i) *Suppose that Γ is smooth, i.e., the outer normal vector depends continuously on $x \in \Gamma$. Provided that $u \in C^2(\overline{\Omega})$ solves (2.18) for a given source term $f \in C(\overline{\Omega})$ and Neumann data $\phi \in C(\Gamma)$, it holds that $u \in H^1(\Omega)$ and*

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} + (\phi ; \gamma v)_{L^2(\Gamma)} \quad \text{for all } v \in H^1(\Omega). \tag{2.21}$$

(ii) *Given $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma)$, the variational formulation*

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} + (\phi ; \gamma v)_{L^2(\Gamma)} \quad \text{for all } v \in H_*^1(\Omega) \tag{2.22}$$

has a unique solution $u \in H_^1(\Omega) := \{v \in H^1(\Omega) \mid \int_{\Omega} v \, dx = 0\}$.*

(iii) *Provided that the data $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma)$ satisfy (2.20), the unique solution $u \in H_*^1(\Omega)$ of (2.22) even solves the **weak form** (2.21). Moreover, it holds that*

$$\begin{aligned} \|u\|_{H^1(\Omega)} &\leq C_1 \left(\sup_{v \in H^1(\Omega) \setminus \{0\}} \frac{(f ; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma) \setminus \{0\}} \frac{(\phi ; w)_{L^2(\Gamma)}}{\|w\|_{H^{1/2}(\Gamma)}} \right) \\ &\leq C_2 (\|f\|_{L^2(\Omega)} + \|\phi\|_{L^2(\Gamma)}) \end{aligned} \tag{2.23}$$

where the constants $C_1, C_2 > 0$ depend only on Ω .

(iv) *Provided that $f \in C(\overline{\Omega})$ and $\phi \in C(\Gamma)$ satisfy (2.20) and that the weak solution $u \in H_*^1(\Omega)$ of (2.21) resp. (2.22) additionally satisfies $u \in C^2(\overline{\Omega})$, then u even solves the strong form (2.18).*

Proof. (i) The variational form (2.21) holds for test functions $v \in C^1(\overline{\Omega})$ according to integration by parts. For fixed u , the left-hand as well as the right-hand side define continuous linear functionals on $H^1(\Omega)$. Thus, (2.21) follows for $v \in H^1(\Omega)$ by density arguments. (ii) According to the Poincaré inequality, it holds that

$$\|\nabla v\|_{L^2(\Omega)}^2 \leq \|v\|_{H^1(\Omega)}^2 \leq (1 + \tilde{C}_P^2) \|\nabla v\|_{L^2(\Omega)}^2 \quad \text{for all } v \in H_*^1(\Omega).$$

Therefore, the left-hand side of (2.22) defines an equivalent scalar product on $H_*^1(\Omega)$. Note that $H_*^1(\Omega)$ is a closed subspace of $H^1(\Omega)$ and hence a Hilbert space. Therefore, (2.22) follows from the Riesz theorem. (iii) For a function $v \in H^1(\Omega)$, we define $\tilde{v} := v - v_\Omega \in H_*^1(\Omega)$, where $v_\Omega \in \mathbb{R}$ denotes the integral mean $v_\Omega := (1/|\Omega|) \int_\Omega v \, dx \in \mathbb{R}$. Note that (2.20) implies that

$$(f; v_\Omega)_{L^2(\Omega)} + (\phi; v_\Omega)_{L^2(\Gamma)} = 0.$$

Thus, (2.22) proves that

$$(\nabla u; \nabla v)_{L^2(\Omega)} = (\nabla u; \nabla \tilde{v})_{L^2(\Omega)} = (f; \tilde{v})_{L^2(\Omega)} + (\phi; \gamma \tilde{v})_{L^2(\Gamma)} = (f; v)_{L^2(\Omega)} + (\phi; \gamma v)_{L^2(\Gamma)},$$

i.e., u even solves (2.21). Plugging-in $u = v$, we see that

$$\|\nabla u\|_{L^2(\Omega)}^2 \leq \sup_{v \in H^1(\Omega) \setminus \{0\}} \frac{(f; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} \|u\|_{H^1(\Omega)} + \sup_{w \in H^{1/2}(\Gamma) \setminus \{0\}} \frac{(\phi; w)_{L^2(\Gamma)}}{\|w\|_{H^{1/2}(\Gamma)}} \|\gamma u\|_{H^{1/2}(\Gamma)},$$

where we have used that $H^{1/2}(\Gamma) = \text{range}(\gamma)$. Note that the $H^{1/2}$ -norm is defined in such a way that $\gamma \in L(H^1(\Omega); H^{1/2}(\Gamma))$ with $\|\gamma u\|_{H^{1/2}(\Gamma)} \leq \|u\|_{H^1(\Omega)}$. Therefore,

$$\|\nabla u\|_{L^2(\Omega)}^2 \leq \|u\|_{H^1(\Omega)} \left(\sup_{v \in H^1(\Omega) \setminus \{0\}} \frac{(f; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma) \setminus \{0\}} \frac{(\phi; w)_{L^2(\Gamma)}}{\|w\|_{H^{1/2}(\Gamma)}} \right).$$

Together with $(1 + \tilde{C}_P^2)^{-1} \|u\|_{H^1(\Omega)}^2 \leq \|\nabla u\|_{L^2(\Omega)}^2$, this proves the first estimate in (2.23). As above, the first supremum may be estimated by $\|f\|_{L^2(\Omega)}$. With the continuous embedding $H^{1/2}(\Gamma) \subset L^2(\Gamma)$, the numerator of the second supremum can be dominated by

$$(\phi; w)_{L^2(\Gamma)} \leq \|\phi\|_{L^2(\Gamma)} \|w\|_{L^2(\Gamma)} \leq \tilde{C} \|\phi\|_{L^2(\Gamma)} \|w\|_{H^{1/2}(\Gamma)}.$$

This provides the upper bound $\tilde{C} \|\phi\|_{L^2(\Gamma)}$ for the second supremum. (iv) As above, we may use integration by parts to see that

$$(f + \Delta u; v)_{L^2(\Omega)} + (\phi - \partial u / \partial n; \gamma v)_{L^2(\Gamma)} = 0 \quad \text{for all } v \in H^1(\Omega).$$

From this, we first conclude $f = -\Delta u$ by use of Theorem 2.1 for test functions $v \in \mathcal{D}(\Omega) \subset H_0^1(\Omega) \subset H^1(\Omega)$. To prove $\phi = \partial u / \partial n$, one proceeds analogously to the remark right after Theorem 2.1. ■

Chapter 3

A Priori Analysis

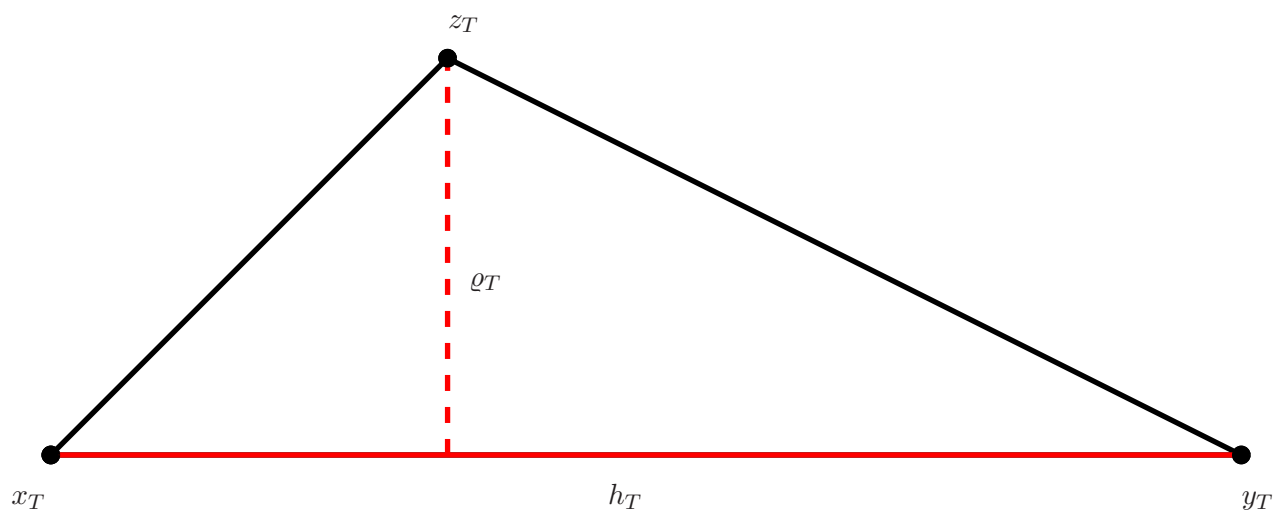


FIGURE 3.1. The diameter h_T of the triangle T is the length of the longest edge (possibly non-unique). The quantity ϱ_T denotes the corresponding height.

3.1 P1-Finite Element Method in 2D (18.10.2017)

A set $T \subset \mathbb{R}^2$ is called a **non-degenerate triangle** provided that there are nodes $x_T, y_T, z_T \in \mathbb{R}^2$ with $T = \text{conv}\{x_T, y_T, z_T\}$ and provided that $|T| > 0$, i.e., T has positive measure. We note that T is in particular bounded and closed, whence compact. We denote by

$$\mathcal{K}_T := \{x_T, y_T, z_T\} \tag{3.1}$$

the **set of nodes** of T and by

$$\mathcal{E}_T := \{ \text{conv}\{x_T, y_T\}, \text{conv}\{y_T, z_T\}, \text{conv}\{z_T, x_T\} \} \tag{3.2}$$

the **set of edges** of T . The **diameter** of T is denoted by

$$h_T := \text{diam}(T) := \max \{ |x - y| \mid x, y \in T \}. \tag{3.3}$$

Moreover, we define the **edge length**

$$h_E := \text{diam}(E) := \max \{|x - y| \mid x, y \in E\} \quad (3.4)$$

for all edges $E \in \mathcal{E}_T$. Clearly, the diameter h_T of a triangle is the length of the longest edge (possibly non unique), i.e., there is some $E \in \mathcal{E}_T$ with $h_T = h_E$. The **height** over the longest edge E of T is denoted by ϱ_T , cf. Figure 3.1. Recall that the measure of the triangle reads

$$|T| = \frac{h_T \varrho_T}{2}. \quad (3.5)$$

The most important example is the **reference triangle**

$$T_{\text{ref}} := \text{conv}\{(0, 0), (1, 0), (0, 1)\} \quad (3.6)$$

which has measure $|T_{\text{ref}}| = 1/2$.

Exercise 14. Give a formal proof that the diameter of a triangle T is the length of one longest edge, i.e., $h_T = \max_{E \in \mathcal{E}_T} h_E$. *Hint:* Use that the convex hull $\text{conv}(M) := \bigcap \{\widehat{M} \subseteq \mathbb{R}^d \mid \widehat{M} \text{ is convex with } M \subseteq \widehat{M}\}$ of a set $M \subseteq \mathbb{R}^d$ is also characterized by $\text{conv}(M) = \{\sum_{j=1}^N \lambda_j x_j \mid N \in \mathbb{N}, x_j \in M, \lambda_j \geq 0 \text{ with } \sum_{j=1}^N \lambda_j = 1\}$. The proof then directly applies to general simplices in \mathbb{R}^d , i.e., $T = \text{conv}\{x_0, \dots, x_d\} \subset \mathbb{R}^d$. \square

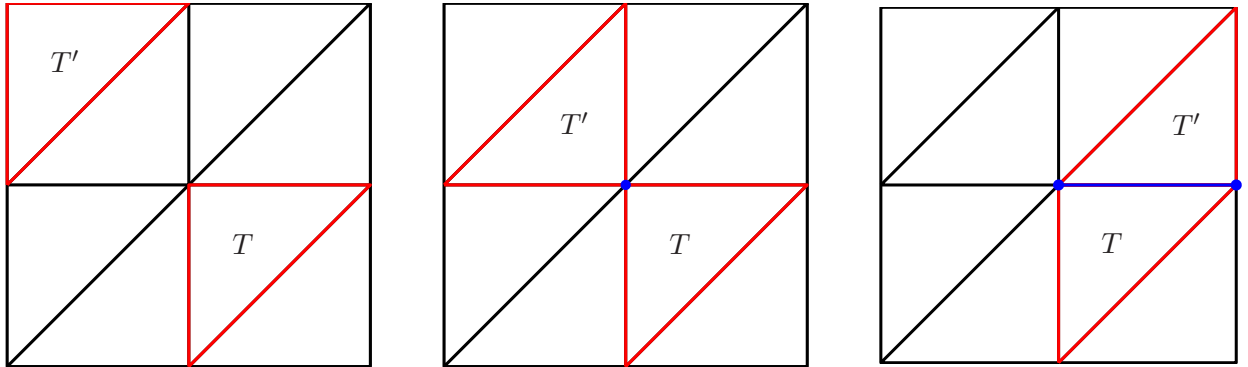


FIGURE 3.2. For a regular triangulation \mathcal{T} , the intersection of two elements $T \neq T'$ is either empty, a joint node, or a joint edge.

Definition. A set \mathcal{T} is a **triangulation** of Ω (consisting of triangles) if and only if

- \mathcal{T} is a finite set of non-degenerate triangles,
- the closure of Ω is covered by \mathcal{T} , i.e., $\overline{\Omega} = \bigcup \mathcal{T}$,
- for all $T, T' \in \mathcal{T}$ with $T \neq T'$, it holds that $|T \cap T'| = 0$, i.e., the overlap is a set of measure zero.

By $\mathcal{K} := \bigcup \{x \in \mathcal{K}_T \mid T \in \mathcal{T}\}$, we then denote the **set of nodes** of the triangulation \mathcal{T} and by $\mathcal{E} := \bigcup \{E \in \mathcal{E}_T \mid T \in \mathcal{T}\}$ the **set of edges** of the triangulation \mathcal{T} . A triangulation of Ω is called **conforming** or **regular (in the sense of Ciarlet)** provided that the intersection of two elements $T, T' \in \mathcal{T}$ with $T \neq T'$ is

- either empty,
- or a joint node, i.e., $T \cap T' = \{z\} = \mathcal{K}_T \cap \mathcal{K}_{T'}$,
- or a joint edge, i.e., $E := T \cap T' \in \mathcal{E}_T \cap \mathcal{E}_{T'}$,

cf. Figure 3.2. According to this regularity assumption, an edge $E \in \mathcal{E}$ with surface measure $|E \cap \Gamma| > 0$ automatically satisfies $E \subseteq \Gamma$, i.e., an edge E is either a boundary edge or an interior edge. Additionally, we always assume that a regular triangulation resolves the boundary conditions: If $\Gamma = \partial\Omega$ is partitioned into Dirichlet and Neumann boundary Γ_D and Γ_N , respectively, each boundary edge $E \in \mathcal{E}$ with $E \subseteq \Gamma$ satisfies

- either $E \subseteq \overline{\Gamma}_D$
- or $E \subseteq \overline{\Gamma}_N$.

With this assumption, we define the (disjoint) sets of boundary edges

$$\mathcal{E}_D := \{E \in \mathcal{E} \mid E \subseteq \overline{\Gamma}_D\} \quad \text{and} \quad \mathcal{E}_N := \{E \in \mathcal{E} \mid E \subseteq \overline{\Gamma}_N\} \quad (3.7)$$

as well as the set of all interior edges

$$\mathcal{E}_\Omega := \mathcal{E} \setminus (\mathcal{E}_D \cup \mathcal{E}_N). \quad (3.8)$$

We finally note that, for each $E \in \mathcal{E}_\Omega$, there are two elements $T, T' \in \mathcal{T}$ with $E = T \cap T'$.

Exercise 15. Let \mathcal{T} be a regular triangulation of Ω and $v : \Omega \rightarrow \mathbb{R}$ such that $v|_T \in C^1(T)$ for all $T \in \mathcal{T}$. Prove that $v \in H^1(\Omega)$ if and only if $v \in C(\Omega)$. □

The following proposition essentially follows from the regularity of the triangulation \mathcal{T} .

v5:
18.10.2017

Proposition 3.1. For a regular triangulation \mathcal{T} of Ω , we define the discrete space

$$\mathcal{S}^1(\mathcal{T}) := \{v_h \in C(\Omega) \mid \forall T \in \mathcal{T} \quad v_h|_T \text{ affine}\} \quad (3.9)$$

of all \mathcal{T} -piecewise affine and globally continuous functions. Then, there holds the following:

- (i) $\mathcal{S}^1(\mathcal{T})$ is an N -dimensional subspace of $H^1(\Omega)$ with $N = \#\mathcal{K}$ the number of nodes.
- (ii) For each node $z \in \mathcal{K}$, there is a unique **hat function**

$$\zeta_z \in \mathcal{S}^1(\mathcal{T}) \quad \text{with} \quad \zeta_z(z') = \delta_{zz'} \quad \text{for all } z' \in \mathcal{K}. \quad (3.10)$$

- (iii) The set $\mathcal{B} := \{\zeta_z \mid z \in \mathcal{K}\}$ is a basis of $\mathcal{S}^1(\mathcal{T})$, the so-called **nodal basis**.

Proof. 1. step. According to the regularity of \mathcal{T} , hat functions ζ_z are automatically continuous on Ω : For each element $T \in \mathcal{T}$, an affine function $v_h : T \rightarrow \mathbb{R}$ is uniquely determined by the nodal values $v_h(z)$ for $z \in \mathcal{K}_T$. Therefore, the \mathcal{T} -piecewise affine hat function ζ_z defined by $\zeta_z(z') = \delta_{zz'}$ is uniquely defined. We now show that $\zeta_z \in C(\Omega)$: If $T, T' \in \mathcal{T}$ are elements with $T \cap T' \neq \emptyset$, regularity of \mathcal{T} implies that either $T = T'$ or $\{z'\} = T \cap T'$ is a joint point or $E = T \cap T'$ is a joint edge. In the latter case, note that the trace on E of the affine function $\zeta_z|_T$ as well as of $\zeta_z|_{T'}$ is

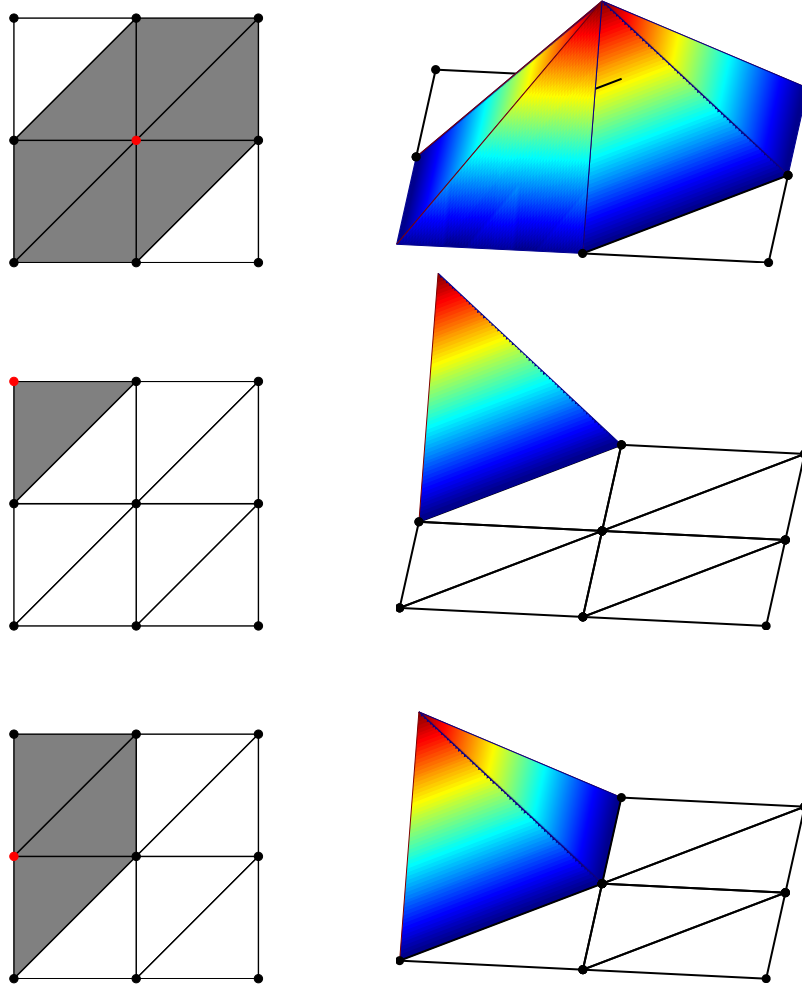


FIGURE 3.3. Examples of $P1$ hat functions ζ_z : The left figures show the mesh as well as the support $\text{supp}(\zeta_z)$ in grey, where the corresponding node $z \in \mathcal{K}$ is indicated in red. The right figures show the plots of the hat functions. Triangles $T \in \mathcal{T}$ with $\zeta_z|_T = 0$ are filled with white.

uniquely defined on the edge E by the nodal values $\zeta_z(x_E)$ and $\zeta_z(y_E)$, where $E = \text{conv}\{x_E, y_E\}$. Therefore the traces of $\zeta_z|_T$ and $\zeta_z|_{T'}$ on E coincide, i.e., ζ_z is continuous on each interior edge.

2. step. The nodal basis \mathcal{B} is a basis of $\mathcal{S}^1(\mathcal{T})$ and $\dim \mathcal{S}^1(\mathcal{T}) = \#\mathcal{K}$: Clearly, the hat functions are linearly independent, $\mathcal{B} \subseteq \mathcal{S}^1(\mathcal{T})$, and $\#\mathcal{B} = \#\mathcal{K}$. Moreover, each function $v_h \in \mathcal{S}^1(\mathcal{T})$ is uniquely defined by the nodal values $v_h(z)$ for $z \in \mathcal{K}$ and can thus be written as the linear combination of the hat functions, i.e., $\mathcal{S}^1(\mathcal{T}) \subseteq \text{span}(\mathcal{B})$.

3. step. The inclusion $\mathcal{S}^1(\mathcal{T}) \subset H^1(\Omega)$ follows from Exercise 15. ■

Remark. Examples for hat functions ζ_z are shown in Figure 3.3. Note that the support $\text{supp}(\zeta_z)$ is always local. This leads to a sparse Galerkin matrix A , i.e., most of the entries of A are zero. □

For a given Dirichlet boundary $\Gamma_D \subseteq \Gamma$, we use the discrete space $\mathcal{S}_D^1(\mathcal{T})$ to discretize the weak form of the mixed boundary value problem. In case of $\Gamma_D = \Gamma$, we consider the space $\mathcal{S}_0^1(\mathcal{T})$.

Corollary 3.2. *Let \mathcal{T} be a regular triangulation of Ω . Then, the space*

$$\mathcal{S}_D^1(\mathcal{T}) := \{v_h \in \mathcal{S}^1(\mathcal{T}) \mid \forall z \in \mathcal{K} \cap \bar{\Gamma}_D \quad v_h(z) = 0\} \quad (3.11)$$

is a finite dimensional subspace of $H_D^1(\Omega)$ of dimension $\#\{z \in \mathcal{K} \mid z \notin \bar{\Gamma}_D\}$. The space

$$\mathcal{S}_0^1(\mathcal{T}) := \{v_h \in \mathcal{S}^1(\mathcal{T}) \mid \forall z \in \mathcal{K} \cap \Gamma \quad v_h(z) = 0\} \quad (3.12)$$

is a finite dimensional subspace of $H_0^1(\Omega)$ of dimension $\#\{z \in \mathcal{K} \mid z \notin \Gamma\}$.

Proof. We only need to show that $v_h|_{\Gamma_D} = 0$ for $v_h \in \mathcal{S}_D^1(\mathcal{T})$. Let $x \in \Gamma_D$. According to the regularity of \mathcal{T} , there is an edge $E \in \mathcal{E}_D$ such that $x \in E$. Since the trace $v_h|_E$ is affine, it is uniquely determined by the nodal values $v_h(x_T) = 0 = v_h(y_T)$, where $E = \text{conv}\{x_T, y_T\}$. Consequently, $v_h|_E = 0$ for all $E \in \mathcal{E}_D$ and hence $v_h \in H_D^1(\Omega)$. In particular, we obtain the claim for $\mathcal{S}_0^1(\mathcal{T})$ in case of $\Gamma_D = \Gamma$. ■

For the discretization of the Neumann problem, we are dealing with $\mathcal{S}_*^1(\mathcal{T})$.

Corollary 3.3. *For a regular triangulation \mathcal{T} of Ω , the space*

$$\mathcal{S}_*^1(\mathcal{T}) := \{v_h \in \mathcal{S}^1(\mathcal{T}) \mid \int_{\Omega} v_h \, dx = 0\} \quad (3.13)$$

is a finite dimensional subspace of $H_^1(\Omega)$ of dimension $\#\mathcal{K} - 1$.*

Proof. Clearly, it holds that $\mathcal{S}_*^1(\mathcal{T}) \subseteq H_*^1(\Omega)$. Note that $I(v_h) := \int_{\Omega} v_h \, dx$ is a linear functional on $\mathcal{S}^1(\mathcal{T})$ with kernel $\mathcal{S}_*^1(\mathcal{T}) = \ker(I)$. Since $\text{rank}(I) = 1$, Linear Algebra yields that $\dim \mathcal{S}_*^1(\mathcal{T}) = \dim \mathcal{S}^1(\mathcal{T}) - 1$. ■

The **P1 Finite Element Method** now consists of using the Galerkin method with the discrete spaces $\mathcal{S}_0^1(\mathcal{T})$, $\mathcal{S}_D^1(\mathcal{T})$, and $\mathcal{S}_*^1(\mathcal{T})$ to approximate the weak solution of the Dirichlet problem, the mixed boundary value problem, and the Neumann problem, respectively. From now on, we shall assume that \mathcal{T} is a regular triangulation of Ω . We start with the **Dirichlet problem**

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

for given data $f \in L^2(\Omega)$. The P1-FEM then reads: Find $u_h \in \mathcal{S}_0^1(\mathcal{T})$ such that

$$(\nabla u_h ; \nabla v_h)_{L^2(\Omega)} = (f ; v_h)_{L^2(\Omega)} \quad \text{for all } v_h \in \mathcal{S}_0^1(\mathcal{T}). \quad (3.14)$$

Second, the **mixed boundary value problem** reads

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \Gamma_D, \\ \partial u / \partial n &= \phi && \text{on } \Gamma_N, \end{aligned}$$

with $\Gamma = \bar{\Gamma}_D \cup \bar{\Gamma}_N$, $\Gamma_D \cap \Gamma_N = \emptyset$, and $|\Gamma_D| > 0$. The data satisfy $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma_N)$. The P1-FEM for the mixed BVP reads: Find $u_h \in \mathcal{S}_D^1(\mathcal{T})$ such that

$$(\nabla u_h ; \nabla v_h)_{L^2(\Omega)} = (f ; v_h)_{L^2(\Omega)} + (\phi ; v_h)_{L^2(\Gamma_N)} \quad \text{for all } v_h \in \mathcal{S}_D^1(\mathcal{T}). \quad (3.15)$$

Finally, we consider the **Neumann problem**

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ \partial u / \partial n &= \phi && \text{on } \Gamma, \end{aligned}$$

where the data $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma)$ are assumed to satisfy $\int_{\Omega} f \, dx + \int_{\Gamma} \phi \, ds = 0$. The P1-FEM for the Neumann problem reads: Find $u_h \in \mathcal{S}_*^1(\mathcal{T})$ such that

$$(\nabla u_h ; \nabla v_h)_{L^2(\Omega)} = (f ; v_h)_{L^2(\Omega)} + (\phi ; v_h)_{L^2(\Gamma)} \quad \text{for all } v_h \in \mathcal{S}_*^1(\mathcal{T}). \quad (3.16)$$

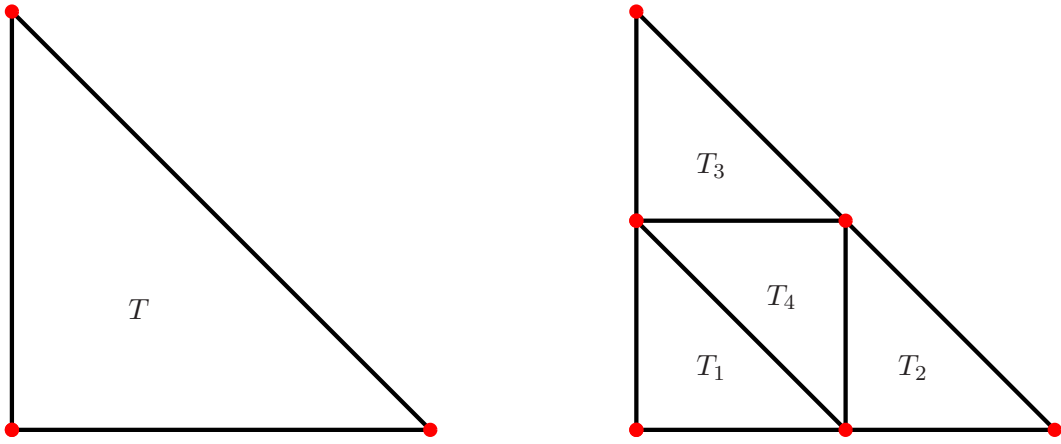


FIGURE 3.4. Red-refinement refines the element $T \in \mathcal{T}^{(\text{old})}$ into 4 similar elements $T_1, \dots, T_4 \in \mathcal{T}^{(\text{new})}$. The new nodes $\mathcal{K}^{(\text{new})} \setminus \mathcal{K}^{(\text{old})}$ are just the edge midpoints for all edges $E \in \mathcal{E}^{(\text{old})}$. In particular, regularity of $\mathcal{T}^{(\text{old})}$ implies regularity of $\mathcal{T}^{(\text{new})}$.

3.2 Approximation Theorem and Bramble-Hilbert Lemma (18.10.2017)

3.2.1 Uniform Mesh-Refinement and Shape Regularity

Let $h \in L^\infty(\Omega)$ and $\varrho \in L^\infty(\Omega)$ denote the **local mesh-width** functions which are defined by

$$h|_T := h_T = \text{diam}(T) \quad \text{and} \quad \varrho|_T := \varrho_T \quad \text{for all } T \in \mathcal{T}. \quad (3.17)$$

Moreover, the quantities

$$\sigma(T) := \frac{h_T}{\varrho_T} \quad \text{and} \quad \sigma(\mathcal{T}) := \|h/\varrho\|_{L^\infty(\Omega)} = \max_{T \in \mathcal{T}} \frac{h_T}{\varrho_T} \geq 1 \quad (3.18)$$

denote the **shape regularity constant** of an element $T \in \mathcal{T}$ resp. the triangulation \mathcal{T} . Note that $|T| = h_T \varrho_T / 2$ so that $2h_T / \varrho_T = h_T^2 / |T|$. The shape regularity constant will affect all error estimates, so that mesh-refinement has to avoid a blow-up of $\sigma(\mathcal{T})$. We say that a regular mesh \mathcal{T} is **γ -shape regular**, if $\sigma(\mathcal{T}) \leq \gamma < \infty$.

For this section, we stick with the so-called **uniform mesh-refinement**: Given a regular triangulation $\mathcal{T}^{(\text{old})}$, we obtain a new triangulation $\mathcal{T}^{(\text{new})}$ as follows: Each element $T \in \mathcal{T}^{(\text{old})}$ is split into 4 similar triangles $T_1, \dots, T_4 \in \mathcal{T}^{(\text{new})}$, cf. Figure 3.4. Therefore, each node $z \in \mathcal{K}^{(\text{new})}$ either belongs to $\mathcal{K}^{(\text{old})}$ or is the midpoint of an edge $E \in \mathcal{E}^{(\text{old})}$. We stress some simple observations:

- The new triangulation $\mathcal{T}^{(\text{new})}$ is also regular.
- The local mesh-width functions satisfy $h^{(\text{new})} = h^{(\text{old})}/2$ and $\varrho^{(\text{new})} = \varrho^{(\text{old})}/2$.
- In particular, the shape regularity constant satisfies that $\sigma(\mathcal{T}^{(\text{old})}) = \sigma(\mathcal{T}^{(\text{new})})$.

Further mesh-refinement strategies are discussed in the following section.

Exercise 16. Let $T = \text{conv}\{z_1, z_2, z_3\}$ be a non-degenerate triangle in \mathbb{R}^2 . Prove that the shape regularity constant h_T/ϱ_T tends to infinity if and only if the smallest angle in T tends to zero. □

Exercise 17. Often, the shape regularity constant is defined as the maximal quotient h_T/r_T , where $r_T > 0$ denotes the maximal radius of a ball $B(x, r_T) := \{y \in \mathbb{R}^2 \mid |x-y| \leq r_T\}$ inscribed in T , i.e., $B(x, r_T) \subseteq T$. Let $T = \text{conv}\{z_1, z_2, z_3\}$ be a non-degenerate triangle in \mathbb{R}^2 . What is the relation between ϱ_T and r_T ? □

3.2.2 Statement and Interpretation of Approximation Theorem

To state our first main result in this section, we need to know that certain Sobolev functions are at least continuous.

Theorem 3.4 (Sobolev). *Let Ω be a Lipschitz domain in \mathbb{R}^d and $m > d/2$. Then, there holds the continuous inclusion $H^m(\Omega) \subseteq C(\overline{\Omega})$.* ■

In particular, for $d = 2, 3$, each Sobolev function $u \in H^2(\Omega)$ is continuous so that evaluation of u at the nodes $z \in \mathcal{K}$ is well-defined. Throughout the remaining section, we assume that \mathcal{T} is a regular triangulation of a bounded Lipschitz domain $\Omega \subset \mathbb{R}^2$. We stress, however, that the same results — even with the same proofs — hold for $d = 3$ as well. As in the previous section, the nodal basis function corresponding to a node $z \in \mathcal{K}$ is denoted by $\zeta_z \in \mathcal{S}^1(\mathcal{T})$.

Theorem 3.5 (Approximation Theorem). *For $u \in H^2(\Omega)$, the nodal interpolant reads*

$$I_h u := \sum_{z \in \mathcal{K}} u(z) \zeta_z \in \mathcal{S}^1(\mathcal{T}). \quad (3.19)$$

For all $T \in \mathcal{T}$, there hold the elementwise error estimates

$$\|u - I_h u\|_{L^2(T)} \leq C \|h^2 D^2 u\|_{L^2(T)} \quad (3.20)$$

and

$$\|\nabla(u - I_h u)\|_{L^2(T)} \leq C \sigma(T) \|h D^2 u\|_{L^2(T)}, \quad (3.21)$$

where the generic constant $C > 0$ is independent of u , \mathcal{T} , and Ω , but depends only on the reference triangle. In particular, this proves for all $\alpha \in \mathbb{R}$ the global error estimates

$$\|h^\alpha(u - I_h u)\|_{L^2(\Omega)} \leq C \|h^{2+\alpha} D^2 u\|_{L^2(\Omega)} \quad (3.22)$$

and

$$\|h^\alpha \nabla(u - I_h u)\|_{L^2(\Omega)} \leq C \sigma(\mathcal{T}) \|h^{1+\alpha} D^2 u\|_{L^2(\Omega)}. \quad (3.23)$$

Before the proof of Theorem 3.5, we discuss the following immediate consequence:

Corollary 3.6. *For $u \in H^2(\Omega) \cap H_D^1(\Omega)$, it holds that $I_h u \in \mathcal{S}_D^1(\mathcal{T})$ and thus*

$$\min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} \leq \|u - I_h u\|_{H^1(\Omega)} \leq C \sigma(\mathcal{T}) \|h D^2 u\|_{L^2(\Omega)}. \quad (3.24)$$

For $u \in H^2(\Omega) \cap H_*^1(\Omega)$, it holds that

$$\begin{aligned} \min_{v_h \in \mathcal{S}_*^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} &= \min_{v_h \in \mathcal{S}^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} \leq \|u - I_h u\|_{H^1(\Omega)} \\ &\leq C \sigma(\mathcal{T}) \|h D^2 u\|_{L^2(\Omega)}. \end{aligned} \quad (3.25)$$

In either case, the constant $C > 0$ depends only on $\text{diam}(\Omega)$.

Proof. Let $C_{\text{apx}} > 0$ denote the constant from the approximation theorem. Then,

$$\|u - I_h u\|_{H^1(\Omega)}^2 = \|u - I_h u\|_{L^2(\Omega)}^2 + \|\nabla(u - I_h u)\|_{L^2(\Omega)}^2 \leq C_{\text{apx}}^2 (\text{diam}(\Omega)^2 + \sigma(\mathcal{T})^2) \|h D^2 u\|_{L^2(\Omega)}^2.$$

Since $\sigma(\mathcal{T}) \geq 1$, we obtain that

$$\|u - I_h u\|_{H^1(\Omega)} \leq C_{\text{apx}} \sigma(\mathcal{T}) (\text{diam}(\Omega)^2 + 1)^{1/2} \|h D^2 u\|_{L^2(\Omega)}.$$

For $u \in H^2(\Omega) \cap H_D^1(\Omega)$, it holds that $u(z) = 0$ for all $z \in \bar{\Gamma}_D$. This implies that $I_h u \in \mathcal{S}_D^1(\mathcal{T})$ and hence (3.24). Before we prove (3.25), note that $I_h u \in \mathcal{S}^1(\mathcal{T})$ does not belong to $\mathcal{S}_*^1(\mathcal{T})$ in general. However, let $\mathbb{P}_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ denote the H^1 -orthogonal projection onto $\mathcal{S}^1(\mathcal{T})$. Since $1 \in \mathcal{S}^1(\mathcal{T})$, it holds that

$$0 = \int_{\Omega} u \, dx = (u ; 1)_{H^1(\Omega)} = (\mathbb{P}_h u ; 1)_{H^1(\Omega)} = \int_{\Omega} \mathbb{P}_h u \, dx \quad \text{for all } u \in H_*^1(\Omega).$$

Therefore, $\mathbb{P}_h u \in \mathcal{S}_*^1(\mathcal{T})$, and the best approximation property of the orthogonal projection \mathbb{P}_h thus implies that

$$\|u - \mathbb{P}_h u\|_{H^1(\Omega)} = \min_{v_h \in \mathcal{S}_*^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} \leq \min_{v_h \in \mathcal{S}_*^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} \leq \|u - \mathbb{P}_h u\|_{H^1(\Omega)}$$

and hence equality. As before, this proves (3.25). ■

Remark. Corollary 3.6 has two important consequences: First, according to Céa's lemma, the Galerkin error is up to a constant the best approximation error. For a smooth exact solution $u \in H^2(\Omega)$, the P1-FEM thus leads (at least and in fact even) to a convergence order $\mathcal{O}(h)$. Second, $C_D^\infty(\bar{\Omega})$ is dense in $H_D^1(\Omega)$ and $C_*^\infty(\bar{\Omega}) := \{v \in C^\infty(\bar{\Omega}) \mid \int_{\Omega} v \, dx = 0\}$ is dense in $H_*^1(\Omega)$. Corollary 3.6 therefore implies convergence of the Galerkin scheme on a dense subspace. The abstract framework provides convergence of the P1-FEM even without any regularity assumptions on u , cf. Proposition 1.7. □

Exercise 18. Use the Poincaré inequality and the Meyers-Serrin theorem to prove that $C_*^\infty(\overline{\Omega})$ is dense in $H_*^1(\Omega)$. □

3.2.3 Bramble-Hilbert Lemma

It now remains to prove the Approximation Theorem 3.5. The proof of which needs three lemmata. The first two lemmata provide the basis for general scaling arguments. We therefore state the results even in a slightly generalized setting.

Definition. For a multiindex $\alpha \in \mathbb{N}_0^d$ and $x \in \mathbb{R}^d$, we define the **monomial** $x^\alpha := \prod_{j=1}^d x_j^{\alpha_j}$, where $|\alpha| := \sum_{j=1}^d \alpha_j$ is the **(total) degree** of α . For a Lipschitz domain $T \subseteq \mathbb{R}^d$, we define

$$\mathcal{P}^m(T) := \{v : T \rightarrow \mathbb{R} \mid v \text{ is linear combination of monomials of degree } \leq m\} \quad (3.26)$$

the space that consists of all **polynomials** of degree less than or equal to $m \in \mathbb{N}$.

Lemma 3.7 (Bramble-Hilbert). For a Lipschitz domain $T \subset \mathbb{R}^d$ and a normed space X , let $A \in L(H^{m+1}(T); X)$ be a linear and continuous operator with $\mathcal{P}^m(T) \subseteq \ker(A)$. Besides the classical continuity estimate

$$\|Av\|_X \leq \|A\| \|v\|_{H^{m+1}(T)} \quad \text{for all } v \in H^{m+1}(T), \quad (3.27)$$

it holds that

$$\|Av\|_X \leq C \|A\| \|D^{m+1}v\|_{L^2(T)} \quad \text{for all } v \in H^{m+1}(T), \quad (3.28)$$

where the constant $C > 0$ depends only on m and T .

Proof. 1. step. Construct an equivalent norm on $H^{m+1}(T)$: Note that $\mathcal{P}^m(T)$ is a finite dimensional space. Let $\Pi : L^2(T) \rightarrow \mathcal{P}^m(T)$ denote the L^2 -orthogonal projection onto $\mathcal{P}^m(T)$. We define

$$\|v\| := \|D^{m+1}v\|_{L^2(T)} + \|\Pi v\|_{L^2(T)} \quad \text{for } v \in H^{m+1}(T).$$

From $\|\Pi v\|_{L^2(T)} \leq \|v\|_{L^2(T)}$, we infer that

$$\|v\| \leq \|D^{m+1}v\|_{L^2(T)} + \|v\|_{L^2(T)} \leq \sqrt{2} \|v\|_{H^{m+1}(T)}.$$

Next, we prove the converse inequality, i.e., there exists a constant $C > 0$ such that

$$\|v\|_{H^{m+1}(T)} \leq C \|v\| \quad \text{for all } v \in H^{m+1}(T).$$

As above, we use the Rellich theorem and argue by contradiction: If the claim is wrong, we find $v_n \in H^{m+1}(T)$ such that $\|v_n\|_{H^{m+1}(T)} > n \|v_n\|$. We define $w_n := v_n / \|v_n\|_{H^{m+1}(T)}$. Note that

$$\|w_n\|_{H^{m+1}(T)} = 1 \quad \text{as well as} \quad \|w_n\| \leq \frac{1}{n}.$$

According to reflexivity, we may thus assume that $w_n \rightharpoonup w \in H^{m+1}(T)$. According to Lemma 2.7, convexity and continuity of $\|\cdot\|$ imply that $\|w\| = 0$. Therefore, it holds that $D^{m+1}w = 0$ as well as $\Pi w = 0$. With the help of Exercise 19, we deduce that $w \in \mathcal{P}^m(T)$ and consequently $\|w\|_{L^2(T)} = \|\Pi w\|_{L^2(T)} = 0$. According to Rellich's theorem, we have $w_n \rightarrow w = 0 \in H^m(T)$. Since $D^{m+1}w_n \rightarrow 0 \in L^2(T)$, we even conclude that $w_n \rightarrow 0 = w \in H^{m+1}(T)$. This however, contradicts $\|w_n\|_{H^{m+1}(T)} = 1$. Altogether, we have shown that $\|\cdot\|$ is an equivalent norm on $H^{m+1}(T)$.

v6:
20.10.2017

2. step. With the norm equivalence constant $C > 0$ of step 1, it holds that

$$\|Av\|_X = \|A(v - \Pi v)\|_X \leq \|A\| \|v - \Pi v\|_{H^{m+1}(T)} \leq C \|A\| \|v - \Pi v\| = C \|A\| \|D^{m+1}v\|_{L^2(T)}$$

for all $v \in H^{m+1}(T)$. ■

Exercise 19. Prove that a function $v \in H^{m+1}(T)$ on a bounded Lipschitz domain $T \subset \mathbb{R}^d$ satisfies $D^{m+1}v = 0$ if and only if $v \in \mathcal{P}^m(T)$. **Hint:** You should use the case $m = 0$ without a proof, cf. Theorem 2.3. □

3.2.4 Scaling Argument and Proof of Approximation Theorem

Lemma 3.8 (Transformation Formula). Let $T, \hat{T} \subset \mathbb{R}^d$ be Lipschitz domains. Let $\Phi(x) := Bx + y$ with regular matrix $B \in \mathbb{R}^{d \times d}$ and vector $y \in \mathbb{R}^d$ be an affine diffeomorphism with $\Phi(\hat{T}) = T$. For $u \in H^m(T)$, it holds that $u \circ \Phi \in H^m(\hat{T})$ with

$$\|D^m(u \circ \Phi)\|_{L^2(\hat{T})} \leq |\det B|^{-1/2} \|B\|_F^m \|D^m u\|_{L^2(T)}, \quad (3.29)$$

where $\|B\|_F$ denotes the Frobenius norm of B . Moreover, for $m = 0$, there even holds equality.

Proof. 1. step. The case $m = 0$: According to the transformation theorem and $D\Phi(x) = B$, it holds that

$$\|u\|_{L^2(T)}^2 = \int_T u^2 dy = \int_{\hat{T}} (u \circ \Phi)^2 |\det D\Phi| dx = |\det B| \|u \circ \Phi\|_{L^2(\hat{T})}^2.$$

2. step. To treat the higher-order case for smooth functions $u \in C^\infty(\bar{T})$, we first prove by induction on m that for all $j_\ell \in \{1, \dots, d\}$, it holds that

$$\partial_{j_1} \cdots \partial_{j_m} (u \circ \Phi)(x) = \sum_{k_1=1}^d \cdots \sum_{k_m=1}^d \partial_{k_1} \cdots \partial_{k_m} u(\Phi(x)) \prod_{\ell=1}^m B_{k_\ell j_\ell}, \quad (3.30)$$

which is the special case of the Faà di Bruno formula (chain rule for partial derivatives): The case $m = 1$ follows from the chain rule $D(u \circ \Phi)(x) = Du(\Phi(x))D\Phi(x) = Du(\Phi(x))B$, where, e.g., $Du(y) = (\partial_1 u, \dots, \partial_d u)(y)$. Therefore,

$$\partial_j (u \circ \Phi)(x) = \sum_{k=1}^d \partial_k u(\Phi(x)) B_{kj}.$$

Assuming that (3.30) holds up to $m \in \mathbb{N}$, we now prove the equality for $m + 1$:

$$\begin{aligned}
 \partial_{j_1} \cdots \partial_{j_{m+1}}(u \circ \Phi)(x) &\stackrel{!}{=} \partial_{j_1} \left(\sum_{k_2=1}^d \cdots \sum_{k_{m+1}=1}^d \partial_{k_2} \cdots \partial_{k_{m+1}} u(\Phi(x)) \prod_{\ell=2}^{m+1} B_{k_\ell j_\ell} \right) \\
 &= \sum_{k_2=1}^d \cdots \sum_{k_{m+1}=1}^d \partial_{j_1} (\partial_{k_2} \cdots \partial_{k_{m+1}} u(\Phi(x))) \prod_{\ell=2}^{m+1} B_{k_\ell j_\ell} \\
 &\stackrel{!}{=} \sum_{k_2=1}^d \cdots \sum_{k_{m+1}=1}^d \sum_{k_1=1}^d \partial_{k_1} \partial_{k_2} \cdots \partial_{k_{m+1}} u(\Phi(x)) B_{k_1 j_1} \prod_{\ell=2}^{m+1} B_{k_\ell j_\ell} \\
 &= \sum_{k_1=1}^d \cdots \sum_{k_{m+1}=1}^d \partial_{k_1} \partial_{k_2} \cdots \partial_{k_{m+1}} u(\Phi(x)) \prod_{\ell=1}^{m+1} B_{k_\ell j_\ell},
 \end{aligned}$$

where we have used the induction hypothesis for m and the initial step $m = 1$. This verifies (3.30).

3. step. We apply the Cauchy inequality to (3.30) to see that

$$\begin{aligned}
 |\partial_{j_1} \cdots \partial_{j_m}(u \circ \Phi)(x)|^2 &\leq \left(\sum_{k_1=1}^d \cdots \sum_{k_m=1}^d |\partial_{k_1} \cdots \partial_{k_m} u(\Phi(x))|^2 \right) \left(\sum_{k_1=1}^d \cdots \sum_{k_m=1}^d \left| \prod_{\ell=1}^m B_{k_\ell j_\ell} \right|^2 \right) \\
 &= \left(\sum_{k_1=1}^d \cdots \sum_{k_m=1}^d |\partial_{k_1} \cdots \partial_{k_m} u(\Phi(x))|^2 \right) \left(\sum_{k_1=1}^d \cdots \sum_{k_m=1}^d \prod_{\ell=1}^m B_{k_\ell j_\ell}^2 \right) \\
 &\stackrel{!}{=} \left(\sum_{k_1=1}^d \cdots \sum_{k_m=1}^d |\partial_{k_1} \cdots \partial_{k_m} u(\Phi(x))|^2 \right) \left(\prod_{\ell=1}^m \sum_{k_\ell=1}^d B_{k_\ell j_\ell}^2 \right),
 \end{aligned}$$

where the last equality follows from another simple induction argument.

4. step. We prove the transformation formula (3.29) for $u \in C^\infty(\bar{T})$:

$$\begin{aligned}
 |\det B| \|D^m(u \circ \Phi)\|_{L^2(\hat{T})}^2 &= \int_{\hat{T}} \sum_{j_1=1}^d \cdots \sum_{j_m=1}^d |\partial_{j_1} \cdots \partial_{j_m}(u \circ \Phi)(x)|^2 |\det D\Phi(x)| dx \\
 &\leq \underbrace{\left(\sum_{j_1=1}^d \cdots \sum_{j_m=1}^d \prod_{\ell=1}^m \sum_{k_\ell=1}^d B_{k_\ell j_\ell}^2 \right)}_{= \prod_{\ell=1}^m \sum_{j_\ell=1}^d \sum_{k_\ell=1}^d B_{k_\ell j_\ell}^2} \underbrace{\left(\int_{\hat{T}} \sum_{k_1=1}^d \cdots \sum_{k_m=1}^d |\partial_{k_1} \cdots \partial_{k_m} u(\Phi(x))|^2 |\det D\Phi(x)| dx \right)}_{= \|D^m u\|_{L^2(T)}^2} \\
 &= \|B\|_F^{2m} \|D^m u\|_{L^2(T)}^2.
 \end{aligned}$$

5. step. We prove the transformation formula (3.29) for general $u \in H^m(T)$: According to the Meyers-Serrin theorem, $C^\infty(\bar{T})$ is a dense subspace of $H^m(T)$. Note that (3.29) implies for $u \in C^\infty(\bar{T})$ the estimate $\|u \circ \Phi\|_{H^m(\hat{T})} \leq C \|u\|_{H^m(T)}$, where $C > 0$ depends only on m and B . Hence, $\Psi u := u \circ \Phi$ extends uniquely to a linear and continuous mapping $\Psi : H^m(T) \rightarrow H^m(\hat{T})$. For $u \in H^m(T)$, choose $(u_n) \subset C^\infty(\bar{T})$ with $u_n \rightarrow u \in H^m(T)$. By continuity of Ψ , it holds that $u_n \circ \Phi = \Psi u_n \rightarrow \Psi u$ in $H^m(\hat{T})$. Moreover, according to step 1, it holds that $u_n \circ \Phi \rightarrow u \circ \Phi \in L^2(\hat{T})$. This implies that $u \circ \Phi = \Psi u \in H^m(\hat{T})$, i.e., the (unique) extension of Ψ from $C^\infty(\bar{T})$ to $H^m(T)$

is, in fact, the composition. Moreover, the left-hand side and the right-hand side of (3.29) depend continuously (with respect to $H^m(T)$) on u . This and (3.29) for $u_n \in C^\infty(\bar{T})$ prove that

$$\begin{aligned} \|D^m(u \circ \Phi)\|_{L^2(\hat{T})} &= \lim_{n \rightarrow \infty} \|D^m(u_n \circ \Phi)\|_{L^2(\hat{T})} \leq \lim_{n \rightarrow \infty} |\det B|^{-1/2} \|B\|_F^m \|D^m u_n\|_{L^2(T)} \\ &= |\det B|^{-1/2} \|B\|_F^m \|D^m u\|_{L^2(T)} \end{aligned}$$

and conclude the proof. \blacksquare

Lemma 3.9. For $\hat{T} = T_{\text{ref}}$ the reference element and $T = \text{conv}\{z_1, z_2, z_3\} \subset \mathbb{R}^2$ being a non-degenerate triangle, we define

$$\Phi_T : T_{\text{ref}} \rightarrow T, \quad \Phi_T(s, t) := z_1 + B \begin{pmatrix} s \\ t \end{pmatrix}, \quad \text{where } B := \begin{pmatrix} z_2 - z_1 & z_3 - z_1 \end{pmatrix} \in \mathbb{R}^{2 \times 2}. \quad (3.31)$$

Then, it holds that $|\det B| = 2|T|$ and

$$h_T/\sqrt{2} \leq \|B\|_F \leq \sqrt{2} h_T \quad \text{as well as} \quad \varrho_T^{-1}/\sqrt{2} \leq \|B^{-1}\|_F \leq \sqrt{2} \varrho_T^{-1}. \quad (3.32)$$

Proof. It holds that

$$\|B\|_F^2 = |z_2 - z_1|^2 + |z_3 - z_1|^2 \leq 2h_T^2.$$

Moreover,

$$|z_3 - z_2| \leq |z_3 - z_1| + |z_2 - z_1| \leq \sqrt{2} (|z_3 - z_1|^2 + |z_2 - z_1|^2)^{1/2} \leq \sqrt{2} \|B\|_F.$$

In particular, $h_T = \max\{|z_2 - z_1|, |z_3 - z_1|, |z_3 - z_2|\} \leq \sqrt{2} \|B\|_F$. The transformation theorem gives

$$\frac{1}{2} |\det B| = |T_{\text{ref}}| |\det B| = \int_{T_{\text{ref}}} |\det D\Phi_T| dx = \int_T dx = |T| > 0.$$

Hence, $0 < |\det B| = 2|T| = h_T \varrho_T$. In particular, B^{-1} as well as ϱ_T^{-1} are well-defined. It holds that

$$B^{-1} = \frac{1}{\det B} \begin{pmatrix} b_{22} & -b_{12} \\ -b_{21} & b_{11} \end{pmatrix} \quad \text{for} \quad B = \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix}.$$

In particular, this proves that

$$\|B^{-1}\|_F = \frac{\|B\|_F}{|\det B|} = \frac{\|B\|_F}{h_T \varrho_T},$$

and the second estimate in (3.32) follows from the first. \blacksquare

Proof of Approximation Theorem 3.5. 1. step. Estimate on the reference element T_{ref} : Let $I_h^{\text{ref}} : H^2(T_{\text{ref}}) \rightarrow \mathcal{P}^1(T_{\text{ref}})$ denote the nodal interpolation operator on the reference element. We consider the operator

$$A := 1 - I_h^{\text{ref}} : H^2(T_{\text{ref}}) \rightarrow H^k(T_{\text{ref}}) \quad \text{for } k = 0, 1$$

and observe that $\mathcal{P}^1(T_{\text{ref}}) \subseteq \ker(A)$. To see that A is continuous, we estimate

$$\|Av\|_{H^k(T_{\text{ref}})} \leq \|v\|_{H^2(T_{\text{ref}})} + \|I_h^{\text{ref}}v\|_{H^k(T_{\text{ref}})}.$$

Let z_1, z_2, z_3 denote the nodes of the reference element. Since all norms on the finite dimensional space $\mathcal{P}^1(T_{\text{ref}})$ are equivalent, we use the Sobolev inequality to see that

$$\|I_h^{\text{ref}}v\|_{H^k(T_{\text{ref}})} \leq C_{\text{norm}} \max_{j=1,\dots,3} |I_h^{\text{ref}}v(z_j)| \leq C_{\text{norm}} \|v\|_{\infty, T_{\text{ref}}} \leq C_{\text{norm}} C_{\text{sobolev}} \|v\|_{H^2(T_{\text{ref}})}.$$

Altogether, we obtain that $\|Av\|_{H^k(T_{\text{ref}})} \leq (1 + C_{\text{norm}} C_{\text{sobolev}}) \|v\|_{H^2(T_{\text{ref}})}$, whence continuity of the operator A . Consequently, the Bramble-Hilbert lemma provides a constant $C_{\text{ref}} > 0$ that depends only on T_{ref} with

$$\|v - I_h^{\text{ref}}v\|_{H^k(T_{\text{ref}})} \leq C_{\text{ref}} \|D^2v\|_{L^2(T_{\text{ref}})} \quad \text{for all } v \in H^2(T_{\text{ref}}) \text{ and } k = 0, 1.$$

2. step. Scaling arguments provide the estimate on each element T : Let $\Phi = \Phi_T$ denote the affine diffeomorphism from Lemma 3.9. Note that $I_h^{\text{ref}}(u \circ \Phi) = (I_h u) \circ \Phi$. Define $v := u \circ \Phi$ and observe that $(u - I_h u) \circ \Phi = (1 - I_h^{\text{ref}})v$. First, we apply the transformation formula to Φ^{-1} ,

$$\begin{aligned} \|D^k(u - I_h u)\|_{L^2(T)} &= \|D^k((v - I_h^{\text{ref}}v) \circ \Phi^{-1})\|_{L^2(T)} \\ &\leq |\det B^{-1}|^{-1/2} \|B^{-1}\|_F^k \|D^k(v - I_h^{\text{ref}}v)\|_{L^2(T_{\text{ref}})} \\ &\leq C_{\text{ref}} |\det B|^{1/2} \|B^{-1}\|_F^k \|D^2v\|_{L^2(T_{\text{ref}})}. \end{aligned}$$

Second, we plug-in $v = u \circ \Phi$ and apply the transformation formula to Φ ,

$$\|D^2v\|_{L^2(T_{\text{ref}})} = \|D^2(u \circ \Phi)\|_{L^2(T_{\text{ref}})} \leq |\det B|^{-1/2} \|B\|_F^2 \|D^2u\|_{L^2(T)}.$$

The combination of the last two estimates proves that

$$\|D^k(u - I_h u)\|_{L^2(T)} \leq C_{\text{ref}} \|B^{-1}\|_F^k \|B\|_F^2 \|D^2u\|_{L^2(T)} \leq C_{\text{ref}} 2^{(k+2)/2} h_T^2 \varrho_T^{-k} \|D^2u\|_{L^2(T)},$$

where we have used the geometric interpretation of $\|B\|_F$ and $\|B^{-1}\|_F$. This proves that

$$\|u - I_h u\|_{L^2(T)} \leq 2C_{\text{ref}} \|h^2 D^2u\|_{L^2(T)} \quad \text{and} \quad \|\nabla(u - I_h u)\|_{L^2(T)} \leq 2^{3/2} C_{\text{ref}} \sigma(\mathcal{T}) \|h D^2u\|_{L^2(T)}.$$

and thus concludes the proof. \blacksquare

Remark. The proof of Theorem 3.5 shows that it is enough to assume $u \in C(\overline{\Omega}) \cap H^2(\mathcal{T})$, where $H^k(\mathcal{T}) := \{u \in L^2(\Omega) \mid \forall T \in \mathcal{T} \quad u|_T \in H^k(T)\}$ for $k \geq 1$. According to the Sobolev inequality, it holds that $H^2(\Omega) \subseteq C(\overline{\Omega}) \cap H^2(\mathcal{T})$. For the *broken Sobolev spaces* $H^k(\mathcal{T})$, we write $D_h^k v$ for the \mathcal{T} -piecewise k -th derivative of v and, in particular, $\nabla_h v = D_h^1 v$ for the \mathcal{T} -piecewise gradient. \square

Remark. We recall the procedure of a scaling argument for proving an estimate. To that end, let $\Phi_T : T_{\text{ref}} \rightarrow T$ be the affine diffeomorphism with linear part B .

- First, transfer the left-hand side from T to T_{ref} :

$$\begin{aligned} \|D^k v\|_{L^2(T)} &= \|D^k(v \circ \Phi_T \circ \Phi_T^{-1})\|_{L^2(T)} \leq |\det B^{-1}|^{-1/2} \|B^{-1}\|_F^k \|D^k(v \circ \Phi_T)\|_{L^2(T_{\text{ref}})} \\ &\simeq |T| \varrho_T^{-k} \|D^k(v \circ \Phi_T)\|_{L^2(T_{\text{ref}})}, \end{aligned}$$

i.e., derivative on the left-hand side give rise to negative powers of ϱ_T .

- Second, prove an appropriate estimate on the reference element T_{ref} .
- Third, transfer the right-hand side from T_{ref} to T :

$$\|D^\ell(w \circ \Phi_T)\|_{L^2(T_{\text{ref}})} \leq |\det B|^{-1/2} \|B\|_F^\ell \|D^\ell w\|_{L^2(T)} \simeq |T|^{-1/2} h_T^\ell \|D^\ell w\|_{L^2(T)},$$

i.e., derivatives on the right-hand side give rise to positive powers of h_T .

Plugging everything together, proves the desired estimate. \square

Note that the heart of the proof of the approximation theorem is the Rellich theorem and thus a compactness argument. The following exercise shows that approximation results are necessarily proved by use of compactness.

Exercise 20. Let X be a Banach space and Y be a normed space with continuous inclusion $Y \subseteq X$. For $h \rightarrow 0$, let X_h be finite dimensional subspaces of X and $I_h \in L(Y; X_h)$ be a continuous and linear operator with

$$\|u - I_h u\|_X \leq C h^\alpha \|u\|_Y \quad \text{for all } u \in Y,$$

where the constants $C, \alpha > 0$ are independent of u and h . Then, the continuous inclusion $Y \subseteq X$ is already compact. \square

v7:
25.10.2017

3.3 Inverse Estimates and Optimality of Approximation Results

Theorem 3.10. For all polynomial degrees $m \in \mathbb{N}$, there exists a constant $C > 0$ such that

$$\|\nabla v_h\|_{L^2(T)} \leq C \|\varrho^{-1} v_h\|_{L^2(T)} \quad \text{for all } v_h \in \mathcal{P}^m(\mathcal{T}) \text{ and all } T \in \mathcal{T}, \quad (3.33)$$

where $\mathcal{P}^m(\mathcal{T}) := \{v_h : \Omega \rightarrow \mathbb{R} \mid \forall T \in \mathcal{T} \quad v_h|_T \in \mathcal{P}^m(T)\}$. With $\nabla_h(\cdot)$ being the \mathcal{T} -piecewise gradient, this implies that

$$\|h^\alpha \nabla_h v_h\|_{L^2(\Omega)} \leq C \sigma(\mathcal{T}) \|h^{\alpha-1} v_h\|_{L^2(\Omega)} \quad \text{for all } v_h \in \mathcal{P}^m(\mathcal{T}) \text{ and } \alpha \in \mathbb{R}. \quad (3.34)$$

The constant $C > 0$ depends only on m , but neither on Ω nor \mathcal{T} .

Proof. The proof is done \mathcal{T} -elementwise and follows from a scaling argument. We start with an abstract observation.

1. step. Let X be a finite dimensional space, $\|\cdot\|_X$ be a norm on X and $|\cdot|_X$ be a seminorm on X . Then, there exists a constant $C > 0$ such that

$$|x|_X \leq C \|x\|_X \quad \text{for all } x \in X :$$

We consider the quotient space X/Y , where $Y := \{x \in X \mid |x|_X = 0\}$. Note that X/Y is finite dimensional and that

$$\|x + Y\|_{X/Y} := \inf_{y \in Y} \|x + y\|_X \quad \text{as well as} \quad |x + Y|_{X/Y} := \inf_{y \in Y} |x + y|_X = |x|_X$$

are norms on the finite dimensional space X/Y . Therefore, there is a norm equivalence constant $C > 0$ such that

$$|x|_X = |x + Y|_{X/Y} \leq C \|x + Y\|_{X/Y} \leq C \|x\|_X \quad \text{for all } x \in X.$$

2. step. There exists a constant $C_{\text{ref}} > 0$ such that

$$\|\nabla w_h\|_{L^2(T_{\text{ref}})} \leq C_{\text{ref}} \|w_h\|_{L^2(T_{\text{ref}})} \quad \text{for all } w_h \in \mathcal{P}^m(T_{\text{ref}}).$$

This follows from the abstract framework for $X = \mathcal{P}^m(T_{\text{ref}})$.

3. step. For each element $T \in \mathcal{T}$, it holds that

$$\|\nabla v_h\|_{L^2(T)} \leq \sqrt{2} C_{\text{ref}} \varrho_T^{-1} \|v_h\|_{L^2(T)} \leq \sqrt{2} C_{\text{ref}} \sigma(\mathcal{T}) h_T^{-1} \|v_h\|_{L^2(T)} \quad \text{for all } v_h \in \mathcal{P}^m(T):$$

Let $\Phi : T_{\text{ref}} \rightarrow T$ be an affine diffeomorphism and $B \in \mathbb{R}^{2 \times 2}$ its linear part. We apply the transformation formula to Φ^{-1} to see that

$$\|\nabla v_h\|_{L^2(T)} \leq |\det B^{-1}|^{-1/2} \|B^{-1}\|_F \|\nabla(v_h \circ \Phi)\|_{L^2(T_{\text{ref}})}.$$

Note that the L^2 -norm can be estimated by step 2 since $v_h \circ \Phi \in \mathcal{P}^m(T_{\text{ref}})$. The application of the transformation formula to Φ proves that

$$\|v_h \circ \Phi\|_{L^2(T_{\text{ref}})} = |\det B|^{-1/2} \|v_h\|_{L^2(T)}.$$

By definition of the shape regularity constant $\sigma(\mathcal{T})$, we obtain that

$$\|\nabla v_h\|_{L^2(T)} \leq C_{\text{ref}} \|B^{-1}\|_F \|v_h\|_{L^2(T)} \leq \sqrt{2} C_{\text{ref}} \varrho_T^{-1} \|v_h\|_{L^2(T)} \leq \sqrt{2} C_{\text{ref}} \sigma(\mathcal{T}) h_T^{-1} \|v_h\|_{L^2(T)},$$

where we have used that $\|B^{-1}\|_F \leq \sqrt{2} \varrho_T^{-1}$. In particular, this proves (3.33).

4. step. Multiplying the estimate in step 3 by h_T^α , we finally sum the (squared) estimates over all $T \in \mathcal{T}$ to obtain (3.34). ■

Exercise 21. Suppose that $|\cdot|_1$ and $|\cdot|_2$ are seminorms on a finite dimensional space X . Prove that the following two claims are equivalent:

- (i) There exists a constant $C > 0$ such that $|x|_1 \leq C |x|_2$ for all $x \in X$.
- (ii) There holds nestedness of the kernels $\{x \in X \mid |x|_2 = 0\} \subseteq \{x \in X \mid |x|_1 = 0\}$.

Hint. Refine the abstract argument in the proof of Theorem 3.11 □

Exercise 22. We consider the \mathcal{T} -piecewise constant functions

$$\mathcal{P}^0(\mathcal{T}) := \{v_h : \Omega \rightarrow \mathbb{R} \mid \forall T \in \mathcal{T} \quad v_h|_T \in \mathcal{P}^0(T)\}.$$

Prove that, for $v \in L^2(\Omega)$, the function $v_{\mathcal{T}} \in \mathcal{P}^0(\mathcal{T})$ defined by $v_{\mathcal{T}}|_T := (1/|T|) \int_T v_h dx$, is the L^2 -best approximation, i.e.,

$$\|v - v_{\mathcal{T}}\|_{L^2(\Omega)} = \min_{v_h \in \mathcal{P}^0(\mathcal{T})} \|v - v_h\|_{L^2(\Omega)}.$$

Use a scaling argument to prove the **Poincaré** inequality

$$\|h^\alpha(v - v_{\mathcal{T}})\|_{L^2(\Omega)} \leq C \|h^{1+\alpha}\nabla_h v\|_{L^2(\Omega)} \quad \text{for all } v \in H^1(\mathcal{T}) \text{ and } \alpha \in \mathbb{R}, \quad (3.35)$$

where the constant C does neither depend on Ω nor on \mathcal{T} , v , or α . \square

Remark. The combination of the inverse estimate (3.34) and the Poincaré inequality (3.35) shows that the powers of h are optimal provided that $h \rightarrow 0$ for the mesh-size. To see this, we proceed as follows: Let $\alpha \in \mathbb{R}$ be the optimal exponent for the Poincaré inequality, i.e.,

$$\|v - v_{\mathcal{T}}\|_{L^2(\Omega)} \leq C_1 \|h^\alpha \nabla_h v\|_{L^2(\Omega)} \quad \text{for all } v \in H^1(\mathcal{T}).$$

According to (3.35), it holds that $\alpha \geq 1$. Let $\beta \in \mathbb{R}$ be the optimal exponent for the inverse estimate, i.e.,

$$\|h^\beta \nabla_h v_h\|_{L^2(\Omega)} \leq C_2 \|v_h\|_{L^2(\Omega)} \quad \text{for all } v_h \in \mathcal{P}^m(\mathcal{T}).$$

According to the inverse estimate (3.34), it holds that $\beta \leq 1$. This, however, implies that

$$C_2^{-1} \|h^\beta \nabla v\|_{L^2(\Omega)} = C_2^{-1} \|h^\beta \nabla_h(v - v_{\mathcal{T}})\|_{L^2(\Omega)} \leq \|v - v_{\mathcal{T}}\|_{L^2(\Omega)} \leq C_1 \|h^\alpha \nabla v\|_{L^2(\Omega)}$$

for all $v \in \mathcal{P}^1(\Omega)$. Consequently, it holds that $\beta \geq \alpha$, whence $\alpha = 1 = \beta$. In particular, we see that order 1 is attained for a non-constant affine function $v \in \mathcal{P}^1(\Omega)$ with $\nabla v \neq 0$. \square

With the same arguments as in the preceding remark, one can prove that the convergence of the P1-FEM is $\mathcal{O}(h)$ and that this order is optimal. We state this observation in the following exercises:

Exercise 23. Assume that \mathcal{T} is a regular triangulation with $\|h\|_{L^\infty(\Omega)} \leq 1$. For a polynomial degree $m \in \mathbb{N}$, it holds that

$$\|h^\alpha D_h^2 v_h\|_{L^2(\Omega)} \leq C \|h^{\alpha-1} v_h\|_{H^1(\mathcal{T})} \quad \text{for all } v_h \in \mathcal{P}^m(\mathcal{T}) \text{ and } \alpha \in \mathbb{R}. \quad (3.36)$$

The constant $C > 0$ depends only on m and the shape regularity $\sigma(\mathcal{T})$ of \mathcal{T} . \square

Exercise 24. Use Exercise 23 and the Approximation Theorem 3.5 to conclude that the optimal exponents $\alpha, \beta, \gamma, \delta > 0$ in the estimates

- $\|u - I_h u\|_{H^1(\Omega)} \leq C_1 \|h^\alpha D^2 u\|_{L^2(\Omega)}$ for all $u \in H^2(\Omega)$,
- $\|h^\beta D_h^2 v_h\|_{L^2(\Omega)} \leq C_2 \|v_h\|_{H^1(\mathcal{T})}$ for all $v_h \in \mathcal{P}^m(\mathcal{T})$,
- $\min_{v_h \in \mathcal{S}^1(\mathcal{T})} \|u - v_h\|_{H^1(\Omega)} \leq C_3 \|h^\gamma D^2 u\|_{L^2(\Omega)}$ for all $u \in H^2(\Omega)$
- $\min_{v_h \in \mathcal{P}^1(\mathcal{T})} \|u - v_h\|_{H^1(\mathcal{T})} \leq C_4 \|h^\delta D_h^2 u\|_{L^2(\Omega)}$ for all $u \in H^2(\mathcal{T})$

satisfy that $\alpha = \beta = \gamma = \delta = 1$. Conclude that, for $u \in H^2(\Omega)$, the order $\mathcal{O}(h)$ for the convergence of the P1-FEM cannot be improved in general. \square

Another important implication of inverse estimates is the proof of bounds for the condition number

$$\text{cond}(A) := \|A\|_2 \|A^{-1}\|_2$$

of the stiffness matrix $A \in \mathbb{R}_{\text{sym}}^{N \times N}$ with respect to the operator norm $\|\cdot\|_2$ induced by the usual Euclidean norm $|\cdot|$ on \mathbb{R}^N .

Theorem 3.11. *Let \mathcal{T} be a regular triangulation. Let $\{\zeta_1, \dots, \zeta_N\}$ be the hat function basis of $\mathcal{S}_D^1(\mathcal{T})$ and $A \in \mathbb{R}_{\text{sym}}^{N \times N}$ be the corresponding stiffness matrix, i.e., $A_{jk} = \int_{\Omega} \nabla \zeta_j \cdot \nabla \zeta_k \, dx$. Then, it holds that*

$$\text{cond}(A) \leq C h_{\min}^{-2}, \quad (3.37)$$

where $h_{\min} := \min_{T \in \mathcal{T}} h_T$ and where the constant $C > 0$ depends only on $\sigma(\mathcal{T})$.

Proof. We have already seen that A is symmetric and positive definite. Hence, the condition number of A reads

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

where λ_{\max} and λ_{\min} denote the maximal and minimal eigenvalue of A , respectively. Moreover, both eigenvalues can be represented by use of the Rayleigh quotient

$$\lambda_{\max} = \max_{x \in \mathbb{R}^N \setminus \{0\}} \frac{x \cdot Ax}{|x|^2} \quad \text{and} \quad \lambda_{\min} = \min_{x \in \mathbb{R}^N \setminus \{0\}} \frac{x \cdot Ax}{|x|^2}.$$

We thus aim to provide constants $C_1, C_2 > 0$ such that

$$C_1^{-1} h_{\min}^2 |x|^2 \leq x \cdot Ax \leq C_2 |x|^2 \quad \text{for all } x \in \mathbb{R}^N. \quad (3.38)$$

This results in $\lambda_{\min} \geq C_1^{-1} h_{\min}^2$ and $\lambda_{\max} \leq C_2$ and hence $\text{cond}(A) \leq C_1 C_2 h_{\min}^{-2}$.

From now on, let $x \in \mathbb{R}^N$ be arbitrary. Note that x corresponds to a numbering $\mathcal{K} \setminus \bar{\Gamma}_D = \{z_1, \dots, z_N\}$ of the free nodes of \mathcal{T} . For z_j , let ζ_j denote the corresponding hat function and $v_h := \sum_{j=1}^N x_j \zeta_j \in \mathcal{S}_D^1(\mathcal{T})$ be the discrete function associated with the coefficient vector $x \in \mathbb{R}^N$. Recall that $x \cdot Ax = \|\nabla v_h\|_{L^2(\Omega)}^2$ by definition of A .

1. step. Proof of lower estimate in (3.38): For a fixed element $T \in \mathcal{T}$ and $\Phi : T_{\text{ref}} \rightarrow T$ the associated affine transformation with linear part B , it holds that

$$\begin{aligned} \|v_h\|_{L^2(T)}^2 &= |\det B^{-1}|^{-1} \|v_h \circ \Phi\|_{L^2(T_{\text{ref}})}^2 \geq 2|T| C_{\text{norm}}^2 \sum_{z \in \mathcal{K}_{\text{ref}}} |v_h \circ \Phi(z)|^2 \\ &= 2|T| C_{\text{norm}}^2 \sum_{z \in \mathcal{K}_T} |v_h(z)|^2, \end{aligned}$$

where $C_{\text{norm}} > 0$ stems from norm equivalence on the finite dimensional space $\mathcal{P}^1(T_{\text{ref}})$. With

$$|T| = \frac{1}{2} h_T \varrho_T = \frac{1}{2} h_T^2 \sigma(T)^{-1} \geq \frac{1}{2} h_T^2 \sigma(\mathcal{T})^{-1},$$

this implies that

$$\begin{aligned} \|v_h\|_{L^2(\Omega)}^2 &= \sum_{T \in \mathcal{T}} \|v_h\|_{L^2(T)}^2 \geq 2 C_{\text{norm}}^2 \sum_{T \in \mathcal{T}} |T| \sum_{z \in \mathcal{K}_T} |v_h(z)|^2 \\ &\geq C_{\text{norm}}^2 \sigma(\mathcal{T})^{-1} \sum_{T \in \mathcal{T}} h_T^2 \sum_{z \in \mathcal{K}_T} |v_h(z)|^2 \\ &\geq C_{\text{norm}}^2 \sigma(\mathcal{T})^{-1} h_{\min}^2 \sum_{T \in \mathcal{T}} \sum_{z \in \mathcal{K}_T} |v_h(z)|^2. \end{aligned}$$

We stress that some nodes $z \in \mathcal{K}$ are counted several times. Therefore,

$$\|v_h\|_{L^2(\Omega)}^2 \geq C_{\text{norm}}^2 \sigma(\mathcal{T})^{-1} h_{\min}^2 \sum_{z \in \mathcal{K}} |v_h(z)|^2 = C_{\text{norm}}^2 \sigma(\mathcal{T})^{-1} h_{\min}^2 |x|^2,$$

since $x_j = v_h(z_j)$ for a free node $z_j \in \mathcal{K} \setminus \bar{\Gamma}_D$ and since $v_h(z) = 0$ for $z \in \mathcal{K} \cap \bar{\Gamma}_D$. Finally, the Friedrichs inequality proves that

$$C_{\text{norm}}^2 \sigma(\mathcal{T})^{-1} h_{\min}^2 |x|^2 \leq \|v_h\|_{L^2(\Omega)}^2 \leq C_F^2 \|\nabla v_h\|_{L^2(\Omega)}^2 = C_F^2 (x \cdot Ax).$$

2. step. Proof of upper estimate in (3.38): We start from the inverse estimate

$$x \cdot Ax = \|\nabla v_h\|_{L^2(\Omega)}^2 \leq C_{\text{inv}}^2 \|\varrho^{-1} v_h\|_{L^2(\Omega)}^2 = C_{\text{inv}}^2 \sum_{T \in \mathcal{T}} \varrho_T^{-2} \|v_h\|_{L^2(T)}^2.$$

As before, norm equivalence on $\mathcal{P}^1(T_{\text{ref}})$ provides a constant $\tilde{C}_{\text{norm}} > 0$ with

$$\|\nabla v_h\|_{L^2(\Omega)}^2 \leq 2 C_{\text{inv}}^2 \tilde{C}_{\text{norm}}^2 \sum_{T \in \mathcal{T}} \varrho_T^{-2} |T| \sum_{z \in \mathcal{K}_T} |v_h(z)|^2 \leq C_{\text{inv}}^2 \tilde{C}_{\text{norm}}^2 \sigma(\mathcal{T}) \sum_{T \in \mathcal{T}} \sum_{z \in \mathcal{K}_T} |v_h(z)|^2.$$

Since $\sigma(\mathcal{T})$ provides a lower bound for the minimal interior angle of an element, each node $z \in \mathcal{K}$ is node of at most $C_{\text{count}} > 0$ elements $T \in \mathcal{T}$. Altogether, this proves that

$$x \cdot Ax \leq C_{\text{inv}}^2 \tilde{C}_{\text{norm}}^2 C_{\text{count}} \sigma(\mathcal{T}) \sum_{z \in \mathcal{K}} |v_h(z)|^2.$$

The proof is concluded by $\sum_{z \in \mathcal{K}} |v_h(z)|^2 = |x|^2$ which we already noted before. ■

Remark. A sequence $(\mathcal{T}_\ell)_{\ell \in \mathbb{N}}$ of meshes is **quasi-uniform** provided that

- the shape regularity constants are uniformly bounded, i.e.,

$$\gamma := \sup_{\ell \in \mathbb{N}} \sigma(\mathcal{T}_\ell) < \infty,$$

- the quotients of maximal and minimal mesh-size are uniformly bounded, i.e

$$q := \sup_{\ell \in \mathbb{N}} \frac{h_{\max}(\mathcal{T}_\ell)}{h_{\min}(\mathcal{T}_\ell)} < \infty.$$

where $h_{\max}(\mathcal{T}_\ell) := \max_{T \in \mathcal{T}_\ell} h_T$ and $h_{\min}(\mathcal{T}_\ell) := \min_{T \in \mathcal{T}_\ell} h_T$, respectively.

In this case, the condition estimate (3.37) cannot be improved, i.e., the power of h is optimal. To realize this, the reader should check the proof of Theorem 3.11 to see that all estimates, including the inverse estimate, are sharp with respect to the powers of h . \square

Remark. For \mathbb{R}^d and $|T| \simeq h_T^d$, step 1 of the proof of Theorem 3.11 shows $h_{\min}^d |x|^2 \lesssim x \cdot Ax$. Moreover, since $|T| \simeq \varrho_T^d \simeq h_T^d$, step 2 yields $x \cdot Ax \lesssim h_{\max}^{d-2}$. Overall, one thus obtains that

$$\text{cond}(A) \lesssim \frac{h_{\max}^{d-2}}{h_{\min}^d} = \left(\frac{h_{\max}}{h_{\min}} \right)^{d-2} h_{\min}^{-2},$$

i.e., the power 2 does not stem from the dimension, but only from the order of the differential operator. \square

3.4 MATLAB Implementation (31.10.2017)

Throughout, the idea will be to represent the Galerkin solution u_h in the nodal basis of the entire discrete space $\mathcal{S}^1(\mathcal{T})$. In particular, we build — in any of the three cases (3.14)–(3.16) — the Galerkin matrix A and the right-hand side b of Theorem 1.4 for the entire nodal basis $\mathcal{B} = \{\zeta_z \mid z \in \mathcal{K}\}$. The so-called assembly is done elementwise, and we stress that this is of (almost) linear complexity with respect to the number of elements. With $N := \#\mathcal{K}$, we assume a numbering $\mathcal{K} = \{z_1, \dots, z_N\}$ of the nodes. The hat functions are numbered accordingly $\mathcal{B} = \{\zeta_1, \dots, \zeta_N\}$.

v8:
27.10.2017

3.4.1 The Data Structure

The usual FEM implementation is done as follows:

- With $N = \#\mathcal{K}$ the number of nodes, `coordinates` is a double $N \times 2$ array that contains the coordinates of the nodes: The k -th row `coordinates(k, :)` gives the coordinates of $z_k \in \mathbb{R}^2$.
- With $M = \#\mathcal{T}$ the number of elements, `elements` is an integer $M \times 3$ array that contains the numbers of the nodes of all triangles: The j -th row `elements(j, :)` provides the numbers of the nodes of the element $T_j = \text{conv}\{x_T, y_T, z_T\}$ with respect to `coordinates`. In particular, this fixes an ordering $\mathcal{T} = \{T_1, \dots, T_M\}$ of the elements.
- With $d = \#\mathcal{E}_D$ the number of Dirichlet edges, `dirichlet` is an integer $d \times 2$ array that contains the numbers of the nodes of all edges on the Dirichlet boundary: The ℓ -th row `dirichlet(l, :)` provides the numbers of the nodes of an edge $E_\ell = \text{conv}\{x_E, y_E\} \in \mathcal{E}_D$ with respect to `coordinates`. Implicitly, we are given a numbering of the Dirichlet edges \mathcal{E}_D .
- With $n = \#\mathcal{E}_N$ the number of Neumann edges, `neumann` is an integer $n \times 2$ array that contains the numbers of the nodes of all edges on the Neumann boundary: The ℓ -th row `neumann(l, :)` provides the numbers of the nodes of an edge $E_\ell = \text{conv}\{x_E, y_E\} \in \mathcal{E}_N$ with respect to `coordinates`.

Conventions on Array `elements`: It is a usual convention to store the nodes of a triangle T counter clockwise. This allows to compute the outer normal on ∂T by a closed formulae, cf. the Triangle $T_2 = \text{conv}\{z_1, z_2, z_5\}$ which is stored in the form `element(2, :) = [5 1 2]` in Figure 3.5 and Figure 3.6. Moreover, the numbering is chosen in such a way that the first edge is always some

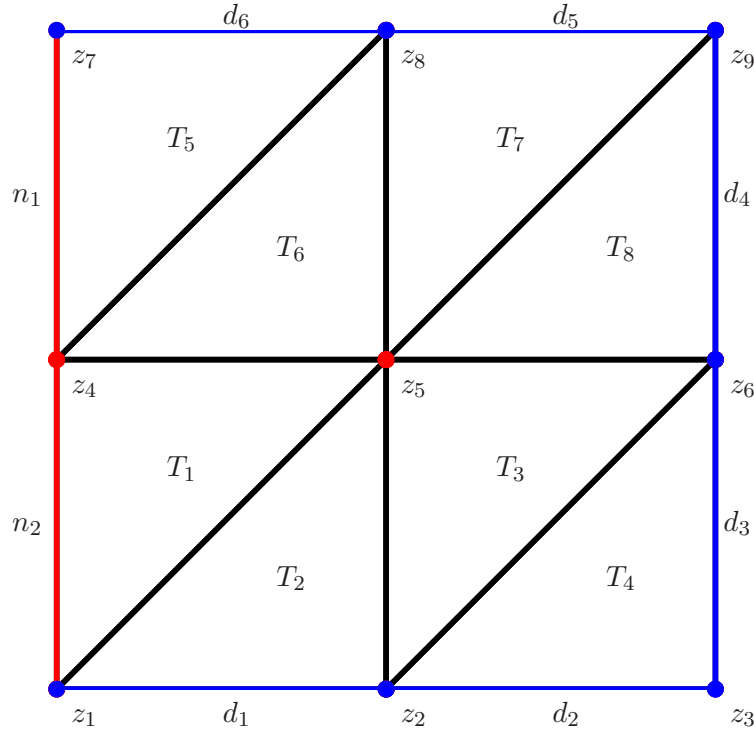


FIGURE 3.5. Example of a regular triangulation $\mathcal{T} = \{T_1, \dots, T_8\}$. The triangulation has nine nodes $\mathcal{K} = \{z_1, \dots, z_9\}$. The Dirichlet boundary $\mathcal{E}_D = \{d_1, \dots, d_6\}$ is indicated by blue color, the Neumann boundary $\mathcal{E}_N = \{n_1, n_2\}$ is indicated by red color. The blue bullets z_ℓ denote nodes, where the Dirichlet data $u_h(z_\ell) = 0$ is prescribed, whereas red bullets denote free nodes $z_\ell \in \mathcal{K} \setminus \overline{\Gamma}_D = \{z_4, z_5\}$, where $u_h(z_\ell)$ has to be computed.

special edge of T , e.g., the longest edge. In our example, the longest edge of T_2 is $E = \text{conv}\{z_1, z_5\}$. This convention is useful for the mesh-refinement, where one has to take care of the angles (i.e., shape regularity).

Conventions on Array `dirichlet`: The nodes are stored in such an order that the tangential vector directed from the first node to the second node is counter clockwise with respect to the boundary $\partial\Omega$, cf. $d_5 = \text{conv}\{z_8, z_9\}$ which is stored in the form `dirichlet(5,:) = [9 8]` in Figure 3.5 and Figure 3.6.

Conventions on Array `neumann`: We apply the same conventions as for the array `dirichlet`.

3.4.2 Building the Right-Hand Side

Let $\mathcal{K} = \{z_1, \dots, z_N\}$ with corresponding hat functions ζ_k . We aim to compute the vector $b \in \mathbb{R}^N$ of the right-hand side of the Galerkin scheme, cf. Theorem 1.4. For fixed $k = 1, \dots, N$, a coefficient b_k of the vector b reads

$$b_k := (f ; \zeta_k)_{L^2(\Omega)} + (\phi ; \zeta_k)_{L^2(\Gamma_N)} = \sum_{T \in \mathcal{T}} \int_T f \zeta_k dx + \sum_{E \in \mathcal{E}_N} \int_E \phi \zeta_k ds,$$

% elements.dat	% coordinates.dat	% dirichlet.dat	% neumann.dat
1 5 4	0.0 0.0	1 2	7 4
5 1 2	0.5 0.0	2 3	4 1
2 6 5	1.0 0.0	3 6	
6 2 3	0.0 0.5	6 9	
4 8 7	0.5 0.5	9 8	
8 4 5	1.0 0.5	8 7	
5 9 8	0.0 1.0		
9 5 6	0.5 1.0		
	1.0 1.0		

FIGURE 3.6. The data files for the regular triangulation \mathcal{T} from Figure 3.5. In MATLAB the command, e.g., `load elements.dat` creates the array `elements` described in Section 3.4.1.

where $\mathcal{E}_N = \emptyset$ in case of the Dirichlet problem. The integrals are usually *not* computed analytically. Instead, one uses numerical quadrature. We shall see below that — according to the Strang lemma — a one point quadrature is sufficient and does not disturb the convergence order of the P1-FEM: For a triangle $T = \text{conv}\{x_T, y_T, z_T\} \in \mathcal{T}$, we use the center of mass $s_T := \frac{1}{3}(x_T + y_T + z_T)$ and approximate

$$\int_T f \zeta_k dx \approx |T| f(s_T) \zeta_k(s_T).$$

According to the definition of ζ_k , it holds that

$$\zeta_k(s_T) = \begin{cases} 0 & \text{if } z_k \notin \{x_T, y_T, z_T\}. \\ 1/3 & \text{else.} \end{cases}$$

For an edge $E = \text{conv}\{x_E, y_E\} \in \mathcal{E}_N$, we consider the midpoint $m_E := \frac{1}{2}(x_E + y_E)$ and approximate

$$\int_E \phi \zeta_k ds \approx h_E \phi(m_E) \zeta_k(m_E),$$

where h_E denotes the length of E . Again, the definition of ζ_k shows that

$$\zeta_k(m_E) = \begin{cases} 0 & \text{if } z_k \notin \{x_E, y_E\}. \\ 1/2 & \text{else.} \end{cases}$$

The pseudo MATLAB code below computes the approximate right-hand side

$$b_k := \sum_{T \in \mathcal{T}} [|T| f(s_T) \zeta_k(s_T)] + \sum_{E \in \mathcal{E}_N} [h_E \phi(m_E) \zeta_k(m_E)] \quad \text{for all } k = 1, \dots, N$$

The assembly is done \mathcal{T} -elementwise resp. \mathcal{E}_N -edgewise since the sums appear to consist of 3 resp. 2 terms only:

```

N = size(coordinates,1);
b = zeros(N,1);
for j = 1:size(elements,1)
    nodes = elements(j,:);

```

```

    b(nodes) = b(nodes) + |Tj|/3*f(sTj);
end
for j = 1:size(neumann,1)
    nodes = neumann(j,:);
    b(nodes) = b(nodes) + hEj/2*phi(mEj);
end

```

Of course, the quadrature is only allowed if the data f and ϕ are (at least \mathcal{T} -piecewise) continuous in Ω and on Γ , respectively. However, this is the usual case in practice.

If $E = \text{conv}\{x_E, y_E\}$ is an edge, its length is simply computed by $h_E = |x_E - y_E|$. In this sense, it only remains to provide a formula for the computation of $|T|$.

Lemma 3.12. *For a triangle $T = \text{conv}\{z_1, z_2, z_3\}$ with nodes $z_1, z_2, z_3 \in \mathbb{R}^2$, there holds*

$$|T| = \frac{1}{2} \left| \det \begin{pmatrix} z_2 - z_1 & z_3 - z_1 \end{pmatrix} \right| = \frac{1}{2} \left| \det \begin{pmatrix} 1 & 1 & 1 \\ z_1 & z_2 & z_3 \end{pmatrix} \right|, \quad (3.39)$$

where the occurring matrices are 2×2 and 3×3 , respectively, i.e., $z_j \in \mathbb{R}^2$ are column vectors.

Proof. Without loss of generality, we may assume that T is non-degenerate. The first equality has already been proved in Lemma 3.9. The second equality follows from simple Linear Algebra

$$\det \begin{pmatrix} 1 & 1 & 1 \\ z_1 & z_2 & z_3 \end{pmatrix} = \det \begin{pmatrix} 1 & 0 & 0 \\ z_1 & z_2 - z_1 & z_3 - z_1 \end{pmatrix} = \det \begin{pmatrix} z_2 - z_1 & z_3 - z_1 \end{pmatrix},$$

which concludes the proof. ■

We remark that the order of the nodes z_1, z_2, z_3 of T matters in the following way: If the nodes are given in counter clockwise order with respect to the boundary ∂T , the determinants in the previous lemma are always non-negative. According to our implementational conventions, we are going to use this observation for the MATLAB implementation given below.

Exercise 25. Let $T = \text{conv}\{z_1, z_2, z_3\}$ be a non-degenerate triangle in \mathbb{R}^2 with center of mass s_T . Let $v_h : T \rightarrow \mathbb{R}$ be an affine function. Prove that

$$\int_T v_h dx = |T| v_h(s_T),$$

which means that the quadrature involved above is exact for affine functions. Conclude that

$$\int_T \zeta_k dx = \frac{|T|}{3},$$

where ζ_k is the hat function corresponding to one of the nodes z_k . □

3.4.3 Building the Galerkin Matrix

Next, we discuss the assembly of the Galerkin matrix $A \in \mathbb{R}_{\text{sym}}^{N \times N}$, where the entries take the form

$$A_{jk} = (\nabla \zeta_k ; \nabla \zeta_j)_{L^2(\Omega)} = \sum_{T \in \mathcal{T}} \int_T \nabla \zeta_k \cdot \nabla \zeta_j \, dx.$$

Obviously, a summand is only nonzero provided that the nodes z_j and z_k corresponding to the hat functions ζ_j and ζ_k are nodes of $T = \text{conv}\{x_T, y_T, z_T\}$. For each element $T \in \mathcal{T}$, we thus have to compute a matrix $A_T \in \mathbb{R}_{\text{sym}}^{3 \times 3}$, which is called **local stiffness matrix**. The pseudo MATLAB code for the computation of A then takes the following form:

```

N = size(coordinates,1);
A = zeros(N,N);
for j = 1:size(elements,1)
    nodes = elements(j,:);
    A(nodes,nodes) = A(nodes,nodes) + A_Tj;
end
    
```

We stress that the gradients $\nabla \zeta_\ell$ are constant on an element $T \in \mathcal{T}$. Therefore, the entries of A_T can be computed analytically without quadrature. This is done in the following lemma:

Lemma 3.13. *Let $T = \text{conv}\{z_1, z_2, z_3\}$ be a non-degenerate triangle in \mathbb{R}^2 . Let $\zeta_1, \zeta_2, \zeta_3$ be the corresponding hat functions on T . We define the matrices $A, M \in \mathbb{R}^{3 \times 3}$ by*

$$A_{jk} = \int_T \nabla \zeta_j \cdot \nabla \zeta_k \, dx \quad \text{and} \quad M = \begin{pmatrix} 1 & 1 & 1 \\ z_1 & z_2 & z_3 \end{pmatrix}. \quad (3.40)$$

Then, M is regular and with $B \in \mathbb{R}^{3 \times 2}$ defined by

$$B := M^{-1} \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad (3.41)$$

it holds that $A = \frac{1}{2} |\det(M)| BB^T$.

Proof. According to Lemma 3.12, it holds that $|T| = |\det(M)|/2$. Therefore, non-degeneracy of T implies that M is regular. For given $x \in \mathbb{R}^2$, let the **barycentric coordinates** $\lambda(x) \in \mathbb{R}^3$ be the unique solution of

$$M\lambda(x) = \begin{pmatrix} 1 \\ x \end{pmatrix}.$$

Note that this is equivalent to $\sum_{j=1}^3 \lambda_j(x) = 1$ and $x = \sum_{j=1}^3 \lambda_j(x) z_j$. In particular, this and the regularity of M imply that $\lambda_j(z_k) = \delta_{jk}$. Moreover, writing $\lambda(x)$ in the form

$$\lambda(x) = \begin{pmatrix} \lambda_1(x) \\ \lambda_2(x) \\ \lambda_3(x) \end{pmatrix} = \begin{pmatrix} c_{11} & c_{12} & c_{13} \\ c_{21} & c_{22} & c_{23} \\ c_{31} & c_{32} & c_{33} \end{pmatrix} \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix} \quad \text{with the matrix } C := M^{-1},$$

we see that each λ_j is affine. This, in fact, implies that $\lambda_j = \zeta_j$. By definition of B , we see that

$$B = C \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} c_{12} & c_{13} \\ c_{22} & c_{23} \\ c_{32} & c_{33} \end{pmatrix} = \begin{pmatrix} \nabla \lambda_1(x) \\ \nabla \lambda_2(x) \\ \nabla \lambda_3(x) \end{pmatrix},$$

i.e., the j -th row of B coincides with $\nabla \lambda_j$. This finally leads to

$$\nabla \lambda_j \cdot \nabla \lambda_k = \sum_{\ell=1}^2 B_{j\ell} B_{k\ell} = (BB^T)_{jk}.$$

Altogether, $A_{jk} = |T| (BB^T)_{jk} = \frac{1}{2} |\det(M)| (BB^T)_{jk}$. This concludes the proof. ■

A MATLAB realization of the previous lemma reads as follows:

```

1 function AT = stima(nodes)
2
3 % STIMA(NODES) returns the local P1 stiffness matrix for a
4 % non-degenerate triangle T which is given by the 3 x 2 array
5 % NODES, i.e. each row of NODES provides the coordinates of
6 % one vertex of T
7
8 M = [1 1 1 ; nodes'];
9 B = M \ [0 0 ; 1 0 ; 0 1];
10 AT = det(M)*B*B'/2;
```

With the MATLAB function `stima`, the assembly of A reads as follows:

```

N = size(coordinates,1);
A = zeros(N,N);
for j = 1:size(elements,1)
    nodes = elements(j,:);
    A(nodes,nodes) = A(nodes,nodes) + stima(coordinates(nodes,:));
end
```

Exercise 26. Let $T = \text{conv}\{z_1, z_2, z_3\}$ be a non-degenerate triangle in \mathbb{R}^2 . Let $\zeta_1, \zeta_2, \zeta_3$ denote the hat functions corresponding to the nodes z_j . Define the local mass matrix $A \in \mathbb{R}_{\text{sym}}^{3 \times 3}$ by

$$A_{jk} = \int_T \zeta_j \zeta_k dx.$$

Give an explicit formula for A . **Hint:** What are the hat functions on the reference element? Provide a formula for $T = T_{\text{ref}}$ first. Then, use the transformation formula to obtain the general result. □

Exercise 27. Lemma 3.12 and Lemma 3.13 as well as Exercise 26 can be generalized to simplices $T = \text{conv}\{z_1, \dots, z_{d+1}\}$ in \mathbb{R}^d which are non-degenerate, i.e., $|T| > 0$. State and prove the extended results. \square

3.4.4 Computation of Galerkin Solution

We stress that we do not have built the Galerkin data according to Theorem 1.4. Instead we have built the matrix $A \in \mathbb{R}_{\text{sym}}^{N \times N}$ and the right-hand side $b \in \mathbb{R}^N$ with respect to the nodal basis $\mathcal{B} = \{\zeta_1, \dots, \zeta_N\}$ of $\mathcal{S}^1(\mathcal{T})$. In case of the **Dirichlet problem** (3.14) or the **mixed boundary value problem** (3.15), the discrete space is $\mathcal{S}_D^1(\mathcal{T})$, and a corresponding basis is $\mathcal{B}_D := \{\zeta_\ell \mid z_\ell \in \mathcal{K} \setminus \bar{\Gamma}_D\}$. Therefore, the Galerkin system according to Theorem 1.4 is simply a subsystem of the linear system built: We have to solve the system only for the so-called **free nodes** $\mathcal{K} \setminus \bar{\Gamma}_D$. A realization of this idea in pseudo MATLAB code reads like this:

```
N = size(coordinates,1);
x = zeros(N,1);
freenodes = setdiff(1:N, unique(dirichlet));
x(freenodes) = A(freenodes,freenodes)\b(freenodes);
```

Here, `unique` provides the list $\{\ell \mid z_\ell \in \bar{\Gamma}_D\}$ of all indices ℓ such that z_ℓ belongs to the Dirichlet boundary Γ_D . Moreover, `setdiff` denotes the difference $\{1, \dots, N\} \setminus \{\ell \mid z_\ell \in \mathcal{K} \cap \bar{\Gamma}_D\}$ and thus provides the indices of the free nodes. — Note that this subsystem is, in fact, the linear system stated in Theorem 1.4. — Altogether,

$$u_h := \sum_{j=1}^N x_j \zeta_j$$

thus provides the Galerkin approximation $u_h \in \mathcal{S}_D^1(\mathcal{T})$ of $u \in H_D^1(\Omega)$, where the subscript D is replaced by the subscript 0 in case of the Dirichlet problem $\Gamma_D = \Gamma$.

For the **Neumann problem** (3.16), we have to enforce the side constraint $\int_\Omega u_h dx = 0 = \int_\Omega v_h dx$ for the Galerkin solution u_h as well as for the test functions v_h . This can be done in two ways, which are stated in the following exercises:

Exercise 28. We adopt the foregoing notations and assume that $A \in \mathbb{R}_{\text{sym}}^{N \times N}$ and $b \in \mathbb{R}^N$ have been computed with respect to the nodal basis \mathcal{B} of $\mathcal{S}^1(\mathcal{T})$. We define an additional vector $c \in \mathbb{R}^N$ by $c_k = \int_\Omega \zeta_k dx$. Prove that the matrix of the linear system

$$\begin{pmatrix} A & c \\ c^T & 0 \end{pmatrix} \begin{pmatrix} x \\ \lambda \end{pmatrix} = \begin{pmatrix} b \\ 0 \end{pmatrix}$$

is regular and thus leads to a unique solution $(x, \lambda) \in \mathbb{R}^N \times \mathbb{R}$. Prove that $u_h = \sum_{j=1}^N x_j \zeta_j$ is the unique Galerkin solution $u_h \in \mathcal{S}_*^1(\mathcal{T})$ of the Neumann problem (3.16). Write a pseudo MATLAB code that realizes this Lagrange multiplier ansatz for the computation of the coefficient vector $x \in \mathbb{R}^N$. \square

One drawback of the preceding ansatz is that the system matrix

$$M := \begin{pmatrix} A & c \\ c^T & 0 \end{pmatrix}$$

is symmetric but *not* positive definite. Since SPD matrices are good in practice, e.g., for the direct solution by use of the Cholesky factorization, in engineering the following ansatz is used:

Exercise 29. We again adopt the foregoing notations and assume that $A \in \mathbb{R}_{\text{sym}}^{N \times N}$ and $b \in \mathbb{R}^N$ have been computed with respect to the nodal basis \mathcal{B} of $\mathcal{S}^1(\mathcal{T})$. We define the restricted system $\tilde{A} \in \mathbb{R}_{\text{sym}}^{(N-1) \times (N-1)}$ and $\tilde{b} \in \mathbb{R}^{N-1}$ by

$$\tilde{A}_{jk} = A_{jk}, \quad \tilde{b}_k = b_k \quad \text{for } j, k = 1, \dots, N-1.$$

Prove that \tilde{A} is symmetric and positive definite, whence regular. With the unique solution $\tilde{x} \in \mathbb{R}^{N-1}$ of $\tilde{A}\tilde{x} = \tilde{b}$, we define

$$u_h := \tilde{u}_h - \frac{1}{|\Omega|} \int_{\Omega} \tilde{u}_h dx, \quad \text{where } \tilde{u}_h := \sum_{j=1}^{N-1} \tilde{x}_j \zeta_j.$$

Prove that u_h is the unique Galerkin solution $u_h \in \mathcal{S}_*^1(\mathcal{T})$ of the Neumann problem (3.16). Write a pseudo MATLAB code which realizes this algorithm and which provides the coefficient vector $x \in \mathbb{R}^N$ of u_h with respect to the nodal basis \mathcal{B} . \square

3.4.5 A First Matlab Implementation

The following function `solveLaplace0` provides a MATLAB implementation of the P1-FEM for the Dirichlet problem and the mixed boundary value problem. We suppose that `f` and `phi` are function handles in MATLAB 7. Since most of the entries of the Galerkin matrix A are zero, we use the data format `sparse` instead of a dense matrix obtained by `zeros`. This reduces the storage requirements from N^2 to $\mathcal{O}(M)$ since the matrix has at most $9M$ non-zero entries, where $N = \#\mathcal{K}$ denotes the number of nodes and $M = \#\mathcal{T}$ is the number of elements. For usual triangulations, it holds that $N = \mathcal{O}(M)$ so that the storage requirements are, in fact, decreased down to linear. Moreover, we do not use the function `stima` but include its code for the assembly of A . This avoids an unnecessary computational overhead which arises from the internal realization of function calls in MATLAB.

```

1 function x = solveLaplace0(coordinates,elements,f,dirichlet,neumann,phi)
2
3 % SOLVELAPLACE0(ELEMENTS,COORDINATES,DIRICHLET,NEUMANN,F,PHI)
4 % solves the mixed boundary value problem for the Laplacian
5 % with Neumann data PHI and zero Dirichlet data. In particular,
6 % the code solves the Dirichlet problem, if NEUMANN is an
7 % empty matrix. The source term F as well as (in case of) the
8 % Neumann data PHI have to be function handles.
9
10 % (c) 2007,2008 by Dirk Praetorius, last modified 07.01.2008

```



```

11 % dirk.praetorius@tuwien.ac.at - http://www.asc.tuwien.ac.at/~dirk
12
13 N = size(coordinates,1);
14 x = zeros(N,1);
15 b = zeros(N,1);
16 A = sparse(N,N);
17
18 %*** Assemble the right-hand side
19 for j = 1:size(elements,1)
20     nodes = elements(j,:);
21     s = [1 1 1]*coordinates(nodes,:)/3;
22     b(nodes) = b(nodes) + det([1 1 1 ; coordinates(nodes,:)'])*f(s)/6;
23 end
24 for j = 1:size(neumann,1)
25     nodes = neumann(j,:);
26     m = [1 1]*coordinates(nodes,:)/2;
27     b(nodes) = b(nodes) + norm([1 -1]*coordinates(nodes,:))*phi(m)/2;
28 end
29
30 %*** Assemble the stiffness matrix
31 for j = 1:size(elements,1)
32     nodes = elements(j,:);
33     M = [1 1 1 ; coordinates(nodes,:)'];
34     B = M \ [0 0 ; 1 0 ; 0 1];
35     A(nodes,nodes) = A(nodes,nodes) + det(M)*B*B'/2;
36 end
37
38 %*** Compute P1-FEM approximation
39 freenodes = setdiff(1:N, unique(dirichlet));
40 x(freenodes) = A(freenodes,freenodes)\b(freenodes);

```

Exercise 30. Prove that the stiffness matrix A has at most $9M$ entries, where $M = \#\mathcal{T}$ denotes the number of elements. □

Exercise 31. Write a MATLAB function `solveLaplace` which realizes the P1-FEM for the Dirichlet problem and for the mixed boundary value problem as well as for the Neumann problem. For the extension of the code of the function `solveLaplace0`, use Exercise 28 or Exercise 29. □

In MATLAB, the visualization of the P1-FEM solution $u_h = \sum_{z \in \mathcal{K}} x_z \zeta_z$ is rather simple. Let $N = \#\mathcal{K}$ be the number of nodes of the triangulation \mathcal{T} . The given MATLAB code computes the column vector $x \in \mathbb{R}^N$ of the basis representation of u_h with respect to the nodal basis \mathcal{B} of $\mathcal{S}^1(\mathcal{T})$. Then, the following MATLAB function generates a surface plot of u_h :

```
trisurf(elements,coordinates(:,1),coordinates(:,2),x,'facecolor','interp')
```

Here, `elements` and `coordinates` are matrices as introduced above.

3.4.6 A Second Matlab Implementation with (Almost) Linear Complexity

On the first glance, the MATLAB implementation of the previous section has linear complexity with respect to the number $M = \#\mathcal{T}$ of elements. However, if we plot the building time of the stiffness matrix A over M , we observe a quadratic dependence instead, cf. Figure 3.7 below. This is due to the internal storage of sparse matrices in MATLAB. A matrix of type `sparse` is stored in the compressed column storage format. To explain this storage format, we consider the sparse matrix

$$A = \begin{pmatrix} 10 & 0 & 0 & -2 & 0 \\ 3 & 9 & 0 & 0 & 3 \\ 0 & 7 & 8 & 0 & 0 \\ 3 & 0 & 8 & 5 & 0 \\ 0 & 8 & 0 & 9 & 13 \\ 0 & 4 & 0 & 2 & -1 \end{pmatrix} \quad (3.42)$$

The probably most simple format to store a sparse matrix, is the so-called **coordinate format** which is also used as the output format of a sparse matrix in MATLAB.

- For the coordinate format, one stores the dimensions M, N of the matrix $A \in \mathbb{R}^{M \times N}$ as well as the number n of nonzero elements together with two index vectors $I, J \in \mathbb{N}^n$ and a vector of entries $a \in \mathbb{R}^n$. For each index $1 \leq \ell \leq n$, it holds that

$$a_\ell = A_{ij} \quad \text{with} \quad i = I_\ell \quad \text{and} \quad j = J_\ell.$$

The ordering of the vectors I, J, a does not matter in the following sense: For the matrix $A \in \mathbb{R}^{6 \times 5}$ from (3.42) holds $n = 16$. It can be stored, for instance, columnwise

$$\begin{aligned} a &= (\quad 10 \quad 3 \quad 3 \quad | \quad 9 \quad 7 \quad 8 \quad 4 \quad | \quad 8 \quad 8 \quad | \quad -2 \quad 5 \quad 9 \quad 2 \quad | \quad 3 \quad 13 \quad -1 \quad), \\ I &= (\quad 1 \quad 2 \quad 4 \quad | \quad 2 \quad 3 \quad 5 \quad 6 \quad | \quad 3 \quad 4 \quad | \quad 1 \quad 4 \quad 5 \quad 6 \quad | \quad 2 \quad 5 \quad 6 \quad), \\ J &= (\quad 1 \quad 1 \quad 1 \quad | \quad 2 \quad 2 \quad 2 \quad 2 \quad | \quad 3 \quad 3 \quad | \quad 4 \quad 4 \quad 4 \quad 4 \quad | \quad 5 \quad 5 \quad 5 \quad), \end{aligned}$$

or rowwise

$$\begin{aligned} a &= (\quad 10 \quad -2 \quad | \quad 3 \quad 9 \quad 3 \quad | \quad 7 \quad 8 \quad | \quad 3 \quad 8 \quad 5 \quad | \quad 8 \quad 9 \quad 13 \quad | \quad 4 \quad 2 \quad -1 \quad), \\ I &= (\quad 1 \quad 1 \quad | \quad 2 \quad 2 \quad 2 \quad | \quad 3 \quad 3 \quad | \quad 4 \quad 4 \quad 4 \quad | \quad 5 \quad 5 \quad 5 \quad | \quad 6 \quad 6 \quad 6 \quad), \\ J &= (\quad 1 \quad 4 \quad | \quad 1 \quad 2 \quad 5 \quad | \quad 2 \quad 3 \quad | \quad 1 \quad 3 \quad 4 \quad | \quad 2 \quad 4 \quad 5 \quad | \quad 2 \quad 4 \quad 5 \quad). \end{aligned}$$

Both combinations describe the same sparse matrix.

An alternative is the so-called **compressed column storage format** (CCS format or Harwell-Boeing format) which is used for the internal storage of sparse matrices in MATLAB.

- For the CCS format, one stores the dimensions M, N of the matrix $A \in \mathbb{R}^{M \times N}$ as well as the number n of nonzero elements together with an index vector $I \in \mathbb{N}^n$ and a vector of entries $a \in \mathbb{R}^n$ which are both ordered columnwise. Moreover, one stores a vector $J \in \mathbb{N}^{N+1}$ such that

$$a_\ell = A_{ij} \quad \text{with} \quad i = I_\ell \quad \text{for all indices} \quad J(j) \leq \ell < J(j+1).$$

By definition $J_{N+1} := n + 1$ so that one does not have to store the number n of non-zero entries explicitly. In the case of the matrix $A \in \mathbb{R}^{6 \times 5}$ from (3.42), there hold

$$\begin{aligned} a &= (10 \ 3 \ 3 \mid 9 \ 7 \ 8 \ 4 \mid 8 \ 8 \mid -2 \ 5 \ 9 \ 2 \mid 3 \ 13 \ -1), \\ I &= (1 \ 2 \ 4 \mid 2 \ 3 \ 5 \ 6 \mid 3 \ 4 \mid 1 \ 4 \ 5 \ 6 \mid 2 \ 5 \ 6), \end{aligned}$$

as well as

$$J = (1 \ 4 \ 8 \ 10 \ 14 \mid 17).$$

We stress again that MATLAB stores a sparse matrix internally by the CCS format, whereas the data is converted on-the-fly into the coordinate format for output only.

Obviously, the update of a sparse matrix which is stored in CCS format, needs some sorting. This is usually of complexity $\mathcal{O}(k \log k)$ if k is the current number of non-zero entries — per update! In the MATLAB implementation of the preceding Section 3.4.5, the sparse matrix A is updated M times, where $M = \#\mathcal{T}$ is the number of elements. Therefore, the assembly of the stiffness matrix A leads to a complexity $\mathcal{O}(M^2)$ (if we neglect the logarithmic factor) which is even observed in praxis. However, MATLAB provides a remedy in terms of the function `sparse`. Let $a \in \mathbb{R}^n$ and $I, J \in \mathbb{N}^n$ be the vectors for the coordinate format of a sparse matrix $A \in \mathbb{R}^{M \times N}$. Then, A can be declared and initialized by use of

$$A = \text{sparse}(I, J, a, M, N)$$

If an index pair $(i, j) = (I_\ell, J_\ell)$ appears twice (or even more), the corresponding entries a_ℓ are added. In particular, the internal realization only needs one sorting of the entries which is only of complexity $\mathcal{O}(n \log n)$ with $n \in \mathbb{N}$ the number of non-zero elements. Note that $n = \mathcal{O}(M + N)$ for a usual sparse matrix $A \in \mathbb{R}^{M \times N}$. The following code is a realization of this possibility, where we use that each element $T \in \mathcal{T}$ leads to 9 entries in the corresponding vectors a, I, J of the coordinate format. Note that the latter observation also allows to avoid the logarithmic term for sorting!

```

1 function x = solveLaplace0(coordinates,elements,f,dirichlet,neumann,phi)
2
3 % SOLVELAPLACE0(ELEMENTS,COORDINATES,DIRICHLET,NEUMANN,F,PHI)
4 % solves the mixed boundary value problem for the Laplacian
5 % with Neumann data PHI and zero Dirichlet data. In particular,
6 % the code solves the Dirichlet problem, if NEUMANN is an
7 % empty matrix. The source term F as well as (in case of) the
8 % Neumann data PHI have to be function handles.
9
10 % (c) 2007,2008 by Dirk Praetorius, last modified 07.01.2008
11 % dirk.praetorius@tuwien.ac.at - http://www.asc.tuwien.ac.at/~dirk
12
13 M = size(elements,1);
14 N = size(coordinates,1);
15
16 x = zeros(N,1);
```

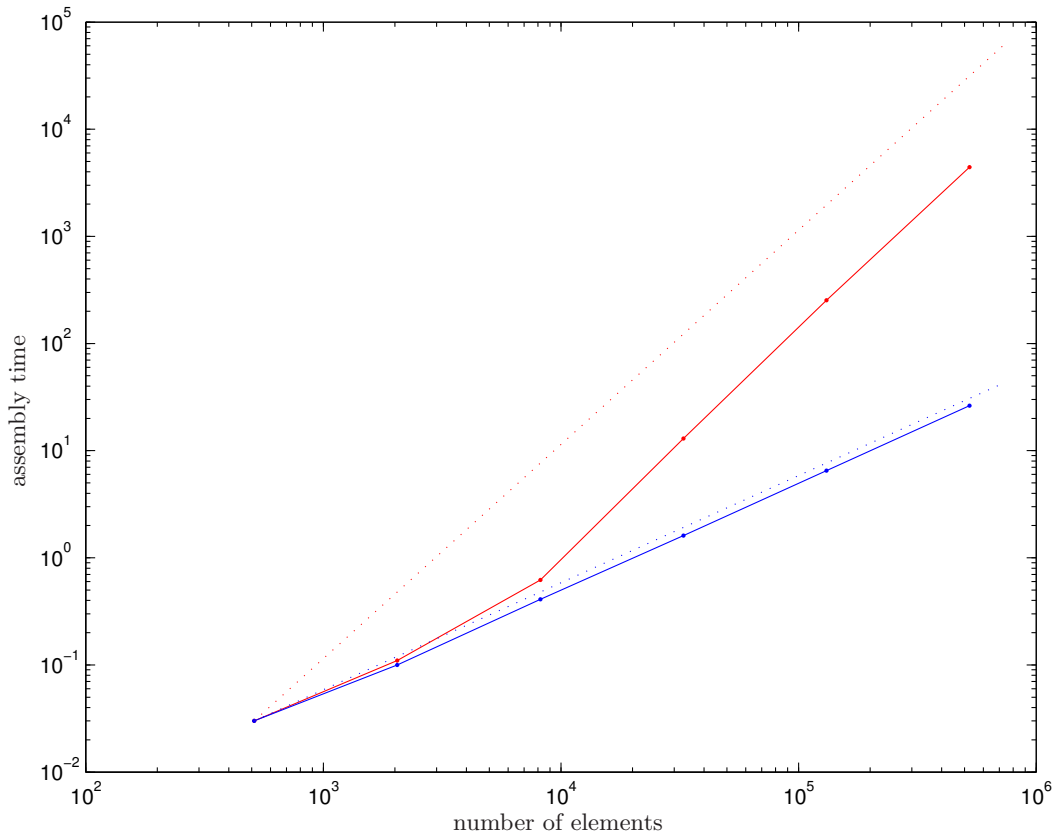


FIGURE 3.7. Assembly time in seconds for the stiffness matrix A by use of the code of Section 3.4.5 (red) and Section 3.4.6, respectively. In the double-logarithmic plot, a dependence $T = M^\alpha$ between assembly time T and number M of elements is visualized by a straight line with slope α . The dotted lines indicate quadratic growth $\alpha = 2$ (red) versus linear growth $\alpha = 1$ (blue).

```

17 b = zeros(N,1);
18 I = zeros(9*M,1);
19 J = zeros(9*M,1);
20 A = zeros(9*M,1);
21
22 %*** Assemble the right-hand side
23 for j = 1:M
24     nodes = elements(j,:);
25     s = [1 1 1]*coordinates(nodes,+)/3;
26     b(nodes) = b(nodes) + det([1 1 1 ; coordinates(nodes,:)'])*f(s)/6;
27 end
28 for j = 1:size(neumann,1)
29     nodes = neumann(j,:);
30     m = [1 1]*coordinates(nodes,+)/2;
31     b(nodes) = b(nodes) + norm([1 -1]*coordinates(nodes,:))*phi(m)/2;
32 end

```

```

33
34 %*** Assemble the stiffness matrix in (almost) linear complexity
35 for j = 1:M
36     nodes = elements(j,:);
37     M = [1 1 1 ; coordinates(nodes,:)'];
38     B = M \ [0 0 ; 1 0 ; 0 1];
39     idx = 9*(j-1)+1:9*j;
40     tmp = [1;1;1]*nodes;
41     I(idx) = reshape(tmp,9,1);
42     J(idx) = reshape(tmp',9,1);
43     A(idx) = det(M)/2*reshape(B*B',9,1);
44 end
45 A = sparse(I,J,A,N,N);
46
47 %*** Compute P1-FEM approximation
48 freenodes = setdiff(1:N, unique(dirichlet));
49 x(freenodes) = A(freenodes,freenodes)\b(freenodes);

```

In Figure 3.7, we compare the building time for the stiffness matrix in the MATLAB code of Section 3.4.5 and Section 3.4.6, respectively. As expected, we observe that the first implementation leads to a quadratic dependence of assembly time and number of elements. In contrast to that, the new implementation of this section provides the optimal complexity in the sense that we observe linear dependence.

Exercise 32. Write MATLAB functions `coordinate2ccs` and `ccs2coordinate` that convert a sparse matrix $A \in \mathbb{R}^{M \times N}$ given by the coordinate format into the CCS format and vice versa. Keep the computational complexity as low as possible. \square

Exercise 33. Write MATLAB functions for the matrix-vector multiplication $b := Ax$, where $A \in \mathbb{R}^{M \times N}$ is a sparse matrix given in coordinate format and CCS format, respectively. Avoid converting to a suitable MATLAB-format but work with the corresponding format vectors directly. You might use the conversion and the usual matrix-vector multiplication for the verification of your code. \square

Exercise 34. Let $f : \mathbb{R}_{>0} \rightarrow \mathbb{R}_{>0}$ be a function of N . For given $N_1 < \dots < N_\ell$, we plot $f(N_j)$ over N in a double logarithmic plot, cf. Figure 3.7. Prove that a dependence $f(N) = CN^\alpha$ for some $\alpha \in \mathbb{R}$ and $C > 0$ leads to a straight line with slope α . \square

3.5 FEM with Data Approximation (31.10.2017)

Now that we have realized that the P1-FEM is of order $\mathcal{O}(h)$, we need to show that the quadrature rules used for our MATLAB implementation are sufficiently accurate. Recall that we are approxi-

mating the right-hand side of the exact P1-FEM

$$F(v) := \int_{\Omega} f v \, dx + \int_{\Gamma_N} \phi v \, ds \quad \text{for } v \in H^1(\Omega) \quad (3.43)$$

by

$$F_h(v_h) := \sum_{T \in \mathcal{T}} |T| f(s_T) v_h(s_T) + \sum_{E \in \mathcal{E}_N} h_E \phi(m_E) v_h(m_E) \quad \text{for } v_h \in \mathcal{S}^1(\mathcal{T}), \quad (3.44)$$

where s_T denotes the center of mass of an element $T \in \mathcal{T}$ and where m_E denotes the midpoint of a Neumann edge $E \in \mathcal{E}_N$. Therefore, our MATLAB code realizes a perturbed P1-FEM and we need to study the convergence of this perturbed scheme.

3.5.1 First Strang Lemma

In this section, we go back to the abstract formulation of Galerkin schemes: Let H be a real Hilbert space with norm $\|\cdot\|_H$. Let $\langle\langle \cdot ; \cdot \rangle\rangle$ be a bilinear form which is assumed to be elliptic and continuous, i.e., it holds that

$$\alpha \|v\|_H^2 \leq \langle\langle v ; v \rangle\rangle \quad \text{as well as} \quad \langle\langle v ; w \rangle\rangle \leq \beta \|v\|_H \|w\|_H \quad \text{for all } v, w \in H. \quad (3.45)$$

Let $F \in H^*$ be a given right-hand side. Then, the Lax-Milgram lemma applies and yields the existence and uniqueness of the solution $u \in H$ of

$$\langle\langle u ; \cdot \rangle\rangle = F \in H^*. \quad (3.46)$$

For a discretization parameter $h > 0$, let X_h be a finite dimensional subspace of H . It is an important property of a Galerkin scheme that it is stable with respect to certain perturbations of the scalar product $\langle\langle \cdot ; \cdot \rangle\rangle$ or the right-hand side F . — For the interpretation, recall that usually the right-hand side $F \in H^*$ as well as the scalar product $\langle\langle \cdot ; \cdot \rangle\rangle$ involve integrals, which are computed numerically by quadrature rules. For a fixed discrete space X_h , this leads to perturbations $F_h \in X_h^*$ and $\langle\langle \cdot ; \cdot \rangle\rangle_h$ of F and $\langle\langle \cdot ; \cdot \rangle\rangle$, respectively. In particular, this gives rise to additional **consistency errors**

$$\|F - F_h\|_{X_h^*} \quad \text{and} \quad \sup_{v_h \in X_h \setminus \{0\}} \frac{\|\langle\langle v_h ; \cdot \rangle\rangle - \langle\langle v_h ; \cdot \rangle\rangle_h\|_{X_h^*}}{\|v_h\|_H}.$$

In practice, the **best approximation error** (or **discretization error**) behaves like

$$\min_{v_h \in X_h} \|u - v_h\|_H = \mathcal{O}(h^\alpha) \quad \text{for } h \rightarrow 0,$$

where the **convergence order** $\alpha > 0$ usually depends on the regularity of the exact solution u . Then, the Céa lemma proves that

$$\|u - \mathbb{G}_h u\|_H = \mathcal{O}(h^\alpha).$$

The following result due to Strang shows that the consistency errors should be at least of the same order, i.e., one needs a sufficiently large order for the quadrature rules. Then, the perturbed Galerkin scheme

$$\langle\langle u_h ; v_h \rangle\rangle_h = F_h(v_h) \quad \text{for all } v_h \in X_h \quad (3.47)$$

still allows for a unique solution $u_h \in X_h$. Moreover the **approximation error** still satisfies

$$\|u - u_h\|_H = \mathcal{O}(h^\alpha).$$

However, the consequence of Strang's lemma even works the other way around: You should avoid to compute integrals exactly (or with high accuracy quadrature rules) since this is usually computationally expensive and since this expense does not pay in the sense of an increased order of convergence. Finally, we note that analytic computation of integrals via antiderivatives, i.e., $\int_a^b f dx = F(b) - F(a)$ for the simple 1D case, necessarily leads to cancellation for small mesh-sizes. These are, however, avoided for numerical integration via Gaussian quadrature rules, since the Gaussian quadrature weights are all positive. In explicit terms, this implies that approximate computation will be numerically more accurate than analytic computation, if the quadrature rules are deliberately chosen.

v9:
03.11.2017

Proposition 3.14 (First Strang Lemma). *Assume that $\langle\langle \cdot ; \cdot \rangle\rangle_h$ is a bilinear form on X_h and that $F_h : X_h \rightarrow \mathbb{R}$ is linear. Then, there holds the following:*

(i) *Assume convergence of $\langle\langle \cdot ; \cdot \rangle\rangle_h$ to $\langle\langle \cdot ; \cdot \rangle\rangle$, i.e.,*

$$\lim_{h \rightarrow 0} E_h = 0 \quad \text{with} \quad E_h := \sup_{v_h, w_h \in X_h \setminus \{0\}} \frac{|\langle\langle v_h ; w_h \rangle\rangle - \langle\langle v_h ; w_h \rangle\rangle_h|}{\|v_h\|_H \|w_h\|_H}. \quad (3.48)$$

Then, the bilinear forms are uniformly elliptic for small h , i.e.,

$$\exists \alpha_0 > 0 \exists h_0 > 0 \forall h \in (0, h_0) \forall v_h \in X_h \quad \alpha_0 \|v_h\|_H^2 \leq \langle\langle v_h ; v_h \rangle\rangle_h. \quad (3.49)$$

In particular, there exist unique $u_h \in X_h$ with $\langle\langle u_h ; \cdot \rangle\rangle_h = F_h \in X_h^$ for sufficiently small $h > 0$.*

(ii) *Provided (3.49), there holds the Céa type estimate*

$$\begin{aligned} C^{-1} \|u - u_h\|_H &\leq \inf_{v_h \in X_h} (\|u - v_h\|_H + \|\langle\langle v_h ; \cdot \rangle\rangle - \langle\langle v_h ; \cdot \rangle\rangle_h\|_{X_h^*}) + \|F - F_h\|_{X_h^*} \\ &\leq (1 + E_h) \min_{v_h \in X_h} \|u - v_h\|_H + E_h \|u\|_H + \|F - F_h\|_{X_h^*} \end{aligned} \quad (3.50)$$

with u being the exact solution of (3.46). The constant $C > 0$ depends only on $\langle\langle \cdot ; \cdot \rangle\rangle$.

Proof. Let $0 < \varepsilon < \alpha$ and $h_0 > 0$ such that

$$\forall h \in (0, h_0) \quad \sup_{v_h \in X_h \setminus \{0\}} \frac{|\langle\langle v_h ; v_h \rangle\rangle - \langle\langle v_h ; v_h \rangle\rangle_h|}{\|v_h\|_H^2} \leq \varepsilon.$$

Then, $\alpha \|v_h\|_H^2 \leq \langle\langle v_h ; v_h \rangle\rangle \leq \langle\langle v_h ; v_h \rangle\rangle_h + |\langle\langle v_h ; v_h \rangle\rangle - \langle\langle v_h ; v_h \rangle\rangle_h| \leq \langle\langle v_h ; v_h \rangle\rangle_h + \varepsilon \|v_h\|_H^2$, whence

$$(\alpha - \varepsilon) \|v_h\|_H^2 \leq \langle\langle v_h ; v_h \rangle\rangle_h,$$

i.e., $\langle\langle \cdot ; \cdot \rangle\rangle_h$ is an elliptic bilinear form on X_h for $h < h_0$. This concludes the proof of (i) with $\alpha_0 := \alpha - \varepsilon > 0$. To prove (ii), let $v_h \in X_h$ be arbitrary. Then,

$$\alpha_0 \|v_h - u_h\|_H^2 \leq \langle\langle v_h - u_h ; v_h - u_h \rangle\rangle_h = \langle\langle v_h ; v_h - u_h \rangle\rangle_h - F_h(v_h - u_h).$$

Together with

$$\langle\langle u - v_h ; v_h - u_h \rangle\rangle = F(v_h - u_h) - \langle\langle v_h ; v_h - u_h \rangle\rangle,$$

we obtain that

$$\begin{aligned} \alpha_0 \|v_h - u_h\|_H^2 &\leq [F(v_h - u_h) - F_h(v_h - u_h)] + [\langle\langle v_h ; v_h - u_h \rangle\rangle_h - \langle\langle v_h ; v_h - u_h \rangle\rangle] \\ &\quad - \langle\langle u - v_h ; v_h - u_h \rangle\rangle \\ &\leq \|v_h - u_h\|_H [\|F - F_h\|_{X_h^*} + \|\langle\langle v_h ; \cdot \rangle\rangle_h - \langle\langle v_h ; \cdot \rangle\rangle\|_{X_h^*} + \beta \|u - v_h\|_H]. \end{aligned}$$

Finally, the combination with a triangle inequality yields that

$$\begin{aligned} \|u - u_h\|_H &\leq \|u - v_h\|_H + \|v_h - u_h\|_H \\ &\leq C [\|F_h - F\|_{X_h^*} + \|\langle\langle v_h ; \cdot \rangle\rangle - \langle\langle v_h ; \cdot \rangle\rangle_h\|_{X_h^*} + \|u - v_h\|_H] \end{aligned}$$

for any $v_h \in X_h$ with $C = 1 + \beta/\alpha_0$. This proves the first estimate in (3.50). To see the second estimate, note that

$$\|\langle\langle v_h ; \cdot \rangle\rangle - \langle\langle v_h ; \cdot \rangle\rangle_h\|_{X_h^*} \leq E_h \|v_h\|_H \leq E_h \|u\|_H + E_h \|u - v_h\|_H.$$

This concludes the proof. ■

Under the assumptions of the Strang lemma, one can even show convergence of the perturbed Galerkin scheme.

Exercise 35. Assume that $\langle\langle \cdot ; \cdot \rangle\rangle_h$ is a symmetric bilinear form on X_h^* and that $F_h \in X_h^*$. We assume convergence of the data in the sense that

$$\lim_{h \rightarrow 0} E_h = 0 = \lim_{h \rightarrow 0} \|F - F_h\|_{X_h^*} \quad \text{with} \quad E_h := \sup_{v_h, w_h \in X_h \setminus \{0\}} \frac{|\langle\langle v_h ; w_h \rangle\rangle - \langle\langle v_h ; w_h \rangle\rangle_h|}{\|v_h\|_H \|w_h\|_H}. \quad (3.51)$$

For sufficiently small $h > 0$, let $u_h \in X_h$ be the unique solutions of the perturbed Galerkin scheme (3.47). Under the approximation assumption

$$\lim_{h \rightarrow 0} \min_{v_h \in X_h} \|v - v_h\|_H = 0 \quad \text{for all } v \in D \quad (3.52)$$

for some dense subspace D of H , there holds convergence

$$\lim_{h \rightarrow 0} \|u - u_h\|_H = 0$$

with u being the exact solution of (3.46). □

3.5.2 Approximation of Volume Forces

For our MATLAB implementation of P1-FEM, we compute the bilinear form $\langle\langle v_h ; w_h \rangle\rangle$ analytically and perturb only the right-hand side. Let F and F_h be given by (3.43)–(3.44), respectively. According to the first Strang lemma 3.14, we only need to show that

$$\|F - F_h\|_{S^1(\mathcal{T})^*} = \mathcal{O}(h)$$

to guarantee that the perturbed P1-FEM is also of order $\mathcal{O}(h)$. We consider the two contributions of the right-hand side separately.

Proposition 3.15. *Let $f \in H^2(\mathcal{T})$ and $F(v) := \int_{\Omega} f v \, dx$ for $v \in H^1(\Omega)$. Let $F_h(v_h) := \sum_{T \in \mathcal{T}} |T| f(s_T) v_h(s_T)$ for $v_h \in \mathcal{S}^1(\mathcal{T})$, where $s_T \in \mathbb{R}^2$ denotes the center of mass of an element $T \in \mathcal{T}$. Then, it holds that*

$$\|F - F_h\|_{\mathcal{S}^1(\mathcal{T})^*} \leq C \|h^2 \nabla f\|_{H^1(\mathcal{T})}, \quad (3.53)$$

where the constant $C > 0$ depends only on T_{ref} , but not on Ω , \mathcal{T} , or f .

Proof. The proof is done elementwise. For $T \in \mathcal{T}$ and $w \in H^1(T)$, we define the integral mean $w_T := |T|^{-1} \int_T w \, dx$. According to the Poincaré inequality, it holds that $\|w - w_T\|_{L^2(T)} \leq C_P h_T \|\nabla w\|_{L^2(T)}$, where the constant $C_P > 0$ is independent of T and w . Moreover, $w \mapsto w_T$ is the L^2 -orthogonal projection onto $\mathcal{P}^0(T)$.

1. step. It holds that

$$\left| \int_T f v_h \, dx - |T| f(s_T) v_h(s_T) \right| \leq C_P^2 h_T^2 \|\nabla f\|_{L^2(T)} \|\nabla v_h\|_{L^2(T)} + \|f_T - f(s_T)\|_{L^2(T)} \|v_h\|_{L^2(T)} :$$

From $\int_T v_h \, dx = |T| v_h(s_T)$, we infer that

$$\begin{aligned} \int_T f v_h \, dx - |T| f(s_T) v_h(s_T) &= (f - f(s_T); v_h)_{L^2(T)} \\ &= (f - f(s_T); v_h - v_{hT})_{L^2(T)} + (f - f(s_T); v_{hT})_{L^2(T)} \\ &= (f - f_T; v_h - v_{hT})_{L^2(T)} + (f_T - f(s_T); v_h)_{L^2(T)} \\ &\leq \|f - f_T\|_{L^2(T)} \|v_h - v_{hT}\|_{L^2(T)} + \|f_T - f(s_T)\|_{L^2(T)} \|v_h\|_{L^2(T)}. \end{aligned}$$

where we have used orthogonality of $(\cdot)_T$ in the last but one step. The Poincaré inequality concludes the proof of step 1.

2. step. It holds that $\|f_T - f(s_T)\|_{L^2(T)} \leq 2C_{\text{ref}} h_T^2 \|D^2 f\|_{L^2(T)}$ with an independent constant $C_{\text{ref}} > 0$, which is obtained from a scaling argument: Let $\Phi : T_{\text{ref}} \rightarrow T$ denote an affine diffeomorphism with linear part $B \in \mathbb{R}^{2 \times 2}$. Note that

$$f_T = \frac{1}{|T|} \int f \, dx = \frac{|\det B|}{|T|} \int_{T_{\text{ref}}} f \circ \Phi \, dx = 2 \int_{T_{\text{ref}}} f \circ \Phi \, dx = \frac{1}{|T_{\text{ref}}|} \int_{T_{\text{ref}}} f \circ \Phi \, dx = (f \circ \Phi)_{T_{\text{ref}}}.$$

Together with $f(s_T) = (f \circ \Phi)(s_{T_{\text{ref}}})$, this yields that

$$\|f_T - f(s_T)\|_{L^2(T)} = |\det B^{-1}|^{-1/2} \|(f \circ \Phi)_{T_{\text{ref}}} - (f \circ \Phi)(s_{T_{\text{ref}}})\|_{L^2(T_{\text{ref}})}.$$

We define $g := f \circ \Phi \in H^2(T_{\text{ref}})$ and consider the operator $A : H^2(T_{\text{ref}}) \rightarrow L^2(T_{\text{ref}})$ defined by $Ag := g_{T_{\text{ref}}} - g(s_{T_{\text{ref}}})$. Then, $\mathcal{P}^1(T_{\text{ref}}) \subseteq \ker A$ and continuity of A follows from the Sobolev inequality

$$\begin{aligned} \|Ag\|_{L^2(T_{\text{ref}})} &\leq \|g_{T_{\text{ref}}}\|_{L^2(T_{\text{ref}})} + |T_{\text{ref}}|^{1/2} |g(s_{T_{\text{ref}}})| \leq \|g\|_{L^2(T_{\text{ref}})} + |T_{\text{ref}}|^{1/2} \|g\|_{\infty, T_{\text{ref}}} \\ &\leq (1 + C_{\text{Sobolev}} |T_{\text{ref}}|^{1/2}) \|g\|_{H^2(T_{\text{ref}})} \end{aligned}$$

Therefore, the Bramble-Hilbert lemma provides a constant $C_{\text{ref}} > 0$ with $\|Ag\|_{L^2(T_{\text{ref}})} \leq C_{\text{ref}} \|D^2 g\|_{L^2(T_{\text{ref}})}$. We conclude the scaling argument by

$$C_{\text{ref}}^{-1} \|(f \circ \Phi)_{T_{\text{ref}}} - (f \circ \Phi)(s_{T_{\text{ref}}})\|_{L^2(T_{\text{ref}})} \leq \|D^2(f \circ \Phi)\|_{L^2(T_{\text{ref}})} \leq |\det B|^{-1/2} \|B\|_F^2 \|D^2 f\|_{L^2(T)},$$

which finally leads to

$$\|f_T - f(s_T)\|_{L^2(T)} \leq 2C_{\text{ref}} h_T^2 \|D^2 f\|_{L^2(T)}.$$

3. step. It holds that $|\int_T f v_h dx - |T|f(s_T)v_h(s_T)| \leq \max\{C_P^2, 2C_{\text{ref}}\} h_T^2 \|\nabla f\|_{H^1(T)} \|v_h\|_{H^1(T)}$: The combination of step 1 and step 2 proves that

$$\begin{aligned} & \left| \int_T f v_h dx - |T|f(s_T)v_h(s_T) \right| \\ & \leq \max\{C_P^2, 2C_{\text{ref}}\} h_T^2 (\|\nabla f\|_{L^2(T)} \|\nabla v_h\|_{L^2(T)} + \|D^2 f\|_{L^2(T)} \|v_h\|_{L^2(T)}). \end{aligned}$$

Note that the brackets contain an \mathbb{R}^2 -scalar product which is estimated with the help of the Cauchy inequality $ab + cd \leq (a^2 + c^2)^{1/2}(b^2 + d^2)^{1/2}$. This concludes the proof of step 3.

4. step. With $C := \max\{C_P^2, 2C_{\text{ref}}\}$, we finally sum over all elements $T \in \mathcal{T}$ to obtain that

$$\begin{aligned} |F(v_h) - F_h(v_h)| & \leq \sum_{T \in \mathcal{T}} \left| \int_T f v_h dx - |T|f(s_T)v_h(s_T) \right| \\ & \leq C \sum_{T \in \mathcal{T}} \|h^2 \nabla f\|_{H^1(T)} \|v_h\|_{H^1(T)} \\ & \leq C \left(\sum_{T \in \mathcal{T}} \|h^2 \nabla f\|_{H^1(T)}^2 \right)^{1/2} \left(\sum_{T \in \mathcal{T}} \|v_h\|_{H^1(T)}^2 \right)^{1/2} \\ & = C \|h^2 \nabla f\|_{H^1(\mathcal{T})} \|v_h\|_{H^1(\Omega)} \end{aligned}$$

by use of the Cauchy inequality. This concludes the proof. \blacksquare

We stress that the proof does not work for $f \in H^1(\mathcal{T})$ since H^1 -functions are in general discontinuous so that the evaluation of f at s_T is not well-defined. However, for $f \in C^1(\mathcal{T})$, everything works well.

Exercise 36. For $f \in C^1(\mathcal{T})$, define $F \in H^1(\Omega)^*$ and $F_h \in \mathcal{S}^1(\mathcal{T})^*$ as in Proposition 3.15. Then, there holds

$$\|F - F_h\|_{\mathcal{S}^1(\mathcal{T})^*} \leq C \|h \nabla f\|_{L^\infty(\Omega)}, \quad (3.54)$$

where the constant $C > 0$ does neither depend on Ω nor \mathcal{T} or f . \square

However, if the volume force only satisfies $f \in H^1(\mathcal{T})$, one can proceed as follows:

Exercise 37. For $f \in H^1(\mathcal{T})$, define $F \in H^1(\Omega)^*$ as in Proposition 3.15 and $F_h \in \mathcal{S}^1(\mathcal{T})^*$ by $F_h(v_h) := \sum_{T \in \mathcal{T}} |T| f_T v_h(s_T)$, where $f_T := |T|^{-1} \int_T f dx$ denotes the integral mean. Then,

$$\|F - F_h\|_{\mathcal{S}^1(\mathcal{T})^*} \leq C \|h^2 \nabla f\|_{L^2(\Omega)}, \quad (3.55)$$

where the constant $C > 0$ does neither depend on Ω nor \mathcal{T} or f . \square

3.5.3 Trace Inequality

The treatment of the Neumann contributions of the right-hand side follows the same techniques as in the proof of Proposition 3.15. However, the analogous arguments lead to norms $\|v_h - v_{hE}\|_{L^2(E)}$ and $\|v_h\|_{L^2(E)}$ in step 1, where $E \in \mathcal{E}_N$ denotes a Neumann edge and where $v_{hE} = h_E^{-1} \int_E v_h ds$ denotes the integral mean over E . These two terms have to be estimated by $\|v_h\|_{H^1(T)}$, where $T \in \mathcal{T}$ is an element with edge E . This needs so-called trace inequalities.

Theorem 3.16 (Trace Inequality). *Let $T = \text{conv}\{z_0, \dots, z_d\} \subset \mathbb{R}^d$ be a simplex in \mathbb{R}^d with $|T| > 0$ and diameter $h_T := \text{diam}(T)$. Let $E = \text{conv}\{z_1, \dots, z_d\}$ denote one particular side of the simplex. Then, for $v \in H^1(T)$, it holds*

$$\|v\|_{L^2(E)}^2 \leq \frac{|E|}{|T|} (\|v\|_{L^2(T)}^2 + \frac{2}{d} h_T \|v \nabla v\|_{L^1(T)}) \leq \frac{|E|}{|T|} (1 + 2 h_T/d) \|v\|_{H^1(T)}^2. \quad (3.56)$$

With the integral means $v_T := |T|^{-1} \int_T v dx$ and $v_E := h_E^{-1} \int_E v ds$, it holds that

$$\|v - v_E\|_{L^2(E)}^2 \leq \|v - v_T\|_{L^2(E)}^2 \leq C \frac{|E| h_T^2}{|T|} \|\nabla v\|_{L^2(T)}^2, \quad (3.57)$$

where $C > 0$ depends only on the reference element T_{ref} and the dimension d .

The proof of the trace inequalities (3.56)–(3.57) is done with the help of the following lemma. In particular, we shall see that both estimates are sharp. Note that, for $d = 2$, it holds that $|E|/|T| \leq 2\varrho_T^{-1}$ and $|E| h_T^2/|T| \leq 2\sigma(T)h_T$.

Lemma 3.17 (Trace Identity). *Let $T = \text{conv}\{z_0, \dots, z_d\} \subset \mathbb{R}^d$ be a simplex in \mathbb{R}^d with $|T| > 0$. Let $E = \text{conv}\{z_1, \dots, z_d\}$ denote one particular side of the simplex. Then,*

$$\frac{1}{|E|} \int_E w ds = \frac{1}{|T|} \int_T w dx + \frac{1}{d|T|} \int_T (x - z_0) \cdot \nabla w dx \quad (3.58)$$

for all $w \in C^1(\overline{T})$.

Proof. We apply the Gauss Divergence Theorem to the function $f(x) := w(x)(x - z_0)$. With $\text{div } f(x) = \nabla w(x) \cdot (x - z_0) + dw(x)$, we obtain that

$$d \int_T w dx + \int_T (x - z_0) \cdot \nabla w(x) dx = \int_T \text{div } f dx = \int_{\partial T} f \cdot n ds.$$

Note that $(x - z_0) \cdot n(x) = 0$ for $x \in \partial T \setminus E$, whereas $(x - z_0) \cdot n(x) = \text{dist}(z_0, H)$, where $H \subset \mathbb{R}^d$ denotes the hyperplane with $E \subseteq H$. Therefore, the boundary integral simplifies to $\int_{\partial T} f \cdot n ds = \text{dist}(z_0, H) \int_E w ds$ and the latter equality reads

$$\frac{1}{|T|} \int_T w dx + \frac{1}{d|T|} \int_T (x - z_0) \cdot \nabla w dx = \frac{\text{dist}(z_0, H)|E|}{d|T|} \frac{1}{|E|} \int_E w ds,$$

which holds for any $w \in C^1(\overline{T})$. The special choice $w = 1$ can be used to determine the w -independent constant $\frac{\text{dist}(z_0, H)|E|}{d|T|} = 1$. This concludes the proof. \blacksquare

v10:

08.11.2017

Remark. Note that Lemma 3.17 holds for any $w \in W^{1,1}(T) := \{w \in L^1(T) \text{ weakly differentiable} \mid \nabla w \in L^1(T)^d\}$, even with the same proof. \square

Proof of Theorem 3.16. According to standard density arguments, it suffices to consider $v \in C^1(\overline{T})$. Plugging $w := v^2 \in C^1(\overline{T})$ into the trace identity (3.58), we see that

$$\frac{1}{|E|} \int_E v^2 ds = \frac{1}{|T|} \int_T v^2 dx + \frac{1}{d|T|} \int_T (x - z_0) \cdot (2v \nabla v) dx.$$

This is rewritten in the form

$$\begin{aligned} \frac{|T|}{|E|} \|v\|_{L^2(E)}^2 &= \|v\|_{L^2(T)}^2 + \frac{2}{d} \int_T (x - z_0) \cdot (v \nabla v) dx \leq \|v\|_{L^2(T)}^2 + \frac{2}{d} h_T \|v \nabla v\|_{L^1(T)} \\ &\leq (1 + 2h_T/d) \|v\|_{H^1(T)}^2 \end{aligned}$$

which proves (3.56). For the proof of (3.57), we simply replace v by $v - v_T$ and apply the Poincaré inequality. This leads to

$$\begin{aligned} \|v - v_T\|_{L^2(E)}^2 &\leq \frac{|E|}{|T|} (\|v - v_T\|_{L^2(T)}^2 + \frac{2}{d} h_T \|v - v_T\|_{L^2(T)} \|\nabla v\|_{L^2(T)}) \\ &\leq \frac{|E|}{|T|} (C_P^2 h_T^2 \|\nabla v\|_{L^2(T)}^2 + \frac{2}{d} C_P h_T^2 \|\nabla v\|_{L^2(T)}^2) \\ &= (C_P^2 + 2C_P/d) \frac{|E| h_T^2}{|T|} \|\nabla v\|_{L^2(T)}^2. \end{aligned}$$

The remaining estimate $\|v - v_E\|_{L^2(E)} \leq \|v - v_T\|_{L^2(E)}$ follows from the L^2 -best approximation property of the integral mean. \blacksquare

3.5.4 Approximation of Neumann Data

Finally, we consider the approximation of the Neumann contribution.

Proposition 3.18. Let $\phi \in C^2(\mathcal{E}_N) := \{\psi \in L^2(\Gamma_N) \mid \forall E \in \mathcal{E}_N \ \psi|_E \in C^2(E)\}$ and $F(v) := \int_{\Gamma_N} \phi v ds$ for $v \in H^1(\Omega)$. Let $F_h(v_h) := \sum_{E \in \mathcal{E}_N} h_E \phi(m_E) v_h(m_E)$ for $v_h \in \mathcal{S}^1(\mathcal{T})$, where $m_E \in \mathbb{R}^2$ denotes the midpoint of a Neumann edge $E \in \mathcal{E}_N$. With the mesh-size function $h \in L^\infty(\Gamma_N)$, $h|_E := h_E = \text{diam}(E)$, it then holds

$$\|F - F_h\|_{\mathcal{S}^1(\mathcal{T})^*} \leq C \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} := \max_{E \in \mathcal{E}_N} (h_E^{3/2} \max\{\|\phi'\|_{L^\infty(E)}, \|\phi''\|_{L^\infty(E)}\}) \quad (3.59)$$

where the constant $C > 0$ depends only on $\sigma(\mathcal{T})$ and $|\Gamma_N|$.

Proof. We aim to follow the lines of the proof of Proposition 3.15. For a Neumann edge $E \in \mathcal{E}_N$ and $w \in L^2(E)$, let $w_E := h_E^{-1} \int_E w ds$ denote the integral mean.

1. step. From $\int_E v_h ds = h_E v_h(m_E)$, we infer that

$$\begin{aligned} \int_E \phi v_h ds - h_E \phi(m_E) v_h(m_E) &= (\phi - \phi(m_E); v_h)_{L^2(E)} \\ &= (\phi - \phi(m_E); v_h - v_{hE})_{L^2(E)} + (\phi - \phi(m_E); v_{hE})_{L^2(E)} \\ &= (\phi - \phi_E; v_h - v_{hE})_{L^2(E)} + (\phi_E - \phi(m_E); v_h)_{L^2(E)} \\ &\leq \|\phi - \phi_E\|_{L^2(E)} \|v_h - v_{hE}\|_{L^2(E)} + \|\phi_E - \phi(m_E)\|_{L^2(E)} \|v_h\|_{L^2(E)}. \end{aligned}$$

where we have simply used orthogonality of $(\cdot)_E$. Therefore, the trace inequalities (3.56)–(3.57) yield that

$$\left| \int_E \phi v_h ds - h_E \phi(m_E) v_h(m_E) \right| \leq C \left(h_E^{1/2} \|\phi - \phi_E\|_{L^2(E)} + h_E^{-1/2} \|\phi_E - \phi(m_E)\|_{L^2(E)} \right) \|v_h\|_{H^1(T)},$$

where $T \in \mathcal{T}$ is an arbitrary element with $E \in \mathcal{E}_T$. The constant $C > 0$ depends only on $\sigma(\mathcal{T})$ and on $|\Gamma_N|$.

2. step. It holds that $\|\phi - \phi_E\|_{L^2(E)} \leq h_E^{3/2} \|\phi'\|_{L^\infty(E)}$: Note that $w := \phi - \phi_E \in C^1(E)$ has necessarily a zero $\zeta \in E$. Therefore, the fundamental theorem of calculus proves that

$$|w(x)| = \left| \int_\zeta^x w' ds \right| \leq h_E^{1/2} \|w'\|_{L^2(E)}.$$

Integration over E thus yields that

$$\|\phi - \phi_E\|_{L^2(E)}^2 = \|w\|_{L^2(E)}^2 = \int_E |w(x)|^2 ds_x \leq h_E^2 \|w'\|_{L^2(E)}^2 = h_E^2 \|\phi'\|_{L^2(E)}^2 \leq h_E^3 \|\phi'\|_{L^\infty(E)}^2.$$

3. step. It holds that $\|\phi_E - \phi(m_E)\|_{L^2(E)} \leq (1/2) h_E^{5/2} \|\phi''\|_{L^\infty(E)}$: Let $p \in \mathcal{P}^1(E)$ be a polynomial on E (with respect to the arclength) such that $\phi(m_E) = p(m_E)$ and $\phi'(m_E) = p'(m_E)$. Then,

$$\|\phi_E - \phi(m_E)\|_{L^2(E)} = h_E^{1/2} |\phi_E - \phi(m_E)| = h_E^{-1/2} \left| \int_E \phi ds - h_E p(m_E) \right| = h_E^{-1/2} \left| \int_E (\phi - p) ds \right|$$

With $w := \phi - p$ and hence $w'' = \phi''$, this implies that

$$\|\phi_E - \phi(m_E)\|_{L^2(E)} \leq h_E^{-1/2} \|w\|_{L^1(E)} \leq \|w\|_{L^2(E)}.$$

Note that w as well as w' have zeros at the edge midpoint m_E . Therefore, the same arguments as in step 2 (with the zero $\zeta = m_E$ and hence integration along a segment of length $h_E/2$) prove that

$$\|w\|_{L^2(E)}^2 \leq \frac{h_E^2}{2} \|w'\|_{L^2(E)}^2 \quad \text{as well as} \quad \|w'\|_{L^2(E)}^2 \leq \frac{h_E^2}{2} \|w''\|_{L^2(E)}^2 = \frac{h_E^2}{2} \|\phi''\|_{L^2(E)}^2.$$

Altogether, we see that

$$\|\phi_E - \phi(m_E)\|_{L^2(E)}^2 \leq \|w\|_{L^2(E)}^2 \leq \frac{h_E^4}{4} \|\phi''\|_{L^2(E)}^2 \leq \frac{h_E^5}{4} \|\phi''\|_{L^\infty(E)}^2.$$

4. step. The combination of the preceding steps proves that

$$\begin{aligned} \left| \int_E \phi v_h ds - h_E \phi(m_E) v_h(m_E) \right| &\leq C h_E^2 (\|\phi'\|_{L^\infty(E)} + \|\phi''\|_{L^\infty(E)}) \|v_h\|_{H^1(T)} \\ &\leq 2C h_E^2 \|\phi'\|_{C^1(E)} \|v_h\|_{H^1(T)} \\ &\leq 2C h_E^{1/2} \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} \|v_h\|_{H^1(T)} \end{aligned}$$

by definition of $\|w\|_{C^1(E)} := \max\{\|w\|_{L^\infty(E)}, \|w'\|_{L^\infty(E)}\}$.

5. step. We obtain the final result by summing over all Neumann edges $E \in \mathcal{E}_N$: For each $E \in \mathcal{E}_N$ we choose an element $T_E \in \mathcal{T}$ with $E \in \mathcal{E}_T$. Note that the element T_E can arise at most 3 times. Therefore,

$$\begin{aligned} |F(v_h) - F_h(v_h)| &\leq 2C \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} \sum_{E \in \mathcal{E}_N} h_E^{1/2} \|v_h\|_{H^1(T_E)} \\ &\leq 2C \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} \left(\sum_{E \in \mathcal{E}_N} h_E \right)^{1/2} \left(\sum_{E \in \mathcal{E}_N} \|v_h\|_{H^1(T_E)}^2 \right)^{1/2} \\ &\leq 2\sqrt{3} C |\Gamma_N|^{1/2} \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} \left(\sum_{T \in \mathcal{T}} \|v_h\|_{H^1(T)}^2 \right)^{1/2} \\ &= 2\sqrt{3} C |\Gamma_N|^{1/2} \|h^{3/2} \phi'\|_{C^1(\mathcal{E}_N)} \|v_h\|_{H^1(\Omega)}. \end{aligned}$$

This concludes the proof. ■

Exercise 38. (i) Extend the MATLAB code `solveLaplace` such that besides the coefficient vector of the Galerkin solution $u_h \in \mathcal{S}^1(\mathcal{T})$ even the energy $\|u_h\|^2 = \|\nabla u_h\|_{L^2(\Omega)}^2$ is returned. The Galerkin orthogonality yields that

$$\|u - u_h\|^2 = \|u\|^2 - \|u_h\|^2.$$

Even if the exact energy $\|u\|^2$ is unknown, it can be extrapolated by use of Aitkin's Δ^2 -method to obtain a good approximation of the error $\|u - u_h\|$.

(ii) Consider the homogenous Dirichlet problems

$$\begin{aligned} -\Delta u &= 1 \quad \text{in } \Omega, \\ u &= 0 \quad \text{on } \Gamma = \partial\Omega, \end{aligned}$$

with Ω being either the square $\Omega = (-1, 1)^2$ or the L -shaped domain $\Omega = (-1, 1)^2 \setminus [0, 1]^2$. Which experimental convergence rates $\|u - u_h\| = \mathcal{O}(h^\alpha)$ are observed? Do you expect that the solutions belong to $H^2(\Omega)$? **Hint:** For a convergent sequence $(x_j)_{j \in \mathbb{N}}$, the Δ^2 -sequence reads

$$y_j = x_j - \frac{(x_{j+1} - x_j)^2}{x_{j+2} - 2x_{j+1} + x_j}.$$

Under certain assumptions on $(x_j)_{j \in \mathbb{N}}$ the sequence $(y_j)_{j \in \mathbb{N}}$ then converges faster to $\lim_{j \rightarrow \infty} x_j$. \square

3.6 Inhomogeneous Dirichlet Data (08.11.2017)

Under the usual assumptions of the mixed boundary value problem of Section 2.3.2, we consider the boundary value problem

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= u_D && \text{on } \Gamma_D, \\ \partial_n u &= \phi && \text{on } \Gamma_N. \end{aligned} \quad (3.60)$$

The only difference to the problem treated above is the fact, that the Dirichlet data u_D might be nontrivial. A function $u \in C^2(\overline{\Omega})$ that solves (3.60) is called **strong solution** of (3.60), and the formulation (3.60) is called the **strong form** of the boundary value problem. A function $u \in H^1(\Omega)$ is **weak solution** of (3.60) provided that

$$\gamma u|_{\Gamma_D} = u_D \quad (3.61a)$$

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} + (\phi ; \gamma v)_{L^2(\Gamma_N)} \quad \text{for all } v \in H_D^1(\Omega). \quad (3.61b)$$

These two equations are referred to as the **weak form** of the boundary value problem (3.60). Note that the variational part (3.61b) of the weak form is the same as for the mixed boundary value problem with homogeneous Dirichlet conditions $u_D = 0$.

The following proposition shows that (3.60) and (3.61) are essentially equivalent and that the weak solution is unique. The unique solvability, however, needs certain assumptions on the Dirichlet data: If (3.61) has a solution $u \in H^1(\Omega)$, then it holds that $\gamma u|_{\Gamma_D} = u_D$, i.e., u_D can be extended from Γ_D to a function $\widehat{u}_D \in H^1(\Omega)$. With the same arguments as above, cf. Exercise 12 on page 18, one shows that

$$H^{1/2}(\Gamma_D) := \{\gamma u|_{\Gamma_D} \mid u \in H^1(\Omega)\} \quad \text{with norm} \quad \|v\|_{H^{1/2}(\Gamma_D)} = \inf \{\|\widehat{v}\|_{H^1(\Omega)} \mid \gamma \widehat{v}|_{\Gamma_D} = v\}$$

is a Hilbert space. Moreover, $H^{1/2}(\Gamma_D)$ is continuously embedded into $L^2(\Gamma_D)$, and the restriction operator $(\cdot)|_{\Gamma_D} : H^{1/2}(\Gamma) \rightarrow H^{1/2}(\Gamma_D)$ is well-defined and continuous.

Proposition 3.19. (i) *Provided that $u \in C^2(\overline{\Omega})$ solves the strong form (3.60), u solves also the weak form (3.61).*

(ii) *Provided that $f \in C(\overline{\Omega})$, $\phi \in C(\overline{\Gamma}_N)$, and $u_D \in C(\overline{\Gamma}_D)$ and that the weak solution $u \in H^1(\Omega)$ of (3.61) additionally satisfies $u \in C^2(\overline{\Omega})$, then u even solves the strong form (3.60).*

(iii) *Let $\widehat{u}_D \in H^1(\Omega)$ be an arbitrary extension of the Dirichlet data $u_D \in H^{1/2}(\Gamma)$. Given $f \in L^2(\Omega)$ and $\phi \in L^2(\Gamma_N)$, there exists a unique $u_0 \in H_D^1(\Omega)$ such that*

$$(\nabla u_0 ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} - (\nabla \widehat{u}_D ; \nabla v)_{L^2(\Omega)} + (\phi ; \gamma v)_{L^2(\Gamma_N)} \quad \text{for all } v \in H_D^1(\Omega). \quad (3.62)$$

(iv) *Under the assumptions of (iii), a function $u \in H^1(\Omega)$ with $\gamma u|_{\Gamma_D} = u_D$ solves the weak form (3.61), if and only if $u_0 := u - \widehat{u}_D \in H_D^1(\Omega)$ solves (3.62).*

(v) *Under the assumptions of (iii), there exists a unique weak solution $u \in H^1(\Omega)$ of (3.61). Contrary to $u_0 \in H_D^1(\Omega)$, however, the function $u \in H^1(\Omega)$ does not depend on the special choice of \widehat{u}_D .*

(vi) The weak solution $u \in H^1(\Omega)$ satisfies

$$\begin{aligned} \|u\|_{H^1(\Omega)} &\leq C_1 \left(\sup_{v \in H_D^1(\Omega) \setminus \{0\}} \frac{(f; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma_N) \setminus \{0\}} \frac{(\phi; w)_{L^2(\Gamma_N)}}{\|w\|_{H^{1/2}(\Gamma_N)}} + \|u_D\|_{H^{1/2}(\Gamma_D)} \right) \\ &\leq C_2 (\|f\|_{L^2(\Omega)} + \|\phi\|_{L^2(\Gamma_N)} + \|u_D\|_{H^{1/2}(\Gamma_D)}) \end{aligned} \quad (3.63)$$

where the constants $C_1, C_2 > 0$ only depend on Ω and Γ_D .

Proof. Note that the variational form (3.61b) does not consider whether u_D is zero or not. Therefore, the same proofs as for the mixed boundary value problem with homogeneous Dirichlet data apply to prove (i) and (ii). To verify (iii), simply note that the left-hand side of (3.62) defines an equivalent scalar product on the Hilbert space $H_D^1(\Omega)$. The right-hand side is linear and continuous on $H_D^1(\Omega)$. Therefore, existence and uniqueness of u_0 follows from the Riesz theorem. (iv) is obvious, and (v) thus an immediate consequence of (iii) and (iv). To prove the stability estimate, we argue as for the homogeneous Dirichlet conditions. With the Friedrichs inequality, we see that

$$\begin{aligned} C_F^{-2} \|u_0\|_{H^1(\Omega)}^2 &\leq \|\nabla u_0\|_{L^2(\Omega)}^2 \\ &= (\nabla u_0; \nabla u_0)_{L^2(\Omega)} \\ &= (f; \nabla u_0)_{L^2(\Omega)} + (\phi; \gamma u_0)_{L^2(\Gamma_N)} - (\nabla \hat{u}_D; \nabla u_0)_{L^2(\Omega)} \\ &\leq \|u_0\|_{H^1(\Omega)} \left(\sup_{v \in H_D^1(\Omega) \setminus \{0\}} \frac{(f; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma_N) \setminus \{0\}} \frac{(\phi; w)_{L^2(\Gamma_N)}}{\|w\|_{H^{1/2}(\Gamma_N)}} + \|\hat{u}_D\|_{H^1(\Omega)} \right) \end{aligned}$$

Second, the triangle inequality gives

$$\begin{aligned} \|u\|_{H^1(\Omega)} &\leq \|\hat{u}_D\|_{H^1(\Omega)} + \|u_0\|_{H^1(\Omega)} \\ &\leq (1 + C_F^2) \left(\sup_{v \in H_D^1(\Omega) \setminus \{0\}} \frac{(f; v)_{L^2(\Omega)}}{\|v\|_{H^1(\Omega)}} + \sup_{w \in H^{1/2}(\Gamma_N) \setminus \{0\}} \frac{(\phi; w)_{L^2(\Gamma_N)}}{\|w\|_{H^{1/2}(\Gamma_N)}} + \|\hat{u}_D\|_{H^1(\Omega)} \right). \end{aligned}$$

Taking the infimum over all \hat{u}_D , we conclude the stability estimate (3.63). \blacksquare

Remark. A first idea for the numerical approximation of the weak solution $u \in H^1(\Omega)$ of (3.60) might be the following:

- Construct an extension $\hat{u}_D \in H^1(\Omega)$ of the Dirichlet data.
- Discretize the variational form (3.62) by P1-FEM to obtain an approximation $u_{0h} \in \mathcal{S}_D^1(\mathcal{T})$ of $u_0 \in H_D^1(\Omega)$.
- Compute $u_h := u_{0h} + \hat{u}_D$ to obtain an approximation of u .

We stress, however, that then $u_h \notin \mathcal{S}^1(\mathcal{T})$ so that a postprocessing or evaluation of u_h is nontrivial. Moreover, we have to compute the scalar product $(\nabla \hat{u}_D; \nabla v_h)$ for discrete functions to build the load vector of the P1-FEM for u_0 . This leads to additional quadrature errors. Finally and most important, it might be hard to compute \hat{u}_D unless the Dirichlet data u_D are rather simple. \square

To overcome the difficulties mentioned in the previous remark, one uses the following approach in practice, which is then called P1-FEM of the weak form (3.61):

- Discretize Dirichlet data $u_D \in H^{1/2}(\Gamma_D)$ by some $u_{Dh} \in \mathcal{S}^1(\mathcal{T}|_{\Gamma_D}) := \{v_h|_{\Gamma_D} \mid v_h \in \mathcal{S}^1(\mathcal{T})\}$.
- Construct extension $\widehat{u}_{Dh} \in \mathcal{S}^1(\mathcal{T})$ with $\widehat{u}_{Dh}|_{\Gamma_D} = u_{Dh}$.
- With \widehat{u}_{Dh} replacing \widehat{u}_D , compute P1-FEM approximation $u_{0h} \in \mathcal{S}_D^1(\mathcal{T})$, cf. (3.62).
- Finally, define $u_h := u_{0h} + \widehat{u}_{Dh} \in \mathcal{S}^1(\mathcal{T})$ as approximation of the weak solution $u \in H^1(\Omega)$.

Note that the discrete solution $u_h \in \mathcal{S}^1(\mathcal{T})$ then belongs to the affine space $\widehat{u}_{Dh} + \mathcal{S}_D^1(\mathcal{T})$. The following result is the corresponding Céa-type lemma:

Lemma 3.20 (Céa lemma, first version). *Let $u \in H^1(\Omega)$ be the weak solution of (3.60). Let $\widehat{u}_{Dh} \in \mathcal{S}^1(\mathcal{T})$ be the approximate Dirichlet data and $u_{Dh} := \widehat{u}_{Dh}|_{\Gamma_D}$. Let $u_h \in \mathcal{S}^1(\mathcal{T})$ be the unique solution of*

$$\begin{aligned} u_h|_{\Gamma_D} &= u_{Dh} \\ (\nabla u_h; \nabla v_h)_{L^2(\Omega)} &= (f; v_h)_{L^2(\Omega)} + (\phi; v_h)_{L^2(\Gamma_N)} \quad \text{for all } v_h \in \mathcal{S}_D^1(\mathcal{T}). \end{aligned} \quad (3.64)$$

Then, u_h is quasioptimal in the sense that there exists a constant $C > 0$ such that

$$C^{-1} \|u - u_h\|_{H^1(\Omega)} \leq \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|u - (v_h + \widehat{u}_{Dh})\|_{H^1(\Omega)} = \min_{\substack{w_h \in \mathcal{S}^1(\mathcal{T}) \\ w_h|_{\Gamma_D} = u_D}} \|u - w_h\|_{H^1(\Omega)}. \quad (3.65)$$

The constant $C > 0$ only depends on Ω and Γ_D .

Proof. Note that the variational formulations (3.61) and (3.64) imply the Galerkin orthogonality

$$(\nabla(u - u_h); \nabla v_h)_{L^2(\Omega)} = 0 \quad \text{for all } v_h \in \mathcal{S}_D^1(\mathcal{T}).$$

We define $u_{0h} := u_h - \widehat{u}_{Dh} \in \mathcal{S}_D^1(\mathcal{T})$ and observe that

$$\begin{aligned} \|\nabla(u - u_h)\|_{L^2(\Omega)}^2 &= (\nabla(u - u_h); \nabla(u - [u_{0h} + \widehat{u}_{Dh}]))_{L^2(\Omega)} \\ &= (\nabla(u - u_h); \nabla(u - [v_h + \widehat{u}_{Dh}]))_{L^2(\Omega)} \\ &\leq \|\nabla(u - u_h)\|_{L^2(\Omega)} \|\nabla(u - [v_h + \widehat{u}_{Dh}])\|_{L^2(\Omega)} \end{aligned}$$

for each $v_h \in \mathcal{S}_D^1(\mathcal{T})$. Next, recall that $\|v\| := \|\nabla v\|_{L^2(\Omega)} + \|\gamma v\|_{L^2(\Gamma_D)}$ provides an equivalent norm on $H^1(\Omega)$, i.e., there are constants $C_1, C_2 > 0$ such that $C_1^{-1}\|v\| \leq \|v\|_{H^1(\Omega)} \leq C_2\|v\|$ for all $v \in H^1(\Omega)$. Consequently,

$$\begin{aligned} C_2^{-1} \|u - u_h\|_{H^1(\Omega)} &\leq \|\nabla(u - u_h)\|_{L^2(\Omega)} + \|\gamma(u - u_h)\|_{L^2(\Gamma_D)} \\ &= \|\nabla(u - u_h)\|_{L^2(\Omega)} + \|u_D - u_{Dh}\|_{L^2(\Gamma_D)} \\ &\leq \|\nabla(u - [v_h + \widehat{u}_{Dh}])\|_{L^2(\Omega)} + \|\gamma(u - [v_h + \widehat{u}_{Dh}])\|_{L^2(\Gamma_D)} \\ &\leq C_1 \|u - (v_h + \widehat{u}_{Dh})\|_{H^1(\Omega)} \end{aligned}$$

for all $v_h \in \mathcal{S}_D^1(\mathcal{T})$. This proves (3.65) with an infimum on the right-hand side. Standard arguments show that this infimum is, in fact, attained. \blacksquare

Exercise 39. Proof that (3.64) has a unique solution $u_h \in \mathcal{S}^1(\mathcal{T})$. □

Remark. Note that Lemma 3.20 is independent of how the Dirichlet data are actually discretized, but the discretization enters the right-hand side, since it constraints the affine space for the minimum in (3.65). Later on, we shall see that appropriate discretization $u_{Dh} = J_h u_D$ by means of the Scott-Zhang projection J_h even guarantees that

$$\|u - u_h\|_{H^1(\Omega)} \leq C \min_{w_h \in \mathcal{S}^1(\mathcal{T})} \|u - w_h\|_{H^1(\Omega)},$$

where the right-hand side is independent of how u_D is actually discretized; see also Exercise 41–42 below. □

Remark. If the Dirichlet data u_D have an extension $\widehat{u}_D \in H^2(\Omega)$ with $\gamma \widehat{u}_D|_{\Gamma_D} = u_D$, then u_D is continuous. We define $\widehat{u}_{Dh} \in \mathcal{S}^1(\mathcal{T})$ nodewise by

$$\widehat{u}_{Dh}(z) = \begin{cases} u_D(z) & \text{for } z \in \overline{\Gamma}_D, \\ 0 & \text{else,} \end{cases}$$

for $z \in \mathcal{K}$. Let $u \in H^1(\Omega)$ denote the weak solution of (3.60) and $u_0 := u - \widehat{u}_D \in H_D^1(\Omega)$. We additionally define $\widetilde{u}_{Dh} \in \mathcal{S}_D^1(\mathcal{T})$ nodewise by

$$\widetilde{u}_{Dh}(z) = \begin{cases} 0 & \text{for } z \in \overline{\Gamma}_D, \\ \widehat{u}_D(z) & \text{else,} \end{cases}$$

for $z \in \mathcal{K}$. Note that the nodal interpolant of \widehat{u}_D reads $I_h \widehat{u}_D = \widehat{u}_{Dh} + \widetilde{u}_{Dh}$ and that $\|\widehat{u}_D - I_h \widehat{u}_D\|_{H^1(\Omega)} = \mathcal{O}(h)$ decays with optimal order. Consequently, we may plug-in $u = \widehat{u}_D + u_0$ into Céa's lemma to observe that

$$\begin{aligned} C^{-1} \|u - u_h\|_{H^1(\Omega)} &\leq \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|u - (\widehat{u}_{Dh} + v_h)\|_{H^1(\Omega)} \\ &= \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|(\widehat{u}_D - I_h \widehat{u}_D) + (u_0 - v_h + \widetilde{u}_{Dh})\|_{H^1(\Omega)} \\ &= \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|(\widehat{u}_D - I_h \widehat{u}_D) + (u_0 - w_h)\|_{H^1(\Omega)} \\ &\leq \|\widehat{u}_D - I_h \widehat{u}_D\|_{H^1(\Omega)} + \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|u_0 - w_h\|_{H^1(\Omega)}. \end{aligned}$$

Conversely, it holds that

$$\begin{aligned} \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|u_0 - w_h\|_{H^1(\Omega)} &= \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|(u - \widehat{u}_D) - w_h\|_{H^1(\Omega)} \\ &= \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|u - (w_h + I_h \widehat{u}_D) - (\widehat{u}_D - I_h \widehat{u}_D)\|_{H^1(\Omega)} \\ &\leq \min_{w_h \in \mathcal{S}_D^1(\mathcal{T})} \|u - (w_h + I_h \widehat{u}_D)\|_{H^1(\Omega)} + \|\widehat{u}_D - I_h \widehat{u}_D\|_{H^1(\Omega)} \\ &\leq \|u - u_h\|_{H^1(\Omega)} + \|\widehat{u}_D - I_h \widehat{u}_D\|_{H^1(\Omega)}. \end{aligned}$$

Therefore, the proposed P1-FEM for the approximation of $u \in H^1(\Omega)$ converges with the same order as the P1-FEM for the approximation of $u_0 \in H_D^1(\Omega)$. □

The inhomogeneous Dirichlet problem allows the proof that the trace operator has a right inverse \mathcal{L} . This inverse is called *lifting operator*.

Exercise 40. Let $\gamma \in L(H^1(\Omega); H^{1/2}(\Gamma))$ denote the trace operator. Prove that there exists a lifting operator $\mathcal{L} \in L(H^{1/2}(\Gamma); H^1(\Omega))$ such that $\gamma\mathcal{L}v = v$ for all $v \in H^{1/2}(\Gamma)$. **Hint.** Consider an appropriate Dirichlet-Problem with inhomogeneous Dirichlet data $v \in H^{1/2}(\Gamma)$ and let $u := \mathcal{L}v \in H^1(\Omega)$ denote the unique solution. \square

The assumptions of the following exercise will be satisfied for the Scott-Zhang projection.

Exercise 41 (Céa lemma, second version). Suppose that there exists a linear projection $P_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ with the following properties

- (i) $\|P_h v\|_{H^1(\Omega)} \leq C_{\text{stab}} \|v\|_{H^1(\Omega)}$ for all $v \in H^1(\Omega)$
- (ii) $P_h v_h = v_h$ for all $v_h \in \mathcal{S}^1(\mathcal{T})$
- (iii) $(P_h v)|_\omega = v|_\omega$ for all $v \in H^1(\Omega)$ with $v|_\omega \in \mathcal{S}^1(\mathcal{T}|_\omega)$ and $\omega \in \{\Gamma, \Gamma_D\}$
- (iv) $(P_h v)|_\omega$ depends only on the trace $v|_\omega$ for all $v \in H^1(\Omega)$ and $\omega \in \{\Gamma, \Gamma_D\}$

Then, for $u \in H^1(\Omega)$ being the solution of (3.60) and $u_{Dh} := (P_h u)|_{\Gamma_D}$, it holds that

$$\min_{\substack{v_h \in \mathcal{S}^1(\mathcal{T}) \\ v_h|_{\Gamma_D} = u_{Dh}}} \|u - v_h\|_{H^1(\Omega)} \leq C \min_{w_h \in \mathcal{S}^1(\mathcal{T})} \|u - w_h\|_{H^1(\Omega)},$$

where $C > 0$ depends only on the stability constant C_{stab} , Ω , and Γ_D . In particular, this implies an unconstrained Céa lemma for the mixed boundary value problem with inhomogeneous Dirichlet data, i.e., under the assumptions of Lemma 3.20 and with $u_{Dh} = (P_h u)|_{\Gamma_D}$, it holds

$$\|u - u_h\|_{H^1(\Omega)} \leq C \min_{w_h \in \mathcal{S}^1(\mathcal{T})} \|u - w_h\|_{H^1(\Omega)}.$$

Hint. Let $w \in H^1(\Omega)$ be the weak solution of $\Delta w = 0$ in Ω subject to the boundary conditions $w = u - u_{Dh}$ on Γ_D and $\partial_n w = 0$ on Γ_N . Define $u_0 := u - w$. Prove that $u_0 = u_{Dh}$ on Γ_D and $\|u - u_0\|_{H^1(\Omega)} \simeq \|u - u_{Dh}\|_{H^{1/2}(\Gamma)}$. Choose $v_h := P_h u_0$. \square

The existence of a Scott-Zhang-type projection is essentially equivalent to the validity of the Céa lemma.

Exercise 42. (a) Suppose that $P_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ satisfies the properties (i)–(iii) of Exercise 41 for $\omega = \Gamma$ only. Then, for all $u \in H^1(\Omega)$ and all $u_{Dh} \in \mathcal{S}^1(\mathcal{T})$, it holds that

$$\min_{\substack{v_h \in \mathcal{S}^1(\mathcal{T}) \\ v_h|_\Gamma = u_{Dh}}} \|u - v_h\|_{H^1(\Omega)} \leq C \left[\min_{w_h \in \mathcal{S}^1(\mathcal{T})} \|u - w_h\|_{H^1(\Omega)} + \|u - u_{Dh}\|_{H^{1/2}(\Gamma)} \right], \quad (3.66)$$

where $C > 0$ depends only on Ω and the stability constant C_{stab} . **Hint.** Argue along the lines of Exercise 41.

(b) Suppose that (3.66) holds true. Then, there exists a linear projection $P_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ which satisfies the properties (i)–(iii) of Exercise 41. **Hint.** For given $u \in H^1(\Omega)$, let $u_{Dh} \in \mathcal{S}^1(\mathcal{T}|_\Gamma)$ be the $H^{1/2}(\Gamma)$ -best approximation of $u|_\Gamma \in H^{1/2}(\Gamma)$. With this, let

$P_h u := u_h \in \mathcal{S}^1(\mathcal{T})$ be the FEM solution of the inhomogeneous Dirichlet problem with discrete Dirichlet data u_{Dh} . □

Remark. Note that, for inhomogeneous Dirichlet data, it holds that

$$\|u - u_h\|^2 \neq \|u\|^2 - \|u_h\|^2$$

in general. Therefore, we cannot proceed as in Exercise 38 to approximate the error. Instead, in academic examples, where u is known, one has to compute

$$\|u - u_h\|^2 = \sum_{T \in \mathcal{T}} \|\nabla u - \nabla u_h\|_{L^2(T)}^2$$

by \mathcal{T} -piecewise numerical quadrature. □

Exercise 43. Write a MATLAB code for the P1-FEM for the mixed boundary value problem (3.60) with inhomogeneous but continuous Dirichlet data u_D . To verify the code, consider the Dirichlet problem

$$\begin{aligned} -\Delta u &= 1 & \text{in } \Omega &= [0, 1]^2, \\ u &= 1 & \text{on } \Gamma. \end{aligned}$$

If u_0 denotes the solution of the corresponding homogeneous problem, then it holds that $u = u_0 + 1$. □

Chapter 4

A Posteriori Analysis

4.1 Introduction (17.11.2014)

We consider the model problem

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= 0 && \text{on } \Gamma_D, \\ \partial u / \partial n &= \phi && \text{on } \Gamma_N. \end{aligned} \tag{4.1}$$

Let $u \in H_D^1(\Omega)$ be the weak solution of (4.1) and $u_h \in \mathcal{S}_D^1(\mathcal{T})$ be the P1-FEM approximation of u . In the previous chapter, we aimed to control the error $\|u - u_h\|_{H^1(\Omega)}$ by a priori knowledge, e.g., regularity of the given data and the exact solution (but essentially independent of the discrete solution u_h). Since u is unknown, in general, the **a priori analysis** provides a qualitative understanding of the FEM, e.g., convergence with certain rates, but the derived bounds are non-computable in practice. In this chapter, we aim to derive *numerically computable* bounds $\eta = \eta(u_h, f, \phi, \mathcal{T})$ for the error $\|u - u_h\|_{H^1(\Omega)}$, which may depend on u_h , the triangulation \mathcal{T} , and the given data f and ϕ (but *not* on the exact solution u). The quantity η is referred to as (**a posteriori**) **error estimator**, and emphasis is laid on the fact that η can be computed algorithmically as soon as the discrete solution $u_h \in \mathcal{S}_D^1(\mathcal{T})$ has been computed. An error estimator η is called **reliable** provided that

$$\|u - u_h\|_{H^1(\Omega)} \leq C_{\text{rel}} \eta. \tag{4.2}$$

Usually, the information η provides, is used to steer a mesh-refinement that leads to a sequence \mathcal{T}_ℓ of regular meshes with nested spaces $\mathcal{S}_D^1(\mathcal{T}_\ell) \subseteq \mathcal{S}_D^1(\mathcal{T}_{\ell+1})$, i.e., $\mathcal{T}_{\ell+1}$ is a certain refinement of \mathcal{T}_ℓ . If η is reliable the (numerically or algorithmically observed) decrease of η to zero implies the convergence of u_h towards u . However, it might (formally) occur that u_h tends to u , while η does not tend to zero. Therefore, an error estimator η is called **efficient** provided that

$$C_{\text{eff}} \eta \leq \|u - u_h\|_{H^1(\Omega)}. \tag{4.3}$$

For an efficient error estimator η , the convergence of u_h to u necessarily implies the convergence of η to zero. Finally, if η is reliable and efficient, we observe for η the same order of convergence as for $\|u - u_h\|_{H^1(\Omega)}$.

The aim of a posteriori error estimates is twofold:

- We want to control the accuracy $\|u - u_h\|_{H^1(\Omega)}$ of a discrete solution u_h and stop the computation if u_h is sufficiently accurate.
- The mesh-refinement should be steered automatically by the algorithm so that we are led to the highest possible accuracy with the lowest number of degrees of freedom.

Remark. Throughout, we allow the cases $\Gamma_D = \Gamma$ as well as $\Gamma_D = \emptyset$. In the latter case, (4.1) becomes the Neumann problem, for which we have to assume the compatibility condition $\int_{\Omega} f \, dx + \int_{\Gamma} \phi \, ds = 0$. Then, $u \in H_*^1(\Omega)$ and, even more important, the test space $H_*^1(\Omega)$ in the weak formulation can equivalently be replaced by the entire space $H^1(\Omega)$. The same holds for the P1-FEM, where $u_h \in \mathcal{S}_*^1(\mathcal{T})$ and where the discrete test is $\mathcal{S}_*^1(\mathcal{T})$ or equivalently $\mathcal{S}^1(\mathcal{T})$. \square

4.2 Scott-Zhang Projection (02.12.2014)

Since H^1 -functions are in general not continuous, nodal interpolation requires additional regularity assumptions. In this section, we aim to provide some quasi-interpolation operator which is well-defined for all $u \in H^1(\Omega)$ and also has the projection property. We start with the following elementary lemma

Lemma 4.1. *For $z \in \mathcal{K}$, choose an edge $E_z \in \mathcal{E}$ with $z \in E_z$. Then, there is a unique dual function $\psi_z \in \mathcal{P}^1(E_z)$ such that*

$$\int_{E_z} \psi_z \zeta_{z'} \, ds = \delta_{zz'} \quad \text{for all } z' \in \mathcal{K}. \quad (4.4)$$

Moreover, it holds that $\|\psi_z\|_{L^\infty(E_z)} \leq C |E_z|^{-1}$ for some generic constant $C > 0$, which is in particular independent of z and \mathcal{T} .

Proof. According to the Riesz theorem, there is a unique function $\widehat{\psi} \in \mathcal{P}^1[0, 1]$ such that

$$\int_0^1 \widehat{\psi} \widehat{\phi} \, dt = \widehat{\phi}(0) \quad \text{for all } \widehat{\phi} \in \mathcal{P}^1[0, 1].$$

Let $\Phi_z : [0, 1] \rightarrow E_z$ be an affine parametrization of the edge E_z with $\Phi_z(0) = z$. We define

$$\psi_z := \frac{1}{|E_z|} \widehat{\psi} \circ \Phi_z^{-1} \in \mathcal{P}^1(E_z).$$

Clearly, $\|\psi_z\|_{L^\infty(E_z)} \leq \|\widehat{\psi}\|_{L^\infty(0,1)} |E_z|^{-1}$. Note that $|\Phi_z'| = |E_z|$ and hence

$$\int_{E_z} \psi_z \zeta_{z'} \, ds = \int_0^1 (\psi_z \circ \Phi_z) (\zeta_{z'} \circ \Phi_z) |\Phi_z'| \, dt = \int_0^1 \widehat{\psi}(t) (\zeta_{z'} \circ \Phi_z)(t) \, dt = \zeta_{z'}(\Phi_z(0)) = \zeta_{z'}(z).$$

This concludes the proof. \blacksquare

Definition. For each node $z \in \mathcal{K}$ of \mathcal{T} , we choose an edge $E_z \in \mathcal{E}$ such that

- $E_z \subseteq \overline{\Gamma}_D$ for $z \in \overline{\Gamma}_D$,

- $E_z \subseteq \Gamma$ for $z \in \Gamma$,
- E_z arbitrary for $z \in \Omega$.

Note that the precise choice is immaterial for the following analysis. For $z \in \mathcal{K}$, let $\psi_z \in \mathcal{P}^1(E_z)$ be the corresponding dual function. Then, the **Scott-Zhang projection** is defined by

$$J_h v := \sum_{z \in \mathcal{K}} \left(\int_{E_z} \psi_z v ds \right) \zeta_z. \quad (4.5)$$

Clearly, $J_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ is well-defined and linear. Our first proposition states that J_h is in fact a projection which preserves discrete boundary data (cf. the properties needed in Exercise 41).

v12:
17.11.2017

Proposition 4.2. *For $v \in H^1(\Omega)$ and $v_h \in \mathcal{S}^1(\mathcal{T})$, the following properties (i)–(iii) are true:*

- (i) $J_h v_h = v_h$.
- (ii) $(J_h v)|_\omega$ depends only on the trace $v|_\omega$ for $\omega \in \{\Gamma, \Gamma_D\}$.
- (iii) $v|_\omega = v_h|_\omega$ implies that $(J_h v)|_\omega = v|_\omega$ for $\omega \in \{\Gamma, \Gamma_D\}$.

Proof. (i) Note that $v_h = \sum_{z' \in \mathcal{K}} v_h(z') \zeta_{z'}$. By choice of ψ_z , this shows

$$\int_{E_z} \psi_z v_h ds = \sum_{z' \in \mathcal{K}} v_h(z') \int_{E_z} \psi_z \zeta_{z'} ds = v_h(z).$$

With this, we deduce

$$J_h v_h = \sum_{z \in \mathcal{K}} \left(\int_{E_z} \psi_z v_h ds \right) \zeta_z = \sum_{z \in \mathcal{K}} v_h(z) \zeta_z = v_h.$$

(ii) follows from the choice of the edges E_z . (iii) We consider only $\omega = \Gamma_D$. We first note that

$$(J_h w)|_{\Gamma_D} = \sum_{z \in \mathcal{K}} \left(\int_{E_z} \psi_z w ds \right) \zeta_z|_{\Gamma_D} = \sum_{z \in \mathcal{K} \cap \bar{\Gamma}_D} \left(\int_{E_z} \psi_z w ds \right) \zeta_z|_{\Gamma_D} \quad \text{for all } w \in H^1(\Omega).$$

For $z \in \mathcal{K} \cap \bar{\Gamma}_D$, it holds that $E_z \subseteq \bar{\Gamma}_D$ and hence $\int_{E_z} \psi_z v_h ds = \int_{E_z} \psi_z v ds$. Together with the last equation and the projection property, we obtain that

$$(J_h v)|_{\Gamma_D} = (J_h v_h)|_{\Gamma_D} = v_h|_{\Gamma_D}.$$

This concludes the proof. ■

Exercise 44. Show that Lemma 4.1 holds for any dimension $d \geq 2$. □

Note that the Scott-Zhang projection $J_h v$ is not defined for general L^2 -functions, since $L^2(T)$ does not provide traces. However, one can define an appropriate variant as follows:

Exercise 45. Construct a linear projection $P_h : L^2(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ which satisfies

- $\|P_h v\|_{L^2(\Omega)} \leq C \|v\|_{L^2(\Omega)}$ for all $v \in L^2(\Omega)$.
- $P_h v_h = v_h$ for all $v_h \in \mathcal{S}^1(\mathcal{T})$.

The constant $C > 0$ may only depend on $\sigma(\mathcal{T})$. **Hint.** Proceed as for the standard Scott-Zhang projection. Instead of an edge E_z , associate with each node $z \in \mathcal{K}$ an arbitrary element $T_z \in \mathcal{T}$ with $z \in T_z$. \square

Next, we aim to show that the Scott-Zhang projection has local stability and approximation properties. Unlike nodal interpolation, this will require appropriate patches.

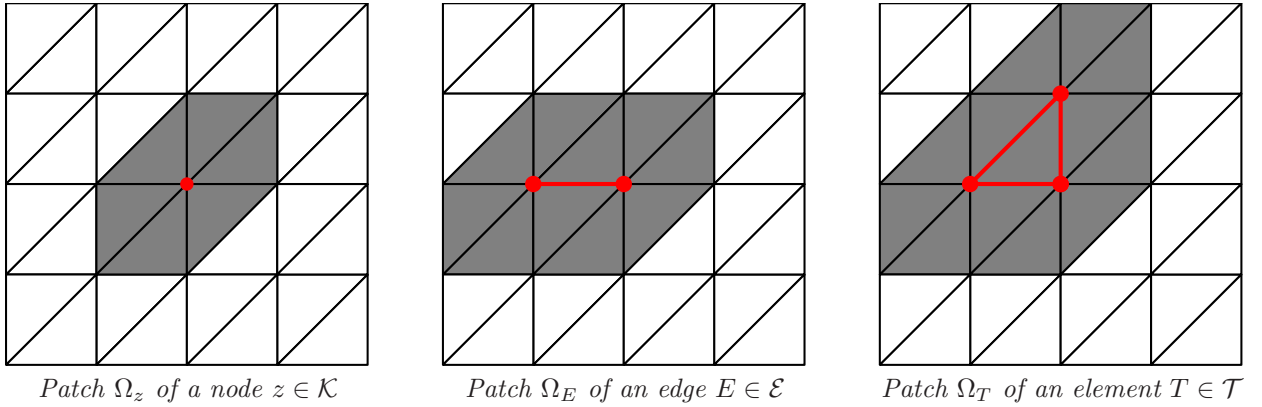


FIGURE 4.1. For the a posteriori analysis, we need three types of patches $\omega \subseteq \Omega$, namely patches of nodes, edges, and elements, respectively. Note that the patch of an edge (or of an element) just is the union of the patches of its nodes.

Definition. For the a posteriori analysis, we need certain unions of elements, called **patches**, cf. Figure 4.1: For a node $z \in \mathcal{K}$, we define

$$\tilde{\Omega}_z := \{T \in \mathcal{T} \mid z \in \mathcal{K}_T\} \quad \text{as well as} \quad \Omega_z := \bigcup \tilde{\Omega}_z := \{x \in \mathbb{R}^2 \mid \exists T \in \tilde{\Omega}_z \quad x \in T\}. \quad (4.6)$$

For an edge $E \in \mathcal{E}$, we define

$$\tilde{\Omega}_E := \{T \in \mathcal{T} \mid \mathcal{K}_E \cap T \neq \emptyset\} = \{T \in \tilde{\Omega}_z \mid z \in \mathcal{K}_E\} \quad \text{as well as} \quad \Omega_E := \bigcup \tilde{\Omega}_E. \quad (4.7)$$

Finally, for an element $T \in \mathcal{T}$, we define

$$\tilde{\Omega}_T := \{T' \in \mathcal{T} \mid \mathcal{K}_T \cap T' \neq \emptyset\} = \{T' \in \tilde{\Omega}_z \mid z \in \mathcal{K}_T\} \quad \text{as well as} \quad \Omega_T := \bigcup \tilde{\Omega}_T. \quad (4.8)$$

The patches Ω_z , Ω_E , and Ω_T are visualized in Figure 4.1.

Lemma 4.3. *There is a constant $C > 0$ which depends only on $\sigma(\mathcal{T})$, such that*

- $\#\tilde{\Omega}_z \leq C$ for all $z \in \mathcal{K}$,
- $\#\tilde{\Omega}_E \leq C$ for all $E \in \mathcal{E}$,

- $\#\tilde{\Omega}_T \leq C$ for all $T \in \mathcal{T}$,

i.e., the number of elements per patch is uniformly bounded. Moreover,

- $\#\{T' \in \mathcal{T} \mid T \in \tilde{\Omega}_{T'}\} \leq C$ for all $T \in \mathcal{T}$,

i.e., an element $T \in \mathcal{T}$ belongs only to finitely many patches.

Proof. Note that $\sigma(\mathcal{T})$ provides a bound for the minimal interior angle of all elements $T \in \mathcal{T}$; see Exercise 16. Consequently, there is a maximal number $C > 0$ of elements in $\tilde{\Omega}_z$, for all nodes $u \in \mathcal{K}$. By definition, there follows $\#\tilde{\Omega}_E \leq 2C$ as well as $\#\tilde{\Omega}_T \leq 3C$. ■

An essential consequence of Lemma 4.3 is that

$$\|v\|_{L^2(\Omega)} \leq \left(\sum_{T \in \mathcal{T}} \|v\|_{L^2(\Omega_T)}^2 \right)^{1/2} \leq C_{\text{patch}} \|v\|_{L^2(\Omega)} \quad \text{for all } v \in L^2(\Omega),$$

where $C_{\text{patch}} > 0$ depends only on $\sigma(\mathcal{T})$. Another consequence of Lemma 4.3 is that the diameter $\text{diam}(\Omega_T)$ of a patch is proportional to $h_T = \text{diam}(T)$. This is stated in the following lemma.

Lemma 4.4. For a regular triangulation, it holds that

- $\text{diam}(\Omega_z) \leq C h_T$ for all $z \in \mathcal{K}$ and $T \in \tilde{\Omega}_z$,
- $\text{diam}(\Omega_E) \leq C h_E \leq C h_T$ for all $E \in \mathcal{E}$ and $T \in \tilde{\Omega}_E$,
- $\text{diam}(\Omega_{T'}) \leq C h_T$ for all $T' \in \mathcal{T}$ and $T \in \tilde{\Omega}_{T'}$.

The constant $C > 0$ depends only on $\sigma(\mathcal{T})$.

Proof. 1. step. Note that $h_T \leq \sigma(\mathcal{T}) \varrho_T \leq \sigma(\mathcal{T}) h_E$ for all $T \in \mathcal{T}$ and all edges $E \in \mathcal{E}_T$.

2. step. Patch of a node $z \in \mathcal{K}$: For $\tilde{\Omega}_z = \{T_1, \dots, T_n\}$, we may choose a numbering such that T_{j-1}, T_j are neighbours, i.e., $T_{j-1} \cap T_j \in \mathcal{E}$. From step 1, we derive $h_{T_{j-1}} \leq \sigma(\mathcal{T}) h_{T_j}$, whence $h_{T'} \leq \sigma(\mathcal{T})^{n-1} h_T$ for all $T, T' \in \tilde{\Omega}_z$. This yields that

$$\text{diam}(\Omega_z) \leq 2 \max_{T' \in \tilde{\Omega}_z} h_{T'} \leq 2\sigma(\mathcal{T})^{n-1} h_T \quad \text{for all } T \in \tilde{\Omega}_z.$$

3. step. Patch of an edge $E \in \mathcal{E}$: With $E = \text{conv}\{z_1, z_2\}$ for some $z_1, z_2 \in \mathcal{K}$, it holds that $\tilde{\Omega}_E = \tilde{\Omega}_{z_1} \cup \tilde{\Omega}_{z_2}$ as well as $\tilde{\Omega}_{z_1} \cap \tilde{\Omega}_{z_2} \neq \emptyset$. Let $T \in \tilde{\Omega}_E$ and $n := \max\{\#\tilde{\Omega}_{z_1}, \#\tilde{\Omega}_{z_2}\}$. Without loss of generality, we may assume $T \in \tilde{\Omega}_{z_1}$. Choose $T' \in \tilde{\Omega}_{z_1} \cap \tilde{\Omega}_{z_2}$. Then,

$$\text{diam}(\Omega_E) \leq \text{diam}(\Omega_{z_1}) + \text{diam}(\Omega_{z_2}) \leq 2\sigma(\mathcal{T})^{n-1} (h_T + h_{T'}) \leq 2\sigma(\mathcal{T})^{n-1} (1 + \sigma(\mathcal{T})^{n-1}) h_T.$$

4. step. Patch of an element $T \in \mathcal{T}$: Simply use the same arguments as in step 3. ■

The Scott-Zhang projection is locally H^1 -stable and has a local first-order approximation property.

Proposition 4.5. For all $T \in \mathcal{T}$, it holds that

$$\|v - J_h v\|_{L^2(T)} + h_T \|\nabla J_h v\|_{L^2(T)} \leq C h_T \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H^1(\Omega). \quad (4.9)$$

The constant $C > 0$ depends only on γ -shape regularity of \mathcal{T} .

The proof requires the following technical lemma which is also valid in any dimension $d \geq 2$ as the proof reveals.

Lemma 4.6 (Generalized Poincaré-Friedrichs inequality). Let $v \in H^1(\Omega)$, $T \in \mathcal{T}$, $T' \in \tilde{\Omega}_T$, and $E' \in \mathcal{E}_{T'}$. Define the integral means $v_T := (1/|T|) \int_T v dx$, $v_{T'} := (1/|T'|) \int_{T'} v dx$, and $v_{E'} := (1/|E'|) \int_{E'} v ds$. Then,

$$\|v_T - v_{T'}\|_{L^2(T)} + \|v_T - v_{E'}\|_{L^2(T)} \leq C h_T \|\nabla v\|_{L^2(\Omega_T)}. \quad (4.10)$$

In particular, this implies that

$$\|v - v_{T'}\|_{L^2(\Omega_T)} + \|v - v_{E'}\|_{L^2(\Omega_T)} \leq C h_T \|\nabla v\|_{L^2(\Omega_T)}. \quad (4.11)$$

In either estimate, the constant $C > 0$ depends only on γ -shape regularity of \mathcal{T} , but is independent of Ω and the shape of Ω_T .

Proof. To ease the notation, let $v_E := (1/|E|) \int_E v ds$ also denote the integral mean over edges. Let T_{ref} denote the reference triangle and $E_{\text{ref}} = [0, 1]$ be the reference edge of T_{ref} .

1. step. For any $w \in H^1(T_{\text{ref}})$, it holds that $|w_{T_{\text{ref}}} - w_{E_{\text{ref}}}| \leq C_{\text{ref}} \|\nabla w\|_{L^2(T_{\text{ref}})}$, where $C_{\text{ref}} > 0$ depends only on T_{ref} : With the trace inequality (3.57) on T_{ref} , we see that

$$\begin{aligned} |w_{T_{\text{ref}}} - w_{E_{\text{ref}}}| &\leq |E_{\text{ref}}|^{-1} \|w - w_{T_{\text{ref}}}\|_{L^1(E_{\text{ref}})} \leq |E_{\text{ref}}|^{-1/2} \|w - w_{T_{\text{ref}}}\|_{L^2(E_{\text{ref}})} \\ &\leq C \frac{h_{T_{\text{ref}}}}{|T_{\text{ref}}|^{1/2}} \|\nabla w\|_{L^2(T_{\text{ref}})} =: C_{\text{ref}} \|\nabla w\|_{L^2(T_{\text{ref}})} \end{aligned}$$

2. step. For all $T \in \mathcal{T}$, $E \in \mathcal{E}_T$, it holds that $\|v_T - v_E\|_{L^2(T)} \leq C h_T \|\nabla v\|_{L^2(T)}$, where $C > 0$ depends only on C_{ref} : Let $\Phi_T : T_{\text{ref}} \rightarrow T$ be an affine diffeomorphism with $\Phi_T(E_{\text{ref}}) = E$. Let $B \in \mathbb{R}^{2 \times 2}$ denote the linear part. Define $w := v \circ \Phi_T \in H^1(T_{\text{ref}})$. Then, the transformation formula gives

$$|v_T - v_E| = |w_{T_{\text{ref}}} - w_{E_{\text{ref}}}| \leq C_{\text{ref}} \|\nabla w\|_{L^2(T_{\text{ref}})} \leq C_{\text{ref}} |\det B|^{-1/2} \|B\|_F \|\nabla v\|_{L^2(T)}.$$

Recall that $|\det B| \simeq |T|$ and $\|B\|_F \simeq h_T$. Therefore,

$$\|v_T - v_E\|_{L^2(T)} = |T|^{1/2} |v_T - v_E| \lesssim h_T \|\nabla v\|_{L^2(T)}.$$

3. step. For all $T, T' \in \mathcal{T}$ with $E := T \cap T' \in \mathcal{E}$, shape regularity and the triangle inequality yield that

$$\begin{aligned} \|v_T - v_{T'}\|_{L^2(T)} &= |T|^{1/2} |v_T - v_{T'}| \lesssim |T|^{1/2} |v_T - v_E| + |T'|^{1/2} |v_E - v_{T'}| \\ &= \|v_T - v_E\|_{L^2(T)} + \|v_{T'} - v_E\|_{L^2(T')}. \end{aligned}$$

v13:
22.11.2017

With Step 2, we thus see that

$$\|v_T - v_{T'}\|_{L^2(T)} \lesssim h_T \|\nabla v\|_{L^2(T)} + h_{T'} \|\nabla v\|_{L^2(T')} \lesssim h_T \|\nabla v\|_{L^2(T \cup T')},$$

where the hidden constant now depends on $C > 0$ from step 2 and from γ -shape regularity of \mathcal{T} .

4. step. If $T \cap T' \neq \emptyset$, there is a minimal $n \in \mathbb{N}$ and elements $T_0, \dots, T_n \in \mathcal{T}$ with $T_0 = T$, $T_j \cap T_{j-1} \in \mathcal{E}$ and $T_j \subseteq \Omega_T$ for all $j = 1, \dots, n$, and $T_n = T'$. Note that n is uniformly bounded in terms of the γ -shape regularity of \mathcal{T} . Iterating the argument from Step 3, we conclude (4.10) with $\bigcup_{j=0}^n T_j \subseteq \Omega_T$. The overall constant then depends on $C > 0$ and γ .

5. step. For each element $T'' \in \tilde{\Omega}_T$, the Poincaré inequality and (4.10) show

$$\begin{aligned} & \|v - v_{T''}\|_{L^2(T'')} + \|v - v_{E'}\|_{L^2(T'')} \\ & \lesssim \|v - v_{T''}\|_{L^2(T'')} + \|v_T - v_{T'}\|_{L^2(T'')} + \|v_T - v_{T''}\|_{L^2(T'')} + \|v_T - v_{E'}\|_{L^2(T'')} \\ & \simeq \|v - v_{T''}\|_{L^2(T'')} + \|v_T - v_{T'}\|_{L^2(T)} + \|v_T - v_{T''}\|_{L^2(T)} + \|v_T - v_{E'}\|_{L^2(T)} \\ & \lesssim h_{T''} \|\nabla v\|_{L^2(T'')} + h_T \|\nabla v\|_{L^2(\Omega_T)} \\ & \lesssim h_T \|\nabla v\|_{L^2(\Omega_T)}. \end{aligned}$$

Summing this estimate over all $T'' \in \tilde{\Omega}_T$, we obtain that

$$\|v - v_{T'}\|_{L^2(\Omega_T)} + \|v - v_{E'}\|_{L^2(\Omega_T)} \lesssim h_T \|\nabla v\|_{L^2(\Omega_T)},$$

where the hidden constants depends only on γ -shape regularity of \mathcal{T} . ■

Proof of Proposition 4.5 (H^1 -stability). For $z \in \mathcal{K}$, let $E_z \subset T_z \in \mathcal{T}$ and $h_z := \text{diam}(T_z)$. Note that $T_z \subseteq \Omega_T$ for $z \in T$. The trace inequality (3.56) yields that

$$\|v\|_{L^2(E_z)}^2 \lesssim h_z^{-1} (\|v\|_{L^2(T_z)}^2 + h_z \|v\|_{L^2(T_z)} \|\nabla v\|_{L^2(T_z)}) \lesssim h_z^{-1} (\|v\|_{L^2(T_z)}^2 + h_z^2 \|\nabla v\|_{L^2(T_z)}^2)$$

With this and Lemma 4.1, we see that

$$\begin{aligned} \left| \int_{E_z} \psi_z v \, ds \right| & \leq \|\psi_z\|_{L^\infty(E_z)} \|v\|_{L^1(E_z)} \lesssim |E_z|^{-1/2} \|v\|_{L^2(E_z)} \\ & \lesssim |E_z|^{-1/2} h_z^{-1/2} (\|v\|_{L^2(T_z)} + h_z \|\nabla v\|_{L^2(T_z)}). \end{aligned}$$

For any hat function, an inverse estimate shows

$$\|\nabla \zeta_z\|_{L^2(T)} \lesssim h_T^{-1} \|\zeta_z\|_{L^2(T)} \leq |T|^{1/2} h_T^{-1}.$$

Together with $|E_z| h_z \simeq |T_z| \simeq |T|$ and $h_z \simeq h_T$, we therefore obtain that, for all $v \in H^1(\Omega)$,

$$\|\nabla J_h v\|_{L^2(T)} \leq \sum_{z \in \mathcal{K} \cap T} \left| \int_{E_z} \psi_z v \, ds \right| \|\nabla \zeta_z\|_{L^2(T)} \lesssim \sum_{z \in \mathcal{K} \cap T} (h_z^{-1} \|v\|_{L^2(T_z)} + \|\nabla v\|_{L^2(T_z)}). \quad (4.12)$$

With the integral mean $v_T := (1/|T|) \int_T v \, dx$ and the projection property $J_h v_T = v_T$, we apply the last estimate for $w := v - v_T$ and see that

$$\|\nabla J_h v\|_{L^2(T)} = \|\nabla J_h (v - v_T)\|_{L^2(T)} \lesssim \sum_{z \in \mathcal{K} \cap T} (h_z^{-1} \|v - v_T\|_{L^2(T_z)} + \|\nabla v\|_{L^2(T_z)}).$$

According to the Poincaré inequality and Lemma 4.6, it holds that for all $z \in \mathcal{K} \cap T$,

$$\|v - v_T\|_{L^2(T_z)} \leq \|v - v_{T_z}\|_{L^2(T_z)} + \|v_{T_z} - v_T\|_{L^2(T_z)} \lesssim h_z \|\nabla v\|_{L^2(\Omega_T)}. \quad (4.13)$$

Combining the last two estimates, we thus conclude $\|\nabla J_h v\|_{L^2(T)} \lesssim \|\nabla v\|_{L^2(\Omega_T)}$. \blacksquare

Proof of Proposition 4.5 (approximation property). We adopt the notation from the proof of local H^1 -stability. Arguing as for (4.12), we see that

$$\|J_h v\|_{L^2(T)} \leq \sum_{z \in \mathcal{K} \cap T} \left| \int_{E_z} \psi_z v ds \right| \|\zeta_z\|_{L^2(T)} \lesssim \sum_{z \in \mathcal{K} \cap T} (\|v\|_{L^2(T_z)} + h_z \|\nabla v\|_{L^2(T_z)}). \quad (4.14)$$

With the integral mean $v_T := (1/|T|) \int_T v dx$ and the projection property $J_h v_T = v_T$, we apply the last estimate for $w := v - v_T$ and see that

$$\begin{aligned} \|v - J_h v\|_{L^2(T)} &= \|(v - v_T) - J_h(v - v_T)\|_{L^2(T)} \\ &\leq \|v - v_T\|_{L^2(T)} + \|J_h(v - v_T)\|_{L^2(T)} \\ &\lesssim h_T \|\nabla v\|_{L^2(T)} + \sum_{z \in \mathcal{K} \cap T} (\|v - v_T\|_{L^2(T_z)} + h_z \|\nabla v\|_{L^2(T_z)}) \end{aligned}$$

Finally, we employ (4.13) and $h_z \simeq h_T$ to conclude $\|v - J_h v\|_{L^2(T)} \lesssim h_T \|\nabla v\|_{L^2(\Omega_T)}$. \blacksquare

The following theorem concludes the main properties of the Scott-Zhang projection:

Theorem 4.7. *The Scott-Zhang projection $J_h : H^1(\Omega) \rightarrow \mathcal{S}^1(\mathcal{T})$ has the following properties (i)–(vii):*

(i) J_h is linear and continuous with respect to the H^1 -norm, i.e.,

$$\|J_h v\|_{H^1(\Omega)} \leq C(1 + \text{diam}(\Omega)) \|v\|_{H^1(\Omega)} \quad \text{for all } v \in H^1(\Omega). \quad (4.15)$$

(ii) J_h is a projection onto $\mathcal{S}^1(\mathcal{T})$, i.e.,

$$J_h v_h = v_h \quad \text{for all } v_h \in \mathcal{S}^1(\mathcal{T}). \quad (4.16)$$

(iii) J_h preserves discrete boundary data, i.e., for $\omega \in \{\Gamma_D, \Gamma\}$ it holds that

$$(J_h v)|_\omega = v|_\omega \quad \text{for all } v \in H^1(\Omega) \text{ with } v|_\omega \in \mathcal{S}^1(\mathcal{T}|_\omega). \quad (4.17)$$

(iv) J_h is locally H^1 -stable, i.e.,

$$\|\nabla J_h v\|_{L^2(T)} \leq C \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H^1(\Omega) \text{ and } T \in \mathcal{T}. \quad (4.18)$$

(v) J_h has a local first-order approximation property, i.e.,

$$\|(1 - J_h)v\|_{L^2(T)} \leq Ch_T \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H^1(\Omega) \text{ and } T \in \mathcal{T}. \quad (4.19)$$

(vi) P_h is quasi-optimal in the sense of the Céa lemma, i.e.,

$$\|(1 - J_h)v\|_{H^1(\Omega)} \leq C(1 + \text{diam}(\Omega)) \min_{v_h \in \mathcal{S}^1(\mathcal{T})} \|v - v_h\|_{H^1(\Omega)} \quad \text{for all } v \in H^1(\Omega). \quad (4.20)$$

(vii) For all $\alpha \in \mathbb{R}$, J_h is quasi-optimal in the sense of

$$\|h^\alpha \nabla(1 - J_h)v\|_{L^2(\Omega)} \leq C \min_{v_h \in \mathcal{S}^1(\mathcal{T})} \|h^\alpha \nabla(v - v_h)\|_{L^2(\Omega)}. \quad (4.21)$$

The constant $C > 0$ in (i)–(vii) depends only on γ -shape regularity of \mathcal{T} .

Proof. (ii)–(v) have already been shown, and (i) is a direct consequence of (vi) and the triangle inequality. (vii) Let $v_h \in \mathcal{S}^1(\mathcal{T})$. With the projection property of J_h and (iv), we see that, for all $T \in \mathcal{T}$,

$$\|\nabla(1 - J_h)v\|_{L^2(T)} = \|\nabla(1 - J_h)(v - v_h)\|_{L^2(T)} \lesssim \|\nabla(v - v_h)\|_{L^2(\Omega_T)}.$$

With γ -shape regularity and hence $h_T \simeq h_{T'}$ for all $T' \subseteq \Omega_T$, we infer

$$\|h^\alpha \nabla(1 - J_h)v\|_{L^2(T)} \lesssim \|h^\alpha \nabla(v - v_h)\|_{L^2(\Omega_T)}.$$

Using the γ -shape regularity again, this results in

$$\begin{aligned} \|h^\alpha \nabla(1 - J_h)v\|_{L^2(\Omega)}^2 &= \sum_{T \in \mathcal{T}} \|h^\alpha \nabla(1 - J_h)v\|_{L^2(T)}^2 \lesssim \sum_{T \in \mathcal{T}} \|h^\alpha \nabla(v - v_h)\|_{L^2(\Omega_T)}^2 \\ &\lesssim \|h^\alpha \nabla(v - v_h)\|_{L^2(\Omega)}^2. \end{aligned}$$

This proves (vii) with an infimum on the right-hand side. Due to finite dimension, this infimum is, in fact, attained. To prove (vi), it remains to estimate the L^2 -part and use $\alpha = 0$ in (vii). With the projection property of J_h and (v), shape regularity yields that

$$\begin{aligned} \|(1 - J_h)v\|_{L^2(\Omega)}^2 &= \sum_{T \in \mathcal{T}} \|(1 - J_h)(v - v_h)\|_{L^2(T)}^2 \lesssim \sum_{T \in \mathcal{T}} h_T^2 \|\nabla(v - v_h)\|_{L^2(\Omega_T)}^2 \\ &\lesssim \text{diam}(\Omega)^2 \|\nabla(v - v_h)\|_{L^2(\Omega)}^2. \end{aligned}$$

Altogether, we thus see that

$$\|(1 - J_h)v\|_{H^1(\Omega)}^2 \lesssim (1 + \text{diam}(\Omega)^2) \|\nabla(v - v_h)\|_{L^2(\Omega)}^2 \lesssim (1 + \text{diam}(\Omega))^2 \|v - v_h\|_{H^1(\Omega)}^2.$$

This concludes the proof of (vi). ■

Remark. Theorem 4.7 holds for any dimension $d \geq 2$ and for any fixed polynomial degree $p \geq 1$. □

One drawback of the Scott-Zhang projection is that it is not positivity conserving, i.e., $v \geq 0$ does not necessarily imply that $J_h v \geq 0$.

Exercise 46. Suppose that \mathcal{T} is a regular triangulation of $\Omega := [0, 1]^2$ into 2 triangles. Find an example of a function $v \in H^1(\Omega)$ with $v \geq 0$ such that there exists some $x \in \Omega$ with $J_h v < 0$.
Hint. Compute the function $\hat{\psi} \in \mathcal{P}^1(0, 1)$ from Lemma 4.1 explicitly. \square

Exercise 47. Extend the approach of Exercise 45 and construct an operator $P_h : L^2(\Omega) \rightarrow \mathcal{S}_D^1(\mathcal{T})$ with the following properties:

(i) $P_h : L^2(\Omega) \rightarrow \mathcal{S}_D^1(\mathcal{T})$ is a well-defined linear projection,

$$P_h v_h = v_h \quad \text{for all } v_h \in \mathcal{S}_D^1(\mathcal{T}).$$

(ii) P_h is locally L^2 -stable, i.e., for all $T \in \mathcal{T}$, it holds that

$$\|(1 - P_h)v\|_{L^2(T)} \leq C \|v\|_{L^2(\Omega_T)} \quad \text{for all } v \in L^2(\Omega).$$

(iii) P_h is locally H_D^1 -stable, i.e., for all $T \in \mathcal{T}$, it holds that

$$\|\nabla(1 - P_h)v\|_{L^2(T)} \leq C \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_D^1(\Omega).$$

(iv) P_h has a local first-order approximation property

$$\|(1 - P_h)v\|_{L^2(T)} \leq Ch_T \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_D^1(\Omega).$$

(v) $P_h : L^2(\Omega) \rightarrow L^2(\Omega)$ as well as $P_h : H_D^1(\Omega) \rightarrow H_D^1(\Omega)$ are bounded linear operators.

(vi) J_h is quasi-optimal in the sense of the Céa lemma, i.e.,

$$\|(1 - P_h)v\|_{H^1(\Omega)} \leq C \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|v - v_h\|_{H^1(\Omega)} \quad \text{for all } v \in H_D^1(\Omega).$$

(vii) For all $\alpha \in \mathbb{R}$, P_h is quasi-optimal in the sense of

$$\|h^\alpha(1 - P_h)v\|_{L^2(\Omega)} \leq C \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|h^\alpha(v - v_h)\|_{L^2(\Omega)} \quad \text{for all } v \in L^2(\Omega).$$

(viii) For all $\alpha \in \mathbb{R}$, P_h is quasi-optimal in the sense of

$$\|h^\alpha \nabla(1 - P_h)v\|_{L^2(\Omega)} \leq C \min_{v_h \in \mathcal{S}_D^1(\mathcal{T})} \|h^\alpha \nabla(v - v_h)\|_{L^2(\Omega)} \quad \text{for all } v \in H_D^1(\Omega).$$

The constant $C > 0$ in (i)–(viii) depends only on γ -shape regularity of \mathcal{T} . **Hint.** Let $\mathcal{K}_F := \mathcal{K} \setminus \Gamma_D$ denote the free nodes, where possibly $\mathcal{K}_F = \mathcal{K}$ for $\Gamma_D = \emptyset$. Then, P_h can be chosen as

$$P_h v = \sum_{z \in \mathcal{K}_F} \left(\int_{T_z} v \psi_z dx \right) \zeta_z$$

with appropriate $T_z \in \mathcal{T}$ and $\psi_z \in \mathcal{P}^1(T_z)$. □

Definition. The Scott-Zhang projection is just a special example of a Clément-type quasi-interpolation operator: We say that an operator $J_h : H_D^1(\Omega) \rightarrow \mathcal{S}_D^1(\mathcal{T})$ is a **Clément-type quasi-interpolation operator** if, for all $v \in H_D^1(\Omega)$ and all $T \in \mathcal{T}$, it holds that

v14:
24.11.2017

- it is locally H^1 -stable

$$\|\nabla(1 - J_h)v\|_{L^2(T)} \leq C \|\nabla v\|_{L^2(\Omega_T)}, \quad (4.22)$$

- and has a local first-order approximation property

$$\|(1 - J_h)v\|_{L^2(T)} \leq Ch_T \|\nabla v\|_{L^2(\Omega_T)}. \quad (4.23)$$

The constant $C > 0$ may only depend on γ -shape regularity of \mathcal{T} (and possibly the shapes of possible patches in \mathcal{T}).

For the a posteriori error analysis, we shall need the following simple consequence.

Lemma 4.8. *Suppose that $J_h : H_D^1(\Omega) \rightarrow \mathcal{S}_D^1(\mathcal{T})$ is a Clément-type operator, i.e., (4.22)–(4.23) hold. Let $T \in \mathcal{T}$ and $E \in \mathcal{E}_T$. Then, it holds that*

$$\|(1 - J_h)v\|_{L^2(E)} \leq Ch_E^{1/2} \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_D^1(\Omega). \quad (4.24)$$

The constant $C > 0$ depends only on γ -shape regularity of \mathcal{T} .

Proof. We apply the trace inequality

$$\|w\|_{L^2(E)}^2 \lesssim h_T^{-1} (\|w\|_{L^2(T)}^2 + h_T \|w\|_{L^2(T)} \|\nabla w\|_{L^2(T)})$$

for $w := (1 - J_h)v \in H_D^1(\Omega)$. With the Clément properties (4.22)–(4.23), this yields that

$$\|(1 - J_h)v\|_{L^2(E)}^2 \lesssim h_T \|\nabla v\|_{L^2(\Omega_T)}^2.$$

Shape regularity and hence $h_T \simeq h_E$ concludes the proof. \blacksquare

The following example is one further *classical* example of a Clément-type operator. The analysis will be left to the reader, but requires the following simple observation:

Exercise 48. Use a scaling argument to show that

$$C^{-1}h_T \|\nabla v_h\|_{L^\infty(T)} \leq \|\nabla v_h\|_{L^2(T)} \leq \frac{h_T}{\sqrt{2}} \|\nabla v_h\|_{L^\infty(T)} \quad \text{for all } v_h \in \mathcal{P}^m(T),$$

where the constant $C > 0$ only depends on $\sigma(\mathcal{T})$ and the polynomial degree $m \in \mathbb{N}_0$. \square

Exercise 49. Let $\mathcal{K}_F := \mathcal{K} \setminus \overline{\Gamma}_D$ denote the free nodes (where possibly $\mathcal{K}_F = \mathcal{K}$ if $\Gamma_D = \emptyset$). Define

$$J_h v := \sum_{z \in \mathcal{K}_F} v_z \zeta_z \quad \text{with} \quad v_z := \frac{1}{|\Omega_z|} \int_{\Omega_z} v \, dx, \quad (4.25)$$

where $\Omega_z \subseteq \overline{\Omega}$ denotes the patch of a node $z \in \mathcal{K}$. Prove that J_h satisfies the following properties:

- $J_h : L^2(\Omega) \rightarrow \mathcal{S}_D^1(\mathcal{T})$ is a well-defined linear operator.

(ii) J_h is locally L^2 -stable, i.e., for all $T \in \mathcal{T}$, it holds that

$$\|(1 - J_h)v\|_{L^2(T)} \leq C \|v\|_{L^2(\Omega_T)} \quad \text{for all } v \in L^2(\Omega).$$

(iii) J_h is locally H_D^1 -stable, i.e., for all $T \in \mathcal{T}$, it holds that

$$\|\nabla(1 - J_h)v\|_{L^2(T)} \leq C \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_D^1(\Omega).$$

(iv) J_h has a local first-order approximation property

$$\|(1 - J_h)v\|_{L^2(T)} \leq Ch_T \|\nabla v\|_{L^2(\Omega_T)} \quad \text{for all } v \in H_D^1(\Omega).$$

(v) $J_h : L^2(\Omega) \rightarrow L^2(\Omega)$ as well as $J_h : H_D^1(\Omega) \rightarrow H_D^1(\Omega)$ are bounded linear operators.

(vi) J_h is positivity preserving, i.e., $J_h v \geq 0$ for all $v \in L^2(\Omega)$ with $v \geq 0$.

(vii) With $\Pi_h : L^2(\Omega) \rightarrow \mathcal{P}^0(\mathcal{T})$ the L^2 -orthogonal projection onto $\mathcal{P}^0(\mathcal{T})$, it holds that $J_h \Pi_h = J_h$.

The constant $C > 0$ depends only on γ -shape regularity of \mathcal{T} . □

Exercise 50. Find a counter example which shows that the operator J_h from Exercise 49 is no projection, i.e., it holds that $J_h v_h \neq v_h$ for some $v_h \in \mathcal{S}_D^1(\mathcal{T})$. □

4.3 Residual-Based Error Estimator (12.12.2014)

Residual-based a posteriori error estimates follow a general strategy. Recall that the weak solution $u \in H_D^1(\Omega)$ of (4.1) solves the variational form

$$(\nabla u ; \nabla v)_{L^2(\Omega)} = (f ; v)_{L^2(\Omega)} + (\phi ; v)_{L^2(\Gamma_N)} \quad \text{for all } v \in H_D^1(\Omega). \quad (4.26)$$

For an approximation $u_h \in \mathcal{S}_D^1(\mathcal{T})$ it is thus natural to define the **residual** $R_h \in H_D^1(\Omega)^*$ by

$$R_h(v) := (f ; v)_{L^2(\Omega)} + (\phi ; v)_{L^2(\Gamma_N)} - (\nabla u_h ; \nabla v)_{L^2(\Omega)}, \quad (4.27)$$

i.e., $R_h = 0$ if and only if $u_h = u$. Let $\|w\| := \|\nabla w\|_{L^2(\Omega)}$ denote the energy norm on $H_D^1(\Omega)$ and

$$\|\Phi\|_* := \sup_{w \in H_D^1(\Omega) \setminus \{0\}} \frac{\Phi(w)}{\|w\|}$$

the induced operator norm on $H_D^1(\Omega)^*$, where we stress that both are equivalent norms on $H_D^1(\Omega)$ and its dual space, respectively. Then, the Riesz theorem and $R_h(v) = (\nabla(u - u_h) ; \nabla v)_{L^2(\Omega)}$ yield

$$\|R_h\|_* = \|u - u_h\|.$$

To derive a reliable error estimator η , we thus need to prove an estimate of the type

$$R_h(v) \leq \widetilde{C}_{\text{rel}} \eta \|v\| \quad \text{for all } v \in H_D^1(\Omega). \quad (4.28)$$

To derive an efficient error estimator η , we need to show

$$R_h(v) \geq \widetilde{C}_{\text{eff}} \eta \|v\| \quad \text{for some } v \in H_D^1(\Omega) \setminus \{0\}, \quad (4.29)$$

where this $v \in H_D^1(\Omega)$ has to be constructed appropriately.

Exercise 51. Prove that reliability (4.2) of an error estimator η is, in fact, equivalent to (4.28). Prove that efficiency (4.3) of η holds if and only if (4.29) holds. \square

So far, our observations did not use that we are dealing with Galerkin schemes. We stress that the Galerkin orthogonality reads

$$R_h(v_h) = 0 \quad \text{for all } v_h \in \mathcal{S}_D^1(\mathcal{T}) \quad (4.30)$$

with respect to the residual R_h . To provide a reliable (and residual-based) error estimator η , we will use some Clément-type operator $J_h : H^1(\Omega) \rightarrow \mathcal{S}_D^1(\Omega)$ in connection with the Galerkin orthogonality (4.30).

Before introducing a first a posteriori error estimator, we introduce the following notational conventions. We define the \mathcal{T} -piecewise resp. \mathcal{E} -piecewise constant mesh-width functions

$$h_{\mathcal{T}}|_T := h_T \quad \text{and} \quad h_{\mathcal{E}}|_T := h_E$$

for elements $T \in \mathcal{T}$ and edges $E \in \mathcal{E}$, respectively. Moreover, we write

$$\|h_{\mathcal{E}}^{1/2} \psi\|_{L^2(\mathcal{E}_*)} := \left(\sum_{E \in \mathcal{E}_*} h_E \|\psi\|_{L^2(E)}^2 \right)^{1/2}$$

for any set $\mathcal{E}_* \subseteq \mathcal{E}$ of edges and any function ψ which belongs to $L^2(E)$ for all $E \in \mathcal{E}_*$. Recall that \mathcal{E}_D and \mathcal{E}_N denote the Dirichlet and Neumann edges of \mathcal{T} , respectively. Moreover, let \mathcal{E}_Ω denote the set of all **interior edges**, i.e., for $E \in \mathcal{E}_\Omega$, there are unique elements $T_E^+, T_E^- \in \mathcal{T}$ with $E = T_E^+ \cap T_E^-$. Finally, for $E \in \mathcal{E}_\Omega$, we define the **jump of the normal derivative** by

$$[[\partial_n u_h]]_E := \frac{\partial u_h}{\partial n_E^+} + \frac{\partial u_h}{\partial n_E^-} \in \mathbb{R}, \quad (4.31)$$

where n_E^\pm denote the outer normal vectors of the elements T_E^\pm on the edge E . Note that $n_E^+ = -n_E^-$ so that the sum in the definition is, in fact, a difference.

Theorem 4.9. *The error estimator*

$$\eta := \left(\|h_{\mathcal{T}} f\|_{L^2(\Omega)}^2 + \|h_{\mathcal{E}}^{1/2} [[\partial_n u_h]]\|_{L^2(\mathcal{E}_\Omega)}^2 + \|h_{\mathcal{E}}^{1/2} (\phi - \partial_n u_h)\|_{L^2(\mathcal{E}_N)}^2 \right)^{1/2} \quad (4.32)$$

satisfies the reliability estimate

$$\|u - u_h\|_{H^1(\Omega)} \leq C \eta, \quad (4.33)$$

where the constant $C > 0$ depends only on γ -shape regularity of \mathcal{T} .

Proof. For all $w \in H_D^1(\Omega)$, elementwise integration by parts proves

$$\begin{aligned}
 R_h(w) &= (f; w)_{L^2(\Omega)} + (\phi; w)_{L^2(\Gamma_N)} - \sum_{T \in \mathcal{T}} (\nabla u_h; \nabla w)_{L^2(T)} \\
 &= (f; w)_{L^2(\Omega)} + \sum_{E \in \mathcal{E}_N} (\phi; w)_{L^2(E)} - \sum_{T \in \mathcal{T}} (\partial_n u_h; w)_{L^2(\partial T)} \\
 &= \sum_{T \in \mathcal{T}} (f; w)_{L^2(T)} + \sum_{E \in \mathcal{E}_N} (\phi - \partial_n u_h; w)_{L^2(E)} - \sum_{E \in \mathcal{E}_\Omega} ([[\partial_n u_h]]; w)_{L^2(E)} \\
 &\leq \sum_{T \in \mathcal{T}} \|f\|_{L^2(T)} \|w\|_{L^2(T)} + \sum_{E \in \mathcal{E}_N} \|\phi - \partial_n u_h\|_{L^2(E)} \|w\|_{L^2(E)} + \sum_{E \in \mathcal{E}_\Omega} \|[[\partial_n u_h]]\|_{L^2(E)} \|w\|_{L^2(E)}.
 \end{aligned}$$

For arbitrary $v \in H_D^1(\Omega)$, we now choose $w = v - J_h v$ and note that $R_h(v) = R_h(w)$ according to the Galerkin orthogonality. Then, we estimate the three sums separately. The approximation property of the Clément-type operator J_h and Lemma 4.3 imply

$$\begin{aligned}
 \sum_{T \in \mathcal{T}} \|f\|_{L^2(T)} \|v - J_h v\|_{L^2(T)} &\lesssim \left(\sum_{T \in \mathcal{T}} \|h_{\mathcal{T}} f\|_{L^2(T)}^2 \right)^{1/2} \left(\sum_{T \in \mathcal{T}} \|\nabla v\|_{L^2(\Omega_T)}^2 \right)^{1/2} \\
 &\lesssim \left(\sum_{T \in \mathcal{T}} \|h_{\mathcal{T}} f\|_{L^2(T)}^2 \right)^{1/2} \left(\sum_{T \in \mathcal{T}} \|\nabla v\|_{L^2(T)}^2 \right)^{1/2} \\
 &= \|h_{\mathcal{T}} f\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)}.
 \end{aligned}$$

For each edge $E \in \mathcal{E}$, we choose an arbitrary element $T_E \in \mathcal{T}$ with $E \in \mathcal{E}_{T_E}$. Let $\mathcal{E}_* \subset \mathcal{E}$ and $\psi \in L^2(E)$ for all $E \in \mathcal{E}_*$. Recall that $\|(1 - J_h)v\|_{L^2(E)} \lesssim h_E^{1/2} \|\nabla v\|_{L^2(\Omega_{T_E})}$. Therefore, the same arguments as before prove

$$\begin{aligned}
 \sum_{E \in \mathcal{E}_*} \|\psi\|_{L^2(E)} \|v - J_h v\|_{L^2(E)} &\lesssim \left(\sum_{E \in \mathcal{E}_*} \|h_{\mathcal{E}}^{1/2} \psi\|_{L^2(E)}^2 \right)^{1/2} \left(\sum_{E \in \mathcal{E}_*} \|\nabla v\|_{L^2(\Omega_{T_E})}^2 \right)^{1/2} \\
 &\lesssim \|h_{\mathcal{E}}^{1/2} \psi\|_{L^2(\mathcal{E}_*)} \|\nabla v\|_{L^2(\Omega)},
 \end{aligned}$$

where we note that an element $T \in \mathcal{T}$ may satisfy $T = T_E$ for at most three edges. Altogether, we now see

$$\begin{aligned}
 R_h(v) &\lesssim \|\nabla v\|_{L^2(\Omega)} \left(\|h_{\mathcal{T}} f\|_{L^2(\Omega)} + \|h_{\mathcal{E}}^{1/2} [[\partial_n u_h]]\|_{L^2(\mathcal{E}_\Omega)} + \|h_{\mathcal{E}}^{1/2} (\phi - \partial_n u_h)\|_{L^2(\mathcal{E}_N)} \right) \\
 &\leq \sqrt{3} \|\nabla v\|_{L^2(\Omega)} \eta.
 \end{aligned}$$

The hidden constant C depends only (on the Clément operator J_h and) on γ -shape regularity of \mathcal{T} . \blacksquare

Remark. Note that we have used $u_h \in \mathcal{S}^1(\mathcal{T})$ in the sense that the elementwise Laplacian satisfies $\Delta u_h|_T = 0$ for all $T \in \mathcal{T}$. For general \mathcal{T} -piecewise polynomials, the same proof applies with $\|h_{\mathcal{T}} f\|_{L^2(\Omega)}$ replaced by $\|h_{\mathcal{T}}(f + \Delta u_h)\|_{L^2(\Omega)}$. \square

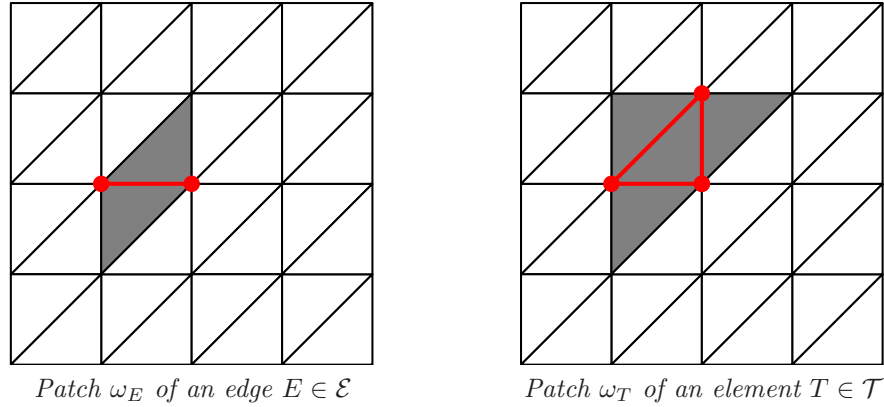


FIGURE 4.2. To prove the efficiency estimate, it suffices to consider smaller patches $\omega_E \subseteq \Omega_E$ and $\omega_T \subseteq \Omega_T$, for edges $E \in \mathcal{E}$ and elements $T \in \mathcal{T}$, respectively. For comparison with the larger patches Ω_E and Ω_T , the reader may consider Figure 4.1 on page 72.

Exercise 52. We consider the mixed boundary value problem

$$\begin{aligned} -\Delta u &= f && \text{in } \Omega, \\ u &= u_D && \text{on } \Gamma_D, \\ \partial_n u &= \phi && \text{on } \Gamma_N. \end{aligned}$$

with inhomogeneous Dirichlet data $u_D \in H^{1/2}(\Gamma_D)$. Let $u \in H^1(\Omega)$ denote the weak solution and $u_h \in \mathcal{S}^1(\mathcal{T}_h)$ the P1-FEM solution for discrete Dirichlet data $u_{Dh} := \hat{u}_{Dh}|_{\Gamma_D}$ with $\hat{u}_{Dh} \in \mathcal{S}^1(\mathcal{T}_h)$. Use the additional problem

$$\begin{aligned} -\Delta w &= 0 && \text{in } \Omega, \\ w &= u_D - u_{Dh} && \text{on } \Gamma_D, \\ \partial_n w &= 0 && \text{on } \Gamma_N. \end{aligned}$$

with weak solution $w \in H^1(\Omega)$ to derive a reliable error estimator for $\|u - u_h\|_{H^1(\Omega)}$.

Hint. Prove that $\|w\|_{H^1(\Omega)} \simeq \|u_D - u_{Dh}\|_{H^{1/2}(\Gamma_D)}$, where the right-hand side is already an a posteriori term. Then, consider the residual $\tilde{R}_h \in H_D^1(\Omega)^*$ corresponding to the function $(u - u_h) - w \in H_D^1(\Omega)$. □

Next, we prove the efficiency of the residual-based error estimator η from (4.32) — at least up to terms of higher order. The efficiency estimate even holds locally with refined patches ω_E and ω_T shown in Figure 4.2: For an interior edge $E \in \mathcal{E}_\Omega$, let $T_E^+, T_E^- \in \mathcal{T}$ be the unique elements with $E = T_E^+ \cap T_E^-$. For a boundary edge $E \in \mathcal{E}_\Gamma$, there is a unique element $T_E \in \mathcal{T}$ with $E \in \mathcal{E}_{T_E}$. We define the refined patch of an edge $E \in \mathcal{E}$ by

$$\omega_E := \begin{cases} T_E^+ \cup T_E^- & \text{for } E \in \mathcal{E}_\Omega, \\ T_E & \text{for } E \in \mathcal{E}_\Gamma. \end{cases} \quad (4.34)$$

Moreover, we define the refined patch of an element $T \in \mathcal{T}$ by

$$\omega_T := \bigcup \{ \omega_E \mid E \in \mathcal{E}_T \}. \quad (4.35)$$

Note that $\omega_E \subseteq \Omega_E$ and $\omega_T \subseteq \Omega_T$, so that Lemma 4.3 and Lemma 4.4 even hold for the refined patches.

Usually, one is interested in error estimators which are localized with respect to the elements or the edges of \mathcal{T} , respectively. For instance, one considers the **element-based residual error estimator**

$$\eta_T := \left(\sum_{T \in \mathcal{T}} \eta_T^2 \right)^{1/2}, \quad (4.36)$$

where

$$\eta_T = \left(h_T^2 \|f\|_{L^2(T)}^2 + h_T \|[\![\partial_n u_h]\!] \|_{L^2(\partial T \cap \Omega)}^2 + h_T \|\phi - \partial_n u_h\|_{L^2(\partial T \cap \Gamma_N)}^2 \right)^{1/2} \quad (4.37)$$

or the **edge-based residual error estimator**

$$\eta_{\mathcal{E}} := \left(\sum_{E \in \mathcal{E}} \eta_E^2 \right)^{1/2}, \quad (4.38)$$

where

$$\eta_E = \begin{cases} \left(h_E^2 \|f\|_{L^2(\omega_E)}^2 + h_E \|[\![\partial_n u_h]\!] \|_{L^2(E)}^2 \right)^{1/2} & \text{for } E \in \mathcal{E}_\Omega, \\ \left(h_E^2 \|f\|_{L^2(\omega_E)}^2 + h_E \|\phi - \partial_n u_h\|_{L^2(E)}^2 \right)^{1/2} & \text{for } E \in \mathcal{E}_N, \\ 0 & \text{for } E \in \mathcal{E}_D. \end{cases} \quad (4.39)$$

Alternatively, one could also define

$$\eta_{\mathcal{T} \cup \mathcal{E}} := \left(\sum_{T \in \mathcal{T}} \eta_T^2 + \sum_{E \in \mathcal{E}} \eta_E^2 \right)^{1/2}, \quad (4.40)$$

where

$$\eta_T = h_T \|f\|_{L^2(T)}, \quad (4.41a)$$

$$\eta_E = \begin{cases} h_E^{1/2} \|[\![\partial_n u_h]\!] \|_{L^2(E)} & \text{for } E \in \mathcal{E}_\Omega, \\ h_E^{1/2} \|\phi - \partial_n u_h\|_{L^2(E)} & \text{for } E \in \mathcal{E}_N, \\ 0 & \text{for } E \in \mathcal{E}_D. \end{cases} \quad (4.41b)$$

Note that $\eta_{\mathcal{T}}$ as well as $\eta_{\mathcal{E}}$ are equivalent to the error estimator η from (4.32): There holds

$$\eta = \eta_{\mathcal{T} \cup \mathcal{E}} \leq \eta_{\mathcal{T}} \leq \sqrt{2} \sigma(\mathcal{T}_h)^{1/2} \eta \quad \text{as well as} \quad \sigma(\mathcal{T})^{-1} \eta \leq \eta_{\mathcal{E}} \leq \sqrt{3} \eta,$$

since $\eta_{\mathcal{T}}$ adds the contributions of interior edges twice and $h_E \leq h_T \leq \sigma(\mathcal{T}_h) h_E$ for each edge $E \in \mathcal{E}_T$, whereas $\eta_{\mathcal{E}}$ adds the element contribution at most three times. The local quantities η_T and η_E can be used to steer an adaptive mesh-refining algorithm. They are therefore called **refinement indicators**. We are going to discuss adaptive mesh-refinement below.

Theorem 4.10. We define $f_{\mathcal{T}} \in \mathcal{P}^0(\mathcal{T})$ by $f_{\mathcal{T}}|_T := |T|^{-1} \int_T f dx$ and $\phi_{\mathcal{E}} \in \mathcal{P}^0(\mathcal{E}_N)$ by $\phi_{\mathcal{E}}|_E := h_E^{-1} \int_E \phi ds$. For each element $T \in \mathcal{T}$, the refinement indicator η_T from (4.37) satisfies

$$\eta_T \leq C \left(\|\nabla(u - u_h)\|_{L^2(\omega_T)}^2 + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(\omega_T)}^2 + \|h_{\mathcal{E}}^{1/2}(\phi - \phi_{\mathcal{E}})\|_{L^2(\partial T \cap \Gamma_N)}^2 \right)^{1/2}. \quad (4.42)$$

Moreover, the error estimator η from (4.32) is efficient in the sense that

$$\eta \leq C \left(\|u - u_h\|_{H^1(\Omega)} + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(\Omega)} + \|h_{\mathcal{E}}^{1/2}(\phi - \phi_{\mathcal{E}})\|_{L^2(\Gamma_N)} \right). \quad (4.43)$$

The constant $C > 0$ only depends on the shape regularity constant $\sigma(\mathcal{T})$.

Remark. For $f \in H^1(\mathcal{T})$ holds $\|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(\Omega)} = \mathcal{O}(h^2)$. For $\phi \in C^1(\mathcal{E}_N)$ holds $\|h_{\mathcal{E}}^{1/2}(\phi - \phi_{\mathcal{E}})\|_{L^2(\Gamma_N)} = \mathcal{O}(h^{3/2})$. Even for $u \in H^2(\Omega)$, the error as well as the error estimator η only satisfy $\|u - u_h\|_{H^1(\Omega)} = \mathcal{O}(h) = \eta$. Therefore, the two terms on the right-hand side are of higher order. \square

Proof of Theorem 4.10. 1. step. Estimate (4.43) is a consequence of (4.42) since

$$\eta \leq \eta_{\mathcal{T}} = \left(\sum_{T \in \mathcal{T}} \eta_T^2 \right)^{1/2} \leq 2C \left(\|u - u_h\|_{H^1(\Omega)} + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(\Omega)} + \|h_{\mathcal{E}}^{1/2}(\phi - \phi_{\mathcal{E}})\|_{L^2(\Gamma_N)} \right).$$

Here, the factor $2 = 4^{1/2}$ appears since each element $T \in \mathcal{T}$ belongs at most to four patches $\omega_{T'}$.

The proof of (4.42) is split into three steps, where we consider each of the three contributions of η_T separately.

2. step. There holds

$$\|h_{\mathcal{T}}f\|_{L^2(T)} \leq C \left(\|\nabla(u - u_h)\|_{L^2(T)} + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(T)} \right): \quad (4.44)$$

For $T \in \mathcal{T}$, we define the **element bubble function**

$$b_T := \prod_{z \in \mathcal{K}_T} \zeta_z \in H_0^1(T) \cap \mathcal{P}^3(T)$$

as product of all three hat functions. It is essential to observe the following estimate

$$\|f_{\mathcal{T}}b_T\|_{L^2(T)} \leq \|f_{\mathcal{T}}b_T^{1/2}\|_{L^2(T)} \leq \|f_{\mathcal{T}}\|_{L^2(T)} \leq C_{\text{ref}} \|f_{\mathcal{T}}b_T\|_{L^2(T)}, \quad (4.45)$$

where the existence of an independent constant $C_{\text{ref}} > 0$ follows from a scaling argument. We stress, however, that $\|f_{\mathcal{T}}b_T\|_{L^2(T)}$ and hence C_{ref} — since $f_{\mathcal{T}}$ is constant on T — can explicitly be computed. In the following, the main idea is to use integration by parts for $v := f_{\mathcal{T}}b_T \in H_0^1(T)$ to show

$$\begin{aligned} C_{\text{ref}}^{-2} \|f_{\mathcal{T}}\|_{L^2(T)}^2 &\leq \|f_{\mathcal{T}}b_T^{1/2}\|_{L^2(T)}^2 = (f_{\mathcal{T}}; v)_{L^2(T)} = (f_{\mathcal{T}} - f; v)_{L^2(T)} + (f - \Delta u_h; v)_{L^2(T)} \\ &= (f_{\mathcal{T}} - f; v)_{L^2(T)} + (\nabla(u - u_h); \nabla v)_{L^2(T)}. \end{aligned}$$

Now, we estimate each of the two scalar products on the right-hand side by use of the Cauchy inequality. Together with $v = f_{\mathcal{T}} b_T \in \mathcal{P}^3(\mathcal{T})$ we observe

$$(f_{\mathcal{T}} - f; v)_{L^2(T)} \leq \|f_{\mathcal{T}} - f\|_{L^2(T)} \|f_{\mathcal{T}} b_T\|_{L^2(T)} \leq \|f_{\mathcal{T}} - f\|_{L^2(T)} \|f_{\mathcal{T}}\|_{L^2(T)}$$

as well as

$$\begin{aligned} (\nabla(u - u_h); \nabla v)_{L^2(T)} &\leq \|\nabla(u - u_h)\|_{L^2(T)} \|\nabla(f_{\mathcal{T}} b_T)\|_{L^2(T)} \\ &\leq C_{\text{inv}} h_T^{-1} \|\nabla(u - u_h)\|_{L^2(T)} \|f_{\mathcal{T}} b_T\|_{L^2(T)} \\ &\leq C_{\text{inv}} h_T^{-1} \|\nabla(u - u_h)\|_{L^2(T)} \|f_{\mathcal{T}}\|_{L^2(T)}. \end{aligned}$$

Altogether, we see

$$h_T \|f_{\mathcal{T}}\|_{L^2(T)} \leq C_{\text{ref}}^2 (h_T \|f_{\mathcal{T}} - f\|_{L^2(T)} + C_{\text{inv}} \|\nabla(u - u_h)\|_{L^2(T)}),$$

which finally results in

$$h_T \|f\|_{L^2(T)} \leq (1 + C_{\text{ref}}^2) (h_T \|f_{\mathcal{T}} - f\|_{L^2(T)} + C_{\text{inv}} \|\nabla(u - u_h)\|_{L^2(T)}).$$

3. step. For an interior edge $E \in \mathcal{E}_{\Omega}$, there holds

$$h_E^{1/2} \|[\![\partial_n u_h]\!] \|_{L^2(E)} \leq C (\|\nabla(u - u_h)\|_{L^2(\omega_E)} + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(\omega_E)}): \quad (4.46)$$

To prove this estimate, we define the **edge bubble function**

$$b_E := \prod_{z \in \mathcal{K}_E} \zeta_z \in H_0^1(\omega_E) \cap \mathcal{P}^2(\mathcal{T}).$$

The essential estimate reads

$$\|b_E\|_{L^2(E)} \leq \|b_E^{1/2}\|_{L^2(E)} \leq h_E^{1/2} \leq C_{\text{ref}} \|b_E\|_{L^2(E)}, \quad (4.47)$$

where the constant $C_{\text{ref}} > 0$ is independent of E . In particular, this provides

$$C_{\text{ref}}^{-2} \|[\![\partial_n u_h]\!] \|_{L^2(E)}^2 \leq \|[\![\partial_n u_h]\!] b_E^{1/2}\|_{L^2(E)}^2 = ([\![\partial_n u_h]\!] ; [\![\partial_n u_h]\!] b_E)_{L^2(E)}.$$

Let $T_E^+, T_E^- \in \mathcal{T}$ be the unique elements with $T_E^+ \cap T_E^- = E$ and $\omega_E = T_E^+ \cup T_E^-$. Note that $v := [\![\partial_n u_h]\!] b_E \in \mathcal{P}^2(T_E^{\pm})$ satisfies $v|_{\partial T_E^{\pm} \setminus E} = 0$. Therefore, integration by parts on T_E^{\pm} proves

$$\begin{aligned} ([\![\partial_n u_h]\!] ; v)_{L^2(E)} &= (\partial_n u_h ; v)_{L^2(\partial T_E^+)} + (\partial_n u_h ; v)_{L^2(\partial T_E^-)} \\ &= (\nabla u_h ; \nabla v)_{L^2(\omega_E)} \\ &= (\nabla(u_h - u) ; \nabla v)_{L^2(\omega_E)} + (f ; v)_{L^2(\omega_E)} \\ &\leq (C_{\text{inv}} \|\nabla(u_h - u)\|_{L^2(\omega_E)} + \|h_{\mathcal{T}} f\|_{L^2(\omega_E)}) \|h_{\mathcal{T}}^{-1} v\|_{L^2(\omega_E)}, \end{aligned}$$

where we have applied the Cauchy inequality and an inverse estimate for $v \in \mathcal{P}^2(\mathcal{T})$. For $T \in \{T_E^+, T_E^-\}$ holds

$$\|v\|_{L^2(T)} = [\![\partial_n u_h]\!]_E \|b_E\|_{L^2(T)} \leq |T|^{1/2} [\![\partial_n u_h]\!]_E \leq \frac{h_T^{1/2}}{\sqrt{2}} \|[\![\partial_n u_h]\!] \|_{L^2(E)},$$

since $|T| \leq \frac{1}{2}h_T h_E$. From this, we infer

$$h_E^{1/2} \|h_{\mathcal{T}}^{-1} v\|_{L^2(\omega_E)} \leq \|h_{\mathcal{T}}^{-1/2} v\|_{L^2(\omega_E)} \leq \|[\partial_n u_h]\|_{L^2(E)}.$$

This finally proves

$$h_E^{1/2} \|[\partial_n u_h]\|_{L^2(E)}^2 \leq C_{\text{ref}}^2 (C_{\text{inv}} \|\nabla(u_h - u)\|_{L^2(\omega_E)} + \|h_{\mathcal{T}} f\|_{L^2(\omega_E)}) \|[\partial_n u_h]\|_{L^2(E)}$$

and we may conclude this step by use of step 2 to dominate $\|h_{\mathcal{T}} f\|_{L^2(\omega_E)}$.

4. step. For $T \in \mathcal{T}$ and a Neumann edge $E \in \mathcal{E}_N \cap \mathcal{E}_T$, it holds

$$\begin{aligned} h_E^{1/2} \|\phi - \partial_n u_h\|_{L^2(E)} \\ \leq C (\|h_{\mathcal{E}}^{1/2} (\phi - \phi_{\mathcal{E}})\|_{L^2(E)} + \|\nabla(u - u_h)\|_{L^2(T)} + \|h_{\mathcal{T}}(f - f_{\mathcal{T}})\|_{L^2(T)}) : \end{aligned} \quad (4.48)$$

We consider again the edge bubble function $b_E \in \mathcal{P}^2(T)$ and note that $b_E|_{\partial T \setminus E} = 0$. With $v := (\phi_{\mathcal{E}} - \partial_n u_h)b_E \in \mathcal{P}^2(T)$, we proceed as in step 3 and obtain

$$C_{\text{ref}}^{-2} \|\phi_{\mathcal{E}} - \partial_n u_h\|_{L^2(E)}^2 \leq (\phi_{\mathcal{E}} - \partial_n u_h; v)_{L^2(E)} = (\phi_{\mathcal{E}} - \phi; v)_{L^2(E)} + (\phi - \partial_n u_h; v)_{L^2(E)}.$$

For the second term, we employ integration by parts to see

$$\begin{aligned} (\phi - \partial_n u_h; v)_{L^2(E)} &= (\partial_n(u - u_h); v)_{L^2(\partial T)} \\ &= (\nabla(u - u_h); \nabla v)_{L^2(T)} - (f; v)_{L^2(T)} \\ &\leq (C_{\text{inv}} \|\nabla(u - u_h)\|_{L^2(T)} + \|h_{\mathcal{T}} f\|_{L^2(T)}) h_E^{-1/2} \|\phi_{\mathcal{E}} - \partial_n u_h\|_{L^2(E)}. \end{aligned}$$

The first term is estimated by the Cauchy inequality directly

$$(\phi_{\mathcal{E}} - \phi; v)_{L^2(E)} \leq \|h_{\mathcal{E}}^{1/2} (\phi_{\mathcal{E}} - \phi)\|_{L^2(E)} \|h_{\mathcal{E}}^{-1/2} v\|_{L^2(E)}.$$

There holds

$$\|v\|_{L^2(E)} = |(\phi_{\mathcal{E}} - \partial_n u_h)|_E \|b_E\|_{L^2(E)} \leq h_E^{1/2} |(\phi_{\mathcal{E}} - \partial_n u_h)|_E = \|\phi_{\mathcal{E}} - \partial_n u_h\|_{L^2(E)}.$$

Altogether, we thus have shown

$$\begin{aligned} h_E^{1/2} \|\phi_{\mathcal{E}} - \partial_n u_h\|_{L^2(E)}^2 \\ \leq C_{\text{ref}}^2 (C_{\text{inv}} \|\nabla(u - u_h)\|_{L^2(T)} + \|h_{\mathcal{T}} f\|_{L^2(T)} + \|h_{\mathcal{E}}^{1/2} (\phi_{\mathcal{E}} - \phi)\|_{L^2(E)}) \|\phi_{\mathcal{E}} - \partial_n u_h\|_{L^2(E)}. \end{aligned}$$

Here, $\|h_{\mathcal{T}} f\|_{L^2(T)}$ is estimated by step 2, and $\phi_{\mathcal{E}}$ on the left-hand side is replaced by ϕ with the help of the triangle inequality. \blacksquare

Exercise 53. Prove that $f_{\mathcal{T}}$ in Theorem 4.10 can be replaced by an arbitrary \mathcal{T} -elementwise polynomial $f_{\mathcal{T}} \in \mathcal{P}^m(\mathcal{T})$. The constant $C > 0$ in (4.42)–(4.43) then additionally depends on the polynomial degree $m \in \mathbb{N}_0$. \square

Remark. With the help of a so-called extension operator that extends a polynomial $p : E \rightarrow \mathbb{R}$ to a polynomial $F_{\text{ext}} p : T \rightarrow \mathbb{R}$, one can show that $\phi_{\mathcal{E}}$ in Theorem 4.10 can be replaced by an arbitrary \mathcal{E}_N -edgewise polynomial (with respect to the arclength). \square

Actually, the the volume residual contribution $\|h_{\mathcal{T}} f\|_{L^2(\Omega)} = \mathcal{O}(h)$ to η can be improved. This is done in the following exercise, where this term is replaced by some higher-order term $\mathcal{O}(h^2)$.

Exercise 54. Let $\Omega_z = \text{supp}(\zeta_z)$ denote the node patch of $z \in \mathcal{K}$. For $f \in L^2(\Omega)$, let $f_z := |\Omega_z|^{-1} \int_{\Omega_z} f dx$ denote the corresponding integral mean. Prove the following claims:

(i) For all inner nodes $z \in \mathcal{K} \setminus \Gamma$, it holds

$$\int_{\Omega_z} f f_z \zeta_z dx \leq C \left(\sum_{\substack{E \in \mathcal{E}_\Omega \\ z \in E}} \|\llbracket \partial_n u_h \rrbracket\|_{L^2(E)}^2 \right)^{1/2} \|h_{\mathcal{T}}^{-1/2} f_z\|_{L^2(\Omega_z)}.$$

(ii) For all inner nodes $z \in \mathcal{K} \setminus \Gamma$ and elements $T \in \mathcal{T}$ with $z \in T$, it holds

$$C^{-1} \|h_{\mathcal{T}} f\|_{L^2(T)}^2 \leq \|h_{\mathcal{T}}(f - f_z)\|_{L^2(\Omega_z)}^2 + \sum_{\substack{E \in \mathcal{E}_\Omega \\ z \in E}} \|h_E^{1/2} \llbracket \partial_n u_h \rrbracket\|_{L^2(E)}^2.$$

(iii) Derive the equivalence

$$\begin{aligned} C^{-1} \eta^2 \leq \tilde{\eta}^2 &:= \|h_{\mathcal{E}}^{1/2} \llbracket \partial_n u_h \rrbracket\|_{L^2(\mathcal{E}_\Omega)}^2 + \|h_{\mathcal{E}}^{1/2}(\phi - \partial_n u_h)\|_{L^2(\mathcal{E}_N)}^2 \\ &+ \sum_{z \in \mathcal{K} \setminus \Omega} \|h_{\mathcal{T}}(f - f_z)\|_{L^2(\Omega_z)}^2 \leq C \eta^2. \end{aligned}$$

(iv) Conclude that the improved error estimator $\tilde{\eta}$ is reliable and efficient.

(v) In what sense is the error estimator $\tilde{\eta}$ improved when compared to η .

The constant $C > 0$ in (i)–(iii) depends only on γ -shape regularity of \mathcal{T} . □

The following MATLAB code computes the vector of the element-based refinement indicators η_T from (4.37). The integral

$$h_T^2 \|f\|_{L^2(T)}^2 = h_T^2 \int_T f^2 dx \approx h_T^2 |T| f(s_T)^2 \simeq |T|^2 f(s_T)^2$$

is computed by 1-point quadrature associated with the center of mass s_T of T . The integral

$$\begin{aligned} h_T \|\phi - \partial_n u_h\|_{L^2(\partial T \cap \Gamma_N)}^2 &= \sum_{E \in \mathcal{E}_T \cap \mathcal{E}_N} h_T \int_E (\phi - \partial_n u_h)^2 ds \\ &\approx \sum_{E \in \mathcal{E}_T \cap \mathcal{E}_N} h_T h_E (\phi(m_E) - (\partial_n u_h)|_E)^2 \\ &\simeq |T| \sum_{E \in \mathcal{E}_T \cap \mathcal{E}_N} (\phi(m_E) - (\partial_n u_h)|_E)^2 \end{aligned}$$

is computed edge-wise by midpoint quadrature.

```

1 function etaR = computeEtaR(x,coordinates,elements,f,dirichlet,neumann,phi)
2
3 % ETAR = COMPUTEETAR(X,COORDINATES,ELEMENTS,F,DIRICHLET,NEUMANN,PHI)
    
```

```

4 % computes the element-based refinement indicators associated with
5 % the residual-based error estimator. ETAR is a column vector, where
6 %  $ETAR(j)^2 = |T_j| * || f ||_{L^2(T_j)}^2$ 
7 %   +  $|T_j|^{1/2} * || \text{jump}(\partial_n u_h) ||_{L^2(\partial T_j \cap \Omega)}^2$ 
8 %   +  $|T_j|^{1/2} * || \phi - \partial_n u_h ||_{L^2(\partial T_j \cap \Gamma_N)}^2$ 
9 % The exact integrals involving F and PHI\lastmodified{11.05.2009}
10 % are integrated by midpoint
11 % quadrature.
12
13 % (c) 2007,2008 by Dirk Praetorius, last modified 08.01.2008
14 % dirk.praetorius@tuwien.ac.at - http://www.asc.tuwien.ac.at/~dirk
15
16 M = size(elements,1);
17 N = size(coordinates,1);
18
19 etaR = zeros(M,1);
20 int = sparse(N,N);
21
22 %*** Compute normal derivatives  $\partial_n T(u_h)$  on all edges
23 for j = 1:M
24     nodes = elements(j,:);
25     B = [1 1 1 ; coordinates(nodes,:)']';
26     G = B \ [0 0 ; 1 0 ; 0 1];
27     grad = G'*x(nodes); % gradient  $\nabla u_h$  on element  $T_j$ 
28     for k = 1:3
29         node1 = nodes(k);
30         node2 = nodes(mod(k,3)+1);
31         normal = [1 -1]*coordinates([node1,node2],:);
32         normal = normal*[0 1;-1 0] / norm(normal);
33         int(node1,node2) = normal*grad;
34     end
35 end
36
37 %*** Delete data in case of Dirichlet edges
38 for j = 1:size(dirichlet,1)
39     nodes = dirichlet(j,:);
40     int(nodes(1),nodes(2)) = 0;
41 end
42
43 %*** Evaluate exact Neumann data on Neumann edges
44 for j = 1:size(neumann,1)
45     nodes = neumann(j,:);
46     m = [1 1]*coordinates(nodes,+)/2;
47     int(nodes(2),nodes(1)) = -phi(m);
48 end

```

```

49
50  %*** Compute residual-based refinement indicators
51  for j = 1:M
52      nodes = elements(j,:);
53      %*** Compute volume contribution by midpoint quadrature
54      sizeT = det([1 1 1 ; coordinates(nodes,:)'])/2;
55      s = [1 1 1]*coordinates(nodes,+)/3;
56      etaR(j) = sizeT^2*f(s)^2;
57      %*** Add edge contributions
58      for k = 1:3
59          node1 = nodes(k);
60          node2 = nodes(mod(k,3)+1);
61          hE = norm([1 -1]*coordinates([node1,node2],:));
62          etaR(j) = etaR(j) + sizeT*(int(node1,node2)+int(node2,node1))^2;
63      end
64  end
65  etaR = sqrt(etaR);

```

Exercise 55. Consider the homogenous Dirichlet problems

$$\begin{aligned} -\Delta u &= 1 && \text{in } \Omega, \\ u &= 0 && \text{on } \Gamma = \partial\Omega, \end{aligned}$$

with Ω being either the square $\Omega = (-1, 1)^2$ or the L -shaped domain $\Omega = (-1, 1)^2 \setminus [0, 1]^2$. Plot error and error estimator over the number of elements. How can one use the plot to see whether an error estimator is reliable and/or efficient? \square

Exercise 56. Note that the computational time of the function `computeEtaR` grows quadratically with the number $M = \#\mathcal{T}$ of elements. This is due to the successive assembly of the sparse matrix `int`. Improve the implementation so that one observes real linear complexity. \square

Exercise 57. For the computation of the residual-based refinement indicators η_T , the given MATLAB codes approximates the exact data f and ϕ for the terms

$$\|h_{\mathcal{T}}f\|_{L^2(T)} \quad \text{and} \quad \|h_{\mathcal{E}}^{1/2}(\phi - \partial_n u_h)\|_{L^2(E)} \quad \text{for } T \in \mathcal{T} \text{ resp. } E \in \mathcal{E}_N$$

by $f|_T \approx f(s_T)$ and $\phi|_E \approx \phi(m_E)$. Here, s_T and m_E denote the center of mass of T and the midpoint of E , respectively. Formally, this leads to an approximation $\tilde{\eta}_R$ of η_R . Prove that, for $f \in H^2(\mathcal{T})$ and $\phi \in C^2(\mathcal{E}_N)$, there holds

$$|\eta_R - \tilde{\eta}_R| \leq C (\|h_{\mathcal{T}}^2 \nabla f\|_{H^1(\mathcal{T})} + \|h_{\mathcal{E}}^{3/2} \phi'\|_{C^1(\mathcal{E}_N)})$$

with a constant $C > 0$ that only depends on $\sigma(\mathcal{T})$. Consequently, the computed estimator $\tilde{\eta}_R$ is, in fact, reliable and efficient up to terms of higher order. \square

4.4 Adaptive Mesh-Refining Algorithm (15.12.2014)

Usually, a posteriori error estimates are not only used to estimate the (unknown) error $\|\nabla(u - u_h)\|_{L^2(\Omega)}$ but even to steer the local mesh-refinement. Let

v15:
29.11.2017

$$\eta := \left(\sum_{T \in \mathcal{T}} \eta(T)^2 \right)^{1/2}$$

be an a posteriori error estimator, where the quantities $\eta(T) := \eta_T$ reflect —at least heuristically— the (unknown) local error $\|\nabla(u - u_h)\|_{L^2(T)}$ for all $T \in \mathcal{T}$. We then aim to refine only the elements $T \in \mathcal{T}$, where $\eta(T)$ is large. Therefore, the quantities $\eta(T)$ are usually called **refinement indicators** (or error indicators). To state our version of an adaptive algorithm, we introduce some additional notation which will be used from now on.

- the index $\ell \in \mathbb{N}_0$ denotes the step of the adaptive algorithm,
- \mathcal{T}_ℓ is the mesh in the ℓ -th step of the adaptive algorithm.
- \mathcal{N}_ℓ and \mathcal{E}_ℓ denote the associated sets of nodes and edges, respectively.
- $U_\ell \in \mathcal{X}_\ell := \mathcal{S}_D^1(\mathcal{T}_\ell)$ denotes the Galerkin solution in the ℓ -th step.
- $h_\ell \in \mathcal{P}^0(\mathcal{T}_\ell)$, $h_\ell|_T := \text{diam}(T)$ is the local mesh-side function.

With this notation, one common strategy is the following: Let $\theta \in (0, 1)$ be the parameter for the adaptive algorithm.

Algorithm 4.11 (Adaptive Mesh-Refinement). **Input:** Initial triangulation \mathcal{T}_0 , tolerance $\tau > 0$, adaptivity parameter $\theta \in (0, 1)$, counter $\ell := 0$.

- (i) Compute discrete solution U_ℓ .
- (ii) Compute refinement indicators $\eta_\ell(T)$ and error estimator $\eta_\ell = \left(\sum_{T \in \mathcal{T}_\ell} \eta_\ell(T)^2 \right)^{1/2}$.
- (iii) Stop computation provided that $\eta_\ell \leq \tau$
- (iv) Choose the minimal set $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell$ of marked elements such that

$$\theta \eta_\ell^2 = \theta \sum_{T \in \mathcal{T}_\ell} \eta_\ell(T)^2 \leq \sum_{T \in \mathcal{M}_\ell} \eta_\ell(T)^2. \quad (4.49)$$

- (v) Generate a new regular mesh $\mathcal{T}_{\ell+1}$, where at least all marked elements have been refined.
- (vi) Update $\ell \mapsto \ell + 1$ and goto (i).

Output: Finite sequence of discrete solutions U_ℓ and corresponding error estimators η_ℓ .

Remark. Clearly, the stopping criterion (iii) is only meaningful if η_ℓ is reliable and if the reliability constant in $\|\nabla(u - U_\ell)\|_{L^2(\Omega)} \leq C_{\text{rel}} \eta_\ell$ is known. In practice, runtime and storage requirements are the limiting quantities for a numerical simulation. Usually, one thus uses rather a maximal runtime or a maximal storage requirement, e.g., the maximal number of elements, as a stopping criterion. Adaptivity is then used to obtain an—in some sense—optimal approximation with respect to these side constraints. \square

Remark. The marking criterion (4.49) was introduced by DÖRFLER (1996). It will be crucial to prove convergence of U_ℓ to the exact solution $u \in H_D^1(\Omega)$ of (4.1). Note that the choice $\theta \rightarrow 0$ in (4.49) leads to highly adapted meshes, whereas $\theta \rightarrow 1$ corresponds to (almost) uniform mesh-refinement. However, for small θ , only a few elements are refined per step. This may result in too many steps in the sense that usually the assembly of the Galerkin data is the most time consuming part of the algorithm. In practice, a good compromise between sufficient mesh-adaption and as few steps in the loop as possible appears to be $\theta \approx 0.25$. \square

Remark. In the beginning of the analysis of adaptive FEM, Babuška proposed the following marking criterion: An element $T \in \mathcal{T}$ is marked for refinement if and only if

$$\eta_T \geq \theta \max \{ \eta_{T'} \mid T' \in \mathcal{T} \}, \quad (4.50)$$

which is called **bulk criterion** in the literature. Convergence (but *not* optimality) of this version of adaptive FEM was proven by MORIN, SIEBERT & VEESER (2008). Very recently, DIENING, KREUZER & STEVENSON (2014) proved the so-called *instance optimality* of the adaptive algorithm for some extended bulk criterion. \square

Before we comment on the local mesh-refinement in step (v) of Algorithm 4.11, we give a simple MATLAB realization of Algorithm 4.11. We use the number of elements $M := \#\mathcal{T}$ and stop the adaptive algorithm when $M \geq M_{\text{max}}$.

```

1 function [x,M,energy,etaR] ...
2     = solveLaplaceAdaptively(coordinates,elements,f,dirichlet,neumann,phi,theta,Mmax)
3
4 ell = 1;
5 while 1
6     M(ell) = size(elements,1);
7
8     %*** Compute discrete solution and cooresponding energy
9     [x,energy(ell)] = solveLaplace(coordinates,elements,f,dirichlet,neumann,phi);
10
11     %*** Compute refinement indicators and error estimator
12     indicators = computeEtaR(x,coordinates,elements,f,dirichlet,neumann,phi);
13     etaR(ell) = norm(indicators);

```

```

14
15     *** Stopping criterion
16     if M(e11) >= Mmax
17         break
18     end
19
20     *** Use Doerfler marking to mark elements for refinement
21     [indicators,idx] = sort(indicators.^2,'descend');
22     sumeta = cumsum(indicators);
23     m = find(theta*sumeta(end)<=sumeta,1);
24     marked = idx(1:m);
25
26     *** Generate a new mesh by RGB-refinement
27     [coordinates,elements,dirichlet,neumann] = ...
28         rgbrefine(coordinates,elements,dirichlet,neumann,marked);
29
30     *** Update counter
31     e11 = e11 + 1;
32 end

```

4.4.1 Red-Green-Blue Refinement

It now remains to discuss the mesh-refinement. Recall that all error estimates are affected by the shape regularity $\sigma(\mathcal{T}_\ell)$ in the sense that the involved constants become unbounded for $\sigma(\mathcal{T}_\ell) \xrightarrow{\ell \rightarrow \infty} \infty$. Therefore, the mesh-refinement has to take care of the interior angles of the elements $T \in \mathcal{T}_\ell$ since $\sigma(\mathcal{T}_\ell)$ tends to infinity if and only if the minimal interior angle of the triangulation tends to zero. We follow the so-called **red-green-blue strategy** (or **RGB-refinement**): This refinement strategy is based on edge-refinement. First, we thus use the following marking rule:

- If an element $T \in \mathcal{T}_\ell$ is marked for refinement, we mark all edges $E \in \mathcal{E}_T$ for refinement.

We now proceed recursively as follows:

- For each element $T \in \mathcal{T}_\ell$, we mark its *longest* edge $E \in \mathcal{E}_T$ for refinement provided that \mathcal{E}_T contains a marked edge.

Each marked edge will be halved, i.e., the midpoint m_E of a marked edge belongs to the new set $\mathcal{K}_{\ell+1}$ of nodes. Finally, we have the following refinement rules, for all $T \in \mathcal{T}_\ell$:

- If no edge in \mathcal{E}_T is marked for refinement, T is not refined, i.e., $T \in \mathcal{T}_{\ell+1}$.
- If all edges in \mathcal{E}_T are marked, we use a **red-refinement** of T , i.e., T is refined uniformly into four similar triangles, cf. Figure 4.3.
- If one edge in \mathcal{E}_T is marked (and hence the longest edge), we use a **green-refinement**, i.e., T is refined into two triangles, cf. Figure 4.4.
- If two edges in \mathcal{E}_T are marked — one of which is, according to the marking rule, the longest edge of T —, we use a **blue-refinement**, i.e., T is split into three triangles, cf. Figure 4.5.



FIGURE 4.3. Red-refinement: If all edges of a triangle $T \in \mathcal{T}_\ell$ are marked (left), T is refined into four similar triangles $T_1, T_2, T_3, T_4 \in \mathcal{T}_{\ell+1}$ (right).



FIGURE 4.4. Green-refinement: If only the longest edge of a triangle $T \in \mathcal{T}_\ell$ is marked (left), T is refined into two new triangles $T_1, T_2 \in \mathcal{T}_{\ell+1}$ (right).

In Figure 4.6, we visualize a simple example for an RGB-refined mesh.

We state the following elementary but important theorem without a proof.

Theorem 4.12. *Let \mathcal{T}_0 be a regular triangulation such that $\varepsilon > 0$ is a lower bound for the smallest angle of a triangle $T \in \mathcal{T}_0$. Let \mathcal{T}_ℓ be a sequence of meshes, where \mathcal{T}_ℓ is obtained by RGB-refinement of the mesh $\mathcal{T}_{\ell-1}$ and where the set $\mathcal{M}_{\ell-1} \subseteq \mathcal{T}_{\ell-1}$ of marked elements is arbitrary. Then, \mathcal{T}_ℓ is regular and the smallest angle of all triangles $T \in \mathcal{T}_\ell$ is bounded from below by $\varepsilon/2$. In particular, there holds*

$$\sup_{\ell \in \mathbb{N}} \sigma(\mathcal{T}_\ell) < \infty, \quad (4.51)$$

which is an equivalent formulation for the fact that the smallest angles of the triangulations \mathcal{T}_ℓ do not tend to zero. ■

The following MATLAB code is an implementation of the RGB mesh-refinement strategy, which additionally takes care of the specification of the domain boundary.

```

1 function [coordinates,newelements,varargout] ...
2     = rgbrefine(coordinates,elements,varargin)
3
4 % [COORDINATES,ELEMENTS [,DIRICHLET] [,ROBIN] [,NEUMANN] ]
5 % = RGBREFINE(COORDINATES,ELEMENTS [,DIRICHLET] [,ROBIN] [,NEUMANN], MARKED)
6 % refines the MARKED elements of a regular triangulation by a
7 % uniform refinement (red refinement). A green-blue closure leads
8 % to a new regular triangulation.
9 %
10 % vector MARKED contains the indices of all elements that will be refined
11 %
```

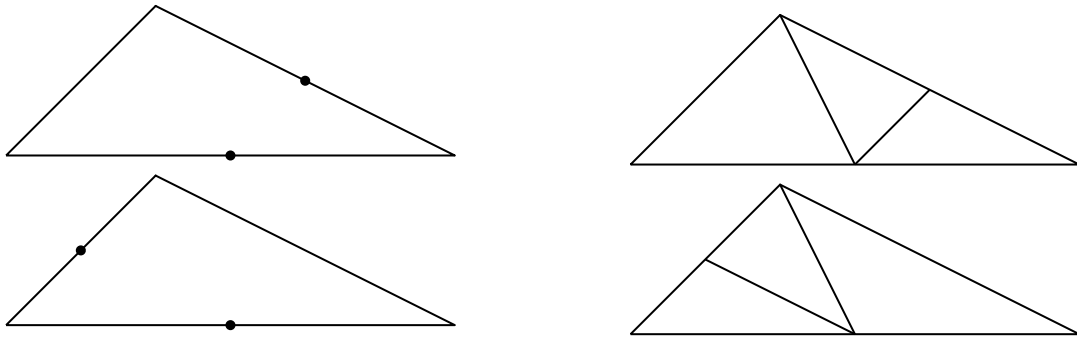


FIGURE 4.5. *Blue-refinement: If besides the longest edge of a triangle $T \in \mathcal{T}_\ell$ just one other edge is marked for refinement (left), T is refined into three new triangles $T_1, T_2, T_3 \in \mathcal{T}_{\ell+1}$ (right).*

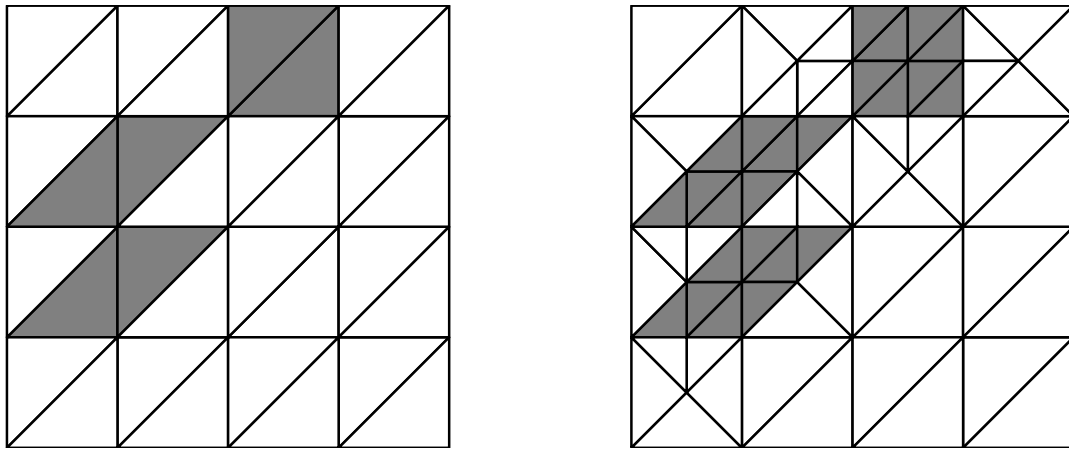


FIGURE 4.6. *The left plot shows an initial mesh \mathcal{T}_ℓ with marked elements coloured in grey. The right plot shows the mesh $\mathcal{T}_{\ell+1}$ obtained by RGB-refinement of the marked elements. The grey elements are obtained by uniform refinement of a marked element $T \in \mathcal{T}_0$.*

```

12 % Optionally, one can provide the specification of boundary conditions,
13 % e.g., Dirichlet, Robin, and/or Neumann boundaries. Then, the refined
14 % boundary conditions are returned in the same order
15
16 % (c) 2007 by Dirk Praetorius, last modified 21.11.2007
17 % dirk.praetorius@tuwien.ac.at - http://www.math.tuwien.ac.at/~dirk
18
19 M = size(elements,1);
20 N = size(coordinates,1);
21 markedelements = varargin{end};
22
23 %*** Sort elements such that the longest edge is always the first edge,
24 %*** i.e. we sort the entries in each row elements(j,:) accordingly.
25

```



```

26 for j = 1:M
27     [hmax,idx] = max(sum((coordinates(elements(j,[2,3,1]),:)- ...
28                         coordinates(elements(j,[1,2,3]),:)).^2'));
29     elements(j,:) = elements(j,[idx,mod(idx,3)+1,mod(idx+1,3)+1]);
30 end
31
32 %*** Introduce numbering of edges, stored in a sparse matrix EDGES:
33 %*** - EDGES(J,K) \neq 0 if and only if nodes J and K connected by edge,
34 %*** - EDGES(J,K) \neq EDGES(K,J) if and only if edge on boundary.
35
36 edges = sparse(N,N);
37 noedges = 0; % number of edges
38 for j = 1:M
39     for k = 1:3
40         a = [elements(j,k),elements(j,mod(k,3)+1)];
41         if edges(a(2),a(1))
42             edges(a(1),a(2)) = edges(a(2),a(1));
43         else
44             noedges = noedges+1;
45             edges(a(1),a(2)) = noedges;
46         end
47     end
48 end
49
50 %*** Transfer marking of elements to marking of edges.
51 %*** If element(j) is marked, we mark all of its edges (red-refinement).
52 %*** - MARKEDEDGES(k) \neq 0 if and only if edge K will be refined.
53
54 element2edges = zeros(M,3);
55 for j = 1:M
56     element2edges(j,:) = diag(edges(elements(j,:),elements(j,[2,3,1])))';
57 end
58 markededges = sparse(noedges,1);
59 markededges(element2edges(markedelements,:)) = ones(3*length(markedelements),1);
60
61 %*** Mark further edges according to green-blue closure:
62 %*** To ensure that the triangles do not degenerate, we always refine
63 %*** the longest edge, i.e. the first edge of an element.
64
65 edge2elements = sparse(N,N);
66 for j = 1:M
67     edge2elements(elements(j,:),elements(j,[2,3,1])) = ...
68         edge2elements(elements(j,:),elements(j,[2,3,1]))+j*eye(3,3);
69     k = j;
70     while k

```

```

71     I = element2edges(k,:);
72     if markededges(I(1))==1 | markededges(I(2:3))==[0;0]
73         k = 0;
74     else
75         markededges(I(1))=1;
76         k = edge2elements(elements(k,2),elements(k,1));
77     end
78 end
79 end
80
81 %*** For each marked edge, its midpoint becomes a new node.
82 %*** We store the number of the new nodes in MARKEDEDGES instead of 1.
83
84 idx = find(markededges);
85 markededges(idx) = N+1:N+length(idx);
86 for j = 1:nnz(markededges)
87     [a,b] = find(idx(j) == edges);
88     coordinates(markededges(idx(j)),:)=(coordinates(a(1),:)+coordinates(b(1),:))/2;
89 end
90
91 %*** Create new elements
92
93 I = reshape(edges(size(edges,1)*(elements(:, [2,3,1])-1)+elements(:, [1,2,3]))),M,3);
94
95 boundaryedges = nonzeros(tril(abs(edges-edges')));
96 newelements = zeros(2*length(idx)-nnz(markededges(boundaryedges))+M,3);
97
98 counter = 0;
99 for j = 1:M
100     RefineEdge = find(markededges(I(j,:)));
101     newnodes=markededges(I(j,RefineEdge))';
102     if size(RefineEdge,1)==3 % red refinement
103         new = [ newnodes([2,3,1]);
104                 [elements(j,1) newnodes(1) newnodes(3)];
105                 [newnodes(1) elements(j,2) newnodes(2)];
106                 [newnodes(3) newnodes(2) elements(j,3)] ];
107     elseif size(RefineEdge,1)==2 % blue refinement
108         new = [ [newnodes(1), elements(j,RefineEdge(2)),newnodes(2)];
109                 [elements(j,5-RefineEdge(2)), ...
110                 elements(j,rem(5-RefineEdge(2),3)+1),newnodes(1)];
111                 [elements(j,rem(RefineEdge(2),3)+1),newnodes(1),newnodes(2)] ];
112     elseif size(RefineEdge,1)==1 % green refinement
113         new = [ [elements(j,[2,3]),newnodes];
114                 [elements(j,[3,1]),newnodes] ];
115     else % no refinement

```

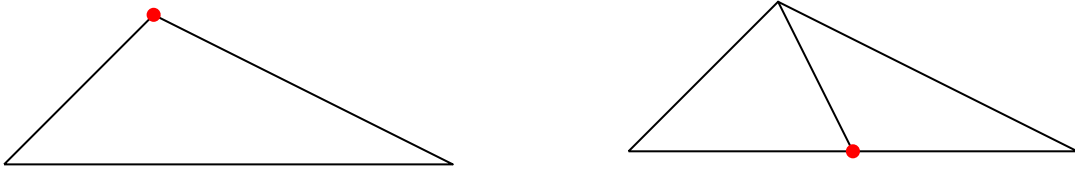


FIGURE 4.7. *Newest vertex bisection: For each element $T \in \mathcal{T}_\ell$, there is one newest vertex (left). The edge opposite to the newest vertex is the reference edge (which is not necessarily the longest edge of T). If the element is marked for refinement, the reference edge is halved and T is split into two sons. The new node becomes the newest vertex of the sons (right).*

```

116     new = elements(j,:);
117     end
118     newelements(counter+1:counter+size(new,1),:) = new;
119     counter = counter + size(new,1);
120 end
121
122 %*** Update boundary conditions
123
124 for j = 1:nargin-3
125     boundary = varargin{j};
126     if ~isempty(boundary)
127         counter = 0;
128         boundarynr = edges(size(edges,1)*(boundary(:,2)-1)+boundary(:,1));
129         for k = 1:size(boundary,1)
130             if markededges(boundarynr(k))
131                 boundary = [ boundary(1:k-1+counter,:);
132                             boundary(k+counter,1),markededges(boundarynr(k));
133                             markededges(boundarynr(k)),boundary(k+counter,2);
134                             boundary(k+1+counter:size(boundary,1),:) ];
135                 counter = counter + 1;
136             end
137         end
138     end
139     varargout(j) = {boundary};
140 end

```

4.4.2 Newest Vertex Bisection

An alternative to red-green-blue refinement is the so-called newest vertex bisection, which is very popular because of two reasons: First, it is easy to implement. Second, mesh-refinement by newest vertex bisection can be characterized by a binary forest, i.e., each element $T \in \mathcal{T}_0$ of the initial mesh is the root of a binary tree. From a mathematical point of view, this makes newest vertex bisection easier to analyze than the red-green-blue strategy.

Newest vertex bisection needs an initialization step on the initial mesh \mathcal{T}_0 :

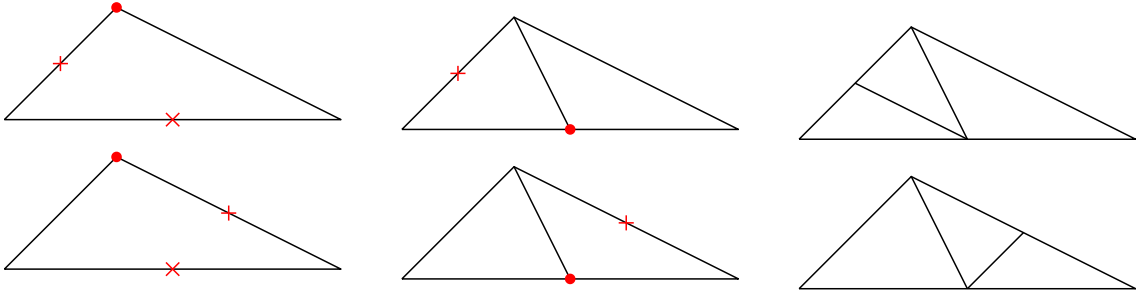


FIGURE 4.8. *Bisec(2)*: If two edges of T are marked for refinement, one is the reference edge (left). This is bisected in a first step, and the second marked edge becomes the reference edge of one son element (middle). In a second step, the corresponding son element is bisected (right).

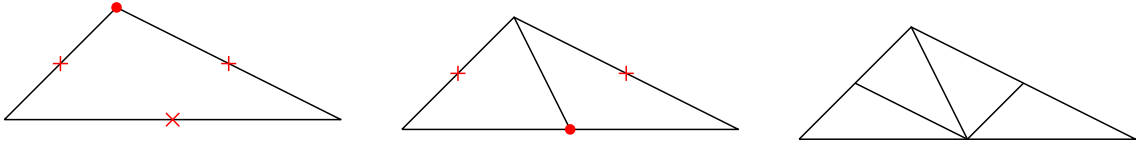


FIGURE 4.9. *Bisec(3)*: If all edges of T are marked for refinement (left), we bisect the reference edge first (middle). The other two edges become the reference edges of the two son elements. In a second step, we bisect both son elements (right).

- For each triangle of the initial mesh \mathcal{T}_0 , we choose an arbitrary (newest) vertex.
- The **reference edge** is always the edge opposite to the newest vertex.

Then, the refinement rule reads as follows:

- A marked element T is refined into two sons T_1, T_2 by halving the reference edge, cf. Figure 4.7.
- The midpoint z_T of the reference edge of T is the newest vertex of T_1 and T_2 , i.e., the other two edges of T now become reference edges of T_1 and T_2 , respectively.

To ensure that this procedure leads to a regular triangulation, the entire refinement strategy reads as follows:

- If an element $T \in \mathcal{T}_\ell$ is marked for refinement, we mark the reference edge $E \in \mathcal{E}_T$ (or even all edges $E \in \mathcal{E}_T$) for refinement.

We then proceed recursively as follows:

- For any element $T \in \mathcal{T}_\ell$, we mark its reference edge $E \in \mathcal{E}_T$ for refinement provided that \mathcal{E}_T contains a marked edge.

Now, we simply do iterated newest vertex bisection until we are led to a regular triangulation. Luckily, only 4 configurations per element can appear, and this iteration has at most two steps, which are realized as for red-green-blue refinement:

- If no edge in \mathcal{E}_T is marked for refinement, T is not refined, i.e., $T \in \mathcal{T}_{\ell+1}$.

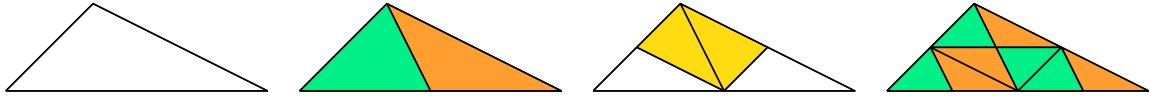


FIGURE 4.10. Refinement by newest vertex bisection only leads to finitely many interior angles for the family of all possible triangulations obtained by arbitrary newest vertex bisections. To see this, we start from a macro element (left), where the bottom edge is the reference edge. Using iterated newest vertex bisection, one observes that only four similarity classes of triangles occur, which are indicated by the coloring. After three steps of bisections (right), no additional similarity class appears.

- If all edges in \mathcal{E}_T are marked, we first bisect T into two sons T_1, T_2 . Note that now the reference edges of T_1 and T_2 are marked for refinement, so that we bisect each son into its sons, cf. Figure 4.9.
- If only one edge in \mathcal{E}_T is marked (and hence the reference edge), we bisect T into two sons, cf. Figure 4.7.
- If two edges in \mathcal{E}_T are marked — one of which is, according to the marking rule, the reference edge of T —, we bisect T into two sons T_1, T_2 . Now, the reference edge of one of these sons is marked. Hence, we bisect this son into its sons, cf. Figure 4.8.

Again, we state the elementary but important theorem that newest vertex bisection does not lead to degenerate triangles.

Theorem 4.13. *Let \mathcal{T}_0 be a regular triangulation. Let \mathcal{T}_ℓ be a sequence of meshes, where \mathcal{T}_ℓ is a refinement of $\mathcal{T}_{\ell-1}$ based on newest vertex bisection and an arbitrary set $\mathcal{M}_{\ell-1} \subseteq \mathcal{T}_{\ell-1}$ of marked elements. Then, \mathcal{T}_ℓ is regular and there holds*

$$\sup_{\ell \in \mathbb{N}} \sigma(\mathcal{T}_\ell) < \infty. \quad (4.52)$$

More precisely, there are only finitely many angles in $\bigcup_{\ell \in \mathbb{N}_0} \mathcal{T}_\ell$. In particular, newest vertex bisection only leads to finitely many similarity classes of triangles and patches.

Proof. The proof follows from Figure 4.10 and some basic geometry on high-school level. ;-) ■

Exercise 58. Modify the MATLAB code of `rgbrefine` to obtain a mesh-refinement based on newest vertex bisection. □

4.5 Convergence of Adaptive FEM (07.12.2017)

In the following, we aim to show that Algorithm 4.11 creates a sequence U_ℓ of discrete solutions which converges to the exact solution $u \in H := H_D^1(\Omega)$. The adaptive algorithm generates a sequence $\mathcal{X}_\ell = \mathcal{S}_D^1(\mathcal{T}_\ell)$ of finite dimensional nested subspaces of H , i.e., $\mathcal{X}_\ell \subsetneq \mathcal{X}_{\ell+1}$ for all $\ell \in \mathbb{N}_0$, since $\mathcal{T}_{\ell+1}$ is some refinement of \mathcal{T}_ℓ . We first stress that the sequence U_ℓ is always convergent to some limit $U_\infty \in H$. However, we even stress that one may in general expect that $U_\infty \neq u$.

Exercise 59. Let \mathcal{X}_ℓ be nested subspaces of a Hilbert space H , i.e., $\mathcal{X}_\ell \subseteq \mathcal{X}_{\ell+1}$ for all $\ell \in \mathbb{N}_0$. Let $\langle \cdot ; \cdot \rangle$ be an elliptic and continuous bilinear form on H with corresponding Galerkin solutions $U_\ell \in \mathcal{X}_\ell$. Prove that the limit $U_\infty := \lim_{\ell \rightarrow \infty} U_\ell$ exists in H . **Hint:** Define \mathcal{X}_∞ as the closure of $\bigcup_{\ell=0}^{\infty} \mathcal{X}_\ell$ in H . Let $U_\infty \in \mathcal{X}_\infty$ be the corresponding Galerkin solution, and prove that U_∞ is the limit of the sequence U_ℓ . \square

Exercise 60. Let $H = H_D^1(\Omega)$ and $\mathcal{X}_\ell = \mathcal{S}_D^1(\mathcal{T}_\ell)$, where the regular initial mesh \mathcal{T}_0 is given and where \mathcal{T}_ℓ is obtained iteratively by uniform refinement of $\mathcal{T}_{\ell-1}$. Prove that $\mathcal{X}_\infty = H$ for the space \mathcal{X}_∞ from Exercise 59. \square

The interpretation of the last exercises is the following: For uniform mesh-refinement, there usually holds $\mathcal{X}_\infty = H$ and thus $u = U_\infty$, i.e., we have convergence of the sequence of discrete solutions U_ℓ towards the unique solution u . However, adaptive mesh-refinement may lead to $\mathcal{X}_\infty \subsetneq H$. Consequently, the question arises whether the adaptive algorithm guarantees $U_\infty = u$ or not. This will be discussed in the following sections.

Throughout the subsequent section, we use the following notation, which is now collected for the convenience of the reader:

- $U_\ell \in \mathcal{X}_\ell := \mathcal{S}_D^1(\mathcal{T}_\ell)$ denotes the Galerkin solution.
- For $T \in \mathcal{T}_\ell$ and some $V \in \mathcal{S}_D^1(\mathcal{T}_\ell)$, $\eta_\ell(T, V)$ denotes the associated refinement indicator, e.g.,

$$\eta_\ell(T, V)^2 = h_T^2 \|f\|_{L^2(T)}^2 + h_T \|[\partial_n V]\|_{L^2(\partial T \cap \Omega)}^2 + h_T \|\phi - \partial_n V\|_{L^2(\partial T \cap \Gamma_N)}^2. \quad (4.53)$$
- For some subset $\mathcal{M} \subseteq \mathcal{T}_\ell$ and $V \in \mathcal{S}_D^1(\mathcal{T}_\ell)$, let $\eta_\ell(\mathcal{M}, V) := (\sum_{T \in \mathcal{M}} \eta_\ell(T, V)^2)^{1/2}$.
- We abbreviate $\eta_\ell(\mathcal{M}) = \eta_\ell(\mathcal{M}, U_\ell)$ and $\eta_\ell = \eta_\ell(\mathcal{T}_\ell)$.

Note that in case of (4.53), η_ℓ is the residual a posteriori error estimator discussed in Section 4.3. We recall some technical terms, proven above for the residual error estimator η_ℓ .

- η_ℓ is **reliable** if

$$\|u - U_\ell\|_H \leq C_{\text{rel}} \eta_\ell. \quad (4.54)$$

- η_ℓ is **efficient** (up to oscillation terms which depend only on \mathcal{T}_ℓ), if

$$\eta_\ell \leq C_{\text{eff}} (\|u - U_\ell\|_H + \text{osc}_\ell), \quad (4.55)$$

where $\text{osc}_\ell := \text{osc}_\ell(\mathcal{T}_\ell)$, $\text{osc}_\ell(\mathcal{M}) := (\sum_{T \in \mathcal{M}} \text{osc}_\ell(T)^2)^{1/2}$ for $\mathcal{M} \subseteq \mathcal{T}_\ell$, and

$$\text{osc}_\ell(T)^2 := h_T^2 \|f - f_{\mathcal{T}}\|_{L^2(T)}^2 + h_T \|\phi - \phi_{\mathcal{E}}\|_{L^2(\partial T \cap \Gamma_N)}^2. \quad (4.56)$$

- The set $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell$ of marked elements is usually assumed to satisfy the **Dörfler marking**

$$\theta \eta_\ell \leq \eta_\ell(\mathcal{M}_\ell) \quad (4.57)$$

for some fixed parameter $\theta \in (0, 1)$.

Exercise 61. Prove that $\|u - U_\ell\|_H$ as well as osc_ℓ are monotonously decreasing for $\ell \rightarrow \infty$. Prove that in case of the residual-based indicators (4.53), there holds $\text{osc}_\ell(T) \leq \eta_\ell(T)$ for all $T \in \mathcal{T}_\ell$, i.e., the error estimator dominates the oscillation terms. \square

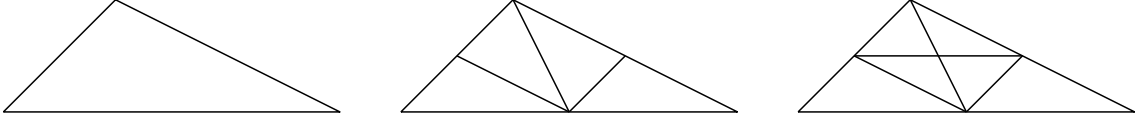


FIGURE 4.11. *Bisec(5) guarantees the inner node property: Let T be marked for refinement (left) and assume that the bottom edge is the reference edge. With five bisections, we pass the configuration of bisec(3) in the middle and end up with an inner node (right).*

4.5.1 Inner Node Property and Strict Error Reduction

The idea of the analysis in this section goes back to DÖRFLER (1996) and has been elaborated by MORIN, NOCHETTO & SIEBERT (2000). Our presentation follows the latter work. The main idea is the so-called discrete local efficiency which is only proven, to the best of our knowledge, for the residual-based error estimator with local contributions (4.53).

Recall that the efficiency (4.55) has been proven even elementwise in the form

$$\eta_\ell(T) \leq C_{\text{eff}} (\|\nabla(u - U_\ell)\|_{L^2(\omega_T)} + \text{osc}_\ell(\tilde{\omega}_T)) \quad \text{for all } T \in \mathcal{T}_\ell,$$

cf. Theorem 4.10. We now make the following assumptions on the refinement of marked elements. Whereas the first assumption is usually satisfied, the second assumption is stronger than what is usually done in practice. The combination of both assumption is called **inner node property** in the literature:

- For each marked element $T \in \mathcal{M}_\ell$, all edges of T are halved, i.e., the midpoints z_E of all edges $E \in \mathcal{E}_T$ belong to $\mathcal{N}_{\ell+1}$.
- For each marked element $T \in \mathcal{M}_\ell$, there is at least one node $z_T \in \mathcal{N}_{\ell+1}$ which lies in the interior of T , cf. Figure 4.11.

Exercise 62. Prove that under the foregoing assumptions, there holds the **discrete local efficiency**

$$\eta_\ell(T) \leq C_{\text{eff}} (\|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\omega_T)}^2 + \text{osc}_\ell(\tilde{\omega}_T)^2)^{1/2} \quad \text{for all } T \in \mathcal{M}_\ell, \quad (4.58)$$

where the exact solution u is replaced by the Galerkin solution $U_{\ell+1}$. The constant depends only on γ -shape regularity of $\mathcal{T}_{\ell+1}$. **Hint:** Replace the element bubble b_T (resp. the edge bubbles b_E) in the proof of Theorem 4.10 by the hat function for the node z_T (resp. z_E) with respect to the refined mesh $\mathcal{T}_{\ell+1}$. \square

Theorem 4.14 (Dörfler '96; Morin, Nochetto & Siebert '00). *Suppose that the set of marked elements \mathcal{M}_ℓ satisfies the Dörfler marking (4.57) for some arbitrary but fixed $\theta \in (0, 1)$. Under the inner node property, there is a constant $q \in (0, 1)$ such that*

$$\|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 \leq q \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \text{osc}_\ell^2 \quad \text{for all } \ell \in \mathbb{N}, \quad (4.59)$$

*which is called **error reduction property** in the literature. The contraction constant q depends only on Ω , the γ -shape regularity of \mathcal{T}_ℓ and $\mathcal{T}_{\ell+1}$, and the marking parameter θ . Moreover, assume vanishing oscillations in the sense that $\lim_{\ell \rightarrow \infty} \text{osc}_\ell^2 = 0$. Then, there holds $\lim_{\ell \rightarrow \infty} U_\ell = u$.*

Proof. Note that the Galerkin orthogonality implies

$$\|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 = \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2.$$

According to reliability and Dörfler marking, we infer

$$\begin{aligned} \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 &= \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 - \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 \\ &\leq C_{\text{rel}}^2 \eta_\ell^2 - \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 \\ &\leq C_{\text{rel}}^2 \theta^{-1} \eta_\ell(\mathcal{M}_\ell)^2 - \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2. \end{aligned}$$

Discrete local efficiency (4.58) and finite overlap of the patches ω_T prove

$$\begin{aligned} \eta_\ell(\mathcal{M}_\ell)^2 &\leq C_{\text{eff}}^2 \sum_{T \in \mathcal{M}_\ell} (\|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\omega_T)}^2 + \text{osc}_\ell(\tilde{\omega}_T)^2) \\ &\leq 4C_{\text{eff}}^2 (\|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 + \text{osc}_\ell^2) \end{aligned}$$

With $C := 4C_{\text{eff}}^2 C_{\text{rel}}^2 \theta^{-1}$ and the Galerkin orthogonality, we obtain

$$\begin{aligned} \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 &\leq C (\|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 + \text{osc}_\ell^2) - \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 \\ &= (C - 1) \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 + C \text{osc}_\ell^2 \\ &= (C - 1) \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 - (C - 1) \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + C \text{osc}_\ell^2. \end{aligned}$$

Collecting $\|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2$ on the left-hand side, we conclude

$$\|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 \leq \frac{C - 1}{C} \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \text{osc}_\ell^2$$

Since $(C - 1)/C = (1 - C^{-1}) < 1$, this yields (4.59). To see that vanishing oscillations imply convergence, let $\alpha_\ell := \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2$. Note that $\alpha_\ell \geq 0$ is bounded. From (4.59), we obtain

$$0 \leq M := \limsup_{\ell \rightarrow \infty} \alpha_\ell = \limsup_{\ell \rightarrow \infty} \alpha_{\ell+1} \leq q \limsup_{\ell \rightarrow \infty} \alpha_\ell + \limsup_{\ell \rightarrow \infty} \text{osc}_\ell = q M.$$

From this, we infer $M = 0$. Hence, $0 \leq \liminf_\ell \alpha_\ell \leq \limsup_\ell \alpha_\ell = 0$ and basic calculus yields $\alpha_\ell \rightarrow 0$ as $\ell \rightarrow \infty$. \blacksquare

Remark. For the convenience of the reader, we collect the arguments of the proof of Theorem 4.14. The proof uses Galerkin orthogonality, reliability of η_ℓ , local discrete efficiency of η_ℓ , and the Dörfler

marking. Generalizations of the proof were done by VEESER (2002, nonlinear model problem, convergence as consequence of energy decay), and CARSTENSEN & HOPPE (2006, for nonconforming and mixed FEM). \square

Remark. As a consequence of Theorem 4.14, Algorithm 4.11 can be extended by some additional marking to ensure vanishing oscillations. For instance, one could fix a contractive constant $\rho \in (0, 1)$. Having obtained $\mathcal{T}_{\ell+1}$ by one step of Algorithm 4.11, we check whether there holds $\text{osc}_{\ell+1} \leq \rho \text{osc}_\ell$ or not. If not, we refine further elements of $\mathcal{T}_{\ell+1}$, where the additional marking is steered by the local oscillations $\text{osc}_\ell(T)$. \square

Remark. DÖRFLER & NOCHETTO (2002) proved that small data oscillations even imply the saturation assumption

$$\|\nabla(u - u_{h/2})\|_{L^2(\Omega)} \leq q \|\nabla(u - u_h)\|_{L^2(\Omega)} \quad (4.60)$$

for some fixed contractive constant $q \in (0, 1)$. Here, $u_h \in \mathcal{S}_D^1(\mathcal{T}_h)$ and $u_{h/2} \in \mathcal{S}_D^1(\mathcal{T}_{h/2})$ are Galerkin solutions with respect to a mesh \mathcal{T}_h and its uniform refinement $\mathcal{T}_{h/2}$. The analysis of DÖRFLER and NOCHETTO is essentially based on the same ideas as above. However, instead of an element $T \in \mathcal{T}_h$, they consider the patch $\tilde{\omega}_z$ of a node (the so-called star). \square

v16:
01.12.2017

4.5.2 Convergence without Inner Node Property

The inner node property is very annoying in the sense that it seemed to be unnecessary in practice. Moreover, the local discrete efficiency seems to be valid only for residual-based error estimators. The following convergence theorem is a result of CASCÓN, KREUZER, NOCHETTO & SIEBERT from 2008, where it is proven that the combined error quantity, which consists of error and error estimator, has a contraction property. We stress two important observations:

- For their analysis, CASCÓN, KREUZER, NOCHETTO, and SIEBERT re-define the **mesh width**

$$h_T := |T|^{1/2} \quad \text{for } T \in \mathcal{T}_\ell, \quad (4.61)$$

whereas we considered $\text{diam}(T)$ before. Note that, however, $|T| \leq \text{diam}(T)^2 \leq 2\sigma(\mathcal{T})|T|$ so that both definition are equivalent for shape regular meshes, and we shall use the new definition in what follows.

- If $T \in \mathcal{T}_\ell$ is refined, each son $T' \in \mathcal{T}_{\ell+1}$ satisfies at least $|T'| \leq |T|/2$, which now results in a strict reduction $h_{T'} \leq h_T/\sqrt{2}$ of the local mesh-width (which fails, in general, for the usual definition $h_T = \text{diam}(T)$). This observation is used in step 2 of the proof of the following theorem.

We note that the analysis holds for general symmetric problems. For non-symmetric problems, the corresponding result has been open until FEISCHL, FÜHRER & PRAETORIUS (2014).

Theorem 4.15 (Cascón, Kreuzer, Nochetto & Siebert '08). *Suppose that the set of marked elements \mathcal{M}_ℓ satisfies the Dörfler marking for some fixed $\theta \in (0, 1)$. Then, there are constants $\kappa > 0$ and $q \in (0, 1)$ which depend only on θ and uniform γ -shape regularity of \mathcal{T}_ℓ for*

all $\ell \in \mathbb{N}_0$, such that

$$(\|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \kappa \eta_{\ell+1}^2)^{1/2} \leq \varrho (\|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \kappa \eta_\ell^2)^{1/2} \text{ for all } \ell \in \mathbb{N}. \quad (4.62)$$

In particular, this implies convergence $\lim_{\ell \rightarrow \infty} \|\nabla(u - U_\ell)\|_{L^2(\Omega)} = 0 = \lim_{\ell \rightarrow \infty} \eta_\ell$.

Proof. 1. step. There holds the following quasi-triangle inequality for the error estimator

$$\eta_\ell(V) \leq \eta_\ell(W) + C_\Delta \|\nabla(V - W)\|_{L^2(\Omega)} \quad \text{for all } V, W \in \mathcal{S}_D^1(\mathcal{T}_\ell) \quad (4.63)$$

with some constant $C_\Delta > 0$ which depends only on $\sigma(\mathcal{T}_\ell)$: From the triangle inequalities in ℓ_2 and L^2 , we infer

$$\begin{aligned} \eta_\ell(V) &= \left[\|h_\ell f\|_{L^2(\Omega)}^2 + \sum_{T \in \mathcal{T}_\ell} h_T (\|[\partial_n V]\|_{L^2(\partial T \cap \Omega)}^2 + \|\phi - \partial_n V\|_{L^2(\partial T \cap \Gamma_N)}^2) \right]^{1/2} \\ &\leq \left[\|h_\ell f\|_{L^2(\Omega)}^2 + \sum_{T \in \mathcal{T}_\ell} h_T (\|[\partial_n W]\|_{L^2(\partial T \cap \Omega)}^2 + \|\phi - \partial_n W\|_{L^2(\partial T \cap \Gamma_N)}^2) \right]^{1/2} \\ &\quad + \left[\sum_{T \in \mathcal{T}_\ell} h_T (\|[\partial_n(V - W)]\|_{L^2(\partial T \cap \Omega)}^2 + \|\partial_n(V - W)\|_{L^2(\partial T \cap \Gamma_N)}^2) \right]^{1/2}. \end{aligned}$$

For fixed $T \in \mathcal{T}_\ell$ and $E \in \mathcal{E}_T$, a scaling argument proves

$$h_T (\|[\partial_n(V - W)]\|_{L^2(E \cap \Omega)}^2 + \|\partial_n(V - W)\|_{L^2(E \cap \Gamma_N)}^2) \lesssim \|\nabla(V - W)\|_{L^2(\omega_E)},$$

where the constant depends only on $\sigma(\mathcal{T}_\ell)$. Consequently, we end up with (4.63).

2. step. There holds an estimator reduction in the sense that there is a constant $\varrho \in (0, 1)$ with

$$\eta_{\ell+1}^2 \leq (1 + \delta)\varrho \eta_\ell^2 + C_\delta \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 \quad \text{for all } \delta > 0, \quad (4.64)$$

where $C_\delta > 0$ depends only on δ and $C_\Delta > 0$. The constant ϱ depends only on θ and the reduction of the mesh-side on marked elements: Let $\Omega_* := \bigcup_{T \in \mathcal{M}_\ell} T$ denote the subdomain of Ω , where the elements are marked. Recall that $h_{T'} \leq h_T/\sqrt{2}$ for all sons $T' \in \mathcal{T}_{\ell+1}$ of a marked element $T \in \mathcal{M}_\ell$. The crucial step is to observe that the error indicators

$$\eta_\ell(T, V)^2 = h_T^2 \|f\|_{L^2(T)}^2 + h_T \|[\partial_n V]\|_{L^2(\partial T \cap \Omega)}^2 + h_T \|\phi - \partial_n V\|_{L^2(\partial T \cap \Gamma_N)}^2.$$

for $h_T = |T|^{1/2}$ lead to

$$\begin{aligned} \eta_{\ell+1}(U_\ell)^2 &= \sum_{\substack{T' \in \mathcal{T}_{\ell+1} \\ T' \subseteq \Omega_*}} \eta_{\ell+1}(T', U_\ell)^2 + \sum_{\substack{T' \in \mathcal{T}_{\ell+1} \\ T' \subseteq \Omega \setminus \Omega_*}} \eta_{\ell+1}(T', U_\ell)^2 \\ &\leq \frac{1}{\sqrt{2}} \sum_{\substack{T \in \mathcal{T}_\ell \\ T \subseteq \Omega_*}} \eta_\ell(T, U_\ell)^2 + \sum_{\substack{T \in \mathcal{T}_\ell \\ T \subseteq \Omega \setminus \Omega_*}} \eta_\ell(T, U_\ell)^2 \\ &= 2^{-1/2} \eta_\ell(\mathcal{M}_\ell)^2 + \eta_\ell(\mathcal{T}_\ell \setminus \mathcal{M}_\ell)^2 \\ &= (2^{-1/2} - 1) \eta_\ell(\mathcal{M}_\ell)^2 + \eta_\ell^2. \end{aligned}$$

By use of the Dörfler marking $\theta\eta_\ell^2 \leq \eta_\ell(\mathcal{M}_\ell)^2$, we thus obtain

$$\eta_{\ell+1}^2(U_\ell) \leq \eta_\ell^2 - (1 - 2^{-1/2})\eta_\ell(\mathcal{M}_\ell)^2 \leq \varrho\eta_\ell^2 \quad \text{with} \quad \varrho := (1 - \theta(1 - 2^{-1/2})).$$

Now, Young's inequality and step 1 conclude

$$\begin{aligned} \eta_{\ell+1}^2 &\leq (1 + \delta)\eta_{\ell+1}(U_\ell)^2 + (1 + \delta^{-1})C_\Delta^2 \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 \\ &\leq (1 + \delta)\varrho\eta_\ell^2 + (1 + \delta^{-1})C_\Delta^2 \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2. \end{aligned}$$

3. step. Proof of contraction property (4.62): Let $\kappa, \delta, \beta > 0$ be constants which are fixed later. Let $\varrho \in (0, 1)$ be the given constant from step 2. We recall the Galerkin orthogonality

$$\|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 = \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2$$

This and the estimator reduction imply

$$\begin{aligned} \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \kappa\eta_{\ell+1}^2 &= \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 - \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 + \kappa\eta_{\ell+1}^2 \\ &\leq \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + (\kappa C_\delta - 1) \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}^2 + \kappa(1 + \delta)\varrho\eta_\ell^2. \end{aligned}$$

Provided that $\kappa C_\delta \leq 1$, we infer

$$\begin{aligned} \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \kappa\eta_{\ell+1}^2 &\leq \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \kappa(1 + \delta)\varrho\eta_\ell^2 \\ &= \|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 - \kappa\beta\eta_\ell^2 + \kappa((1 + \delta)\varrho + \beta)\eta_\ell^2. \end{aligned}$$

Reliability $\|\nabla(u - U_\ell)\|_{L^2(\Omega)} \leq C_{\text{rel}}\eta_\ell$ finally leads to

$$\begin{aligned} \|\nabla(u - U_{\ell+1})\|_{L^2(\Omega)}^2 + \kappa\eta_{\ell+1}^2 &\leq (1 - \kappa\beta C_{\text{rel}}^{-2})\|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \kappa((1 + \delta)\varrho + \beta)\eta_\ell^2 \\ &\leq \max\{1 - \kappa\beta C_{\text{rel}}^{-2}, (1 + \delta)\varrho + \beta\} (\|\nabla(u - U_\ell)\|_{L^2(\Omega)}^2 + \kappa\eta_\ell^2). \end{aligned}$$

It remains to choose the constants κ, δ, β so that $q^2 := \max\{1 - \kappa\beta C_{\text{rel}}^{-2}, (1 + \delta)\varrho + \beta\} \in (0, 1)$:

- Choose $\delta > 0$ such that $(1 + \delta)\varrho < 1$.
- Choose $\kappa > 0$ such that $\kappa C_\delta \leq 1$.
- Choose $\beta > 0$ such that $(1 + \delta)\kappa + \beta < 1$.

This implies $q \in (0, 1)$ and concludes the proof. ■

Remark. We again collect the main arguments of the preceding proof, namely a certain quasi-triangle inequality of the estimator (4.63) and a strict reduction $\eta_{\ell+1}(\text{sons}(\mathcal{M}_\ell), U_\ell) \leq \kappa\eta_\ell(\mathcal{M}_\ell, U_\ell)$ with some $\kappa \in (0, 1)$ based on the strict reduction of the local mesh-width for marked elements. Besides this, only Galerkin orthogonality and Dörfler marking are used. Therefore, the proof works for a quite general class of symmetric problems and a variety of error estimators. The original work of CASCÓN, KREUZER, NOCHETTO & SIEBERT (2008) considers linear second order symmetric and elliptic problems in divergence form and H^1 -conforming finite element spaces with fixed polynomial degree. Finally, we stress that the proof also works for higher dimensions $d \geq 2$, where $h_T = |T|^{-1/d}$. For 2D, the usual definition $h_T := \text{diam}(T)$ is sufficient if marked elements are refined, e.g., by red-refinement or bisec(3), since then all edges are bisected. □

Exercise 63. Prove the following variants of Young's inequality, for all $a, b \in \mathbb{R}$ and $\delta > 0$,

- $ab \leq \frac{a^2}{2\delta} + \frac{\delta b^2}{2}$.
- $(a + b)^2 \leq (1 + \delta^{-1})a^2 + (1 + \delta)b^2$. □

Exercise 64. Prove that the estimator reduction (4.64) with $C_\delta = (1 + \delta^{-1})C_\Delta^2$ is equivalent to $\eta_{\ell+1} \leq \varrho \eta_\ell + C_\Delta \|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)}$. □

Exercise 65. Suppose that an error estimator η_ℓ satisfies the estimator reduction (4.64) and that the discrete spaces are nested, i.e., $\mathcal{S}_D^1(\mathcal{T}_\ell) \subseteq \mathcal{S}_D^1(\mathcal{T}_{\ell+1})$ for all $\ell \in \mathbb{N}_0$. Prove that there holds $\lim_{\ell \rightarrow \infty} \eta_\ell = 0$. **Hint:** Use that there always holds convergence $U_\ell \xrightarrow{\ell \rightarrow \infty} U_\infty$ so that $\|\nabla(U_{\ell+1} - U_\ell)\|_{L^2(\Omega)} \xrightarrow{\ell \rightarrow \infty} 0$, cf. Exercise 59. □

Exercise 66. Prove that adaptive FEM based on the residual error estimator with the usual definition of $h_T := \text{diam}(T)$ instead of (4.61) leads to R -linear convergence $\eta_{\ell+k} \leq C q^k \eta_\ell$ for all $k, \ell \in \mathbb{N}_0$. The constants $C > 0$ and $0 < q < 1$ depend only on θ and the uniform γ -shape regularity of \mathcal{T}_ℓ for all $\ell \in \mathbb{N}_0$. **Hint:** Use Theorem 4.15 and consider the Dörfler marking. □

4.5.3 Optimality of Adaptive FEM

In this section, we briefly comment on the optimality of Algorithm 4.11. Here, optimality is understood with respect to convergence $\|\nabla(u - U_\ell)\|_{L^2(\Omega)} = \mathcal{O}(N_\ell^{-\alpha})$, where N_ℓ denotes the number of elements in \mathcal{T}_ℓ and where $\alpha > 0$ denotes the convergence rate. The analysis goes back to BINEV, DAHMEN & DEVORE (2004) and STEVENSON (2005) who extended Algorithm 4.11 by some additional mesh-coarsening step. STEVENSON (2007) removed the coarsening and proved the same optimal rate $\alpha > 0$, but his proof still relied on the inner node property to ensure local discrete efficiency. The work by CASCÓN, KREUZER, NOCHETTO & SIEBERT (2008) removed the inner node property as has been shown in the previous section and generalized the theory from the 2D Poisson model problem to general linear second-order symmetric and elliptic PDEs. Finally, FEISCHL, FÜHRER & PRAETORIUS (2014) provided a frame to treat general linear second-order elliptic PDEs.

The following theorem is the main result and the current state of the art. Details on the precise assumptions as well as a proof will be given in a following lecture.

Theorem 4.16. *Employ newest vertex bisection for mesh-refinement. Let \mathcal{T}_0 be an initial triangulation, whose newest vertices are properly chosen. Let $\tilde{\mathcal{T}}_\ell$ be an arbitrary sequence of successively refined triangulations, i.e., $\tilde{\mathcal{T}}_{\ell+1} = \text{refine}(\tilde{\mathcal{T}}_\ell, \tilde{\mathcal{M}}_\ell)$ with arbitrary $\tilde{\mathcal{M}}_\ell \subseteq \tilde{\mathcal{T}}_\ell$ for all $\ell \in \mathbb{N}_0$ and $\tilde{\mathcal{T}}_0 := \mathcal{T}_0$. Let $s > 0$ and suppose that the corresponding error estimator $\tilde{\eta}_\ell$ decays with rate s , i.e.,*

$$\sup_{\ell \in \mathbb{N}_0} (\#\tilde{\mathcal{T}}_\ell)^s \tilde{\eta}_\ell < \infty. \tag{4.65}$$

Let $C_{\text{eff}}, C_{\text{rel}} > 0$ denote the constants of efficiency and (discrete local) reliability of the error estimator. Fix a parameter $\theta \in (0, 1)$ with

$$\theta < \theta_* := \frac{1}{1 + C_{\text{rel}}^2 C_{\text{eff}}^2} > 0 \quad (4.66)$$

and let the sets \mathcal{M}_ℓ of marked elements satisfy the Dörfler marking $\theta \eta_\ell^2 \leq \eta_\ell(\mathcal{M}_\ell)^2$ with (up to a fixed multiplicative constant) minimal cardinality. Then, it follows that

$$\sup_{\ell \in \mathbb{N}_0} (\#\mathcal{T}_\ell)^s \eta_\ell < \infty, \quad (4.67)$$

i.e., the adaptive algorithm guarantees (at least) the same algebraic convergence rate $s > 0$. In particular, the triangulations \mathcal{T}_ℓ are (asymptotically) optimal with respect to the degrees of freedom. ■

In practice, the assembly of the Galerkin data, the computation of the discrete solution, the computation of all refinement indicators, and the (local) mesh-refinement can be done in linear complexity $\mathcal{O}(\#\mathcal{T}_\ell)$, for each step $\ell \in \mathbb{N}_0$ of the adaptive algorithm. To determine a set of minimal cardinality $\mathcal{M}_\ell \subseteq \mathcal{T}_\ell$ which satisfies the Dörfler marking (4.49), we need to sort the refinement indicators. This would require $\mathcal{O}(\#\mathcal{T}_\ell \log(\#\mathcal{T}_\ell))$ operations. Based on binning, STEVENSON (2007) proposed a realization of the Dörfler marking (4.49) which only requires $\mathcal{O}(\#\mathcal{T}_\ell)$ operations, and the set \mathcal{M}_ℓ has minimal cardinality up to a factor 2. Then, each step of the adaptive algorithm has linear complexity $\mathcal{O}(\#\mathcal{T}_\ell)$.

Since the adaptive mesh U_ℓ cannot be computed without the full adaptive history $U_j \in \mathcal{S}_D^1(\mathcal{T}_j)$ for all $j = 0, \dots, \ell$, the overall computational cost (and hence the computational time) for the ℓ -th step result in $\mathcal{O}(\sum_{j=0}^{\ell} \#\mathcal{T}_j)$ operations. The following corollary from the PhD thesis of FEISCHL (2015) proves that the adaptive algorithm is even optimal with respect to the computational cost. We note that the assumptions are the same as for Theorem 4.16, but the claim is stronger.

Corollary 4.17 (Feischl '15). *Employ newest vertex bisection for mesh-refinement. Let \mathcal{T}_0 be an initial triangulation, whose newest vertices are properly chosen. Let $\tilde{\mathcal{T}}_\ell$ be an arbitrary sequence of successively refined triangulations, i.e., $\tilde{\mathcal{T}}_{\ell+1} = \text{refine}(\tilde{\mathcal{T}}_\ell, \tilde{\mathcal{M}}_\ell)$ with arbitrary $\tilde{\mathcal{M}}_\ell \subseteq \tilde{\mathcal{T}}_\ell$ for all $\ell \in \mathbb{N}_0$ and $\tilde{\mathcal{T}}_0 := \mathcal{T}_0$. Let $s > 0$ and suppose that the corresponding error estimator $\tilde{\eta}_\ell$ decays with rate s , i.e.,*

$$\sup_{\ell \in \mathbb{N}_0} (\#\tilde{\mathcal{T}}_\ell)^s \tilde{\eta}_\ell < \infty. \quad (4.68)$$

Let $C_{\text{eff}}, C_{\text{rel}} > 0$ denote the constants of efficiency and (discrete local) reliability of the error estimator. Fix a parameter $\theta \in (0, 1)$ with

$$\theta < \theta_* := \frac{1}{1 + C_{\text{rel}}^2 C_{\text{eff}}^2} > 0 \quad (4.69)$$

and let the sets \mathcal{M}_ℓ of marked elements satisfy the Dörfler marking $\theta \eta_\ell^2 \leq \eta_\ell(\mathcal{M}_\ell)^2$ with (up

to a fixed multiplicative constant) minimal cardinality. Then, it follows that

$$\sup_{\ell \in \mathbb{N}_0} \left(\sum_{j=0}^{\ell} \#\mathcal{T}_j \right)^s \eta_{\ell} < \infty, \quad (4.70)$$

i.e., the adaptive algorithm guarantees (at least) the same algebraic convergence rate $s > 0$ with respect to the computational cost. In particular, the triangulations \mathcal{T}_{ℓ} are (asymptotically) optimal with respect to the computational cost.

Proof. From Theorem 4.16 it follows that

$$\sup_{\ell \in \mathbb{N}_0} (\#\mathcal{T}_{\ell})^s \eta_{\ell} < \infty$$

and hence

$$\#\mathcal{T}_{\ell} \lesssim \eta_{\ell}^{-1/s} \quad \text{for all } \ell \in \mathbb{N}_0.$$

From Theorem 4.15, we obtain $0 < q < 1$ as well as linear convergence

$$\begin{aligned} \eta_{\ell+n} &\simeq (\|\nabla(u - U_{\ell+n})\|_{L^2(\Omega)}^2 + \kappa \eta_{\ell+n}^2)^{1/2} \leq q^n (\|\nabla(u - U_{\ell})\|_{L^2(\Omega)}^2 + \kappa \eta_{\ell}^2)^{1/2} \\ &\simeq q^n \eta_{\ell} \quad \text{for all } \ell, n \in \mathbb{N}_0. \end{aligned}$$

Hence, the geometric series implies that

$$\sum_{j=0}^{\ell} \#\mathcal{T}_j \lesssim \sum_{j=0}^{\ell} \eta_j^{-1/s} \lesssim \sum_{j=0}^{\ell} (q^{1/s})^{\ell-j} \eta_{\ell}^{-1/s} \lesssim \eta_{\ell}^{-1/s} \quad \text{for all } \ell \in \mathbb{N}_0.$$

This concludes the proof. ■

Chapter 5

Mixed Problems

5.1 Abstract Analysis of Petrov-Galerkin Schemes (13.01.2014)

Recall that for a continuous linear operator $T \in L(X, Y)$, the **adjoint operator** $T^* : Y^* \rightarrow X^*$ is formally defined by

$$T^*y^* \in X^* \quad \text{with} \quad (T^*y^*)(x) := y^*(Tx) \quad \text{for all } y^* \in Y^* \text{ and } x \in X. \quad (5.1)$$

It is an easy application of the Hahn-Banach extension theorem that $T^* \in L(Y^*, X^*)$ even with the same operator norm $\|T\| = \|T^*\|$. We start this section with some easy, but later on important, observations.

Lemma 5.1. *Let X and Y be normed spaces and $T \in L(X, Y)$. Then, T is an isomorphism between X and $\text{range}(T)$ if and only if*

$$\tau := \inf_{x \in X \setminus \{0\}} \frac{\|Tx\|_Y}{\|x\|_X} > 0. \quad (5.2)$$

In this case, there holds $\|T^{-1} : \text{range}(T) \rightarrow X\| = 1/\tau$. Moreover, the $\text{range}(T)$ is closed provided that X is a Banach space.

Proof. Clearly, $T^{-1} : \text{range}(T) \rightarrow X$ is well-defined (and hence an isomorphism in the sense of Linear Algebra) if and only if T is injective. If T is not injective, there exists some $x \neq 0$ with $Tx = 0$, and hence it follows $\tau = 0$. In particular, $\tau > 0$ implies that T is injective. By elementary calculations, we see

$$\begin{aligned} \tau &= \inf_{x \in X \setminus \{0\}} \frac{\|Tx\|_Y}{\|x\|_X} = \inf_{y \in \text{range}(T) \setminus \{0\}} \frac{\|y\|_Y}{\|T^{-1}y\|_X} = \frac{1}{\sup_{y \in \text{range}(T) \setminus \{0\}} \frac{\|T^{-1}y\|_X}{\|y\|_Y}} \\ &= \frac{1}{\|T^{-1} : \text{range}(T) \rightarrow X\|} \end{aligned}$$

Hence, $\tau > 0$ implies $\|T^{-1} : \text{range}(T) \rightarrow X\| = 1/\tau < \infty$, i.e., T^{-1} is even continuous. The same calculation proves that well-posedness and continuity of T^{-1} imply $\tau > 0$. Finally, suppose that X

is a Banach space and $\tau > 0$. Then, $\text{range}(T)$ is a Banach space as well and hence, in particular, a closed subspace of Y . ■

Exercise 67. For each operator $T \in L(X, Y)$ between normed spaces X and Y holds

$$\overline{\text{range}(T)} = (\ker T^*)_{\circ} := \{y \in Y \mid \forall y^* \in \ker T^* \quad y^*(y) = 0\}. \quad (5.3)$$

Hint: The inclusion $\text{range}(T) \subseteq (\ker T^*)_{\circ}$ can be shown directly, which leads to $\overline{\text{range}(T)} \subseteq (\ker T^*)_{\circ} = (\ker T^*)_{\circ}$. The converse inclusion follows by use of the Hahn-Banach separation theorem. □

According to the Hahn-Banach extension theorem, the **Hahn-Banach embedding**

$$I_X : X \rightarrow X^{**}, \quad (I_X x)(x^*) := x^*(x) \quad \text{for } x \in X \text{ and } x^* \in X^* \quad (5.4)$$

is an isometric linear operator, whence injective and continuous. A normed space X is **reflexive** provided that I_X is also surjective and thus an isometric isomorphism between X and X^{**} . We stress that

- reflexive spaces are, in particular, complete and thus Banach spaces,
- finite dimensional spaces are reflexive,
- all Hilbert spaces are reflexive,
- closed subspaces of reflexive spaces are also reflexive.

All of these facts are simple exercises left to the reader.

Theorem 5.2. Let X and Y be reflexive Banach spaces over \mathbb{R} , and $T \in L(X, Y^*)$. Then, T is an isomorphism if and only if the following two conditions hold:

- **inf-sup condition** $\tau := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{(Tx)(y)}{\|x\|_X \|y\|_Y} > 0$,
- **non-degeneracy condition** $\forall y \in Y \setminus \{0\} \exists x \in X \quad (Tx)(y) \neq 0$.

In this case, there holds $\|T^{-1}\| = 1/\tau$ for the operator norm of the inverse. The combination of inf-sup condition and non-degeneracy condition is called **LBB condition** in the literature, named after Ladyshenskaja, Babuška, and Brezzi.

Proof. 1. step. According to Lemma 5.1, $\tau > 0$ is equivalent to $T : X \rightarrow \text{range}(T)$ being an isomorphism with closed range. According to Exercise 67, it thus follows $\text{range}(T) = \overline{\text{range}(T)} = (\ker T^*)_{\circ}$. Hence, it only remains to show that the non-degeneracy condition is equivalent to $\ker T^* = \{0\}$ and hence $\text{range}(T) = (\ker T^*)_{\circ} = Y^*$: Recall that $T^* \in L(Y^{**}, X^*)$. According to reflexivity of Y , each $y^{**} \in Y^{**}$ allows for some $y \in Y$ with $y^{**} = I_Y y$. For all $x \in X$, it holds

$$(Tx)(y) = (I_Y y)(Tx) = y^{**}(Tx) = (T^* y^{**})(x).$$

2. step. We first show that the non-degeneracy condition (ND) implies $\ker(T^*) = \{0\}$: Let $y^{**} \in Y^{**}$ with $T^*y^{**} = 0$. Choose $y \in Y$ with $i_Y y = y^{**}$. For all $x \in X$, it follows $0 = (T^*y^{**})(x) = (Tx)(y)$. Because of ND, this implies $y = 0$ and hence $y^{**} = i_Y y = 0$.

3. step. To show that $\ker(T^*) = \{0\}$ implies ND, we show the contraposition: Suppose that ND fails. Then, there exists $y \in Y \setminus \{0\}$ such that $0 = (Tx)(y)$ for all $x \in X$. Define $y^{**} = i_Y y \neq 0$. Then, $0 = (Tx)(y) = (T^*y^{**})(x)$ proves that $T^*y^{**} = 0$, i.e., $\ker(T^*) \neq \{0\}$. ■

The following simple exercise proves that the assumptions on X in Theorem 5.2 are sharp.

Exercise 68. Let X be a normed space and Y be a reflexive Banach space over \mathbb{R} . Let $T \in L(X, Y^*)$ be an isomorphism. Prove that X is also a reflexive Banach space. **Hint:** It is known that a Banach space Z is reflexive, if and only if Z^* is reflexive. Moreover, Z is reflexive, if and only if each bounded sequence has a weakly convergent subsequence (i.e., the unit ball of Z is weakly compact). □

We now turn to continuous bilinear forms $a : X \times Y \rightarrow \mathbb{R}$ on normed spaces X and Y . So far, we only considered weak formulations of the type: Find $x \in X$ such that

$$a(x, \cdot) = x^* \in X^*, \tag{5.5}$$

where $a(\cdot, \cdot)$ is a continuous bilinear form on $X = Y$. For the classical Galerkin scheme, we assumed that $a(\cdot, \cdot)$ is even elliptic. Note that the last theorem provides a mathematical framework for weak formulations of the following type: Find $x \in X$ such that

$$a(x, \cdot) = y^* \in Y^*, \tag{5.6}$$

where $a(\cdot, \cdot)$ now is a continuous bilinear form $a : X \times Y \rightarrow \mathbb{R}$. In the literature, this approach is named after Petrov-Galerkin.

Corollary 5.3. Let X and Y be real Banach spaces, where Y is reflexive. Let $a : X \times Y \rightarrow \mathbb{R}$ be bilinear and continuous. Then, the following statements (i)–(ii) are equivalent:

- (i) For each $y^* \in Y^*$, exists a unique $x \in X$ with $a(x, \cdot) = y^*$.
- (ii) The bilinear form satisfies the **LBB condition**:

- **inf-sup condition** $\alpha := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{a(x, y)}{\|x\|_X \|y\|_Y} > 0$,
- **non-degeneracy condition** $\forall y \in Y \setminus \{0\} \exists x \in X \quad a(x, y) \neq 0$.

In this case, it holds

$$\alpha \|x\|_X \leq \|y^*\|_{Y^*} \leq \|a\| \|x\|_X, \tag{5.7}$$

where $\|a\| := \sup_{\substack{x \in X \setminus \{0\} \\ y \in Y \setminus \{0\}}} \frac{a(x, y)}{\|x\|_X \|y\|_Y}$ denotes the continuity bound of $a(\cdot, \cdot)$.

Proof. We associate with $a(\cdot, \cdot)$ the operator $T \in L(X, Y^*)$ given by $Tx = a(x, \cdot)$. Note that (i) is equivalent to the fact that T is an isomorphism. According to Theorem 5.2, the latter is characterized by the LBB condition for T which, in fact, coincides with that for $a(\cdot, \cdot)$. For given $y^* \in Y^*$ and $x \in X$ with $a(x, \cdot) = y^* \in Y^*$, it holds $Tx = y^*$. With $\|T : X \rightarrow Y^*\| = \|a\|$, we see $\|y^*\|_{Y^*} \leq \|a\| \|x\|_X$. With $x = T^{-1}y^*$ and $\|T^{-1} : Y^* \rightarrow X\| \leq 1/\alpha$, we derive $\|x\|_X \leq \|y^*\|_{Y^*}/\alpha$. This concludes the proof. \blacksquare

One important difference to the elliptic framework now is, that we may not simply replace X and Y by discrete spaces X_h and Y_h , respectively. Instead, Corollary 5.3 states that we need to satisfy the inf-sup condition and the non-degeneracy condition not only for the pairing (X, Y) of continuous spaces, but also for any pairing (X_h, Y_h) of discrete spaces. To underline this, note that

$$T = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$$

is an isomorphism on $Y = X = \mathbb{R}^3$. For $X_h = Y_h = \mathbb{R}^2$ and the canonical embedding, i.e., $x \in \mathbb{R}^2$ is identified with $(x, 0) \in \mathbb{R}^3$, the restricted matrix is

$$T_h = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

which is clearly singular. We finally note that in the discrete setting the inf-sup condition and the non-degeneracy condition are equivalent.

Proposition 5.4. *Let X and Y be real Banach spaces with $\dim X < \infty$ and $\dim Y < \infty$. Let $a : X \times Y \rightarrow \mathbb{R}$ be bilinear. Then, there holds the following:*

- (i) *The inf-sup condition $\alpha := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{a(x, y)}{\|x\|_X \|y\|_Y} > 0$ implies $\dim X \leq \dim Y$.*
- (ii) *The non-degeneracy condition $(\forall y \in Y \setminus \{0\}) \exists x \in X \ a(x, y) \neq 0$ implies $\dim Y \leq \dim X$.*
- (iii) *For $\dim X = \dim Y$, the inf-sup condition is satisfied if and only if the non-degeneracy condition is satisfied.*

Proof. We define the operators $A_1 \in L(X, Y^*)$ and $A_2 \in L(Y, X^*)$ by $A_1x := a(x, \cdot)$ and $A_2y := a(\cdot, y)$. According to Linear Algebra, finite dimension implies

$$\begin{aligned} \dim X &= \dim \ker(A_1) + \dim \text{range}(A_1) \leq \dim \ker(A_1) + \dim Y^* = \dim \ker(A_1) + \dim Y, \\ \dim Y &= \dim \ker(A_2) + \dim \text{range}(A_2) \leq \dim \ker(A_2) + \dim X^* = \dim \ker(A_2) + \dim X. \end{aligned}$$

1. step. If $\dim X > \dim Y$, we obtain $\dim \ker(A_1) > 0$. Hence, there exists $x \in X \setminus \{0\}$ with $A_1x = 0$. This implies $a(x, y) = 0$ for all $y \in Y$ and hence $\alpha = 0$ for the inf-sup constant. By contraposition, this shows that the inf-sup condition implies $\dim \ker(A_1) = 0$ and hence $\dim X \leq \dim Y$. This proves (i).

2. step. If $\dim Y > \dim X$, we obtain $\dim \ker(A_2) > 0$. Hence, there exists $y \in Y \setminus \{0\}$ with $A_2y = 0$. This implies $a(x, y) = 0$ for all $x \in X$, and hence the non-degeneracy condition fails. By

contraposition, this shows that the non-degeneracy condition implies $\dim \ker(A_2) = 0$ and hence $\dim Y \leq \dim X$. This proves (ii).

3. step. For (i), we have shown that the inf-sup condition implies injectivity of A_1 . Since $\dim X = \dim Y = \dim Y^*$, this proves that A_1 is bijective. By finite dimension and hence compactness of the unit spheres of X and Y as well as continuity of $a(\cdot, \cdot)$, the converse implication is clear, i.e., A_1 is bijective if and only if the inf-sup condition holds. For (ii), we have shown that the non-degeneracy condition implies injectivity of A_2 . Since $\dim Y = \dim X = \dim X^*$, this proves that A_2 is bijective. The converse implication is obvious, i.e., A_2 is bijective if and only if the non-degeneracy conditions holds. To conclude (iii), we only have to show that bijectivity of A_1 and A_2 are equivalent. To that end, let $\{x_1, \dots, x_n\} \subset X$ and $\{y_1, \dots, y_n\} \subset Y$ be bases. We define the matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{A}_{jk} := a(x_k, y_j)$ and note that $(A_1 x_k)(y_j) = a(x_k, y_j) = \mathbf{A}_{jk}$ as well as $(A_1 y_j)(x_k) = a(x_k, y_j) = \mathbf{A}_{jk}$. Therefore, \mathbf{A} is the Petrov-Galerkin matrix corresponding to A_1 and its transpose \mathbf{A}^T is the Petrov-Galerkin matrix corresponding to A_2 . Therefore, Linear Algebra proves the equivalence

$$A_1 \text{ is bijective} \iff \mathbf{A} \text{ is regular} \iff \mathbf{A}^T \text{ is regular} \iff A_2 \text{ is bijective}$$

This concludes the proof. ■

Exercise 69. Prove that a bilinear form $a : X \times Y \rightarrow \mathbb{R}$ on normed spaces X and Y is continuous if and only if $\|a\| := \sup_{\substack{x \in X \setminus \{0\} \\ y \in Y \setminus \{0\}}} \frac{a(x, y)}{\|x\|_X \|y\|_Y} < \infty$. □

The following exercise states the quasi optimality of Petrov-Galerkin schemes. We stress, however, that the quasi-optimality constant depends on the discrete inf-sup condition.

Exercise 70 (Céa's Lemma for Petrov-Galerkin Schemes). We consider the weak form (5.6) with a continuous bilinear form $a : X \times Y \rightarrow \mathbb{R}$ on Banach spaces X and Y . Let $y^* \in Y^*$. Let X_h and Y_h be finite dimensional subspaces of X resp. Y with $\dim X_h = \dim Y_h$. We assume the

- **discrete inf-sup condition** $\alpha_h := \inf_{x_h \in X_h \setminus \{0\}} \sup_{y_h \in Y_h \setminus \{0\}} \frac{a(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y} > 0,$

Then, there is a unique $x_h \in X_h$ with

$$a(x_h, \cdot) = y^* \in Y_h^*. \tag{5.8}$$

If $x \in X$ solves the weak form (5.6), we have quasi optimality

$$\|x - x_h\|_X \leq (1 + \|a\|/\alpha_h) \min_{v_h \in X_h} \|x - v_h\|_X, \tag{5.9}$$

where $\|a\| := \sup_{\substack{x \in X \setminus \{0\} \\ y \in Y \setminus \{0\}}} \frac{a(x, y)}{\|x\|_X \|y\|_Y}$ denotes the continuity bound of $a(\cdot, \cdot)$. \square

A simple observation is that the LBB theory allows a generalization of the Lax-Milgram lemma to the case of reflexive Banach spaces.

Exercise 71 (Lax-Milgram Lemma for Reflexive Spaces). Let $a : X \times X \rightarrow \mathbb{R}$ be a continuous and elliptic bilinear form on the reflexive Banach space X . Prove that $a(\cdot, \cdot)$ satisfies the inf-sup condition

$$\tau := \inf_{x \in X \setminus \{0\}} \sup_{y \in X \setminus \{0\}} \frac{a(x, y)}{\|x\|_X \|y\|_X} > 0$$

as well as the non-degeneracy condition

$$\forall y \in X \setminus \{0\} \exists x \in X \quad a(x, y) \neq 0.$$

For each given right-hand side $x^* \in X^*$, the weak form (5.5) thus has a unique solution $x \in X$. \square

Another observation is that for reflexive spaces, it is immaterial whether the LBB condition is stated for the first or the second component.

Exercise 72. Let X, Y be reflexive Banach spaces and $a : X \times Y \rightarrow \mathbb{R}$ be a continuous bilinear form. Prove that the following statements (i)–(ii) are equivalent:

(i) The bilinear form satisfies the **LBB condition for the first argument**:

- $\alpha_1 := \inf_{x \in X \setminus \{0\}} \sup_{y \in Y \setminus \{0\}} \frac{a(x, y)}{\|x\|_X \|y\|_Y} > 0,$
- $\forall y \in Y \setminus \{0\} \exists x \in X \quad a(x, y) \neq 0.$

(ii) The bilinear form satisfies the **LBB condition for the second argument**:

- $\alpha_2 := \inf_{y \in Y \setminus \{0\}} \sup_{x \in X \setminus \{0\}} \frac{a(x, y)}{\|x\|_X \|y\|_Y} > 0,$
- $\forall x \in X \setminus \{0\} \exists y \in Y \quad a(x, y) \neq 0.$

Moreover, in this case there holds $\alpha_1 = \alpha_2$. \square

5.2 Abstract Analysis of Mixed Formulations (12.01.2015)

Instead of the general mixed formulation (5.6), we consider linear problems with side constraints in the following. These arise, for instance, for the Stokes problem. [Before we focus on the abstract solution theory, we explain why these problems are called *saddle point problems*: Plotting a function](#)

$f : \mathbb{R}^2 \rightarrow \mathbb{R}$ over the two-dimensional plane, we call a point (x, y) saddle point of f if the real function $f(x+t, y)$ has a minimum at $t=0$ and the function $f(x, y+t)$ has a maximum for $t=0$. This is, what is stated in the following proposition for the so-called *Lagrange functional*.

Proposition 5.5. *Let $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ be bilinear forms on normed spaces X and Y . Assume that $a(\cdot, \cdot)$ is positive semidefinite, i.e., $a(x, x) \geq 0$ and symmetric. Then, given $(x^*, y^*) \in X^* \times Y^*$, $(x, y) \in X \times Y$ is a solution of the saddle point problem*

$$\begin{aligned} a(x, \cdot) + b(\cdot, y) &= x^* \in X^* \\ b(x, \cdot) &= y^* \in Y^*. \end{aligned} \quad (5.10)$$

if and only if the Lagrange functional $\mathcal{L}(v, w) := \frac{1}{2} a(v, v) - x^(v) + b(v, w) - y^*(w)$ satisfies*

$$\mathcal{L}(x, w) \leq \mathcal{L}(x, y) \leq \mathcal{L}(v, y) \quad \text{for all } (v, w) \in X \times Y, \quad (5.11)$$

i.e., (x, y) is a saddle point of $\mathcal{L}(\cdot, \cdot)$. In this case, the first estimate in (5.11) holds with equality.

Proof. First, assume that $(x, y) \in X \times Y$ is a solution of the saddle point problem (5.10). For $w \in Y$, the second equality in (5.10) implies

$$\mathcal{L}(x, y) - \mathcal{L}(x, w) = b(x, y - w) - y^*(y - w) = 0.$$

This proves the lower estimate of (5.11) even with equality. For $v \in X$, symmetry of $a(\cdot, \cdot)$ and the first equality in (5.10) prove

$$\mathcal{L}(v, y) - \mathcal{L}(x, y) = \frac{1}{2} a(x - v, x - v) + \underbrace{a(x, v - x) - x^*(v - x) + b(v - x, y)}_{=0} \geq 0,$$

and we obtain the upper estimate. Altogether, (x, y) is a saddle point of the Lagrange functional. The proof of the converse implication follows from a classical argument from the calculus of variations: Let $(x, y) \in X \times Y$ satisfy (5.11). For fixed $v \in X$, the real function $f(t) := \mathcal{L}(x + tv, y)$ has a global minimum at $t=0$. There holds

$$f(t) = \frac{1}{2} a(x, x) - x^*(x) + b(x, y) - y^*(y) + \frac{t^2}{2} a(v, v) + t\{a(x, v) - x^*(v) + b(v, y)\}.$$

Hence $0 = f'(0) = a(x, v) - x^*(v) + b(v, y)$ for all $v \in X$. This proves the first equality in (5.10). To prove the second equality, consider, for fixed $w \in Y$, the real function $g(t) := \mathcal{L}(x, y + tw)$ which has a global maximum at $t=0$. There holds

$$g(t) = \frac{1}{2} a(x, x) - x^*(x) + b(x, y) - y^*(y) + t\{b(x, w) - y^*(w)\}$$

and thus $0 = g'(0) = b(x, w) - y^*(w)$ for all $w \in Y$, i.e., $b(x, \cdot) = y^* \in Y^*$. ■

The following theorem of Brezzi provides existence and uniqueness of the solution of saddle point problems.

Theorem 5.6 (Brezzi). *Let X be a Hilbert space and Y be a reflexive Banach space. Let $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ be continuous bilinear forms. We define $X_0 := \{x \in X \mid b(x, \cdot) = 0 \in Y^*\}$ and assume*

- $\alpha := \inf_{v \in X_0 \setminus \{0\}} \frac{a(v, v)}{\|v\|_X^2} > 0$, i.e., $a(\cdot, \cdot)$ is elliptic on X_0 ,
- $\beta := \inf_{y \in Y \setminus \{0\}} \sup_{x \in X \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y} > 0$.

Then, for any $(x^*, y^*) \in X^* \times Y^*$, there is a unique solution $(x, y) \in X \times Y$ of

$$\begin{aligned} a(x, \cdot) + b(\cdot, y) &= x^* \in X^* \\ b(x, \cdot) &= y^* \in Y^*. \end{aligned} \quad (5.12)$$

Moreover, we have the stability estimates

$$\|x\|_X \leq \frac{1}{\alpha} \|x^*\|_{X^*} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha}\right) \|y^*\|_{Y^*} \quad (5.13)$$

and

$$\|y\|_Y \leq \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha}\right) \left(\|x^*\|_{X^*} + \frac{\|a\|}{\beta} \|y^*\|_{Y^*}\right) \quad (5.14)$$

Remark. (i) Note that one can identify $X^* \times Y^* = (X \times Y)^*$ as follows: For $x^* \in X^*$ and $y^* \in Y^*$, the definition $z^*(x, y) := x^*(x) + y^*(y)$ yields $z^* \in (X \times Y)^*$. Conversely, $z^* \in (X \times Y)^*$ gives rise to $x^*(x) := z^*(x, 0)$ and $y^*(y) := z^*(0, y)$ with $(x^*, y^*) \in X^* \times Y^*$.

(ii) If we define operators $A_1 \in L(X, X^*)$, $B_1 \in L(X, Y^*)$, and $B_2 \in L(Y, X^*)$ by

$$A_1 x := a(x, \cdot), \quad B_1 x := b(x, \cdot), \quad \text{and} \quad B_2 y := b(\cdot, y),$$

Equation (5.12) can be written in the form

$$\begin{pmatrix} A_1 & B_2 \\ B_1 & 0 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x^* \\ y^* \end{pmatrix}. \quad (5.15)$$

In this form, the Brezzi theorem states that this operator matrix is an isomorphism from $X \times Y$ to $X^* \times Y^* = (X \times Y)^*$ and so fits into the abstract framework given above.

(iii) We stress that the original proof of Brezzi works for reflexive Banach spaces X and Y . Therein, it is proved directly that the operator matrix from (5.15) satisfies the inf-sup condition as well as the non-degeneracy condition. Our stronger assumption that X is not only a reflexive Banach space, but even a Hilbert space, reduces the technical difficulties and leads to a much simpler proof. \square

Sketch of Proof of Theorem 5.6. Let $(x, y) \in X \times Y$. With the orthogonal decomposition $X = X_0 \oplus X_0^\perp$, we write $x = x_1 + x_2$ with $x_1 \in X_0$ and $x_2 \in X_0^\perp$. Note that (5.12) is equivalent to the following three identities:

- $b(x_2, \cdot) = y^* \in Y^*$,
- $a(x_1, \cdot) = x^* - a(x_2, \cdot) \in X_0^*$,
- $b(\cdot, y) = x^* - a(x_1 + x_2, \cdot) \in X^*$.

For the proof of Theorem 5.6 we are going to show that these three equations — proved in the stated order — admit unique solutions $x_2 \in X_0^\perp$, $x_1 \in X_0$, and $y \in Y^*$. This proves existence and uniqueness of the solution $(x, y) = (x_1 + x_2, y) \in X \times Y$ of (5.12). ■

The main ingredient of the proof of Theorem 5.6 is the closed range theorem:

Theorem 5.7 (Banach's Closed Range Theorem). *For an operator $T \in L(X, Y)$ between Banach spaces X and Y , the following is equivalent:*

- (i) $\text{range}(T) \subseteq Y$ is closed,
- (ii) $\text{range}(T) = (\ker T^*)^\circ = \{y \in Y \mid \forall y^* \in \ker T^* \quad y^*(y) = 0\}$,
- (iii) $\text{range}(T^*) \subseteq X^*$ is closed,
- (iv) $\text{range}(T^*) = (\ker T)^\circ = \{x^* \in X^* \mid \forall x \in \ker T \quad x^*(x) = 0\}$. ■

Proof of Theorem 5.6. The essential steps of the proof are based on operator arguments for the operators defined by $B_1x := b(x, \cdot)$ and $B_2y := b(\cdot, y)$. We are going to consider the four operators

$$\begin{aligned} B_1 &\in L(X, Y^*), & B_1^* &\in L(Y^{**}, X^*), \\ B_2 &\in L(Y, X^*), & B_2^* &\in L(X^{**}, Y^*). \end{aligned}$$

More precisely, the first three steps state the essential observations about these operators, whereas the remaining proof follows the line of the sketch given before.

1. step. B_2 is injective with closed range and $\|B_2^{-1} : \text{range}(B_2) \rightarrow Y\| = 1/\beta$, which follows from Lemma 5.1 and

$$\beta = \inf_{y \in Y \setminus \{0\}} \frac{\|B_2y\|_{X^*}}{\|y\|_Y}.$$

2. step. There holds $B_2 = B_1^*I_Y$, which follows from

$$(B_2y)(x) = b(x, y) = (B_1x)(y) = (I_Yy)(B_1x) = (B_1^*I_Yy)(x) \quad \text{for all } x \in X, y \in Y.$$

3. step. Since Y is reflexive, B_1^* is injective with closed range $\text{range}(B_1^*) = \text{range}(B_2)$. Moreover, the closed range theorem even proves

$$\text{range}(B_2) = \text{range}(B_1^*) = (\ker B_1)^\circ = (X_0)^\circ \quad \text{as well as} \quad \text{range}(B_1) = (\ker B_1^*)^\circ = Y^*.$$

4. step. There is a unique $x_2 \in X_0^\perp$ with $b(x_2, \cdot) = y^* \in Y^*$: According to step 3, there is at least one $x \in X$ with $b(x, \cdot) = B_1x = y^*$. The decomposition $x = x_1 + x_2$ with $x_1 \in X_0$ and $x_2 \in X_0^\perp$ proves $b(x_2, \cdot) = b(x, \cdot) = y^* \in Y^*$, which concludes existence. To prove uniqueness, let $\tilde{x}_2 \in X_0^\perp$ with $b(\tilde{x}_2, \cdot) = y^* \in Y^*$. Then, $b(x_2 - \tilde{x}_2, \cdot) = 0 \in Y^*$, whence $x_2 - \tilde{x}_2 \in \ker B_1 = X_0$. From $x_2 - \tilde{x}_2 \in X_0^\perp$, we thus obtain $x_2 = \tilde{x}_2$.

5. step. There is a unique element $x_1 \in X_0$ with $a(x_1, \cdot) = x^* - a(x_2, \cdot) \in X_0^*$ which immediately follows from the Lax-Milgram lemma and the observation that X_0 is a closed subspace of a Hilbert space and hence a Hilbert space as well.

6. step. There is a unique element $y \in Y$ with $b(\cdot, y) = x^* - a(x, \cdot)$, where $x := x_1 + x_2 \in X$: By construction in step 5, there holds

$$x^* - a(x, \cdot) \in (X_0)^\circ = \{v^* \in X^* \mid \forall v \in X_0 \quad v^*(v) = 0\}.$$

According to step 1 and step 3, B_2 is injective with $\text{range}(B_2) = (X_0)^\circ$. Thus, there is a unique $y \in Y$ with $b(\cdot, y) = B_2 y = x^* - a(x, \cdot)$.

7. step. There holds $\|x_2\|_X \leq \|y^*\|_{Y^*}/\beta$: From $x_2 \in X_0^\perp$ follows $(x_2; \cdot)_X \in (X_0)^\circ = \text{range}(B_2)$. Thus, we may choose $\tilde{y} \in Y$ with $B_2 \tilde{y} = (x_2; \cdot)_X$. From $\|B_2^{-1} : (X_0)^\circ \rightarrow Y\| = 1/\beta$, we infer $\|\tilde{y}\|_Y \leq \|(x_2; \cdot)_X\|_{X^*}/\beta = \|x_2\|_X/\beta$. Together with $b(x_2, \cdot) = y^*$, we conclude

$$\|x_2\|_X^2 = (x_2; x_2)_X = (B_2 \tilde{y})(x_2) = b(x_2, \tilde{y}) = y^*(\tilde{y}) \leq \|y^*\|_{Y^*} \|\tilde{y}\|_Y \leq \frac{\|y^*\|_{Y^*}}{\beta} \|x_2\|_X.$$

8. step. There holds $\|x_1\|_X \leq \alpha^{-1} (\|x^*\|_{X^*} + \|a\| \|x_2\|_X)$: Note that $A_1 \in L(X_0, X_0^*)$ is an isomorphism with $\|A_1^{-1} : X_0^* \rightarrow X_0\| \leq 1/\alpha$. From $A_1 x_1 = a(x_1, \cdot) = x^* - a(x_2, \cdot)$, we thus infer

$$\|x_1\|_X \leq \frac{1}{\alpha} \|x^* - a(x_2, \cdot)\|_{X_0^*} \leq \frac{1}{\alpha} (\|x^*\|_{X^*} + \|a\| \|x_2\|_X).$$

9. step. The triangle inequality leads to

$$\|x\|_X \leq \|x_1\|_X + \|x_2\|_X \leq \frac{1}{\alpha} \|x^*\|_{X^*} + \left(\frac{\|a\|}{\alpha} + 1\right) \|x_2\|_X \leq \frac{1}{\alpha} \|x^*\|_{X^*} + \frac{1}{\beta} \left(\frac{\|a\|}{\alpha} + 1\right) \|y^*\|_{Y^*}.$$

10. step. It finally remains to dominate $\|y\|_Y$, where $B_2 y = b(\cdot, y) = x^* - a(x, \cdot) \in (X_0)^\circ$. We use $\|B_2^{-1} : (X_0)^\circ \rightarrow Y\| = 1/\beta$ to see

$$\begin{aligned} \|y\|_Y &\leq \frac{1}{\beta} \|x^* - a(x, \cdot)\|_{X_0^*} \leq \frac{1}{\beta} \|x^*\|_{X^*} + \frac{\|a\|}{\beta} \|x\|_X \\ &\leq \frac{1}{\beta} \|x^*\|_{X^*} + \frac{\|a\|}{\beta} \frac{1}{\alpha} \|x^*\|_{X^*} + \frac{\|a\|}{\beta^2} \left(1 + \frac{\|a\|}{\alpha}\right) \|y^*\|_{Y^*} \\ &= \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha}\right) \left(\|x^*\|_{X^*} + \frac{\|a\|}{\beta} \|y^*\|_{Y^*}\right). \end{aligned}$$

This concludes the proof. ■

Remark. (i) Let $B_1 \in L(X, Y^*)$ and $B_2 \in L(Y, X^*)$ be defined as in the proof of Theorem 5.6. In the proof, we have seen that $\beta > 0$ implies surjectivity of B_1 . We note that even the converse implication holds, i.e.,

$$\beta := \inf_{y \in Y \setminus \{0\}} \sup_{x \in X \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y} > 0 \iff B_1 \text{ is surjective.} \quad (5.16)$$

Suppose that B_1 is surjective. As in step 3 of the preceding proof, the closed range theorem proves that B_1^* is injective with closed range. Moreover, $B_2 = B_1^* I_Y$ proves that B_2 is injective with closed

$\text{range}(B_2) = \text{range}(B_1^*) = (\ker B_1)^\circ = (X_0)^\circ$, i.e., $B_2 : Y \rightarrow \text{range}(B_2)$ is continuous and bijective between the Banach spaces Y and $\text{range}(B_2) \subseteq X^*$. According to the open mapping theorem, $B_2 : Y \rightarrow \text{range}(B_2)$ even is an isomorphism, i.e., $\beta^{-1} = \|B_2 : Y \rightarrow \text{range}(B_2)\| < \infty$, whence $\beta > 0$.

(ii) Altogether, the two main assumptions on $a(\cdot, \cdot)$ and $b(\cdot, \cdot)$ can equivalently be stated as follows:

- The bilinear form $a(\cdot, \cdot)$ is elliptic on $X_0 = \ker B_1$.
- The operator $B_1 \in L(X, Y^*)$ is surjective.

We hope that the reader may keep this (abstract) formulation in mind much easier. For the statement of Theorem 5.6, we used the definition of α and β instead, to provide the stability estimates (5.13)–(5.14) with explicit constants. \square

Going through the proof of Theorem 5.6, one realizes that ellipticity of $a(\cdot, \cdot)$ on X_0 is only used to provide a unique $x_1 \in X_0$ with $a(x_1, \cdot) = x_0^* \in X_0^*$ in step 5. To prove unique existence of x_1 , it is, however, sufficient to assume that the operator $A_1 : X_0 \rightarrow X_0^*$ defined by $A_1 x := a(x, \cdot)$ is an isomorphism. This is done in the following exercise.

Exercise 73. Let $X, Y, a(\cdot, \cdot)$, and $b(\cdot, \cdot)$ be as in Theorem 5.6. Then, the following statements are equivalent:

- (i) For all $(x^*, y^*) \in X^* \times Y^*$, there exists a unique solution $(x, y) \in X \times Y$ of the saddle point problem (5.12).
- (ii) The bilinear forms $a(\cdot, \cdot)$ and $b(\cdot, \cdot)$ satisfy the following three assumptions:

- $\alpha := \inf_{v \in X_0 \setminus \{0\}} \sup_{w \in X_0 \setminus \{0\}} \frac{a(v, w)}{\|v\|_X \|w\|_X} > 0$,
- $\forall w \in X_0 \setminus \{0\} \exists v \in X_0 \quad a(v, w) \neq 0$,
- $\beta := \inf_{y \in Y \setminus \{0\}} \sup_{x \in X \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y} > 0$.

The first two assumptions state that $A_1 : X_0 \rightarrow X_0^*$ is an isomorphism, cf. Theorem 5.2. The assumption on β is the same as in the above statement of the Brezzi theorem. \square

The following corollary provides the relation between saddle point problems and the abstract Petrov-Galerkin scheme from Section 5.1.

Corollary 5.8. *Suppose that X is a Hilbert space, Y is a reflexive Banach space, and $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ are continuous bilinear forms. Then, $Z := X \times Y$ is a reflexive Banach space, and $c((x, y), (\tilde{x}, \tilde{y})) := a(x, \tilde{x}) + b(\tilde{x}, y) + b(x, \tilde{y})$ defines a continuous bilinear form $c : Z \times Z \rightarrow \mathbb{R}$. Moreover, for $(x, y) \in X \times Y$ and $(x^*, y^*) \in X^* \times Y^*$, the saddle point problem (5.12) is equivalent to*

$$c((x, y), (\tilde{x}, \tilde{y})) = x^*(\tilde{x}) + y^*(\tilde{y}) \quad \text{for all } (\tilde{x}, \tilde{y}) \in X \times Y. \quad (5.17)$$

Finally, the following three statements are equivalent:

(i) $a(\cdot, \cdot)$ and $b(\cdot, \cdot)$ satisfy the assumptions of the Brezzi theorem, i.e.,

- $\alpha := \inf_{v \in X_0 \setminus \{0\}} \sup_{w \in X_0 \setminus \{0\}} \frac{a(v, w)}{\|v\|_X \|w\|_X} > 0$ with $X_0 := \{x \in X \mid b(x, \cdot) = 0 \in Y^*\}$,
- $\forall w \in X_0 \setminus \{0\} \exists v \in X_0 \quad a(v, w) \neq 0$,
- $\beta := \inf_{y \in Y \setminus \{0\}} \sup_{x \in X \setminus \{0\}} \frac{b(x, y)}{\|x\|_X \|y\|_Y} > 0$.

(ii) $c(\cdot, \cdot)$ satisfies the LBB conditions

- $\gamma := \inf_{z \in Z \setminus \{0\}} \sup_{w \in Z \setminus \{0\}} \frac{c(z, w)}{\|z\|_Z \|w\|_Z} > 0$,
- $\forall w \in Z \setminus \{0\} \exists z \in Z \quad c(z, w) \neq 0$.

(iii) For all $(x^*, y^*) \in X^* \times Y$, the variational formulation (5.17) has a unique solution $(x, y) \in X \times Y$.

In particular, it holds $\|c\| \leq \|a\| + 2\|b\|$ for the corresponding norms and there exists a constant $C > 0$ such that

$$\gamma \geq C \left[\frac{1}{\alpha} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \left(1 + \frac{\|a\|}{\beta} \right) \right]^{-1}. \quad (5.18)$$

Proof. 1. step. Since X and Y are reflexive, their closed unit balls $B_X \subset X$ and $B_Y \subset Y$ are weakly compact. According to the Tychonov theorem, $B_X \times B_Y$ and hence B_Z are weakly compact as well. Consequently, Z is reflexive. Moreover, it is obvious that $c(\cdot, \cdot)$ is bilinear and continuous with $\|c\| \leq \|a\| + 2\|b\|$.

2. step. Summing the equations of (5.12), we obtain the variational form (5.17). Testing (5.17) with test functions of the type $(\tilde{x}, 0)$ or $(0, \tilde{y})$, we see that (5.12) and (5.17) are, in fact, equivalent.

3. step. The equivalence of (ii) and (iii) is stated in Corollary 5.3. The equivalence of (i) and (iii) follows from step 2 and Exercise 73.

4. step. It remains to prove (5.18): From (5.13)–(5.14), we obtain

$$\begin{aligned} \|x\|_X + \|y\|_Y &\leq \frac{1}{\alpha} \|x^*\|_{X^*} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \|y^*\|_{Y^*} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \left(\|x^*\|_{X^*} + \frac{\|a\|}{\beta} \|y^*\|_{Y^*} \right) \\ &= \left[\frac{1}{\alpha} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \right] \|x^*\|_{X^*} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \left(1 + \frac{\|a\|}{\beta} \right) \|y^*\|_{Y^*} \\ &\leq \left[\frac{1}{\alpha} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \left(1 + \frac{\|a\|}{\beta} \right) \right] [\|x^*\|_{X^*} + \|y^*\|_{Y^*}]. \end{aligned}$$

With the operator $Tz := c(z, \cdot)$, this proves that the solution operator $T^{-1} : X^* \times Y^* \rightarrow X \times Y$ has operator norm $\|T^{-1}\| \leq C \left[\frac{1}{\alpha} + \frac{1}{\beta} \left(1 + \frac{\|a\|}{\alpha} \right) \frac{1}{\beta} \left(1 + \frac{\|a\|}{\beta} \right) \right]$, where $C > 0$ depends only on the norms chosen on $Z = X \times Y$ and $Z^* = X^* \times Y^*$. According to Theorem 5.2, it holds $\|T^{-1}\| = 1/\gamma$. This concludes the proof. ■

Exercise 74. Give a direct proof that $c(\cdot, \cdot)$ from Corollary 5.8 satisfies the LBB condition, i.e., prove directly that (i) implies (ii). **Hint.** For $(x, y) \neq 0$ use the orthogonal decomposition $x = x_1 + x_2 \in X_0 + X_0^\perp$ and estimate $\|x_1\|_X$, $\|x_2\|_X$, and $\|y\|_Y$ separately. \square

Corollary 5.8 together with Exercise 70 provides a solvability theory and the Céa lemma for Galerkin discretizations of saddle point problems.

Corollary 5.9 (Céa Lemma for Saddle Point Problems, Version I). Let $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ be continuous bilinear forms on a Hilbert space X and a reflexive Banach space Y . Given $(x^*, y^*) \in X^* \times Y^*$, let $(x, y) \in X \times Y$ be a solution of the saddle point problem (5.12). Let $X_h \subset X$ and $Y_h \subset Y$ be finite dimensional subspaces and define $X_{0h} := \{x_h \in X_h \mid b(x_h, \cdot) = 0 \in Y_h^*\}$. Suppose that

- $\alpha_h := \inf_{v_h \in X_{0h} \setminus \{0\}} \sup_{w_h \in X_{0h} \setminus \{0\}} \frac{a(v_h, w_h)}{\|v_h\|_X \|w_h\|_X} > 0$,
- $\beta_h := \inf_{y_h \in Y_h \setminus \{0\}} \sup_{x_h \in X_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y} > 0$.

Then, there is a unique solution $(x_h, y_h) \in X_h \times Y_h$ of the discrete saddle point problem

$$\begin{aligned} a(x_h, \cdot) + b(\cdot, y_h) &= x^* \in X_h^*, \\ b(x_h, \cdot) &= y^* \in Y_h^*, \end{aligned} \quad (5.19)$$

and there holds

$$\|x - x_h\|_X + \|y - y_h\|_Y \leq C \left(\min_{\tilde{x}_h \in X_h} \|x - \tilde{x}_h\|_X + \min_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y \right)$$

The constant $C > 0$ depends only on $(\|a\| + \|b\|)/\gamma_h$ with $\gamma_h := \left[\frac{1}{\alpha_h} + \frac{1}{\beta_h} \left(1 + \frac{\|a\|}{\alpha_h} \right) \frac{1}{\beta_h} \left(1 + \frac{\|a\|}{\beta_h} \right) \right]$.

Proof. The existence and uniqueness of (x_h, y_h) follows from the abstract Brezzi theorem; see Corollary 5.8. For Petrov-Galerkin schemes, the constant in the Céa lemma depends only on the quotient of the continuity bound and the discrete inf-sup constant; see Exercise 70. Both constants have been estimated in Corollary 5.8. \blacksquare

Remark. The Galerkin discretization of saddle point problems is structurally much more difficult than for problems of the Lax-Milgram lemma:

(i) Note that $X_{0h} \not\subseteq X_0 := \{v \in X \mid b(v, \cdot) = 0 \in Y^*\}$. There may be even no relation between X_0 and X_{0h} besides the trivial $X_0 \cap X_h \subseteq X_{0h}$. In particular, there is no relation between α and α_h even if $a(\cdot, \cdot)$ is elliptic on X_0 .

(ii) However, if $a(\cdot, \cdot)$ is already elliptic on X , i.e., $\tau := \inf_{x \in X \setminus \{0\}} \frac{a(x, x)}{\|x\|_X^2} > 0$ this implies $\alpha \geq \tau$ and $\alpha_h \geq \tau$ for the continuous and discrete inf-sup constant of $a(\cdot, \cdot)$.

(iii) Moreover, $\beta > 0$ from the continuous formulation does not imply $\beta_h > 0$ for the discrete formulation. Below, we introduce Fortin's criterium which provides some help on this matter.

(iv) Finally, we recall that $\beta_h > 0$ implies necessarily $\dim Y_h \leq \dim X_h$; see Proposition 5.4. \square

Exercise 75. For a matrix $A \in \mathbb{R}^{m \times n}$ holds $\ker(A^T) = (\text{range} A)^\perp$ as well as $\text{range}(A^T) = (\ker A)^\perp$, where $(\cdot)^\perp$ denotes the orthogonal complement with respect to the usual Euclidean product in \mathbb{R}^m resp. \mathbb{R}^n . \square

The following two exercises consider the discretization of the mixed problem (5.12). We stress that a linear system similar to the one here, also appeared for the discretization of the Neumann problem, where we had to realize the linear side constraint $\int_\Omega u_h dx = 0$.

Exercise 76. Let $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ be continuous bilinear forms on a Hilbert space X and a reflexive Banach space Y . We replace X and Y by finite dimensional subspaces X_h and Y_h , respectively. Show that the computation of a discrete solution $(x_h, y_h) \in X_h \times Y_h$ of

$$\begin{aligned} a(x_h, \cdot) + b(\cdot, y_h) &= x^* \in X_h^*, \\ b(x_h, \cdot) &= y^* \in Y_h^*, \end{aligned} \tag{5.20}$$

is equivalent to the solution of a linear system with a matrix of the type $M := \begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix}$. \square

Exercise 77. Let $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{m \times n}$, and $M := \begin{pmatrix} A & B^T \\ B & 0 \end{pmatrix}$. Assume that A is positive definite on the kernel of B . Prove that M is regular if and only if $\text{range}(B) = \mathbb{R}^m$. \square

We conclude this section with an improved Céa lemma for saddle point problems; cf. Corollary 5.9.

Theorem 5.10 (Céa Lemma for Saddle Point Problems, Version II). *Let $a : X \times X \rightarrow \mathbb{R}$ and $b : X \times Y \rightarrow \mathbb{R}$ be continuous bilinear forms on a Hilbert space X and a reflexive Banach space Y . Given $(x^*, y^*) \in X^* \times Y^*$, let $(x, y) \in X \times Y$ be a solution of the saddle point problem (5.12). Let $X_h \subset X$ and $Y_h \subset Y$ be finite dimensional subspaces and define $X_{0h} := \{x_h \in X_h \mid b(x_h, \cdot) = 0 \in Y_h^*\}$. Suppose that*

- $\alpha_h := \inf_{v_h \in X_{0h} \setminus \{0\}} \frac{a(v_h, v_h)}{\|v_h\|_X^2} > 0$,
- $\beta_h := \inf_{y_h \in Y_h \setminus \{0\}} \sup_{x_h \in X_h \setminus \{0\}} \frac{b(x_h, y_h)}{\|x_h\|_X \|y_h\|_Y} > 0$.

Then, there is a unique solution $(x_h, y_h) \in X_h \times Y_h$ of the discrete saddle point problem

$$\begin{aligned} a(x_h, \cdot) + b(\cdot, y_h) &= x^* \in X_h^*, \\ b(x_h, \cdot) &= y^* \in Y_h^*, \end{aligned} \tag{5.21}$$

and there holds

$$\|x - x_h\|_X \leq \left(1 + \frac{\|a\|}{\alpha_h}\right) \left(1 + \frac{\|b\|}{\beta_h}\right) \min_{\tilde{x}_h \in X_h} \|x - \tilde{x}_h\|_X + \frac{\|b\|}{\alpha_h} \min_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y \quad (5.22)$$

and

$$\|y - y_h\|_Y \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \min_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y + \frac{\|a\|}{\beta_h} \|x - x_h\|_X. \quad (5.23)$$

Sketch of Proof of Theorem 5.10. The unique existence of a discrete solution $(x_h, y_h) \in X_h \times Y_h$ follows from the Brezzi Theorem 5.6 applied for $X_h \times Y_h$. The quasioptimality is proven in three steps:

- First, we prove estimate (5.23).
- Second, we prove quasioptimality of $\|x - x_h\|_X$ with respect to the affine space $Z_h := \{\tilde{x}_h \in X_h \mid b(\tilde{x}_h, \cdot) = y^* \in Y_h^*\}$.
- In a final step, we estimate the bestapproximation error with respect to Z_h by the bestapproximation error with respect to the entire discrete space X_h which then leads to (5.22).

This general concept even works for nonlinear problems with linear side constraint. ■

Proof. We first note the Galerkin orthogonality, which now reads

$$\begin{aligned} a(x - x_h, \cdot) + b(\cdot, y - y_h) &= 0 \in X_h^*, \\ b(x - x_h, \cdot) &= 0 \in Y_h^*, \end{aligned} \quad (5.24)$$

1. step. There holds

$$\|y - y_h\|_Y \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \|y - \tilde{y}_h\|_Y + \frac{\|a\|}{\beta_h} \|x - x_h\|_X \quad \text{for all } \tilde{y}_h \in Y_h :$$

According to the definition of β_h , there holds

$$\beta_h \| \tilde{y}_h - y_h \|_Y \leq \sup_{\tilde{x}_h \in X_h \setminus \{0\}} \frac{b(\tilde{x}_h, \tilde{y}_h - y_h)}{\|x_h\|_X}.$$

With the Galerkin orthogonality, the nominator may be written as

$$\begin{aligned} b(\tilde{x}_h, \tilde{y}_h - y_h) &= -(a(x - x_h, \tilde{x}_h) + b(\tilde{x}_h, y - y_h)) + b(\tilde{x}_h, \tilde{y}_h - y_h) \\ &= -a(x - x_h, \tilde{x}_h) + b(\tilde{x}_h, \tilde{y}_h - y) \end{aligned}$$

Therefore, continuity of $a(\cdot, \cdot)$ and $b(\cdot, \cdot)$ lead to

$$\beta_h \| \tilde{y}_h - y_h \|_Y \leq \|a\| \|x - x_h\|_X + \|b\| \| \tilde{y}_h - y \|_Y.$$

Altogether, a triangle inequality $\|y - y_h\|_Y \leq \|y - \tilde{y}_h\|_Y + \| \tilde{y}_h - y_h \|_Y$ yields step 1.

2. step. With the affine space $Z_h := \{\tilde{x}_h \in X_h \mid b(\tilde{x}_h, \cdot) = y^* \in Y_h^*\}$, there holds

$$\|x - x_h\|_X \leq \left(1 + \frac{\|a\|}{\alpha_h}\right) \|x - z_h\|_X + \frac{\|b\|}{\alpha_h} \|y - \tilde{y}_h\|_Y \quad \text{for all } z_h \in Z_h \text{ and } \tilde{y}_h \in Y_h :$$

Since $x_h, z_h \in Z_h$, there holds $x_h - z_h \in X_{0h}$. According to the definition of α_h , we see

$$\alpha_h \|x_h - z_h\|_X^2 \leq a(x_h - z_h, x_h - z_h) = a(x_h - x, x_h - z_h) + a(x - z_h, x_h - z_h).$$

For the first term, the Galerkin orthogonality implies

$$a(x_h - x, x_h - z_h) = b(x_h - z_h, y - y_h) = b(x_h - z_h, \tilde{y}_h - y_h) + b(x_h - z_h, y - \tilde{y}_h),$$

where the first summand $b(x_h - z_h, \tilde{y}_h - y_h) = 0$ drops out by use of $x_h - z_h \in X_{0h}$. By continuity of $a(\cdot, \cdot)$ and $b(\cdot, \cdot)$, we see

$$\alpha_h \|x_h - z_h\|_X \leq \|a\| \|x - z_h\|_X + \|b\| \|y - \tilde{y}_h\|_Y.$$

Again, a triangle inequality $\|x - x_h\|_X \leq \|x - z_h\|_X + \|x_h - z_h\|_X$ yields step 2.

3. step. There holds

$$\|x - z_h\|_X \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \|x - \tilde{x}_h\|_X \quad \text{for all } \tilde{x}_h \in X_h \text{ and some } z_h \in Z_h \text{ depending on } \tilde{x}_h :$$

We define $W_h := (X_{0h})^\perp \subseteq X_h$ and consider the operators $B_1 \in L(W_h, Y_h^*)$ and $B_2 \in L(Y_h, W_h^*)$ defined by $B_1 w_h := b(w_h, \cdot)$ and $B_2 y_h := b(\cdot, y_h)$. Note that

$$0 < \beta_h = \inf_{\tilde{y}_h \in Y_h \setminus \{0\}} \sup_{\tilde{x}_h \in X_h \setminus \{0\}} \frac{b(\tilde{x}_h, \tilde{y}_h)}{\|\tilde{x}_h\|_X \|\tilde{y}_h\|_Y} = \inf_{\tilde{y}_h \in Y_h \setminus \{0\}} \sup_{w_h \in W_h \setminus \{0\}} \frac{b(w_h, \tilde{y}_h)}{\|w_h\|_X \|\tilde{y}_h\|_Y}.$$

According to Lemma 5.1, the operator B_2 is injective with closed range and $1/\beta_h = \|B_2^{-1}\| : \text{range}(B_2) \rightarrow Y_h$. From this, we derive that $B_1 = B_2^* \circ I_{Y_h}$ is surjective due to $\text{range}(B_1) = \text{range}(B_2^*) = (\ker B_2)^\circ = Y_h^*$. Note that by definition of $W_h := (X_{0h})^\perp \subseteq X_h$, the operator B_1 is injective and thus an isomorphism between W_h and Y_h^* . In particular, this yields bijectivity of B_2 as well as

$$\|B_1^{-1}\| = \|I_{Y_h}^{-1}(B_2^*)^{-1}\| = \|(B_2^{-1})^*\| = \|B_2^{-1}\| = 1/\beta_h.$$

In particular, there is a unique element $w_h \in W_h$ with $b(w_h, \cdot) = B_1 w_h = b(x - \tilde{x}_h, \cdot) \in Y_h^*$ and there holds $\|w_h\|_X \leq \beta_h^{-1} \|b(x - \tilde{x}_h, \cdot)\|_{X^*} \leq (\|b\|/\beta_h) \|x - \tilde{x}_h\|_X$. The element $z_h := \tilde{x}_h + w_h \in X_h$ satisfies $b(z_h, \cdot) = b(x, \cdot) = y^* \in Y_h^*$ and thus $z_h \in Z_h$. Now, we finally see

$$\|x - z_h\|_X \leq \|x - \tilde{x}_h\|_X + \|w_h\|_X \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \|x - \tilde{x}_h\|_X.$$

This concludes step 3.

4. step. The proof of (5.23) follows by finite dimension: Note that step 1 implies

$$\|y - y_h\|_Y \leq \left(1 + \frac{\|b\|}{\beta_h}\right) \inf_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y + \frac{\|a\|}{\beta_h} \|x - x_h\|_X,$$

and it only remains to see that the infimum is, in fact, attained: To that end, choose an infimizing sequence (y_k) in Y_h , i.e.

$$\lim_{k \rightarrow \infty} \|y - y_k\|_Y = \inf_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y.$$

According to the triangle inequality, there holds $\|y_k\|_Y \leq \|y\|_Y + \|y - y_k\|_Y$, i.e. the sequence (y_k) is a bounded sequence in the finite dimensional space Y_h . Thus, the Bolzano-Weierstrass theorem yields the existence of a convergent subsequence (y_{k_ℓ}) with limit $y_0 \in Y_h$. By continuity, we conclude

$$\inf_{\tilde{y}_h \in Y_h} \|y - \tilde{y}_h\|_Y = \lim_{\ell \rightarrow \infty} \|y - y_{k_\ell}\|_Y = \|y - y_0\|_Y.$$

5. step. The proof of (5.22) now follows from a combination of step 2 and step 3: For arbitrary $\tilde{x}_h \in X_h$, choose $z_h \in Z_h$ by use of step 3. Let $\tilde{y}_h \in Y_h$ and be arbitrary. We then infer

$$\begin{aligned} \|x - x_h\|_X &\leq \left(1 + \frac{\|a\|}{\alpha_h}\right) \|x - z_h\|_X + \frac{\|b\|}{\alpha_h} \|y - \tilde{y}_h\|_Y \\ &\leq \left(1 + \frac{\|a\|}{\alpha_h}\right) \left(1 + \frac{\|b\|}{\beta_h}\right) \|x - \tilde{x}_h\|_X + \frac{\|b\|}{\alpha_h} \|y - \tilde{y}_h\|_Y. \end{aligned}$$

Now, we take the infimum over \tilde{x}_h and \tilde{y}_h and note that, according to finite dimension, this infimum is attained by independent minima. ■

Bibliography

- [Bra] Dietrich Braess: *Finite elements. Theory, fast solvers, and applications in elasticity theory*, Cambridge University Press, Cambridge, 2007.
- [McL] William McLean: *Strongly elliptic systems and boundary integral equations*, Cambridge University Press, Cambridge, 2000.