

Numerics of Partial Differential Equations: Stationary Problems

Lecture Notes

Michael Feischl and Dirk Praetorius

October 21, 2020

Chapter 1

Introduction

1.1 Strong Form and Variational Form

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

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

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 $C^k(\bar{\Omega}) := \{w|_{\bar{\Omega}} \mid w \in C^k(\mathbb{R}^d)\}$. 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

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 $\|u\|_H = \|F\|_{H^}$. ■*

The proof of this theorem can be found in each textbook of functional analysis.

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

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 .

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

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

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 $\langle\langle \cdot ; \cdot \rangle\rangle$, 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

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

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 vector space over \mathbb{R} 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 (1.22) 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. ■

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 $M \subseteq X$. □

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 $\langle \cdot ; \cdot \rangle$ 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:

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

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 $a(u, \cdot) = F \in H^*$.

□

Chapter 2

Sobolev Spaces and Poisson Problem

2.1 Sobolev Spaces on Domains

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

The following main theorem of calculus, which will not be proven in this lecture, 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 . ■*

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 8. 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 8. (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 $\partial_j u, \widetilde{\partial_j u} \in L^1_{\text{loc}}(\Omega)$ 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 $u \in L^1_{\text{loc}}(\Omega)$ 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

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 and only if each bounded set $S \subseteq X$ is mapped to a pre-compact set $A(S) \subseteq Y$, i.e., $\overline{A(S)} \subseteq Y$ is compact.

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 2.6. *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

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

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 $x_n \rightharpoonup x \in X$. For $\alpha := \liminf_n f(x_n) = \infty$, (2.6) is trivial. We thus may assume $\alpha < \infty$. Let $\beta > \alpha$ and define $\alpha_n := \max\{\beta, f(x_n)\} \rightarrow \beta$. 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 $H^1(\Omega)$.

Proposition 2.7. *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 : By definition of continuity, there exists an open neighborhood $O \subseteq H^1(\Omega)$ of 0 such that $|v|_{H^1} \leq 1$ for all $v \in O$. Without loss of generality, we may choose a

radius $r > 0$ sufficiently small such that $\overline{B_r(0)} \subset O \subset H^1(\Omega)$ for the closed ball with radius r and center zero. This implies

$$|v|_{H^1} = \frac{1}{r} \|v\|_{H^1(\Omega)} \left| r \frac{v}{\|v\|_{H^1(\Omega)}} \right|_{H^1} \leq \frac{1}{r} \|v\|_{H^1(\Omega)}.$$

This proves existence of $C_1 := 1/r$.

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 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 Banach-Alaoglou 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 2.6 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 2.6 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 2.8 (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 2.7, 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

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

and from the boundedness of Ω . ■

Corollary 2.9 (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 $D\Phi = \lambda^{-1}\mathbf{I}$ 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 2.10 (Meyers-Serrin). *For each integer order $m \in \mathbb{N}$, $C^\infty(\bar{\Omega})$ and, in particular, $C^\infty(\Omega) \cap H^m(\Omega)$ are dense subspaces of $H^m(\Omega)$.* ■

Theorem 2.11 (Trace Operator). *There is a unique operator $\gamma \in L(H^1(\Omega); L^2(\Gamma))$ such that $\gamma v = v|_{\Gamma}$ for all $v \in C^1(\bar{\Omega})$, 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 2.11, we can prove that the integration by parts formula also holds for Sobolev functions $u, v \in H^1(\Omega)$.

Corollary 2.12 (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 2.13 (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 2.7. 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

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

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 2.14 (Friedrichs Inequality). *It holds that*

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

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 2.15. *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 9. 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 10. 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

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 2.16. (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

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

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

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

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 9. Then, the main proposition reads as follows:

Proposition 2.17. (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

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

(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

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

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. ■

2.3.3 Neumann Problem

Finally, we consider the Neumann problem

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

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

$$\int_{\Omega} u \, dx = 0. \quad (2.19)$$

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

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

to allow for the existence of (strong) solutions. Recall the trace norm $\|\cdot\|_{H^{1/2}(\Gamma)}$ from Exercise 9.

Proposition 2.18. (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

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

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

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

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

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

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

3.1 P1-Finite Element Method in 2D

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

$$\mathcal{K}_T := \{x_T, y_T, z_T\} \quad (3.1)$$

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

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

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

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

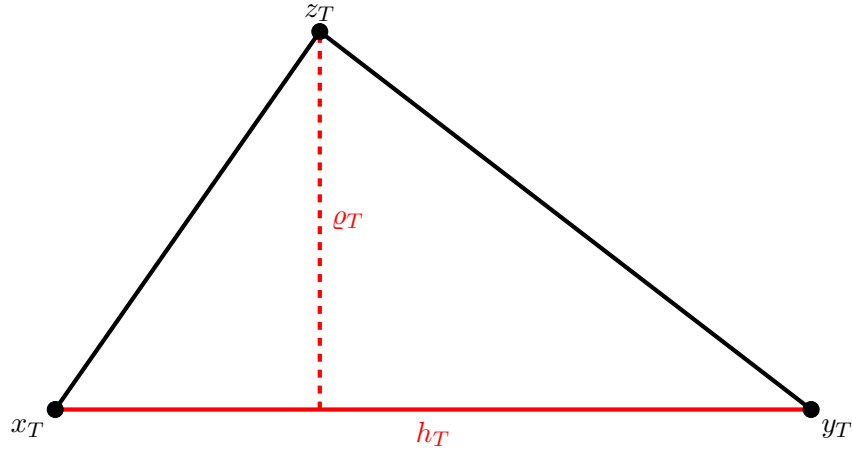


FIGURE 3.1. The diameter h_T of the triangle T is the length of the longest edge (possibly non unique). The quantity ϱ_T denotes the corresponding height.

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 11. 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 $\text{conv}(M) = \left\{ \sum_{j=1}^N \lambda_j x_j \mid N \in \mathbb{N}, x_j \in M, \lambda_j \geq 0 \text{ with } \sum_{j=1}^N \lambda_j = 1 \right\}$. 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

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

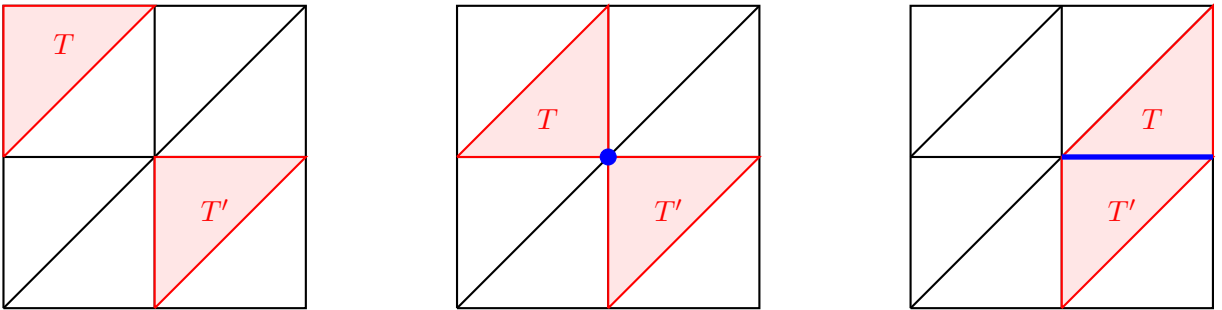


FIGURE 3.2. For a regular triangulation \mathcal{T} , the intersection of two elements $T \neq T'$ is either empty, a joint node, or a joint edge.

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

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

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

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 12. 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 $v \in C(\Omega)$. \square

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

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

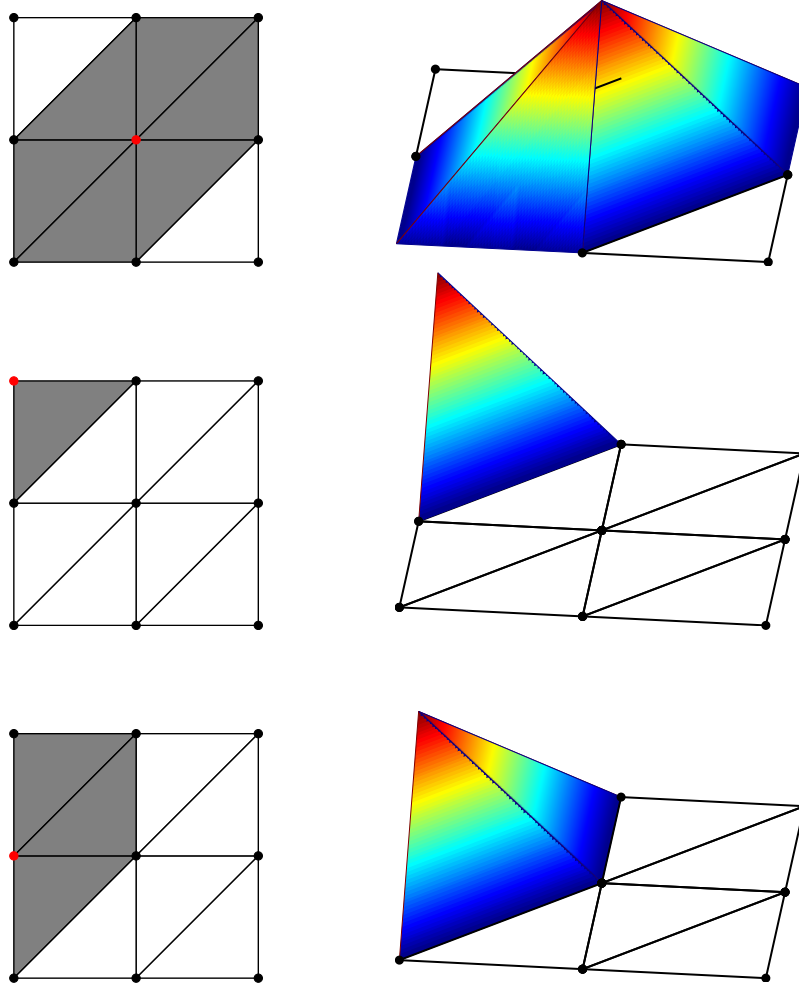


FIGURE 3.3. Examples of $P1$ hat functions ζ_z : The left figures show the mesh as well as the support $\text{supp}(\zeta_z)$ in grey, where the corresponding node $z \in \mathcal{K}$ is indicated in red. The right figures show the plots of the hat functions. Triangles $T \in \mathcal{T}$ with $\zeta_z|_T = 0$ 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 12. ■

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. \square

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**

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

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

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

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

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

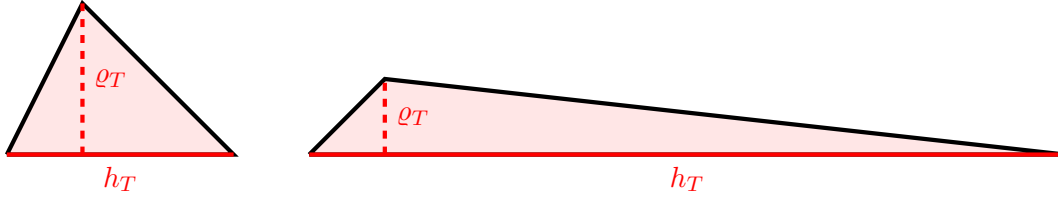


FIGURE 3.4. two triangles with small (left) and large (right) shape regularity constant $\sigma(T) := h_T/\varrho_T$

Finally, we consider the **Neumann problem**

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

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

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

3.2 Approximation Theorem and Bramble-Hilbert Lemma

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} , see Fig. 3.4. 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 3.5. 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.

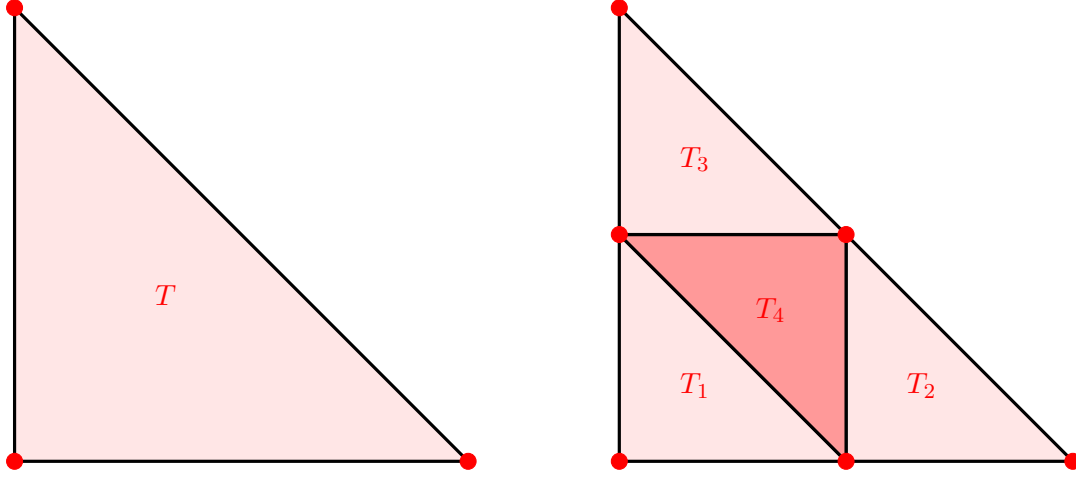


FIGURE 3.5. Red-refinement refines the element $T \in \mathcal{T}^{(\text{old})}$ into 4 similar elements $T_1, \dots, T_4 \in \mathcal{T}^{(\text{new})}$. The new nodes $\mathcal{K}^{(\text{new})} \setminus \mathcal{K}^{(\text{old})}$ are just the edge midpoints for all edges $E \in \mathcal{E}^{(\text{old})}$. In particular, regularity of $\mathcal{T}^{(\text{old})}$ implies regularity of $\mathcal{T}^{(\text{new})}$.

Exercise 13. 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. \square

Exercise 14. 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 ? \square

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 $H^m(\Omega) \subseteq C(\overline{\Omega})$. \blacksquare

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

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

This theorem will be shown later. First, 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). \blacksquare

Remark. Corollary 3.6 has two important consequences: First, according to C ea’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. \square

Exercise 15. Use the Poincar  inequality and the Meyers-Serrin theorem to prove that $C_*^\infty(\bar{\Omega})$ is dense in $H_*^1(\Omega)$. \square

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 2.6, 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 16, 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 $H^{m+1}(T)$.

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 16. 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 $u \in C^\infty(\overline{T})$, 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. ■

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

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

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

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

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

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

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} ,

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

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

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

The combination of the last two estimates proves that

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

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

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

and thus concludes the proof. \blacksquare

Remark. The proof of Theorem 3.5 shows that it is enough to assume $u \in C(\bar{\Omega}) \cap H^2(\mathcal{T})$, 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 $H^2(\Omega) \subseteq C(\bar{\Omega}) \cap H^2(\mathcal{T})$. 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 17. 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

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.

Appendix A

Some Facts from Functional Analysis

In this appendix we collect some results from introductory functional analysis courses which are used throughout. We stick with the case of vector spaces over \mathbb{R} .

A.1 Main Theorems from Functional Analysis

Theorem A.1 (Hahn-Banach Extension Theorem). *Let $p : X \rightarrow \mathbb{R}$ be a sublinear functional on a linear space X , i.e. $p(x + y) \leq p(x) + p(y)$ and $p(\lambda x) = \lambda p(x)$ for all $x, y \in X$ and $\lambda \geq 0$. If Y is a subspace of X and $f : Y \rightarrow \mathbb{R}$ is a linear functional with $f \leq p$ on Y , there is a linear extension $F : X \rightarrow \mathbb{R}$ with $F|_Y = f$ and $F \leq p$ on X . ■*

If X is a normed space and $f \in Y^*$, one may choose $p(x) = \|x\|_X \|f\|_{X^*}$ to prove the extension theorem for continuous linear functionals.

Corollary A.2. *If Y is the subspace of a normed space X and $f \in Y^*$, there is an extension $F \in X^*$ with $F|_Y = f$ and $\|F\|_{X^*} = \|f\|_{Y^*}$. ■*

One then considers the subspace $Y := \text{span}\{x\}$ and $f(\lambda x) = \lambda \|x\|_X$ to derive the following corollary:

Corollary A.3. *If X is a normed space and $x \in X$, there is a linear functional $f \in X^*$ with $\|f\|_{X^*} = 1$ and $f(x) = \|x\|_X = \sup_{\|f\|_{X^*}=1} |f(x)|$. ■*

Theorem A.4 (Hahn-Banach Separation Theorem). *Let X be a normed space, and let A and B be convex, nonempty subsets of X with $A \cap B = \emptyset$.*

- (i) If A is open, there is a linear functional $f \in X^*$ and a scalar $\lambda \in \mathbb{R}$ such that $f(x) < \lambda \leq f(y)$ for all $x \in A$ and $y \in B$.*
- (ii) If A is compact and B is closed, there is a linear functional $f \in X^*$ and scalars $\lambda_1, \lambda_2 \in \mathbb{R}$ such that $f(x) \leq \lambda_1 < \lambda_2 \leq f(y)$ for all $x \in A$ and $y \in B$. ■*

If Y is a subspace of X , one can use (ii) to characterize the closure \bar{Y} of Y in X . The proof only needs that each bounded linear functional $f \in Y^*$ is trivial, i.e. $f|_Y = 0$.

Corollary A.5. *Let Y be a subspace of the normed space X . Then, $x \in X$ satisfies $x \in \overline{Y}$ if and only if $f(x) = 0$ for all $f \in X^*$ with $f|_Y = 0$.*

Proof. For $x \in \overline{Y}$ and $f \in X^*$ with $f|_Y = 0$, continuity yields $f(x) = 0$. The converse implication is proven by contradiction: We assume that $x \notin \overline{Y}$ and choose $f \in X^*$ such that $f(x) < \lambda \leq f(y)$ for all $y \in Y$ and some fixed $\lambda \in \mathbb{R}$. Using that Y is a vector space, we infer that $\lambda \leq f(\pm y) = -f(\mp y) \leq -\lambda$ and thus $f(y) \in [\lambda, -\lambda]$ for all $y \in Y$. As bounded linear functionals are trivial, we obtain $f|_Y = 0$. According to our assumptions, this implies $f(x) = 0$ and thus contradicts $f(x) < \lambda \leq f(0) = 0$. ■

The following corollary is an immediate consequence of the last one.

Corollary A.6. *Let Y be a subspace of the normed space X . Then, Y is dense in X if and only if each functional $f \in X^*$ with $f|_Y = 0$ is trivial, i.e., $f = 0 \in X^*$.* ■

For an operator $T \in L(X; Y)$, one defines $(T^*y^*)(x) := y^*(Tx)$ for all $y^* \in Y^*$ and $x \in X$. From the continuity of T , we see that $T^*y^* \in X^*$, and obviously $T^* : Y^* \rightarrow X^*$ is a linear operator. From the corollary of the Hahn-Banach extension theorem, we derive for the operator norm

$$\begin{aligned} \|T^*\| &= \sup_{\|y^*\|_{Y^*}=1} \|T^*y^*\|_{X^*} = \sup_{\|y^*\|_{Y^*}=1} \sup_{\|x\|_X=1} (T^*y^*)(x) \\ &= \sup_{\|x\|_X=1} \sup_{\|y^*\|_{Y^*}=1} (y^*)(Tx) = \sup_{\|x\|_X=1} \|Tx\|_Y = \|T\|, \end{aligned}$$

i.e. there holds $T^* \in L(Y^*; X^*)$ with operator norm $\|T^*\| = \|T\|$. The operator T^* is called the **adjoint operator** of T .

Theorem A.7 (Banach Closed Range Theorem). *For an operator $T \in L(X; Y)$ between Banach spaces X and Y and $T^* \in L(Y^*; X^*)$ its adjoint, the following is pairwise equivalent:*

- (i) $\text{range}(T)$ is a closed subspace of Y .
- (ii) $\text{range}(T) = (\ker T^*)^\circ := \{y \in Y \mid \forall y^* \in \ker(T^*) \quad y^*(y) = 0\}$.
- (iii) $\text{range}(T^*)$ is a closed subspace of X^* .
- (iv) $\text{range}(T^*) = (\ker T)^\circ := \{x^* \in X^* \mid \forall x \in \ker(T) \quad x^*(x) = 0\}$. ■

A.2 Hilbert Spaces

A space X is called **Hilbert space** if it is a Banach space whose norm is induced by a scalar product.

Theorem A.8. *Let Y be the closed subspace of a Hilbert space X and $Y^\perp := \{x \in X \mid \forall y \in Y \quad (x; y)_X = 0\}$ the orthogonal complement. Then, there holds $X = Y \oplus Y^\perp$ in the sense of the linear algebra, i.e. every element $x \in X$ has a unique decomposition $x = y + y^\perp$ with some $y \in Y$ and $y^\perp \in Y^\perp$.* ■

With the orthogonal decomposition $X = Y \oplus Y^\perp$, one can define a projection $\pi_Y : X \rightarrow Y$ by $x = y + y^\perp \mapsto y$.

Corollary A.9. *Let Y be the closed subspace of a Hilbert space X . Then, there is a unique linear operator $\Pi : X \rightarrow Y$ with $\Pi|_Y = \text{id}$ and $\ker(\Pi) = Y^\perp$, which is called **orthogonal projection** onto Y . This projection is continuous with operator norm $\|\Pi\| = 1$ and symmetric, i.e. $(x ; y)_X = (\Pi x ; y)_X$ for all $x \in X$ and $y \in Y$. Moreover, the orthogonal projection is the solution operator for the best approximation problem, $\|x - \Pi x\|_X = \min_{y \in Y} \|x - y\|_X$. ■*

The dual space X^* of a Hilbert space X has a straight-forward representation, and one can somehow identify X with X^* .

Theorem A.10 (Riesz). *For a Hilbert space X , the **Riesz mapping** $I_X : X \rightarrow X^*$, $I_X x := (x ; \cdot)_X \in X^*$, is an isometric isomorphism. ■*