# Numerical Methods for Partial Differential Equations

Volker John

Summer Semester 2013

# Contents

# Chapter 1

# Some Partial Differential Equations From Physics

**Remark 1.1** *Contents.* This chapter introduces some partial differential equations (pde's) from physics to show the importance of this kind of equations and to motivate the application of numerical methods for their solution. □

## 1.1 The Heat Equation

**Remark 1.2** *Derivation.* The derivation of the heat equation follows (Wladimirow, 1972, p. 39). Let $\mathbf{x} = (x_1, x_2, x_3)^T \in \Omega \subset \mathbb{R}^3$, where $\Omega$ is a domain, $t \in \mathbb{R}$, and consider the following physical quantities

- $u(t, \mathbf{x})$ – temperature at time $t$ and at the point $\mathbf{x}$ with unit $[K]$,
- $\rho(t, \mathbf{x})$ – density of the considered species with unit $[kg/m^3]$,
- $c(t, \mathbf{x})$ – specific heat capacity of the species with unit $[J/(kg\ K)] = [W\ s/(kg\ K)]$,
- $k(t, \mathbf{x})$ – thermal conductivity of the species with unit $[W/(m\ K)]$,
- $F(t, \mathbf{x})$ – intensity of heat sources or sinks with unit $[W/m^3]$.

Consider the heat equilibrium in an arbitrary volume $V \subset \Omega$ and in an arbitrary time interval $(t, t + \Delta t)$. First, there are sources or sinks of heat: heat can enter or leave $V$ through the boundary $\partial V$, or heat can be produced or absorbed in $V$. Let $\mathbf{n}(\mathbf{x})$ be the unit outer normal at $\mathbf{x} \in \partial V$. Due to Fourier's law , one finds that the heat

$$Q_1 = \int_t^{t+\Delta t} \int_{\partial V} k \frac{\partial u}{\partial \mathbf{n}}(t, \mathbf{s})\ d\mathbf{s}\ dt = \int_t^{t+\Delta t} \int_{\partial V} (k\nabla u \cdot \mathbf{n})(t, \mathbf{s})\ d\mathbf{s}\ dt,\ [J],$$

enters through $\partial V$ into $V$. One obtains with partial integration (Gaussian theorem)

$$Q_1 = \int_t^{t+\Delta t} \int_V \nabla \cdot (k\nabla u)(t, \mathbf{x})\ d\mathbf{x}\ dt.$$

In addition, the heat

$$Q_2 = \int_t^{t+\Delta t} \int_V F(t, \mathbf{x})\ d\mathbf{x}\ dt,\ [W\ s] = [J],$$

is produced in $V$.

Second, a law for the change of the temperature in $V$ has to be derived. Using a Taylor expansion, on gets that the temperature at $\mathbf{x}$ changes in $(t, t + \Delta t)$ by

$$u(t + \Delta t, \mathbf{x}) - u(t, \mathbf{x}) = \frac{\partial u}{\partial t}(t, \mathbf{x})\Delta t + \mathcal{O}((\Delta t)^2).$$

Now, a linear ansatz is utilized, i.e.,

$$u(t + \Delta t, \mathbf{x}) - u(t, \mathbf{x}) = \frac{\partial u}{\partial t}(t, \mathbf{x})\Delta t.$$

With this ansatz, one has that for the change of the temperature in $V$ and for arbitrary $\Delta t$, the heat

$$Q_3 = \int_t^{t+\Delta t} \int_V c\rho \frac{u(t + \Delta t, \mathbf{x}) - u(t, \mathbf{x})}{\Delta t} \, d\mathbf{x} \, dt = \int_t^{t+\Delta t} \int_V c\rho \frac{\partial u}{\partial t}(t, \mathbf{x}) \, d\mathbf{x} \, dt$$

is needed. This heat has to be equal to the heat sources, i.e., it holds $Q_3 = Q_2 + Q_1$, from what follows that

$$\int_t^{t+\Delta t} \int_V \left[ c\rho \frac{\partial u}{\partial t} - \nabla \cdot (k\nabla u) - F \right] (t, \mathbf{x}) \, d\mathbf{x} \, dt = 0.$$

Since the volume $V$ was chosen to be arbitrary and $\Delta t$ was arbitrary as well, the term in the integral has to vanish. One obtains the so-called heat equation

$$c\rho \frac{\partial u}{\partial t} - \nabla \cdot (k\nabla u) = F \quad \text{in } (0, T) \times \Omega.$$

At this point of modeling one should check if the equation is dimensionally correct. One finds that all terms have the unit $[W/m^3]$.

For a homogeneous species, $c$, $\rho$, and $k$ are positive constants. Then, the heat equation simplifies to

$$\frac{\partial u}{\partial t} - \varepsilon^2 \Delta u = f \quad \text{in } (0, T) \times \Omega, \tag{1.1}$$

with $\varepsilon^2 = k/(c\rho)$, $[m^2/s]$ and $f = F/(c\rho)$, $[K/s]$. To obtain a well-posed problem, (1.1) has to be equipped with an initial condition $u(0, \mathbf{x})$ and appropriate boundary conditions on $(0, T)\partial\Omega$. □

**Remark 1.3** *Boundary conditions.* For the theory and the numerical simulation of partial differential equations, the choice of boundary conditions is of utmost importance. For the heat equation (1.1), one can prescribe the following types of boundary conditions:

- Dirichlet[1] condition: The temperature $u(t, \mathbf{x})$ at a part of the boundary is prescribed

  $$u = g_1 \text{ on } (0, T) \times \partial\Omega_D$$

  with $\partial\Omega_D \subset \partial\Omega$. In the context of the heat equation, the Dirichlet condition is also called essential boundary conditions.

- Neumann[2] condition: The heat flux is prescribed at a part of the boundary

  $$-k \frac{\partial u}{\partial \mathbf{n}} = g_2 \text{ on } (0, T) \times \partial\Omega_N$$

  with $\partial\Omega_N \subset \partial\Omega$. This boundary condition is a so-called natural boundary condition for the heat equation.

- Mixed boundary condition, Robin[3] boundary condition: At the boundary, there is a heat exchange according to Newton's law

  $$k \frac{\partial u}{\partial \mathbf{n}} + h(u - u_{\text{env}}) = 0 \text{ on } (0, T) \times \partial\Omega_m,$$

  with $\partial\Omega_m \subset \partial\Omega$, the heat exchange coefficient $h$, $[W/(m^2 K^2)]$, and the temperature of the environment $u_{\text{env}}$.

---

[1] Johann Peter Gustav Lejeune Dirichlet (1805 –1859)
[2] Carl Gottfried Neumann (1832 – 1925)
[3] Gustave Robin (1855 – 1897)

□

**Remark 1.4** *The stationary case.* An important special case is that the temperature is constant in time $u(t, \mathbf{x}) = u(\mathbf{x})$. Then, one obtains the stationary heat equation

$$- \varepsilon^2 \Delta u = f \quad \text{in } \Omega. \tag{1.2}$$

This equation is called Poisson[4] equation. Its homogeneous form, i.e., with $f(\mathbf{x}) = 0$, is called Laplace[5] equation. Solution of the Laplace equation are called harmonic functions. The Poisson equation is the simplest partial differential equation. The most part of this lecture will consider numerical methods for solving this equation.

□

**Remark 1.5** *Another application of the Poisson equation.* The stationary distribution of an electric field with charge distribution $f(\mathbf{x})$ satisfies also the Poisson equation (1.2).

□

**Remark 1.6** *Non-dimensional equations.* The application of numerical methods relies on equations for functions without physical units, the so-called non-dimensional equations. Let

- $L$ – a characteristic length scale of the problem, $[m]$,
- $U$ – a characteristic temperature scale of the problem, $[K]$,
- $T^*$ – a characteristic time scale of the problem, $[s]$.

If the new coordinates and functions are denoted with a prime, one gets with the transformations

$$\mathbf{x}' = \frac{\mathbf{x}}{L}, \quad u' = \frac{u}{U}, \quad t' = \frac{t}{T^*}$$

from (1.1) the non-dimensional equations

$$\frac{\partial}{\partial t'} (U u') \frac{\partial t'}{\partial t} - \varepsilon^2 \sum_{i=1}^{d} \frac{\partial}{\partial x_i'} \left( \frac{\partial}{\partial x_i'} (U u') \frac{\partial x_i'}{\partial x_i} \right) \frac{\partial x_i'}{\partial x_i} = f \quad \text{in } \left( 0, \frac{T}{T^*} \right) \times \Omega' \quad \Longleftrightarrow$$

$$\frac{U}{T^*} \frac{\partial u'}{\partial t'} - \frac{\varepsilon^2 U}{L^2} \sum_{i=1}^{d} \frac{\partial^2 u'}{\partial (x_i')^2} = f \quad \text{in } \left( 0, \frac{T}{T^*} \right) \times \Omega'.$$

Usually, one denotes the non-dimensional functions like the dimensional functions, leading to

$$\frac{\partial u}{\partial t} - \frac{\varepsilon^2 T^*}{L^2} \Delta u = \frac{T^*}{U} f \quad \text{in } \left( 0, \frac{T}{T^*} \right) \times \Omega.$$

For the analysis, one sets $L = 1m$, $U = 1K$, and $T^* = 1s$ which yields

$$\frac{\partial u}{\partial t} - \varepsilon^2 \Delta u = f \quad \text{in } (0, T) \times \Omega, \tag{1.3}$$

with a non-dimensional temperature diffusion $\varepsilon^2$ and a non-dimensional right hand side $f(t, \mathbf{x})$.

The same approach can be applied to the stationary equation (1.2) and one gets

$$- \varepsilon^2 \Delta u = f \quad \text{in } \Omega, \tag{1.4}$$

with the non-dimensional temperature diffusion $\varepsilon^2$ and the non-dimensional right hand side $f(\mathbf{x})$.

□

---

[4]Siméon Denis Poisson (1781 – 1840)
[5]Pierre Simon Laplace (1749 – 1829)

**Remark 1.7** *A standard approach for solving the instationary equation.* The heat equation (1.3) is an initial value problem with respect to time and a boundary value problem with respect to space. Numerical methods for solving initial value problems were topic of Numerical Mathematics 2.

A standard approach for solving the instationary problem consists in using a so-called one-step $\theta$-scheme for discretizing the temporal derivative. Consider two consecutive discrete times $t_n$ and $t_{n+1}$ with $\tau = t_{n+1} - t_n$. Then, the application of a one-step $\theta$-scheme yields for the solution at $t_{n+1}$

$$\frac{u_{n+1} - u_n}{\tau} - \theta\varepsilon^2\Delta u_{n+1} - (1-\theta)\varepsilon^2\Delta u_n = \theta f_{n+1} + (1-\theta)f_n,$$

where the subscript at the functions denotes the time level. This equation is equivalent to

$$u_{n+1} - \tau\theta\varepsilon^2\Delta u_{n+1} = u_n + \tau(1-\theta)\varepsilon^2\Delta u_n + \tau\theta f_{n+1} + \tau(1-\theta)f_n. \qquad (1.5)$$

For $\theta = 0$, one obtains the forward Euler scheme, for $\theta = 0.5$ the Crank–Nicolson scheme (trapezoidal rule), and for $\theta = 1$ the backward Euler scheme.

Given $u_n$, (1.5) is a boundary value problem for $u_{n+1}$. That means, one has to solve in each discrete time a boundary value problem. For this reason, this lecture will concentrate on the numerical solution of boundary value problems. □

**Example 1.8** *Demonstrations with the code* MooNMD JOHN AND MATTHIES (2004).

- Consider the Poisson equation (1.4) in $\Omega = (0,1)^2$ with $\varepsilon = 1$. The right hand side and the Dirichlet boundary conditions are chosen such that $u(x,y) = \sin(\pi x)\sin(\pi y)$ is the prescribed solution, see Figure 1.1 Hence, this solution satisfies homogeneous Dirichlet boundary conditions. Denote by $u_h(x,y)$ the computed solution, where $h$ indicates the refinement of a mesh in $\Omega$. The errors obtained on successively refined meshes with the simplest finite element method are presented in Table 1.1.



Figure 1.1: Solution of the two-dimensional example of Example 1.8.

One can observe in Table 1.1 that $\|u - u_h\|_{L^2(\Omega)}$ converges with second order and $\|\nabla(u - u_h)\|_{L^2(\Omega)}$ converges with first order. A main topic of the numerical analysis of discretizations for partial differential equations consists in showing that the computed solution converges to the solution of an appropriate continuous problem in appropriate norms. In addition, to prove a certain order of convergence (in the asymptotic regime) is of interest.

Table 1.1: Example 1.8, two-dimensional example.

| $h$ | degrees of freedom | $\|u - u_h\|_{L^2(\Omega)}$ | $\|\nabla(u - u_h)\|_{L^2(\Omega)}$ |
|---|---|---|---|
| 1/4 | 25 | 8.522e-2 | 8.391e-1 |
| 1/8 | 81 | 2.256e-2 | 4.318e-1 |
| 1/16 | 289 | 5.726e-3 | 2.175e-1 |
| 1/32 | 1089 | 1.437e-3 | 1.089e-1 |
| 1/64 | 4225 | 3.596e-4 | 5.451e-2 |
| 1/128 | 16641 | 8.993e-5 | 2.726e-2 |
| 1/256 | 66049 | 2.248e-5 | 1.363e-2 |
| 1/512 | 263169 | 5.621e-6 | 6.815e-3 |

- Consider the Poisson equation (1.4) in $\Omega = (0,1)^3$ with $\varepsilon = 1$ and $f = 0$. At $z = 1$ the temperature profile should be $u(x, y, 1) = 16x(1 - x)y(1 - y)$ and at the opposite wall should be cooled $u(x, y, 0) = 0$. At all other walls, there should be an undisturbed temperature flux $\frac{\partial u}{\partial \mathbf{n}}(x, y, z) = 0$. A approximation of the solution computed with a finite element method is presented in Figure 1.2.



Figure 1.2: Contour lines of the solution of the three-dimensional example of Example 1.8.

The analytical solution is not known in this example (or it maybe hard to compute). It is important for applications that one obtains, e.g., good visualizations of the solution or approximate values for quantities of interest. One knows by the general theory that the computed solution converges to the solution of the continuous problem in appropriate norms and one hopes that the computed solution is already sufficiently close.

□

## 1.2 The Diffusion Equation

**Remark 1.9** *Derivation.* Diffusion is the transport of a species caused by the movement of particles. Instead of Fourier's law, Newton's law for the particle flux through $\partial V$ per time unit is used

$$dQ = -D\nabla u \cdot \mathbf{n} \, d\mathbf{s}$$

with

- $u(t, \mathbf{x})$ – particle density, concentration with unit $[mol/m^3]$,
- $D(t, \mathbf{x})$ – diffusion coefficient with unit $[m^2/s]$.

The derivation of the diffusion equation proceeds in the same way as for the heat equation. It has the form

$$c\frac{\partial u}{\partial t} - \nabla \cdot (D\nabla u) + qu = F \quad \text{in } (0,T) \times \Omega, \qquad (1.6)$$

where
- $c(t,\mathbf{x})$ – is the porosity of the species, $[\cdot]$,
- $q(t,\mathbf{x})$ – is the absorption coefficient of the species with unit $[1/s]$,
- $F(t,\mathbf{x})$ – describes sources and sinks, $[mol/(s\ m^3)]$.

The porosity and the absorption coefficient are positive functions. To obtain a well posed problem, an initial condition and boundary conditions are necessary.

If the concentration is constant in time, one obtains

$$-\nabla \cdot (D\nabla u) + qu = F \quad \text{in } \Omega. \qquad (1.7)$$

Hence, the diffusion equation possesses a similar form as the heat equation. □

## 1.3 The Navier–Stokes Equations

**Remark 1.10** *Generalities.* The Navier[6]–Stokes[7] equations are the fundamental equations of fluid dynamics. In this section, a viscous fluid (with internal friction) with constant density (incompressible) will be considered. □

**Remark 1.11** *Conservation of mass.* The first basic principle of the flow of an incompressible fluid is the conservation of mass. Let $V$ be an arbitrary volume. Then, the change of fluid in $V$ satisfies

$$-\underbrace{\frac{\partial}{\partial t}\int_V \rho\, d\mathbf{x}}_{\text{change}} = \underbrace{\int_{\partial V} \rho\mathbf{v}\cdot\mathbf{n}\, d\mathbf{s}}_{\text{flux through the boundary of } V} = \int_V \nabla\cdot(\rho\mathbf{v})\, d\mathbf{x},$$

where
- $\mathbf{v}(t,\mathbf{x})$ – velocity $(v_1, v_2, v_3)^T$ at time $t$ and at point $\mathbf{x}$ with unit $[m/s]$,
- $\rho$ – density of the fluid, $[kg/m^3]$.

Since $V$ is arbitrary, the terms in the volume integrals have to be the same. One gets the so-called continuity equation

$$\rho_t + \nabla \cdot (\rho\mathbf{v}) = 0 \quad \text{in } (0,T) \times \Omega.$$

Since $\rho$ is constant, one obtains the first equation of the Navier–Stokes equation, the so-called incompressibility constraint,

$$\nabla \cdot \mathbf{v} = 0 \quad \text{in } (0,T) \times \Omega. \qquad (1.8)$$

□

**Remark 1.12** *Conservation of linear momentum.* The second equation of the Navier–Stokes equations represents Newton's second law of motion

$$\text{net force} = \text{mass} \times \text{acceleration}.$$

It states that the rate of change of the linear momentum must be equal to the net force acting on a collection of fluid particles.

---

[6]Claude Louis Marie Henri Navier (1785 - 1836)
[7]George Gabriel Stokes (1819 - 1903)

The forces acting on an arbitrary volume $V$ are given by

$$F_V = \underbrace{\int_{\partial V} -P\mathbf{n}\,d\mathbf{s}}_{\text{outer pressure}} + \underbrace{\int_{\partial V} \mathbb{S}'\mathbf{n}\,d\mathbf{s}}_{\text{friction}} + \underbrace{\int_V \rho\mathbf{g}\,d\mathbf{x}}_{\text{gravitation}},$$

where
- $S'(t,\mathbf{x})$ – stress tensor with unit $[N/m^2]$,
- $P(t,\mathbf{x})$ – the pressure with unit $[N/m^2]$,
- $\mathbf{g}(t,\mathbf{x})$ – standard gravity (directed), $[m/s^2]$.

The pressure possesses a negative sign since it is directed into $V$, whereas the stress acts outwardly.

The integral on $\partial V$ can be transformed into an integral on $V$ with integration by parts. One obtains the force per unit volume

$$-\nabla P + \nabla \cdot \mathbb{S}' + \rho\mathbf{g}.$$

On the basis of physical considerations (Landau and Lifschitz, 1966, p. 53), one uses the following ansatz for the stress tensor

$$\mathbb{S}' = \eta\left(\nabla\mathbf{v} + \nabla\mathbf{v}^T - \frac{2}{3}(\nabla\cdot\mathbf{v})\mathbb{I}\right) + \zeta(\nabla\cdot\mathbf{v})\mathbb{I},$$

where
- $\eta$ – first order viscosity of the fluid, $[kg/(m\ s)]$,
- $\zeta$ – second order viscosity of the fluid, $[kg/(m\ s)]$,
- $\mathbb{I}$ – unit tensor.

For Newton's second law of motion one considers the movement of particles with velocity $\mathbf{v}(t,\mathbf{x}(t))$. One obtains the following equation

$$\underbrace{-\nabla P + \nabla\cdot\mathbb{S}' + \rho\mathbf{g}}_{\text{force per unit volume}} = \underbrace{\rho}_{\text{mass per unit volume}} \underbrace{\frac{d\mathbf{v}(t,\mathbf{x}(t))}{dt}}_{\text{acceleration}}$$

$$= \rho\left(\mathbf{v}_t + (\mathbf{v}\cdot\nabla)\mathbf{v}\right).$$

The second formula was obtained with the chain rule. The detailed form of the second term is

$$(\mathbf{v}\cdot\nabla)\mathbf{v} = \begin{pmatrix} v_1(v_1)_x + v_2(v_1)_y + v_3(v_1)_z \\ v_1(v_2)_x + v_2(v_2)_y + v_3(v_2)_z \\ v_1(v_3)_x + v_2(v_3)_y + v_3(v_3)_z \end{pmatrix}.$$

If both viscosities are constant, one gets

$$\frac{\partial\mathbf{v}}{\partial t} - \nu\Delta\mathbf{v} + (\mathbf{v}\cdot\nabla)\mathbf{v} - \frac{\nabla P}{\rho} = \mathbf{g} + \frac{1}{\rho}\left(\frac{\eta}{3} + \zeta\right)\nabla(\nabla\cdot\mathbf{v}),$$

where $\nu = \eta/\rho, [m^2/s]$ is the kinematic viscosity. The second term on the right hand side vanishes because of the incompressibility constraint (1.8).

One obtains the dimensional Navier–Stokes equations

$$\frac{\partial\mathbf{v}}{\partial t} - \nu\Delta\mathbf{v} + (\mathbf{v}\cdot\nabla)\mathbf{v} - \frac{\nabla P}{\rho} = \mathbf{g}, \quad \nabla\cdot\mathbf{v} = 0 \ \text{ in } (0,T)\times\Omega.$$

$\square$

**Remark 1.13** *Non-dimensional Navier–Stokes equations.* The final step in the modeling process is the derivation of non-dimensional equations. Let

- $L$ – a characteristic length scale of the problem, $[m]$,
- $U$ – a characteristic velocity scale of the problem, $[m/s]$,
- $T^*$ – a characteristic time scale of the problem, $[s]$.

Denoting here the old coordinates with a prime, one obtains with the transformations

$$\mathbf{x} = \frac{\mathbf{x}'}{L}, \quad \mathbf{u} = \frac{\mathbf{v}}{U}, \quad t = \frac{t'}{T^*}$$

the non-dimensional equations

$$\frac{L}{UT^*}\partial_t \mathbf{u} - \frac{\nu}{UL}\Delta \mathbf{u} + (\mathbf{u} \cdot \nabla)\mathbf{u} + \nabla p = \mathbf{f}, \quad \nabla \cdot \mathbf{u} = 0 \;\; \text{in } (0,T) \times \Omega,$$

with the redefined pressure and the new right hand side

$$p(t,\mathbf{x}) = \frac{P}{\rho U^2}(t,\mathbf{x}), \quad \mathbf{f}(t,\mathbf{x}) = \frac{L\mathbf{g}}{U^2}(t,\mathbf{x}).$$

This equation has two dimensionless characteristic parameters: the Strouhal[8] number $St$ and the Reynolds [9] number $Re$

$$St := \frac{L}{UT^*}, \quad Re := \frac{UL}{\nu}.$$

Setting $T^* = L/U$, one obtains the form of the incompressible Navier–Stokes equations which can be found in the literature

$$\begin{aligned}
\frac{\partial \mathbf{u}}{\partial t} - Re^{-1}\Delta \mathbf{u} + (\mathbf{u} \cdot \nabla)\mathbf{u} + \nabla p &= \mathbf{f} &&\text{in } (0,T) \times \Omega, \\
\nabla \cdot \mathbf{u} &= 0 &&\text{in } [0,T) \times \Omega.
\end{aligned}$$

$\square$

**Remark 1.14** *About the incompressible Navier–Stokes equations.* The Navier–Stokes equations are not yet understood completely. For instance, the existence of an appropriately defined classical solution for $\Omega \subset \mathbb{R}^3$ is not clear. This problem is among the so-called millennium problems of mathematics Fefferman (2000) and its answer is worth one million dollar. Also the numerical methods for solving the Navier–Stokes equations are by far not developed sufficiently well as it is required by many applications, e.g. for turbulent flows in weather prediction. $\square$

**Remark 1.15** *Slow flows.* Am important special case is the case of slow flows which lead to a stationary (independent of time) flow field. In this case, the first term in the in the momentum balance equation vanish. In addition, if the flow is very slow, the nonlinear term can be neglected. One gets the so-called Stokes equations

$$\begin{aligned}
-Re^{-1}\Delta \mathbf{u} + \nabla p &= \mathbf{f} &&\text{in } \Omega, \\
\nabla \cdot \mathbf{u} &= 0 &&\text{in } \Omega.
\end{aligned}$$

$\square$

## 1.4 Classification of Second Order Partial Differential Equations

**Definition 1.16 Quasi-linear and linear second order partial differential equation.** Let $\Omega \subset \mathbb{R}^d$, $d \in \mathbb{N}$. A quasi-linear second order partial differential

---

[8]Čeněk Strouhal (1850 – 1923)
[9]Osborne Reynolds (1842 - 1912)

equation defined on $\Omega$ has the form

$$\sum_{j,k=1}^{d} a_{jk}(\mathbf{x})\partial_j\partial_k u + F\left(\mathbf{x}, u, \partial_1 u, \ldots, \partial_d u\right) = 0 \qquad (1.9)$$

or in nabla notation

$$\nabla \cdot A(\mathbf{x})\nabla u + \tilde{F}\left(\mathbf{x}, u, \partial_1 u, \ldots, \partial_d u\right) = 0.$$

A linear second order partial differential equation has the form

$$\sum_{j,k=1}^{d} a_{jk}(\mathbf{x})\partial_j\partial_k u + \mathbf{b}(\mathbf{x}) \cdot \nabla u + c(\mathbf{x})u = F(\mathbf{x}).$$

$\square$

**Remark 1.17** *The matrix of the second order operator.* If $u(\mathbf{x})$ is sufficiently regular, then the application of the Schwarz'[10] theorem yields $\partial_j\partial_k u(\mathbf{x}) = \partial_k\partial_j u(\mathbf{x})$. It follows that equation (1.9) contains the coefficient $\partial_j\partial_k u(\mathbf{x})$ twice, namely in $a_{jk}(\mathbf{x})$ and $a_{kj}(\mathbf{x})$. For definiteness, one requires that

$$a_{jk}(\mathbf{x}) = a_{kj}(\mathbf{x}).$$

Now, one can write the coefficient of the second order derivative with the symmetric matrix

$$A(\mathbf{x}) = \begin{pmatrix} a_{11}(\mathbf{x}) & \cdots & a_{1d}(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ a_{d1}(\mathbf{x}) & \cdots & a_{dd}(\mathbf{x}) \end{pmatrix}.$$

All eigenvalues of this matrix are real and the classification of quasi-linear second order partial differential equations is based on these eigenvalues. $\square$

**Definition 1.18 Classification of quasi-linear second order partial differential equation.** On a subset $\tilde{\Omega} \subset \Omega$ let $\alpha$ be the number of positive eigenvalues of $A(\mathbf{x})$, $\beta$ be the number of negative eigenvalues, and $\gamma$ be the multiplicity of the eigenvalue zero. The quasi-linear second order partial differential equation (1.9) is said to be of type $(\alpha, \beta, \gamma)$ on $\tilde{\Omega}$. It is called to be
• elliptic on $\tilde{\Omega}$ if it is of type $(d, 0, 0) = (0, d, 0)$,
• hyperbolic on $\tilde{\Omega}$, if its type is $(d-1, 1, 0) = (1, d-1, 0)$,
• parabolic on $\tilde{\Omega}$, if it is of type $(d-1, 0, 1) = (0, d-1, 1)$.
In the case of linear partial differential equations, one speaks of a parabolic equation if in addition to the requirement from above it holds that

$$\mathrm{rank}(A(\mathbf{x}), \mathbf{b}(\mathbf{x})) = d$$

in $\tilde{\Omega}$. $\square$

**Remark 1.19** *Other cases.* Definition 1.18 does not cover all possible cases. However, the other cases are only of little interest in practice. $\square$

**Example 1.20** *Types of second order partial differential equations.*
• For the Poisson equation (1.4) one has $a_{ii} = -\varepsilon^2 < 0$ and $a_{ij} = 0$ for $i \neq j$. It follows that all eigenvalues of $A$ are negative and the Poisson equation is an elliptic partial differential equation. The same reasoning can be applied to the stationary diffusion equation (1.7).

---
[10]Hermann Amandus Schwarz (1843 – 1921)

- In the heat equation (1.3) there is besides the spatial derivatives also the temporal derivative. The derivative in time has to be taken into account in the definition of the matrix $A$. Since this derivative is only of first order, one obtains in $A$ a zero row and a zero column. One has, e.g., $a_{ii} = -\varepsilon^2 < 0, i = 2, \ldots, d + 1$, $a_{11} = 0$, and $a_{ij} = 0$ for $i \neq j$. It follows that one eigenvalue is zero and the others have the same sign. The vector of the first order term has the form $\mathbf{b} = (1, 0, \ldots, 0)^T \in \mathbb{R}^{d+1}$, where the one comes from the $\partial_t u(t, \mathbf{x})$. Now, one can see immediately that $(A, \mathbf{b})$ possesses full column rank. Hence, (1.3) is a parabolic partial differential equation.
- An example for a hyperbolic partial differential equation is the wave equation

$$\partial_t^2 u - \varepsilon^2 \Delta u = f \quad \text{in } (0, T) \times \Omega.$$

$\square$

## 1.5 Literature

**Remark 1.21** *Some books about the topic of this class.* Books about finite difference methods are
- Samarskij (1984), classic book, the English version is Samarskii (2001)
- LeVeque (2007)

Much more books can be found about finite element methods
- Ciarlet (2002), classic text,
- Strang and Fix (2008), classic text,
- Braess (2001), very popular book in Germany,
- Brenner and Scott (2008), rather abstract treatment, from the point of view of functional analysis,
- Ern and Guermond (2004), modern comprehensive book,
- Grossmann and Roos (2007)
- Šolín (2006), written by somebody who worked a lot in the implementation of the methods,
- Goering et al. (2010), introductory text, good for beginners,
- Deuflhard and Weiser (2012), strong emphasis on adaptive methods
- Dziuk (2010).

These lists are not complete.

These lectures notes are based in some parts on lecture notes from Sergej Rjasanow (Saarbrücken) and Manfred Dobrowolski (Würzburg). $\square$

# Chapter 2

# Finite Difference Methods for Elliptic Equations

**Remark 2.1** *Model problem.* The model problem in this chapter is the Poisson equation with Dirichlet boundary conditions

$$
\begin{aligned}
-\Delta u &= f &&\text{in } \Omega, \\
u &= g &&\text{on } \partial\Omega,
\end{aligned}
\tag{2.1}
$$

where $\Omega \subset \mathbb{R}^2$. This chapter follows in wide parts Samarskij (1984). □

## 2.1 Basics on Finite Differences

**Remark 2.2** *Grid.* This section considers the one-dimensional case. Consider the interval $[0, 1]$ which is decomposed by an equidistant grid

$$
\begin{aligned}
x_i &= ih, \quad i = 0, \ldots, n, \quad h = 1/n, \ -\text{ nodes}, \\
\omega_h &= \{x_i \ : \ i = 0, \ldots, n\} \ -\text{ grid}.
\end{aligned}
$$

□

**Definition 2.3 Grid function.** A vector $\mathbf{u}_h = (u_0, \ldots, u_n)^T \in \mathbb{R}^{n+1}$ which assigns every grid point a function value is called grid function. □

**Definition 2.4 Finite differences.** Let $v(x)$ be a sufficiently smooth function and denote by $v_i = v(x_i)$, where $x_i$ are the nodes of the grid. The following quotients are called

$$
\begin{aligned}
v_{x,i} &= \frac{v_{i+1} - v_i}{h} \ -\text{ forward difference}, \\
v_{\overline{x},i} &= \frac{v_i - v_{i-1}}{h} \ -\text{ backward difference}, \\
v_{\mathring{x},i} &= \frac{v_{i+1} - v_{i-1}}{2h} \ -\text{ central difference}, \\
v_{\overline{x}x,i} &= \frac{v_{i+1} - 2v_i + v_{i-1}}{h^2} \ -\text{ second order difference},
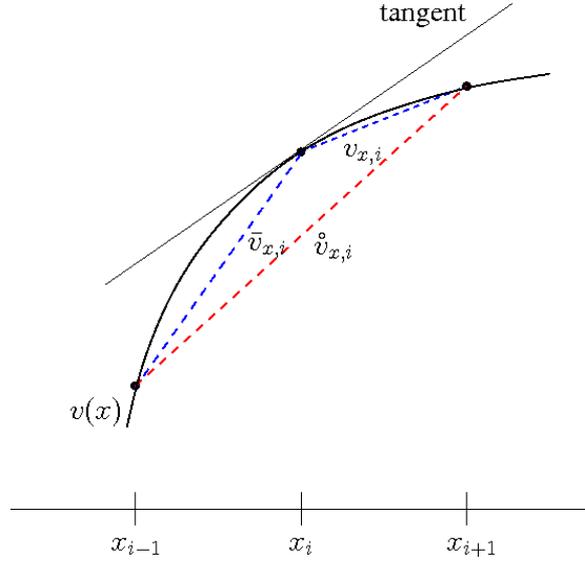\end{aligned}
$$

see Figure 2.1. □

Figure 2.1: Illustration of the finite differences.

**Remark 2.5** *Some properties of the finite differences.* It is (*exercise*)

$$v_{\mathring{x},i} = \frac{1}{2}(v_{x,i} + v_{\overline{x},i}), \quad v_{\overline{x}x,i} = (v_{\overline{x},i})_{x,i}.$$

Using the Taylor series expansion for $v(x)$ at the node $x_i$, one gets (*exercise*)

$$
\begin{aligned}
v_{x,i} &= v'(x_i) + \frac{1}{2}hv''(x_i) + \mathcal{O}\left(h^2\right), \\
v_{\overline{x},i} &= v'(x_i) - \frac{1}{2}hv''(x_i) + \mathcal{O}\left(h^2\right), \\
v_{\mathring{x},i} &= v'(x_i) + \mathcal{O}\left(h^2\right), \\
v_{\overline{x}x,i} &= v''(x_i) + \mathcal{O}\left(h^2\right).
\end{aligned}
$$

$\square$

**Definition 2.6 Consistent difference operator.** Let $L$ be a differential operator. The difference operator $L_h : \mathbb{R}^{n+1} \to \mathbb{R}^{n+1}$ is called consistent with $L$ of order $k$ if

$$\max_{0 \le i \le n} |(Lu)(x_i) - (L_h u_h)_i| = \|(Lu)(x_i) - (L_h u_h)_i\|_{\infty, \omega_h} = \mathcal{O}\left(h^k\right)$$

for all sufficiently smooth functions $u(x)$. $\square$

**Example 2.7** *Consistency orders.* The order of consistency measures the quality of approximation of $L$ by $L_h$.

The difference operators $v_{x,i}, v_{\overline{x},i}, v_{\mathring{x},i}$ are consistent to $L = \frac{d}{dx}$ with order $1, 1$, and $2$, respectively. The operator $v_{\overline{x}x,i}$ is consistent of second order to $L = \frac{d^2}{dx^2}$, see Remark 2.5. $\square$

**Example 2.8** *Approximation of a more complicated differential operator by difference operators.* Consider the differential operator

$$Lu = \frac{d}{dx}\left(k(x)\frac{du}{dx}\right),$$

14

where $k(x)$ is assumed to be continuously differentiable. Define the difference operator $L_h$ as follows

$$
\begin{aligned}
(L_h u_h)_i &= (a u_{\overline{x},i})_{x,i} = \frac{1}{h}\Big(a(x_{i+1})u_{\overline{x},i}(x_{i+1}) - a(x_i)u_{\overline{x},i}(x_i)\Big) \\
&= \frac{1}{h}\left(a_{i+1}\frac{u_{i+1} - u_i}{h} - a_i\frac{u_i - u_{i-1}}{h}\right),
\end{aligned}
$$

where $a$ is a grid function which has to be determined appropriately. One gets with the product rule

$$
(Lu)_i = k'(x_i)(u')_i + k(x_i)(u'')_i
$$

and a Taylor series expansion for $u_{i-1}$, $u_{i+1}$

$$
(L_h u_h)_i = \frac{a_{i+1} - a_i}{h}(u')_i + \frac{a_{i+1} + a_i}{2}(u'')_i + \frac{h(a_{i+1} - a_i)}{6}(u''')_i + \mathcal{O}\left(h^2\right).
$$

Thus, the difference of the differential operator and the difference operator is

$$
\begin{aligned}
(Lu)_i - (L_h u_h)_i &= \left(k'(x_i) - \frac{a_{i+1} - a_i}{h}\right)(u')_i + \left(k(x_i) - \frac{a_{i+1} + a_i}{2}\right)(u'')_i \\
&\quad - \frac{h(a_{i+1} - a_i)}{6}(u''')_i + \mathcal{O}\left(h^2\right).
\end{aligned} \tag{2.2}
$$

In order to define $L_h$ such that it is consistent of second order to $L$, one has to satisfy the following two conditions

$$
\frac{a_{i+1} - a_i}{h} = k'(x_i) + \mathcal{O}\left(h^2\right), \qquad \frac{a_{i+1} + a_i}{2} = k(x_i) + \mathcal{O}\left(h^2\right).
$$

From the first requirement, it follows that $a_{i+1} - a_i = \mathcal{O}(h)$. Hence, the third term in the consistency error equation (2.2) is of order $\mathcal{O}\left(h^2\right)$. Possible choices for the grid function are (*exercise*)

$$
a_i = \frac{k_i + k_{i-1}}{2}, \quad a_i = k\left(x_i - \frac{h}{2}\right), \quad a_i = (k_i k_{i-1})^{1/2}.
$$

Note that the 'natural' choice, $a_i = k_i$, leads only to first order consistency. (*exercise*) $\qquad\square$

## 2.2 Finite Difference Approximation of the Laplacian in Two Dimensions

**Remark 2.9** *The five point stencil.* The Laplacian in two dimensions is defined by

$$
\Delta u(\mathbf{x}) = \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = \partial_x^2 u + \partial_y^2 u = u_{xx} + u_{yy}, \quad \mathbf{x} = (x, y).
$$

The simplest approximation uses for both second order derivatives the second order differences. One obtains the so-called five point stencil and the approximation

$$
\Delta u \approx \Lambda u = u_{\overline{x}x} + u_{\overline{y}y} = \frac{u_{i+1,j} - 2u_{ij} + u_{i-1,j}}{h_x^2} + \frac{u_{i,j+1} - 2u_{ij} + u_{i,j-1}}{h_y^2}, \tag{2.3}
$$

see Figure 2.2. From the consistency order of the second order difference it follows immediately that $\Lambda u$ approximates the Laplacian of order $\mathcal{O}\left(h_x^2 + h_y^2\right)$. $\qquad\square$
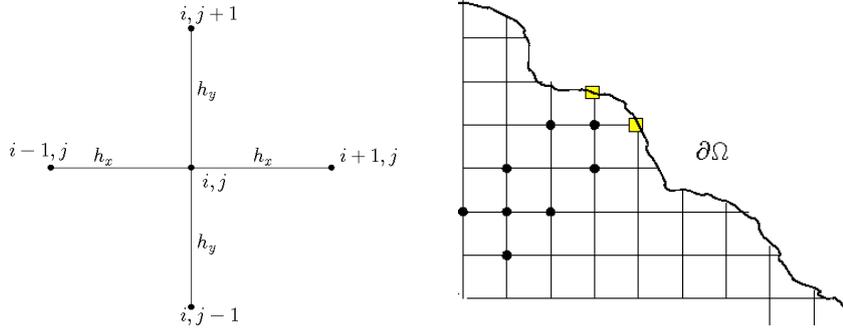
Figure 2.2: Five point stencils.

**Remark 2.10** *The five point stencil on curvilinear boundaries.* There is a difficulty if the five point stencil is used in domains with curvilinear boundaries. The approximation of the second derivative requires three function values in each coordinate direction

$$(x - h_x^-, y), (x, y), (x + h_x^+, y),$$
$$(x, y - h_y^-), (x, y), (x, y + h_y^+),$$

see Figure 2.3. A guideline of defining the approximation is that the five point stencil is recovered in the case $h_x^- = h_x^+$. A possible approximation of this type is

$$\frac{\partial^2 u}{\partial x^2} \approx \frac{1}{\overline{h}_x} \left( \frac{u(x + h_x^+, y) - u(x, y)}{h_x^+} - \frac{u(x, y) - u(x - h_x^-, y)}{h_x^-} \right) \qquad (2.4)$$

with $\overline{h}_x = (h_x^+ + h_x^-)/2$. Using a Taylor series expansion, one finds that the error of this approximation is

$$\frac{\partial^2 u}{\partial x^2} - \frac{1}{\overline{h}_x} \left( \frac{u(x + h_x^+, y) - u(x, y)}{h_x^+} - \frac{u(x, y) - u(x - h_x^-, y)}{h_x^-} \right)$$
$$= -\frac{1}{3}(h_x^+ - h_x^-) \frac{\partial^3 u}{\partial x^3} + \mathcal{O}\left( \overline{h}_x^2 \right).$$

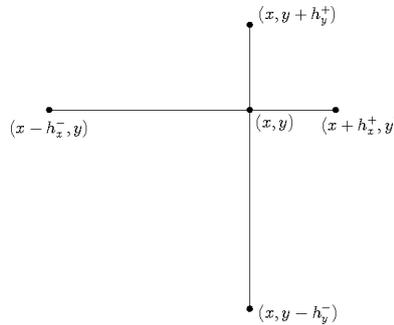For $h_x^+ \neq h_x^-$, this approximation is of first order.



Figure 2.3: Sketch to Remark 2.10.

A different way consists in using

$$\frac{\partial^2 u}{\partial x^2} \approx \frac{1}{\tilde{h}_x} \left( \frac{u(x + h_x^+, y) - u(x, y)}{h_x^+} - \frac{u(x, y) - u(x - h_x^-, y)}{h_x^-} \right)$$

with $\tilde{h}_x = \max\{h_x^+, h_x^-\}$. However, this approximation possesses only the order zero, i.e., there is actually no approximation.

Altogether, there is a loss of order of consistency in this situation. $\square$

**Example 2.11** *The Dirichlet problem.* Consider the Poisson equation which is equipped with Dirichlet boundary conditions (2.1). First, $\mathbb{R}^2$ is decomposed by a grid with rectangular mesh cells $x_i = ih_x, y_j = jh_y$, $h_x, h_y > 0$, $i, j \in \mathbb{Z}$. Denote by

$$
\begin{array}{rcll}
w_h^\circ & = & \{\circ\} & \text{inner nodes, five point stencil completely in } \Omega, \\
w_h^* & = & \{*\} & \text{inner nodes that are close to the boundary,} \\
\gamma_h & = & \{*\} & \text{boundary nodes,} \\
\omega_h & = & w_h^\circ \cup w_h^* & \text{inner nodes,} \\
\omega_h \cup \gamma_h & & & \text{grid,}
\end{array}
$$

see Figure 2.4.



Figure 2.4: Different types of nodes in the grid.

The finite difference approximation of problem (2.1) which will be studied in the following consists in finding a mesh function $u(\mathbf{x})$ such that

$$
\begin{array}{rcll}
-\Lambda u(\mathbf{x}) & = & \phi(\mathbf{x}) & \mathbf{x} \in w_h^\circ, \\
-\Lambda^* u(\mathbf{x}) & = & \phi(\mathbf{x}) & \mathbf{x} \in w_h^*, \\
u(\mathbf{x}) & = & g(\mathbf{x}) & \mathbf{x} \in \gamma_h,
\end{array}
\tag{2.5}
$$

where $\phi(\mathbf{x})$ is a grid function that approximates $f(\mathbf{x})$ and $\Lambda^*$ is an approximation of the Laplacian for nodes that are close to the boundary, e.g., defined by (2.4). The discrete problem is a large sparse linear system of equations. The most important questions are:

- Which properties possesses the solution of (2.5)?
- Converges the solution of (2.5) to the solution of the Poisson problem and if yes, with which order?

$\square$

## 2.3 The Discrete Maximum Principle for a Finite Difference Operator

**Remark 2.12** *Contents of this section.* Solutions of the Laplace equation, i.e., of (2.1) with $f(\mathbf{x}) = 0$, fulfill so-called maximum principles. This section shows, that the finite difference approximation of an operator, where the five point stencil of the Laplacian is a special case, satisfies a discrete analog of one of the maximum principles. $\square$

**Theorem 2.13 Maximum principles for harmonic functions.** *Let $\Omega \subset \mathbb{R}^d$ be a bounded domain and $u \in C^2(\Omega) \cap C(\overline{\Omega})$ harmonic in $\Omega$, i.e. $u(\mathbf{x})$ solves the Laplace equation $-\Delta u = 0$ in $\Omega$.*

- *Weak maximum principle. It holds*

$$\max_{\mathbf{x} \in \overline{\Omega}} u(\mathbf{x}) = \max_{\mathbf{x} \in \partial\Omega} u(\mathbf{x}).$$

  *That means, $u(\mathbf{x})$ takes it maximal value at the boundary.*
- *Strong maximum principle. If $\Omega$ is connected and if the maximum is taken in $\Omega$ (note that $\Omega$ is open), i.e., $u(\mathbf{x}_0) = \max_{\mathbf{x} \in \overline{\Omega}} u(\mathbf{x})$ for a point $\mathbf{x}_0 \in \Omega$, then $u(\mathbf{x})$ is constant*

$$u(\mathbf{x}) = \max_{\mathbf{x} \in \overline{\Omega}} u(\mathbf{x}) = u(\mathbf{x}_0) \quad \forall\ \mathbf{x} \in \overline{\Omega}.$$

**Proof:** See the literature, e.g., (Evans, 2010, p. 27, Theorem 4) or the class on the theory of partial differential equations. ∎

**Remark 2.14** *Interpretation of the maximum principle.* The Laplace equation models the temperature distribution of a heated body without heat sources in $\Omega$. Then, the weak maximum principle just states that the temperature in the interior of the body cannot be higher than the highest temperature at the boundary.

There are maximum principles also for more complicated operators than the Laplacian, e.g., see Evans (2010).

Since the solution of the partial differential equation will be only approximated by a discretization like a finite difference method, one has to expect that basic physical properties are satisfied by the numerical solution also only approximately. However, in applications, it is often very important that such properties are satisfied exactly. □

**Remark 2.15** *The difference equation.* In this section, a difference equation of the form

$$a(\mathbf{x})u(\mathbf{x}) = \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})u(\mathbf{y}) + F(\mathbf{x}),\ \mathbf{x} \in \omega_h \cup \gamma_h, \tag{2.6}$$

will be considered. In (2.6), for each node $\mathbf{x}$, the set $S(\mathbf{x})$ is the set of all nodes on which the sum has to be performed, $\mathbf{x} \notin S(\mathbf{x})$. That means, $a(\mathbf{x})$ describes the contribution of the finite difference scheme of a node $\mathbf{x}$ to itself and $b(\mathbf{x}, \mathbf{y})$ describes the contributions from the neighbors.

It will be assumed that the grid $\omega_h$ is connected, i.e., for all $\mathbf{x}_a, \mathbf{x}_e \in \omega_h$ exist $\mathbf{x}_1, \ldots, \mathbf{x}_m \in \omega_h$ with $\mathbf{x}_1 \in S(\mathbf{x}_a), \mathbf{x}_2 \in S(\mathbf{x}_1), \ldots, \mathbf{x}_e \in S(\mathbf{x}_m)$. E.g., the situation depicted in Figure 2.5 is not allowed.



Figure 2.5: Grid which is not allowed in Section 2.3.

It will be assumed that the coefficients $a(\mathbf{x})$ and $b(\mathbf{x}, \mathbf{y})$ satisfy the following conditions:

$$\begin{aligned} a(\mathbf{x}) &> 0,\ b(\mathbf{x}, \mathbf{y}) > 0,\ \forall\ \mathbf{x} \in \omega_h, \forall\ \mathbf{y} \in S(\mathbf{x}), \\ a(\mathbf{x}) &= 1,\ b(\mathbf{x}, \mathbf{y}) = 0\ \forall\ \mathbf{x} \in \gamma_h\ \text{(Dirichlet boundary condition)}. \end{aligned}$$

The values of the Dirichlet boundary condition are incorporated in (2.6) into the function $F(\mathbf{x})$. □

**Example 2.16** *Five point stencil for approximating the Laplacian.* Inserting the approximation of the Laplacian with the five point stencil (2.3) for $\mathbf{x} = (x, y) \in \omega_h^\circ$ into the scheme (2.6) gives

$$
\frac{2(h_x^2 + h_y^2)}{h_x^2 h_y^2} u(x, y) = \left[ \frac{1}{h_x^2} u(x + h_x, y) + \frac{1}{h_x^2} u(x - h_x, y) \right.
$$
$$
\left. + \frac{1}{h_y^2} u(x, y + h_y) + \frac{1}{h_y^2} u(x, y - h_y) \right] + \phi(x, y).
$$

It follows that

$$
a(\mathbf{x}) = \frac{2(h_x^2 + h_y^2)}{h_x^2 h_y^2},
$$
$$
b(\mathbf{x}, \mathbf{y}) \in \{h_x^{-2}, h_y^{-2}\},
$$
$$
S(\mathbf{x}) = \{(x - h_x, y), (x + h_x, y), (x, y - h_y), (x, y + h_y)\}.
$$

For inner nodes that are close to the boundary, only the one-dimensional case (2.4) will be considered for simplicity. Let $x + h_x^+ \in \gamma_h$, then it follows by inserting (2.4) into (2.6)

$$
\frac{1}{\overline{h}_x} \left( \frac{1}{h_x^+} + \frac{1}{h_x^-} \right) u(x, y) = \frac{u(x - h_x^-, y)}{\overline{h}_x h_x^-} + \underbrace{\frac{u(x + h_x^+, y)}{\overline{h}_x h_x^+}}_{\text{on } \gamma_h \to F(x)} + \phi(x),
$$

where $a(x) = \frac{1}{\overline{h}_x} \left( \frac{1}{h_x^+} + \frac{1}{h_x^-} \right)$, $b(x, y) = \frac{1}{\overline{h}_x h_x^-}$ und $S(x) = \{(x - h_x^-, y)\}$. □

**Remark 2.17** *Reformulation of the difference scheme.* Scheme (2.6) can be reformulated in the form

$$
d(\mathbf{x})u(\mathbf{x}) = \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})\big(u(\mathbf{y}) - u(\mathbf{x})\big) + F(\mathbf{x}) \tag{2.7}
$$

with $d(\mathbf{x}) = a(\mathbf{x}) - \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})$. □

**Example 2.18** *Five point stencil for approximating the Laplacian.* Using the five point stencil for approximating the Laplacian, form (2.7) of the scheme is obtained with

$$
d(\mathbf{x}) = \frac{2(h_x^2 + h_y^2)}{h_x^2 h_y^2} - \frac{2}{h_x^2} - \frac{2}{h_y^2} = 0
$$

for $\mathbf{x} \in \omega_h^\circ$.

The coefficients $a(\mathbf{x})$ and $b(\mathbf{x}, \mathbf{y})$ are the weights of the finite difference stencil for approximating the Laplacian. A minimal condition for consistency is that this approximation vanishes for constant functions. It follows that also for the nodes $\mathbf{x} \in \omega_h^*$ it is $a(\mathbf{x}) = \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})$. However, as it was shown in Example 2.16, in this case the contributions from the neighbors on $\gamma_h$ are included in the scheme (2.6) in $F(x)$. Hence, one obtains for nodes that are close to the boundary

$$
d(\mathbf{x}) = \underbrace{\sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})}_{=a(\mathbf{x})} - \sum_{\mathbf{y} \in S(\mathbf{x}), \mathbf{y} \notin \gamma_h} b(\mathbf{x}, \mathbf{y}) = \sum_{\mathbf{y} \in S(\mathbf{x}), \mathbf{y} \in \gamma_h} b(\mathbf{x}, \mathbf{y}). \tag{2.8}
$$

In the one-dimensional case, one has, by the definition of $\overline{h}_x$ and with $h_x^- = h_x \geq h_x^+$,

$$
\begin{aligned}
d(x) & = \frac{1}{\overline{h}_x}\left(\frac{1}{h_x^+} + \frac{1}{h_x^-}\right) - \frac{1}{\overline{h}_x h_x^-} = \frac{1}{\overline{h}_x h_x^+} = \frac{2}{h_x h_x^+ + h_x^+ h_x^+} \\
& \geq \frac{2}{h_x h_x + h_x h_x} = \frac{1}{h_x h_x} > 0.
\end{aligned}
$$

$\square$

**Lemma 2.19 Discrete maximum principle (DMP).** *Let $u(\mathbf{x}) \neq const$ on $\omega_h$ and $d(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \omega_h$. Then, it follows from*

$$
L_h u(\mathbf{x}) := d(\mathbf{x})u(\mathbf{x}) - \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y})\big(u(\mathbf{y}) - u(\mathbf{x})\big) \leq 0 \tag{2.9}
$$

*(or $L_h u(\mathbf{x}) \geq 0$, respectively) on $\omega_h$ that $u(\mathbf{x})$ does not possess a positive maximum (or negative minimum, respectively) on $\omega_h$.*

**Proof:** The proof is performed by contradiction. Let $L_h u(\mathbf{x}) \leq 0$ for all $\mathbf{x} \in \omega_h$ and assume that $u(\mathbf{x})$ has a positive maximum on $\omega_h$ at $\overline{\mathbf{x}}$, i.e., $u(\overline{\mathbf{x}}) = \max_{\mathbf{x} \in \omega_h} u(\mathbf{x}) > 0$. Then, the idea of the proof consists in showing that with these assumptions there is a node $\tilde{\mathbf{x}} \in \omega_h$ with $L_h u(\tilde{\mathbf{x}}) > 0$.

For the node $\overline{\mathbf{x}}$ it holds that

$$
L_h u(\overline{\mathbf{x}}) = d(\overline{\mathbf{x}})u(\overline{\mathbf{x}}) - \sum_{\mathbf{y} \in S(\overline{\mathbf{x}})} \underbrace{b(\overline{\mathbf{x}}, \mathbf{y})}_{>0} \underbrace{\big(u(\mathbf{y}) - u(\overline{\mathbf{x}})\big)}_{\leq 0 \text{ by definition of } \overline{\mathbf{x}}} \geq d(\overline{\mathbf{x}})u(\overline{\mathbf{x}}) \geq 0.
$$

Hence, it follows that $L_h u(\overline{\mathbf{x}}) = 0$ and, in particular, that $d(\overline{\mathbf{x}}) = 0$. All terms in the sum are non-positive. Consequently, if the sum should be zero, all terms have to be zero, too. Since it was assumed that $b(\overline{\mathbf{x}}, \mathbf{y})$ is positive, it must also hold

$$
u(\mathbf{y}) = u(\overline{\mathbf{x}}) \; \forall \; \mathbf{y} \in S(\mathbf{x}).
$$

From the assumption $u(\mathbf{x}) \neq const$ it follows that there exists a node $\hat{\mathbf{x}} \in \omega_h$ with $u(\overline{\mathbf{x}}) > u(\hat{\mathbf{x}})$. Because the grid is connected, there is a path $\overline{\mathbf{x}}, \mathbf{x}_1, \ldots, \mathbf{x}_m, \hat{\mathbf{x}}$ such that

$$
\begin{aligned}
\mathbf{x}_1 &\in S(\overline{\mathbf{x}}), \quad u(\mathbf{x}_1) = u(\overline{\mathbf{x}}), \\
\mathbf{x}_2 &\in S(\mathbf{x}_1), \quad u(\mathbf{x}_2) = u(\mathbf{x}_1) = u(\overline{\mathbf{x}}), \\
&\cdots \\
\hat{\mathbf{x}} &\in S(\mathbf{x}_m), \quad u(\mathbf{x}_m) = u(\mathbf{x}_{m-1}) = \ldots = u(\overline{\mathbf{x}}) > u(\hat{\mathbf{x}}).
\end{aligned}
$$

For the last node $\mathbf{x}_m$, for which $u(\mathbf{x})$ has the same value as for $\overline{\mathbf{x}}$, it holds that

$$
L_h u(\mathbf{x}_m) \geq \underbrace{d(\mathbf{x}_m)}_{\geq 0} \underbrace{u(\mathbf{x}_m)}_{>0} - \underbrace{b(\mathbf{x}_m, \hat{\mathbf{x}})}_{>0} \underbrace{\big(u(\hat{\mathbf{x}}) - u(\mathbf{x}_m)\big)}_{<0} > 0.
$$

Hence, the node $\mathbf{x}_m$ is the wanted node $\tilde{\mathbf{x}}$. ∎

**Corollary 2.20 Non-negativity of the grid function.** *Let $u(\mathbf{x}) \geq 0$ for $\mathbf{x} \in \gamma_h$ and $L_h u(\mathbf{x}) \geq 0$ on $\omega_h$. Then, the grid function $u(\mathbf{x})$ is non-negative for all $\mathbf{x} \in \omega_h \cup \gamma_h$.*

**Proof:** Assume there is a node $\overline{\mathbf{x}} \in \omega_h$ with $u(\overline{\mathbf{x}}) < 0$. Then, the grid function has a negative minimum on $\omega_h$, which is a contradiction to the discrete maximum principle. ∎

**Corollary 2.21 Unique solution of the discrete Laplace equation with homogeneous Dirichlet boundary conditions.** *The Laplace equation $L_h u(\mathbf{x}) = 0$ possesses only the trivial solution $u(\mathbf{x}) = 0$ for $\mathbf{x} \in \omega_h \cup \gamma_h$.*

**Proof:** The statement of the corollary follows by applying Corollary 2.20 and its analog for the non-positivity of the grid function if $u(\mathbf{x}) \leq 0$ for $\mathbf{x} \in \gamma_h$ and $L_h u(\mathbf{x}) \leq 0$ on $\omega_h$. Note that in the definition $L_h u(\mathbf{x}) = 0$ contains also the boundary values, which are homogeneous Dirichlet. ∎

**Corollary 2.22 Comparison lemma.** *Let*

$$
\begin{aligned}
L_h u(\mathbf{x}) &= f(\mathbf{x}) \quad \text{for } \mathbf{x} \in \omega_h; \quad u(\mathbf{x}) = g(\mathbf{x}) \text{ for } \mathbf{x} \in \gamma_h, \\
L_h \overline{u}(\mathbf{x}) &= \overline{f}(\mathbf{x}) \quad \text{for } \mathbf{x} \in \omega_h; \quad \overline{u}(\mathbf{x}) = \overline{g}(\mathbf{x}) \text{ for } \mathbf{x} \in \gamma_h,
\end{aligned}
$$

*with $|f(\mathbf{x})| \leq \overline{f}(\mathbf{x})$ and $|g(\mathbf{x})| \leq \overline{g}(\mathbf{x})$. Then it is $|u(\mathbf{x})| \leq \overline{u}(\mathbf{x})$ for all $\mathbf{x} \in \omega_h \cup \gamma_h$.*

**Proof:** Exercise. ∎

**Remark 2.23** *Remainder of this section.* The remaining corollaries presented in this section will be applied in the stability proof in Section 2.4. In this proof, the homogeneous problem (right hand side vanishes) and the problem with homogeneous Dirichlet boundary conditions will be analyzed separately. □

**Corollary 2.24 Homogeneous problem.** *For the solution of the problem*

$$
\begin{aligned}
L_h u(\mathbf{x}) &= 0, \quad &\mathbf{x} \in \omega_h, \\
u(\mathbf{x}) &= g(\mathbf{x}), \quad &\mathbf{x} \in \gamma_h,
\end{aligned}
$$

*with $d(\mathbf{x}) = 0$ for all $\mathbf{x} \in \omega_h^\circ$, it holds that*

$$
\|u\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \|g\|_{l^\infty(\gamma_h)}.
$$

**Proof:** Consider the problem

$$
\begin{aligned}
L_h \overline{u}(\mathbf{x}) &= 0, \quad &\mathbf{x} \in \omega_h, \\
\overline{u}(\mathbf{x}) &= \overline{g}(\mathbf{x}) = const = \|g\|_{l^\infty(\gamma_h)}, \quad &\mathbf{x} \in \gamma_h.
\end{aligned}
$$

It is $\overline{u}(\mathbf{x}) = \|g\|_{l^\infty(\gamma_h)} = const$, since for inner nodes that are not close to the boundary it holds that

$$
L_h \overline{u}(\mathbf{x}) = \underbrace{d(\mathbf{x})}_{=0} \overline{u}(\mathbf{x}) - \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y}) \underbrace{\left(\overline{u}(\mathbf{y}) - \overline{u}(\mathbf{x})\right)}_{=0} = 0.
$$

By definition of the problem, $L_h$ vanishes for constant functions. With the same arguments as in Example 2.18, one can derive the representation (2.8) for inner nodes that are close to the boundary. Inserting (2.8) into (2.9) and using in addition $\overline{u}(\mathbf{x}) = \overline{u}(\mathbf{y})$ yields

$$
L_h \overline{u}(\mathbf{x}) = d(\overline{\mathbf{x}})u(\overline{\mathbf{x}}) = \sum_{\mathbf{y} \in S(\mathbf{x}), \mathbf{y} \in \gamma_h} b(\mathbf{x}, \mathbf{y})\overline{u}(\mathbf{x}) = \sum_{\mathbf{y} \in S(\mathbf{x}), \mathbf{y} \in \gamma_h} b(\mathbf{x}, \mathbf{y})\overline{u}(\mathbf{y}).
$$

This expression is exactly the contribution of the nodes on $\gamma_h$ that are included in $F(\mathbf{x})$ in the scheme (2.6), see also Example 2.16. That means, the finite difference equation is also satisfied by the nodes that are close to the wall.

Now, the statement of the corollary follows by the application of Corollary 2.22, since $\overline{u}(\mathbf{x}) \geq |u(\mathbf{x})|$. ∎

**Corollary 2.25 Problem with homogeneous boundary condition.** *For the solution of the problem*

$$
\begin{aligned}
L_h u(\mathbf{x}) &= f(\mathbf{x}), \quad &\mathbf{x} \in \omega_h, \\
u(\mathbf{x}) &= 0, \quad &\mathbf{x} \in \gamma_h,
\end{aligned}
$$

*with $d(\mathbf{x}) > 0$ for all $\mathbf{x} \in \omega_h$, it is*

$$
\|u\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \left\|D^{-1}f\right\|_{l^\infty(\omega_h)}
$$

*with $D = diag(d(\mathbf{x}))$ for $\mathbf{x} \in \omega_h$.*

**Proof:** Consider the grid function

$$\overline{f}(\mathbf{x}) \;=\; |f(\mathbf{x})| \geq f(\mathbf{x}) \quad \forall \; \mathbf{x} \in \omega_h.$$

From the discrete maximum principle it follows that the solution of the problem

$$\begin{aligned}
L_h \overline{u}(\mathbf{x}) &= \overline{f}(\mathbf{x}), & \mathbf{x} &\in \omega_h, \\
\overline{u}(\mathbf{x}) &= 0, & \mathbf{x} &\in \gamma_h,
\end{aligned}$$

is non-negative, i.e., it holds $\overline{u}(\mathbf{x}) \geq 0$ for $\mathbf{x} \in \omega_h \cup \gamma_h$. Define the node $\overline{\mathbf{x}}$ by the condition

$$\overline{u}(\overline{\mathbf{x}}) = \|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)}\,.$$

In $\overline{\mathbf{x}}$, it is

$$L_h \overline{u}(\overline{\mathbf{x}}) = d(\overline{\mathbf{x}})\overline{u}(\overline{\mathbf{x}}) - \sum_{\mathbf{y}\in S(\overline{\mathbf{x}})} \underbrace{b(\overline{\mathbf{x}},\mathbf{y})}_{>0} \underbrace{\big(\overline{u}(\mathbf{y}) - \overline{u}(\overline{\mathbf{x}})\big)}_{\leq 0} = |f(\overline{\mathbf{x}})|\,,$$

from what follows that

$$\overline{u}(\overline{\mathbf{x}}) \leq \frac{|f(\overline{\mathbf{x}})|}{d(\overline{\mathbf{x}})} \leq \max_{\mathbf{x}\in\omega_h} \frac{|f(\mathbf{x})|}{d(\mathbf{x})} = \max_{\mathbf{x}\in\omega_h} \left| \frac{f(\mathbf{x})}{d(\mathbf{x})} \right| = \left\| D^{-1} f \right\|_{l^\infty(\omega_h)}\,.$$

Since $u(\mathbf{x}) \leq \overline{u}(\overline{\mathbf{x}})$ for all $\mathbf{x} \in \omega_h \cup \gamma_h$ because of Corollary 2.22, the statement of the corollary is proved. ∎

**Corollary 2.26 Another problem with homogeneous boundary condition.**
*Consider*

$$\begin{aligned}
L_h u(\mathbf{x}) &= f(\mathbf{x}), & \mathbf{x} &\in \omega_h, \\
u(\mathbf{x}) &= 0, & \mathbf{x} &\in \gamma_h,
\end{aligned}$$

*with $f(\mathbf{x}) = 0$ for all $\mathbf{x} \in \omega_h^\circ$. With respect to the finite difference scheme it will be assumed that $d(\mathbf{x}) = 0$ for all $\mathbf{x} \in \omega_h^\circ$, and $d(\mathbf{x}) > 0$ for all $\mathbf{x} \in \omega_h^*$. Then the following estimate is valid*

$$\|u\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \left\| D^+ f \right\|_{l^\infty(\omega_h)}$$

*with $D^+ = diag(0, d(\mathbf{x})^{-1})$. The zero entries appear for $\mathbf{x} \in \omega_h^\circ$ and the entries $d(\mathbf{x})^{-1}$ for $\mathbf{x} \in \omega_h^*$.*

**Proof:** Let $\overline{f}(\mathbf{x}) = |f(\mathbf{x})|$, $\mathbf{x} \in \omega_h$, and $\overline{g}(\mathbf{x}) = 0, \mathbf{x} \in \gamma_h$. The solution $\overline{u}(\mathbf{x})$ is non-negative, $\overline{u}(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \omega_h \cup \gamma_h$, see the proof of Corollary 2.25. Define $\overline{\mathbf{x}}$ by

$$\overline{u}(\overline{\mathbf{x}}) = \|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)}\,.$$

One can choose $\overline{\mathbf{x}} \in \omega_h^*$, because if $\overline{\mathbf{x}} \in \omega_h^\circ$, then it holds that

$$\underbrace{d(\overline{\mathbf{x}})}_{=0}\overline{u}(\overline{\mathbf{x}}) - \sum_{\mathbf{y}\in S(\overline{\mathbf{x}})} \underbrace{b(\overline{\mathbf{x}},\mathbf{y})}_{>0} \underbrace{\big(\overline{u}(\mathbf{y}) - \overline{u}(\overline{\mathbf{x}})\big)}_{\leq 0} = \overline{f}(\overline{\mathbf{x}}) = 0,$$

i.e. $\overline{u}(\overline{\mathbf{x}}) = \overline{u}(\mathbf{y})$ for all $\mathbf{y} \in S(\mathbf{x})$. Let $\hat{\mathbf{x}} \in \omega_h^*$ and $\overline{\mathbf{x}}, \mathbf{x}_1, \ldots, \mathbf{x}_m, \hat{\mathbf{x}}$ be a connection with $\mathbf{x}_i \notin \omega_h^*$, $i = 1, \ldots, m$. For $\mathbf{x}_m$ it holds analogously that

$$\overline{u}(\mathbf{x}_m) = \|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)} = \overline{u}(\mathbf{y}) \; \forall \; \mathbf{y} \in S(\mathbf{x}_m).$$

Hence, it follow in particular that $\overline{u}(\hat{\mathbf{x}}) = \|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)}$ such that one can choose $\overline{\mathbf{x}} = \hat{\mathbf{x}}$. It follows that

$$\underbrace{d(\hat{\mathbf{x}})}_{>0} \underbrace{\overline{u}(\hat{\mathbf{x}})}_{=\|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)}} - \sum_{\mathbf{y}\in S(\hat{\mathbf{x}})} \underbrace{b(\hat{\mathbf{x}},\mathbf{y})}_{>0} \underbrace{\big(\overline{u}(\mathbf{y}) - \overline{u}(\hat{\mathbf{x}})\big)}_{\leq 0} = \overline{f}(\hat{\mathbf{x}}).$$

Since all terms in the sum over $\mathbf{x} \in \omega_h$ are non-negative, it follows, using also Corollary 2.22, that

$$\|u\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \|\overline{u}\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \frac{f(\hat{\mathbf{x}})}{d(\hat{\mathbf{x}})} \leq \frac{\overline{f}(\hat{\mathbf{x}})}{d(\hat{\mathbf{x}})} \leq \left\| D^+ f \right\|_{l^\infty(\omega_h)}\,.$$

∎

## 2.4 Stability and Convergence of the Finite Difference Approximation of the Poisson Problem with Dirichlet Boundary Conditions

**Remark 2.27** *Decomposition of the solution.* A short form to write (2.5) is

$$L_h u(\mathbf{x}) = f(\mathbf{x}), \ \mathbf{x} \in \omega_h, \quad u(\mathbf{x}) = g(\mathbf{x}), \ \mathbf{x} \in \gamma_h.$$

The solution of (2.5) can be decomposed into

$$u(\mathbf{x}) = u_1(\mathbf{x}) + u_2(\mathbf{x}),$$

with

$$
\begin{aligned}
L_h u_1(\mathbf{x}) &= f(\mathbf{x}), \ \mathbf{x} \in \omega_h, \quad u_1(\mathbf{x}) = 0, \ \mathbf{x} \in \gamma_h \ \text{(homogeneous boundary cond.)}, \\
L_h u_2(\mathbf{x}) &= 0, \ \mathbf{x} \in \omega_h, \quad u_2(\mathbf{x}) = g(\mathbf{x}), \ \mathbf{x} \in \gamma_h \ \text{(homogeneous right hand side)}.
\end{aligned}
$$

□

### Stability with Respect to the Boundary Condition

**Remark 2.28** *Stability with respect to the boundary condition.* From Corollary 2.24 it follows that

$$\|u_2\|_{l^\infty(\omega_h)} \le \|g\|_{l^\infty(\gamma_h)}. \tag{2.10}$$

□

### Stability with Respect to the Right Hand Side

**Remark 2.29** *Decomposition of the right hand side.* The right hand side will be decomposed into

$$f(\mathbf{x}) = f^\circ(\mathbf{x}) + f^*(\mathbf{x})$$

with

$$f^\circ(\mathbf{x}) = \begin{cases} f(\mathbf{x}), & \mathbf{x} \in \omega_h^\circ \\ 0, & \mathbf{x} \in \omega_h^* \end{cases}, \quad f^*(\mathbf{x}) = f(\mathbf{x}) - f^\circ(\mathbf{x}).$$

Since the considered finite difference scheme is linear, also the function $u_1(\mathbf{x})$ can be decomposed into

$$u_1(\mathbf{x}) = u_1^\circ(\mathbf{x}) + u_1^*(\mathbf{x})$$

with

$$
\begin{aligned}
L_h u_1^\circ(\mathbf{x}) &= f^\circ(\mathbf{x}), \ \mathbf{x} \in \omega_h, \quad u_1^\circ(\mathbf{x}) = 0, \ \mathbf{x} \in \gamma_h, \\
L_h u_1^*(\mathbf{x}) &= f^*(\mathbf{x}), \ \mathbf{x} \in \omega_h, \quad u_1^*(\mathbf{x}) = 0, \ \mathbf{x} \in \gamma_h.
\end{aligned}
$$

□

**Remark 2.30** *Estimate for the inner nodes.* Let $B((0,0), R)$ be a circle with center $(0,0)$ and radius $R$, which is chosen such that $R \ge \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \Omega$. Consider the function

$$\overline{u}(\mathbf{x}) = \alpha \left( R^2 - x^2 - y^2 \right) \quad \text{with } \alpha > 0,$$

which takes only non-negative values for $(x,y) \in \Omega$. Applying the definition of the five point stencil, it follows that

$$
\begin{aligned}
\Lambda \overline{u}(\mathbf{x}) &= -\alpha \Lambda (x^2 + y^2 - R^2) \\
&= -\alpha \left( \frac{(x+h_x)^2 - 2x^2 + (x-h_x)^2}{h_x^2} + \frac{(y+h_y)^2 - 2y^2 + (y-h_y)^2}{h_y^2} \right) \\
&= -4\alpha =: -\overline{f}(\mathbf{x}), \ \mathbf{x} \in \omega_h^\circ,
\end{aligned}
$$

and

$$\Lambda^* \overline{u}(\mathbf{x}) = -\alpha \left[ \frac{1}{\overline{h}_x} \left( \frac{(x + h_x^+)^2 - x^2}{h_x^+} - \frac{x^2 - (x - h_x^-)^2}{h_x^-} \right) \right.$$

$$+ \frac{1}{\overline{h}_y} \left( \frac{(y + h_y^+)^2 - y^2}{h_y^+} - \frac{y^2 - (y - h_y^-)^2}{h_y^-} \right) \right]$$

$$= -\alpha \left( \frac{h_x^+ + h_x^-}{\overline{h}_x} + \frac{h_y^+ + h_y^-}{\overline{h}_y} \right) =: -\overline{f}(\mathbf{x}), \ \mathbf{x} \in \omega_h^*.$$

Hence, $\overline{u}(\mathbf{x})$ is the solution of the problem

$$\begin{aligned} L_h \overline{u}(\mathbf{x}) &= \overline{f}(\mathbf{x}), & \mathbf{x} \in \omega_h, \\ \overline{u}(\mathbf{x}) &= \alpha \left( R^2 - x^2 - y^2 \right) \geq 0, & \mathbf{x} \in \gamma_h. \end{aligned}$$

It is $\overline{u}(\mathbf{x}) \geq 0$ for all $\mathbf{x} \in \gamma_h$. Choosing $\alpha = \frac{1}{4} \|f^\circ\|_{l^\infty(\omega_h)}$, one obtains

$$\begin{aligned} \overline{f}(\mathbf{x}) &= 4\alpha = \|f^\circ\|_{l^\infty(\omega_h)} \geq |f^\circ(\mathbf{x})|, \ \mathbf{x} \in \omega_h^\circ, \\ \overline{f}(\mathbf{x}) &\geq 0 = |f^\circ(\mathbf{x})| \ \mathbf{x} \in \omega_h^*. \end{aligned}$$

Now, Lemma 2.22 (Comparison Lemma) can be applied, which leads to

$$\|u_1^\circ\|_{l^\infty(\omega_h)} \leq \|\overline{u}\|_{l^\infty(\omega_h)} \leq \alpha R^2 = \frac{R^2}{4} \|f^\circ\|_{l^\infty(\omega_h)}. \tag{2.11}$$

One gets the final lower or equal estimate because $(0,0)$ does not need to belong to $\Omega$ or $\omega_h$. $\qquad \square$

**Remark 2.31** *Estimate for the nodes that are close to the boundary.* Corollary 2.26 can be applied to estimate $u_1^*(\mathbf{x})$. For $\mathbf{x} \in \omega_h^\circ$ it is $d(\mathbf{x}) = 0$, see Example 2.18. For $\mathbf{x} \in \omega_h^*$ one has

$$d(\mathbf{x}) = \sum_{\mathbf{y} \in S(\mathbf{x}), \mathbf{y} \in \gamma_h} b(\mathbf{x}, \mathbf{y}) \geq \frac{1}{h^2}$$

with $h = \max\{h_x, h_y\}$, since all terms are of the form

$$\frac{1}{\overline{h}_x h_x^+}, \ \frac{1}{\overline{h}_x h_x^-}, \ \frac{1}{\overline{h}_y h_y^+}, \ \frac{1}{\overline{h}_y h_y^-},$$

see Example 2.18. One obtains

$$\|u_1^*\|_{l^\infty(\omega_h)} \leq \|D^+ f^*\|_{l^\infty(\omega_h)} \leq h^2 \|f^*\|_{l^\infty(\omega_h)}. \tag{2.12}$$

$\square$

**Lemma 2.32 Stability estimate** *The solution of the discrete Dirichlet problem (2.5) satisfies*

$$\|u\|_{l^\infty(\omega_h \cup \gamma_h)} \leq \|g\|_{l^\infty(\gamma_h)} + \frac{R^2}{4} \|\phi\|_{l^\infty(\omega_h^\circ)} + h^2 \|\phi\|_{l^\infty(\omega_h^*)} \tag{2.13}$$

*with $R \geq \|\mathbf{x}\|_2$ for all $\mathbf{x} \in \Omega$ and $h = \max\{h_x, h_y\}$, i.e., the solution $u(\mathbf{x})$ can be bounded in the norm $\|\cdot\|_{l^\infty(\omega_h \cup \gamma_h)}$ by the data of the problem.*

**Proof:** The statement of the lemma is obtained by combining the estimates (2.10), (2.11), and (2.12). $\qquad \blacksquare$

**Convergence**

**Theorem 2.33 Convergence.** *Let $u(\mathbf{x})$ be the solution of the Poisson equation (2.1) and $u_h(\mathbf{x})$ be the finite difference approximation given by the solution of (2.5). Then, it is*
$$\|u - u_h\|_{l^\infty(\omega_h \cup \gamma_h)} \le Ch^2$$
*with $h = \max\{h_x, h_y\}$.*

**Proof:** The error in the node $(x_i, y_j)$ is defined by $e_{ij} = u(x_i, y_j) - u_h(x_i, y_j)$. With
$$-\Lambda u(x_i, y_j) = -\Delta u(x_i, y_j) + \mathcal{O}\left(h^2\right) = f(x_i, y_j) + \mathcal{O}\left(h^2\right),$$
one obtains by subtracting the finite difference equation, the following problem for the error
$$\begin{aligned}
-\Lambda e(\mathbf{x}) &= \psi(\mathbf{x}), & \mathbf{x} \in w_h^\circ, \ \psi(\mathbf{x}) = \mathcal{O}\left(h^2\right), \\
-\Lambda^* e(\mathbf{x}) &= \psi(\mathbf{x}), & \mathbf{x} \in w_h^*, \ \psi(\mathbf{x}) = \mathcal{O}(1), \\
e(\mathbf{x}) &= 0, & \mathbf{x} \in \gamma_h,
\end{aligned}$$
where $\psi(\mathbf{x})$ is the consistency error, see Section 2.2. Applying the stability estimate (2.13) to this problem, one obtains immediately
$$\|e\|_{l^\infty(\omega_h \cup \gamma_h)} \le \frac{R^2}{4}\|\psi\|_{l^\infty(\omega_h^\circ)} + h^2\|\psi\|_{l^\infty(\omega_h^*)} = \mathcal{O}\left(h^2\right).$$
∎

## 2.5 An Efficient Solver for the Dirichlet Problem in the Rectangle

**Remark 2.34** *Contents of this section.* This section considers the Poisson equation (2.1) in the special case $\Omega = (0, l_x) \times (0, l_y)$. In this case, a modification of the difference stencil in a neighborhood of the boundary of the domain is not needed. The convergence of the finite difference approximation was already established in Theorem 2.33. Applying this approximation results in a large linear system of equations $A\mathbf{u} = \mathbf{f}$ which has to be solved. This section presents an approach for solving this system in the case of a rectangular domain in an almost optimal way. □

**Remark 2.35** *The considered problem and its approximation.* The considered continuous problem consists in solving
$$\begin{aligned}
-\Delta u &= f & \text{in } \Omega = (0, l_x) \times (0, l_y), \\
u &= g & \text{on } \partial\Omega,
\end{aligned}$$
and the corresponding discrete problem in solving
$$\begin{aligned}
-\Lambda u(\mathbf{x}) &= \phi(\mathbf{x}), & \mathbf{x} \in \omega_h, \\
u(\mathbf{x}) &= g(\mathbf{x}), & \mathbf{x} \in \gamma_h,
\end{aligned}$$
where the discrete Laplacian is of the form (for simplicity of notation, the subscript $h$ is omitted)
$$\Lambda u = \frac{u_{i+1,j} - 2u_{ij} + u_{i-1,j}}{h_x^2} + \frac{u_{i,j+1} - 2u_{ij} + u_{i,j-1}}{h_y^2} =: \Lambda_x u + \Lambda_y u, \quad (2.14)$$
with $h_x = l_x/n_x, h_y = l_y/n_y, i = 0, \ldots, n_x, j = 0, \ldots, n_y$, see Figure 2.6. □
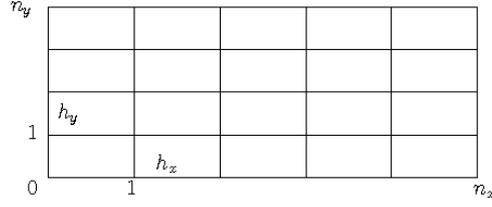
Figure 2.6: Grid for the Dirichlet problem in the rectangular domain.

**Remark 2.36** *The linear system of equations.* The difference scheme (2.14) is equivalent to a linear system of equations $A\mathbf{u} = \mathbf{f}$.

For assembling the matrix and the right hand side of the system, usually a lexicographical enumeration of the nodes of the grid is used. The nodes are called enumerated lexicographically if the node $(i_1, j_1)$ has a smaller number than the node $(i_2, j_2)$, if for the corresponding coordinates it is

$$y_1 < y_2 \ \ \text{or} \ \ (y_1 = y_2) \wedge (x_1 < x_2).$$

Using this lexicographical enumeration of the nodes, one obtains for the inner nodes a system of the form

$$
\begin{aligned}
A &= \text{BlockTriDiag}(C, B, C) \in \mathbb{R}^{(n_x-1)(n_y-1) \times (n_x-1)(n_y-1)}, \\
B &= \text{TriDiag}\left(-\frac{1}{h_x^2}, \frac{2}{h_x^2} + \frac{2}{h_y^2}, -\frac{1}{h_x^2}\right) \in \mathbb{R}^{(n_x-1) \times (n_x-1)}, \\
C &= \text{Diag}\left(-\frac{1}{h_y^2}\right) \in \mathbb{R}^{(n_x-1) \times (n_x-1)}, \\
\mathbf{f} &= \begin{cases}
\phi(\mathbf{x}), & \mathbf{x} \in \omega_h^\circ, \\
\phi(\mathbf{x}) + \dfrac{g(x \pm h_x, y)}{h_x^2}, & i \in \{1, n_x - 1\}; j \notin \{1, n_y - 1\}, \\
\phi(\mathbf{x}) + \dfrac{g(x, y \pm h_y)}{h_y^2}, & i \notin \{1, n_x - 1\}; j \in \{1, n_y - 1\}, \\
\phi(\mathbf{x}) + \dfrac{g(x \pm h_x, y)}{h_x^2} + \dfrac{g(x, yx \pm h_y)}{h_y^2}, & i \in \{1, n_x - 1\}; j \in \{1, n_y - 1\}.
\end{cases}
\end{aligned}
$$

The last line of the right hand side vector is for inner nodes which are situated in corner points of $\omega_h^\circ$. In this approach, the known Dirichlet boundary values are already substituted into the system and they appear in the right hand side vector. The matrices $B$ and $C$ possess some modifications for nodes which have a neighbor on the boundary.

The linear system of equations has the following properties:

- high dimension: $N = (n_x - 1)(n_y - 1) \sim 10^3 \cdots 10^7$,
- sparse: per row and column of the matrix there are only 3, 4, or 5 non-zero entries,
- symmetric: hence, all eigenvalues are real,
- positive definite: all eigenvalues are positive. It holds that

$$
\begin{aligned}
\lambda_{\min} &= \lambda_{(1,1)} \sim \pi^2 \left(\frac{1}{l_x^2} + \frac{1}{l_y^2}\right) = \mathcal{O}(1), \\
\lambda_{\max} &= \lambda_{(n_x-1, n_y-1)} \sim \pi^2 \left(\frac{1}{h_x^2} + \frac{1}{h_y^2}\right) = \mathcal{O}\left(h^{-2}\right) \qquad (2.15)
\end{aligned}
$$

with $h = \max\{h_x, h_y\}$, see Remark 2.37 below.

- high condition number: For the spectral condition number of a symmetric and positive definite matrix it is

$$\kappa_2(A) = \frac{\lambda_{\max}}{\lambda_{\min}} = \mathcal{O}\left(h^{-2}\right).$$

Since the dimension of the matrix is large, iterative solvers are an appropriate approach for solving the linear system of equations. The main costs for iterative solvers are the matrix-vector multiplications (often one per iteration). The cost of one matrix-vector multiplication is for sparse matrices proportional to the number of unknowns. Hence, an optimal solver is given if the number of operations for solving the linear system of equations is proportional to the number of unknowns. It is known that the number of iterations of many iterative solvers depends on the condition number of the matrix:

- (damped) Jacobi method, SOR, SSOR. The number of iteration is proportional to $\kappa_2(A)$. That means, if the grid is refined once, $h \to h/2$, then the number of unknowns is increased by around the factor 4 in two dimensions and also the number of iterations increases by a factor of around 4. Altogether, for one refinement step, the total costs increase by a factor of around 16.
- (preconditioned) conjugate gradient (PCG) method. The number of iterations is proportional to $\sqrt{\kappa_2(A)}$, see the corresponding theorem from the class Numerical Mathematics II. Hence, the total costs increase by a factor of around 8 if the grid is refined once.
- multigrid methods. For multigrid methods, the number of iterations is constant. Hence, the total costs are proportional to the number of unknowns and these methods are optimal. However, the implementation of multigrid methods is involved.

$\square$

**Remark 2.37** *An eigenvalue problem.* The derivation of an alternative direct solver is based on the eigenvalues and eigenvectors of the discrete Laplacian. It is possible to computed these quantities only in special situations, e.g., if the Poisson problem with Dirichlet boundary conditions is considered, the domain is rectangular, and the Laplacian is approximated with the five point stencil.

Consider the following eigenvalue problem

$$\begin{aligned} -\Lambda v(\mathbf{x}) &= \lambda v(\mathbf{x}), \quad \mathbf{x} \in \omega_h, \\ v(\mathbf{x}) &= 0, \qquad \mathbf{x} \in \gamma_h. \end{aligned}$$

The solution of this problem is sought in product form (separation of variables)

$$v_{ij}^{(\mathbf{k})} = v_i^{(k_x),x} v_j^{(k_y),y}, \quad \mathbf{k} = (k_x, k_y)^T.$$

It is

$$\Lambda v_{ij}^{(\mathbf{k})} = \Lambda_x v_i^{(k_x),x} v_j^{(k_y),y} + v_i^{(k_x),x} \Lambda_y v_j^{(k_y),y} = -\lambda_{\mathbf{k}} v_i^{(k_x),x} v_j^{(k_y),y}$$

with $i = 0, \ldots, n_x$, $j = 0, \ldots, n_y$ refers to the nodes and $k_x = 1, \ldots, n_x - 1$, $k_y = 1, \ldots, n_y - 1$ refers to the eigenvalues. Note that the number of eigenvalues is equal to the number of inner nodes, i.e. it is $(n_x - 1)(n_y - 1)$. In this ansatz, also a splitting of the eigenvalues in a contribution from the $x$ coordinate and a contribution from the $y$ coordinate is included. From the boundary condition it follows that

$$v_0^{(k_x),x} = v_{n_x}^{(k_x),x} = v_0^{(k_y),y} = v_{n_y}^{(k_y),y} = 0.$$

Now, the eigenvalue problem can be split

$$\frac{\Lambda_x v_i^{(k_x),x}}{v_i^{(k_x),x}} + \lambda_{k_x}^{(x)} = -\frac{\Lambda_y v_j^{(k_y),y}}{v_j^{(k_y),y}} - \lambda_{k_y}^{(y)}$$

with $\lambda_{\mathbf{k}} = \lambda_{k_x}^{(x)} + \lambda_{k_y}^{(y)}$. Both sides of this equation have to be constant since one of them depends only on $i$, i.e., on $x$, and the other only on $j$, i.e., on $y$. The splitting of $\lambda_{\mathbf{k}}$ can be chosen such that the constant is zero. Then, one gets

$$\Lambda_x v_i^{(k_x),x} + \lambda_{k_x}^{(x)} v_i^{(k_x),x} = 0, \quad \Lambda_y v_j^{(k_y),y} + \lambda_{k_y}^{(y)} v_j^{(k_y),y} = 0.$$

The solution of these eigenvalue problems is well known (exercise)

$$
\begin{aligned}
v_i^{(k_x),x} &= \sqrt{\frac{2}{l_x}} \sin\left(\frac{k_x \pi i}{n_x}\right), \quad \lambda_{k_x}^{(x)} = \frac{4}{h_x^2} \sin^2\left(\frac{k_x \pi}{2 n_x}\right), \\
v_j^{(k_y),y} &= \sqrt{\frac{2}{l_y}} \sin\left(\frac{k_y \pi j}{n_y}\right), \quad \lambda_{k_y}^{(y)} = \frac{4}{h_y^2} \sin^2\left(\frac{k_y \pi}{2 n_y}\right).
\end{aligned}
$$

It follows that the solution of the full eigenvalue problem is

$$v_{ij}^{(\mathbf{k})} = \frac{2}{\sqrt{l_x l_y}} \sin\left(\frac{k_x \pi i}{n_x}\right) \sin\left(\frac{k_y \pi j}{n_y}\right), \quad \lambda_{\mathbf{k}} = \frac{4}{h_x^2} \sin^2\left(\frac{k_x \pi}{2 n_x}\right) + \frac{4}{h_y^2} \sin^2\left(\frac{k_y \pi}{2 n_y}\right)$$

with $i = 0, \ldots, n_x, j = 0, \ldots, n_y$ and $k_x = 1, \ldots, n_x - 1, k_y = 1, \ldots, n_y - 1$. Using a Taylor series expansion, one obtains now the asymptotic behavior of the eigenvalues as given in (2.15). Note that because of the splitting of the eigenvalues into the directional contributions, the number of individual terms for computing the eigenvalues is only $\mathcal{O}(n_x + n_y)$.

Since the matrix corresponding to $\Lambda$ is symmetric, the eigenvectors are orthogonal with respect to the Euclidean vector product. They become orthonormal with respect to the weighted Euclidean vector product

$$\langle u, v \rangle = h_x h_y \sum_{\mathbf{x} \in \omega_h \cup \gamma_h} u(\mathbf{x}) v(\mathbf{x}) = h_x h_y \sum_{i=0}^{n_x} \sum_{j=0}^{n_y} u_{ij}(\mathbf{x}) v_{ij}(\mathbf{x}), \quad h_x = \frac{l_x}{n_x}, h_y = \frac{l_y}{n_y},$$
(2.16)

i.e., then it is

$$\langle v^{(\mathbf{k})}, v^{(\mathbf{m})} \rangle = \delta_{\mathbf{k},\mathbf{m}}.$$

This property can be checked by using the relation

$$\sum_{i=0}^{n} \sin^2\left(\frac{i\pi}{n}\right) = \frac{n}{2}, \quad n > 1.$$

The norm induced by the weighted Euclidean vector product is given by

$$\|v\|_h = \langle v, v \rangle^{1/2} = \left( h_x h_y \sum_{i=0}^{n_x} \sum_{j=0}^{n_y} v_{ij}^2(\mathbf{x}) \right)^{1/2}.$$
(2.17)

The weights are such that this norm is for constants (almost) independent of the mesh, i.e.,

$$\|1\|_h = (h_x h_y (n_x + 1)(n_y + 1))^{1/2} = \left( l_x l_y \frac{n_x + 1}{n_x} \frac{n_y + 1}{n_y} \right)^{1/2} \approx (l_x l_y)^{1/2}.$$

$\square$

**Remark 2.38** *Solver based on the eigenvalues and eigenvectors.* One uses the ansatz

$$f(\mathbf{x}) = \sum_{\mathbf{k}} f_{\mathbf{k}} v^{(\mathbf{k})}(\mathbf{x})$$

with the Fourier coefficients

$$f_{\mathbf{k}} = \langle f, v^{(\mathbf{k})} \rangle = \frac{2h_x h_y}{\sqrt{l_x l_y}} \sum_{i=0}^{n_x} \sum_{j=0}^{n_y} f_{ij} \sin\left(\frac{k_x \pi i}{n_x}\right) \sin\left(\frac{k_y \pi j}{n_y}\right), \quad \mathbf{k} = (k_x, k_y).$$

The solution $u(\mathbf{x})$ of (2.14) is sought in the form

$$u(\mathbf{x}) = \sum_{\mathbf{k}} u_{\mathbf{k}} v^{(\mathbf{k})}(\mathbf{x})$$

with unknown coefficients $u_{\mathbf{k}}$. With this ansatz, one obtains

$$\Lambda u = \sum_{\mathbf{k}} u_{\mathbf{k}} \Lambda v^{(\mathbf{k})} = \sum_{\mathbf{k}} u_{\mathbf{k}} \lambda_{\mathbf{k}} v^{(\mathbf{k})}.$$

Since the eigenfunctions form a basis of the space of the grid functions, a comparison of the coefficients with the right hand side gives

$$u_{\mathbf{k}} = \frac{f_{\mathbf{k}}}{\lambda_{\mathbf{k}}}$$

or, for each component,

$$u_{ij} = \sum_{\mathbf{k}} \frac{f_{\mathbf{k}}}{\lambda_{\mathbf{k}}} v_{ij}^{(\mathbf{k})} = \frac{2h_x h_y}{\sqrt{l_x l_y}} \sum_{k_x=1}^{n_x-1} \sum_{k_y=1}^{n_y-1} \frac{f_{\mathbf{k}}}{\lambda_{\mathbf{k}}} \sin\left(\frac{k_x \pi i}{n_x}\right) \sin\left(\frac{k_y \pi j}{n_y}\right),$$

$i = 0, \ldots, n_x$, $j = 0, \ldots, n_y$.

It is possible to implement this approach with the Fast Fourier Transform (FFT) with

$$\mathcal{O}\left(n_x n_y \log_2 n_x + n_x n_y \log_2 n_y\right) = \mathcal{O}\left(N \log_2 N\right), \quad N = (n_x - 1)(n_y - 1),$$

operations. Hence, this method is almost optimal. $\qquad \square$

## 2.6 A Higher Order Discretizations

**Remark 2.39** *Contents.* The five point stencil is a second order discretization of the Laplacian. In this section, a discretization of higher order will be studied. In these studies, only the case of a rectangular domain $\Omega = (0, l_x) \times (0, l_y)$ and Dirichlet boundary conditions will be considered. $\qquad \square$

**Remark 2.40** *Derivation of a fourth order approximation.* Let $u(\mathbf{x})$ be the solution of the Poisson equation (2.1) and assume that $u(\mathbf{x})$ is sufficiently smooth. It is

$$Lu(\mathbf{x}) = \Delta u(\mathbf{x}) = L_x u(\mathbf{x}) + L_y u(\mathbf{x}), \quad L_\alpha u := \frac{\partial^2 u}{\partial x_\alpha^2}.$$

Let the five point stencil be represented by the following operator

$$\Lambda u = \Lambda_x u + \Lambda_y u.$$

Applying a Taylor series expansion, one finds that

$$\Lambda u - \Delta u = \frac{h_x^2}{12} L_x^2 u + \frac{h_y^2}{12} L_y^2 u + \mathcal{O}\left(h^4\right). \tag{2.18}$$

From the equation $-Lu = f$ it follows that

$$L_x^2 u = -L_x f - L_x L_y u, \quad L_y^2 u = -L_y f - L_y L_x u.$$

Inserting these expressions into (2.18) gives

$$\Lambda u - \Delta u = -\frac{h_x^2}{12} L_x f - \frac{h_y^2}{12} L_y f - \frac{h_x^2 + h_y^2}{12} L_x L_y u + \mathcal{O}\left(h^4\right). \qquad (2.19)$$

The operator $L_x L_y = \frac{\partial^4}{\partial x^2 \partial y^2}$ can be approximated as follows

$$L_x L_y u \approx \Lambda_x \Lambda_y u = u_{\bar{x}x\bar{y}y}.$$

The difference operator in this approximation requires nine points, see Figure 2.7

$$\begin{aligned}
\Lambda_x \Lambda_y u \quad = \quad & \frac{1}{h_x^2 h_y^2} \Big( u_{i+1,j+1} - 2u_{i,j+1} + u_{i-1,j+1} - 2u_{i+1,j} + 4u_{ij} \\
& -2u_{i-1,j} + u_{i+1,j-1} - 2u_{i,j-1} + u_{i-1,j-1} \Big).
\end{aligned}$$
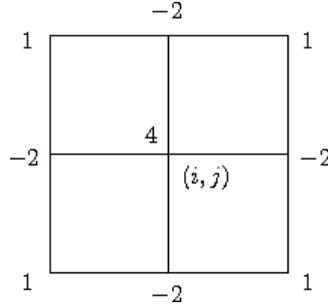
Therefore it is called nine point stencil.



Figure 2.7: The nine point stencil.

One checks, as usual by using a Taylor series expansion, that this approximation is of second order

$$L_x L_y u - \Lambda_x \Lambda_y u = \mathcal{O}\left(h^2\right).$$

Inserting this expansion into (2.19) and using the partial differential equation shows that the difference equation

$$-\left(\Lambda + \frac{h_x^2 + h_y^2}{12} \Lambda_x \Lambda_y\right) u = \left(f + \frac{h_x^2}{12} L_x f + \frac{h_y^2}{12} L_y f\right)$$

is a fourth order approximation of the differential equation (2.1). In addition, one can replace the derivatives of $f(\mathbf{x})$ also by finite differences

$$L_x f = \Lambda_x f + \mathcal{O}\left(h_x^2\right), \quad L_y f = \Lambda_y f + \mathcal{O}\left(h_y^2\right).$$

Finally, one obtains a finite difference equation $-\Lambda' u = \phi$ with

$$\Lambda' = \Lambda_x + \Lambda_y + \frac{h_x^2 + h_y^2}{12} \Lambda_x \Lambda_y, \quad \phi = f + \frac{h_x^2}{12} \Lambda_x f + \frac{h_y^2}{12} \Lambda_y f.$$

$\square$

**Remark 2.41** *On the convergence of the fourth order approximation.* The finite difference problem with the higher order approximation property can be written with the help of the second order differences. Since the convergence proof is based on the five point stencil, the following lemma considers this stencil. It will be proved that one can estimate the values of the grid function by the second order differences. This result will be used in the convergence proof for the fourth order approximation.

$\square$

**Lemma 2.42 Embedding lemma.** *Let*

$$\omega_h = \{(ih_x, jh_y) \ : \ i = 1, \ldots, n_x - 1, j = 1, \ldots, n_y - 1\},$$

*and let $y$ be a grid function on $\omega_h \cup \gamma_h$ with $y(\mathbf{x}) = 0$ for $\mathbf{x} \in \gamma_h$. Then, the following estimate holds*

$$\|y\|_{l^\infty(\omega_h \cup \gamma_h)} \le M \, \|Ay\|_h \, ,$$

*with $M = \frac{\max\{l_x^2, l_y^2\}}{2\sqrt{l_x l_y}}$, $A$ is the matrix obtained by using the five point stencil $\Lambda = \Lambda_x + \Lambda_y$ for approximating the second derivatives, and the norm on the right hand side is defined in (2.17).*

**Proof:** Let $\{v_{ij}^{\mathbf{k}}\}$, $\mathbf{k} = (k_x, k_y)$, be the orthonormal basis with

$$v_{ij}^{\mathbf{k}} = \frac{2}{\sqrt{l_x l_y}} \sin\left(\frac{k_x \pi i}{n_x}\right) \sin\left(\frac{k_y \pi j}{n_y}\right)$$

which was derived in Remark 2.37. Then, there is a unique representation of the grid function $y = \sum_{\mathbf{k}} y_{\mathbf{k}} v^{\mathbf{k}}$ and it holds with (2.16)

$$Ay = \sum_{\mathbf{k}} y_{\mathbf{k}} \lambda_{\mathbf{k}} v^{\mathbf{k}}, \quad \|Ay\|_h^2 = \frac{1}{h_x h_y} \sum_{\mathbf{k}} y_{\mathbf{k}}^2 \lambda_{\mathbf{k}}^2.$$

It follows for $\mathbf{x} \in \omega_h$, because of $|\sin(x)| \le 1$ for all $x \in \mathbb{R}$, that

$$|y(\mathbf{x})| = \left| \sum_{\mathbf{k}} y_{\mathbf{k}} v^{\mathbf{k}}(\mathbf{x}) \right| \le \sum_{\mathbf{k}} |y_{\mathbf{k}}| \left| v^{\mathbf{k}}(\mathbf{x}) \right| \le \sum_{\mathbf{k}} |y_{\mathbf{k}}| \max_{\mathbf{k}} \left| v^{\mathbf{k}}(\mathbf{x}) \right| \le \frac{2}{\sqrt{l_x l_y}} \sum_{\mathbf{k}} |y_{\mathbf{k}}| \, .$$

Applying the Cauchy–Schwarz inequality for sums gives

$$
\begin{aligned}
|y(\mathbf{x})|^2 &\le& \frac{4}{l_x l_y} \left( \sum_{\mathbf{k}} |y_{\mathbf{k}}| \right)^2 \\
&=& \frac{4}{l_x l_y} \left( \sum_{\mathbf{k}} |\lambda_{\mathbf{k}} y_{\mathbf{k}}| \frac{1}{\lambda_{\mathbf{k}}} \right)^2 \\
&\le& \frac{4}{l_x l_y} \sum_{\mathbf{k}} \lambda_{\mathbf{k}}^2 y_{\mathbf{k}}^2 \sum_{\mathbf{k}} \frac{1}{\lambda_{\mathbf{k}}^2} \\
&=& \frac{4}{l_x l_y} \|Ay\|_h^2 \sum_{\mathbf{k}} \frac{1}{\lambda_{\mathbf{k}}^2}.
\end{aligned}
\tag{2.20}
$$

Now, one has to estimate the last sum. It is already known that

$$\lambda_{\mathbf{k}} = \frac{4}{h_x^2} \sin^2\left(\frac{k_x \pi}{2 n_x}\right) + \frac{4}{h_y^2} \sin^2\left(\frac{k_y \pi}{2 n_y}\right), \quad k_x = 1, \ldots, n_x - 1, \ k_y = 1, \ldots, n_y - 1.$$

Setting $l = \max\{l_x, l_y\}$ and $h_\alpha = l_\alpha / n_\alpha, \phi_\alpha = \frac{k_\alpha \pi}{2 n_\alpha} \in (0, \pi/2)$, $\alpha \in \{x, y\}$, leads to

$$\lambda_{\mathbf{k}} = \frac{k_x^2 \pi^2}{l_x^2} \left(\frac{\sin \phi_x}{\phi_x}\right)^2 + \frac{k_y^2 \pi^2}{l_y^2} \left(\frac{\sin \phi_y}{\phi_y}\right)^2 \ge 4 \left(\frac{k_x^2}{l_x^2} + \frac{k_y^2}{l_y^2}\right) \ge \frac{4}{l^2} \left(k_x^2 + k_y^2\right).$$

In performing this estimate, it was used that the function $\sin(\phi)/\phi$ is monotonically decreasing on $(0, \pi/2)$, see Figure 2.8, and that

$$\frac{\sin \phi}{\phi} \ge \frac{\sin(\pi/2)}{\pi/2} = \frac{2}{\pi} \ \forall \ \phi \in (0, \pi/2).$$

Let $G = \{(x, y) \ : \ x > 0, y > 0, x^2 + y^2 > 1\}$ be the first quadrant of the complex plane without the part that belongs to the unit circle, see Figure 2.9. The function $\left(k_x^2 + k_y^2\right)^{-2}$
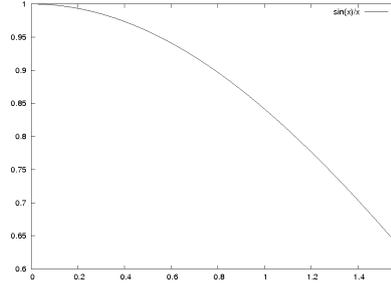
Figure 2.8: The function $\sin(\phi)/\phi$.

has its smallest value in the square $[k_x - 1, k_x] \times [k_y - 1, k_y]$ in the point $(k_x, k_y)$. Using the lower estimate of $\lambda_{\mathbf{k}}$, one obtains

$$
\begin{aligned}
\sum_{\mathbf{k}, \mathbf{k} \neq (1,1)} \frac{1}{\lambda_{\mathbf{k}}^2} \quad &\leq \quad \frac{l^4}{16} \sum_{\mathbf{k}, \mathbf{k} \neq (1,1)} \left(k_x^2 + k_y^2\right)^{-2} \\
&= \quad \frac{l^4}{16} \sum_{\mathbf{k}, \mathbf{k} \neq (1,1)} \underbrace{\left(k_x^2 + k_y^2\right)^{-2}}_{\text{smallest value in square}} \underbrace{\int_{k_x-1}^{k_x} \int_{k_y-1}^{k_y} dx dy}_{=1} \\
&= \quad \frac{l^4}{16} \sum_{\mathbf{k}, \mathbf{k} \neq (1,1)} \int_{k_x-1}^{k_x} \int_{k_y-1}^{k_y} \left(k_x^2 + k_y^2\right)^{-2} dx dy \\
&\leq \quad \frac{l^4}{16} \int_G \left(x^2 + y^2\right)^{-2} dx dy \\
&\overset{\text{polar coord.}}{=} \quad \frac{l^4}{16} \int_1^\infty \int_0^{\pi/2} \frac{\rho}{\rho^4} d\phi d\rho = \frac{l^4}{16} \frac{\pi}{2} \left( -\frac{\rho^2}{2} \Big|_{\rho=1}^{\rho=\infty} \right) = \frac{\pi l^4}{64}.
\end{aligned}
$$

For performing this computation, one has to exclude $\rho \to 0$.



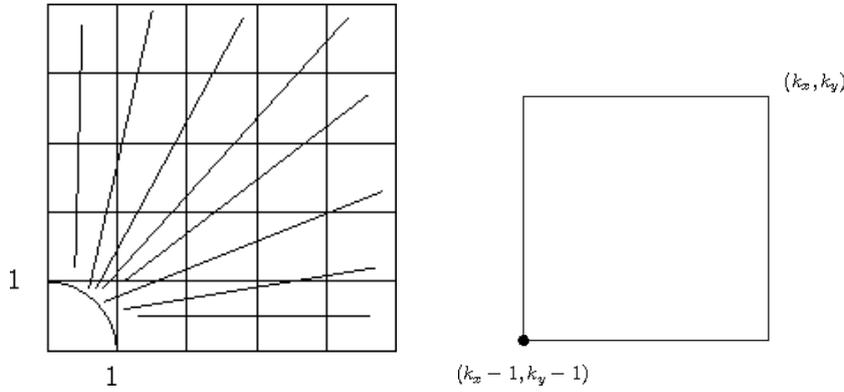Figure 2.9: Illustration to the proof of Lemma 2.42.

For $\lambda_{(1,1)}$ it is

$$
\begin{aligned}
\lambda_{(1,1)} \quad &= \quad \frac{4}{h_x^2} \sin^2\left(\frac{\pi}{2n_x}\right) + \frac{4}{h_y^2} \sin^2\left(\frac{\pi}{2n_y}\right) = \frac{4}{h_x^2} \sin^2\left(\frac{h_x \pi}{2l_x}\right) + \frac{4}{h_y^2} \sin^2\left(\frac{h_y \pi}{2l_y}\right) \\
&= \quad \frac{\pi^2}{l_x^2} \left(\frac{2l_x}{h_x \pi}\right)^2 \sin^2\left(\frac{h_x \pi}{2l_x}\right) + \frac{\pi^2}{l_y^2} \left(\frac{2l_y}{h_y \pi}\right)^2 \sin^2\left(\frac{h_y \pi}{2l_y}\right) \\
&\geq \quad \frac{\pi^2}{l_x^2} \frac{8}{\pi^2} + \frac{\pi^2}{l_y^2} \frac{8}{\pi^2} \geq \frac{16}{l^2}.
\end{aligned} \tag{2.21}
$$

For this estimate, the following relations and the monotonicity of $\sin(x)/x$, see Figure 2.8, were used

$$h_\alpha \le \frac{l_\alpha}{2}, \quad \phi_\alpha = \frac{h_\alpha \pi}{2l_\alpha} \le \frac{\pi}{4}, \quad \left(\frac{\sin \phi_\alpha}{\phi_\alpha}\right)^2 \ge \left(\frac{\sin(\pi/4)}{\pi/4}\right)^2 = \frac{8}{\pi^2}.$$

Collecting all estimates gives

$$\sum_{\mathbf{k}} \frac{1}{\lambda_{\mathbf{k}}^2} = \lambda_{(1,1)}^{-2} + \sum_{\mathbf{k},\mathbf{k}\ne(1,1)} \frac{1}{\lambda_{\mathbf{k}}^2} \le \frac{l^4}{256} + \frac{\pi l^4}{64} \le \frac{l^4}{16}.$$

Inserting this estimate into (2.20), the final estimate has the form

$$\|y\|_{l^\infty (\omega_h \cup \gamma_h)} \le \frac{2}{\sqrt{l_x l_y}} \|Ay\|_h \frac{l^2}{4} =: M \|Ay\|_h.$$

$\blacksquare$

**Theorem 2.43 Convergence of the higher order finite difference scheme.**
*The finite difference scheme*

$$\begin{aligned} -\Lambda' u(\mathbf{x}) &= \phi(\mathbf{x}), \quad \mathbf{x} \in \omega_h^\circ, \\ u(\mathbf{x}) &= g(\mathbf{x}), \quad \mathbf{x} \in \gamma_h, \end{aligned}$$

*with*

$$\Lambda' = \Lambda_x + \Lambda_y + \frac{h_x^2 + h_y^2}{12} \Lambda_x \Lambda_y, \quad \phi = f + \frac{h_x^2}{12} \Lambda_x f + \frac{h_y^2}{12} \Lambda_y f,$$

*converges of fourth order.*

**Proof:** Analogously as in the proof of Theorem 2.33, one finds that the following equation holds for the error $e = u(x_i, y_j) - u_{ij}$:

$$\begin{aligned} -\Lambda' e(\mathbf{x}) &= \psi(\mathbf{x}), \quad \psi = \mathcal{O}\left(h^4\right), \mathbf{x} \in \omega_h, \\ e(\mathbf{x}) &= 0, \quad \mathbf{x} \in \gamma_h. \end{aligned}$$

Let $\Omega_h$ be the vector space of grid functions, which are non-zero only in the interior, i.e., at the nodes from $\omega_h$, and which vanish on $\gamma_h$. Let $A_\alpha y = -\Lambda_\alpha y$, $y \in \Omega_h$, $\alpha \in \{x, y\}$. The operators $A_\alpha : \Omega_h \to \Omega_h$ are linear and they have the following properties:

- They are symmetric and positive definite, i.e., $A_\alpha = A_\alpha^* > 0$, where $A_\alpha^*$ is the adjoint (transposed) of $A_\alpha$, and $(A_\alpha u, v) = (u, A_\alpha v)$, $\forall\, u, v \in \Omega_h$.
- They are elliptic, i.e., $(A_\alpha u, u) \ge \lambda_1^{(\alpha)}(u, u)$, $\forall u \in \Omega_h$, with

$$\lambda_1^{(\alpha)} = \frac{4}{h_\alpha^2} \sin^2\left(\frac{\pi h_\alpha}{2l_\alpha}\right) \ge \frac{8}{l_\alpha^2},$$

  see (2.21).
- They are bounded, i.e., it holds $(A_\alpha u, u) \le \lambda_{n_\alpha - 1}^{(\alpha)}(u, u)$ with

$$\lambda_{n_\alpha - 1}^{(\alpha)} = \frac{4}{h_\alpha^2} \sin^2\left(\frac{k_\alpha \pi}{2n_\alpha}\right) \le \frac{4}{h_\alpha^2}$$

  and $\|A_\alpha\|_2 \le 4/h_\alpha^2$, since the spectral norm of a symmetric positive definite matrix is the largest eigenvalue.
- They are commutative, i.e., $A_x A_y = A_y A_x$.
- It holds $A_x A_y = (A_x A_y)^*$.

The error equation on $\omega_h$ is given by

$$A_x e + A_y e - (\kappa_x + \kappa_y) A_x A_y e = \psi \quad \text{with} \quad \kappa_\alpha = \frac{h_\alpha^2}{12}. \tag{2.22}$$

33

Using the boundedness of the operators, one finds for all $v \in \Omega_h$ that

$$
\begin{aligned}
(\kappa_x A_x A_y v + \kappa_y A_x A_y v, v) &= ((\kappa_x A_x) A_y v, v) + ((\kappa_y A_y) A_x v, v) \\
&\leq \frac{h_x^2}{12} \frac{4}{h_x^2} (A_y, v) + \frac{h_y^2}{12} \frac{4}{h_y^2} (A_x v, v) \\
&= \frac{1}{3} ((A_x + A_y) v, v).
\end{aligned}
$$

Now, it follows for all $v \in \Omega_h$ that

$$
\begin{aligned}
(\underbrace{(A_x + A_y - (\kappa_x + \kappa_y) A_x A_y)}_{=A'} v, v) &= ((A_x + A_y) v, v) - (\kappa_x A_x A_y v + \kappa_y A_x A_y v, v) \\
&\geq \frac{2}{3} ((A_x + A_y) v, v) \geq 0.
\end{aligned}
$$

The matrices on both sides of this inequality are symmetric and because the matrix on the lower estimate is positive definite, also the matrix at the upper estimate is positive definite. Since matrices commute since the order of applying the finite differences in $x$ and $y$ direction does not matter. Using these properties, one gets (*exercise*)

$$
\left\| \frac{2}{3} (A_x + A_y) e \right\|_h \leq \|A' e\|_h = \|\psi\|_h,
$$

where the last equality follows from (2.22). The application of Lemma 2.42 to the error gives

$$
\begin{aligned}
\|e(\mathbf{x})\|_{l^\infty (\omega_h \cup \gamma_h)} &\leq \frac{l^2}{2\sqrt{l_x l_y}} \|(\Lambda_x + \Lambda_y) e\|_h \leq \frac{3l^2}{4\sqrt{l_x l_y}} \|A' e\|_h = \frac{3l^2}{4\sqrt{l_x l_y}} \|\psi\|_h \\
&\leq \frac{3l^2}{4\sqrt{l_x l_y}} (h_x h_y (n_x + 1)(n_y + 1))^{1/2} \|\psi\|_{l^\infty (\omega_h \cup \gamma_h)} = \mathcal{O}\left(h^4\right).
\end{aligned}
$$

∎

**Remark 2.44** *On the discrete maximum principle.* Reformulation of the finite difference scheme $-\Lambda' u = \phi$ in the form studied for the discrete maximum principle gives for $u_{ij}$

$$
\begin{aligned}
a(\mathbf{x}) u(\mathbf{x}) &= \sum_{\mathbf{y} \in S(\mathbf{x})} b(\mathbf{x}, \mathbf{y}) u(\mathbf{y}) + \phi(\mathbf{x}), \\
a(\mathbf{x}) &= \frac{2}{h_x^2} + \frac{2}{h_y^2} - \frac{1}{12} \left(h_x^2 + h_y^2\right) \frac{4}{h_x^2 h_y^2} = \frac{5}{3} \left(\frac{1}{h_x^2} + \frac{1}{h_y^2}\right) > 0, \\
b(\mathbf{x}, \mathbf{y}) &= \frac{1}{h_x^2} - \frac{1}{12} \left(h_x^2 + h_y^2\right) \frac{2}{h_x^2 h_y^2} = \frac{1}{6} \left(\frac{5}{h_x^2} - \frac{1}{h_y^2}\right), \ i \pm 1, j, \\
& \hspace{5cm} \text{(left, right node)} \\
b(\mathbf{x}, \mathbf{y}) &= \frac{1}{6} \left(-\frac{1}{h_x^2} + \frac{5}{h_y^2}\right), \ i, j \pm 1, \ \text{(bottom, top node)} \\
b(\mathbf{x}, \mathbf{y}) &= \frac{1}{12} \left(\frac{1}{h_x^2} + \frac{1}{h_y^2}\right), \ i \pm 1, j \pm 1, \ \text{(other neighbors)}.
\end{aligned}
$$

Hence, the assumptions for the discrete maximum principle, see Remark 2.15, are satisfied only if

$$
\frac{1}{\sqrt{5}} < \frac{h_x}{h_y} < \sqrt{5}.
$$

Consequently, the ratio of the grid widths has to be bounded and it has to be of order one. In this case, one speaks of an isotropic grid. □

## 2.7 Summary

**Remark 2.45** *Summary.*

- Finite difference methods are the simplest approach for discretizing partial differential equations. The derivatives are just approximated by difference quotients.
- They are very popular in the engineering community.
- One large drawback are the difficulties in approximating domains which are not of tensor-product type. However, in the engineering communities, a number of strategies have been developed to deal with this issue in practice.
- Another drawback arises from the point of view of numerical analysis. The numerical analysis of finite difference methods is mainly based on Taylor series expansions. For this tool to be applicable, one has to assume a high regularity of the solution. These assumptions are generally not realistic.
- In Numerical Mathematics, one considers often other schemes then finite difference methods. However, there are problems, where finite difference methods can compete with other discretizations, like finite element methods.

□

# Chapter 3

# Introduction to Sobolev Spaces

**Remark 3.1** *Contents.* Sobolev spaces are the basis of the theory of weak or variational forms of partial differential equations. A very popular approach for discretizing partial differential equations, the finite element method, is based on variational forms. In this chapter, a short introduction into Sobolev spaces will be given. Recommended literature are the books Adams (1975); Adams and Fournier (2003) and Evans (2010). □

## 3.1  Elementary Inequalities

**Lemma 3.2 Inequality for strictly monotonically increasing function.** *Let $f : \mathbb{R}_+ \cup \{0\} \to \mathbb{R}$ be a continuous and strictly monotonically increasing function with $f(0) = 0$ and $f(x) \to \infty$ for $x \to \infty$. Then, for all $a, b \in \mathbb{R}_+ \cup \{0\}$ it is*

$$ab \leq \int_0^a f(x) \ dx + \int_0^b f^{-1}(y) \ dy,$$

*where $f^{-1}(y)$ is the inverse of $f(x)$.*

**Proof:** Since $f(x)$ is strictly monotonically increasing, the inverse function exists. The proof is based on a geometric argument, see Figure 3.1.



Figure 3.1: Sketch to the proof of Lemma 3.2.

Consider the interval $(0, a)$ on the $x$-axis and the interval $(0, b)$ on the $y$-axis. Then, the area of the corresponding rectangle is given by $ab$, $\int_0^a f(x) \ dx$ is the area below the curve, and $\int_0^b f^{-1}(y) \ dy$ is the area between the positive $y$-axis and the curve. From Figure 3.1, the inequality follows immediately. The equal sign holds only iff $f(a) = b$. ∎

**Remark 3.3** *Young's inequality.* Young's inequality

$$ab \le \frac{\varepsilon}{2}a^2 + \frac{1}{2\varepsilon}b^2 \quad \forall\, a,b,\varepsilon \in \mathbb{R}_+ \tag{3.1}$$

follows from Lemma 3.2 with $f(x) = \varepsilon x$, $f^{-1}(y) = \varepsilon^{-1}y$. It is also possible to derive this inequality from the binomial theorem. For proving the generalized Young inequality

$$ab \le \frac{\varepsilon^p}{p}a^p + \frac{1}{q\varepsilon^q}b^q, \quad \forall\, a,b,\varepsilon \in \mathbb{R}_+ \tag{3.2}$$

with $p^{-1}+q^{-1}=1, p,q \in (1,\infty)$, one chooses $f(x) = x^{p-1}$, $f^{-1}(y) = y^{1/(p-1)}$ and applies Lemma 3.2 with intervals where the upper bounds are given by $\varepsilon a$ and $\varepsilon^{-1}b$.
□

**Remark 3.4** *Cauchy–Schwarz inequality.* The Cauchy[1]–Schwarz[2] inequality (for vectors, for sums)

$$|(\mathbf{x},\mathbf{y})| \le \|\mathbf{x}\|_2 \|\mathbf{y}\|_2 \ \forall\, \mathbf{x},\mathbf{y} \in \mathbb{R}^n, \tag{3.3}$$

where $(\cdot,\cdot)$ is the Euclidean product and $\|\cdot\|_2$ the Euclidean norm, is well known. One can prove this inequality with the help of Young's inequality.

First, it is clear that the Cauchy–Schwarz inequality is correct if one of the vectors is the zero vector. Now, let $\mathbf{x},\mathbf{y}$ with $\|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1$. One obtains with the triangle inequality and Young's inequality (3.1)

$$|(\mathbf{x},\mathbf{y})| = \left|\sum_{i=1}^n x_i y_i\right| \le \sum_{i=1}^n |x_i|\,|y_i| \le \frac{1}{2}\sum_{i=1}^n |x_i|^2 + \frac{1}{2}\sum_{i=1}^n |y_i|^2 = 1.$$

Hence, the Cauchy–Schwarz inequality is correct for $\mathbf{x},\mathbf{y}$. Last, one considers arbitrary vectors $\tilde{\mathbf{x}} \neq \mathbf{0}, \tilde{\mathbf{y}} \neq \mathbf{0}$. Now, one can utilize the homogeneity of the Cauchy–Schwarz inequality. From the validity of the Cauchy–Schwarz inequality for $\mathbf{x}$ and $\mathbf{y}$, one obtains by a scaling argument

$$\left|(\underbrace{\|\tilde{\mathbf{x}}\|_2^{-1}\,\tilde{\mathbf{x}}}_{\mathbf{x}}, \underbrace{\|\tilde{\mathbf{y}}\|_2^{-1}\,\tilde{\mathbf{y}}}_{\mathbf{y}})\right| \le 1$$

Both vectors $\mathbf{x},\mathbf{y}$ have the Euclidean norm 1, hence

$$\frac{1}{\|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{y}}\|_2}|(\tilde{\mathbf{x}},\tilde{\mathbf{y}})| \le 1 \quad \Longleftrightarrow \quad |(\tilde{\mathbf{x}},\tilde{\mathbf{y}})| \le \|\tilde{\mathbf{x}}\|_2 \|\tilde{\mathbf{y}}\|_2.$$

The generalized Cauchy–Schwarz inequality or Hölder inequality

$$|(\mathbf{x},\mathbf{y})| \le \left(\sum_{i=1}^n |x_i|^p\right)^{1/p} \left(\sum_{i=1}^n |y_i|^q\right)^{1/q}$$

with $p^{-1}+q^{-1}=1, p,q \in (1,\infty)$, can be proved in the same way with the help of the generalized Young inequality.
□

**Definition 3.5** *Lebesgue spaces.* The space of functions which are Lebesgue integrable on $\Omega$ to the power of $p \in [1,\infty)$ is denoted by

$$L^p(\Omega) = \left\{f \ : \ \int_\Omega |f|^p(\mathbf{x})\, d\mathbf{x} < \infty\right\},$$

---

[1] Augustin Louis Cauchy (1789 – 1857)
[2] Hermann Amandus Schwarz (1843 – 1921)

which is equipped with the norm

$$\|f\|_{L^p(\mathbf{x})} = \left(\int_\Omega |f|^p(\mathbf{x})\ d\mathbf{x}\right)^{1/p}.$$

For $p = \infty$, this space is

$$L^\infty(\Omega) = \{f\ :\ |f(\mathbf{x})| < \infty \text{ almost everywhere in } \Omega\}$$

with the norm

$$\|f\|_{L^\infty(\Omega)} = \text{ess sup}_{\mathbf{x}\in\Omega}|f(\mathbf{x})|.$$

$\square$

**Lemma 3.6 Hölder's inequality.** *Let $p^{-1} + q^{-1} = 1, p, q \in [1, \infty]$. If $u \in L^p(\Omega)$ and $v \in L^q(\Omega)$, then it is $uv \in L^1(\Omega)$ and it holds that*

$$\|uv\|_{L^1(\Omega)} \leq \|u\|_{L^p(\Omega)}\|v\|_{L^q(\Omega)}. \tag{3.4}$$

*If $p = q = 2$, then this inequality is also known as Cauchy–Schwarz inequality*

$$\|uv\|_{L^1(\Omega)} \leq \|u\|_{L^2(\Omega)}\|v\|_{L^2(\Omega)}. \tag{3.5}$$

**Proof:** $p, q \in (1, \infty)$. First, one has to show that $|uv(\mathbf{x})|$ can be estimated from above by an integrable function. Setting in the generalized Young inequality (3.2) $\varepsilon = 1$, $a = |u(\mathbf{x})|$, and $b = |v(\mathbf{x})|$ gives

$$|u(\mathbf{x})v(\mathbf{x})| \leq \frac{1}{p}|u(\mathbf{x})|^p + \frac{1}{q}|v(\mathbf{x})|^q.$$

Since the right hand side of this inequality is integrable, by assumption, it follows that $uv \in L^1(\Omega)$. In addition, Hölder's inequality is proved for the case $\|u\|_{L^p(\Omega)} = \|v\|_{L^q(\Omega)} = 1$ using this inequality

$$\int_\Omega |u(\mathbf{x})v(\mathbf{x})|\ d\mathbf{x} \leq \frac{1}{p}\int_\Omega |u(\mathbf{x})|^p\ d\mathbf{x} + \frac{1}{q}\int_\Omega |v(\mathbf{x})|^q\ d\mathbf{x} = 1.$$

The general inequality follows, for the case that both functions do not vanish almost everywhere, with the same homogeneity argument as used for proving the Cauchy–Schwarz inequality of sums. In the case that one of the functions vanishes almost everywhere, (3.4) is trivially satisfied.

$p = 1, q = \infty$. It is

$$\int_\Omega |u(\mathbf{x})v(\mathbf{x})|\ d\mathbf{x} \leq \int_\Omega |u(\mathbf{x})|\, \text{ess sup}_{\mathbf{x}\in\Omega}|v(\mathbf{x})|\ d\mathbf{x} = \|u\|_{L^1(\Omega)}\|v\|_{L^\infty(\Omega)}.$$

$\blacksquare$

## 3.2 Weak Derivative and Distributions

**Remark 3.7** *Contents.* This section introduces a generalization of the derivative which is needed for the definition of weak or variational problems. For an introduction to the topic of this section, see, e.g., Haroske and Triebel (2008)

Let $\Omega \subset \mathbb{R}^d$ be a domain with boundary $\Gamma = \partial\Omega$, $d \in \mathbb{N}$, $\Omega \neq \emptyset$. A domain is always an open set. $\square$

**Definition 3.8 The space $C_0^\infty(\Omega)$.** The space of infinitely often differentiable real functions with compact (closed and bounded) support in $\Omega$ is denoted by $C_0^\infty(\Omega)$

$$C_0^\infty(\Omega) = \{v\ :\ v \in C^\infty(\Omega),\ \text{supp}(v) \subset \Omega\},$$

where

$$\text{supp}(v) = \overline{\{\mathbf{x}\in\Omega\ :\ v(\mathbf{x})\neq 0\}}.$$

$\square$

**Definition 3.9 Convergence in $C_0^\infty(\Omega)$.** The sequence of functions $\{\phi_n(\mathbf{x})\}_{n=1}^\infty$, $\phi_n \in C_0^\infty(\Omega)$ for all $n$, is said to convergence to the zero functions if and only if

a) $\exists K \subset \Omega, K$ compact (closed and bounded) with $\text{supp}(\phi_n) \subset K$ for all $n$,

b) $D^{\boldsymbol{\alpha}}\phi_n(\mathbf{x}) \to 0$ for $n \to \infty$ on $K$ for all multi-indices $\alpha = (\alpha_1, \ldots, \alpha_d)$, $|\alpha| = \alpha_1 + \ldots + \alpha_d$.

It is

$$\lim_{n\to\infty} \phi_n(\mathbf{x}) = \phi(\mathbf{x}) \quad \Longleftrightarrow \quad \lim_{n\to\infty}(\phi_n(\mathbf{x}) - \phi(\mathbf{x})) = 0.$$

$\square$

**Definition 3.10 Weak derivative.** Let $f, F \in L^1_{\text{loc}}(\Omega)$. ($L^1_{\text{loc}}(\Omega)$: for each compact subset $\Omega' \subset \Omega$ it holds

$$\int_{\Omega'} |u(\mathbf{x})| \ d\mathbf{x} < \infty \ \forall \ u \in L^1_{\text{loc}}(\Omega).)$$

If for all functions $g \in C_0^\infty(\Omega)$ it holds that

$$\int_\Omega F(\mathbf{x})g(\mathbf{x}) \ d\mathbf{x} = (-1)^{|\boldsymbol{\alpha}|} \int_\Omega f(\mathbf{x})D^{\boldsymbol{\alpha}}g(\mathbf{x}) \ d\mathbf{x},$$

then $F(\mathbf{x})$ is called weak derivative of $f(\mathbf{x})$ with respect to the multi-index $\boldsymbol{\alpha}$. $\square$

**Remark 3.11** *On the weak derivative.*

- One uses the same notations for the derivative as in the classical case : $F(\mathbf{x}) = D^{\boldsymbol{\alpha}}f(\mathbf{x})$.
- If $f(\mathbf{x})$ is classically differentiable on $\Omega$, then the classical derivative is also the weak derivative.
- The assumptions on $f(\mathbf{x})$ and $F(\mathbf{x})$ are such that the integrals in the definition of the weak derivative are well defined. In particular, since the test functions vanish in a neighborhood of the boundary, the behavior of $f(\mathbf{x})$ and $F(\mathbf{x})$ if $\mathbf{x}$ approaches the boundary is not of importance.
- The main aspect of the weak derivative is due to the fact that the (Lebesgue) integral is not influenced from the values of the functions on a set of (Lebesgue) measure zero. Hence, the weak derivative is uniquely defined only up to a set of measure zero. It follows that $f(\mathbf{x})$ might be not classically differentiable on a set of measure zero, e.g., in a point, but it can still be weakly differentiable.
- The weak derivative is uniquely determined, in the sense described above.

$\square$

**Example 3.12** *Weak derivative.* The weak derivative of the function $f(x) = |x|$ is

$$f'(x) = \begin{cases} -1 & x < 0 \\ 0 & x = 0 \\ 1 & x > 0 \end{cases}$$

In $x = 0$, one can use also any other real number. The proof of this statement follows directly from the definition and it is left as an exercise. $\square$

**Definition 3.13 Distribution.** A continuous linear functional defined on $C_0^\infty(\Omega)$ is called distribution. The set of all distributions is denoted by $(C_0^\infty(\Omega))'$.

Let $u \in C_0^\infty(\Omega)$ and $\psi \in (C_0^\infty(\Omega))'$, then the following notation is used for the application of the distribution to the function

$$\psi(u(\mathbf{x})) = \langle \psi, u \rangle \in \mathbb{R}.$$

$\square$

**Remark 3.14** *On distributions.* Distributions are a generalization of functions. They assign each function from $C_0^\infty(\Omega)$ a real number. □

**Example 3.15** *Regular distribution.* Let $u(\mathbf{x}) \in L_{\mathrm{loc}}^1(\Omega)$. Then, a distribution is defined by

$$\int_\Omega u(\mathbf{x})\phi(\mathbf{x})\,d\mathbf{x} = \langle \psi, \phi \rangle,\ \forall \phi \in C_0^\infty(\Omega).$$

This distribution will be identified with $u(\mathbf{x}) \in L_{\mathrm{loc}}^1(\Omega)$.

Distributions with such an integral representation are called regular, otherwise they are called singular. □

**Example 3.16** *Dirac distribution.* Let $\boldsymbol{\xi} \in \Omega$ fixed, then

$$\langle \delta_{\boldsymbol{\xi}}, \phi \rangle = \phi(\boldsymbol{\xi})\ \forall\ \phi \in C_0^\infty(\Omega)$$

defines a singular distribution, the so-called Dirac distribution or $\delta$-distribution. It is denoted by $\delta_{\boldsymbol{\xi}} = \delta(\mathbf{x} - \boldsymbol{\xi})$. □

**Definition 3.17** *Derivatives of distributions.* Let $\phi \in (C_0^\infty(\Omega))'$ be a distribution. The distribution $\psi \in (C_0^\infty(\Omega))'$ is called derivative in the sense of distributions or distributional derivative of $\phi$ if

$$\langle \psi, u \rangle = (-1)^{|\boldsymbol{\alpha}|} \langle \phi, D^{\boldsymbol{\alpha}} u \rangle\ \forall u \in C_0^\infty(\Omega),$$

$\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_d),\ \alpha_j \geq 0, j = 1, \ldots, d,\ |\boldsymbol{\alpha}| = \alpha_1 + \ldots + \alpha_d$. □

**Remark 3.18** *On derivatives of distributions.* Each distribution has derivatives in the sense of distributions of arbitrary order.

If the derivative in the sense of distributions $D^\alpha u(\mathbf{x})$ of $u(\mathbf{x}) \in L_{\mathrm{loc}}^1(\Omega)$ is a regular distribution, then also the weak derivative of $u(\mathbf{x})$ exists and both derivatives are identified. □

## 3.3 Lebesgue Spaces and Sobolev Spaces

**Remark 3.19** *On the spaces $L^p(\Omega)$.* These spaces were introduced in Definition 3.5.
- The elements of $L^p(\Omega)$ are, strictly speaking, equivalence classes of functions which are different only on a set of Lebesgue measure zero.
- The spaces $L^p(\Omega)$ are Banach spaces (complete normed spaces). A space $X$ is complete, if each so-called Cauchy sequence $\{u_n\}_{n=0}^\infty \in X$, i.e., for all $\varepsilon > 0$ there is an index $n_0(\varepsilon)$ such that for all $i, j > n_0(\varepsilon)$

$$\|u_i - u_j\|_X < \varepsilon.$$

converges and the limit is an element of $X$.
- The space $L^2(\Omega)$ becomes a Hilbert spaces with the inner product

$$(f, g) = \int_\Omega f(\mathbf{x})g(\mathbf{x})\,d\mathbf{x}, \quad \|f\|_{L^2} = (f, f)^{1/2}, \quad f, g \in L^2(\Omega).$$

- The dual space of a space $X$ is the space of all bounded linear functionals defined on $X$. Let $\Omega$ be a domain with sufficiently smooth boundary $\Gamma$. of the Lebesgue spaces $L^p(\Omega)$, $p \in [1, \infty]$, then

$$
\begin{aligned}
(L^p(\Omega))' &= L^q(\Omega)\ \text{ with }\ p, q \in (1, \infty),\ \frac{1}{p} + \frac{1}{q} = 1, \\
(L^1(\Omega))' &= L^\infty(\Omega), \\
(L^\infty(\Omega))' &\neq L^1(\Omega).
\end{aligned}
$$

The spaces $L^1(\Omega)$, $L^\infty(\Omega)$ are not reflexive, i.e., the dual space of the dual space is not the original space again.

$\square$

**Definition 3.20 Sobolev[3] spaces.** Let $k \in \mathbb{N} \cup \{0\}$ and $p \in [1, \infty]$, then the Sobolev space $W^{k,p}(\Omega)$ is defined by

$$W^{k,p}(\Omega) := \{u \in L^p(\Omega) \ : \ D^{\boldsymbol{\alpha}} u \in L^p(\Omega) \ \forall \ \boldsymbol{\alpha} \text{ with } |\boldsymbol{\alpha}| \le k\}.$$

This space is equipped with the norm

$$\|u\|_{W^{k,p}(\Omega)} := \sum_{|\boldsymbol{\alpha}| \le k} \|D^{\boldsymbol{\alpha}} u\|_{L^p(\Omega)}. \tag{3.6}$$

$\square$

**Remark 3.21** *On the spaces $W^{k,p}(\Omega)$.*
- Definition 3.20 has the following meaning. From $u \in L^p(\Omega)$, $p \in [1, \infty)$, it follows in particular that $u \in L^1_{\mathrm{loc}}(\Omega)$, such that $u(\mathbf{x})$ defines (represents) a distribution. Then, all derivatives $D^{\boldsymbol{\alpha}} u$ exist in the sense of distributions. The statement $D^{\boldsymbol{\alpha}} u \in L^p(\Omega)$ means that the distribution $D^{\boldsymbol{\alpha}} u \in (C_0^\infty(\Omega))'$ can be represented by a function from $L^p(\Omega)$.
- One can add elements from $W^{k,p}(\Omega)$ and one can multiply them with real numbers. The result is again a function from $W^{k,p}(\Omega)$. With this property, the space $W^{k,p}(\Omega)$ becomes a vector space (linear space). It is straightforward to check that (3.6) is a norm. (*exercise*)
- It is $D^{\boldsymbol{\alpha}} u(\mathbf{x}) = u(\mathbf{x})$ for $\boldsymbol{\alpha} = (0, \ldots, 0)$ and $W^{0,p}(\Omega) = L^p(\Omega)$.
- The spaces $W^{k,p}(\Omega)$ are Banach spaces.
- Sobolev spaces have for $p \in [1, \infty)$ a countable basis $\{\varphi_n(\mathbf{x})\}_{n=1}^\infty$ (Schauder basis), i.e., each element $u(\mathbf{x})$ can be written in the form

$$u(\mathbf{x}) = \sum_{n=1}^\infty u_n \varphi_n(\mathbf{x}), \quad u_n \in \mathbb{R} \ n = 1, \ldots, \infty.$$

- Sobolev spaces are uniformly convex for $p \in (1, \infty)$, i.e., for each $\varepsilon \in (0, 2]$ (note that the largest distance in the ball is equal to 2) there is a $\delta(\varepsilon) > 0$ such that for all $u, v \in W^{k,p}(\Omega)$ with $\|u\|_{W^{k,p}(\Omega)} = \|v\|_{W^{k,p}(\Omega)} = 1$, and $\|u - v\|_{W^{k,p}(\Omega)} > \varepsilon$ it holds that $\left\|\frac{u+v}{2}\right\|_{W^{k,p}(\Omega)} \le 1 - \delta(\varepsilon)$, see Figure 3.2 for an illustration.
- Sobolev spaces are reflexive for $p \in (1, \infty)$.
- On can show that $C^\infty(\Omega)$ is dense in $W^{k,p}(\Omega)$, e.g., see (Alt, 1999, Satz 1.21, Satz 2.10) or (Adams, 1975, Lemma 3.15). With this property, one can characterize the Sobolev spaces $W^{k,p}(\Omega)$ as completion of the functions from $C^\infty(\Omega)$ with respect to the norm (3.6). For domains with smooth boundary, one can even show that $C^\infty(\overline{\Omega})$ is dense in $W^{k,p}(\Omega)$.
- The Sobolev space $H^k(\Omega) = W^{k,2}(\Omega)$ is a Hilbert space with the inner product

$$(u, v)_{H^k(\Omega)} = \sum_{|\boldsymbol{\alpha}| \le k} \int_\Omega D^{\boldsymbol{\alpha}} u(\mathbf{x}) D^{\boldsymbol{\alpha}} v(\mathbf{x}) \ d\mathbf{x}$$

and the norm $\|u\|_{H^k(\Omega)} = (u, u)^{1/2}$.

$\square$

---

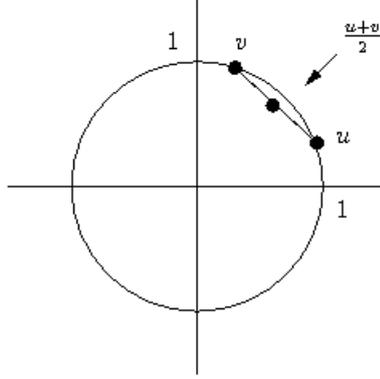[3]Sergei Lvovich Sobolev (1908 – 1989)

Figure 3.2: Illustration of the uniform convexity of Sobolev spaces.

**Definition 3.22 The space $W_0^{k,p}(\Omega)$.** The Sobolev space $W_0^{k,p}(\Omega)$ is defined as the completion of $C_0^\infty(\Omega)$ in the norm of $W^{k,p}(\Omega)$

$$W_0^{k,p}(\Omega) = \overline{C_0^\infty(\Omega)}^{\|\cdot\|_{W^{k,p}(\Omega)}}.$$

$\square$

## 3.4 The Trace of a Function from a Sobolev Space

**Remark 3.23** *Motivation.* This class considers boundary value problems for partial differential equations. In the theory of weak or variational solutions, the solution of the partial differential equation is searched in an appropriate Sobolev space. Then, for the boundary value problem, this solution has to satisfy the boundary condition. However, since the boundary of a domain is a manifold of dimension $(d-1)$, and consequently it has Lebesgue measure zero, one has to clarify how a function from a Sobolev space is defined on this manifold. This definition will be presented in this section. $\square$

**Definition 3.24 Boundary of class $C^{k,\alpha}$.** A bounded domain $\Omega \subset \mathbb{R}^d$ and its boundary $\Gamma$ are of class $C^{k,\alpha}$, $0 \le \alpha \le 1$ if for all $\mathbf{x}_0 \in \Gamma$ there is a ball $B(\mathbf{x}_0, r)$ with $r > 0$ and a bijective map $\psi : B(\mathbf{x}_0, r) \to D \subset \mathbb{R}^d$ such that
1) $\psi(B(\mathbf{x}_0, r) \cap \Omega) \subset \mathbb{R}_+^d$,
2) $\psi(B(\mathbf{x}_0, r) \cap \Gamma) \subset \partial\mathbb{R}_+^d$,
3) $\psi \in C^{k,\alpha}(B(\mathbf{x}_0, r)), \psi^{-1} \in C^{k,\alpha}(D)$, are Hölder continuous.
That means, $\Gamma$ is locally the graph of a function with $d-1$ arguments. (A function $u(\mathbf{x})$ is Hölder continuous if

$$\|u\|_{C^{k,\alpha}(\Omega)} = \sum_{|\boldsymbol{\alpha}| \le k} \|D^{\boldsymbol{\alpha}} u\|_{C(\overline{\Omega})} + \sum_{|\boldsymbol{\alpha}| = k} [D^{\boldsymbol{\alpha}} u]_{C^{0,\alpha}(\overline{\Omega})}$$

with

$$[D^{\boldsymbol{\alpha}} u]_{C^{0,\alpha}(\overline{\Omega})} = \sup_{\mathbf{x},\mathbf{y} \in \Omega} \left\{ \frac{|u(\mathbf{x}) - u(\mathbf{y})|}{|\mathbf{x} - \mathbf{y}|^\alpha} \right\}$$

is finite.) $\square$

**Remark 3.25** *Lipschitz boundary.* It will be generally assumed that the boundary of $\Omega$ is of class $C^{0,1}$. That means, the map is Lipschitz[4] continuous. Such a boundary

---
[4] Rudolf Otto Sigismund Lipschitz (1832 – 1903)

is simply called Lipschitz boundary and the domain is called Lipschitz domain. An important feature of a Lipschitz boundary is that the outer normal vector is defined almost everywhere at the boundary and it is almost everywhere continuous.  □

**Example 3.26** *On Lipschitz domains.*
- Domains with Lipschitz boundary are, for example, balls or polygonal domains in two dimensions where the domain is always on one side of the boundary.
- A domain which is not a Lipschitz domain is a circle with a slit

$$\Omega = \{(x,y) \ : \ x^2 + y^2 < 1\} \setminus \{(x,y) \ : \ x \geq 0, y = 0\}.$$

At the slit, the domain is on both sides of the boundary.

- In three dimension, a polyhedral domain is not not necessarily a Lipschitz domain. For instance, if the domain is build of two bricks which are laying on each other like in Figure 3.3, then the boundary is not Lipschitz continuous where the edge of one brick meets the edge of the other brick.
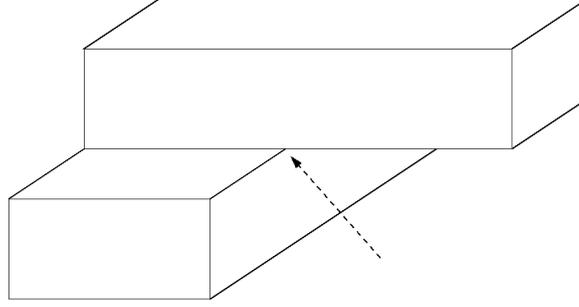
  □



Figure 3.3: Polyhedral domain in three dimensions which is not Lipschitz continuous (at the corner where the arrow points to).

**Theorem 3.27 Trace theorem.** *Let $\Omega \subset \mathbb{R}^d$, $d \geq 2$, with a Lipschitz boundary. Then, there is exactly one linear and continuous operator $\gamma : W^{1,p}(\Omega) \to L^p(\Gamma)$, $p \in [1,\infty)$, which gives for functions $u \in C(\overline{\Omega}) \cap W^{1,p}(\Omega)$ the classical boundary values*

$$\gamma u(\mathbf{x}) = u(\mathbf{x}), \ \mathbf{x} \in \Gamma, \ \forall \, u \in C(\overline{\Omega}) \cap W^{1,p}(\Omega),$$

*i.e., $\gamma u(\mathbf{x}) = u(\mathbf{x})|_{\mathbf{x} \in \Gamma}$.*

**Proof:** The proof can be found in the literature, e.g., in Adams (1975); Adams and Fournier (2003). ∎

**Remark 3.28** *On the trace.* The operator $\gamma$ is called trace or trace operator.
- Since a linear and continuous operator is bounded, there is a constant $C > 0$ with

$$\|\gamma u\|_{L^p(\Gamma)} \leq C \, \|u\|_{W^{1,p}(\Omega)} \ \forall \ u \in W^{1,p}(\Omega)$$

  or

$$\|\gamma\|_{\mathcal{L}(W^{1,p}(\Omega), L^p(\Gamma))} \leq C.$$

- By definition of the trace, one gets for $u \in C(\overline{\Omega})$ the classical boundary values. By the density of $C^\infty(\overline{\Omega})$ in $W^{1,p}(\Omega)$ for domains with smooth boundary, it follows that $C(\overline{\Omega})$ is also dense in $W^{1,p}(\Omega)$ such that for all $u \in W^{1,p}(\Omega)$ there is a sequence $\{u_n\}_{n=1}^{\infty} \in C(\overline{\Omega})$ with $u_n \to u$ in $W^{1,p}(\Omega)$. Then, the trace of $u$ is defined to be $\gamma u = \lim_{k\to\infty}(\gamma u_k)$.

- It is

$$
\begin{aligned}
\gamma u(\mathbf{x}) &= 0 \quad \forall\, u \in W_0^{1,p}(\Omega), \\
\gamma D^{\boldsymbol{\alpha}} u(\mathbf{x}) &= 0 \quad \forall\, u \in W_0^{k,p}(\Omega), |\boldsymbol{\alpha}| \le k-1. \tag{3.7}
\end{aligned}
$$

$\qquad\qquad\square$

## 3.5 Sobolev Spaces with Non-Integer and Negative Exponents

**Remark 3.29** *Motivation.* Sobolev spaces with non-integer and negative exponents are important in the theory of variational formulations of partial differential equations.

Let $\Omega \subset \mathbb{R}^d$ be a domain and $p \in (1,\infty)$ mit $p^{-1} + q^{-1} = 1$. $\qquad\square$

**Definition 3.30 The space $W^{-k,q}(\Omega)$.** The space $W^{-k,q}(\Omega), k \in \mathbb{N} \cup \{0\}$, contains distributions which are defined on $W^{k,p}(\Omega)$

$$
W^{-k,q}(\Omega) = \left\{ \varphi \in (C_0^\infty(\Omega))' \; : \; \|\varphi\|_{W^{-k,q}} < \infty \right\}
$$

with

$$
\|\varphi\|_{W^{-k,q}} = \sup_{u \in C_0^\infty(\Omega), u \ne 0} \frac{\langle \varphi, u \rangle}{\|u\|_{W^{k,p}(\Omega)}}.
$$

$\qquad\qquad\square$

**Remark 3.31** *On the spaces $W^{-k,p}(\Omega)$.*
- It is $W^{-k,q}(\Omega) = \left[ W_0^{k,p}(\Omega) \right]'$, i.e., $W^{-k,q}(\Omega)$ can be identified with the dual space of $W_0^{k,p}(\Omega)$. In particular it is $H^{-1}(\Omega) = \left( H_0^1(\Omega) \right)'$.
- It is

$$
\ldots \subset W^{2,p}(\Omega) \subset W^{1,p}(\Omega) \subset L^p(\Omega) \subset W^{-1,q}(\Omega) \subset W^{-2,q}(\Omega) \ldots
$$

$\qquad\qquad\square$

**Definition 3.32 Sobolev–Slobodeckij space.** Let $s \in \mathbb{R}$, then the Sobolev–Slobodeckij or Sobolev space $H^s(\Omega)$ is defined as follows:
- $s \in \mathbb{Z}$. $H^s(\Omega) = W^{s,2}(\Omega)$.
- $s > 0$ with $s = k + \sigma$, $k \in \mathbb{N} \cup \{0\}$, $\sigma \in (0,1)$. The space $H^s(\Omega)$ contains all functions $u$ for which the following norm is finite:

$$
\|u\|_{H^s(\Omega)}^2 = \|u\|_{H^k(\Omega)}^2 + |u|_{k+\sigma}^2 ,
$$

with

$$
\begin{aligned}
(u,v)_{H^s(\Omega)} &= (u,v)_{H^k} + (u,v)_{k+\sigma}, \quad |u|_{k+\sigma}^2 = (u,u)_{k+\sigma}, \\
(u,v)_{k+\sigma} &= \sum_{|\boldsymbol{\alpha}|=k} \int_\Omega \int_\Omega \frac{(D^{\boldsymbol{\alpha}} u(\mathbf{x}) - D^{\boldsymbol{\alpha}} u(\mathbf{y}))\,(D^{\boldsymbol{\alpha}} v(\mathbf{x}) - D^{\boldsymbol{\alpha}} v(\mathbf{y}))}{\|\mathbf{x} - \mathbf{y}\|_2^{d+2\sigma}} \; d\mathbf{x} d\mathbf{y},
\end{aligned}
$$

- $s < 0$. $H^s(\Omega) = \left[ H_0^{-s}(\Omega) \right]'$ with $H_0^{-s}(\Omega) = \overline{C_0^\infty(\Omega)}^{\|\cdot\|_{H^{-s}(\Omega)}}$.

$\qquad\qquad\square$

## 3.6 Theorem on Equivalent Norms

**Definition 3.33 Equivalent norms.** Two norms $\|\cdot\|_1$ and $\|\cdot\|_2$ on the linear space $X$ are said to be equivalent if there are constants $C_1$ and $C_2$ such that

$$C_1 \|u\|_1 \le \|u\|_2 \le C_2 \|u\|_1 \ \forall \, u \in X.$$

$\square$

**Remark 3.34** *On equivalent norms.* Many important properties, like continuity or convergence, do not change if an equivalent norm is considered. $\square$

**Theorem 3.35 Equivalent norms in $W^{k,p}(\Omega)$.** *Let $\Omega \subset \mathbb{R}^d$ be a domain with Lipschitz boundary $\Gamma$, $p \in [1,\infty]$, and $k \in \mathbb{N}$. Let $\{f_i\}_{i=1}^l$ be a system with the following properties:*
*1) $f_i : W^{k,p}(\Omega) \to \mathbb{R}_+ \cup \{0\}$ is a semi norm,*
*2) $\exists C_i > 0$ with $0 \le f_i(v) \le C_i \|v\|_{W^{k,p}(\Omega)}$, $\forall \, v \in W^{k,p}(\Omega)$,*
*3) $f_i$ is a norm on the polynomials of degree $k-1$, i.e., if for $v \in P_{k-1} = \left\{\sum_{|\boldsymbol{\alpha}| \le k-1} C_{\boldsymbol{\alpha}} x^{\boldsymbol{\alpha}}\right\}$ it holds that $f_i(v) = 0$, $i = 1, \ldots, l$, then it is $v \equiv 0$.*
*Then, the norm $\|\cdot\|_{W^{k,p}(\Omega)}$ defined in (3.6) and the norm*

$$\|u\|'_{W^{k,p}(\Omega)} := \left( \sum_{i=1}^l f_i^p(u) + |u|_{W^{k,p}(\Omega)}^p \right)^{1/p} \quad with$$

$$|u|_{W^{k,p}(\Omega)} = \left( \sum_{|\boldsymbol{\alpha}|=k} \int_\Omega |D^{\boldsymbol{\alpha}} u(\mathbf{x})|^p \ d\mathbf{x} \right)^{1/p}$$

*are equivalent.*

**Remark 3.36** *On semi norms.* For a semi norm $f_i(\cdot)$, one cannot conclude from $f_i(v) = 0$ that $v = 0$. The third assumptions however states, that this conclusion can be drawn for all polynomials up to a certain degree. $\square$

**Example 3.37** *Equivalent norms in Sobolev spaces.*
- The following norms are equivalent to the standard norm in $W^{1,p}(\Omega)$:

$$\text{a)} \quad \|u\|'_{W^{1,p}(\Omega)} = \left( \left| \int_\Omega u \ d\mathbf{x} \right|^p + |u|_{W^{1,p}(\Omega)}^p \right)^{1/p},$$

$$\text{b)} \quad \|u\|'_{W^{1,p}(\Omega)} = \left( \left| \int_\Gamma u \ d\mathbf{s} \right|^p + |u|_{W^{1,p}(\Omega)}^p \right)^{1/p},$$

$$\text{c)} \quad \|u\|'_{W^{1,p}(\Omega)} = \left( \int_\Gamma |u|^p \ d\mathbf{s} + |u|_{W^{1,p}(\Omega)}^p \right)^{1/p}.$$

- In $W^{k,p}(\Omega)$ it is

$$\|u\|'_{W^{k,p}(\Omega)} = \left( \sum_{i=0}^{k-1} \int_\Gamma \left| \frac{\partial^i u}{\partial \mathbf{n}^i} \right|^p \ d\mathbf{s} + |u|_{W^{k,p}(\Omega)}^p \right)^{1/p}$$

equivalent to the standard norm. Here, $\mathbf{n}$ denotes the outer normal on $\Gamma$ with $\|\mathbf{n}\|_2 = 1$.

- In the case $W_0^{k,p}(\Omega)$, one does not need the regularity of the boundary. It is

$$\|u\|'_{W_0^{k,p}(\Omega)} = |u|_{W^{k,p}(\Omega)},$$

i.e., in the spaces $W_0^{k,p}(\Omega)$ the standard semi norm is equivalent to the standard norm.

In particular, it is for $u \in H_0^1(\Omega)$ $(k = 1, p = 2)$

$$C_1 \|u\|_{H^1(\Omega)} \leq \|\nabla u\|_{L^2(\Omega)} \leq C_2 \|u\|_{H^1(\Omega)}.$$

It follows that there is a constant $C > 0$ such that

$$\|u\|_{L^2(\Omega)} \leq C \|\nabla u\|_{L^2(\Omega)} \quad \forall \, u \in H_0^1(\Omega). \tag{3.8}$$

$\square$

## 3.7 Some Inequalities in Sobolev Spaces

**Remark 3.38** *Motivation.* This section presents a generalization of the last part of Example 3.37. It will be shown that for inequalities of type (3.8) it is not necessary that the trace vanishes on the complete boundary.

Let $\Omega \subset \mathbb{R}^d$ be a bounded domain with boundary $\Gamma$ and let $\Gamma_1 \subset \Gamma$ with $\text{meas}_{\mathbb{R}^{d-1}}(\Gamma_1) = \int_{\Gamma_1} d\mathbf{s} > 0$.

One considers the space

$$
\begin{aligned}
V_0 &= \left\{ v \in W^{1,p}(\Omega) \; : \; v|_{\Gamma_1} = 0 \right\} \subset W^{1,p}(\Omega) \text{ if } \Gamma_1 \subset \Gamma, \\
V_0 &= W_0^{1,p}(\Omega) \text{ if } \Gamma_1 = \Gamma
\end{aligned}
$$

with $p \in [1, \infty)$. $\square$

**Lemma 3.39 Friedrichs[5] inequality, Poincaré[6] inequality, Poincaré–Friedrichs inequality.** *Let $p \in [1, \infty)$ and $\text{meas}_{\mathbb{R}^{d-1}}(\Gamma_1) > 0$. Then it is for all $u \in V_0$*

$$\int_\Omega |u(\mathbf{x})|^p \; d\mathbf{x} \leq C_P \int_\Omega \|\nabla u(\mathbf{x})\|_2^p \; d\mathbf{x}, \tag{3.9}$$

*where $\|\cdot\|_2$ is the Euclidean vector norm.*

**Proof:** The inequality will be proved with the theorem on equivalent norms, Theorem 3.35. Let $f_1(u) \; : \; W^{1,p}(\Omega) \to \mathbb{R}_+ \cup \{0\}$ with

$$f_1(u) = \left( \int_{\Gamma_1} |u(\mathbf{s})|^p \; d\mathbf{s} \right)^{1/p}.$$

This functions has the following properties:
1) $f_1(u)$ is a semi norm.
2) It is

$$
\begin{aligned}
0 \quad &\leq \quad f_1(u) = \left( \int_{\Gamma_1} |u(\mathbf{s})|^p \; d\mathbf{s} \right)^{1/p} \leq \left( \int_{\Gamma} |u(\mathbf{s})|^p \; d\mathbf{s} \right)^{1/p} \\
&= \quad \|u\|_{L^p(\Gamma)} = \|\gamma u\|_{L^p(\Gamma)} \leq C \|u\|_{W^{1,p}(\Omega)}.
\end{aligned}
$$

The last estimate follows from the continuity of the trace operator.

---

[5] Friedrichs
[6] Poincaré

3) Let $v \in P_0$, i.e., $v$ is a constant. Then, one obtains from

$$0 = f_1(v) = \left(\int_{\Gamma_1} |v(\mathbf{s})|^p \ d\mathbf{s}\right)^{1/p} = |v| \left(\mathrm{meas}_{\mathbb{R}^{d-1}}(\Gamma_1)\right)^{1/p},$$

that $|v| = 0$.

Hence, all assumptions of Theorem 3.35 are satisfied. That means, there are two constants $C_1$ and $C_2$ with

$$C_1 \underbrace{\left(\int_{\Gamma_1} |u(\mathbf{s})|^p \ d\mathbf{s} + \int_{\Omega} \|\nabla u(\mathbf{x})\|_2^p \ d\mathbf{x}\right)^{1/p}}_{\|u\|'_{W^{1,p}(\Omega)}} \leq \|u\|_{W^{1,p}(\Omega)} \leq C_2 \|u\|'_{W^{1,p}(\Omega)} \ \forall \ u \in W^{1,p}(\Omega).$$

In particular, it follows that

$$\int_{\Omega} |u(\mathbf{x})|^p \ d\mathbf{x} + \int_{\Omega} \|\nabla u(\mathbf{x})\|_2^p \ d\mathbf{x} \leq C_2^p \left(\int_{\Gamma_1} |u(\mathbf{s})|^p \ d\mathbf{s} + \int_{\Omega} \|\nabla u(\mathbf{x})\|_2^p \ d\mathbf{x}\right)$$

or

$$\int_{\Omega} |u(\mathbf{x})|^p \ d\mathbf{x} \leq C_P \left(\int_{\Gamma_1} |u(\mathbf{s})|^p \ d\mathbf{s} + \int_{\Omega} \|\nabla u(\mathbf{x})\|_2^p \ d\mathbf{x}\right)$$

with $C_P = C_2^p$. Since $u \in V_0$ vanishes on $\Gamma_1$, the statement of the lemma is proved. ∎

**Remark 3.40** *On the Poincaré–Friedrichs inequality.* In the space $V_0$ becomes $|\cdot|_{W^{1,p}}$ a norm which is equivalent to $\|\cdot\|_{W^{1,p}(\Omega)}$. The classical Poincaré–Friedrichs inequality is given for $\Gamma_1 = \Gamma$ and $p = 2$

$$\|u\|_{L^2} \leq C_P \|\nabla u\|_{L^2} \ \forall \ u \in H_0^1(\Omega),$$

where the constant depends only on the diameter of the domain $\Omega$. □

**Lemma 3.41 Another inequality of Poincaré–Friedrichs type.** *Let $\Omega' \subset \Omega$ with $\mathrm{meas}_{\mathbb{R}^d}(\Omega') = \int_{\Omega'} \ d\mathbf{x} > 0$, then for all $u \in W^{1,p}(\Omega)$ it is*

$$\int_{\Omega} |u(\mathbf{x})|^p \ d\mathbf{x} \leq C \left(\left|\int_{\Omega'} u(\mathbf{x}) \ d\mathbf{x}\right|^p + \int_{\Omega} \|\nabla u(\mathbf{x})\|_2^p \ d\mathbf{x}\right).$$

**Proof:** *Exercise.* ∎

## 3.8 The Gaussian Theorem

**Remark 3.42** *Motivation.* The Gaussian theorem is the generalization of the integration by parts from calculus. This operation is very important for the theory of weak or variational solutions of partial differential equations. One has to study, under which conditions on the regularity of the domain and of the functions it is well defined. □

**Theorem 3.43 Gaussian theorem.** *Let $\Omega \subset \mathbb{R}^d, d \geq 2$, be a bounded domain with Lipschitz boundary $\Gamma$. Then, the following identity holds for all $u \in W^{1,1}(\Omega)$*

$$\int_{\Omega} \partial_i u(\mathbf{x}) \ d\mathbf{x} = \int_{\Gamma} u(\mathbf{s}) \mathbf{n}_i(\mathbf{s}) \ d\mathbf{s}, \tag{3.10}$$

*where $\mathbf{n}$ is the unit outer normal vector on $\Gamma$.*

**Proof:** sketch. First of all, one proves the statement for functions from $C^1(\overline{\Omega})$. This proof is somewhat longer and it is referred to the literature, e.g., Evans (2010).

The space $C^1(\overline{\Omega})$ is dense in $W^{1,1}(\Omega)$, see Remark 3.21. Hence, for all $u \in W^{1,1}(\Omega)$ there is a sequence $\{u_n\}_{n=1}^{\infty} \in C^1(\overline{\Omega})$ with

$$\lim_{n \to \infty} \|u - u_n\|_{W^{1,1}(\Omega)} = 0$$

and (3.10) holds for all functions $u_n(\mathbf{x})$. It will be shown that the limit of the left hand side converges to the left hand side of (3.10) and the limit of the right hand side converges to the right hand side of (3.10).

From the convergence in $\|\cdot\|_{W^{1,1}(\Omega)}$, one has in particular

$$\lim_{n \to \infty} \int_{\Omega} \partial_i u_n(\mathbf{x}) \, d\mathbf{x} = \int_{\Omega} \partial_i u(\mathbf{x}) \, d\mathbf{x}.$$

On the other hand, the continuity of the trace operator gives

$$\lim_{n \to \infty} \|u - u_n\|_{L^1(\Gamma)} \leq C \lim_{n \to \infty} \|u - u_n\|_{W^{1,1}} = 0,$$

from what follows that

$$\lim_{n \to \infty} \int_{\Gamma} u_n(\mathbf{s}) \, d\mathbf{s} = \int_{\Gamma} u(\mathbf{s}) \, d\mathbf{s}.$$

Since for a Lipschitz boundary, the normal $\mathbf{n}$ is almost everywhere continuous, it is

$$\lim_{n \to \infty} \int_{\Gamma} u_n(\mathbf{s})\mathbf{n}_i(\mathbf{s}) \, d\mathbf{s} = \int_{\Gamma} u(\mathbf{s})\mathbf{n}_i(\mathbf{s}) \, d\mathbf{s}.$$

Thus, the limits lead to (3.10). ∎

**Corollary 3.44 Vector field.** *Let the conditions of Theorem 3.43 on the domain $\Omega$ be satisfied and let $\mathbf{u} \in \left(W^{1,1}(\Omega)\right)^d$ be a vector field. Then it is*

$$\int_{\Omega} \nabla \cdot \mathbf{u}(\mathbf{x}) \, d\mathbf{x} = \int_{\Gamma} \mathbf{u}(\mathbf{s}) \cdot \mathbf{n}(\mathbf{s}) \, d\mathbf{s}.$$

**Proof:** The statement follows by adding (3.10) from $i = 1$ to $i = d$. ∎

**Corollary 3.45 Integration by parts.** *Let the conditions of Theorem 3.43 on the domain $\Omega$ be satisfied. Consider $u \in W^{1,p}(\Omega)$ and $v \in W^{1,q}(\Omega)$ with $p \in (1, \infty)$ and $\frac{1}{p} + \frac{1}{q} = 1$. Then it is*

$$\int_{\Omega} \partial_i u(\mathbf{x})v(\mathbf{x}) \, d\mathbf{x} = \int_{\Gamma} u(\mathbf{s})v(\mathbf{s})\mathbf{n}_i(\mathbf{s}) \, d\mathbf{s} - \int_{\Omega} u(\mathbf{x})\partial_i v(\mathbf{x}) \, d\mathbf{x}.$$

**Proof:** *exercise.* ∎

**Corollary 3.46 First Green[7]'s formula.** *Let the conditions of Theorem 3.43 on the domain $\Omega$ be satisfied, then it is*

$$\int_{\Omega} \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} = \int_{\Gamma} \frac{\partial u}{\partial \mathbf{n}}(\mathbf{s})v(\mathbf{s}) \, d\mathbf{s} - \int_{\Omega} \Delta u(\mathbf{x})v(\mathbf{x}) \, d\mathbf{x}$$

*for all $u \in H^2(\Omega)$ and $v \in H^1(\Omega)$.*

**Proof:** From the definition of the Sobolev spaces it follows that the integrals are well defined. Now, the proof follows the proof of Corollary 3.45, where one has now to sum over the components. ∎

---

[7]Georg Green (1793 – 1841)

**Remark 3.47** *On the first Green's formula.* The first Green's formula is the formula of integrating by parts once. The boundary integral can be equivalently written in the form

$$\int_\Gamma \nabla u(\mathbf{s}) \cdot \mathbf{n}(\mathbf{s}) v(\mathbf{s}) \, d\mathbf{s}.$$

The formula of integrating by parts twice is called second Green's formula. $\square$

**Corollary 3.48 Second Green's formula.** *Let the conditions of Theorem 3.43 on the domain $\Omega$ be satisfied, then one has*

$$\int_\Omega \left( \Delta u(\mathbf{x}) v(\mathbf{x}) - \Delta v(\mathbf{x}) u(\mathbf{x}) \right) \, d\mathbf{x} = \int_\Gamma \left( \frac{\partial u}{\partial \mathbf{n}}(\mathbf{s}) v(\mathbf{s}) - \frac{\partial v}{\partial \mathbf{n}}(\mathbf{s}) u(\mathbf{s}) \right) \, d\mathbf{s}$$

*for all $u, v \in H^2(\Omega)$.*

## 3.9  Sobolev Imbedding Theorems

**Remark 3.49** *Motivation.* This section studies the question which Sobolev spaces are subspaces of other Sobolev spaces. With this property, called imbedding, it is possible to estimate the norm of a function in the subspace by the norm in the larger space. $\square$

**Lemma 3.50 Imbedding of Sobolev spaces with same integration power $p$ and different orders of the derivative.** *Let $\Omega \subset \mathbb{R}^d$ be a domain with $p \in [1, \infty)$ and $k \leq m$, then it is $W^{m,p}(\Omega) \subset W^{k,p}(\Omega)$.*

**Proof:** The statement of this lemma follows directly from the definition of Sobolev spaces, see Definition 3.20. ∎

**Lemma 3.51 Imbedding of Sobolev spaces with the same order of the derivative $k$ and different integration powers.** *Let $\Omega \subset \mathbb{R}^d$ be a bounded domain, $k \geq 0$, and $p, q \in [1, \infty]$ with $q > p$. Then it is $W^{k,q}(\Omega) \subset W^{k,p}(\Omega)$.*

**Proof:** *exercise.* ∎

**Remark 3.52** *Imbedding of Sobolev spaces with the same order of the derivative $k$ and the same integration power $p$ in imbedded domains.* Let $\Omega \subset \mathbb{R}^d$ be a domain with sufficiently smooth boundary $\Gamma$, $k \geq 0$, and $p \in [1, \infty]$. Then there is a map $E : W^{k,p}(\Omega) \to W^{k,p}(\mathbb{R}^d)$, the so-called (simple) extension, with
- $Ev|_\Omega = v$,
- $\|Ev\|_{W^{k,p}(\mathbb{R}^d)} \leq C \|v\|_{W^{k,p}(\Omega)}$, with $C > 0$,

e.g., see (Adams, 1975, Chapter IV) for details. Likewise, the natural restriction $e : W^{k,p}(\mathbb{R}^d) \to W^{k,p}(\Omega)$ can be defined and it is $\|ev\|_{W^{k,p}(\Omega)} \leq \|v\|_{W^{k,p}(\mathbb{R}^d)}$. $\square$

**Theorem 3.53 A Sobolev inequality.** *Let $\Omega \subset \mathbb{R}^d$ be a bounded domain with Lipschitz boundary $\Gamma$, $k \geq 0$, and $p \in [1, \infty)$ with*

$$\begin{aligned} k &\geq d \quad && \text{for } p = 1, \\ k &> d/p \quad && \text{for } p > 1. \end{aligned}$$

*Then there is a constant $C$ such that for all $u \in W^{k,p}(\Omega)$ it follows that $u \in C_B(\Omega)$, where*

$$C_B(\Omega) = \{ v \in C(\Omega) \ : \ v \text{ is bounded} \},$$

*and it is*

$$\|u\|_{C_B(\Omega)} = \|u\|_{L^\infty(\Omega)} \leq C \|u\|_{W^{k,p}(\Omega)}. \tag{3.11}$$

**Proof:** See literature, e.g., Adams (1975); Adams and Fournier (2003). ∎

**Remark 3.54** *On the Sobolev inequality.* The Sobolev inequality states that each function with sufficiently many weak derivatives (the number depends on the dimension of $\Omega$ and the integration power) can be considered as a continuous and bounded function in $\Omega$. One says that $W^{k,p}(\Omega)$ is imbedded in $C_B(\Omega)$. It is

$$C\left(\overline{\Omega}\right) \subset C_B(\Omega) \subset C(\Omega).$$

Consider $\Omega = (0,1)$ and $f_1(x) = 1/x$ and $f_2(x) = \sin(1/x)$. Then, $f_1 \in C(\Omega)$, $f_1 \notin C_B(\Omega)$ and $f_2 \in C_B(\Omega)$, $f_2 \notin C(\overline{\Omega})$.

Of course, it is possible to apply this theorem to weak derivatives of functions. Then, one obtains imbeddings like $W^{k,p}(\Omega) \to C_B^s(\Omega)$ for $(k-s)p > d, p > 1$. A comprehensive overview on imbeddings can be found in Adams (1975); Adams and Fournier (2003). □

**Example 3.55** $H^1(\Omega)$ *in one dimension.* Let $d = 1$ and $\Omega$ be a bounded interval. Then, each function from $H^1(\Omega)$ $(k = 1, p = 2)$ is continuous and bounded in $\Omega$. □

**Example 3.56** $H^1(\Omega)$ *in higher dimensions.* The functions from $H^1(\Omega)$ are in general not continuous for $d \geq 2$. This property will be shown with the following example.
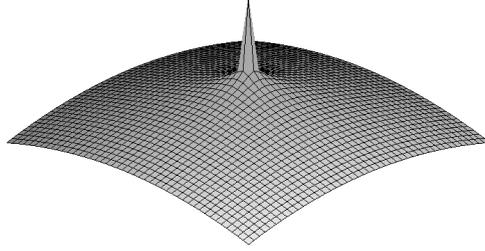


Figure 3.4: The function $f(\mathbf{x})$ of Example 3.56 for $d = 2$.

Let $\Omega = \{\mathbf{x} \in \mathbb{R}^d \ : \ \|\mathbf{x}\|_2 < 1/2\}$ and $f(\mathbf{x}) = \ln|\ln\|\mathbf{x}\|_2|$, see Figure 3.4. For $\|\mathbf{x}\|_2 < 1/2$ it is $|\ln\|\mathbf{x}\|_2| = -\ln\|\mathbf{x}\|_2$ and one gets for $\mathbf{x} \neq \mathbf{0}$

$$\partial_i f(\mathbf{x}) = -\frac{1}{\ln\|\mathbf{x}\|_2}\frac{1}{\|\mathbf{x}\|_2}\frac{x_i}{\|\mathbf{x}\|_2} = -\frac{x_i}{\|\mathbf{x}\|_2^2 \ln\|\mathbf{x}\|_2}.$$

For $p \leq d$, one obtains

$$\left|\frac{\partial f}{\partial x_i}(\mathbf{x})\right|^p = \underbrace{\left|\frac{x_i}{\|\mathbf{x}\|_2}\right|^p}_{\leq 1}\underbrace{\left|\frac{1}{\|\mathbf{x}\|_2 \ln\|\mathbf{x}\|_2}\right|^p}_{\geq e} \leq \left|\frac{1}{\|\mathbf{x}\|_2 \ln\|\mathbf{x}\|_2}\right|^d.$$

The estimate of the second factor can be obtained, e.g., with a discussion of the curve. Using now spherical coordinates, $\rho = e^{-t}$ and $S^{d-1}$ is the unit sphere, yields

$$
\begin{aligned}
\int_\Omega |\partial_i f(\mathbf{x})|^p \, d\mathbf{x} &\leq \int_\Omega \frac{d\mathbf{x}}{\|\mathbf{x}\|_2^d |\ln\|\mathbf{x}\|_2|^d} = \int_{S^{d-1}}\int_0^{1/2}\frac{\rho^{d-1}}{\rho^d |\ln\rho|^d} \, d\rho d\omega \\
&= \operatorname{meas}\left(S^{d-1}\right)\int_0^{1/2}\frac{d\rho}{\rho |\ln\rho|^d} = -\operatorname{meas}\left(S^{d-1}\right)\int_\infty^{\ln 2}\frac{dt}{t^d} < \infty,
\end{aligned}
$$

because of $d \geq 2$.

It follows that $\partial_i f \in L^p(\Omega)$ with $p \leq d$. Analogously, one proves that $f \in L^p(\Omega)$ with $p \leq d$. Altogether, one has $f \in W^{1,p}(\Omega)$ with $p \leq d$. However, it is $f \notin L^\infty(\Omega)$. This example shows that the condition $k > d/p$ for $p > 1$ is sharp.

In particular, it was proved for $p = 2$ that from $f \in H^1(\Omega)$ in general it does not follow that $f \in C(\Omega)$. $\qquad\square$

**Example 3.57** *The assumption of a Lipschitz boundary.* Also the assumption that $\Omega$ is a Lipschitz domain is of importance.

Consider $\Omega = \{(x,y) \in \mathbb{R}^2 \ : \ 0 < x < 1, \ |y| < x^r, r > 1\}$, see Figure 3.5 for $r = 2$.



Figure 3.5: Domain of Example 3.57.

For $u(x,y) = x^{-\varepsilon/p}$ with $0 < \varepsilon < r$ it is

$$\partial_x u = x^{-\varepsilon/p-1}\left(-\frac{\varepsilon}{p}\right) = C(\varepsilon,p)x^{-\varepsilon/p-1}, \ \partial_y u = 0.$$

It follows that

$$
\begin{aligned}
\sum_{|\boldsymbol{\alpha}|=1}\int_\Omega |D^{\boldsymbol{\alpha}} u(x,y)|^p \ dxdy &= C(\varepsilon,p)\int_\Omega x^{-\varepsilon-p} \ dxdy \\
&= C(\varepsilon,p)\int_0^1 x^{-\varepsilon-p}\left(\int_{-x^r}^{x^r} dy\right) dx \\
&= \tilde{C}(\varepsilon,p)\int_0^1 x^{-\varepsilon-p+r} \ dx.
\end{aligned}
$$

This value is finite for $-\varepsilon - p + r > -1$ or for $p < 1 + r - \varepsilon$, respectively. If one chooses $r \geq \varepsilon > 0$, then it is $u \in W^{1,p}(\Omega)$. But for $\varepsilon > 0$ the function $u(\mathbf{x})$ is not bounded in $\Omega$, i.e., $u \notin L^\infty(\Omega)$.

The unbounded values of the function are compensated in the integration by the fact that the neighborhood of the singular point $(0,0)$ possesses a small measure. $\qquad\square$

# Chapter 4

# The Ritz Method and the Galerkin Method

**Remark 4.1** *Contents.* This chapter studies variational or weak formulations of boundary value problems of partial differential equations in Hilbert spaces. The existence and uniqueness of an appropriately defined weak solution will be discussed. The approximation of this solution with the help of finite-dimensional spaces is called Ritz method or Galerkin method. Some basic properties of this method will be proved.

In this chapter, a Hilbert space $V$ will be considered with inner product $a(\cdot, \cdot)$ : $V \times V \to \mathbb{R}$ and norm $\|v\|_V = a(v, v)^{1/2}$. □

## 4.1 The Theorems of Riesz and Lax–Milgram

**Theorem 4.2 Representation theorem of Riesz.** *Let $f \in V'$ be a continuous and linear functional, then there is a uniquely determined $u \in V$ with*

$$a(u, v) = f(v) \quad \forall \, v \in V. \tag{4.1}$$

*In addition, $u$ is the unique solution of the variational problem*

$$F(v) = \frac{1}{2}a(v, v) - f(v) \to \min \ \forall \, v \in V. \tag{4.2}$$

**Proof:** First, the existence of a solution $u$ of the variational problem will be proved. Since $f$ is continuous, it holds

$$|f(v)| \le c \, \|v\|_V \quad \forall \, v \in V,$$

from what follows that

$$F(v) \ge \frac{1}{2} \, \|v\|_V^2 - c \, \|v\|_V \ge -\frac{1}{2}c^2,$$

where in the last estimate the necessary criterion for a local minimum of the expression of the first estimate is used. Hence, the function $F(\cdot)$ is bounded from below and

$$d = \inf_{v \in V} F(v)$$

exists.

Let $\{v_k\}_{k \in \mathbb{N}}$ be a sequence with $F(v_k) \to d$ for $k \to \infty$. A straightforward calculation (parallelogram identity in Hilbert spaces) gives

$$\|v_k - v_l\|_V^2 + \|v_k + v_l\|_V^2 = 2 \, \|v_k\|_V^2 + 2 \, \|v_l\|_V^2 \,.$$

Using the linearity of $f(\cdot)$ and $d \leq F(v)$ for all $v \in V$, one obtains

$$\|v_k - v_l\|_V^2$$

$$= \quad 2\|v_k\|_V^2 + 2\|v_l\|_V^2 - 4\left\|\frac{v_k + v_l}{2}\right\|_V^2 - 4f(v_k) - 4f(v_l) + 8f\left(\frac{v_k + v_l}{2}\right)$$

$$= \quad 4F(v_k) + 4F(v_l) - 8F\left(\frac{v_k + v_l}{2}\right)$$

$$\leq \quad 4F(v_k) + 4F(v_l) - 8d \to 0$$

for $k, l \to \infty$. Hence $\{v_k\}_{k\in\mathbb{N}}$ is a Cauchy sequence. Because $V$ is a complete space, there exists a limit $u$ of this sequence with $u \in V$. Because $F(\cdot)$ is continuous, it is $F(u) = d$ and $u$ is a solution of the variational problem.

In the next step, it will be shown that each solution of the variational problem (4.2) is also a solution of (4.1). It is

$$\Phi(\varepsilon) \quad = \quad F(u + \varepsilon v) = \frac{1}{2}a(u + \varepsilon v, u + \varepsilon v) - f(u + \varepsilon v)$$

$$= \quad \frac{1}{2}a(u, u) + \varepsilon a(u, v) + \frac{\varepsilon^2}{2}a(v, v) - f(u) - \varepsilon f(v).$$

If $u$ is a minimum of the variational problem, then the function $\Phi(\varepsilon)$ has a local minimum at $\varepsilon = 0$. The necessary condition for a local minimum leads to

$$0 = \Phi'(0) = a(u, v) - f(v) \quad \text{for all } v \in V.$$

Finally, the uniqueness of the solution will be proved. It is sufficient to prove the uniqueness of the solution of the equation (4.1). If the solution of (4.1) is unique, then the existence of two solutions of the variational problem (4.2) would be a contradiction to the fact proved in the previous step. Let $u_1$ and $u_2$ be two solutions of the equation (4.1). Computing the difference of both equations gives

$$a(u_1 - u_2, v) = 0 \quad \text{for all } v \in V.$$

This equation holds, in particular, for $v = u_1 - u_2$. Hence, $\|u_1 - u_2\|_V = 0$, such that $u_1 = u_2$. ∎

**Definition 4.3 Bounded bilinear form, coercive bilinear form, $V$-elliptic bilinear form.** Let $b(\cdot, \cdot) : V \times V \to \mathbb{R}$ be a bilinear form on the Banach space $V$. Then it is bounded if

$$|b(u, v)| \leq M \|u\|_V \|v\|_V \quad \forall\, u, v \in V, M > 0, \tag{4.3}$$

where the constant $M$ is independent of $u$ and $v$. The bilinear form is coercive or $V$-elliptic if

$$b(u, u) \geq m \|u\|_V^2 \quad \forall\, u \in V, m > 0, \tag{4.4}$$

where the constant $m$ is independent of $u$. □

**Remark 4.4** *Application to an inner product.* Let $V$ be a Hilbert space. Then the inner product $a(\cdot, \cdot)$ is a bounded and coercive bilinear form, since by the Cauchy–Schwarz inequality

$$|a(u, v)| \leq \|u\|_V \|v\|_V \quad \forall\, u, v \in V,$$

and obviously $a(u, u) = \|u\|_V^2$. Hence, the constants can be chosen to be $M = 1$ and $m = 1$.

Next, the representation theorem of Riesz will be generalized to the case of coercive and bounded bilinear forms. □

**Theorem 4.5 Theorem of Lax–Milgram.** *Let $b(\cdot, \cdot) : V \times V \to \mathbb{R}$ be a bounded and coercive bilinear form on the Hilbert space $V$. Then, for each bounded linear functional $f \in V'$ there is exactly one $u \in V$ with*

$$b(u, v) = f(v) \quad \forall\, v \in V. \tag{4.5}$$

**Proof:** One defines linear operators $T, T' : V \to V$ by

$$a(Tu, v) = b(u, v) \ \forall \ v \in V, \quad a(T'u, v) = b(v, u) \ \forall \ v \in V. \tag{4.6}$$

Since $b(u, \cdot)$ and $b(\cdot, u)$ are continuous linear functionals on $V$, it follows from Theorem 4.2 that the elements $Tu$ and $T'u$ exist and they are defined uniquely. Because the operators satisfy the relation

$$a(Tu, v) = b(u, v) = a(T'v, u) = a(u, T'v), \tag{4.7}$$

$T'$ is called adjoint operator of $T$. Setting $v = Tu$ in (4.6) and using the boundedness of $b(\cdot, \cdot)$ yields

$$\|Tu\|_V^2 = a(Tu, Tu) = b(u, Tu) \le M \|u\|_V \|Tu\|_V \implies \|Tu\|_V \le M \|u\|_V$$

for all $u \in V$. Hence, $T$ is bounded. Since $T$ is linear, it follows that $T$ is continuous. Using the same argument, one shows that $T'$ is also bounded and continuous.

Define the bilinear form

$$d(u, v) := a(TT'u, v) = a(T'u, T'v) \quad \forall \ u, v \in V,$$

where (4.7) was used. Hence, this bilinear form is symmetric. Using the coercivity of $b(\cdot, \cdot)$ and the Cauchy–Schwarz inequality gives

$$m^2 \|v\|_V^4 \le b(v, v)^2 = a(T'v, v)^2 \le \|v\|_V^2 \|T'v\|_V^2 = \|v\|_V^2 \, a(T'v, T'v) = \|v\|_V^2 \, d(v, v).$$

Applying now the boundedness of $a(\cdot, \cdot)$ and of $T'$ yields

$$m^2 \|v\|_V^2 \le d(v, v) = a(T'v, T'v) = \|T'v\|_V^2 \le M \|v\|_V^2. \tag{4.8}$$

Hence, $d(\cdot, \cdot)$ is also coercive and, since it is symmetric, it defines an inner product on $V$. From (4.8) one has that the norm induced by $d(v, v)^{1/2}$ is equivalent to the norm $\|v\|_V$. From Theorem 4.2 it follows that there is a exactly one $w \in V$ with

$$d(w, v) = f(v) \quad \forall \ v \in V.$$

Inserting $u = T'w$ into (4.5) gives with (4.6)

$$b(T'w, v) = a(TT'w, v) = d(w, v) = f(v) \quad \forall \ v \in V,$$

hence $u = T'w$ is a solution of (4.5).

The uniqueness of the solution is proved analogously as in the symmetric case. ∎

# 4.2 Weak Formulation of Boundary Value Problems

**Remark 4.6** *Model problem.* Consider the Poisson equation with homogeneous Dirichlet boundary conditions

$$\begin{aligned} -\Delta u &= f &\text{in } \Omega \subset \mathbb{R}^d, \\ u &= 0 &\text{on } \partial\Omega. \end{aligned} \tag{4.9}$$

□

**Definition 4.7 Weak formulation of** (4.9)**.** Let $f \in L^2(\Omega)$. A weak formulation of (4.9) consists in finding $u \in V = H_0^1(\Omega)$ such that

$$a(u, v) = (f, v) \quad \forall \ v \in V \tag{4.10}$$

with

$$a(u, v) = (\nabla u, \nabla v) = \int_\Omega \nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \ d\mathbf{x}$$

and $(\cdot, \cdot)$ is the inner product in $L^2(\Omega)$.

□

**Remark 4.8** *On the weak formulation.*
- The weak formulation is also called variational formulation.
- As usual in mathematics, 'weak' means that something holds for all appropriately chosen test functions.
- Formally, one obtains the weak formulation by multiplying the strong form of the equation (4.9) with the test function, by integrating the equation on $\Omega$, and applying integration by parts. Because of the Dirichlet boundary condition, on can use as test space $H_0^1(\Omega)$ and therefore the integral on the boundary vanishes.
- The ansatz space for the solution and the test space are defined such that the arising integrals are well defined.
- The weak formulation reduces the necessary regularity assumptions for the solution by the integration and the transfer of derivatives to the test function. Whereas the solution of (4.9) has to be in $C^2(\overline{\Omega})$, the solution of (4.10) has to be only in $H_0^1(\Omega)$. The latter assumption is much more realistic for problems coming from applications.
- The regularity assumption on the right hand side can be relaxed to $f \in H^{-1}(\Omega)$.
$\square$

**Theorem 4.9 Existence and uniqueness of the weak solution.** *Let* $f \in L^2(\Omega)$. *There is exactly one solution of* (4.10).

**Proof:** Because of the Poincaré inequality (3.9), there is a constant $c$ with

$$\|v\|_{L^2(\Omega)} \le c \, \|\nabla v\|_{L^2(\Omega)} \quad \forall \, v \in H_0^1(\Omega).$$

It follows for $v \in H_0^1(\Omega) \subset H^1(\Omega)$ that

$$
\begin{aligned}
\|v\|_{H^1(\Omega)} &= \left( \|v\|_{L^2(\Omega)}^2 + \|\nabla v\|_{L^2(\Omega)}^2 \right)^{1/2} \le \left( c \, \|\nabla v\|_{L^2(\Omega)}^2 + \|\nabla v\|_{L^2(\Omega)}^2 \right)^{1/2} \\
&\le C \, \|\nabla v\|_{L^2(\Omega)} \le C \, \|v\|_{H^1(\Omega)}.
\end{aligned}
$$

Hence, $a(\cdot, \cdot)$ is an inner product on $H_0^1(\Omega)$ with the induced norm

$$\|v\|_{H_0^1(\Omega)} = a(v, v)^{1/2},$$

which is equivalent to the norm $\|\cdot\|_{H^1(\Omega)}$.

Define for $f \in L^2(\Omega)$ the linear functional

$$\tilde{f}(v) := \int_\Omega f(\mathbf{x}) v(\mathbf{x}) \, d\mathbf{x} \quad \forall \, v \in H_0^1(\Omega).$$

Applying the Cauchy–Schwarz inequality (3.5) and the Poincaré inequality (3.9)

$$\left| \tilde{f}(v) \right| = |(f, v)| \le \|f\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \le c \, \|f\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)} = c \, \|f\|_{L^2(\Omega)} \|v\|_{H_0^1(\Omega)}$$

shows that this functional is continuous on $H_0^1(\Omega)$. Applying the representation theorem of Riesz, Theorem 4.2, gives the existence and uniqueness of the weak solution of (4.10). In addition, $u(\mathbf{x})$ solves the variational problem

$$F(v) = \frac{1}{2} \|\nabla v\|_2^2 - \int_\Omega f(\mathbf{x}) v(\mathbf{x}) \, d\mathbf{x} \to \min \quad \text{for all } v \in H_0^1(\Omega).$$

$\blacksquare$

**Example 4.10** *A more general elliptic problem.* Consider the problem

$$
\begin{aligned}
-\nabla \cdot (A(\mathbf{x}) \nabla u) + c(\mathbf{x}) u &= f \quad \text{in } \Omega \subset \mathbb{R}^d, \\
u &= 0 \quad \text{on } \partial\Omega,
\end{aligned}
\tag{4.11}
$$

with $A(\mathbf{x}) \in \mathbb{R}^{d \times d}$ for each point $\mathbf{x} \in \Omega$. It will be assumed that the coefficients $a_{i,j}(\mathbf{x})$ and $c(\mathbf{x}) \geq 0$ are bounded, $f \in L^2(\Omega)$, and that the matrix (tensor) $A(\mathbf{x})$ is for all $\mathbf{x} \in \Omega$ uniformly elliptic, i.e., there are positive constants $m$ and $M$ such that

$$m \|\mathbf{y}\|_2^2 \leq \mathbf{y}^T A(\mathbf{x})\mathbf{y} \leq M \|\mathbf{y}\|_2^2 \quad \forall\, \mathbf{y} \in \mathbb{R}^d,\ \forall\, \mathbf{x} \in \Omega.$$

The weak form of (4.11) is obtained in the usual way by multiplying (4.11) with test functions $v \in H_0^1(\Omega)$, integrating on $\Omega$, and applying integration by parts: Find $u \in H_0^1(\Omega)$, such that

$$a(u,v) = f(v) \quad \forall\, v \in H_0^1(\Omega)$$

with

$$a(u,v) = \int_\Omega \left( \nabla u(\mathbf{x})^T A(\mathbf{x}) \nabla v(\mathbf{x}) + c(\mathbf{x})u(\mathbf{x})v(\mathbf{x}) \right)\, d\mathbf{x}.$$

This bilinear form is bounded (*exercise*). The coercivity of the bilinear form is proved by using the uniform ellipticity of $A(\mathbf{x})$ and the non-negativity of $c(\mathbf{x})$:

$$
\begin{aligned}
a(u,u) &= \int_\Omega \nabla u(\mathbf{x})^T A(\mathbf{x}) \nabla u(\mathbf{x}) + c(\mathbf{x})u(\mathbf{x})u(\mathbf{x})\, d\mathbf{x} \\
&\geq \int_\Omega m \nabla u(\mathbf{x})^T \nabla u(\mathbf{x})\, d\mathbf{x} = m \|u\|_{H_0^1(\Omega)}^2.
\end{aligned}
$$

Applying the Theorem of Lax–Milgram, Theorem 4.5, gives the existence and uniqueness of a weak solution of (4.11).

If the tensor is not symmetric, $a_{ij}(\mathbf{x}) \neq a_{ji}(\mathbf{x})$ for one pair $i, j$, then the solution cannot be characterized as the solution of a variational problem. $\qquad\square$

## 4.3 The Ritz Method and the Galerkin Method

**Remark 4.11** *Idea of the Ritz method.* Let $V$ be a Hilbert space with the inner product $a(\cdot, \cdot)$. Consider the problem

$$F(v) = \frac{1}{2}a(v,v) - f(v) \rightarrow \min, \tag{4.12}$$

where $f : V \to \mathbb{R}$ is a bounded linear functional. As already proved in Theorem 4.2, there is a unique solution $u \in V$ of this variational problem which is also the unique solution of the equation

$$a(u,v) = f(v) \quad \forall\, v \in V. \tag{4.13}$$

For approximating the solution of (4.12) or (4.13) with a numerical method, it will be assumed that $V$ has a countable orthonormal basis (Schauder basis). Then, there are finite-dimensional subspaces $V_1, V_2, \ldots \subset V$ with $\dim V_k = k$, which has the following property: for each $u \in V$ and each $\varepsilon > 0$ there is a $K \in \mathbb{N}$ and a $u_k \in V_k$ with

$$\|u - u_k\|_V \leq \varepsilon \quad \forall\, k \geq K. \tag{4.14}$$

Note that it is not required that there holds an inclusion of the form $V_k \subset V_{k+1}$.

The Ritz approximation of (4.12) and (4.13) is defined by: Find $u_k \in V_k$ with

$$a(u_k, v_k) = f(v_k) \quad \forall\, v_k \in V_k. \tag{4.15}$$

$\qquad\square$

**Lemma 4.12 Existence and uniqueness of a solution of** (4.15). *There exists exactly one solution of* (4.15).

**Proof:** Finite-dimensional subspaces of Hilbert spaces are Hilbert spaces as well. For this reason, one can apply the representation theorem of Riesz, Theorem 4.2, to (4.15) which gives the statement of the lemma. In addition, the solution of (4.15) solves a minimization problem on $V_k$. ∎

**Lemma 4.13 Best approximation property.** *The solution of* (4.15) *is the best approximation of $u$ in $V_k$, i.e., it is*

$$\|u - u_k\|_V = \inf_{v_k \in V_k} \|u - v_k\|_V .\tag{4.16}$$

**Proof:** Since $V_k \subset V$, one can use the test functions from $V_k$ in the weak equation (4.13). Then, the difference of (4.13) and (4.15) gives the orthogonality, the so-called Galerkin orthogonality,

$$a(u - u_k, v_k) = 0 \quad \forall \, v_k \in V_k.\tag{4.17}$$

Hence, the error $u - u_k$ is orthogonal to the space $V_k$: $u - u_k \perp V_k$. That means, $u_k$ is the orthogonal projection of $u$ onto $V_k$ with respect of the inner product of $V$.

Let now $w_k \in V_k$ be an arbitrary element, then it follows with the Galerkin orthogonality (4.17) and the Cauchy–Schwarz inequality that

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

Since $w_k \in V_k$ was arbitrary, also $v_k \in V_k$ is arbitrary. If $\|u - u_k\|_V > 0$, division by $\|u - u_k\|_V$ gives the statement of the lemma. If $\|u - u_k\|_V = 0$, the statement of the lemma is trivially true. ∎

**Theorem 4.14 Convergence of the Ritz approximation.** *The Ritz approximation converges*

$$\lim_{k \to \infty} \|u - u_k\|_V = 0.$$

**Proof:** The best approximation property (4.16) and property (4.14) give

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

for each $\varepsilon > 0$ and $k \geq K(\varepsilon)$. Hence, the convergence is proved. ∎

**Remark 4.15** *Formulation of the Ritz method as linear system of equations.* One can use an arbitrary basis $\{\phi_i\}_{i=1}^k$ of $V_k$ for the computation of $u_k$. First of all, the equation for the Ritz approximation (4.15) is satisfied for all $v_k \in V_k$ if and only if it is satisfied for each basis function $\phi_i$. This statement follows from the linearity of both sides of the equation with respect to the test function and from the fact that each function $v_k \in V_k$ can be represented as linear combination of the basis functions. Let $v_k = \sum_{i=i}^k \alpha_i \phi_i$, then from (4.15) it follows that

$$a(u_k, v_k) = \sum_{k=1}^k \alpha_i a(u_k, \phi_i) = \sum_{k=1}^k \alpha_i f(\phi_i) = f(v_k).$$

This equation is satisfied if $a(u_k, \phi_i) = f(\phi_i)$, $i = 1, \ldots, k$. On the other hand, if (4.15) holds then it holds in particular for each basis function $\phi_i$.

Then, one uses as ansatz for the solution also a linear combination of the basis functions

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

with unknown coefficients $u^j \in \mathbb{R}$. Using as test functions now the basis functions yields

$$\sum_{j=1}^{k} a(u^j \phi_j, \phi_i) = \sum_{j=1}^{k} a(\phi_j, \phi_i) u^j = f(\phi_i), \quad i = 1, \ldots, k.$$

This equation is equivalent to the linear system of equations $A\mathbf{u} = \mathbf{f}$, where

$$A = (a_{ij})_{i,j=1}^{k} = a(\phi_j, \phi_i)_{i,j=1}^{k}$$

is called stiffness matrix. Note that the order of the indices is different for the entries of the matrix and the arguments of the inner product. The right hand side is a vector of length $k$ with the entries $f_i = f(\phi_i)$, $i = 1, \ldots, k$.

Using the one-to-one mapping between the coefficient vector $(v^1, \ldots, v^k)^T$ and the element $v_k = \sum_{i=1}^{k} v^i \phi_i$, one can show that the matrix $A$ is symmetric and positive definite (*exercise*)

$$A = A^T \iff a(v, w) = a(w, v) \quad \forall\, v, w \in V_k,$$
$$\mathbf{x}^T A \mathbf{x} > 0 \text{ for } \mathbf{x} \neq \mathbf{0} \iff a(v, v) > 0 \quad \forall\, v \in V_k, v \neq 0.$$

$\square$

**Remark 4.16** *The case of a bounded and coercive bilinear form.* If $b(\cdot, \cdot)$ is bounded and coercive, but not symmetric, it is possible to approximate the solution of (4.5) with the same idea as for the Ritz method. In this case, it is called Galerkin method. The discrete problem consists in finding $u_k \in V_k$ such that

$$b(u_k, v_k) = f(v_k) \quad \forall\, v_k \in V_k. \tag{4.18}$$

$\square$

**Lemma 4.17 Existence and uniqueness of a solution of** (4.18)**.** *There is exactly one solution of* (4.18)*.*

**Proof:** The statement of the lemma follows directly from the Theorem of Lax–Milgram, Theorem 4.5. ∎

**Remark 4.18** *On the discrete solution.* The discrete solution is not the orthogonal projection into $V_k$ in the case of a bounded and coercive bilinear form, which is not the inner product of $V$. $\square$

**Lemma 4.19 Lemma of Cea, error estimate.** *Let $b : V \times V \to \mathbb{R}$ be a bounded and coercive bilinear form on the Hilbert space $V$ and let $f \in V'$ be a bounded linear functional. Let $u$ be the solution of (4.5) and $u_k$ be the solution of (4.18), then the following error estimate holds*

$$\|u - u_k\|_V \leq \frac{M}{m} \inf_{v_k \in V_k} \|u - v_k\|_V, \tag{4.19}$$

*where the constants $M$ and $m$ are given in (4.3) and (4.4).*

**Proof:** Considering the difference of the continuous equation (4.5) and the discrete equation (4.18), one obtains the error equation

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

which is also called Galerkin orthogonality. With (4.4), the Galerkin orthogonality, and (4.3) it follows that

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

from what the statement of the lemma follows immediately. ∎

**Remark 4.20** *On the best approximation error.* It follows from estimate (4.19) that the error is bounded by a multiple of the best approximation error, where the factor depends on properties of the bilinear form $b(\cdot, \cdot)$. Thus, concerning error estimates for concrete finite-dimensional spaces, the study of the best approximation error will be of importance. □

**Remark 4.21** *The corresponding linear system of equations.* The corresponding linear system of equations is derived analogously to the symmetric case. The system matrix is still positive definite but not symmetric. □

**Remark 4.22** *Choice of the basis.* The most important issue of the Ritz and Galerkin method is the choice of the spaces $V_k$, or more concretely, the choice of an appropriate basis $\{\phi_i\}_{i=1}^{k}$ that spans the space $V_k$. From the point of view of numerics, there are the requirements that it should be possible to compute the entries $a_{ij}$ of the stiffness matrix efficiently and that the matrix $A$ should be sparse. □

# Chapter 5

# Finite Element Methods

## 5.1 Finite Element Spaces

**Remark 5.1** *Mesh cells, faces, edges, vertices.* A mesh cell $K$ is a compact poly-hedron in $\mathbb{R}^d$, $d \in \{2,3\}$, whose interior is not empty. The boundary $\partial K$ of $K$ consists of $m$-dimensional linear manifolds (points, pieces of straight lines, pieces of planes), $0 \leq m \leq d-1$, which are called $m$-faces. The 0-faces are the vertices of the mesh cell, the 1-faces are the edges, and the $(d-1)$-faces are just called faces.
□

**Remark 5.2** *Finite dimensional spaces defined on $K$.* Let $s \in \mathbb{N}$. Finite element methods use finite dimensional spaces $P(K) \subset C^s(K)$ which are defined on $K$. In general, $P(K)$ consists of polynomials. The dimension of $P(K)$ will be denoted by $\dim P(K) = N_K$.
□

**Example 5.3** *The space $P(K) = P_1(K)$.* The space consisting of linear polynomials on a mesh cell $K$ is denoted by $P_1(K)$:

$$P_1(K) = \left\{ a_0 + \sum_{i=1}^{d} a_i x_i \; : \; \mathbf{x} = (x_1, \ldots, x_d)^T \in K \right\}.$$

There are $d+1$ unknown coefficients $a_i$, $i = 0, \ldots, d$, such that $\dim P_1(K) = N_K = d+1$.
□

**Remark 5.4** *Linear functionals defined on $P(K)$.* For the definition of finite elements, linear functional which are defined on $P(K)$ are of importance.

Consider linear and continuous functionals $\Phi_{K,1}, \ldots, \Phi_{K,N_K} \; : \; C^s(K) \to \mathbb{R}$ which are linearly independent. There are different types of functionals which can be utilized in finite element methods:

- point values: $\Phi(v) = v(\mathbf{x})$, $\mathbf{x} \in K$,
- point values of a first partial derivative: $\Phi(v) = \partial_i v(\mathbf{x})$, $\mathbf{x} \in K$,
- point values of the normal derivative on a face $E$ of $K$: $\Phi(v) = \nabla v(\mathbf{x}) \cdot \mathbf{n}_E$, $\mathbf{n}_E$ is the outward pointing unit normal vector on $E$,
- integral mean values on $K$: $\Phi(v) = \frac{1}{|K|} \int_K v(\mathbf{x}) \, d\mathbf{x}$,
- integral mean values on faces $E$: $\Phi(v) = \frac{1}{|E|} \int_E v(\mathbf{s}) \, d\mathbf{s}$.

The smoothness parameter $s$ has to be chosen in such a way that the functionals $\Phi_{K,1}, \ldots, \Phi_{K,N_K}$ are continuous. If, e.g., a functional requires the evaluation of a partial derivative or a normal derivative, then one has to choose at least $s = 1$. For the other functionals given above, $s = 0$ is sufficient.
□

**Definition 5.5 Unisolvence of $P(K)$ with respect to the functionals $\Phi_{K,1}$,**
$\ldots, \Phi_{K,N_K}$. The space $P(K)$ is called unisolvent with respect to the functionals $\Phi_{K,1}, \ldots, \Phi_{K,N_K}$ if there is for each $\mathbf{a} \in \mathbb{R}^{N_K}$, $\mathbf{a} = (a_1, \ldots, a_{N_K})^T$, exactly one $p \in P(K)$ with

$$\Phi_{K,i}(p) = a_i, \quad 1 \le i \le N_K.$$

$\square$

**Remark 5.6** *Local basis.* Unisolvence means that for each vector $\mathbf{a} \in \mathbb{R}^{N_K}$, $\mathbf{a} = (a_1, \ldots, a_{N_K})^T$, there is exactly one element in $P(K)$ such that $a_i$ is the image of the $i$-th functional, $i = 1, \ldots, N_K$.

Choosing in particular the Cartesian unit vectors for $\mathbf{a}$, then it follows from the unisolvence that a set $\{\phi_{K,i}\}_{i=1}^{N_K}$ exists with $\phi_{K,i} \in P(K)$ and

$$\Phi_{K,i}(\phi_{K,j}) = \delta_{ij}, \quad i, j = 1, \ldots, N_K.$$

Consequently, the set $\{\phi_{K,i}\}_{i=1}^{N_K}$ forms a basis of $P(K)$. This basis is called local basis. $\square$

**Remark 5.7** *Transform of an arbitrary basis to the local basis.* If an arbitrary basis $\{p_i\}_{i=1}^{N_K}$ of $P(K)$ is known, then the local basis can be computed by solving a linear system of equations. To this end, represent the local basis in terms of the known basis

$$\phi_{K,j} = \sum_{k=1}^{N_K} c_{jk} p_k, \quad c_{jk} \in \mathbb{R}, \; j = 1, \ldots, N_K,$$

with unknown coefficients $c_{jk}$. Applying the definition of the local basis leads to the linear system of equations

$$\Phi_{K,i}(\phi_{K,j}) = \sum_{k=1}^{N_K} c_{jk} a_{ik} = \delta_{ij}, \quad i, j = 1, \ldots, N_K, \quad a_{ik} = \Phi_{K,i}(p_k).$$

Because of the unisolvence, the matrix $A = (a_{ij})$ is non-singular and the coefficients $c_{jk}$ are determined uniquely. $\square$

**Example 5.8** *Local basis for the space of linear functions on the reference triangle.*
Consider the reference triangle $\hat{K}$ with the vertices $(0,0)$, $(1,0)$, and $(0,1)$. A linear space on $\hat{K}$ is spanned by the functions $1, \hat{x}, \hat{y}$. Let the functionals be defined by the values of the functions in the vertices of the reference triangle. Then, the given basis is not a local basis because the function 1 does not vanish at the vertices.

Consider first the vertex $(0,0)$. A linear basis function $a\hat{x} + b\hat{y} + c$ which has the value 1 in $(0,0)$ and which vanishes in the other vertices has to satisfy the following set of equations

$$\begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}.$$

The solution is $a = -1, b = -1, c = 1$. The two other basis functions of the local basis are $\hat{x}$ and $\hat{y}$, such that the local basis has the form $\{1 - \hat{x} - \hat{y}, \hat{x}, \hat{y}\}$. $\square$

**Remark 5.9** *Triangulation, grid, mesh, grid cell.* For the definition of global finite element spaces, a decomposition of the domain $\Omega$ into polyhedrons $K$ is needed. This decomposition is called triangulation $\mathcal{T}^h$ and the polyhedrons $K$ are called mesh cells. The union of the polyhedrons is called grid or mesh.

A triangulation is called regular, see the definition in Ciarlet Ciarlet (1978), if:
- It holds $\overline{\Omega} = \cup_{K \in \mathcal{T}^h} K$.

- Each mesh cell $K \in \mathcal{T}^h$ is closed and the interior $\mathring{K}$ is non-empty.
- For distinct mesh cells $K_1$ and $K_2$ there holds $\mathring{K}_1 \cap \mathring{K}_2 = \emptyset$.
- For each $K \in \mathcal{T}^h$, the boundary $\partial K$ is Lipschitz-continuous.
- The intersection of two mesh cells is either empty or a common $m$-face, $m \in \{0, \ldots, d-1\}$.

$\square$

**Remark 5.10** *Global and local functionals.* Let $\Phi_1, \ldots, \Phi_N : C^s(\overline{\Omega}) \to \mathbb{R}$ continuous linear functionals of the same types as given in Remark 5.4. The restriction of the functionals to $C^s(K)$ defines local functionals $\Phi_{K,1}, \ldots, \Phi_{K,N_K}$, where it is assumed that the local functionals are unisolvent on $P(K)$. The union of all mesh cells $K_j$, for which there is a $p \in P(K_j)$ with $\Phi_i(p) \neq 0$, will be denoted by $\omega_i$. $\square$

**Example 5.11** *On subdomains $\omega_i$.* Consider the two-dimensional case and let $\Phi_i$ be defined as nodal value of a function in $\mathbf{x} \in K$. If $\mathbf{x} \in \mathring{K}$, then $\omega_i = K$. In the case that $\mathbf{x}$ is on a face of $K$ but not in a vertex, then $\omega_i$ is the union of $K$ and the other mesh cell whose boundary contains this face. Last, if $\mathbf{x}$ is a vertex of $K$, then $\omega_i$ is the union of all mesh cells which possess this vertex, see Figure 5.1. $\square$
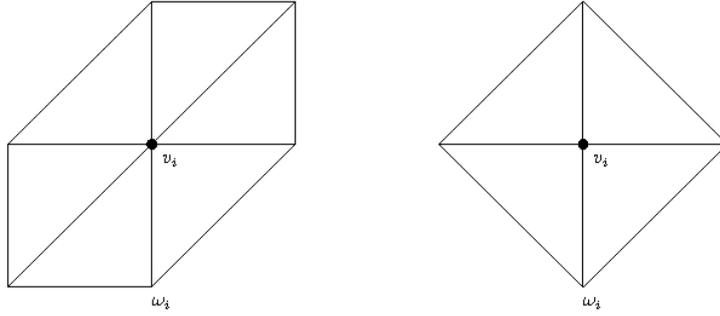


Figure 5.1: Subdomains $\omega_i$.

**Definition 5.12 Finite element space, global basis.** A function $v(\mathbf{x})$ defined on $\Omega$ with $v|_K \in P(K)$ for all $K \in \mathcal{T}^h$ is called continuous with respect to the functional $\Phi_i : \Omega \to \mathbb{R}$ if

$$\Phi_i(v|_{K_1}) = \Phi_i(v|_{K_2}), \quad \forall K_1, K_2 \in \omega_i.$$

The space

$$S = \left\{ v \in L^\infty(\Omega) : v|_K \in P(K) \text{ and } v \text{ is continuous with respect to} \right.$$
$$\left. \Phi_i, i = 1, \ldots, N \right\}$$

is called finite element space.

The global basis $\{\phi_j\}_{j=1}^N$ of $S$ is defined by the condition

$$\phi_j \in S, \quad \Phi_i(\phi_j) = \delta_{ij}, \quad i, j = 1, \ldots, N.$$

$\square$

**Example 5.13** *Piecewise linear global basis function.* Figure 5.2 shows a piecewise linear global basis function in two dimensions. Because of its form, such a function is called hat function. $\square$
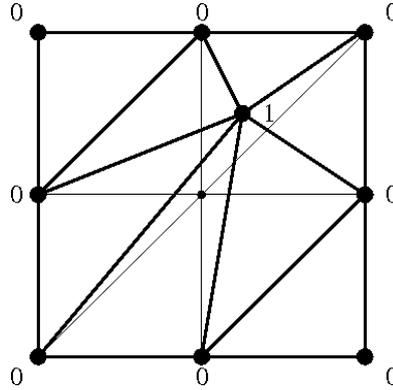
Figure 5.2: Piecewise linear global basis function (boldface lines), hat function.

**Remark 5.14** *On global basis functions.* A global basis function coincides on each mesh cell with a local basis function. This property implies the uniqueness of the global basis functions.

For many finite element spaces it follows from the continuity with respect to $\{\Phi_i\}_{i=1}^N$, the continuity of the finite element functions themselves. Only in this case, one can speak of values of finite element functions on $m$-faces with $m < d$. □

**Definition 5.15 Parametric finite elements.** Let $\hat{K}$ be a reference mesh cell with the local space $P(\hat{K})$, the local functionals $\hat{\Phi}_1, \ldots, \hat{\Phi}_{\hat{N}}$, and a class of bijective mappings $\{F_K \; : \; \hat{K} \to K\}$. A finite element space is called a parametric finite element space if:
- The images $\{K\}$ of $\{F_K\}$ form the set of mesh cells.
- The local spaces are given by

$$P(K) = \left\{ p \; : \; p = \hat{p} \circ F_K^{-1}, \hat{p} \in \hat{P}(\hat{K}) \right\}. \tag{5.1}$$

- The local functionals are defined by

$$\Phi_{K,i}(v(\mathbf{x})) = \hat{\Phi}_i \left( v(F_K(\hat{\mathbf{x}})) \right), \tag{5.2}$$

where $\hat{\mathbf{x}} = (\hat{x}_1, \ldots, \hat{x}_d)^T$ are the coordinates of the reference mesh cell and it holds $\mathbf{x} = F_K(\hat{\mathbf{x}})$.

□

**Remark 5.16** *Motivations for using parametric finite elements.* Definition 5.12 of finite elements spaces is very general. For instance, different types of mesh cells are allowed. However, as well the finite element theory as the implementation of finite element methods become much simpler if only parametric finite elements are considered. □

## 5.2 Finite Elements on Simplices

**Definition 5.17** *$d$-simplex.* **A $d$-simplex $K \subset \mathbb{R}^d$ is the convex hull of $(d+1)$ points $\mathbf{a}_1, \ldots, \mathbf{a}_{d+1} \in \mathbb{R}^d$ which form the vertices of $K$.** □

**Remark 5.18** *On $d$-simplices.* It will be always assumed that the simplex is not degenerated, i.e., its $d$-dimensional measure is positive. This property is equivalent

to the non-singularity of the matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1,d+1} \\ a_{21} & a_{22} & \dots & a_{2,d+1} \\ \vdots & \vdots & \ddots & \vdots \\ a_{d1} & a_{d2} & \dots & a_{d,d+1} \\ 1 & 1 & \dots & 1 \end{pmatrix},$$

where $\mathbf{a}_i = (a_{1i}, a_{2i}, \dots, a_{di})^T$, $i = 1, \dots, d+1$.

For $d = 2$, the simplices are the triangles and for $d = 3$ they are the tetrahedrons.
$\square$

**Definition 5.19 Barycentric coordinates.** Since $K$ is the convex hull of the points $\{\mathbf{a}_i\}_{i=1}^{d+1}$, the parametrization of $K$ with a convex combination of the vertices reads as follows

$$K = \left\{ \mathbf{x} \in \mathbb{R}^d \ : \ \mathbf{x} = \sum_{i=1}^{d+1} \lambda_i \mathbf{a}_i, \ 0 \le \lambda_i \le 1, \ \sum_{i=1}^{d+1} \lambda_i = 1 \right\}.$$

The coefficients $\lambda_1, \dots, \lambda_{d+1}$ are called barycentric coordinates of $\mathbf{x} \in K$. $\square$

**Remark 5.20** *On barycentric coordinates.* From the definition it follows that the barycentric coordinates are the solution of the linear system of equations

$$\sum_{i=1}^{d+1} a_{ji} \lambda_i = x_j, \quad 1 \le j \le d, \quad \sum_{i=1}^{d+1} \lambda_i = 1.$$

Since the system matrix is non-singular, see Remark 5.18, the barycentric coordinates are determined uniquely.

The barycentric coordinates of the vertex $\mathbf{a}_i$, $i = 1, \dots, d+1$, of the simplex is $\lambda_i = 1$ and $\lambda_j = 0$ if $i \ne j$. Since $\lambda_i(\mathbf{a}_j) = \delta_{ij}$, the barycentric coordinate $\lambda_i$ can be identified with the linear function which has the value 1 in the vertex $\mathbf{a}_i$ and which vanishes in all other vertices $\mathbf{a}_j$ with $j \ne i$.

The barycenter of the simplex is given by

$$S_K = \frac{1}{d+1} \sum_{i=1}^{d+1} \mathbf{a}_i = \sum_{i=1}^{d+1} \frac{1}{d+1} \mathbf{a}_i.$$

Hence, its barycentric coordinates are $\lambda_i = 1/(d+1)$, $i = 1, \dots, d+1$. $\square$

**Remark 5.21** *Simplicial reference mesh cells.* A commonly used reference mesh cell for triangles and tetrahedrons is the unit simplex

$$\hat{K} = \left\{ \hat{\mathbf{x}} \in \mathbb{R}^d \ : \ \sum_{i=1}^{d} \hat{x}_i \le 1, \ \hat{x}_i \ge 0, \ i = 1, \dots, d \right\},$$

see Figure 5.3. The class $\{F_K\}$ of admissible mappings are the bijective affine mappings

$$F_K \hat{\mathbf{x}} = B\hat{\mathbf{x}} + \mathbf{b}, \quad B \in \mathbb{R}^{d \times d}, \ \det(B) \ne 0, \ \mathbf{b} \in \mathbb{R}^d.$$

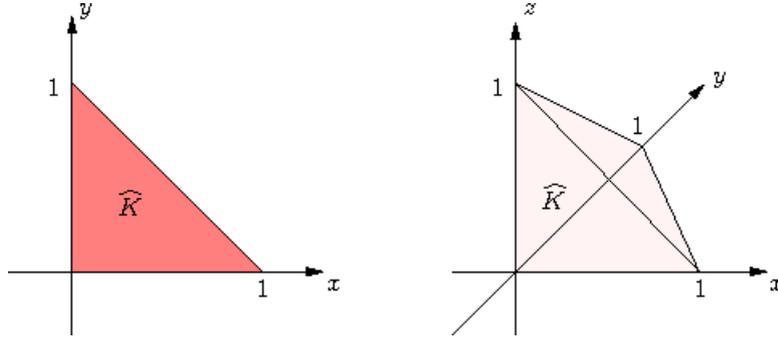The images of these mappings generate the set of the non-degenerated simplices $\{K\} \subset \mathbb{R}^d$. $\square$

Figure 5.3: The unit simplices in two and three dimensions.

**Definition 5.22 Affine family of simplicial finite elements.** Given a simplicial reference mesh cell $\hat{K}$, affine mappings $\{F_K\}$, and an unisolvent set of functionals on $\hat{K}$. Using (5.1) and (5.2), one obtains a local finite element space on each non-degenerated simplex. The set of these local spaces is called affine family of simplicial finite elements. □

**Definition 5.23 Polynomial space $P_k$.** Let $\mathbf{x} = (x_1, \ldots, x_d)^T$, $k \in \mathbb{N} \cup \{0\}$, and $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_d)^T$. Then, the polynomial space $P_k$ is given by

$$P_k = \text{span}\left\{ \prod_{i=1}^{d} x_i^{\alpha_i} = \mathbf{x}^{\boldsymbol{\alpha}} \; : \; \alpha_i \in \mathbb{N} \cup \{0\} \;\; \text{for} \;\; i = 1, \ldots, d, \; \sum_{i=1}^{d} \alpha_i \leq k \right\}.$$

□

**Remark 5.24** *Lagrangian finite elements.* In all examples given below, the linear functionals on the reference mesh cell $\hat{K}$ are the values of the polynomials with the same barycentric coordinates as on the general mesh cell $K$. Finite elements whose linear functionals are values of the polynomials on certain points in $K$ are called Lagrangian finite elements. □

**Example 5.25** $P_0$ : *piecewise constant finite element.* The piecewise constant finite element space consists of discontinuous functions. The linear functional is the value of the polynomial in the barycenter of the mesh cell, see Figure 5.4. It is $\dim P_0(K) = 1$. □
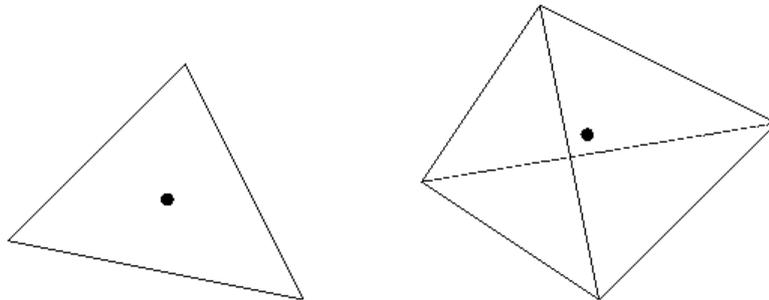


Figure 5.4: The finite element $P_0(K)$.

**Example 5.26** $P_1$ : *conforming piecewise linear finite element.* This finite element space is a subspace of $C(\overline{\Omega})$. The linear functionals are the values of the function in the vertices of the mesh cells, see Figure 5.5. It follows that $\dim P_1(K) = d + 1$.
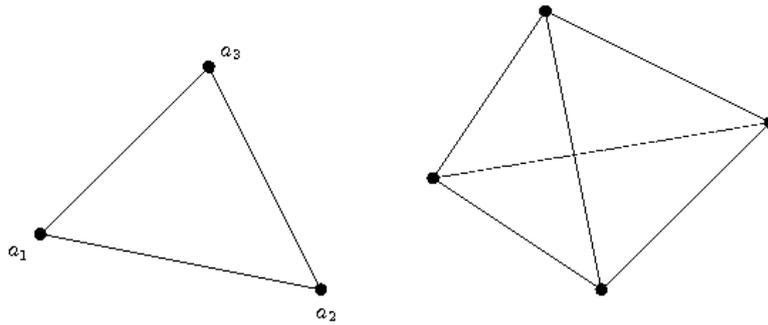
Figure 5.5: The finite element $P_1(K)$.

The local basis for the functionals $\{\Phi_i(v) = v(\mathbf{a}_i),\ i = 1, \ldots, d+1\}$, is $\{\lambda_i\}_{i=1}^{d+1}$ since $\Phi_i(\lambda_j) = \delta_{ij}$, see Remark 5.20. Since a local basis exists, the functionals are unisolvent with respect to the polynomial space $P_1(K)$.

Now, it will be shown that the corresponding finite element space consists of continuous functions. Let $K_1, K_2$ be two mesh cells with the common face $E$ and let $v \in P_1(= S)$. The restriction of $v_{K_1}$ on $E$ is a linear function on $E$ as well as the restriction of $v_{K_2}$ on $E$. It has to be shown that both linear functions are identical. A linear function on the $(d-1)$-dimensional face $E$ is uniquely determined with $d$ linearly independent functionals which are defined on $E$. These functionals can be chosen to be the values of the function in the $d$ vertices of $E$. The functionals in $S$ are continuous, by the definition of $S$. Thus, it must hold that both restrictions on $E$ have the same values in the vertices of $E$. Hence, it is $v_{K_1}|_E = v_{K_2}|_E$ and the functions from $P_1$ are continuous. □

**Example 5.27** $P_2$ : *conforming piecewise quadratic finite element.* This finite element space is also a subspace of $C(\overline{\Omega})$. It consists of piecewise quadratic functions. The functionals are the values of the functions in the $d+1$ vertices of the mesh cell and the values of the functions in the centers of the edges, see Figure 5.6. Since each vertex is connected to each other vertex, there are $\sum_{i=1}^{d} i = d(d+1)/2$ edges. Hence, it follows that $\dim P_2(K) = (d+1)(d+2)/2$.
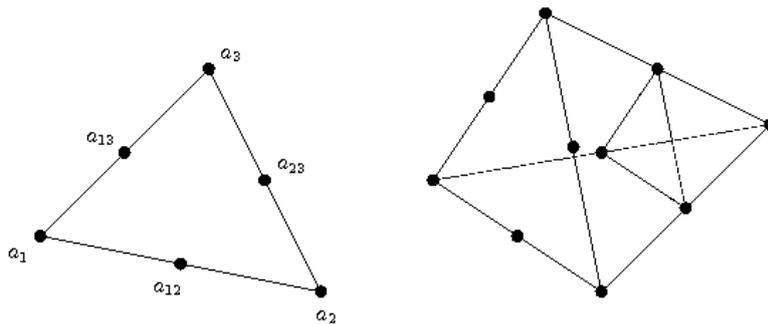


Figure 5.6: The finite element $P_2(K)$.

The part of the local basis which belongs to the functionals $\{\Phi_i(v) = v(\mathbf{a}_i),\ i = 1, \ldots, d+1\}$, is given by

$$\{\phi_i(\lambda) = \lambda_i(2\lambda_i - 1), \quad i = 1, \ldots, d+1\}.$$

Denote the center of the edges between the vertices $\mathbf{a}_i$ and $\mathbf{a}_j$ by $\mathbf{a}_{ij}$. The corre-

sponding part of the local basis is given by

$$\{\phi_{ij} = 4\lambda_i\lambda_j, \quad i,j = 1,\ldots,d+1, \ i < j\}.$$

The unisolvence follows from the fact that there exists a local basis. The continuity of the corresponding finite element space is shown in the same way as for the $P_1$ finite element. The restriction of a quadratic function in a mesh cell to a face $E$ is a quadratic function on that face. Hence, the function on $E$ is determined uniquely with $d(d+1)/2$ linearly independent functionals on $E$.

The functions $\phi_{ij}$ are called in two dimensions edge bubble functions. $\qquad\square$

**Example 5.28** $P_3$ : *conforming piecewise cubic finite element.* This finite element space consists of continuous piecewise cubic functions. It is a subspace of $C(\overline{\Omega})$. The functionals in a mesh cell $K$ are defined to be the values in the vertices $((d+1)$ values), two values on each edge (dividing the edge in three parts of equal length) $(2\sum_{i=1}^{d} i = d(d+1)$ values), and the values in the barycenter of the 2-faces of $K$, see Figure 5.7. Each 2-face of $K$ is defined by three vertices. If one considers for each vertex all possible pairs with other vertices, then each 2-face is counted three times. Hence, there are $(d+1)(d-1)d/6$ 2-faces. The dimension of $P_3(K)$ is given by

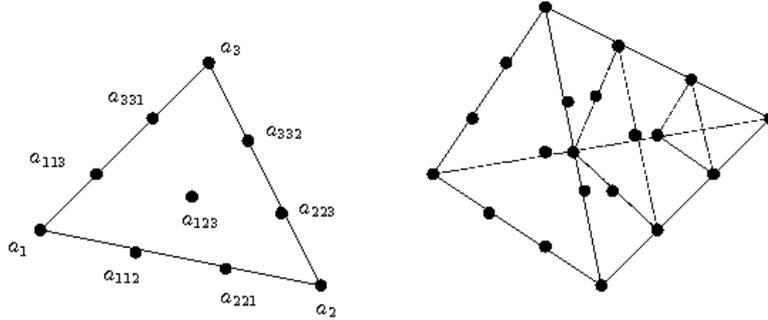$$\dim P_3(K) = (d+1) + d(d+1) + \frac{(d-1)d(d+1)}{6} = \frac{(d+1)(d+2)(d+3)}{6}.$$



Figure 5.7: The finite element $P_3(K)$.

For the functionals

$$\left\{
\begin{array}{rcll}
\Phi_i(v) & = & v(\mathbf{a}_i), \ i = 1,\ldots,d+1, & \text{(vertex)}, \\
\Phi_{iij}(v) & = & v(\mathbf{a}_{iij}), \ i,j = 1,\ldots,d+1, i \neq j, & \text{(point on edge)}, \\
\Phi_{ijk}(v) & = & v(\mathbf{a}_{ijk}), \ i = 1,\ldots,d+1, i < j < k & \text{(point on 2-face)}
\end{array}
\right\},$$

the local basis is given by

$$\left\{
\begin{array}{rcl}
\phi_i(\lambda) & = & \dfrac{1}{2}\lambda_i(3\lambda_i - 1)(3\lambda_i - 2), \\[2mm]
\phi_{iij}(\lambda) & = & \dfrac{9}{2}\lambda_i\lambda_j(3\lambda_i - 1), \\[2mm]
\phi_{ijk}(\lambda) & = & 27\lambda_i\lambda_j\lambda_k
\end{array}
\right\}.$$

In two dimensions, the function $\phi_{ijk}(\lambda)$ is called cell bubble function. $\qquad\square$

**Example 5.29** *Cubic Hermite element.* The finite element space is a subspace of $C(\overline{\Omega})$, its dimension is $(d+1)(d+2)(d+3)/6$ and the functionals are the values of the function in the vertices of the mesh cell ($(d+1)$ values), the value of the barycenter at the 2-faces of $K$ ($(d+1)(d-1)d/6$ values), and the partial derivatives at the vertices ($d(d+1)$ values), see Figure 5.8. The dimension is the same as for the $P_3$ element. Hence, the local polynomials can be defined to be cubic.
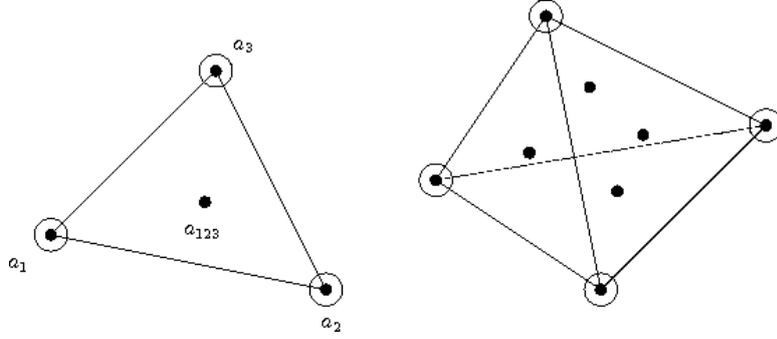


Figure 5.8: The cubic Hermite element.

This finite element does not define an affine family in the strict sense, because the functionals for the partial derivatives $\hat{\Phi}_i(\hat{v}) = \partial_i \hat{v}(\mathbf{0})$ on the reference cell are mapped to the functionals $\Phi_i(v) = \partial_{\mathbf{t}_i} v(\mathbf{a})$, where $\mathbf{a} = F_K(\mathbf{0})$ and $\mathbf{t}_i$ are the directions of edges which are adjacent to $\mathbf{a}$, i.e., $\mathbf{a}$ is an end point of this edge. This property suffices to control all first derivatives. On has to take care of this property in the implementation of this finite element.

Because of this property, one can use the derivatives in the direction of the edges as functionals

$$
\begin{array}{llll}
\Phi_i(v) & = & v(\mathbf{a}_i), & \text{(vertices)}\\
\Phi_{ij}(v) & = & \nabla v(\mathbf{a}_i) \cdot (\mathbf{a}_j - \mathbf{a}_i), \ i,j = 1,\ldots,d-1, i \neq j, & \text{(directional derivative)}\\
\Phi_{ijk}(v) & = & v(\mathbf{a}_{ijk}), \ i < j < k, & \text{(2-faces)}
\end{array}
$$

with the corresponding local basis

$$
\begin{array}{lll}
\phi_i(\lambda) & = & -2\lambda_i^3 + 3\lambda_i^2 - 7\lambda_i \sum_{j<k,j\neq i,k\neq i} \lambda_j \lambda_k,\\
\phi_{ij}(\lambda) & = & \lambda_i \lambda_j (2\lambda_i - \lambda_j - 1),\\
\phi_{ijk}(\lambda) & = & 27\lambda_i \lambda_j \lambda_k.
\end{array}
$$

The proof of the unisolvence can be found in the literature.

Here, the continuity of the functions will be shown only for $d = 2$. Let $K_1, K_2$ be two mesh cells with the common edge $E$ and the unit tangential vector $\mathbf{t}$. Let $V_1, V_2$ be the end points of $E$. The restriction $v|_{K_1}, v|_{K_2}$ to $E$ satisfy four conditions

$$
v|_{K_1}(V_i) = v|_{K_2}(V_i), \quad \partial_{\mathbf{t}} v|_{K_1}(V_i) = \partial_{\mathbf{t}} v|_{K_2}(V_i), \ i = 1, 2.
$$

Since both restrictions are cubic polynomials and four conditions have to be satisfied, their values coincide on $E$.

The cubic Hermite finite element possesses an advantage in comparison with the $P_3$ finite element. For $d = 2$, it holds for a regular triangulation $\mathcal{T}^h$ that

$$
\#(K) \approx 2\#(V), \quad \#(E) \approx 2\#(V),
$$

where $\#(\cdot)$ denotes the number of triangles, nodes, and edges, respectively. Hence, the dimension of $P_3$ is approximately $7\#(V)$, whereas the dimension of the cubic

Hermite element is approximately $5\#(V)$. This difference comes from the fact that both spaces are different. The elements of both spaces are continuous functions, but for the functions of the cubic Hermite finite element, in addition, the first derivatives are continuous at the nodes. That means, these two spaces are different finite element spaces whose degree of the local polynomial space is the same (cubic). One can see at this example the importance of the functionals for the definition of the global finite element space. □

**Example 5.30** $P_1^{\mathrm{nc}}$ : *nonconforming linear finite element, Crouzeix–Raviart finite element Crouzeix and Raviart (1973).* This finite element consists of piecewise linear but discontinuous functions. The functionals are given by the values of the functions in the barycenters of the faces such that $\dim P_1^{\mathrm{nc}}(K) = (d+1)$. It follows from the definition of the finite element space, Definition 5.12, that the functions from $P_1^{\mathrm{nc}}$ are continuous in the barycenter of the faces

$$P_1^{\mathrm{nc}} = \Big\{ v \in L^2(\Omega) \; : \; v|_K \in P_1(K), \; v(\mathbf{x}) \text{ is continuous at the barycenter}$$
$$\text{of all faces} \Big\}. \tag{5.3}$$

Equivalently, the functionals can be defined to be the integral mean values on the faces and then the global space is defined to be

$$P_1^{\mathrm{nc}} = \Bigg\{ v \in L^2(\Omega) \; : \; v|_K \in P_1(K),$$
$$\int_E v|_K \, d\mathbf{s} = \int_E v|_{K'} \, d\mathbf{s} \; \forall \, E \in \mathcal{E}(K) \cap \mathcal{E}(K') \Bigg\}, \tag{5.4}$$

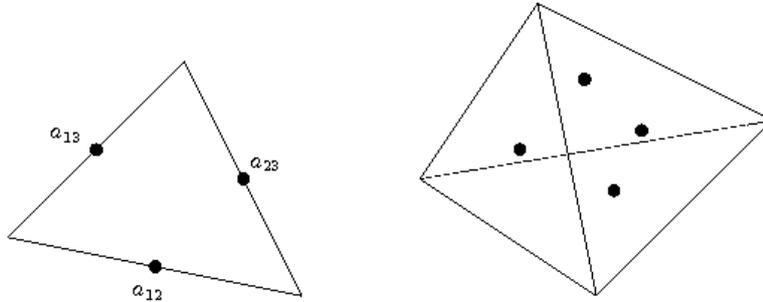where $\mathcal{E}(K)$ is the set of all $(d-1)$ dimensional faces of $K$.



Figure 5.9: The finite element $P_1^{\mathrm{nc}}$.

For the description of this finite element, one defines the functionals by

$$\Phi_i(v) = v(\mathbf{a}_{i-1,i+1}) \text{ for } d = 2, \quad \Phi_i(v) = v(\mathbf{a}_{i-2,i-1,i+1}) \text{ for } d = 3,$$

where the points are the barycenters of the faces with the vertices that correspond to the indices. This system is unisolvent with the local basis

$$\phi_i(\lambda) = 1 - d\lambda_i, \quad i = 1, \ldots, d+1.$$

□

## 5.3 Finite Elements on Parallelepipeds

**Remark 5.31** *Reference mesh cells, reference map.* On can find in the literature two reference cells: the unit cube $[0,1]^d$ and the large unit cube $[-1,1]^d$. It does

not matter which reference cell is chosen. Here, the large unit cube will be used: $\hat{K} = [-1, 1]^d$. The class of admissible reference maps $\{F_K\}$ consists of bijective affine mappings of the form

$$F_K \hat{\mathbf{x}} = B\hat{\mathbf{x}} + \mathbf{b}, \quad B \in \mathbb{R}^{d \times d}, \ \mathbf{b} \in \mathbb{R}^d.$$

If $B$ is a diagonal matrix, then $\hat{K}$ is mapped to $d$-rectangles.

The class of mesh cells which are obtained in this way is not sufficient to triangulate general domains. If one wants to use more general mesh cells than parallelepipeds, then the class of admissible reference maps has to be enlarged, see Section 5.4. □

**Definition 5.32 Polynomial space $Q_k$.** Let $\mathbf{x} = (x_1, \ldots, x_d)^T$ and denote by $\boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_d)^T$ a multi-index. Then, the polynomial space $Q_k$ is given by

$$Q_k = \text{span}\left\{ \prod_{i=1}^d x_i^{\alpha_i} = \mathbf{x}^{\boldsymbol{\alpha}} \ : \ 0 \leq \alpha_i \leq k \ \text{ for } \ i = 1, \ldots, d \right\}.$$

□

**Example 5.33 $Q_1$ vs. $P_1$.** The space $Q_1$ consists of all polynomials which are $d$-linear. Let $d = 2$, then it is

$$Q_1 = \text{span}\{1, x, y, xy\},$$

whereas

$$P_1 = \text{span}\{1, x, y\}.$$

□

**Remark 5.34** *Finite elements on $d$-rectangles.* For simplicity of presentation, the examples below consider $d$-rectangles. In this case, the finite elements are just tensor products of one-dimensional finite elements. In particular, the basis functions can be written as products of one-dimensional basis functions. □

**Example 5.35 $Q_0$** *: piecewise constant finite element.* Similarly to the $P_0$ space, the space $Q_0$ consists of piecewise constant, discontinuous functions. The functional is the value of the function in the barycenter of the mesh cell $K$ and it holds $\dim Q_0(K) = 1$. □

**Example 5.36 $Q_1$** *: conforming piecewise $d$-linear finite element.* This finite element space is a subspace of $C(\overline{\Omega})$. The functionals are the values of the function in the vertices of the mesh cell, see Figure 5.10. Hence, it is $\dim Q_1(K) = 2^d$.

The one-dimensional local basis functions, which will be used for the tensor product, are given by

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

With these functions, e.g., the basis functions in two dimensions are computed by

$$\hat{\phi}_1(\hat{x})\hat{\phi}_1(\hat{y}), \ \hat{\phi}_1(\hat{x})\hat{\phi}_2(\hat{y}), \ \hat{\phi}_2(\hat{x})\hat{\phi}_1(\hat{y}), \ \hat{\phi}_2(\hat{x})\hat{\phi}_2(\hat{y}).$$

The continuity of the functions of the finite element space $Q_1$ is proved in the same way as for simplicial finite elements. It is used that the restriction of a function from $Q_k(K)$ to a face $E$ is a function from the space $Q_k(E)$, $k \geq 1$. □
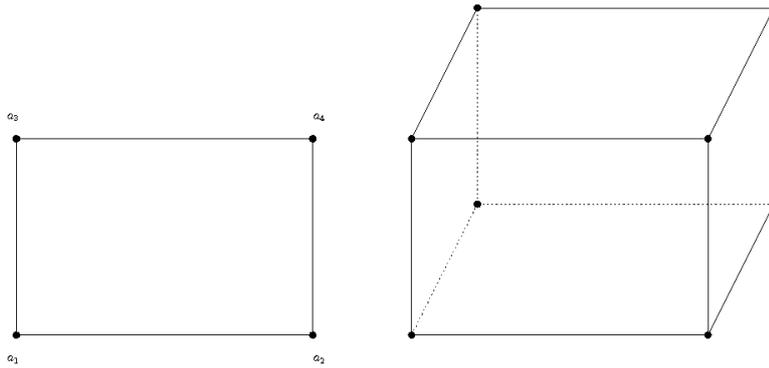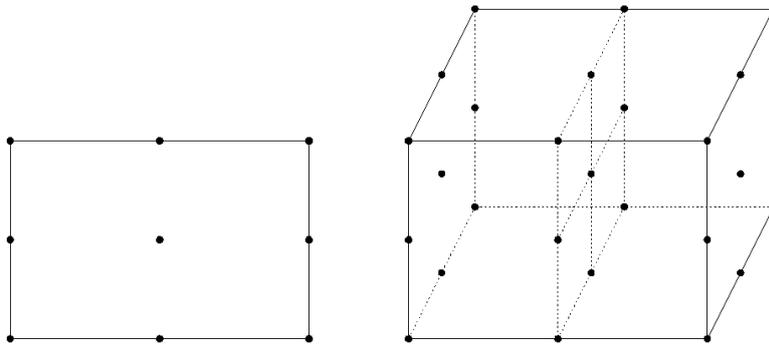
Figure 5.10: The finite element $Q_1$.



Figure 5.11: The finite element $Q_2$.

**Example 5.37** $Q_2$ : *conforming piecewise d-quadratic finite element.* It holds that $Q_2 \subset C(\overline{\Omega})$. The functionals in one dimension are the values of the function at both ends of the interval and in the center of the interval, see Figure 5.11. In $d$ dimensions, they are the corresponding values of the tensor product of the intervals. It follows that $\dim Q_2(K) = 3^d$.

The one-dimensional basis function on the reference interval are defined by

$$\hat{\phi}_1(\hat{x}) = -\frac{1}{2}\hat{x}(1-\hat{x}), \quad \hat{\phi}_2(\hat{x}) = (1-\hat{x})(1+\hat{x}), \quad \hat{\phi}_3(\hat{x}) = \frac{1}{2}(1+\hat{x})\hat{x}.$$

The basis function $\prod_{i=1}^d \hat{\phi}_2(\hat{x}_i)$ is called cell bubble function. $\qquad\square$

**Example 5.38** $Q_3$ : *conforming piecewise d-quadratic finite element.* This finite element space is a subspace of $C(\overline{\Omega})$. The functionals on the reference interval are given by the values at the end of the interval and the values at the points $\hat{x} = -1/3$, $\hat{x} = 1/3$. In multiple dimensions, it is the corresponding tensor product, see Figure 5.12. The dimension of the local space is $\dim Q_3(K) = 4^d$.

The one-dimensional basis functions in the reference interval are given by

$$
\begin{aligned}
\hat{\phi}_1(\hat{x}) &= -\frac{1}{16}(3\hat{x}+1)(3\hat{x}-1)(\hat{x}-1), \\
\hat{\phi}_2(\hat{x}) &= \frac{9}{16}(\hat{x}+1)(3\hat{x}-1)(\hat{x}-1), \\
\hat{\phi}_3(\hat{x}) &= -\frac{9}{16}(\hat{x}+1)(3\hat{x}+1)(\hat{x}-1), \\
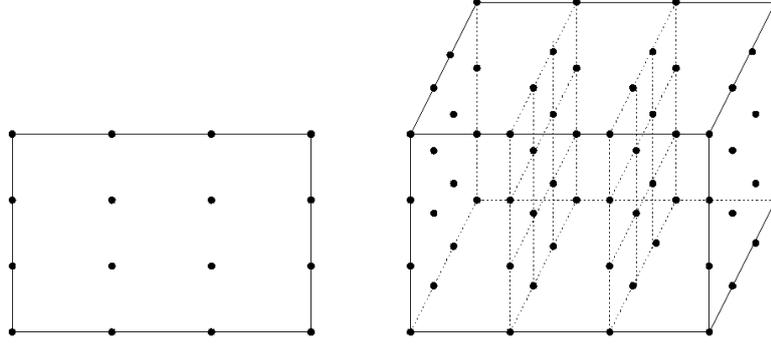\hat{\phi}_4(\hat{x}) &= \frac{1}{16}(3\hat{x}+1)(3\hat{x}-1)(\hat{x}+1).
\end{aligned}
$$

Figure 5.12: The finite element $Q_3$.

□

**Example 5.39** $Q_1^{\mathrm{rot}}$ : *rotated nonconforming element of lowest order, Rannacher–Turek element Rannacher and Turek (1992):* This finite element space is a generalization of the $P_1^{\mathrm{nc}}$ finite element to quadrilateral and hexahedral mesh cells. It consists of discontinuous functions which are continuous at the barycenter of the faces. The dimension of the local finite element space is $\dim Q_1^{\mathrm{rot}}(K) = 2d$. The space on the reference mesh cell is defined by

$$
\begin{aligned}
Q_1^{\mathrm{rot}}\left(\hat{K}\right) &= \{\hat{p} \,:\, \hat{p} \in \mathrm{span}\{1, \hat{x}, \hat{y}, \hat{x}^2 - \hat{y}^2\}\} && \text{for } d = 2, \\
Q_1^{\mathrm{rot}}\left(\hat{K}\right) &= \{\hat{p} \,:\, \hat{p} \in \mathrm{span}\{1, \hat{x}, \hat{y}, \hat{z}, \hat{x}^2 - \hat{y}^2, \hat{y}^2 - \hat{z}^2\}\} && \text{for } d = 3.
\end{aligned}
$$

Note that the transformed space

$$
Q_1^{\mathrm{rot}}(K) = \{p = \hat{p} \circ F_K^{-1}, \hat{p} \in Q_1^{\mathrm{rot}}(\hat{K})\}
$$

contains polynomials of the form $ax^2 - by^2$, where $a, b$ depend on $F_K$.
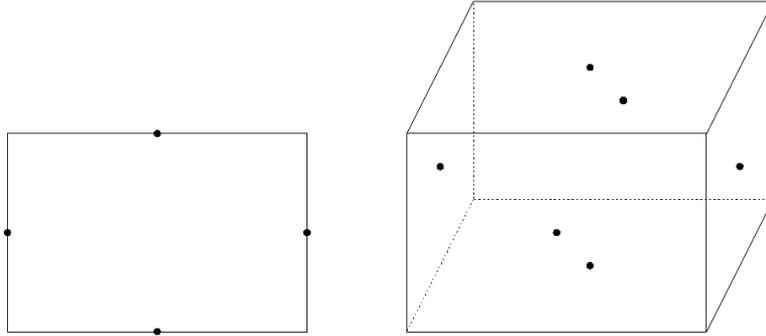


Figure 5.13: The finite element $Q_1^{\mathrm{rot}}$.

For $d = 2$, the local basis on the reference cell is given by

$$
\begin{aligned}
\phi_1(\hat{x}, \hat{y}) &= -\frac{3}{8}(\hat{x}^2 - \hat{y}^2) - \frac{1}{2}\hat{y} + \frac{1}{4}, \\
\phi_2(\hat{x}, \hat{y}) &= \frac{3}{8}(\hat{x}^2 - \hat{y}^2) + \frac{1}{2}\hat{x} + \frac{1}{4}, \\
\phi_3(\hat{x}, \hat{y}) &= -\frac{3}{8}(\hat{x}^2 - \hat{y}^2) + \frac{1}{2}\hat{y} + \frac{1}{4}, \\
\phi_4(\hat{x}, \hat{y}) &= \frac{3}{8}(\hat{x}^2 - \hat{y}^2) - \frac{1}{2}\hat{x} + \frac{1}{4}.
\end{aligned}
$$

72

Analogously to the Crouzeix–Raviart finite element, the functionals can be defined as point values of the functions in the barycenters of the faces, see Figure 5.13, or as integral mean values of the functions at the faces. Consequently, the finite element spaces are defined in the same way as (5.3) or (5.4), with $P_1^{\mathrm{nc}}(K)$ replaced by $Q_1^{\mathrm{rot}}(K)$.

In the code MooNMD John and Matthies (2004), the mean value oriented $Q_1^{\mathrm{rot}}$ finite element space is implemented fro two dimensions and the point value oriented $Q_1^{\mathrm{rot}}$ finite element space for three dimensions. For $d = 3$, the integrals on the faces of mesh cells, whose equality is required in the mean value oriented $Q_1^{\mathrm{rot}}$ finite element space, involve a weighting function which depends on the particular mesh cell $K$. The computation of these weighting functions for all mesh cells is an additional computational overhead. For this reason, Schieweck (Schieweck, 1997, p. 21) suggested to use for $d = 3$ the simpler point value oriented form of the $Q_1^{\mathrm{rot}}$ finite element. $\qquad\square$

## 5.4 Parametric Finite Elements on General $d$-Dimensional Quadrilaterals

**Remark 5.40** *Parametric mappings.* The image of an affine mapping of the reference mesh cell $\hat{K} = [-1, 1]^d$, $d \in \{2, 3\}$, is a parallelepiped. If one wants to consider finite elements on general $q$-quadrilaterals, then the class of admissible reference maps has to be enlarged.

The simplest parametric finite element on quadrilaterals in two dimensions uses bilinear mappings. Let $\hat{K} = [-1, 1]^2$ and let

$$F_K(\hat{\mathbf{x}}) = \begin{pmatrix} F_K^1(\hat{\mathbf{x}}) \\ F_K^2(\hat{\mathbf{x}}) \end{pmatrix} = \begin{pmatrix} a_{11} + a_{12}\hat{x} + a_{13}\hat{y} + a_{14}\hat{x}\hat{y} \\ a_{21} + a_{22}\hat{x} + a_{23}\hat{y} + a_{24}\hat{x}\hat{y} \end{pmatrix}, F_K^i \in Q_1, \ i = 1, 2,$$

be a bilinear mapping from $\hat{K}$ on the class of admissible quadrilaterals. A quadrilateral $K$ is called admissible if
- the length of all edges of $K$ is larger than zero,
- the interior angles of $K$ are smaller than $\pi$, i.e. $K$ is convex.

This class contains, e.g., trapezoids and rhombi. $\qquad\square$

**Remark 5.41** *Parametric finite element functions.* The functions of the local space $P(K)$ on the mesh cell $K$ are defined by $p = \hat{p} \circ F_K^{-1}$. These functions are in general rational functions. However, using $d$-linear mappings, then the restriction of $F_K$ on an edge of $\hat{K}$ is an affine map. For instance, in the case of the $Q_1$ finite element, the functions on $K$ are linear functions on each edge of $K$ for this reason. It follows that the functions of the corresponding finite element space are continuous, see Example 5.26. $\qquad\square$

## 5.5 Transform of Integrals

**Remark 5.42** *Motivation.* The transform of integrals from the reference mesh cell to mesh cells of the grid and vice versa is used as well for analysis as for the implementation of finite element methods. This section provides an overview of the most important formulae for transforms.

Let $\hat{K} \subset \mathbb{R}^d$ be the reference mesh cell, $K$ be an arbitrary mesh cell, and $F_K : \hat{K} \to K$ with $\mathbf{x} = F_K(\hat{\mathbf{x}})$ be the reference map. It is assumed that the reference map is a continuous differentiable one-to-one map. The inverse map is

denoted by $F_K^{-1} : K \to \hat{K}$. For the integral transforms, the derivatives (Jacobians) of $F_K$ and $F_K^{-1}$ are needed

$$DF_K(\hat{\mathbf{x}})_{ij} = \frac{\partial x_i}{\partial \hat{x}_j}, \quad DF_K^{-1}(\mathbf{x})_{ij} = \frac{\partial \hat{x}_i}{\partial x_j}, \quad i,j = 1,\dots,d.$$

$\square$

**Remark 5.43** *Integral with a function without derivatives.* This integral transforms with the standard rule of integral transforms

$$\int_K v(\mathbf{x})\ d\mathbf{x} = \int_{\hat{K}} \hat{v}(\hat{\mathbf{x}})\ |\det DF_K(\hat{\mathbf{x}})|\ d\hat{\mathbf{x}}, \tag{5.5}$$

where $\hat{v}(\hat{\mathbf{x}}) = v(F_K(\hat{\mathbf{x}}))$. $\square$

**Remark 5.44** *Transform of derivatives.* Using the chain rule, one obtains

$$
\begin{aligned}
\frac{\partial v}{\partial x_i}(\mathbf{x}) &= \sum_{j=1}^d \frac{\partial \hat{v}}{\partial \hat{x}_j}(\hat{\mathbf{x}}) \frac{\partial \hat{x}_j}{\partial x_i} = \nabla_{\hat{\mathbf{x}}} \hat{v}(\hat{\mathbf{x}}) \cdot \left( \left( DF_K^{-1}(\mathbf{x}) \right)^T \right)_i \\
&= \nabla_{\hat{\mathbf{x}}} \hat{v}(\hat{\mathbf{x}}) \cdot \left( \left( DF_K^{-1}(F_K(\hat{\mathbf{x}})) \right)^T \right)_i, \tag{5.6} \\
\frac{\partial \hat{v}}{\partial \hat{x}}(\hat{\mathbf{x}}) &= \sum_{j=1}^d \frac{\partial v}{\partial x_j}(\mathbf{x}) \frac{\partial x_j}{\partial \hat{x}_i} = \nabla v(\mathbf{x}) \cdot \left( \left( DF_K(\hat{\mathbf{x}}) \right)^T \right)_i \\
&= \nabla v(\mathbf{x}) \cdot \left( \left( DF_K(F_K^{-1}(\mathbf{x})) \right)^T \right)_i. \tag{5.7}
\end{aligned}
$$

The index $i$ denotes the $i$-th row of a matrix. Derivatives on the reference mesh cell are marked with a symbol on the operator. $\square$

**Remark 5.45** *Integrals with a gradients.* Using the rule for transforming integrals and (5.6) gives

$$
\begin{aligned}
&\int_K \mathbf{b}(\mathbf{x}) \cdot \nabla v(\mathbf{x})\ d\mathbf{x} \\
&= \int_{\hat{K}} \mathbf{b}\left(F_K(\hat{\mathbf{x}})\right) \cdot \left[ \left(DF_K^{-1}\right)^T \left(F_K(\hat{\mathbf{x}})\right) \right] \nabla_{\hat{\mathbf{x}}} \hat{v}(\hat{\mathbf{x}})\ |\det DF_K(\hat{\mathbf{x}})|\ d\hat{\mathbf{x}}. \tag{5.8}
\end{aligned}
$$

Similarly, one obtains

$$
\begin{aligned}
&\int_K \nabla v(\mathbf{x}) \cdot \nabla w(\mathbf{x})\ d\mathbf{x} \\
&= \int_{\hat{K}} \left[ \left(DF_K^{-1}\right)^T \left(F_K(\hat{\mathbf{x}})\right) \right] \nabla_{\hat{\mathbf{x}}} \hat{v}(\hat{\mathbf{x}}) \cdot \left[ \left(DF_K^{-1}\right)^T \left(F_K(\hat{\mathbf{x}})\right) \right] \nabla_{\hat{\mathbf{x}}} \hat{w}(\hat{\mathbf{x}}) \\
&\quad \times |\det DF_K(\hat{\mathbf{x}})|\ d\hat{\mathbf{x}}. \tag{5.9}
\end{aligned}
$$

$\square$

**Remark 5.46** *Integral with the divergence.* Integrals of the following type are important for the Navier–Stokes equations

$$
\begin{aligned}
\int_K \nabla \cdot v(\mathbf{x}) q(\mathbf{x})\ d\mathbf{x} &= \int_K \sum_{i=1}^d \frac{\partial v_i}{\partial x_i}(\mathbf{x}) q(\mathbf{x})\ d\mathbf{x} \\
&= \int_{\hat{K}} \sum_{i=1}^d \left[ \left( \left(DF_K^{-1}(F_K(\hat{\mathbf{x}}))\right)^T \right)_i \cdot \nabla_{\hat{\mathbf{x}}} \hat{v}_i(\hat{\mathbf{x}}) \right] \hat{q}(\hat{\mathbf{x}})\ |\det DF_K(\hat{\mathbf{x}})|\ d\hat{\mathbf{x}} \\
&= \int_{\hat{K}} \left[ \left(DF_K^{-1}(F_K(\hat{\mathbf{x}}))\right)^T : D_{\hat{\mathbf{x}}} \mathbf{v}(\hat{\mathbf{x}}) \right] \hat{q}(\hat{\mathbf{x}})\ |\det DF_K(\hat{\mathbf{x}})|\ d\hat{\mathbf{x}}. \tag{5.10}
\end{aligned}
$$

In the derivation, (5.6) was used. $\square$

**Example 5.47** *Affine transform.* The most important class of reference maps are affine transforms

$$\mathbf{x} = B\hat{\mathbf{x}} + \mathbf{b}, \quad B \in \mathbb{R}^{d \times d}, \mathbf{b} \in \mathbb{R}^d,$$

where the invertible matrix $B$ and the vector $\mathbf{b}$ are constants. It follows that

$$\hat{\mathbf{x}} = B^{-1}(\mathbf{x} - \mathbf{b}) = B^{-1}\mathbf{x} - B^{-1}\mathbf{b}.$$

In this case, there are

$$DF_K = B, \quad DF_K^{-1} = B^{-1}, \quad \det DF_K = \det(B).$$

One obtains for the integral transforms from (5.5), (5.8), (5.9), and (5.10)

$$\int_K v(\mathbf{x}) \, d\mathbf{x} \;=\; |\det(B)| \int_{\hat{K}} \hat{v}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}, \tag{5.11}$$

$$\int_K \mathbf{b}(\mathbf{x}) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x} \;=\; |\det(B)| \int_{\hat{K}} \mathbf{b}\,(F_K(\hat{\mathbf{x}})) \cdot B^{-T}\nabla_{\hat{\mathbf{x}}}\hat{v}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}, \tag{5.12}$$

$$\int_K \nabla v(\mathbf{x}) \cdot \nabla w(\mathbf{x}) \, d\mathbf{x} \;=\; |\det(B)| \int_{\hat{K}} B^{-T}\nabla_{\hat{\mathbf{x}}}\hat{v}(\hat{\mathbf{x}}) \cdot B^{-T}\nabla_{\hat{\mathbf{x}}}\hat{w}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}, \tag{5.13}$$

$$\int_K \nabla \cdot v(\mathbf{x})q(\mathbf{x}) \, d\mathbf{x} \;=\; |\det(B)| \int_{\hat{K}} \left[B^{-T} : D_{\hat{\mathbf{x}}}\mathbf{v}(\hat{\mathbf{x}})\right] \hat{q}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}. \tag{5.14}$$

$\square$

# Chapter 6

# Interpolation

**Remark 6.1** *Motivation.* Variational forms of partial differential equations use functions in Sobolev spaces. The solution of these equations shall be approximated with the Ritz method in finite dimensional spaces, the finite element spaces. The best possible approximation of an arbitrary function from the Sobolev space by a finite element function is a factor in the upper bound for the finite element error, e.g., see the Lemma of Cea, estimate (4.19).

This section studies the approximation quality of finite element spaces. Estimates are proved for interpolants of functions. Interpolation estimates are of course upper bounds for the best approximation error and they can serve as factors in finite element error estimates. □

## 6.1 Interpolation in Sobolev Spaces by Polynomials

**Lemma 6.2 Unique determination of a polynomial with integral conditions.** *Let $\Omega$ be a bounded domain in $\mathbb{R}^d$ with Lipschitz boundary. Let $m \in \mathbb{N} \cup \{0\}$ be given and let for all derivatives with multi-index $\boldsymbol{\alpha}$, $|\boldsymbol{\alpha}| \leq m$, a value $a_{\boldsymbol{\alpha}} \in \mathbb{R}$ be given. Then, there is a uniquely determined polynomial $p \in P_m(\Omega)$ such that*

$$\int_\Omega \partial_{\boldsymbol{\alpha}} p(\mathbf{x}) \; d\mathbf{x} = a_{\boldsymbol{\alpha}}, \quad |\boldsymbol{\alpha}| \leq m. \tag{6.1}$$

**Proof:** Let $p \in P_m(\Omega)$ be an arbitrary polynomial. It has the form

$$p(\mathbf{x}) = \sum_{|\boldsymbol{\beta}| \leq m} b_{\boldsymbol{\beta}} \mathbf{x}^{\boldsymbol{\beta}}.$$

Inserting this representation into (6.1) leads to a linear system of equations $M\mathbf{b} = \mathbf{a}$ with

$$M = (M_{\boldsymbol{\alpha}\boldsymbol{\beta}}), \; M_{\boldsymbol{\alpha}\boldsymbol{\beta}} = \int_\Omega \partial_{\boldsymbol{\alpha}} \mathbf{x}^{\boldsymbol{\beta}} \; d\mathbf{x}, \; \mathbf{b} = (b_{\boldsymbol{\beta}}), \; \mathbf{a} = (a_{\boldsymbol{\alpha}}),$$

for $|\boldsymbol{\alpha}|, |\boldsymbol{\beta}| \leq m$. Since $M$ is a squared matrix, the linear system of equations possesses a unique solution if and only if $M$ is non-singular.

The proof is performed by contradiction. Assume that $M$ is singular. Then there exists a non-trivial solution of the homogeneous system. That means, there is a polynomial $q \in P_m(\Omega) \setminus \{0\}$ with

$$\int_\Omega \partial_{\boldsymbol{\alpha}} q(\mathbf{x}) \; d\mathbf{x} = 0 \text{ for all } |\boldsymbol{\alpha}| \leq m.$$

The polynomial $q(\mathbf{x})$ has the representation $q(\mathbf{x}) = \sum_{|\boldsymbol{\beta}| \leq m} c_{\boldsymbol{\beta}} \mathbf{x}^{\boldsymbol{\beta}}$. Now, one can choose a $c_{\boldsymbol{\beta}} \neq 0$ with maximal value $|\boldsymbol{\beta}|$. Then, it is $\partial_{\boldsymbol{\beta}} q(\mathbf{x}) = C c_{\boldsymbol{\beta}} = const \neq 0$, where $C > 0$ comes

from the differentiation rule for polynomials, which is a contradiction to the vanishing of the integral for $\partial_{\boldsymbol{\beta}} q(\mathbf{x})$. ∎

**Remark 6.3** *To Lemma 6.2.* Lemma 6.2 states that a polynomial is uniquely determined if a condition on the integral on $\Omega$ is prescribed for each derivative. □

**Lemma 6.4 Poincaré-type inequality.** *Denote by $D^k v(\mathbf{x})$, $k \in \mathbb{N} \cup \{0\}$, the total derivative of order $k$ of a function $v(\mathbf{x})$, e.g., for $k = 1$ the gradient of $v(\mathbf{x})$. Let $\Omega$ be convex and be included into a ball of radius $R$. Let $k, l \in \mathbb{N} \cup \{0\}$ with $k \leq l$ and let $p \in \mathbb{R}$ with $p \in [1, \infty]$. Assume that $v \in W^{l,p}(\Omega)$ satisfies*

$$\int_{\Omega} \partial_{\boldsymbol{\alpha}} v(\mathbf{x}) \, d\mathbf{x} = 0 \text{ for all } |\boldsymbol{\alpha}| \leq l - 1,$$

*then it holds the estimate*

$$\left\| D^k v \right\|_{L^p(\Omega)} \leq C R^{l-k} \left\| D^l v \right\|_{L^p(\Omega)},$$

*where the constant $C$ does not depend on $\Omega$ and on $v(\mathbf{x})$.*

**Proof:** There is nothing to prove if $k = l$. In addition, it suffices to prove the lemma for $k = 0$ and $l = 1$, since the general case follows by applying the result to $\partial_{\boldsymbol{\alpha}} v(\mathbf{x})$. Only the case $p < \infty$ will be discussed here in detail.

Since $\Omega$ is assumed to be convex, the integral mean value theorem can be written in the form

$$v(\mathbf{x}) - v(\mathbf{y}) = \int_0^1 \nabla v(t\mathbf{x} + (1-t)\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \, dt, \quad \mathbf{x}, \mathbf{y} \in \Omega.$$

Integration with respect to $\mathbf{y}$ yields

$$v(\mathbf{x}) \int_{\Omega} d\mathbf{y} - \int_{\Omega} v(\mathbf{y}) \, d\mathbf{y} = \int_{\Omega} \int_0^1 \nabla v(t\mathbf{x} + (1-t)\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \, dt \, d\mathbf{y}.$$

It follows from the assumption that the second integral on the left hand side vanishes. Hence, one gets

$$v(\mathbf{x}) = \frac{1}{|\Omega|} \int_{\Omega} \int_0^1 \nabla v(t\mathbf{x} + (1-t)\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \, dt \, d\mathbf{y}.$$

Now, taking the absolute value on both sides, using that the absolute value of an integral is estimated from above by the integral of the absolute value, applying the Cauchy–Schwarz inequality for vectors and the estimate $\|\mathbf{x} - \mathbf{y}\|_2 \leq 2R$ yields

$$
\begin{aligned}
|v(\mathbf{x})| &= \frac{1}{|\Omega|} \left| \int_{\Omega} \int_0^1 \nabla v(t\mathbf{x} + (1-t)\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y}) \, dt \, d\mathbf{y} \right| \\
&\leq \frac{1}{|\Omega|} \int_{\Omega} \int_0^1 |\nabla v(t\mathbf{x} + (1-t)\mathbf{y}) \cdot (\mathbf{x} - \mathbf{y})| \, dt \, d\mathbf{y} \\
&\leq \frac{2R}{|\Omega|} \int_{\Omega} \int_0^1 \|\nabla v(t\mathbf{x} + (1-t)\mathbf{y})\|_2 \, dt \, d\mathbf{y}. \quad (6.2)
\end{aligned}
$$

Then (6.2) is raised to the power $p$ and then integrated with respect to $\mathbf{x}$. One obtains with Hölder's inequality (3.4), with $p^{-1} + q^{-1} = 1 \implies p/q - p = p(1/q - 1) = -1$, that

$$
\begin{aligned}
\int_{\Omega} |v(\mathbf{x})|^p \, d\mathbf{x} &\leq \frac{CR^p}{|\Omega|^p} \int_{\Omega} \left( \int_{\Omega} \int_0^1 \|\nabla v(t\mathbf{x} + (1-t)\mathbf{y})\|_2 \, dt \, d\mathbf{y} \right)^p d\mathbf{x} \\
&\leq \frac{CR^p}{|\Omega|^p} \int_{\Omega} \left[ \underbrace{\left( \int_{\Omega} \int_0^1 1^q \, dt \, d\mathbf{y} \right)^{p/q}}_{|\Omega|^{p/q}} \right. \\
&\quad \left. \times \left( \int_{\Omega} \int_0^1 \|\nabla v(t\mathbf{x} + (1-t)\mathbf{y})\|_2^p \, dt \, d\mathbf{y} \right) \right] d\mathbf{x} \\
&= \frac{CR^p}{|\Omega|} \int_{\Omega} \left( \int_{\Omega} \int_0^1 \|\nabla v(t\mathbf{x} + (1-t)\mathbf{y})\|_2^p \, dt \, d\mathbf{y} \right) d\mathbf{x}.
\end{aligned}
$$

Applying the theorem of Fubini allows the commutation of the integration

$$\int_\Omega |v(\mathbf{x})|^p \ d\mathbf{x} \le \frac{CR^p}{|\Omega|} \int_0^1 \int_\Omega \left( \int_\Omega \|\nabla v(t\mathbf{x} + (1-t)\mathbf{y})\|_2^p \ d\mathbf{y} \right) d\mathbf{x} \ dt.$$

Using the integral mean value theorem in one dimension gives that there is a $t_0 \in [0,1]$, such that

$$\int_\Omega |v(\mathbf{x})|^p \ d\mathbf{x} \le \frac{CR^p}{|\Omega|} \int_\Omega \left( \int_\Omega \|\nabla v(t_0\mathbf{x} + (1-t_0)\mathbf{y})\|_2^p \ d\mathbf{y} \right) d\mathbf{x}.$$

The function $\|\nabla v(\mathbf{x})\|_2^p$ will be extended to $\mathbb{R}^d$ by zero and the extension will be also denoted by $\|\nabla v(\mathbf{x})\|_2^p$. Then, it is

$$\int_\Omega |v(\mathbf{x})|^p \ d\mathbf{x} \le \frac{CR^p}{|\Omega|} \int_\Omega \left( \int_{\mathbb{R}^d} \|\nabla v(t_0\mathbf{x} + (1-t_0)\mathbf{y})\|_2^p \ d\mathbf{y} \right) d\mathbf{x}. \qquad (6.3)$$

Let $t_0 \in [0, 1/2]$. Since the domain of integration is $\mathbb{R}^d$, a substitution of variables $t_0\mathbf{x} + (1-t_0)\mathbf{y} = \mathbf{z}$ can be applied and leads to

$$\int_{\mathbb{R}^d} \|\nabla v(t_0\mathbf{x} + (1-t_0)\mathbf{y})\|_2^p \ d\mathbf{y} = \frac{1}{1-t_0} \int_{\mathbb{R}^d} \|\nabla v(\mathbf{z})\|_2^p \ d\mathbf{z} \le 2 \|\nabla v\|_{L^p(\Omega)}^p,$$

since $1/(1-t_0) \le 2$. Inserting this expression into (6.3) gives

$$\int_\Omega |v(\mathbf{x})|^p \ d\mathbf{x} \le 2CR^p \|\nabla v\|_{L^p(\Omega)}^p.$$

If $t_0 > 1/2$ then one changes the roles of $\mathbf{x}$ and $\mathbf{y}$, applies the theorem of Fubini to change the sequence of integration, and uses the same arguments.

The estimate for the case $p = \infty$ is also based on (6.2).  ■

**Remark 6.5** *On Lemma 6.4.* The Lemma 6.4 proves an inequality of Poincaré-type. It says that it is possible to estimate the $L^p(\Omega)$ norm of a lower derivative of a function $v(\mathbf{x})$ by the same norm of a higher derivative if the integral mean values of some lower derivatives vanish.

An important application of Lemma 6.4 is in the proof of the Bramble–Hilbert lemma. The Bramble–Hilbert lemma considers a continuous linear functional which is defined on a Sobolev space and which vanishes for all polynomials of degree less or equal than $m$. It states that the value of the functional can be estimated by the Lebesgue norm of the $(m+1)$th total derivative of the functions from this Sobolev space.  □

**Theorem 6.6** *Bramble–Hilbert lemma.* *Let $m \in \mathbb{N} \cup \{0\}$, $m \ge 0$, $p \in [1, \infty]$, and $F : W^{m+1,p}(\Omega) \to \mathbb{R}$ be a continuous linear functional, and let the conditions of Lemma 6.2 and 6.4 be satisfied. Let*

$$F(p) = 0 \quad \forall \ p \in P_m(\Omega),$$

*then there is a constant $C(\Omega)$, which is independent of $v(\mathbf{x})$ and $F$, such that*

$$|F(v)| \le C(\Omega) \left\| D^{m+1}v \right\|_{L^p(\Omega)} \quad \forall \ v \in W^{m+1,p}(\Omega).$$

**Proof:** Let $v \in W^{m+1,p}(\Omega)$. It follows from Lemma 6.2 that there is a polynomial from $P_m(\Omega)$ with

$$\int_\Omega \partial_{\boldsymbol{\alpha}}(v+p)(\mathbf{x}) \ d\mathbf{x} = 0 \text{ for } |\boldsymbol{\alpha}| \le m.$$

Lemma 6.4 gives, with $l = m+1$ and considering each term in $\|\cdot\|_{W^{m+1,p}(\Omega)}$ individually, the estimate

$$\|v+p\|_{W^{m+1,p}(\Omega)} \le C(\Omega) \left\| D^{m+1}(v+p) \right\|_{L^p(\Omega)} = C(\Omega) \left\| D^{m+1}v \right\|_{L^p(\Omega)}.$$

From the vanishing of $F$ for $p \in P_m(\Omega)$ and the continuity of $F$ it follows that

$$|F(v)| = |F(v+p)| \le c \, \|v+p\|_{W^{m+1,p}(\Omega)} \le C(\Omega) \left\| D^{m+1} v \right\|_{L^p(\Omega)}.$$

∎

**Remark 6.7** *Strategy for estimating the interpolation error.* The Bramble–Hilbert lemma will be used for estimating the interpolation error for an affine family of finite elements. The strategy is as follows:
- Show first the estimate on the reference mesh cell $\hat{K}$.
- Transform the estimate on an arbitrary mesh cell $K$ to the reference mesh cell $\hat{K}$.
- Apply the estimate on $\hat{K}$.
- Transform back to $K$.

One has to study what happens if the transforms are applied to the estimate. □

**Remark 6.8** *Assumptions, definition of the interpolant.* Let $\hat{K} \subset \mathbb{R}^d, d \in \{2,3\}$, be a reference mesh cell (compact polyhedron), $\hat{P}(\hat{K})$ a polynomial space of dimension $N$, and $\hat{\Phi}_1, \ldots, \hat{\Phi}_N \;:\; C^s(\hat{K}) \to \mathbb{R}$ continuous linear functionals. It will be assumed that the space $\hat{P}(\hat{K})$ is unisolvent with respect to these functionals. Then, there is a local basis $\hat{\phi}_1, \ldots, \hat{\phi}_N \in \hat{P}(\hat{K})$.

Consider $\hat{v} \in C^s(\hat{K})$, then the interpolant $I_{\hat{K}} \hat{v} \in \hat{P}(\hat{K})$ is defined by

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

The operator $I_{\hat{K}}$ is a continuous and linear operator from $C^s(\hat{K})$ to $\hat{P}(\hat{K})$. From the linearity it follows that $I_{\hat{K}}$ is the identity on $\hat{P}(\hat{K})$

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

□

**Example 6.9** *Interpolation operators.*
- Let $\hat{K} \subset \mathbb{R}^d$ be an arbitrary reference cell, $\hat{P}(\hat{K}) = P_0(\hat{K})$, and

$$\hat{\Phi}(\hat{v}) = \frac{1}{\left| \hat{K} \right|} \int_{\hat{K}} \hat{v}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}.$$

The functional $\hat{\Phi}$ is continuous on $C^0(\hat{K})$ since

$$\left| \hat{\Phi}(\hat{v}) \right| \le \frac{1}{\left| \hat{K} \right|} \int_{\hat{K}} |\hat{v}(\hat{\mathbf{x}})| \; d\hat{\mathbf{x}} \le \frac{\left| \hat{K} \right|}{\left| \hat{K} \right|} \max_{\hat{\mathbf{x}} \in \hat{K}} |\hat{v}(\hat{\mathbf{x}})| = \|\hat{v}\|_{C^0(\hat{K})}.$$

For the constant function $1 \in P_0(\hat{K})$ it is $\hat{\Phi}(1) = 1 \neq 0$. Hence, $\{\hat{\phi}\} = \{1\}$ is the local basis and the space is unisolvent with respect to $\hat{\Phi}$. The operator

$$I_{\hat{K}} \hat{v}(\hat{\mathbf{x}}) = \hat{\Phi}(\hat{v}) \hat{\phi}(\hat{\mathbf{x}}) = \frac{1}{\left| \hat{K} \right|} \int_{\hat{K}} \hat{v}(\hat{\mathbf{x}}) \, d\hat{\mathbf{x}}$$

is an integral mean value operator, i.e., each continuous function on $\hat{K}$ will be approximated by a constant function whose value equals the integral mean value, see Figure 6.1
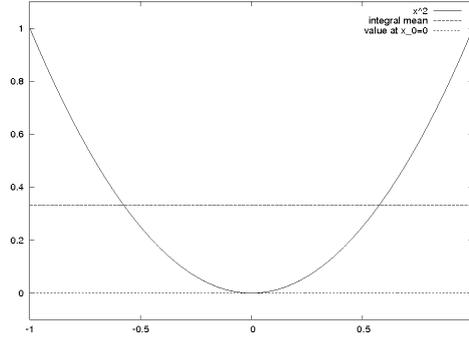
Figure 6.1: Interpolation of $x^2$ in $[-1, 1]$ by a $P_0$ function with the integral mean value and with the value of the function at $x_0 = 0$.

- It is possible to define $\hat{\Phi}(\hat{v}) = \hat{v}(\hat{\mathbf{x}}_0)$ for an arbitrary point $\hat{\mathbf{x}}_0 \in \hat{K}$. This functional is also linear and continuous in $C^0(\hat{K})$. The interpolation operator $I_{\hat{K}}$ defined in this way interpolates each continuous function by a constant function whose value is equal to the value of the function at $\hat{\mathbf{x}}_0$, see also Figure 6.1.
  Interpolation operators which are defined by using values of functions, are called Lagrangian interpolation operators.

This example demonstrates that the interpolation operator $I_{\hat{K}}$ depends on $\hat{P}(\hat{K})$ and on the functionals $\hat{\Phi}_i$. □

**Theorem 6.10 Interpolation error estimate on a reference mesh cell.** *Let $P_m(\hat{K}) \subset \hat{P}(\hat{K})$ and $p \in [1, \infty]$ with $(m + 1 - s)p > d$. Then there is a constant $C$ that is independent of $\hat{v}(\hat{\mathbf{x}})$ such that*

$$\left\| \hat{v} - I_{\hat{K}} \hat{v} \right\|_{W^{m+1,p}(\hat{K})} \leq C \left\| D^{m+1} \hat{v} \right\|_{L^p(\hat{K})} \quad \forall \, \hat{v} \in W^{m+1,p}(\hat{K}). \tag{6.4}$$

**Proof:** Because of the Sobolev imbedding, Theorem 3.53, ($\lambda = 0, j = s, m$ (of Sobolev imbedding) $= m + 1 - s$) it holds that

$$W^{m+1,p}(\hat{K}) \to C^s(\hat{K})$$

if $(m+1-s)p > d$. That means, the interpolation operator is well defined in $W^{m+1,p}(\hat{K})$. From the identity of the interpolation operator in $P_m(\hat{K})$, the triangle inequality, the boundedness of the interpolation operator (it is a linear and continuous operator mapping $C^s(\hat{K}) \to \hat{P}(\hat{K}) \subset W^{m+1,p}(\hat{K})$), and the Sobolev imbedding, one obtains for $\hat{q} \in P_m(\hat{K})$

$$
\begin{aligned}
\left\| \hat{v} - I_{\hat{K}} \hat{v} \right\|_{W^{m+1,p}(\hat{K})} &= \left\| \hat{v} + \hat{q} - I_{\hat{K}}(\hat{v} + \hat{q}) \right\|_{W^{m+1,p}(\hat{K})} \\
&\leq \left\| \hat{v} + \hat{q} \right\|_{W^{m+1,p}(\hat{K})} + \left\| I_{\hat{K}}(\hat{v} + \hat{q}) \right\|_{W^{m+1,p}(\hat{K})} \\
&\leq \left\| \hat{v} + \hat{q} \right\|_{W^{m+1,p}(\hat{K})} + c \left\| \hat{v} + \hat{q} \right\|_{C^s(\hat{K})} \\
&\leq c \left\| \hat{v} + \hat{q} \right\|_{W^{m+1,p}(\hat{K})}.
\end{aligned}
$$

Choosing $\hat{q}(\hat{\mathbf{x}})$ in Lemma 6.2 such that

$$\int_{\hat{K}} \partial_{\boldsymbol{\alpha}} (\hat{v} + \hat{q}) \, d\hat{\mathbf{x}} = 0 \quad \forall \, |\boldsymbol{\alpha}| \leq m,$$

the assumptions of Lemma 6.4 are satisfied. It follows that

$$\left\| \hat{v} + \hat{q} \right\|_{W^{m+1,p}(\hat{K})} \leq c \left\| D^{m+1}(\hat{v} + \hat{q}) \right\|_{L^p(\hat{K})} = c \left\| D^{m+1} \hat{v} \right\|_{L^p(\hat{K})}.$$

■

**Remark 6.11** *On Theorem 6.10.*

- One can construct examples which show that the Sobolev imbedding is not valid if $(m + 1 - s)p > d$ is not satisfied. In the case $(m + 1 - s)p \le d$, the statement of Theorem 6.10 is not true.

  Consider the interpolation of continuous functions $(s = 0)$ with piecewise linear elements $(m = 1)$ in Sobolev spaces that are also Hilbert spaces $(p = 2)$. Then $(m+1-s)p = 4$ and it follows that the theorem can be applied only for $d \in \{2, 3\}$. For piecewise constant finite elements, the statement of the theorem is true only for $d = 1$.
- The theorem requires only that $P_m(\hat{K}) \subset \hat{P}(\hat{K})$. This requirement does not exclude that $\hat{P}(\hat{K})$ contains polynomials of higher degree, too. However, this property is not utilized and also not needed if the other assumptions of the theorem are satisfied.

$\square$

**Remark 6.12** *Assumptions on the triangulation.* For deriving the interpolation error estimate for arbitrary mesh cells $K$, and finally for the finite element space, one has to study the properties of the affine mapping from $K$ to $\hat{K}$ and of the back mapping.

Consider an affine family of finite elements whose mesh cells are generated by affine mappings

$$F_K \hat{\mathbf{x}} = B\hat{\mathbf{x}} + \mathbf{b},$$

where $B$ is a non-singular $d \times d$ matrix and $\mathbf{b}$ is a $d$ vector.

Let $h_K$ be the diameter of $K = F_K(\hat{K})$, i.e., the largest distance of two points that are contained in $K$. The images $\{K = F_K(\hat{K})\}$ are assumed to satisfy the following conditions:
- $K \subset \mathbb{R}^d$ is contained in a ball of radius $C_R h_K$,
- $K$ contains a ball of radius $C_R^{-1} h_K$,

where the constant $C_R$ is independent of $K$. Hence, it follows for all $K$ that

$$\frac{\text{radius of circumcircle}}{\text{radius of inscribed circle}} \le C_R^2.$$

A triangulation with this property is called a quasi-uniform triangulation. $\square$

**Lemma 6.13 Estimates of matrix norms.** *For each matrix norm $\|\cdot\|$ one has the estimates*

$$\|B\| \le ch_K, \quad \|B^{-1}\| \le ch_K^{-1},$$

*where the constants depend on the matrix norm and on $C_R$.*

**Proof:** Since $\hat{K}$ is a Lipschitz domain with polyhedral boundary, it contains a ball $B(\hat{\mathbf{x}}_0, r)$ with $\hat{\mathbf{x}}_0 \in \hat{K}$ and some $r > 0$. Hence, $\hat{\mathbf{x}}_0 + \hat{\mathbf{y}} \in \hat{K}$ for all $\|\hat{\mathbf{y}}\|_2 = r$. It follows that the images

$$\mathbf{x}_0 = B\hat{\mathbf{x}}_0 + \mathbf{b}, \quad \mathbf{x} = B(\hat{\mathbf{x}}_0 + \hat{\mathbf{y}}) + \mathbf{b} = \mathbf{x}_0 + B\hat{\mathbf{y}}$$

are contained in $K$. Since the triangulation is assumed to be quasi-uniform, one obtains for all $\hat{\mathbf{y}}$

$$\|B\hat{\mathbf{y}}\|_2 = \|\mathbf{x} - \mathbf{x}_0\|_2 \le C_R h_K.$$

Now, it holds for the spectral norm that

$$\|B\|_2 = \sup_{\hat{\mathbf{z}} \ne \mathbf{0}} \frac{\|B\hat{\mathbf{z}}\|_2}{\|\hat{\mathbf{z}}\|_2} = \frac{1}{r} \sup_{\|\hat{\mathbf{z}}\|_2 = r} \|B\hat{\mathbf{z}}\|_2 \le \frac{C_R}{r} h_K.$$

An estimate of this form, with a possible different constant, holds also for all other matrix norms since all matrix norms are equivalent.

The estimate for $\|B^{-1}\|$ proceeds in the same way with interchanging the roles of $K$ and $\hat{K}$. $\blacksquare$

**Theorem 6.14 Local interpolation estimate.** *Let an affine family of finite elements be given by its reference cell $\hat{K}$, the functionals $\{\hat{\Phi}_i\}$, and a space of polynomials $\hat{P}(\hat{K})$. Let all assumptions of Theorem 6.10 be satisfied. Then, for all $v \in W^{m+1,p}(K)$ there is a constant $C$, which is independent of $v(\mathbf{x})$ such that*

$$\left\| D^k(v - I_K v) \right\|_{L^p(K)} \le C h_K^{m+1-k} \left\| D^{m+1} v \right\|_{L^p(K)}, \quad k \le m+1. \tag{6.5}$$

**Proof:** The idea of the proof consists in transforming left hand side of (6.5) to the reference cell, using the interpolation estimate on the reference cell and transforming back.

*i).* Denote the elements of the matrices $B$ and $B^{-1}$ by $b_{ij}$ and $b_{ij}^{(-1)}$, respectively. Since $\|B\|_\infty = \max_{i,j} |b_{ij}|$ is also a matrix norm, it holds that

$$|b_{ij}| \le C h_K, \quad \left| b_{ij}^{(-1)} \right| \le C h_K^{-1}. \tag{6.6}$$

Using element-wise estimates for the matrix $B$ (Leibniz formula for determinants), one obtains

$$|\det B| \le C h_K^d, \quad \left| \det B^{-1} \right| \le C h_K^{-d}. \tag{6.7}$$

*ii).* The next step consists in proving that the transformed interpolation operator is equal to the natural interpolation operator on $K$. The latter one is given by

$$I_K v = \sum_{i=1}^N \Phi_{K,i}(v) \phi_{K,i}, \tag{6.8}$$

where $\{\phi_{K,i}\}$ is the basis of the space

$$P(K) = \{ p \ : \ K \to \mathbb{R} \ : \ p = \hat{p} \circ F_K^{-1}, \hat{p} \in \hat{P}(\hat{K}) \},$$

which satisfies $\Phi_{K,i}(\phi_{K,j}) = \delta_{ij}$. The functionals are defined by

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

Hence, it follows with $v = \hat{\phi}_j \circ F_K^{-1}$ from the condition on the local basis on $\hat{K}$ that

$$\Phi_{K,i}(\hat{\phi}_j \circ F_K^{-1}) = \hat{\Phi}_i(\hat{\phi}_j) = \delta_{ij},$$

i.e., the local basis on $K$ is given by $\phi_{K,j} = \hat{\phi}_j \circ F_K^{-1}$. Using (6.8), one gets

$$
\begin{aligned}
I_{\hat{K}} \hat{v} &= \sum_{i=1}^N \hat{\Phi}_i(\hat{v}) \hat{\phi}_i = \sum_{i=1}^N \Phi_{K,i}(\underbrace{\hat{v} \circ F_K^{-1}}_{=v}) \phi_{K,i} \circ F_K = \left( \sum_{i=1}^N \Phi_{K,i}(v) \phi_{K,i} \right) \circ F_K \\
&= I_K v \circ F_K.
\end{aligned}
$$

Hence, $I_{\hat{K}} \hat{v}$ is transformed correctly.

*iii).* One obtains with the chain rule

$$\frac{\partial v(\mathbf{x})}{\partial \mathbf{x}_i} = \sum_{j=1}^d \frac{\partial \hat{v}(\hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}_j} b_{ji}^{(-1)}, \quad \frac{\partial \hat{v}(\hat{\mathbf{x}})}{\partial \hat{\mathbf{x}}_i} = \sum_{j=1}^d \frac{\partial v(\mathbf{x})}{\partial \mathbf{x}_j} b_{ji}.$$

It follows with (6.6) that (with each derivative one obtains an additional factor of $B$ or $B^{-1}$, respectively)

$$\left\| D_{\mathbf{x}}^k v(\mathbf{x}) \right\|_2 \le C h_K^{-k} \left\| D_{\hat{\mathbf{x}}}^k \hat{v}(\hat{\mathbf{x}}) \right\|_2, \quad \left\| D_{\hat{\mathbf{x}}}^k \hat{v}(\hat{\mathbf{x}}) \right\|_2 \le C h_K^k \left\| D_{\mathbf{x}}^k v(\mathbf{x}) \right\|_2.$$

One gets with (6.7)

$$\int_K \left\| D_{\mathbf{x}}^k v(\mathbf{x}) \right\|_2^p \, d\mathbf{x} \le C h_K^{-kp} |\det B| \int_{\hat{K}} \left\| D_{\hat{\mathbf{x}}}^k \hat{v}(\hat{\mathbf{x}}) \right\|_2^p \, d\hat{\mathbf{x}} \le C h_K^{-kp+d} \int_{\hat{K}} \left\| D_{\hat{\mathbf{x}}}^k \hat{v}(\hat{\mathbf{x}}) \right\|_2^p \, d\hat{\mathbf{x}}$$

and

$$\int_{\hat{K}} \left\| D_{\hat{\mathbf{x}}}^k \hat{v}(\hat{\mathbf{x}}) \right\|_2^p \, d\hat{\mathbf{x}} \le C h_K^{kp} \left| \det B^{-1} \right| \int_K \left\| D_{\mathbf{x}}^k v(\mathbf{x}) \right\|_2^p \, d\mathbf{x} \le C h_K^{kp-d} \int_K \left\| D_{\mathbf{x}}^k v(\mathbf{x}) \right\|_2^p \, d\mathbf{x}.$$

Using now the interpolation estimate on the reference cell (6.4) yields

$$\left\| D_{\hat{\mathbf{x}}}^{k}(\hat{v} - I_{\hat{K}}\hat{v}) \right\|_{L^{p}(\hat{K})}^{p} \leq C \left\| D_{\hat{\mathbf{x}}}^{m+1}\hat{v} \right\|_{L^{p}(\hat{K})}^{p}, \quad 0 \leq k \leq m+1.$$

It follows that

$$
\begin{aligned}
\left\| D_{\mathbf{x}}^{k}(v - I_{K}v) \right\|_{L^{p}(K)}^{p} &\leq Ch_{K}^{-kp+d} \left\| D_{\hat{\mathbf{x}}}^{k}(\hat{v} - I_{\hat{K}}\hat{v}) \right\|_{L^{p}(\hat{K})}^{p} \\
&\leq Ch_{K}^{-kp+d} \left\| D_{\hat{\mathbf{x}}}^{m+1}\hat{v} \right\|_{L^{p}(\hat{K})}^{p} \\
&\leq Ch_{K}^{(m+1-k)p} \left\| D_{\mathbf{x}}^{m+1}v \right\|_{L^{p}(K)}^{p}.
\end{aligned}
$$

Taking the $p$-th root proves the statement of the theorem. ∎

**Remark 6.15** *On estimate* (6.5).
- Note that the power of $h_{K}$ does not depend on $p$ and $d$.
- Consider a quasi-uniform triangulation and define

$$h = \max_{K \in \mathcal{T}^{h}} \{h_{K}\}.$$

Then, one obtains by summing over all mesh cells an interpolation estimate for the global finite element space

$$
\begin{aligned}
\left\| D^{k}(v - I_{h}v) \right\|_{L^{p}(\Omega)} &= \left( \sum_{K \in \mathcal{T}^{h}} \left\| D^{k}(v - I_{K}v) \right\|_{L^{p}(K)}^{p} \right)^{1/p} \\
&\leq \left( \sum_{K \in \mathcal{T}^{h}} ch_{K}^{(m+1-k)p} \left\| D^{m+1}v \right\|_{L^{p}(K)}^{p} \right)^{1/p} \\
&\leq ch^{(m+1-k)} \left\| D^{m+1}v \right\|_{L^{p}(\Omega)}. \tag{6.9}
\end{aligned}
$$

For linear finite elements $P_{1}$ $(m = 1)$ it is, in particular,

$$\left\| v - I_{h}v \right\|_{L^{p}(\Omega)} \leq ch^{2} \left\| D^{2}v \right\|_{L^{p}(\Omega)}, \quad \left\| \nabla(v - I_{h}v) \right\|_{L^{p}(\Omega)} \leq ch \left\| D^{2}v \right\|_{L^{p}(\Omega)},$$

if $v \in W^{2,p}(\Omega)$. □

**Corollary 6.16 Finite element error estimate.** *Let $u(\mathbf{x})$ be the solution of the model problem (4.9) with $u \in H^{m+1}(\Omega)$ and let $u^{h}(\mathbf{x})$ be the solution of the corresponding finite element problem. Consider a family of quasi-uniform triangulations and let the finite element spaces $V^{h}$ contain polynomials of degree $m$. Then, the following finite element error estimate holds*

$$\left\| \nabla(u - u^{h}) \right\|_{L^{2}(\Omega)} \leq ch^{m} \left\| D^{m+1}u \right\|_{L^{2}(\Omega)} = ch^{m} |u|_{H^{m+1}(\Omega)}. \tag{6.10}$$

**Proof:** The statement follows by combining Lemma 4.13 (for $V = H_{0}^{1}(\Omega)$) and (6.9)

$$\left\| \nabla(u - u^{h}) \right\|_{L^{2}(\Omega)} \leq \inf_{v^{h} \in V^{h}} \left\| \nabla(u - v^{h}) \right\|_{L^{2}(\Omega)} \leq \left\| \nabla(u - I_{h}u) \right\|_{L^{2}(\Omega)} \leq ch^{m} |u|_{H^{m+1}(\Omega)}.$$

∎

**Remark 6.17** *To* (6.10). Note that Lemma 4.13 provides only information about the error in the norm on the left-hand side of (6.10), but not in other norms. □

## 6.2 Inverse Estimate

**Remark 6.18** *On inverse estimates.* The approach for proving interpolation error estimates can be uses also to prove so-called inverse estimates. In contrast to interpolation error estimates, a norm of a higher order derivative of a finite element function will be estimated by a norm of a lower order derivative of this function. One obtains as penalty a factor with negative powers of the diameter of the mesh cell. □

**Theorem 6.19 Inverse estimate.** *Let $0 \leq k \leq l$ be natural numbers and let $p, q \in [1, \infty]$. Then there is a constant $C_{\mathrm{inv}}$, which depends only on $k, l, p, q, \hat{K}, \hat{P}(\hat{K})$ such that*

$$\left\| D^l v^h \right\|_{L^q(K)} \leq C_{\mathrm{inv}} h_K^{(k-l) - d(p^{-1} - q^{-1})} \left\| D^k v^h \right\|_{L^p(K)} \quad \forall \, v^h \in P(K). \quad (6.11)$$

**Proof:** In the first step, (6.11) is shown for $h_{\hat{K}} = 1$ and $k = 0$ on the reference mesh cell. Since all norms are equivalent in finite dimensional spaces, one obtains

$$\left\| D^l \hat{v}^h \right\|_{L^q(\hat{K})} \leq \left\| \hat{v}^h \right\|_{W^{l,q}(\hat{K})} \leq C \left\| \hat{v}^h \right\|_{L^p(\hat{K})} \quad \forall \, \hat{v}^h \in \hat{P}(\hat{K}).$$

If $k > 0$, then one sets

$$\tilde{P}(\hat{K}) = \left\{ \partial_{\boldsymbol{\alpha}} \hat{v}^h \; : \; \hat{v}^h \in \hat{P}(\hat{K}), |\boldsymbol{\alpha}| = k \right\},$$

which is also a space consisting of polynomials. The application of the first estimate of the proof to $\tilde{P}(\hat{K})$ gives

$$
\begin{aligned}
\left\| D^l \hat{v}^h \right\|_{L^q(\hat{K})} &= \sum_{|\boldsymbol{\alpha}|=k} \left\| D^{l-k} \left( \partial_{\boldsymbol{\alpha}} \hat{v}^h \right) \right\|_{L^q(\hat{K})} \leq C \sum_{|\boldsymbol{\alpha}|=k} \left\| \partial_{\boldsymbol{\alpha}} \hat{v}^h \right\|_{L^p(\hat{K})} \\
&= C \left\| D^k \hat{v}^h \right\|_{L^p(\hat{K})}.
\end{aligned}
$$

This estimate is transformed to an arbitrary mesh cell $K$ analogously as for the interpolation error estimates. From the estimates for the transformations, one obtains

$$
\begin{aligned}
\left\| D^l v^h \right\|_{L^q(K)} &\leq C h_K^{-l + d/q} \left\| D^l \hat{v}^h \right\|_{L^q(\hat{K})} \leq C h_K^{-l+d/q} \left\| D^k \hat{v}^h \right\|_{L^p(\hat{K})} \\
&\leq C_{\mathrm{inv}} h_K^{k-l+d/q-d/p} \left\| D^k v^h \right\|_{L^p(K)}.
\end{aligned}
$$

■

**Remark 6.20** *On the proof.* The crucial point in the proof was the equivalence of all norms in finite dimensional spaces. Such a property does not exist in infinite dimensional spaces. □

**Corollary 6.21 Global inverse estimate.** *Let $p = q$ and let $\mathcal{T}^h$ be a regular triangulation of $\Omega$, then*

$$\left\| D^l v^h \right\|_{L^{p,h}(\Omega)} \leq C_{\mathrm{inv}} h^{k-l} \left\| D^k v^h \right\|_{L^{p,h}(\Omega)},$$

*where*

$$\|\cdot\|_{L^{p,h}(\Omega)} = \left( \sum_{K \in \mathcal{T}^h} \|\cdot\|_{L^p(K)}^p \right)^{1/p}.$$

**Remark 6.22** *On $\|\cdot\|_{L^{p,h}(\Omega)}$.* The cell wise definition of the norm is important for $l \geq 2$ since in this case finite element functions generally do not possess the regularity for the global norm to be well defined. It is also important for $l \geq 1$ and non-conforming finite element functions. □

## 6.3  Interpolation of Non-Smooth Functions

**Remark 6.23** *Motivation.* The interpolation theory of Section 6.1 requires that the interpolation operator is continuous on the Sobolev space to which the function belongs that should be interpolated. But if one, e.g., wants to interpolate discontinuous functions with continuous, piecewise linear elements, then Section 6.1 does not provide estimates.

A simple remedy seems to be first to apply some smoothing operator to the function to be interpolated and then to interpolate the smoothed function. However, this approach leads to difficulties at the boundary of $\Omega$ and it will not be considered further.

There are two often used interpolation operators for non-smooth functions. The interpolation operator of Clément (1975) is defined for functions from $L^1(\Omega)$ and it can be generalized to more or less all finite elements. The interpolation operator of Scott and Zhang (1990) is more special. It has the advantage that it preserves homogeneous Dirichlet boundary conditions. Here, only the interpolation operator of Clément, for linear finite elements, will be considered.

Let $\mathcal{T}^h$ be a regular triangulation of the polyhedral domain $\Omega \subset \mathbb{R}^d, d \in \{2,3\}$, with simplices $K$. Denote by $P_1$ the space of continuous, piecewise linear finite elements on $\mathcal{T}^h$. □

**Remark 6.24** *Construction of the interpolation Operator of Clément.* For each vertex $V_i$ of the triangulation, the union of all grid cells which possess $V_i$ as vertex will be denoted by $\omega_i$, see Figure 5.1.

The interpolation operator of Clément is defined with the help of local $L^2(\omega_i)$ projections. Let $v \in L^1(\Omega)$ and let $P_1(\omega_i)$ be the space of continuous piecewise linear finite elements on $\omega_i$. Then, the local $L^2(\omega_i)$ projection of $v \in L^1(\omega_i)$ is the solution $p_i \in P_1(\omega_i)$ of

$$\int_{\omega_i} (v - p_i)(\mathbf{x}) q(\mathbf{x})\ d\mathbf{x} = 0 \quad \forall\ q \in P_1(\omega_i) \tag{6.12}$$

or equivalently of

$$(v - p_i, q)_{L^2(\omega_i)} = 0 \quad \forall\ q \in P_1(\omega_i).$$

Then, the Clément interpolation operator is defined by

$$P_{\mathrm{Cle}}^h v(\mathbf{x}) = \sum_{i=1}^N p_i(V_i)\phi_i^h(\mathbf{x}), \tag{6.13}$$

where $\{\phi_i^h\}_{i=1}^N$ is the standard basis of the global finite element space $P_1$. Since $P_{\mathrm{Cle}}^h v(\mathbf{x})$ is a linear combination of basis functions of $P_1$, it defines a map $P_{\mathrm{Cle}}^h : L^1(\Omega) \to P_1$. □

**Theorem 6.25 Interpolation estimate.** *Let $k, l \in \mathbb{N} \cup \{0\}$ and $q \in \mathbb{R}$ with $k \le l \le 2$, $1 \le q \le \infty$ and let $\omega_K$ be the union of all subdomains $\omega_i$ that contain the mesh cell $K$, see Figure 6.2. Then it holds for all $v \in W^{l,q}(\omega_K)$ the estimate*

$$\left\| D^k(v - P_{\mathrm{Cle}}^h v) \right\|_{L^q(K)} \le Ch^{l-k} \left\| D^l v \right\|_{L^q(\omega_K)}, \tag{6.14}$$

*with $h = diam(\omega_K)$, where the constant $C$ is independent of $v(\mathbf{x})$ and $h$.*

**Proof:** The statement of the lemma is obvious in the case $k = l = 2$ since it is $D^2 P_{\mathrm{Cle}}^h v(\mathbf{x})|_K = 0$.
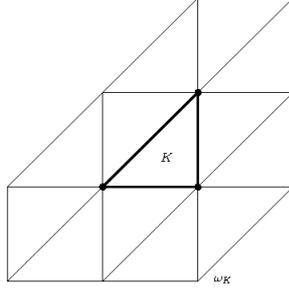
Figure 6.2: A subdomain $\omega_K$.

Let $k \in \{0, 1\}$. Because the $L^2(\Omega)$ projection gives an element with best approximation, one gets with (6.12)

$$P_{\text{Cle}}^h p = p \quad \text{in } K \quad \forall\, p \in P_1(\omega_K). \tag{6.15}$$

One says that $P_{\text{Cle}}^h$ is a consistent operator.

The next step consists in proving the stability of $P_{\text{Cle}}^h$. One obtains with the inverse inequality (6.11)

$$\|p\|_{L^\infty(\omega_i)} \le ch^{-d/2} \|p\|_{L^2(\omega_i)} \quad \text{for all } p \in P_1(\omega_i).$$

The inverse inequality and definition (6.12) of the local $L^2$ projection with the test function $q = p_i$ gives

$$\|p_i\|_{L^\infty(\omega_i)}^2 \le ch^{-d} \|p_i\|_{L^2(\omega_i)}^2 \le ch^{-d} \|v\|_{L^1(\omega_i)} \|p_i\|_{L^\infty(\omega_i)}.$$

Dividing by $\|p_i\|_{L^\infty(\omega_i)}$ and applying Hölder's inequality, one obtains for $p^{-1} = 1 - q^{-1}$ (*exercise*)

$$|p_i(V_i)| \quad \le \quad ch^{-d/q} \|v\|_{L^q(\omega_i)} \tag{6.16}$$

for all $V_i \in K$. From the regularity of the triangulation, it follows for the basis functions that (inverse estimate)

$$\left\| D^k \phi_i \right\|_{L^\infty(K)} \le ch^{-k}, \quad k = 0, 1. \tag{6.17}$$

Using the triangle inequality, combining (6.16) and (6.17) yields the stability of $P_{\text{Cle}}^h$

$$
\begin{aligned}
\left\| D^k P_{\text{Cle}}^h v \right\|_{L^q(K)} &\le& \sum_{V_i \in K} |p_i(V_i)| \left\| D^k \phi_i \right\|_{L^q(K)} \\
&\le& c \sum_{V_i \in K} h^{-d/q} \|v\|_{L^q(\omega_i)} \left\| D^k \phi_i \right\|_{L^\infty(K)} \|1\|_{L^q(K)} \\
&\le& c \sum_{V_i \in K} h^{-d/q} \|v\|_{L^q(\omega_i)} h^{-k} h^{d/q} \\
&=& ch^{-k} \|v\|_{L^q(\omega_K)}.
\end{aligned}
\tag{6.18}
$$

The remainder of the proof follows the proof of the interpolation error estimate for the polynomial interpolation, Theorem 6.10, apart from the fact that a reference cell is not used for the Clément interpolation operator. Using Lemma 6.2 and 6.4, one can find a polynomial $p \in P_1(\omega_K)$ with

$$\left\| D^j(v - p) \right\|_{L^q(\omega_K)} \le ch^{l-j} \left\| D^l v \right\|_{L^q(\omega_K)}, \quad 0 \le j \le l \le 2. \tag{6.19}$$

With (6.15), the triangle inequality, $\|\cdot\|_{L^q(K)} \le \|\cdot\|_{L^q(\omega_K)}$, (6.18), and (6.19), one obtains

$$
\begin{aligned}
\left\| D^k \left( v - P_{\mathrm{Cle}}^h v \right) \right\|_{L^q(K)} &= \left\| D^k \left( v - p + P_{\mathrm{Cle}}^h p - P_{\mathrm{Cle}}^h v \right) \right\|_{L^q(K)} \\
&\le \left\| D^k (v - p) \right\|_{L^q(K)} + \left\| D^k P_{\mathrm{Cle}}^h (v - p) \right\|_{L^q(K)} \\
&\le \left\| D^k (v - p) \right\|_{L^q(\omega_K)} + c h^{-k} \left\| v - p \right\|_{L^q(\omega_K)} \\
&\le c h^{l-k} \left\| D^l v \right\|_{L^q(\omega_K)} + c h^{-k} h^l \left\| D^l v \right\|_{L^q(\omega_K)} \\
&= c h^{l-k} \left\| D^l v \right\|_{L^q(\omega_K)}.
\end{aligned}
$$

∎

**Remark 6.26** *Uniform meshes.*
- If all mesh cells in $\omega_K$ are of the same size, then one can replace $h$ by $h_K$ in the interpolation error estimate (6.14). This property is given in many cases.
- If one assumes that the number of mesh cells in $\omega_K$ is bounded uniformly for all considered triangulations, the global interpolation estimate

$$
\left\| D^k (v - P_{\mathrm{Cle}}^h v) \right\|_{L^q(\Omega)} \le C h^{l-k} \left\| D^l v \right\|_{L^q(\Omega)}, \quad 0 \le k \le l \le 2,
$$

follows directly from (6.14).

□

# Chapter 7

# Finite Element Methods for Second Order Elliptic Equations

## 7.1 General Convergence Theorems

**Remark 7.1** *Motivation.* In Section 5.1, non-conforming finite element methods have been introduced, i.e., methods where the finite element space $V^h$ is not a subspace of $V$, which is the space in the definition of the continuous variational problem. The property $V^h \not\subset V$ is given for the Crouzeix–Raviart and the Rannacher–Turek element. Another case of non-conformity is given if the domain does not possess a polyhedral boundary and one has to apply some approximation of the boundary.

For non-conforming methods, the finite element approach is not longer a Ritz method. Hence, the convergence proof from Theorem 4.14 cannot be applied in this case. The abstract convergence theorem, which will be proved in this section, allows the numerical analysis of complex finite element methods. □

**Remark 7.2** *Notations, Assumptions.* Let $\{h > 0\}$ be a set of mesh widths and let $S^h, V^h$ normed spaces of functions which are defined on domains $\{\Omega^h \subset \mathbb{R}^d\}$. It will be assumed that the space $S^h$ has a finite dimension and that $S^h$ and $V^h$ possess a common norm $\|\cdot\|_h$. In the application of the abstract theory, $S^h$ will be a finite element space and $V^h$ is defined such that the restriction and/or extension of the solution of the continuous problem to $\Omega^h$ is contained in $V^h$. Strictly speaking, the modified solution of the continuous problem does not solve the given problem any longer. Hence, it is consequent that the continuous problem does not appear explicitly in the abstract theory.

Given the bilinear forms

$$
\begin{aligned}
a^h &: \quad S^h \times S^h \to \mathbb{R}, \\
\tilde{a}^h &: \quad (S^h + V^h) \times (S^h + V^h) \to \mathbb{R}.
\end{aligned}
$$

Let the bilinear form $a^h$ be regular in the sense that there is a constant $m > 0$, which is independent of $h$, such that for each $v^h \in S^h$ there is a $w^h \in S^h$ with $\|w^h\|_h = 1$ such that

$$
m \|v^h\|_h \leq a^h(v^h, w^h). \tag{7.1}
$$

This condition is equivalent to the requirement that the stiffness matrix $A$ with the entries $a_{ij} = a^h(\phi_j, \phi_i)$, where $\{\phi_i\}$ is a basis of $S^h$, is uniformly non-singular, i.e.,

its regularity is independent of $h$. For the second bilinear form, only its boundedness will be assumed

$$\tilde{a}^h(u,v) \le M \left\| u \right\|_h \left\| v \right\|_h \quad \forall \ u, v \in S^h + V^h. \tag{7.2}$$

Let the linear functionals $\{f^h(\cdot)\} : S^h \to \mathbb{R}$ be given. Then, the following discrete problems will be considered: Find $u^h \in S^h$ with

$$a^h(u^h, v^h) = f^h(v^h) \quad \forall \, v^h \in S^h. \tag{7.3}$$

Because the stiffness matrix is assumed to be non-singular, there is a unique solution of (7.3). $\qquad\square$

**Theorem 7.3 Abstract finite element error estimate.** *Let the conditions (7.1) and (7.2) be satisfied and let $u^h$ be the solution of (7.3). Then the following error estimate holds for each $\tilde{u} \in V^h$*

$$
\begin{aligned}
\left\| \tilde{u} - u^h \right\|_h \quad &\le \quad c \inf_{v^h \in S^h} \left\{ \left\| \tilde{u} - v^h \right\|_h + \sup_{w^h \in S^h} \frac{\left| \tilde{a}^h(v^h, w^h) - a^h(v^h, w^h) \right|}{\left\| w^h \right\|_h} \right\} \\
&\quad + c \sup_{w^h \in S^h} \frac{\left| \tilde{a}^h(\tilde{u}, w^h) - f^h(w^h) \right|}{\left\| w^h \right\|_h}
\end{aligned}
\tag{7.4}
$$

*with $c = c(m, M)$.*

**Proof:** Because of (7.1) there is for each $v^h \in S^h$ a $w^h \in S^h$ with $\left\| w^h \right\|_h = 1$ and

$$m \left\| u^h - v^h \right\|_h \le a^h(u^h - v^h, w^h).$$

Using the definition of $u^h$ from (7.3), one obtains

$$m \left\| u^h - v^h \right\|_h \le f^h(w^h) - a^h(v^h, w^h) + \tilde{a}^h(v^h, w^h) + \tilde{a}^h(\tilde{u} - v^h, w^h) - \tilde{a}^h(\tilde{u}, w^h).$$

From (7.2) and $\left\| w^h \right\|_h = 1$ it follows that

$$\tilde{a}^h(\tilde{u} - v^h, w^h) \le M \left\| \tilde{u} - v^h \right\|_h.$$

Rearranging the terms appropriately and using $\left\| w^h / \left\| w^h \right\|_h \right\|_h = 1$ gives

$$
\begin{aligned}
m \left\| u^h - v^h \right\|_h \quad &\le \quad M \left\| \tilde{u} - v^h \right\|_h + \sup_{w^h \in S^h} \frac{\left| \tilde{a}^h(v^h, w^h) - a^h(v^h, w^h) \right|}{\left\| w^h \right\|_h} \\
&\quad + \sup_{w^h \in S^h} \frac{\left| \tilde{a}^h(\tilde{u}, w^h) - f^h(w^h) \right|}{\left\| w^h \right\|_h}.
\end{aligned}
$$

Applying the triangle inequality

$$\left\| \tilde{u} - u^h \right\|_h \le \left\| \tilde{u} - v^h \right\|_h + \left\| u^h - v^h \right\|_h$$

and inserting the estimate from above gives (7.4). $\qquad\blacksquare$

**Remark 7.4** *To Theorem 7.3.* An important special case of this theorem is the case that the stiffness matrix is uniformly positive definite, i.e., the condition

$$m \left\| v^h \right\|_h^2 \le a^h(v^h, v^h) \quad \forall \, v^h \in S^h \tag{7.5}$$

is satisfied. Dividing (7.5) by $\left\| v^h \right\|_h$ reveals that condition (7.1) is implied by (7.5).

If the continuous problem is also defined with the bilinear form $\tilde{a}^h(\cdot, \cdot)$, then

$$\sup_{w^h \in S^h} \frac{\left|\tilde{a}^h(v^h, w^h) - a^h(v^h, w^h)\right|}{\|w^h\|_h}$$

can be considered as consistency error of the bilinear forms and the term

$$\sup_{w^h \in S^h} \frac{\left|\tilde{a}^h(\tilde{u}, w^h) - f^h(w^h)\right|}{\|w^h\|_h}$$

as consistency error of the right-hand sides. $\qquad\qquad\qquad\square$

**Theorem 7.5 First Strang lemma** *Let $S^h$ be a conform finite element space, i.e., $S^h \subset V$, with $\|\cdot\|_h = \|\cdot\|_V$ and let the space $V^h$ be independent of $h$. Consider a continuous problem of the form*

$$\tilde{a}^h(u, v) = f(v) \quad \forall\, v \in V,$$

*then the following error estimate holds.*

$$
\begin{aligned}
\left\|u - u^h\right\|_V &\leq c \inf_{v^h \in S^h} \left\{ \left\|u - v^h\right\|_V + \sup_{w^h \in S^h} \frac{\left|\tilde{a}^h(v^h, w^h) - a^h(v^h, w^h)\right|}{\|w^h\|_V} \right\} \\
&\quad + c \sup_{w^h \in S^h} \frac{\left|f(w^h) - f^h(w^h)\right|}{\|w^h\|_V}.
\end{aligned}
$$

**Proof:** The statement of this theorem follows directly from Theorem 7.3. $\qquad\blacksquare$

## 7.2 Linear Finite Element Method on Non-Polyhedral Domains

**Remark 7.6** *The continuous problem.* The abstract theory will be applied to the linear finite element method for the solution of second order elliptic partial differential equations.

Let $\Omega \subset \mathbb{R}^d$, $d \in \{2, 3\}$, be a bounded domain with Lipschitz boundary, which does not need to be polyhedral. Let

$$Lu = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \partial\Omega, \tag{7.6}$$

where the operator is given by

$$Lu = -\nabla \cdot (A\nabla u)$$

with

$$A(\mathbf{x}) = (a_{ij}(\mathbf{x}))_{i,j=1}^d, \quad a_{ij} \in W^{1,p}(\Omega), p > d, \tag{7.7}$$

It will be assumed that there are two positive real numbers $m, M$ such that

$$m \|\boldsymbol{\xi}\|_2^2 \leq \boldsymbol{\xi}^T A(\mathbf{x})\boldsymbol{\xi} \leq M \|\boldsymbol{\xi}\|_2^2 \quad \forall\, \boldsymbol{\xi} \in \mathbb{R}^d, \mathbf{x} \in \overline{\Omega}. \tag{7.8}$$

From the Sobolev inequality it follows that $a_{ij} \in L^\infty(\Omega)$. With

$$a(u, v) = \int_\Omega (A(\mathbf{x})\nabla u(\mathbf{x})) \cdot \nabla v(\mathbf{x})\, d\mathbf{x}$$

and the Cauchy–Schwarz inequality, one obtains

$$|a(u, v)| \leq \|A\|_{L^\infty(\Omega)} \int_\Omega |\nabla u(\mathbf{x}) \cdot \nabla v(\mathbf{x})|\ d\mathbf{x} \leq c \|\nabla u\|_{L^2(\Omega)} \|\nabla v\|_{L^2(\Omega)}$$

for all $u, v \in H_0^1(\Omega)$. In addition, it follows that

$$m \left\|\nabla u\right\|_{L^2(\Omega)}^2 \le a(u, u) \quad \forall \, u \in H_0^1(\Omega).$$

Hence, the bilinear form is bounded and elliptic. Using the Theorem of Lax–Milgram, Theorem 4.5, it follows that there es a unique weak solution $u \in H_0^1(\Omega)$ of (7.6) with

$$a(u, v) = f(v) \quad \forall \, v \in H_0^1(\Omega).$$

$\square$

**Remark 7.7** *The finite element problem.* Let $\mathcal{T}^h$ be a regular triangulation consisting of simplices $\{K\}$ such that the vertices of the simplices belong to $\overline{\Omega}$, see Figure 7.1, and define $\overline{\Omega^h} = \cup_{K \in \mathcal{T}^h} K$.



Figure 7.1: Approximation of the boundary by the finite element mesh.

The space of continuous and piecewise linear functions that vanish at the boundary of $\Omega^h$ will be denoted by $P_1$. It will be assumed that for the data of the problem $a_{ij}(\mathbf{x}), f(\mathbf{x})$ there exist extensions $\tilde{a}_{ij}(\mathbf{x}), \tilde{f}(\mathbf{x})$ to a larger domain $\tilde{\Omega} \supset \overline{\Omega^h}$ with

$$\left\|\tilde{a}_{ij}\right\|_{W^{1,p}(\tilde{\Omega})} \le c \left\|a_{ij}\right\|_{W^{1,p}(\Omega)}, \quad \left\|\tilde{f}\right\|_{L^2(\tilde{\Omega})} \le c \left\|f\right\|_{L^2(\Omega)}. \tag{7.9}$$

In addition, it will be assumed that the coefficients $\tilde{a}_{ij}(\mathbf{x})$ satisfy the ellipticity condition (7.8) on $\tilde{\Omega}$.

Obviously, $f(\mathbf{x})$ can be continued simply by zero. The extensions of $a_{ij}(\mathbf{x})$ have to be weakly differentiable. It is possible to show that such extensions exist, see the literature.

The finite element method is defined as follows: Find $u^h \in P_1$ with

$$a^h(u^h, v^h) = f^h(v^h) \quad \forall \, v^h \in P_1,$$

where

$$a^h(u^h, v^h) = \int_{\Omega^h} \left(\tilde{A}(\mathbf{x}) \nabla u^h(\mathbf{x})\right) \cdot \nabla v^h(\mathbf{x}) \, d\mathbf{x}, \quad f^h(v^h) = \int_{\Omega^h} \tilde{f}(\mathbf{x}) v^h(\mathbf{x}) \, d\mathbf{x}.$$

In practice, it might be hard to apply the method in this form. From the existence of the extension operators for $a_{ij}(\mathbf{x})$ it is not yet clear how to compute them. On the other hand, in practice often the coefficients $a_{ij}(\mathbf{x})$ are constant or at least piecewise constant. In these case, the extension is trivial. As remedy in the general case, one can use quadrature rules whose nodes are situated within $\overline{\Omega}$, see the literature. $\square$

**Remark 7.8** *Goal of the analysis, further assumptions.* The goal consists in proving the linear convergence of the finite element method in $\left\|\cdot\right\|_h = \left\|\cdot\right\|_{H^1(\Omega^h)}$. In the analysis, one has to pay attention to the fact that in general neither holds $\Omega^h \subset \Omega$ nor $\Omega \subset \Omega^h$. It will be assumed that there is an extension $\tilde{u} \in H^2(\tilde{\Omega})$ of $u(\mathbf{x})$ with

$$\left\|\tilde{u}\right\|_{H^2(\tilde{\Omega})} \le c \left\|u\right\|_{H^2(\Omega)}. \tag{7.10}$$

In addition, it will be assumed that $\Omega^h$ is a sufficiently good approximation of $\Omega$ in the following sense

$$\max_{\mathbf{x} \in \partial\Omega^h} \operatorname{dist}(\mathbf{x}, \partial\Omega) \le ch^2. \tag{7.11}$$

One can show that (7.11) is satisfied for $d = 2$ if the boundary of $\Omega$ is piecewise $C^2$ and the corners of $\Omega$ are vertices of the triangulation. In this case, one can rotate the coordinate system locally such that the distance between $\partial\Omega$ and $\partial\Omega^h$ can be represented as the error of a one-dimensional interpolation problem with continuous, piecewise linear finite elements. Using error estimates for this kind of problem, e.g., see Goering et al. (2010), one can estimate the error by $ch^2$. For three-dimensional domains, with piecewise $C^\infty$ boundary, one needs in addition a smoothness assumption for the edges. $\square$

**Lemma 7.9 Estimate of a function on the difference of the domains.** *Let the condition (7.11) be satisfied. Then, for all $v \in W^{1,1}(\Omega)$ it holds the estimate*

$$\int_{\Omega_s} |v(\mathbf{x})| \ d\mathbf{x} \le ch^2 \int_{\Omega} (|v(\mathbf{x})| + \|\nabla v(\mathbf{x})\|_2) \ d\mathbf{x}, \tag{7.12}$$

*where $\Omega_s$ is the set $\Omega \setminus \Omega^h$ or $\Omega^h \setminus \Omega$.*

**Proof:** At the beginning, a one-dimensional estimate will be shown. Let $f \in C^1([0,1])$, then one obtains with the fundamental theorem of calculus

$$f(x) = \int_y^x f'(\xi) \ d\xi + f(y) \quad \forall \, x, y \in [0,1].$$

It follows that

$$|f(x)| \le \int_0^1 |f'(\xi)| \ d\xi + |f(y)| \,.$$

Integrating this inequality with respect to $y$ in $[0,1]$ and with respect to $x$ in $[0,a]$ with $a \in (0,1]$ yields

$$\int_0^a |f(x)| \ dx \le a \int_0^1 |f'(\xi)| \ d\xi + a \int_0^1 |f(y)| \ dy = a \int_0^1 \left( |f(x)| + |f'(x)| \right) \ dx. \tag{7.13}$$

Consider the case $\Omega_s = \Omega \setminus \Omega^h$. Since $\Omega$ has a Lipschitz boundary, it can be shown that $\partial\Omega$ can be covered with a finite number of open sets $U_1, \ldots, U_N$. After a rotation of the coordinate system, one can represent $\partial\Omega \cap U_i$ as a Lipschitz continuous function $g_i(\mathbf{y}')$ of $(d-1)$ arguments $\mathbf{y}' = (y_1, \ldots, y_{d-1}) \in U_i' \subset \mathbb{R}^{d-1}$.

In the next step of the proof, sets will be constructed whose union covers the difference $\Omega \setminus \Omega^h$. Let

$$S_{i,\sigma} = \left\{ (\mathbf{y}', y_d) \ : \ g_i(\mathbf{y}') - \sigma < y_d < g_i(\mathbf{y}'), \ \mathbf{y}' \in U_i' \right\}, \quad i = 1, \ldots, N,$$

see Figure 7.2. Then, using (7.11) it is $(\Omega \setminus \Omega^h) \cap U_i \subset S_{i,c_1 h^2}$, where $c_1$ depends on $g_i(\mathbf{y}')$ but not on $h$. In addition, there is a $\sigma_0$ such that $S_{i,\sigma_0} \subset \Omega$ for all $i$.

The transform of (7.13) to the interval $[0, \sigma_0]$ gives for sufficiently small $h$, such that $c_1 h^2 \le 1$,

$$\int_0^{c_1 h^2} |f(x)| \ dx \le ch^2 \int_0^{\sigma_0} \left( |f(x)| + |f'(x)| \right) \ dx.$$

For $v \in C^1(\overline{\Omega})$, one applies this estimate to the rotated function $v(\mathbf{y}', x)$

$$
\begin{aligned}
\int_{S_{i,c_1 h^2}} |v(\mathbf{y})| \ d\mathbf{y} &= \int_{U_i'} \int_0^{c_1 h^2} |v(\mathbf{y}', x)| \ dx \ d\mathbf{y}' \\
&\le ch^2 \int_{U_i'} \int_0^{\sigma_0} \left( |\partial_x v(\mathbf{y}', x)| + |v(\mathbf{y}', x)| \right) \ dx \ d\mathbf{y}' \\
&\le ch^2 \int_{\Omega} \left( |\partial_{y_d} v(\mathbf{y})| + |v(\mathbf{y})| \right) \ d\mathbf{y},
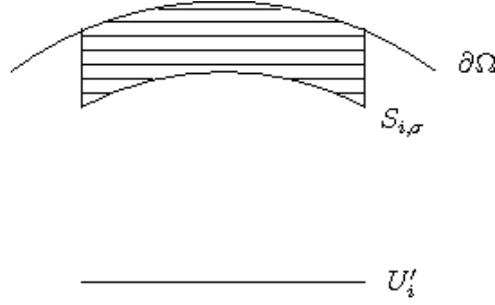\end{aligned}
$$

Figure 7.2: $S_{i,\sigma}$.

where in the first step the theorem of Fubini was used. Taking the sum over $i$ proves the lemma for functions from $C^1(\overline{\Omega})$. Since $C^1(\overline{\Omega})$ is dense in $W^{1,1}(\Omega)$, the statement of the lemma holds also for $v \in W^{1,1}(\Omega)$.

The case $\Omega_s = \Omega^h \setminus \Omega$ is proved analogously. ∎

**Theorem 7.10 Error estimate, linear convergence.** *Let the assumptions* (7.7), *(7.8), (7.9), (7.10), and (7.11) be satisfied. Then, it holds the error estimate*

$$\left\| \tilde{u} - u^h \right\|_{H^1(\Omega^h)} \le ch \left\| u \right\|_{H^2(\Omega)},$$

*where $c$ does not depend on $u$, $f$, and $h$.*

**Proof:** For proving the error estimate, the abstract error estimate, Theorem 7.3, is used with $S^h = P_1$, $V^h = H^1(\Omega^h)$, $\|\cdot\|_h = \|\cdot\|_{H^1(\Omega^h)}$, and

$$a^h(u,v) = \tilde{a}^h(u,v) = \int_{\Omega^h} \left( \tilde{A}(\mathbf{x}) \nabla u(\mathbf{x}) \right) \cdot \nabla v(\mathbf{x}) \, d\mathbf{x}.$$

With this choice of $a^h(\cdot,\cdot)$ and $\tilde{a}^h(\cdot,\cdot)$, the middle term in the abstract error estimate (7.4) vanishes. Setting in the abstract error estimate $v^h = I_h \tilde{u}$, one obtains with the interpolation error estimate (6.5) and (7.10)

$$\left\| \tilde{u} - I_h \tilde{u} \right\|_{H^1(\Omega^h)} \le ch \left\| D^2 \tilde{u} \right\|_{L^2(\Omega^h)} \le ch \left\| u \right\|_{H^2(\Omega)}. \tag{7.14}$$

It remains to estimate the last term of (7.4).

The regularity and the boundedness of $a^h(\cdot, \cdot)$ can be proved easily using the ellipticity and the boundedness of the coefficients $\tilde{a}_{ij}(\mathbf{x})$.

The estimate of the last term of (7.4) starts with integration by parts

$$a^h(\tilde{u}, w^h) = \int_{\Omega^h} \left( \tilde{A}(\mathbf{x}) \nabla \tilde{u}(\mathbf{x}) \right) \cdot \nabla w^h(\mathbf{x}) \, d\mathbf{x} = \int_{\Omega^h} g(\mathbf{x}) w^h(\mathbf{x}) \, d\mathbf{x}$$

with $g(\mathbf{x}) = -\nabla \cdot (\tilde{A} \nabla \tilde{u})(\mathbf{x})$. Because of $g(\mathbf{x}) = \tilde{f}(\mathbf{x}) = f(\mathbf{x})$ in $\Omega$ it is

$$a^h(\tilde{u}, w^h) - f^h(w^h) = \int_{\Omega^h \setminus \Omega} \left( g(\mathbf{x}) - \tilde{f}(\mathbf{x}) \right) w^h(\mathbf{x}) \, d\mathbf{x}.$$

Using the extension of $w^h(\mathbf{x})$ by zero on $\Omega \setminus \Omega^h$, one obtains with (7.12), and noting that in general $\Omega^h \not\subset \Omega$,

$$
\begin{aligned}
\int_{\Omega^h \setminus \Omega} \left| w^h(\mathbf{x}) \right|^2 \, d\mathbf{x} &\leq ch^2 \int_{\Omega} \left( \left\| \nabla w^h(\mathbf{x}) \right\|_2^2 + \left| w^h(\mathbf{x}) \right|^2 \right) \, d\mathbf{x} \\
&\leq ch^2 \int_{\Omega^h} \left( \left\| \nabla w^h(\mathbf{x}) \right\|_2^2 + \left| w^h(\mathbf{x}) \right|^2 \right) \, d\mathbf{x} = ch^2 \left\| w^h \right\|_{H^1(\Omega^h)}^2.
\end{aligned}
$$

Applying the Cauchy–Schwarz inequality and the triangle inequality yields

$$
\begin{aligned}
\left| a^h(\tilde{u}, w^h) - f^h(w^h) \right| &\leq \left\| g - \tilde{f} \right\|_{L^2(\Omega^h \setminus \Omega)} \left\| w^h \right\|_{L^2(\Omega^h \setminus \Omega)} \\
&\leq ch \left( \|g\|_{L^2(\tilde{\Omega})} + \left\| \tilde{f} \right\|_{L^2(\tilde{\Omega})} \right) \left\| w^h \right\|_{H^1(\Omega^h)},
\end{aligned}
$$

where $\tilde{\Omega}$ was introduced in Remark 7.7. Now, a bound for $\|g\|_{L^2(\tilde{\Omega})}$ is needed. Using the product rule and the triangle inequality, one gets

$$\left\| \nabla \cdot (\tilde{A} \nabla \tilde{u}) \right\|_{L^2(\tilde{\Omega})} \leq \left\| \sum_{i,j=1}^d \tilde{a}_{ij} \frac{\partial^2 \tilde{u}}{\partial x_i \partial x_j} \right\|_{L^2(\tilde{\Omega})} + \left\| \left( \nabla \cdot \tilde{A} \right) \cdot \nabla \tilde{u} \right\|_{L^2(\tilde{\Omega})}.$$

Because of the Sobolev imbedding $W^{1,p}(\tilde{\Omega}) \to L^\infty(\tilde{\Omega})$ for $p > d$, Theorem 3.53, it follows that $\left\| \tilde{A} \right\|_{L^\infty(\tilde{\Omega})} \leq c$. One obtains for the first term

$$\left\| \sum_{i,j=1}^d \tilde{a}_{ij} \frac{\partial^2 \tilde{u}}{\partial x_i \partial x_j} \right\|_{L^2(\tilde{\Omega})} \leq c \left\| D^2 \tilde{u} \right\|_{L^2(\tilde{\Omega})}.$$

The estimate of the second term uses Hölders inequality (*exercise*)

$$\left\| \left( \nabla \cdot \tilde{A} \right) \cdot \nabla \tilde{u} \right\|_{L^2(\tilde{\Omega})} \leq \left\| \nabla \cdot \tilde{A} \right\|_{L^p(\tilde{\Omega})}^2 \|\nabla \tilde{u}\|_{L^{2p/(p-2)}(\tilde{\Omega})}^2 \leq c \|\nabla \tilde{u}\|_{L^{2p/(p-2)}(\tilde{\Omega})}^2.$$

Using a Sobolev inequality, e.g., see Adams (1975), one obtains the estimate

$$\|\nabla \tilde{u}\|_{L^{2p/(p-2)}(\tilde{\Omega})} \leq c \|\tilde{u}\|_{H^2(\tilde{\Omega})}.$$

Inserting all estimates, one obtains with (7.9) and (7.10)

$$
\begin{aligned}
\left| a^h(\tilde{u}, w^h) - f^h(w^h) \right| &\leq ch \left( \|\tilde{u}\|_{H^2(\tilde{\Omega})} + \left\| \tilde{f} \right\|_{L^2(\tilde{\Omega})} \right) \left\| w^h \right\|_{H^1(\Omega^h)} \\
&\leq ch \left( \|u\|_{H^2(\Omega)} + \|f\|_{L^2(\Omega)} \right) \left\| w^h \right\|_{H^1(\Omega^h)} \\
&\leq ch \|u\|_{H^2(\Omega)} \left\| w^h \right\|_{H^1(\Omega^h)}.
\end{aligned}
$$

In the final step of this estimate, one uses the representation of $f(\mathbf{x})$ from (7.6), for which one can perform estimates that are analog to the estimates of $g(\mathbf{x})$.

The proof of the linear convergence is finished by using (7.4), (7.14), and the last estimate. ∎

## 7.3 Finite Element Method with the Nonconforming Crouzeix–Raviart Element

**Remark 7.11** *Assumptions and discrete problem.* The nonconforming Crouzeix–Raviart finite element $P_1^{\mathrm{nc}}$ was introduced in Example 5.30. To simplify the presentation, it will be restricted here on the two-dimensional case. In addition, to avoid the estimate of the error coming from approximating the domain, it will be assumed that $\Omega$ is a convex domain with polygonal boundary.

Let $\mathcal{T}^h$ be a regular triangulation of $\Omega$ with triangles. Let $P_1^{\mathrm{nc}}$ (nc – non conforming) be denote the finite element space of piecewise linear functions which are continuous at the midpoints of the edges. This space is nonconforming if it is applied for the discretization of a second order elliptic equation since the continuous problem is given in $H_0^1(\Omega)$ and the functions of $H_0^1(\Omega)$ do not possess jumps. The functions of $P_1^{\mathrm{nc}}$ have generally jumps, see Figure 7.3, and they are not weakly differentiable. In addition, the space is also nonconforming with respect to the boundary condition, which is not satisfied exactly. The functions from $P_1^{\mathrm{nc}}$ vanish in the midpoint of the edges at the boundary. However, in the other points at the boundary, their value is generally not equal to zero.



Figure 7.3: Function from $P_1^{\mathrm{nc}}$.

The bilinear form

$$a(u,v) = \int_\Omega (A(\mathbf{x})\nabla u(\mathbf{x})) \cdot \nabla v(\mathbf{x}) \; d\mathbf{x}$$

will be extended to $H_0^1(\Omega) + P_1^{\mathrm{nc}}$ by

$$a^h(u,v) = \sum_{K \in \mathcal{T}^h} \int_K (A(\mathbf{x})\nabla u(\mathbf{x})) \cdot \nabla v(\mathbf{x}) \; d\mathbf{x} \quad \forall \, u,v \in H_0^1(\Omega) + P_1^{\mathrm{nc}}.$$

Then the nonconforming finite element method is given by: Find $u^h \in P_1^{\mathrm{nc}}$ with

$$a^h(u^h, v^h) = (f, v^h) \quad \forall \, v^h \in P_1^{\mathrm{nc}}.$$

The goal of this section consists in proving the linear convergence with respect to $h$ in the energy norm $\|\cdot\|_h = \left(a^h(\cdot,\cdot)\right)^{1/2}$. It will be assumed that the solution of the continuous problem (7.6) is smooth, i.e., that $u \in H^2(\Omega)$, that $f \in L^2(\Omega)$, and that the coefficients $a_{ij}(\mathbf{x})$ are weakly differentiable with bounded derivatives. $\quad\Box$

**Remark 7.12** *The error equation.* The first step of proving an error estimate consists in deriving an equation for the error. To this end, multiply the continuous

problem (7.6) with a test function from $v^h \in P_1^{\mathrm{nc}}$, integrate the product on $\Omega$, and apply integration by parts on each triangle. This approach gives

$$
\begin{aligned}
(f, v^h) &= -\sum_{K \in \mathcal{T}^h} \int_K \nabla \cdot (A(\mathbf{x}) \nabla u(\mathbf{x})) \, v^h(\mathbf{x}) \, d\mathbf{x} \\
&= \sum_{K \in \mathcal{T}^h} \int_K (A(\mathbf{x}) \nabla u(\mathbf{x})) \cdot \nabla v^h(\mathbf{x}) \, d\mathbf{x} \\
&\quad - \sum_{K \in \mathcal{T}^h} \int_{\partial K} (A(s) \nabla u(s)) \cdot \mathbf{n}_K(s) v^h(s) \, ds \\
&= a^h(u, v^h) - \sum_{K \in \mathcal{T}^h} \int_{\partial K} (A(s) \nabla u(s)) \cdot \mathbf{n}_K(s) v^h(s) \, ds,
\end{aligned}
$$

where $\mathbf{n}_K$ is the unit outer normal at the edges of the triangles. Subtracting the finite element equation, one obtains

$$
a^h(u - u^h, v^h) = \sum_{K \in \mathcal{T}^h} \int_{\partial K} (A(s) \nabla u(s)) \cdot \mathbf{n}_K(s) v^h(s) \, ds \quad \forall \, v^h \in P_1^{\mathrm{nc}}. \qquad (7.15)
$$

$\square$

**Lemma 7.13 Estimate of the right-hand side of the error equation** (7.15).
*Assume that $u \in H^2(\Omega)$ and $a_{ij} \in W^{1,\infty}(\Omega)$, then it is*

$$
\left| \sum_{K \in \mathcal{T}^h} \int_{\partial K} A(s) \nabla u(s) \cdot \mathbf{n}_K(s) v^h(s) \, ds \right| \le ch \, \|u\|_{H^2(\Omega)} \, \|v^h\|_h .
$$

**Proof:** Every edge of the triangulation which is in $\Omega$ appears exactly twice in the boundary integrals on $\partial K$. The corresponding unit normals possess opposite signs. One can choose for each edge one unit normal and then one can write the integrals in the form

$$
\sum_E \int_E \left[\left[ (A(s) \nabla u(s)) \cdot \mathbf{n}_E(s) v^h(s) \right]\right]_E \, ds = \sum_E \int_E (A(s) \nabla u(s)) \cdot \mathbf{n}_E(s) \left[\left[ v^h \right]\right]_E (s) \, ds,
$$

where the sum is taken over all edges $\{E\}$. Here, $\left[\left[ v^h \right]\right]_E$ denotes the jump of $v^h$

$$
\left[\left[ v^h \right]\right]_E (s) = \begin{cases} v^h|_{K_1}(s) - v^h|_{K_2}(s) & s \in E \subset \Omega, \\ v^h(s) & s \in E \subset \partial\Omega, \end{cases}
$$

where $\mathbf{n}_E$ is directed from $K_1$ to $K_2$ or it is the outer normal on $\partial\Omega$. For writing the integrals in this form, it was used that $\nabla u(s)$, $A(s)$, and $\mathbf{n}_E(s)$ are almost everywhere continuous, such that these functions can be written as factor in front of the jumps. Because of the continuity condition for the functions from $P_1^{\mathrm{nc}}$ and the homogeneous Dirichlet boundary condition, it is for all $v^h \in P_1^{\mathrm{nc}}$ that $\left[\left[ v^h \right]\right]_E (P) = 0$ for the midpoints $P$ of all edges. From the linearity of the functions on the edges, it follows that

$$
\int_E \left[\left[ v^h \right]\right]_E (s) \, ds = 0 \quad \forall \, E. \qquad (7.16)
$$

Let $E$ be an arbitrary edge in $\Omega$ which belongs to the triangles $K_1$ and $K_2$. The next goal consists in proving the estimate

$$
\begin{aligned}
&\left| \int_E (A(s) \nabla u(s)) \cdot \mathbf{n}_E(s) \left[\left[ v^h \right]\right]_E (s) \, ds \right| \\
&\le \quad ch \, \|u\|_{H^2(K_1)} \left( \left\| \nabla v^h \right\|_{L^2(K_1)} + \left\| \nabla v^h \right\|_{L^2(K_2)} \right).
\end{aligned} \qquad (7.17)
$$

To this end, one uses a reference configuration $\left(\hat{K}_1, \hat{K}_2, \hat{E}\right)$, where $\hat{K}_1$ is the unit triangle and $\hat{K}_2$ is the triangle which one obtains by reflecting the unit triangle at the $y$-axis. The common edge $\hat{E}$ is the interval $(0, 1)$ on the $y$-axis. The unit normal on $\hat{E}$ will be chosen to be the Cartesian unit vector $\mathbf{e}_x$, see Figure 7.4. This reference configuration can be transformed to $(K_1, K_2, E)$ by a map which is continuous and on both triangles $\hat{K}_i$ affine. For this map one, can prove the same properties for the transform as proved in Chapter 6.



Figure 7.4: Reference configuration.

Using (7.16), the Cauchy–Schwarz inequality, and the trace theorem, one obtains for an arbitrary constant $\alpha \in \mathbb{R}$

$$
\begin{aligned}
\int_{\hat{E}} \left(\hat{A}(\hat{s})\nabla \hat{u}(\hat{s})\right) \cdot \mathbf{e}_x \left[\!\left[\hat{v}_1^h\right]\!\right]_{\hat{E}} \, d\hat{s} &= \int_{\hat{E}} \left(\left(\hat{A}(\hat{s})\nabla \hat{u}(\hat{s})\right) \cdot \mathbf{e}_x - \alpha\right) \left[\!\left[\hat{v}_1^h\right]\!\right]_{\hat{E}} \, d\hat{s} \\
&\leq c \left\|\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x - \alpha\right\|_{H^1(\hat{K}_1)} \left\|\left[\!\left[\hat{v}_1^h\right]\!\right]_{\hat{E}}\right\|_{L^2(\hat{E})} \quad (7.18)
\end{aligned}
$$

In particular, one can choose $\alpha$ such that

$$
\int_{\hat{E}} \left(\left(\hat{A}(\hat{s})\nabla \hat{u}(\hat{s})\right) \cdot \mathbf{e}_x - \alpha\right) \, d\hat{s} = 0.
$$

The $L^2(\Omega)$ term in the first factor of the right-hand side of (7.18) can be bounded using the estimate from Lemma 6.4 for $k = 0$ and $l = 1$

$$
\begin{aligned}
&\left\|\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x - \alpha\right\|_{H^1(\hat{K}_1)} \\
&\leq c \left(\left\|\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x - \alpha\right\|_{L^2(\hat{K}_1)} + \left\|\nabla\left(\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x - \alpha\right\|_{L^2(\hat{K}_1)}\right) \\
&\leq c \left\|\nabla\left(\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x - \alpha\right)\right\|_{L^2(\hat{K}_1)} \\
&= c \left\|\nabla\left(\left(\hat{A}\nabla \hat{u}\right) \cdot \mathbf{e}_x\right)\right\|_{L^2(\hat{K}_1)}.
\end{aligned}
$$

To estimate the second factor, in the first step, the trace theorem is applied

$$
\begin{aligned}
\left\|\left[\!\left[\hat{v}_1^h\right]\!\right]_{\hat{E}}\right\|_{L^2(\hat{E})} &\leq c \left(\left\|\hat{v}^h\right\|_{H^1(\hat{K}_1)} + \left\|\hat{v}^h\right\|_{H^1(\hat{K}_2)}\right) \\
&\leq c \left(\left\|\nabla \hat{v}^h\right\|_{L^2(\hat{K}_1)} + \left\|\nabla \hat{v}^h\right\|_{L^2(\hat{K}_2)}\right).
\end{aligned}
$$

The second estimate follows from the equivalence of all norms in finite dimensional spaces. To apply this argument, one has to prove that the terms in the last line are in fact norms. Let the terms in the last line be zero, then it follows that $\hat{v}^h = c_1$ in $\hat{K}_1$ and $\hat{v}^h = c_2$ in $\hat{K}_2$. Because $\hat{v}^h$ is continuous in the midpoint of $\hat{E}$, one finds that $c_1 = c_2$ and consequently that $\left[\!\left[\hat{v}^h\right]\!\right]_{\hat{E}} = 0$. Hence, also the left hand side of the estimate is zero. It follows that the right-hand side of this estimate defines a norm in the quotient space of $P_1^{\mathrm{nc}}$ with respect to $\left[\!\left[\hat{v}^h\right]\!\right]_{\hat{E}} = 0$.

Altogether, one obtains for the reference configuration

$$\left| \int_{\hat{E}} \left( \hat{A}(\hat{s}) \nabla u(\hat{s}) \right) \cdot \mathbf{e}_x \left[\!\left[ \hat{v}_1^h \right]\!\right]_{\hat{E}} d\hat{s} \right|$$

$$\leq \quad c \left\| \nabla \left( \left( \hat{A} \nabla u \right) \cdot \mathbf{e}_x \right) \right\|_{L^2(\hat{K}_1)} \left( \left\| \nabla \hat{v}^h \right\|_{L^2(\hat{K}_1)} + \left\| \nabla \hat{v}^h \right\|_{L^2(\hat{K}_2)} \right).$$

This estimate has to be transformed to the triple $(K_1, K_2, E)$. In this step, one gets for the integral on the edge the factor $c$ ($ch$ for $\nabla$ and $ch^{-1}$ for $d\hat{s}$). For the product of the norms on the right-hand side, one obtains the factor $ch$ ($ch$ for the first factor and $c$ for the second factor). In addition, one uses that $A(s)$ and all first order derivatives of $A(s)$ are bounded to estimated the first term on the right-hand side (*exercise*). In summary, (7.17) is proved.

The statement of the lemma follows by summing over all edges and by applying on the right-hand side the Cauchy–Schwarz inequality. ∎

**Theorem 7.14 Finite element error estimate.** *Let the assumptions of Lemma 7.13 be satisfied, then it holds the following error estimate*

$$\left\| u - u^h \right\|_h^2 \leq ch \left\| u \right\|_{H^2(\Omega)} \left\| u - u^h \right\|_h + ch^2 \left\| u \right\|_{H^2(\Omega)}^2 .$$

**Proof:** Applying Lemma 7.13, it follows from the error equation (7.15) that

$$\left| a^h(u - u^h, v^h) \right| \leq ch \left\| u \right\|_{H^2(\Omega)} \left\| v^h \right\|_h \quad \forall \ v^h \in P_1^{\mathrm{nc}}.$$

Let $I_h \ : \ H_0^1(\Omega) \to P_1^{\mathrm{nc}}$ be an interpolation operator with optimal interpolation order in $\left\| \cdot \right\|_h$. Then, one obtains with the Cauchy–Schwarz inequality, the triangle inequality, and the interpolation estimate

$$
\begin{aligned}
\left\| u - u^h \right\|_h^2 \ &= \ a^h(u - u^h, u - u^h) = a^h(u - u^h, u - I_h u) + a^h(u - u^h, I_h u - u^h) \\
&\leq \ \left| a^h(u - u^h, u - I_h u) \right| + ch \left\| u \right\|_{H^2(\Omega)} \left\| I_h u - u^h \right\|_h \\
&\leq \ \left\| u - u^h \right\|_h \left\| u - I_h u \right\|_h + ch \left\| u \right\|_{H^2(\Omega)} \left( \left\| I_h u - u \right\|_h + \left\| u - u^h \right\|_h \right) \\
&\leq \ ch \left\| u - u^h \right\|_h \left\| u \right\|_{H^2(\Omega)} + ch \left\| u \right\|_{H^2(\Omega)} \left( h \left\| u \right\|_{H^2(\Omega)} + \left\| u - u^h \right\|_h \right).
\end{aligned}
$$

∎

**Remark 7.15** *To the error estimate.* If $h$ is sufficiently small, than the second term of the error estimate is of higher order and this term can be absorbed into the constant of the first term. One obtains the asymptotic error estimate

$$\left\| u - u^h \right\|_h \leq ch \left\| u \right\|_{H^2(\Omega)} .$$

□

# 7.4 $L^2(\Omega)$ Error Estimates

**Remark 7.16** *Motivation.* A method is called quasi-optimal in a given norm, if the order of the method is the same as the optimal approximation order. Already for one dimension, one can show that at most linear convergence in $H^1(\Omega)$ can be achieved for the best approximation in $P_1$. This statement can be already verified with the function $v(x) = x^2$. Hence, all considered methods so far are quasi-optimal in the energy norm.

However, the best approximation error in $L^2(\Omega)$ is of one order higher than the best approximation error in $H^1(\Omega)$. A natural question is if finite element methods

converge also of higher order with respect to the error in $L^2(\Omega)$ than with respect to the error in the energy norm.

In this section it will be shown that one can obtain for finite element methods a higher order of convergence in $L^2(\Omega)$ than in $H^1(\Omega)$. However, there are more restrictive assumptions to prove this property in comparison with the convergence prove for the energy norm. □

**Remark 7.17** *Model problem.* Let $\Omega \subset \mathbb{R}^d$, $d \in \{2,3\}$, be a convex polyhedral domain with Lipschitz boundary. The model problem has the form

$$-\Delta u = f \ \text{ in } \Omega, \quad u = 0 \ \text{ on } \partial\Omega. \tag{7.19}$$

For proving an error estimate in $L^2(\Omega)$, the regularity of the solution of (7.19) plays an essential role. □

**Definition 7.18** $m$**-regular differential operator.** Let $L$ be a second order differential operator. This operator is called $m$-regular, $m \geq 2$, if for all $f \in H^{m-2}(\Omega)$ the solutions of $Lu = f$ in $\Omega$, $u = 0$ on $\partial\Omega$, are in the space $H^m(\Omega)$ and the following estimate holds

$$\|u\|_{H^m(\Omega)} \leq c \|f\|_{H^{m-2}(\Omega)} + c \|u\|_{H^1(\Omega)}. \tag{7.20}$$

□

**Remark 7.19** *On the m-regularity.*
- The definition is formulated in a way that it can be applied also if the solution of the problem is not unique.
- For the Laplacian, the term $\|u\|_{H^1(\Omega)}$ can be estimated by $\|f\|_{L^2(\Omega)}$ such that with (7.20) one obtains (*exercise*)

$$\|u\|_{H^2(\Omega)} \leq c \|f\|_{L^2(\Omega)}.$$

- Many regularity results can be found in the literature. Loosely speaking, they say that regularity is given if the data of the problem (coefficients of the operator, boundary of the domain) are sufficiently regular. For instance, an elliptic operator in divergence form ($\Delta = \nabla \cdot \nabla$) is 2-regular if the coefficients are from $W^{1,p}(\Omega)$, $p \geq 1$, and if $\partial\Omega$ is a $C^2$ boundary. Another important result is the 2-regularity of the Laplacian on a convex domain. A comprehensive overview on regularity results can be found in Grisvard (1985).

□

**Remark 7.20** *Variational form and finite element formulation of the model problem.* The variational form of (7.19) is: Find $u \in H_0^1(\Omega)$ with

$$(\nabla u, \nabla v) = (f, v) \quad \forall \, v \in H_0^1(\Omega).$$

The $P_1$ finite element space, with zero boundary conditions, will be used for the discretization. Then, the finite element problem reads as follows: Find $u^h \in P_1$ such that

$$(\nabla u^h, \nabla v^h) = (f, v^h) \quad \forall \, v^h \in P_1. \tag{7.21}$$

□

**Theorem 7.21 Finite element error estimates.** *Let $u(\mathbf{x})$ be the solution of (7.19), let (7.19) be 2-regular, and let $u^h(\mathbf{x})$ be the solution of (7.21). Then, the following error estimates hold*

$$
\begin{aligned}
\left\|\nabla(u - u^h)\right\|_{L^2(\Omega)} &\leq ch \|f\|_{L^2(\Omega)}, \\
\left\|u - u^h\right\|_{L^2(\Omega)} &\leq ch^2 \|f\|_{L^2(\Omega)}.
\end{aligned}
$$

**Proof:** With the error estimate in $H^1(\Omega)$, Corollary 6.16, and the 2-regularity, one obtains

$$\left\|\nabla(u-u^h)\right\|_{L^2(\Omega)} \leq ch\,\|u\|_{H^2(\Omega)} \leq ch\,\|f\|_{L^2(\Omega)}.$$

For proving the $L^2(\Omega)$ error estimate, let $w \in H_0^1(\Omega)$ be the unique solution of the so-called dual problem

$$(\nabla v, \nabla w) = (u-u^h, v) \quad \forall\, v \in H_0^1(\Omega).$$

For a symmetric differential operator, the dual problem has the same form like the original (primal) problem. Hence, the dual problem is also 2-regular and it holds the estimate

$$\|w\|_{H^2(\Omega)} \leq c\left\|u-u^h\right\|_{L^2(\Omega)}.$$

For performing the error estimate, the Galerkin orthogonality of the error is utilized

$$(\nabla(u-u^h), \nabla v^h) = (\nabla u, \nabla v^h) - (\nabla u^h, \nabla v^h) = (f, v^h) - (f, v^h) = 0$$

for all $v^h \in P_1$. Now, the error $u-u^h$ is used as test function $v$ in the dual problem. Let $I_h w$ be the interpolant of $w$ in $P_1$. Using the Galerkin orthogonality, the interpolation estimate, and the regularity of $w$, one obtains

$$
\begin{aligned}
\left\|u-u^h\right\|_{L^2(\Omega)}^2 &= (\nabla(u-u^h), \nabla w) = (\nabla(u-u^h), \nabla(w-I_h w)) \\
&\leq \left\|\nabla(u-u^h)\right\|_{L^2(\Omega)} \|\nabla(w-I_h w)\|_{L^2(\Omega)} \\
&\leq ch\,\|w\|_{H^2(\Omega)} \left\|\nabla(u-u^h)\right\|_{L^2(\Omega)} \\
&\leq ch\,\left\|u-u^h\right\|_{L^2(\Omega)} \left\|\nabla(u-u^h)\right\|_{L^2(\Omega)}.
\end{aligned}
$$

Finally, division by $\left\|u-u^h\right\|_{L^2(\Omega)}$ and the application of the already known error estimate for $\left\|\nabla(u-u^h)\right\|_{L^2(\Omega)}$ are used for completing the proof of the theorem. ∎

# Chapter 8

# Outlook

**Remark 8.1** *Outlook to forthcoming classes.* This class provided an introduction to numerical methods for solving partial differential equations and the numerical analysis of these methods. There are many further aspects that will be covered in forthcoming classes.

*Further aspects for elliptic problems.*
- Adaptive methods and a posteriori error estimators. It will be shown how it is possible to estimate the error of the computed solution only using known quantities and in which ways one can decide where it makes sense to refine the mesh and where not. (Numerical Mathematics IV)
- Multigrid methods. Multigrid methods are for certain classes of problems optimal solvers. (probably Numerical Mathematics IV)
- Numerical analysis of problems with other boundary conditions or taking into account quadrature rules.

*Time-dependent problems.* As mentioned in Remark 1.7, standard approaches for the numerical solution of time-dependent problems are based on solving stationary problems in each discrete time.
- The numerical analysis of discretizations of time-dependent problems has some new aspects, but also many tools from the analysis of steady-state problems are used. (Numerical Mathematics IV)

*Convection-diffusion equations.* Convection-diffusion equations are of importance in many applications. Generally, the convection (first order differential operator) dominates the diffusion (second order differential operator).
- In the convection-dominated regime, the Galerkin method as presented in this class does not work. One needs new ideas for discretizations and these new discretizations create new challenges for the numerical analysis. (Numerical Mathematics IV)

*Problems with more than one unknown function.* The fundamental equation of fluid dynamics, the Navier–Stokes equations, Section 1.3, belong to this class.
- It will turn out that the discretization of the Navier–Stokes equations requires special care in the choice of the finite element spaces. The numerical analysis becomes rather involved. (special class)

$\square$

# Bibliography

Adams, R. A., 1975: *Sobolev spaces*. Academic Press [A subsidiary of Harcourt Brace Jovanovich, Publishers], New York-London, xviii+268 pp., pure and Applied Mathematics, Vol. 65.

Adams, R. A. and J. J. F. Fournier, 2003: *Sobolev spaces*, Pure and Applied Mathematics (Amsterdam), Vol. 140. 2d ed., Elsevier/Academic Press, Amsterdam, xiv+305 pp.

Alt, H., 1999: *Lineare Funktionalanalysis. Eine anwendungsorientierte Einführung.* 3d ed., Springer Berlin.

Braess, D., 2001: *Finite elements.* 2d ed., Cambridge University Press, Cambridge, xviii+352 pp., theory, fast solvers, and applications in solid mechanics, Translated from the 1992 German edition by Larry L. Schumaker.

Brenner, S. C. and L. R. Scott, 2008: *The mathematical theory of finite element methods*, Texts in Applied Mathematics, Vol. 15. 3d ed., Springer, New York, xviii+397 pp., doi:10.1007/978-0-387-75934-0, URL `http://dx.doi.org/10.1007/978-0-387-75934-0`.

Ciarlet, P. G., 1978: *The finite element method for elliptic problems.* North-Holland Publishing Co., Amsterdam, xix+530 pp., studies in Mathematics and its Applications, Vol. 4.

Ciarlet, P. G., 2002: *The finite element method for elliptic problems*, Classics in Applied Mathematics, Vol. 40. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, xxviii+530 pp., doi:10.1137/1.9780898719208, URL `http://dx.doi.org/10.1137/1.9780898719208`, reprint of the 1978 original [North-Holland, Amsterdam; MR0520174 (58 #25001)].

Clément, P., 1975: Approximation by finite element functions using local regularization. *Rev. Française Automat. Informat. Recherche Opérationnelle Sér. RAIRO Analyse Numérique*, **9 (R-2)**, 77–84.

Crouzeix, M. and P.-A. Raviart, 1973: Conforming and nonconforming finite element methods for solving the stationary Stokes equations. I. *Rev. Française Automat. Informat. Recherche Opérationnelle Sér. Rouge*, **7 (R-3)**, 33–75.

Deuflhard, P. and M. Weiser, 2012: *Adaptive numerical solution of PDEs.* de Gruyter Textbook, Walter de Gruyter & Co., Berlin, xii+421 pp.

Dziuk, G., 2010: *Theorie und Numerik partieller Differentialgleichungen.* Walter de Gruyter GmbH & Co. KG, Berlin, x+319 pp., doi:10.1515/9783110214819, URL `http://dx.doi.org/10.1515/9783110214819`.

Ern, A. and J.-L. Guermond, 2004: *Theory and practice of finite elements*, Applied Mathematical Sciences, Vol. 159. Springer-Verlag, New York, xiv+524 pp.

Evans, L. C., 2010: *Partial differential equations*, Graduate Studies in Mathematics, Vol. 19. 2d ed., American Mathematical Society, Providence, RI, xxii+749 pp.

Fefferman, C., 2000: Existence & smoothness of the Navier-Stokes equations. Http://www.claymath.org/millennium/Navier-Stokes_Equations/, http://www.claymath.org/millennium/Navier-Stokes_Equations/.

Goering, H., R. Hans-Görg, and L. Tobiska, 2010: *Die Finite-Elemente-Methode fr Anfänger.* 4th ed., Wiley-VCH, Berlin, ix + 219 pp.

Grisvard, P., 1985: *Elliptic problems in nonsmooth domains*, Monographs and Studies in Mathematics, Vol. 24. Pitman (Advanced Publishing Program), Boston, MA, xiv+410 pp.

Grossmann, C. and H.-G. Roos, 2007: *Numerical treatment of partial differential equations.* Universitext, Springer, Berlin, xii+591 pp., doi:10.1007/978-3-540-71584-9, URL `http://dx.doi.org/10.1007/978-3-540-71584-9`, translated and revised from the 3rd (2005) German edition by Martin Stynes.

Haroske, D. D. and H. Triebel, 2008: *Distributions, Sobolev spaces, elliptic equations.* EMS Textbooks in Mathematics, European Mathematical Society (EMS), Zürich, x+294 pp.

John, V. and G. Matthies, 2004: MooNMD—a program package based on mapped finite element methods. *Comput. Vis. Sci.*, **6 (2-3)**, 163–169, doi:10.1007/s00791-003-0120-1, URL `http://dx.doi.org/10.1007/s00791-003-0120-1`.

Landau, L. and E. Lifschitz, 1966: *Lehrbuch der theoretischen Physik*, Vol. VI, Hydrodynamik. Akademie-Verlag Berlin.

LeVeque, R. J., 2007: *Finite difference methods for ordinary and partial differential equations.* Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, xvi+341 pp., doi:10.1137/1.9780898717839, URL `http://dx.doi.org/10.1137/1.9780898717839`, steady-state and time-dependent problems.

Rannacher, R. and S. Turek, 1992: Simple nonconforming quadrilateral Stokes element. *Numer. Methods Partial Differential Equations*, **8 (2)**, 97–111, doi:10.1002/num.1690080202, URL `http://dx.doi.org/10.1002/num.1690080202`.

Samarskii, A. A., 2001: *The theory of difference schemes*, Monographs and Textbooks in Pure and Applied Mathematics, Vol. 240. Marcel Dekker Inc., New York, xviii+761 pp., doi:10.1201/9780203908518, URL `http://dx.doi.org/10.1201/9780203908518`.

Samarskij, A., 1984: *Theorie der Differenzenverfahren*, Mathematik und ihre Anwendungen in Physik und Technik, Vol. 40. Akademische Verlagsgesellschaft Geest & Portig K.-G., Leipzig.

Schieweck, F., 1997: *Parallele Lösung der stationären inkompressiblen Navier-Stokes Gleichungen.* Otto-von-Guericke-Universität Magdeburg, Fakultät für Mathematik, habilitation.

Scott, L. R. and S. Zhang, 1990: Finite element interpolation of nonsmooth functions satisfying boundary conditions. *Math. Comp.*, **54 (190)**, 483–493, doi:10.2307/2008497, URL `http://dx.doi.org/10.2307/2008497`.

Šolín, P., 2006: *Partial differential equations and the finite element method.* Pure and Applied Mathematics (New York), Wiley-Interscience [John Wiley & Sons], Hoboken, NJ, xviii+472 pp.

Strang, G. and G. Fix, 2008: *An analysis of the finite element method.* 2d ed., Wellesley-Cambridge Press, Wellesley, MA, x+402 pp.

Wladimirow, W. S., 1972: *Gleichungen der mathematischen Physik.* VEB Deutscher Verlag der Wissenschaften, Berlin.

# Index