

2 Second-Order Elliptic Equations

2.1 Characteristic Features

The prototypical representative of elliptic PDEs is the POISSON equation

$$-\Delta u = f, \quad (2.1)$$

which will serve as our elliptic model problem. This section is devoted to the most important properties of this problem, which shall be preserved under any 'good' numerical discretisation scheme.

Boundary Conditions

2.1.1 Definition (Boundary Conditions) Let $g, a : \partial\Omega \rightarrow \mathbb{R}$. A boundary condition of the form

(a) $u = g \text{ on } \partial\Omega$

is called DIRICHLET *boundary condition* or *boundary condition of the first kind*,

(b) $\frac{\partial u}{\partial n} = g \text{ on } \partial\Omega$

is called NEUMANN *boundary condition* or *boundary condition of the second kind*,

(c) $\frac{\partial u}{\partial n} + au = g \text{ on } \partial\Omega$

is called ROBIN *boundary condition* or *boundary condition of the third kind*.

These boundary conditions are said to be homogeneous if $g \equiv 0$, otherwise inhomogeneous.

Existence of Strong Solutions

The *classical* or *strong formulation* of Poisson's equation with DIRICHLET boundary conditions reads: find $u \in C^2(\Omega) \cap C(\bar{\Omega})$ such that

$$-\Delta u = f \quad \text{in } \Omega \quad (2.2a)$$

$$u = g \quad \text{on } \partial\Omega. \quad (2.2b)$$

Proving existence of strong solutions of (2.2) on general domains is difficult and not exactly elegant. It also requires (often too) strong assumptions on the regularity of the domain Ω and the data f . We'll skip the discussion and refer to the literature on PDEs—the keyword is GREEN's functions.

Uniqueness of Strong Solutions

The classical uniqueness proof for strong solutions of POISSON's equation equipped with DIRICHLET boundary conditions relies on a *maximum principle* for elliptic equations:

2.1.2 Theorem (Elliptic Maximum Principle) *Let*

$$L = a_{11} \frac{\partial^2}{\partial x_1^2} + 2a_{12} \frac{\partial^2}{\partial x_1 \partial x_2} + a_{22} \frac{\partial^2}{\partial x_2^2} + a_1 \frac{\partial}{\partial x_1} + a_2 \frac{\partial}{\partial x_2}$$

be elliptic (note that $a \equiv 0$) and $u \in C^2(\Omega) \cap C(\bar{\Omega})$. Then

$$Lu \leq 0 \quad \text{in } \Omega \quad \Rightarrow \quad \max_{x \in \bar{\Omega}} u(x) \leq \max_{x \in \partial\Omega} u(x),$$

i.e. the solution u assumes its maximum on the boundary.

Now the uniqueness of strong solutions to (2.2) is obtained as follows:

We consider the difference between two proposed solutions u_1, u_2 and show that $u_1 = u_2$ on $\bar{\Omega}$.

Strong vs Weak Solutions

A solution of the strong formulation of POISSON's equation with homogeneous DIRICHLET boundary conditions

$$\begin{aligned} -\Delta u &= f \quad \text{in } \Omega \\ u &= 0 \quad \text{on } \partial\Omega \end{aligned} \tag{P}$$

has to be twice continuously differentiable in the interior of the domain (so that the Laplacian can be applied to it) and it must be continuous up to the boundary (so that it actually approaches the boundary values from the interior).

However, asking for two continuous derivatives of the solution u is often too strong a condition. For example, if material parameters suddenly change across an interface that cuts through the domain Ω , then this will normally affect the smoothness of the solution u as well, for instance in the form of a 'kink' (a discontinuity in a first derivative).

It would therefore be desirable to have a relaxed formulation, which generalises the notion of solutions to the POISSON-DIRICHLET problem (P). The central idea behind *variational* or *weak formulations*:

Inner product with a test function v and integrate by parts to transfer derivatives of u onto v ; thus reducing the required number of derivatives on u .

The divergence theorem implies the identity

$$\int_{\Omega} (\operatorname{div} F)v \, dx = \int_{\partial\Omega} (F \cdot n)v \, ds - \int_{\Omega} F \cdot \nabla v \, dx$$

which generalises integration by parts to higher dimensions, where $F : \Omega \rightarrow \mathbb{R}^d$ is a sufficiently regular vector field and $v : \Omega \rightarrow \mathbb{R}$ a scalar function on a domain $\Omega \subset \mathbb{R}^d$.

Applied to (P), we obtain

$$\begin{aligned} \int_{\Omega} -\Delta u \cdot v \, dx &= \int_{\Omega} f \cdot v \, dx \\ &= \int_{\Omega} +\nabla \cdot (\nabla u) \cdot v \, dx \\ &= \int_{\Omega} (-\nabla u \cdot n) v \, ds - \int_{\Omega} (-\nabla u) \cdot \nabla v \, dx \\ &= \int_{\partial\Omega} \frac{\partial}{\partial n} u v \, ds + \int_{\Omega} \nabla u \cdot \nabla v \, dx \quad \leftarrow \text{dirichlet: first term drops} \end{aligned}$$

and if v vanishes on the boundary $\partial\Omega$ as well, then

$$\boxed{\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx}$$

Note that all the second derivatives have now disappeared. It is not even necessary to ask for *continuous* first

derivatives of u and v . Instead, we require the following assumptions to ensure that the weak-formulation is well-defined:

1. u, v are once differentiable
2. $\nabla u \cdot \nabla v$ is integrable
3. a. f, g is integrable
b. or, can more general source terms be used, e.g. δ -functions?

Consequently, the spaces C^k of continuous(ly differentiable) functions are not well suited for weak formulations. Instead, we will use the so-called LEBESGUE spaces L^p and SOBOLEV spaces H^k and $W^{k,p}$, which contain functions that meet certain integrability conditions.

2.1.3 Definition (LEBESGUE Space of Square-Integrable Functions) For a domain $\Omega \subset \mathbb{R}^d$ we define

$$\|u\|_{L^2(\Omega)} = \left(\int_{\Omega} |u(x)|^2 dx \right)^{1/2}$$

\uparrow abs. value for scalar fields, vector norm for vector fields

The set

$$L^2(\Omega) = \{ u : \Omega \rightarrow \mathbb{R} \mid \|u\|_{L^2(\Omega)} < \infty \}$$

is called the LEBESGUE space of order 2.

2.1.4 Theorem (L^2 is a HILBERT Space) The LEBESGUE space L^2 of square-integrable functions is a HILBERT space with the scalar product

$$(u, v)_{L^2} = \int_{\Omega} u v dx \quad \leftarrow \text{scalar prod's}$$

$$(\vec{u}, \vec{v})_{L^2} = \int_{\Omega} \vec{u} \cdot \vec{v} dx \quad \leftarrow \text{vector prod's}$$

Exercise: 2. Since $|\int_{\Omega} \nabla u \cdot \nabla v dx| \leq \|\nabla u\|_{L^2} \|\nabla v\|_{L^2}$, and given that the gradients of u, v exist, it is sufficient to require $\nabla u, \nabla v \in L^2(\Omega)$.

3.a. Similarly, $|\int_{\Omega} f v dx| \leq \|f\|_{L^2} \|v\|_{L^2}$, it is sufficient that $f, v \in L^2(\Omega)$ for well-definedness

2.1.5 Remark The space $L^2(\Omega)$ is actually quite different from spaces of continuous functions, such as $C(\Omega)$:

- L^2 -functions are not necessarily continuous

- In LEBESGUE spaces, we usually cannot assign point values to functions. The expression $u(x)$ is not meaningful. Consider for instance the two L^2 -functions on $[-1, 1]$

$$u_1(x) \equiv 1 \quad \text{and} \quad u_2(x) = \begin{cases} 1 & \text{for } x \neq 0 \\ -3 & \text{for } x = 0 \end{cases}$$

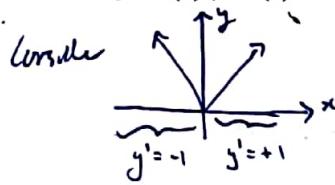
These two functions are not distinguishable in $L^2([-1, 1])$ since

$$\|u_1 - u_2\|_{L^2}^2 = \int_{-1}^1 \left\{ \begin{cases} 0, & x \neq 0 \\ 4, & x = 0 \end{cases} \right. dx = 0$$

$$\int_{-1}^1 = \lim_{\epsilon \rightarrow 0} \int_{-1-\epsilon}^{1+\epsilon} (10dx + \int_{-1-\epsilon}^{1+\epsilon} 4x dx) = \lim_{\epsilon \rightarrow 0} 4(2\epsilon)$$

Thus there is no difference between $u_1(x), u_2(x)$ in $L^2([-1, 1])$

even though $u_1(0) \neq u_2(0)$.



Now, the derivative at $x=0$ is not defined - but that doesn't matter! $\{x=0\}$ is a set of measure zero, and thus information at this single point will not affect the weak solution.

2.1.6 Definition (L^2 -based SOBOLEV Spaces) For a domain $\Omega \subset \mathbb{R}^d$ and $k \in \mathbb{N}_0$ we define the SOBOLEV space $H^k(\Omega)$ as the set of all functions $u \in L^2(\Omega)$, of which all (weak) partial derivatives up to and including order k are in $L^2(\Omega)$ as well.

2.1.7 Theorem (H^k is a HILBERT Space) The L^2 -based SOBOLEV spaces H^k of square-integrable functions with square-integrable derivatives up to order k are HILBERT spaces with the scalar product

$$(u, v)_{H^k} = \int_{\Omega} uv dx + \int_{\Omega} \nabla u \cdot \nabla v dx + \int_{\Omega} \nabla^2 u : \nabla^2 v dx + \dots \quad (\text{up to order } k, \text{ where } \nabla^k \text{ is a } k\text{-tensor and} \\ : \text{ is the Frobenius norm of the pointwise product of the tensors})$$

There are a few more mathematical technicalities involved when it comes to the definition of the homogeneous DIRICHLET boundary conditions of our example problem ((P)). We will not go into further detail here, but refer to the trace theorem for SOBOLEV spaces, which can be found in the literature on functional analysis.

2.1.8 Definition (H_0^1 and H^{-1}) The space of all functions $u \in H^1(\Omega)$ with $u|_{\partial\Omega} = 0$ is denoted by $H_0^1(\Omega)$.

The dual space of $H_0^1(\Omega)$ is denoted by $H^{-1}(\Omega)$. That is,

$$H^{-1}(\Omega) = (H_0^1(\Omega))^* = \{ f : H_0^1(\Omega) \rightarrow \mathbb{R} \mid f \text{ is linear and continuous} \}.$$

For the *duality pairing* of these two spaces, i.e. for $f \in H^{-1}(\Omega)$ and a function $v \in H_0^1(\Omega)$ we normally use the symmetric notation $\langle f, v \rangle_{H^{-1}, H_0^1}$ instead of $f(v)$.

Overall, we have $H_0^1(\Omega) \subset L^2(\Omega) \subset H^{-1}$. The space H^{-1} is actually bigger than L^2 and also contains abstract elements that are no functions on Ω . Depending on the dimension of Ω , this could for instance be objects like delta-'functions':

2.1.9 Example (Source Terms in H^{-1}) **L^2 -functions** If the source term f of the PDE (P) is a given L^2 -function, then f is also an element of H^{-1} and the notation of the duality pairing reduces to the L^2 scalar product: $\langle f, v \rangle_{H^{-1}, H_0^1} = \int_{\Omega} f v \, dx \leftarrow \langle f, v \rangle_{L^2}$

delta-'function' If $f = \delta_0$, i.e. the delta-'function' centred at the point $x = 0$

$$\langle f, v \rangle_{H^{-1}, H_0^1} = \int_{\Omega} \delta_0 v \, dx \leftarrow \begin{array}{l} \text{physicist's notation (we won't use this)} \\ \text{mathematical notation} \end{array}$$

$$= v(0) \leftarrow \text{mathematical notation (note that this only holds for sufficiently Cts., etc. domains, as generally } L^2 \text{ func. aren't point-wise unique, as we've seen)}$$

Note that this delta-'function' requires us to evaluate the test functions v at a point. Whenever this is not possible for all test functions in $H_0^1(\Omega)$ —and this will depend on the dimension of Ω —the delta-'function' does not belong to the space $H^{-1}(\Omega)$.

1. if $u, v \in H_0^1(\Omega)$, then $\nabla u, \nabla v$ exist

2. if $u, v \in H_0^1(\Omega)$, then $\nabla u, \nabla v$ are in $L^2(\Omega)^d$ \leftarrow d -dimensional vector w each component in $L^2(\Omega)$

3. f need not be in $L^2(\Omega)$, rather $H^{-1}(\Omega)$ is sufficient.

We can finally write down a weak formulation of (P): given $f \in H^{-1}(\Omega)$, find a function $u \in H_0^1(\Omega)$ such that for all test functions $v \in H_0^1(\Omega)$

$$\boxed{\int_{\Omega} \nabla u \cdot \nabla v \, dx = \langle f, v \rangle_{H^{-1}, H_0^1}.}$$

(P')

Due to the very different form of the problems (P) and (P'), the classical analytical toolkit required to derive these properties for strong solutions is very different from the functional analytical techniques for weak solutions.

Existence and Uniqueness of Weak Solutions

Two famous results from functional analysis immediately give the existence and uniqueness of weak solutions of (P'): the LAX-MILGRAM lemma together with the POINCARÉ inequality. In our numerical analysis, we will make use of these two results over and over again.

2.1.10 Lemma (POINCARÉ inequality) *There exists a constant $C > 0$ such that for all $u \in H_0^1(\Omega)$*

$$\|u\|_{L^2} \leq C \|\nabla u\|_{L^2}. \quad (2.3)$$

It is important to note that this inequality only holds for functions in H_0^1 , not for all functions in H^1 ! A simple counterexample is given by $u = K > 0$ (const. positive function)

2.1.11 Theorem (LAX-MILGRAM Lemma) *Let V be a HILBERT space with scalar product $(\cdot, \cdot)_V$ and norm $\|\cdot\|_V = \sqrt{(\cdot, \cdot)_V}$. If $B : V \times V \rightarrow \mathbb{R}$ is a bilinear form (linear in each argument)*

- continuous

$$\exists C > 0 \text{ s.t. } \forall u, v \in V \times V, \quad B(u, v) \leq C \|u\|_V \|v\|_V$$

- and coercive

$$\exists c > 0 \text{ s.t. } \forall u \in V, \quad B(u, u) \geq c \|u\|_V^2$$

bilinear form and $f \in V^*$, then the problem

$$B(u, v) = (f, v)_V, \quad \forall v \in V$$

admits a unique solution $u \in V$. This solution satisfies the a priori estimate

$$\|u\|_V \leq \frac{1}{c} \|f\|_{V^*}. \quad (2.4)$$

2.1.12 Example Consider the particular special case where the space V is the n -dimensional Euclidean space \mathbb{R}^n ,

$$B(u, v) = v^T A \cdot u \quad f = b^T$$

where A is a symmetric positive definite $n \times n$ matrix and $b \in \mathbb{R}^n$ a given (column) vector.

We are solving $v^T A u = v^T b$ for all $v \in \mathbb{R}^n$. Rearranging, $(v, A u - b) = 0 \Leftrightarrow A u - b = 0$.

$$\hookrightarrow \text{Continuity: } |v^T A u| \leq \|A u\|_2 \|v\|_2 \leq \|A\|_2 \|u\|_2 \|v\|_2 = \lambda_{\max} \|u\|_2 \|v\|_2$$

$$\hookrightarrow \text{Coercivity: } |v^T A u| \geq \lambda_{\min} \|u\|_2^2 \quad [\text{N.B. For a matrix, coercivity} \Rightarrow \text{positive definiteness,}_{13} \text{but this is not true in general Hilbert spaces}]$$

Then, Lax-Milgram gives us:

$$\hookrightarrow \text{solution to } v^T A u = v^T b \text{ for } v \in \mathbb{R}^n \text{ has a unique soln, and } u \text{ satisfies the a priori estimate } \|u\|_2 \leq \frac{1}{\lambda_{\min}} \|b\|_2$$

- $B(u, v) = v^T A u : |B(u, v)| \leq \|v\| \|A u\|$
 $\leq \|A\| \|u\| \|v\|$
 operator norm of matrix
 thus $|B(u, v)| \leq C \|u\| \|v\|$ where $C = \|A\|$

- $B(u, u) = u^T A u \geq 0 \leftarrow \text{sym. pos. def}$
 Let $A = V D V^T$ where $V V^T = I$, D diagonal matrix of eigenvalues $\begin{bmatrix} \lambda_1 & \dots & \lambda_n \end{bmatrix}$
 $\Rightarrow B(u, u) = u^T V D V^T u \geq \lambda_{\min} u^T u = \lambda_{\min} \|u\|^2$

- Thus, we have: $\|u\|_2 \leq \frac{1}{\lambda_{\min}} \|b\|_2$

2.1.13 Corollary (Existence and Uniqueness of Solutions to (P')) Given $f \in H^{-1}(\Omega)$, there is a unique solution $u \in H_0^1(\Omega)$ such that for all test functions $v \in H_0^1(\Omega)$

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \langle f, v \rangle_{H^{-1}, H_0^1}. \quad (\text{P}')$$

Proof. Let $B(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx$ Cauchy-Schwarz

$$\text{Continuity: } |B(u, v)| \leq \int_{\Omega} \|\nabla u\|_2 \|\nabla v\|_2 \, dx \leq \|u\|_{H_0^1} \|v\|_{H_0^1} \Rightarrow B(u, v) \leq \|u\|_{H_0^1} \|v\|_{H_0^1}$$

Coercivity: $|B(u, u)| = \|\nabla u\|_2^2 \geq \frac{1}{c} \|u\|_2^2$. *Loosely, we write that $B(u, u) \geq c \|u\|_{H_0^1}^2$ for some c*

$$\|u\|_{H_0^1}^2 = \|u\|_2^2 + \|\nabla u\|_2^2 \geq (1 + \frac{1}{c}) \|u\|_2^2$$

[N.B. Coercivity for Laplacian is stronger than positive definiteness, $B(u, u) \geq c \|u\|_{H_0^1}^2$
 as $\|u\|_{H_0^1} \leq \frac{1}{c} \|f\|_{H^{-1}}$ where $B(u, u) \geq c \|u\|_{H_0^1}^2$ necessitates that $\lambda_{\min} \geq c$,]
 (i.e. $\lambda_{\min} \not\rightarrow 0$)

□

Note that this proof of existence and uniqueness is equally valid for strong solutions of (P), since

all strong solutions are also weak solutions.

Continuous Dependence on Data

One more result follows immediately from the LAX-MILGRAM lemma: the continuous dependence of the unique solution u on the data f . If u is a strong or weak solution of the model problem

$$-\Delta u = f \quad \text{in } \Omega \quad u = 0 \quad \text{on } \partial\Omega$$

and \tilde{u} a strong or weak solution of the perturbed problem

$$-\Delta \tilde{u} = \tilde{f} \quad \text{in } \Omega \quad \tilde{u} = 0 \quad \text{on } \partial\Omega,$$

then $u - \tilde{u}$ solves the problem

$$-\Delta(u - \tilde{u}) = f - \tilde{f} \quad \text{in } \Omega \quad u - \tilde{u} = 0 \quad \text{on } \partial\Omega.$$

The estimate (2.4) from the LAX-MILGRAM lemma now gives

$$\|u - \tilde{u}\|_{H_0^1} \leq \frac{1}{C} \|f - \tilde{f}\|_{H^{-1}}$$

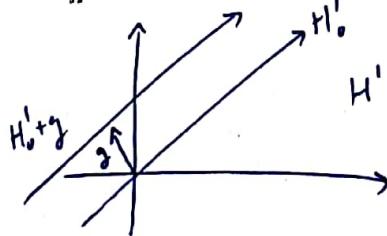
Clearly, the perturbed solution \tilde{u} must approach the unperturbed solution u in the H_0^1 -norm as $\tilde{f} \rightarrow f$ in $H^{-1}(\Omega)$ (or in $L^2(\Omega)$, if both f 's happen to be L^2 -functions).

What about perturbations in the boundary data? So far, we have only considered homogeneous boundary conditions, in fact, without loss of generality. If $g \in H^1(\Omega)$, then its restriction to the boundary $\partial\Omega$ is a function that is less regular by 'half a derivative': $g|_{\partial\Omega} \in H^{1/2}(\partial\Omega)$. Every H^1 -function in Ω has boundary values (aka 'trace') in $H^{1/2}(\partial\Omega)$ and every $H^{1/2}$ -function on $\partial\Omega$ is the trace of a H^1 -function in Ω^* . Consequently, if boundary data $g \in H^{1/2}(\Omega)$ are prescribed,

$$\begin{aligned} & \text{unamerican!} \\ & -\Delta u = f \quad \text{in } \Omega \quad u = g \quad \text{on } \partial\Omega \end{aligned} \tag{2.5}$$

then this function g can be extended to a H^1 -function on the entire domain. The weak formulation of (2.5) reads: find a solution u from the affine subspace $g + H_0^1(\Omega)$ of $H^1(\Omega)$ such that for all test functions $v \in H_0^1(\Omega)$

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \langle f, v \rangle_{H^{-1}, H_0^1}. \tag{2.6}$$



*Traces of H^2 -functions are in $H^{3/2}(\partial\Omega)$, traces of H^3 -functions in $H^{5/2}(\partial\Omega)$ etc. If $u \in H^1(\Omega)$, then its gradient is in $L^2(\Omega)$ and traces of the gradient, such as the normal derivative $\partial_n u$ on the boundary, are in $H^{-1/2}(\partial\Omega)$.

This problem is equivalent to the problem of finding a solution $w = u - g \in H_0^1(\Omega)$ such that for all test functions $v \in H_0^1(\Omega)$

$$\int_{\Omega} \nabla w \cdot \nabla v \, dx = \langle f, v \rangle_{H^{-1}, H_0^1} - \int_{\Omega} \nabla g \cdot \nabla v \, dx, \quad (2.7)$$

another problem with homogeneous boundary conditions, but a modified right hand side instead. The inequality (2.4) from the LAX-MILGRAM lemma implies that w and hence also u depends continuously on the boundary data g as well.

Regularity

In the most general setting, weak solutions to an elliptic, linear, second-order PDE lie in the space $H^1(\Omega)$. In general, such solutions need not be continuous and higher derivatives may not exist. They may not even be bounded. Certain geometries and boundary conditions are particularly notorious for generating unbounded peaks in the solution.

We first introduce LEBESGUE and SOBOLEV spaces with exponents not necessarily equal to two.

2.1.14 Definition (General LEBESGUE Spaces) For a domain $\Omega \subset \mathbb{R}^d$ we define

$$\begin{aligned} \|u\|_{L^p(\Omega)} &= \left(\int_{\Omega} |u(x)|^p \, dx \right)^{1/p} && \text{for } 1 \leq p < \infty \\ \|u\|_{L^\infty(\Omega)} &= \text{ess sup}_{x \in \Omega} |u(x)| = \inf \{ c \geq 0 \mid |u(x)| \leq c \text{ almost everywhere in } \Omega \}. \end{aligned}$$

For $1 \leq p \leq \infty$ the set

$$L^p(\Omega) = \{ u : \Omega \rightarrow \mathbb{R} \mid \|u\|_{L^p(\Omega)} < \infty \}$$

is called a LEBESGUE space of order p .

2.1.15 Definition (General SOBOLEV Spaces) For a domain $\Omega \subset \mathbb{R}^d$ and $k \in \mathbb{N}_0$ we define the SOBOLEV space $W^{k,p}(\Omega)$ as the set of all functions $u \in L^p(\Omega)$, of which all (weak) partial derivatives up to and including order k are in $L^p(\Omega)$ as well.

In many cases, the data of a PDE and its domain are smoother than in the most general setting. For example, we often have a proper function $f \in L^2(\Omega)$ instead of just $f \in H^{-1}(\Omega)$ on the right hand side of the PDE. Our hope is that then this higher regularity of f will propagate through to the solution u so that u will have one extra derivative as well. Provided that the domain is not too irregular, this is indeed the case:

2.1.16 Theorem (H^2 -Regularity) Let Ω be either a domain with C^2 -boundary or a convex polygon. Let $u \in H_0^1(\Omega)$ be the unique weak solution to the elliptic problem

$$Lu = f \quad \text{in } \Omega \quad u = 0 \quad \text{on } \partial\Omega$$

where the linear operator L has coefficients

$$a_{ij} \in C^1(\bar{\Omega}), \quad a_i, a \in L^\infty(\Omega) \quad (i, j \in \{1, 2\})$$

and $f \in L^2(\Omega)$.

Then $u \in H^2(\Omega)$ and there exists a constant $C > 0$ such that

$$\|u\|_{H^2} \leq C\|f\|_{L^2}. \quad (2.8)$$

If the boundary, the coefficients of the differential operator and the data are even smoother, then even higher regularity may be deduced for u . SOBOLEV embeddings can be applied to investigate what other spaces such as L^p or C^k the solution belongs to.

In practical problems, however, domains often possess corners or edges. Then a prediction of H^2 -regularity may already be the optimal smoothness result that the theory allows for. Domains that have corners and that are also non-convex are a classical source of trouble. Around *re-entrant corners* of the domain, i.e. corners where the interior angle is larger than π , solutions often tend to develop singularities that prevent them from being H^2 even if the data are arbitrarily smooth, as the following example shall demonstrate.

2.1.17 Example (Corner Singularities) On the circular sector, described in polar coordinates by $\Omega = \{ (r, \vartheta) \in \mathbb{R}^2 \mid 0 < r < 1 \text{ and } 0 < \vartheta < \frac{3\pi}{2} \}$, we consider the problem

$$-\Delta u = 0 \quad \text{in } \Omega$$

with boundary conditions

$$u(r, \vartheta) = 0 \quad \text{for } \vartheta \in \left\{ 0, \frac{3\pi}{2} \right\} \quad u(r, \vartheta) = \sin \frac{2\vartheta}{3} \quad \text{for } r = 1.$$

The unique strong and weak solution is given by

$$u(r, \vartheta) = r^{2/3} \sin \frac{2\vartheta}{3}$$

and possesses a singularity at the origin. In this example $u \in C^2(\Omega) \cap C(\bar{\Omega})$, but $u \notin C^1(\bar{\Omega})$ and $u \notin H^2(\Omega)$. This result can be generalised to circular sectors with angles $0 < \vartheta < \omega$, where $\omega \in]0, 2\pi]$. Then

$$u(r, \vartheta) = r^{\pi/\omega} \sin \frac{\vartheta\pi}{\omega}$$

solves the homogeneous POISSON equation with appropriate boundary conditions. As soon as $\omega > \pi$, a corner singularity arises in the solution. The extreme case $\omega = 2\pi$ describes a circular disk with a crack. Such problems are of great importance in mechanical and civil engineering.

2.2 Finite Differences for POISSON's Equation

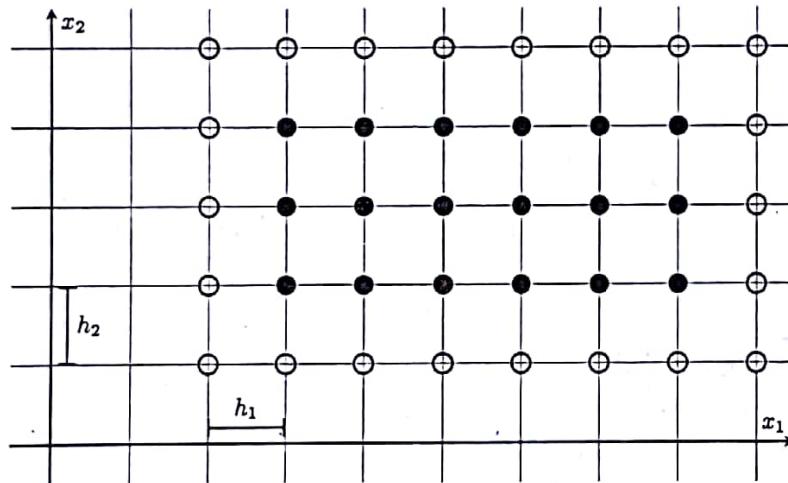
The finite difference method yields a discrete approximation of the *strong formulation* of a PDE. A problem is discretised in three steps:

1. mesh Ω with a (normally Cartesian) grid Ω^h
2. approximate all derivatives with difference quotients
3. set up a system of equations $L^h u^h = f^h$ for the unknown function values u^h on Ω^h

We consider the prototypical elliptic boundary value problem

$$-\Delta u = f \quad \text{in } \Omega \quad u = g \quad \text{on } \partial\Omega \quad (2.9)$$

on a rectangular domain $\Omega \subset \mathbb{R}^2$ that is aligned with the coordinate axes and discretised by a Cartesian grid Ω^h . We assume that Ω^h is equidistant in each direction, with a constant grid spacing of $h_1 > 0$ in x_1 -direction and a constant grid spacing of $h_2 > 0$ in x_2 -direction. For the discrete domain with the boundary points included, we use the notation $\bar{\Omega}^h$.



Approximating Derivatives with Difference Quotients

Our objective is to find an approximation to $-\Delta u(x_1, x_2)$, where $(x_1, x_2) \in \Omega^h$. To this end, we approximate the first partial derivatives at two intermediate points by the central difference quotients

$$\begin{aligned} \frac{\partial u}{\partial x_1} \left(x_1 - \frac{h_1}{2}, x_2 \right) &\approx \frac{u(x_1, x_2) - u(x_1 - h_1, x_2)}{h_1} \\ \frac{\partial u}{\partial x_1} \left(x_1 + \frac{h_1}{2}, x_2 \right) &\approx \frac{u(x_1 + h_1, x_2) - u(x_1, x_2)}{h_1} \end{aligned}$$

and another divided difference using these two quotients yields

$$\frac{\partial^2 u}{\partial x_1^2}(x_1, x_2) \approx \frac{u(x_1 - h_1, x_2) - 2u(x_1, x_2) + u(x_1 + h_1, x_2)}{h_1^2}.$$

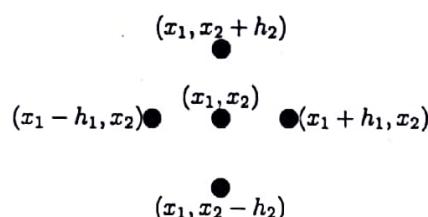
Hence,

$$-\frac{\partial^2 u}{\partial x_1^2}(x_1, x_2) \approx \frac{-u(x_1 - h_1, x_2) + 2u(x_1, x_2) - u(x_1 + h_1, x_2)}{h_1^2} \quad (2.10)$$

and analogously we obtain

$$-\frac{\partial^2 u}{\partial x_2^2}(x_1, x_2) \approx \frac{-u(x_1, x_2 - h_2) + 2u(x_1, x_2) - u(x_1, x_2 + h_2)}{h_2^2} \quad (2.11)$$

The last two equations show that the discrete Laplacian evaluated at the grid point (x_1, x_2) depends on the function values at the point (x_1, x_2) itself plus the four neighbouring grid points $(x_1 \pm h_1, x_2 \pm h_2)$. This dependency pattern is referred to as a *5-point stencil*.



For a computational implementation of the finite difference method, it is desirable to write the difference equations for the unknown function values of u^h on the grid points of Ω^h in matrix form.

If the points in Ω^h are ordered row-wise from bottom left to top right, then (2.10) leads to

$$\frac{1}{h_1^2} \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & \ddots & \ddots & \ddots \\ & & & & 2 & -1 \\ & & & & -1 & 2 & -1 \\ & & & & & \ddots & \ddots & \ddots \\ & & & & & & -1 & 2 \\ & & & & & & & -1 \end{pmatrix} \begin{pmatrix} u_{1,1} \\ u_{1,2} \\ \vdots \\ u_{1,n} \\ u_{2,1} \\ \vdots \\ u_{2,n} \\ \vdots \\ u_{n,1} \\ u_{n,2} \\ \vdots \\ u_{n,n} \end{pmatrix} + \frac{1}{h_1^2} \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ -u_{1,-1} + 2u_{1,1} - u_{1,3} \\ \vdots \\ -u_{n,-1} + 2u_{n,1} - u_{n,3} \end{pmatrix}.$$

Similarly, (2.11) applied to all interior grid points results in

$$\frac{1}{h_2^2} \begin{pmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{pmatrix} \begin{pmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{pmatrix} + \frac{1}{h_2^2} \begin{pmatrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{pmatrix}$$

The complete linear system that represents the discrete counterpart of (2.9) is now given by

$$(L_1^h + L_2^h)u^h = f^h - (l_1^h + l_2^h) \quad (2.12)$$

where the entries of the vector f^h are the function values of the right hand side f evaluated on each grid point.

It is important to note that the discrete negative Laplacian $L^h = L_1^h + L_2^h$ is represented by a *sparse matrix*. If L^h is an $N \times N$ matrix, it contains a total number of N^2 entries. The 5-point stencil of the central differencing scheme reveals that out of the N entries in each row, at most 5 can be non-zero. For points near the boundary, the stencil includes known boundary values. Therefore, the corresponding rows of L^h have 3 or 4 non-zero entries only. Rows corresponding to grid points further away from the boundary contain 5 non-zero entries.

Consequently, only $O(N)$ out of the N^2 entries of L^h are non-zero. For this purpose, software packages for scientific computing usually offer a special data type for sparse matrices where only the non-zero entries and their indices are stored, but not all the $O(N^2)$ zero entries. Since PDE problems lead to systems with very large N ($N \sim 10^7 - 10^9$ for most industrial problems), the memory requirements for the full matrix would be excessive. Furthermore, the multiplication of a sparse matrix with a vector involves only $O(N)$ multiplications of non-zero numbers. If the sparse matrix had been stored fully, a computer would still carry out all N^2 multiplications of matrix entries with vector entries, despite the fact that almost all of these multiplications give zero. Sparse data formats hence allow for an economical use of both memory and CPU resources.

Fundamental Notions of the Numerical Analysis of PDEs

Three closely related notions describe a numerical scheme $T^h(u^h) = 0$ that discretises a problem $T(u) = 0$ (e.g. $T(u) = Lu - f$):

- (1) *Consistency*: Is the discretisation scheme T^h a good approximation of the exact problem T ?

(2) *Stability*: Is the discrete problem $T^h(u^h) = 0$ well-posed and does a small residual $T^h(v^h)$ imply a small error $v^h - u$?

(3) *Convergence*: Is the discrete solution u^h a good approximation of the exact solution u ?

2.2.1 Definition (Consistency) The numerical scheme T^h is said to be

1. *consistent*, if for every solution of $T(u) = 0$

$$\|T^h(u) - T(u)\| \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

2. *consistent of order $O(h^p)$* , if additionally

$$\|T^h(u) - T(u)\| = O(h^p) \quad \text{as } h \rightarrow 0.$$

Note that since $T(u) = 0$, we could have equally written $\|T^h(u) - T(u)\| = \|T^h(u)\|$, but the above notation is probably more illustrative.

2.2.2 Definition (Stability) The numerical scheme T^h is said to be *stable* (with respect to h) if there are constants $h_0 > 0$ and $C > 0$ such that $T^h(u^h) = 0$ has a unique solution for all $h \in]0, h_0]$ and if additionally the stability inequality

$$\|u^h - v^h\| \leq C \|T^h(u^h) - T^h(v^h)\|$$

holds for all discrete functions v^h and all $h \in]0, h_0]$.

Note that since $T^h(u^h) = 0$, we could have equally written $\|T^h(u^h) - T^h(v^h)\| = \|T^h(v^h)\|$, but the above notation is probably more illustrative.

2.2.3 Definition (Convergence) The numerical scheme T^h is said to be

1. *convergent*, if there exists a constant $h_0 > 0$ such that $T^h(u^h) = 0$ has a unique solution for all $h \in]0, h_0]$ and if with the solution of $T(u) = 0$

$$\|u^h - u\| \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

2. *convergent of order $O(h^p)$* , if additionally

$$\|u^h - u\| = O(h^p) \quad \text{as } h \rightarrow 0.$$

2.2.4 Theorem (Consistency \wedge Stability \Rightarrow Convergence) If the scheme T^h is

- (a) *consistent and stable*, then T^h is also *convergent*;
- (b) *consistent of order p and stable*, then T^h is also *convergent of order p* .

Proof. For a scheme that is stable, we have existence and uniqueness of a discrete solution u^h on sufficiently fine grids along with the estimate

$$\|u^h - u\| \leq C \|T^h(u^h) - T^h(u)\| = \|T^h(u)\|$$

where the right hand side $\rightarrow 0$ (or $= O(h^p)$) if the scheme is consistent (of order $O(h^p)$). \square

Consistency of Finite Difference Methods

So far, we have simply written down an ad hoc approximation of (2.9). We shall now (i) confirm that the approximations made do in fact lead to a consistent approximation and (ii) also find a more general, systematic approach for deriving a finite difference scheme. This approach relies on TAYLOR expansions.

For POISSON's equation over a one-dimensional domain Ω , we use the finite difference approximation

$$-u''(x) \approx \frac{-u(x-h) + 2u(x) - u(x+h)}{h^2}.$$

For a solution $u \in C^4(\bar{\Omega})$, TAYLOR expansion yields

$$u(x \pm h) = u(x) \pm hu'(x) + \frac{h^2}{2}u''(x) \pm \frac{h^3}{6}u'''(x) + \frac{h^4}{24}u^{(4)}(\xi_{\pm})$$

where the fourth derivative is evaluated at some point $\xi_+ \in]x, x+h[$ or $\xi_- \in]x-h, x[$, respectively. For the difference between the exact second derivative and its finite difference approximation, the so-called *truncation error*, we obtain

$$\begin{aligned} \left| -u''(x) - \frac{-u(x-h) + 2u(x) - u(x+h)}{h^2} \right| &= \left| -u''(x) - \frac{-\frac{h^2}{2}u''(x) - \frac{h^4}{24}u^{(4)}(\xi_-) - \frac{h^2}{2}u''(x) - \frac{h^4}{24}u^{(4)}(\xi_+)}{h^2} \right| \\ &= \left| \frac{h^2}{24}u^{(4)}(\xi_-) + \frac{h^2}{24}u^{(4)}(\xi_+) \right| \\ &\leq \frac{h^2}{12} \max_{[x-h, x+h]} |u^{(4)}|. \end{aligned}$$

Note that $-u''(x) = f(x)$ and hence we have derived the following result:

$$\|T^h(u)\|_{C(\bar{\Omega})} = \|L^h u\|_{\Omega^h} - f^h\|_{\infty} \leq \frac{h^2}{12} \max_{\bar{\Omega}} |u^{(4)}|$$

i.e. second-order consistency, provided that the exact solution u of the problem is in $C^4(\bar{\Omega})$.

For a 2D domain, we follow the exact same steps, once for the partial derivatives in x_1 -direction and once for the partial derivatives in x_2 -direction. Then we add up the two truncation errors and obtain

2.2.5 Lemma (Consistency on Equidistant Grids) *Let $\Omega^h \subset \mathbb{R}^2$ be a rectangular grid with constant grid spacing h_1 and h_2 in x_1 - and x_2 -direction, respectively. Then the finite difference discretisation (2.12) for the POISSON-DIRICHLET problem is 2nd order consistent, provided that the solution u of the continuous problem is in $C^4(\bar{\Omega})$:*

$$\|L^h u\|_{\Omega^h} - f^h\|_{\infty} \leq \frac{h^2}{6} \max_{\bar{\Omega}} \{ |\partial_{x_1}^4 u|, |\partial_{x_2}^4 u| \} \quad (2.13)$$

with $h = \max \{ h_1, h_2 \}$.

Stability of Finite Difference Methods

For the finite difference method on equidistant grids, the discrete matrix possesses a lot of structure with repeating patterns. For the negative Laplacian on $\Omega = [0, 1]$ with DIRICHLET boundary conditions, we have

$$L^h = \frac{1}{h^2} \begin{pmatrix} 2 & -1 & & \\ -1 & 2 & -1 & \\ & \ddots & \ddots & \\ & & -1 & 2 \end{pmatrix}$$

and its eigenvalues can be computed analytically:

$$\lambda_i = \frac{4}{h^2} \sin^2 \left(\frac{N-i}{N} \frac{\pi}{2} \right) \quad i = 1, \dots, N-1.$$

Here $N = 1/h$ is the number of subintervals.

First of all, we observe that all eigenvalues are positive. Consequently, the linear system of the discrete problem has a unique solution.

Furthermore, with any discrete function v^h we derive

$$\begin{aligned} L^h(u^h - v^h) &= f^h - L^h v^h \\ \Rightarrow u^h - v^h &= (L^h)^{-1}(f^h - L^h v^h) \\ \Rightarrow |u^h - v^h| &\leq \|(L^h)^{-1}\| |(f^h - L^h v^h)| \leq \frac{1}{\lambda_{\min}} |(f^h - L^h v^h)| \end{aligned}$$

where

$$\lambda_{\min} = \lambda_{N-1} = \frac{4}{h^2} \sin^2 \left(\frac{1}{N} \frac{\pi}{2} \right). \quad (2.14)$$

On fine grids with small h and large N ,

$$\lambda_{\min} \approx \frac{4}{h^2} \left(\frac{1}{N} \frac{\pi}{2} \right)^2 = \pi^2, \quad (2.15)$$

so the eigenvalues remain bounded away from zero and we have derived the stability inequality

$$|u^h - v^h| \leq C |(f^h - L^h v^h)| = C \|T^h(v^h)\|.$$

If we use the maximum norm $|\cdot|_\infty$ instead of the Euclidean norm $|\cdot|$, the constant C in this inequality may change.

On a rectangular domain in 2D, one may use separation of variables and then apply the same arguments.

2.2.6 Theorem (Stability Inequality) *For the discrete POISSON-DIRICHLET problem (2.12) we have the stability inequality*

$$|u^h - v^h|_\infty \leq C |L^h v^h - f^h|_\infty \quad (2.16)$$

for all grid functions v^h and with a constant C that is independent of h .

Convergence of Finite Difference Methods

2.2.7 Theorem (Convergence on Equidistant Grids) Let $\Omega^h \subset \mathbb{R}^2$ be a rectangular grid with constant grid spacing h_1 and h_2 in x_1 - and x_2 -direction, respectively. Then the finite difference discretisation (2.12) for the POISSON-DIRICHLET problem is 2nd order convergent, provided that the solution u to the continuous problem is in $C^4(\bar{\Omega})$.

Proof. Second-order consistency & stability \Rightarrow second-order convergence. \square

2.2.8 Remark (C^4 -Regularity up to the Boundary) The assumption that the analytical solution belongs to the space $C^4(\bar{\Omega})$ is not normally satisfied, as such a high regularity of the solution usually requires a sufficiently regular domain, with no corners. Even for the problem

$$-\Delta u = 1 \quad \text{in } \Omega \quad u = 0 \quad \text{on } \partial\Omega$$

with C^∞ -data, $u \notin C^4(\bar{\Omega})$ if Ω is, for example, the square $[0, 1]^2$: assuming that u is four times continuously differentiable up to the boundary, then the PDE prescribes

$$-\Delta u|_{x=0} = 1$$

but the boundary condition implies

$$-\Delta u|_{x=0} = 0$$

Consequently, even with the perfectly smooth data in this example, the corners in the domain do not even allow for second derivatives that are continuous up to the boundary.

Discrete Maximum Principle

We will now have a closer look at the structure of the discretised POISSON-DIRICHLET (and other elliptic) problems. The better we understand the properties of the 'big linear system' $L^h u^h = f^h$, the better we will understand what characteristic features of the continuous problem $Lu = f$ are preserved under a 'suitable' discretisation scheme. Furthermore, we will later use all the information that we can possibly gather on the matrix L^h to select a numerical method for solving the discrete problem that is guaranteed to converge, and that additionally exploits all the structure of L^h to compute u^h as efficiently as somehow possible.

2.2.9 Definition (Diagonal Dominance) If in the i^{th} row of a matrix $A \in \mathbb{R}^{n \times n}$ the absolute value of the diagonal entry is

- greater than or equal to the sum of the absolute values of the off-diagonal terms

$$|A_{ii}| \geq \sum_{j \neq i} |A_{ij}|$$

then we say that A is *weakly diagonally dominant* in this row;

- greater than the sum of the absolute values of the off-diagonal terms

$$|A_{ii}| > \sum_{j \neq i} |A_{ij}|$$

then we say that A is *strictly diagonally dominant* in this row.

A matrix A with the properties that

- (i) A is weakly diagonally dominant in all rows
- (ii) A is strictly diagonally dominant in at least one row
- (iii) for all rows i_0 there exists a chain of indices $i_0 \rightarrow i_1 \rightarrow \dots \rightarrow i_s$ to a strictly diagonally dominant row i_s such that all $a_{i_{l-1} i_l} \neq 0$ ($l = 1, \dots, s$)
- (iv) for any two rows i_0, i_s there exists a chain of indices $i_0 \rightarrow i_1 \rightarrow \dots \rightarrow i_s$ such that all $a_{i_{l-1} i_l} \neq 0$ and all $a_{i_l i_{l-1}} \neq 0$ ($l = 1, \dots, s$)

is called *weakly chained diagonally dominant*. If, instead of (iii), there even holds that

- (iv) for any two rows i_0, i_s there exists a chain of indices $i_0 \rightarrow i_1 \rightarrow \dots \rightarrow i_s$ such that all $a_{i_{l-1} i_l} \neq 0$ and all $a_{i_l i_{l-1}} \neq 0$ ($l = 1, \dots, s$)

then A is called *irreducibly diagonally dominant*.

The two chain properties describe how data from the right hand side and information from each component of the solution propagate through the linear system. As can be seen from the conditions (iii) and (iv), they describe the structure of the sparsity pattern of A .

With the weaker chain property (iii), information from row i_s is referred to in row i_{s-1} . Then the equation in row i_{s-2} refers to the component i_{s-1} of the solution, and hence indirectly also to i_s . Finally, row i_0 directly or indirectly depends on the components i_1, i_2, \dots, i_s of the solution and the right hand side. It is not required for (iii) that conversely row i_s also depends on row i_0 .

This is the difference to the stronger property (iv). Here, information is shared globally and all rows directly or indirectly refer to themselves and all other rows.

2.2.10 Definition (Monotone Matrix) A matrix $A \in \mathbb{R}^{n \times n}$ is said to be (*inverse-*)*monotone* if for all vectors $u \in \mathbb{R}^n$

$$Au \succeq 0 \Rightarrow u \succeq 0$$

Recall that a matrix A is monotone if and only if A^{-1} exists and all its entries are positive or zero: $(A^{-1})_{ij} \geq 0, \forall i, j = 1, \dots, n$.

2.2.11 Definition (*Z*-, *L*- and *M*-Matrices) A matrix $A \in \mathbb{R}^{n \times n}$ is called

- *Z-matrix* or *L_0 -matrix* if all off-diagonal entries are negative or zero:

$$a_{ij} \leq 0 \quad \forall i \neq j$$

- *L-matrix* if all off-diagonal entries are negative or zero and all diagonal entries are positive:

$$a_{ij} \leq 0 \quad \forall i \neq j \quad \text{and} \quad a_{ii} > 0 \quad \forall i$$

- *M-matrix*, if it is a monotone *Z*-matrix.

There are many characterisations of *M*-matrices. The following sufficient condition is the most illustrative one in the context of discretised elliptic PDEs:

Weakly chained diagonal dominance & L-matrix \Rightarrow monotone matrix

2.2.12 Lemma (*M*-Criterion) *If $A \in \mathbb{R}^{n \times n}$ is a weakly chained diagonally dominant L-matrix, then A is also monotone and hence an M-matrix.*

2.2.13 Lemma (Discrete Laplacian is an M-Matrix) *The discretised elliptic operator L^h in the POISSON-DIRICHLET problem (2.12) is an M-matrix, independent of h .*

Proof. The matrix L^h is

- strongly diagonally dominant in all rows corresponding to grid points adjacent to the boundary
- weakly, but not strongly diagonally dominant in all rows corresponding to the other grid points 'further away' from the boundary

The matrix L^h is irreducibly diagonally dominant[†], and in particular it is weakly chained diagonally dominant[‡]. Furthermore, all off-diagonal entries are non-negative while all diagonal entries are positive and hence L^h is an L-matrix. The *M*-criterion now implies that L^h is an M-matrix, as asserted. \square

2.2.14 Theorem (Discrete Elliptic Maximum Principle) *Let the discretised elliptic operator $L^h \in \mathbb{R}^{n \times n}$ be a weakly chained diagonally dominant L-matrix. Then*

$$L^h u^h \leq 0 \quad \text{in } \Omega^h \quad \Rightarrow \quad \max_{x \in \Omega^h} u^h(x) \leq \max_{x \in \partial \Omega^h} u^h(x),$$

i.e. the discrete solution u^h assumes its maximum on the boundary.

Furthermore, the strong connectivity (irreducibility) of the system matrix L^h reveals that even a local change in only one component of the right hand side f^h immediately affects the entire solution u^h globally. This phenomenon of immediate global propagation of information through the entire domain is characteristic for elliptic equations.

2.2.15 Remark (Advantages and Disadvantages of Finite Difference Methods) \oplus Finite difference methods are easy to set up and implement on equidistant, rectangular grids.

- \oplus In these cases, the discrete matrices possess a lot of structure, which can be exploited for efficient solution algorithms of the resulting linear systems.
- \oplus Some finite difference methods accurately reflect characteristic features of an elliptic PDE, e.g. well-posedness, maximum principle, symmetry, positive definiteness, global propagation of information.
- \ominus The equations become very complicated on non-equidistant grids or more complex domains. Then much of the structure of the discrete matrices is lost as well.
- \ominus Finite difference methods only approximate some point values of the solution, they do not return a function that is defined over the entire domain.
- \ominus C^4 regularity up to the boundary is required to guarantee 2nd order convergence for the POISSON equation. This is extremely unrealistic in practice.

[†]"One can walk from any interior grid point to any other interior grid point, taking only steps that are covered by the 5-point stencil."

[‡]"One can walk from any interior grid point to a point adjacent to the boundary where the matrix is strictly diagonally dominant."

2.3 Finite Elements for POISSON's Equation

The finite element method follows an approach which is completely different from the finite difference method. Instead of the strong formulation of a PDE, finite element discretisations start from the *weak formulation* of a PDE. For a linear elliptic equation in weak form, a solution $u \in V$ satisfies an expression of the form

$$B(u, v) = (f, v), \quad \forall v \in V \quad (2.17)$$

with a bilinear form $B : V \times V \rightarrow \mathbb{R}$.

There are two common ways to arrive at a weak form of a given problem, one of which we already know:

- (Weighted) residual method (Galerkin Method) :
 - Multiply a PDE in strong form by a test fun. & I.B.P. to get weak form
- Energy minimization method (Ritz Method) :
 - Less general, but more physical; PDE's tend to arise from a minimization of some "energy", e.g. Lagrangian mechanics, Fermat's principle of least time, etc.
 - e.g. $E(u) = \int_{\Omega} |\nabla u|^2 dx \xrightarrow{\min} (D_u E(u)) = \int_{\Omega} 2 \nabla u \cdot \nabla v dx = 0 \quad \forall v \rightarrow \sqrt{B(u, u)} = \sqrt{\int_{\Omega} 2 \nabla u \cdot \nabla v dx} = \text{"(natural) energy norm"}$

We formulate a number of objectives that a 'good' finite element method should attain:

- Well-posedness/stability
- Maximum principle
- Work on irregular domains
- Weaker regularity requirements on exact soln (i.e. $\in C^4$ for $O(h^2)$ convergence)
- The numerical soln should be in the same fun. space as the true soln (e.g. both H_0^1)

The central idea of finite element methods is:

- Replace the infinite-dimensional fun. space V by a finite-dimensional subspace $V_h \subset V$. (Galerkin approx.)

- Find a solution $u_h \in V_h$ s.t. $v \in V_h$, $B(u_h, v_h) = (f, v_h)$ (Galerkin Eqns.)

$$[\text{N.B.: } B(e_h, v_h) = B(u_h - \bar{u}, v_h) = B(u_h, v_h) - B(\bar{u}, v_h) = (f, v_h) - (f, v_h) = 0,$$

using \uparrow bilinearity

i.e. the error is orthogonal to all test functions v_h (Galerkin Orthogonality)

→ u_h is the best approximation to the true soln \bar{u} from V_h

→ the distance measured in the energy norm of the problem.]

To obtain a discrete problem

$$L^h u^h = f^h$$

from (2.17), we proceed in three steps:

1st Step Choose an N -dimensional subspace V^h and a basis $(\phi_i^h)_{i=1}^N$

2nd Step Write

$$u^h = \sum_{j=1}^N u_j^h \phi_j^h$$

↑
coeff. (piecewise-linear) basis vector

3rd Step Substitute u^h in (2.17) to obtain the GALERKIN equations

$$\sum_{j=1}^N b(\phi_j^h, \phi_i^h) u_j^h = \langle f, \phi_i^h \rangle \quad \forall i = 1, \dots, N$$

(only have to test "orthogonality" against all basis functions ϕ_i^h , not the infinity of all possible v^h)

Linear Finite Elements

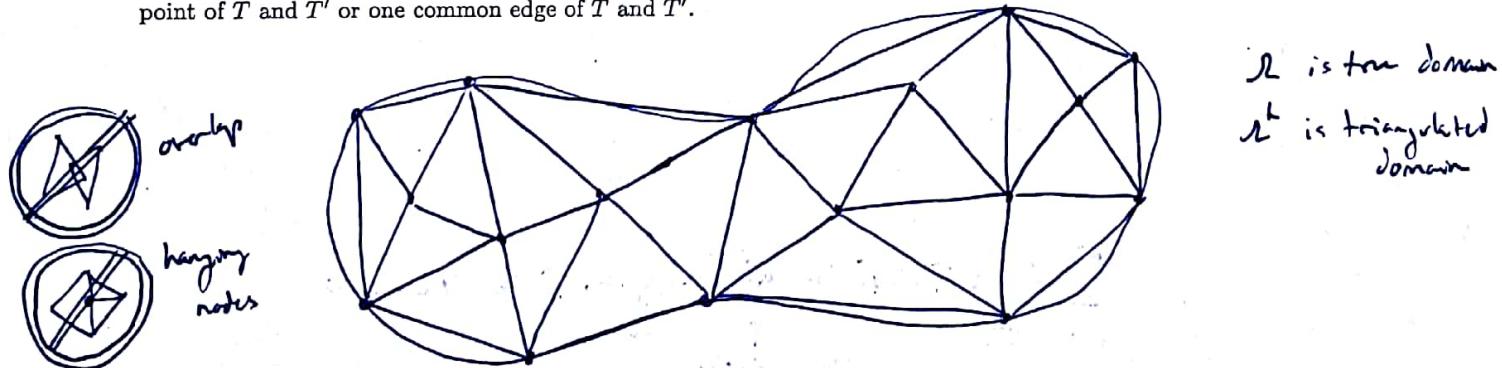
To give an overview of a discretisation with finite elements for a simple example, we re-visit the elliptic model problem

$$-\Delta u = f \quad \text{in } \Omega \quad u = 0 \quad \text{on } \partial\Omega \quad (2.18)$$

in its weak form with a given right hand side $f \in L^2(\Omega)$

$$\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f v \, dx, \quad \forall v \in H_0^1(\Omega). \quad (2.19)$$

First of all, we have to fix a suitable subspace $V^h \subset H_0^1(\Omega)$. To that end, we approximate the domain Ω with a union of closed triangles T_i from a *regular triangulation* $\mathcal{T}^h = \{T_i \mid i = 1, \dots, n_T\}$. 'Regular' means in this context: for any two triangles from $T, T' \in \mathcal{T}^h$ the intersection $T \cap T'$ is either empty, one common corner point of T and T' or one common edge of T and T' .



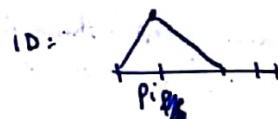
For the mesh size parameter we define $h = \max_{T \in \mathcal{T}^h} \text{diam } T$.

On the discrete domain $\bar{\Omega}^h = \bigcup_{T \in \mathcal{T}^h} T$ we now introduce the space of linear finite elements:

$$V^h = \left\{ v^h \in C(\bar{\Omega}^h) \mid \forall T \in \mathcal{T}^h: v^h|_T \in P_1(T), v^h = 0 \text{ on } \partial\Omega^h \right\}$$

A basis of this space is given by the functions $\phi_i^h \in V^h$ defined by

$$\phi_i^h(p_j) = \delta_{ij} \quad \forall \text{ interior nodes } p_i \text{ of } T^h, j=1 \dots N$$



Moving on to the second step, we may now decompose any function $v^h \in V^h$ (including the discrete solution u^h) in this basis as

$$v^h = \sum_{j=1}^N v^h(p_j) \phi_j^h.$$

Obviously, v^h is fully determined by its function values $v^h(p_j)$ on the interior grid points p_j . Therefore, we also collect these nodal values in a column vector

$$\vec{v}^h = \begin{pmatrix} v^h(p_1) \\ \vdots \\ v^h(p_N) \end{pmatrix}$$

for an alternative representation of v^h .

Thirdly and lastly, the GALERKIN equations for the model problem discretised with linear finite elements read

$$\sum_{j=1}^N \left(\int_{\Omega} \nabla \phi_j^h \cdot \nabla \psi^h \, dx \right) u_j^h = \int_{\Omega} f \cdot \psi^h \, dx \quad \forall i=1 \dots N \quad (2.20)$$

Note that due to the linearity of (2.19) in v , it is sufficient to only use the N basis functions ϕ_i as test functions instead of all infinitely many functions $v^h \in V^h$.

The GALERKIN equations (2.20) can be written in matrix form

$$\begin{matrix} K^h \vec{u}^h & = & \vec{f}^h \\ \uparrow \text{stiffness} & & \uparrow \text{load} \\ \text{matrix} & & \text{vector} \end{matrix}$$

with $k_{ij}^h = \int_0^1 \nabla \phi_j \cdot \nabla \phi_i^h \, dx$

$$f_i^h = \int_0^1 f \phi_i^h \, dx$$

[N.B. here, $k_{ij}^h = k_{ji}^h$, but in general the "shape" basis func.
 ϕ_j will be different from the "test" basis func ϕ_i .
For Galerkin methods, however, they are the same.]

2.3.1 Example (Linear Finite Elements in 1D) In one dimension, the model problem reads: find $u \in H_0^1([0, 1])$ such that for all $v \in H_0^1([0, 1])$:

$$\int_0^1 u' v' \, dx = \int_0^1 f v \, dx.$$

We discretise this problem on the equidistant grid

$$0, h, 2h, 3h, \dots, (N-1)h, 1$$

with N subintervals and grid spacing $h = 1/N$.

Practical Implementation

To evaluate or approximate some integrals that arise in a finite-element discretisation, we often need quadrature formulae, i.e. numerical approximations of integrals. Quadrature formulae are of the form

$$\int_{\Omega} g(x) dx \approx |\Omega| \sum_{i=1}^n w_i g(x^i)$$

where w_i are the weights and x^i the nodes of the quadrature formula.

2.3.2 Theorem (Quadrature Formulae on Intervals)

	Sketch	Nodes	Weights	Error
Midpoint Rule		c	1	$\rightarrow O(h^3)$ if $g \in C^2$ \rightarrow exact for $\deg \leq 1$ polys
Trapezoidal Rule		a b c	$\frac{1}{2}$ $\frac{1}{2}$	(same c = midpoint)
SIMPSON's Rule		a b b c	$\frac{1}{6}$ $\frac{2}{3}$ $\frac{1}{6}$	$\rightarrow O(h^5)$ if $g \in C^4$ \rightarrow exact for $\deg \leq 3$ polys.

At this stage, it is not quite clear yet which quadrature formula should be chosen for what problem. In our convergence analysis of finite element methods, we will derive a rule that will tell us how accurately we have to solve the integrals in the stiffness matrix and on the right hand side. Clearly, we would not want to worsen the convergence rate of the finite element method through too large integration errors. On another hand, we would not want to integrate 'too accurately' if there is a quadrature rule of lower order that already achieves the same convergence rate with less computational effort.

The midpoint rule and the trapezoidal rule extend naturally to higher spatial dimensions:

2.3.3 Theorem (Quadrature Formulae on Triangles)

	Sketch	Nodes	Weights	Error
Midpoint Rule		a	1	$O(h^2)$ if $f \in C^2$ exact for polys of deg ≤ 1
Trapezoidal Rule		a b c	$\frac{1}{3}$ $\frac{1}{3}$ $\frac{1}{3}$	$O(h^3)$ if $f \in C^3$ exact for polys of deg ≤ 2
Cubic Rule		a b c d e f	$\left\{ \begin{array}{l} 3/60 \\ 3/60 \\ 8/60 \\ 8/60 \\ 27/60 \end{array} \right.$	$O(h^5)$ if $f \in C^4$ exact for polys of deg ≤ 3

In practice, it is often easiest to assemble the stiffness matrix, the mass matrix, the right hand side and any other required discrete operators element-by-element. Note that we can split the integration over all of Ω^h into integrals over one triangle only:

We collect the (very few) nonzero contributions from a single element in an *element stiffness matrix* or *element mass matrix*.

2.3.4 Example (Element Mass Matrix) Let T be a triangle and $\phi_1^h, \phi_2^h, \phi_3^h$ the three hat functions that are nonzero on T .

2.3.5 Example (Element Stiffness Matrix) If T has the corner points p_1, p_2, p_3 with corresponding hat functions $\phi_1^h, \phi_2^h, \phi_3^h$ such that $\phi_i^h(p_j) = \delta_{ij}$, then

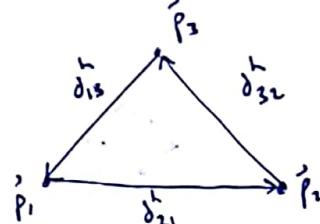
$$K_T^h = \frac{1}{4|T|} \begin{pmatrix} d_{32}^h & \\ d_{13}^h & \\ d_{21}^h & \end{pmatrix} \begin{pmatrix} d_{32}^h & d_{13}^h & d_{21}^h \\ d_{13}^h & d_{32}^h & d_{13}^h \\ d_{21}^h & d_{13}^h & d_{21}^h \end{pmatrix} = \frac{1}{4|T|} \begin{pmatrix} d_{32}^h \cdot d_{32}^h & d_{32}^h \cdot d_{13}^h & d_{32}^h \cdot d_{21}^h \\ d_{13}^h \cdot d_{32}^h & d_{13}^h \cdot d_{13}^h & d_{13}^h \cdot d_{21}^h \\ d_{21}^h \cdot d_{32}^h & d_{21}^h \cdot d_{13}^h & d_{21}^h \cdot d_{21}^h \end{pmatrix}$$

with the triangle edge vectors

$$d_{32}^h = p_3 - p_2 \quad d_{13}^h = p_1 - p_3 \quad d_{21}^h = p_2 - p_1$$

and the triangle area

$$|T| = \frac{1}{2} \|d_{13}^h \times d_{21}^h\| = \dots$$



Proof. An affine function u , i.e. a function of the form

$$u(x) = \alpha x_1 + \beta x_2 + \gamma$$

that is defined by its values on the three vertices $a = p_1, b = p_2, c = p_3 \in \mathbb{R}^2$ of a non-degenerate triangle T (they are not all on one line)

$$u(a) = u_a \quad u(b) = u_b \quad u(c) = u_c$$

satisfies the equation

$$\begin{aligned} u(x) &= \frac{b_1 c_2 - c_1 b_2 + (b_2 - c_2)x_1 + (c_1 - b_1)x_2}{2|T|} u_a \\ &+ \frac{c_1 a_2 - a_1 c_2 + (c_2 - a_2)x_1 + (a_1 - c_1)x_2}{2|T|} u_b \\ &+ \frac{a_1 b_2 - b_1 a_2 + (a_2 - b_2)x_1 + (b_1 - a_1)x_2}{2|T|} u_c. \end{aligned} \tag{2.21}$$

The gradient of u is constant:

$$\nabla u(x) = \frac{1}{4|T|} \begin{pmatrix} (b_2 - c_2)u_a + (c_2 - a_2)u_b + (a_2 - b_2)u_c \\ (c_1 - b_1)u_a + (a_1 - c_1)u_b + (b_1 - a_1)u_c \end{pmatrix}. \tag{2.22}$$

If u is a hat function (restricted to the triangle T), then one of u_a, u_b, u_c is one while the other two are zero. We may now evaluate the integrals

$$\int_T \nabla \phi_i^h \cdot \nabla \phi_j^h \, dx$$

e.g. with the midpoint rule (which is exact for constant functions) for all nine possible combinations of $i, j \in \{1, 2, 3\}$ and this yields the above element stiffness matrix. \square

Software packages for finite element computations typically store the mesh data in three arrays:

- an array $P \in \mathbb{R}^{2 \times n_P}$ with the coordinates of all grid points, such that the k -th column of P defines coordinates p_k of the k -th point ($k = 1, \dots, n_P$)

- an array $E \in \{1, \dots, n_P\}^{2 \times n_E}$ with the indices of points on the boundary *edges*, such that the k -th edge segment $e_{1k} \rightarrow e_{2k}$ coincides with the edge of a triangle ($k = 1, \dots, n_E$)
- an array $T \in \{1, \dots, n_P\}^{3 \times n_T}$ defining the triangulation, such that t_{1k}, t_{2k}, t_{3k} are the indices of the corner points of the k -th *triangle* in anticlockwise order ($k = 1, \dots, n_T$).

Whenever DIRICHLET conditions are imposed on (parts of) the boundary, we want to eliminate these from the system of equations. This is easily achieved by means of projection matrices $P_f \in \{0, 1\}^{N \times n_P}$ and $P_D \in \{0, 1\}^{(n_P - N) \times n_P}$. P_f projects a vector $\bar{u}^h \in \mathbb{R}^{n_P}$ onto its N components $\bar{u}^h \in \mathbb{R}^N$ that are actual degrees of freedom, i.e. points in the interior and boundary points on which no DIRICHLET conditions are imposed. P_D projects a vector $\bar{u}^h \in \mathbb{R}^{n_P}$ onto its components, for which boundary values are already prescribed.

If $\bar{K}^h, \bar{M}^h \in \mathbb{R}^{n_P \times n_P}$ are the stiffness and mass matrices and \bar{f}^h the load vector on all of $\bar{\Omega}^h$, and $g^h \in \mathbb{R}^{n_P - N}$ are the DIRICHLET boundary values, then

$$\begin{aligned}\bar{u}^h &= P_f^T \bar{u}^h + P_D^T g^h \\ \bar{K}^h \bar{u}^h &= \bar{K}^h P_f^T \bar{u}^h + \bar{K}^h P_D^T g^h \\ \bar{M}^h \bar{u}^h &= \bar{M}^h P_f^T \bar{u}^h + \bar{M}^h P_D^T g^h\end{aligned}$$

Recall that the test functions vanish on that part of the boundary, where DIRICHLET conditions apply. Therefore, we also have to delete the $n_P - N$ equations for these DIRICHLET points. For example, the finite element discretisation of the problem

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

reads

$$K^h \bar{u}^h = f^h - k_D^h \quad (2.23)$$

with the (reduced) stiffness matrix and load vector

$$\begin{aligned}K^h &= P_f^T \bar{K}^h P_f \\ f^h &= P_f^T \bar{f}^h\end{aligned}$$

and stiffness terms from any inhomogeneous DIRICHLET boundary values

$$k_D^h = P_f^T \bar{K}^h P_D^T g^h$$

As above, the boundary values from g can be added to the solution \bar{u}^h by setting

$$\bar{u}^h = P_f^T \bar{u}^h + P_D^T g^h$$

General GALERKIN Methods

Besides subspaces of piecewise linear functions, there are many more types of finite-dimensional subspaces V^h of a function space V that one could use for a GALERKIN approximation:

Spaces of Piecewise Polynomial Functions Such spaces V^h are defined based on a mesh of simple polytopes (called *cells*), e.g. intervals in 1D, triangles or quadrilaterals in 2D, tetrahedra or hexahedra in 3D. V^h is the set of all functions that belong to a certain class of polynomials when restricted to one cell, and which often have to meet additional conditions on continuity across cell interfaces. One may choose basis functions for V^h that only have local support.

- (+) domains of complex shape
- (+) all matrices are guaranteed sparse
- (+) (piecewise) polynomial splines have good approximation properties
- (+) polynomials are easy to integrate (theoretically & practically) w/ strong analytic bounds, etc.

GALERKIN approximations with these discrete spaces of locally polynomial functions are called *finite element methods*.

Spaces of Polynomials An alternative approach would be to consider spaces of (globally) polynomial functions over a (hyper-)rectangular domain. On the unit square $\Omega = [0, 1]^2$ these spaces are of the form

$$V^h = \left\{ v^h : \Omega \rightarrow \mathbb{R} \mid v^h(x) = \sum_{i,j=0}^m a_{ij} x_1^i x_2^j \right\}.$$

Taking tensor products of the monomials, i.e. functions like $x_1^i x_2^j$ as basis vectors is not practicable, since no sparsity can be expected for the mass or stiffness matrix and—even worse—since the stiffness matrix is a HILBERT-type matrix. The condition of HILBERT matrices grows exponentially with their dimension N . Even for small values such as $N = 10$, the 10×10 HILBERT matrix already possesses a condition number of $> 10^{13}$. Numerical solutions to that poorly conditioned linear systems are worthless.

Tensor products of polynomials that are orthogonal in $L^2([0, 1]^2)$, such as the CHEBYCHEV or LEGENDRE polynomials lead to well-conditioned mass and stiffness matrices with condition numbers 1 and $O(N)$.

- (+) mass matrix is diagonal
- (-) stiffness matrix is not sparse (assuming derivatives are not orthogonal, which is typically the case)
- (-) Restricted to rectangular domains or unions thereof \square

GALERKIN approximations with spaces of CHEBYCHEV or LEGENDRE polynomials are sometimes classified as *spectral methods*.

Spaces of Trigonometric Functions Again on the unit square $\Omega = [0, 1]^2$, these spaces contain truncated FOURIER series. To approximate a problem with homogeneous DIRICHLET boundary conditions, the space

$$V^h = \left\{ v^h : \Omega \rightarrow \mathbb{R} \mid v^h(x) = \sum_{i,j=0}^n a_{ij} \sin(i\pi x_1) \sin(j\pi x_2) \right\}$$

is an admissible choice with basis vectors $\sin(i\pi x_1) \sin(j\pi x_2)$.

- Ⓐ FFT provides fast solves to discrete GALERKIN equations in complexity $O(N \log N)$
- Ⓑ Gibbs' phenomenon near discontinuity in function derivative
- Ⓒ restricted to rectangular domains

GALERKIN approximations with trigonometric functions are called *spectral methods*.

2.3.6 Remark (PETROV-GALERKIN Methods) The space used for approximating the solution u^h and the space of test functions do not necessarily have to be the same V^h . The generalisation of GALERKIN methods with a space U^h of shape functions and another space V^h of test functions is known as a PETROV-GALERKIN method.

2.3.7 Definition (Finite Element) A *finite element* (in the sense of CIARLET) is defined as a triple (T, P, L) consisting of

- a bounded closed subset $T \subset \mathbb{R}^d$ with nonempty interior and piecewise smooth boundary
- a finite-dimensional space $P = P(T)$ of functions (normally polynomials) over T
- a set of degrees of freedom (or node functionals) $L = L(T)$ that forms a basis for the dual space P^* .

Recall that the dual space is the set of all linear and bounded functionals. Most commonly, these degrees of freedom are set by means of pointwise evaluations of function values or derivatives (and mappings $p \mapsto p(x_0)$, or $p \mapsto \nabla p(x_0)$ are indeed bounded linear functionals on P).

The property that the degrees of freedom must form a basis of P^* , i.e. that they uniquely determine every function in P , is also referred to as *unisolvence*.

2.3.8 Definition (LAGRANGE and HERMITE Elements) If the degrees of freedom of a finite element (T, P, L) only evaluate function values of polynomials in P , then this element is called a LAGRANGE element.

If the degrees of freedom also include evaluations of derivatives of polynomials in P , then this element is called an HERMITE element.

2.3.9 Definition (Interpolation with Finite Elements) Let (T, P, L) be a finite element. We introduce an interpolation operator

$$I^h : H^m(T) \rightarrow P(T), v \mapsto I^h v$$

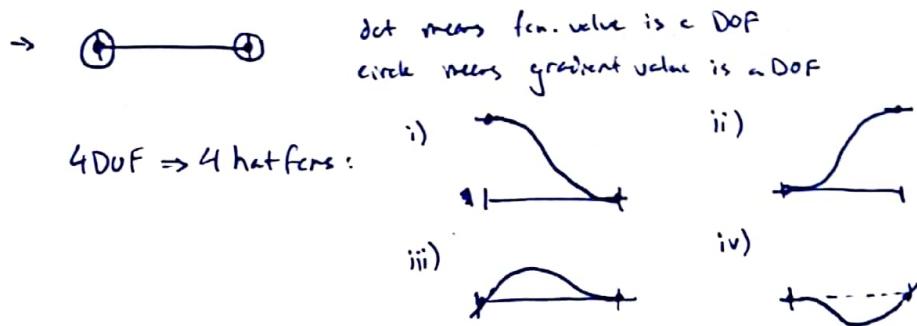
where the interpolant $I^h v$ is defined by

$$\ell(v) = \ell(I^h v), \quad \forall (\text{finitely many}) \ell \in L.$$

By combining the patches from (T, P, L) to a larger domain and composite functions, we obtain a finite-element space V^h . This space may possibly impose further restrictions on continuity or boundary conditions. This could be implemented by equating the degrees of freedom on neighbouring elements or prescribing their values on the boundary, respectively.

2.3.10 Definition (Conforming and Nonconforming Elements) A finite-element space V^h is said to be a *conforming* approximation of a function space V if $V^h \subset V$, otherwise the approximation is called *nonconforming*.

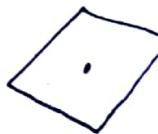
2.3.11 Example (Higher-Order Elements in 1D)



All possible combinations of these give all cubics on the sub-interval;
"classical B-splines"

N.B. Checking conforming vs. non-conforming: show that basis func & sol'n are both in the same space (e.g. H_0^1). Non-conforming: sufficient to provide a counter example

2.3.12 Example (Common Finite Elements in 2D) Piecewise Constant

T	P	L
triangle	$P_0(T)$, dim $P_0(T) = 1$	
quadrilateral	$Q_0(T)$, dim $Q_0(T) = 1$	

$$\operatorname{div} F(u) = \int f \cdot \hat{n}$$

$$\sum_{T_j} \int_T \operatorname{div} F(u) \cdot \hat{n} dx = \sum_{T_j} f \cdot \hat{n} \quad \forall f \in V$$

$$\sum_{T_j} \int_T \operatorname{div} F(u) \phi_i \cdot \hat{n} dx = \sum_{T_j} f \phi_i \cdot \hat{n} \quad \forall i = 1 \dots n_T$$

$$\sum_{j=1}^{n_T} \int_{\partial T_j} \operatorname{div} F(u) \phi_i \cdot \hat{n} ds = \sum_{j=1}^{n_T} \int_{\partial T_j} \phi_i F(u) \cdot \hat{n} ds \quad \leftarrow \text{all integrals vanish except } j=i$$

$$= \int_{\partial T_i} F(u) \cdot \hat{n} ds = \int_{\partial T_i} f \cdot \hat{n} \quad \forall i = 1 \dots n_T$$

"Finite Volume Method"

Piecewise Linear

T	P	L
-----	-----	-----

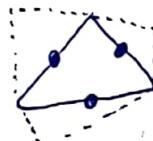
triangle $P_1(T)$, $\dim P_1(T) = 3$



$$P_k = \left\{ \sum_{i+j=0}^k a_{ij} x_1^i x_2^j \right\}$$

$$Q_k = \left\{ \sum_{i+j=0}^k a_{ij} x_1^i x_2^j \right\}$$

triangle $P_1(T)$, $\dim P_1(T) = 3$



Piecewise Quadratic

T	P	L
-----	-----	-----

triangle $P_2(T)$, $\dim P_2(T) = 6$



triangle $P_2(T)$, $\dim P_2(T) = 6$



quadrilateral $Q_1(T)$, $\dim Q_1(T) = 4$



All of these elements possess natural generalisations on tetrahedrons or hexahedrons for three-dimensional problems.

Consistency

A discretisation with finite elements generally introduces errors from three distinct sources:

- ① Galerkin approximation error (= interpolation + geometric approximation error)
- ② quadrature error in load term (from $\int f_i \phi_i d\Omega$)
- ③ Conformity/consistency error

2.3.13 Theorem (BRAMBLE-HILBERT Lemma) Let $T \subset \mathbb{R}^d$ be a domain with LIPSCHITZ boundary and let $F : H^{k+1}(T) \rightarrow \mathbb{R}$ be a bounded and sublinear functional that vanishes for all polynomials of degree $\leq k$:

- $|F(v)| \leq c_1 \|v\|_{H^{k+1}(T)}$ for all $v \in H^{k+1}(T)$
- $|F(u) + F(v)| \leq |F(u)| + |F(v)|$ for all $u, v \in H^{k+1}(T)$
- $F(p) = 0$ for all $p \in P_k(T)$.

Then there exists a constant $c > 0$ such that

$$F(v) \leq c \|\nabla^{k+1} v\|_{L^2(T)}.$$

Proof. This result can be derived with a few technicalities on polynomial projections and a generalised version of POINCARÉ's inequality. The full proof is given on pp 224-225 in the book C GROSSMANN, HG ROOS and M STYNES: *Numerical Treatment of Partial Differential Equations*. Springer, 2007. \square

2.3.14 Theorem (Interpolation Error on Simplices) Let (T, P, L) be a LAGRANGE or HERMITE element, where $T \subset \mathbb{R}^d$ is a d -simplex (interval, triangle, tetrahedron, ...) with inradius r_T and diameter h_T and where $P = P_k(T)$. We denote the corresponding interpolation operator by $I^h : H^{k+1}(T) \rightarrow P_k(T)$.

Then there exists an interpolation constant $c > 0$ depending only on the dimension d and the polynomial degree k such that for all functions $v \in H^{k+1}(T)$ and all $i \in \{0, \dots, k+1\}$

$$\|\nabla^i (v - I^h v)\|_{L^2(T)} \leq c \frac{h_T^{k+1}}{r_T^i} \|\nabla^{k+1} v\|_{L^2(T)}.$$



Proof. We map a simplex T in the domain to a reference simplex \hat{T} , apply the BRAMBLE-HILBERT lemma and map back to T . The mapping between \hat{T} and T is affine:

$$F_T(\hat{x}) = A_T \hat{x} + b_T.$$

The change of variables leads to factors of $\det A_T$, $\|A_T\|$ and $\|A_T^{-1}\|$ in the above norms, which can be estimated in terms of h_T or r_T , respectively. The full details can be found pp 225-226 in the book of GROSSMANN, ROOS, STYNES. \square

2.3.15 Corollary (Interpolation Error with Finite Elements in 2D) Let $\Omega \subset \mathbb{R}^2$ be a polygonal domain and \mathcal{T}^h a triangulation on Ω which satisfies the uniform (in h) shape regularity condition

$$\max_{T \in \mathcal{T}^h} \frac{h_T}{r_T} \leq c. \quad \Leftrightarrow \text{minimum angle condition: } \alpha \geq \alpha^* > 0 \text{ as } h \rightarrow 0 \quad (2.24)$$

Then there exists an interpolation constant $c > 0$ depending only on the regularity constant from (2.24), the dimension d and the polynomial degree k such that

$$\|\nabla^i (v - I^h v)\|_{L^2(\Omega)} \leq ch^{k+1-i} \|\nabla^{k+1} v\|_{L^2(\Omega)}.$$

Proof.

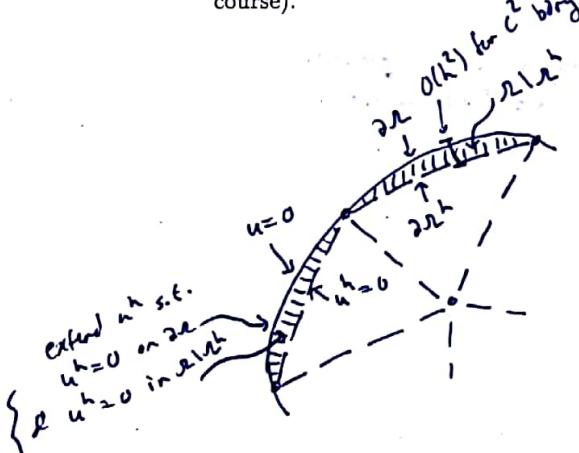
$$\begin{aligned}
 \|\nabla^i(v - I^h v)\|_{L^2(\Omega)}^2 &= \int_{\Omega} |\nabla^i(v - I^h v)|^2 dx \\
 &= \sum_{T \in \mathcal{T}^h} \int_T |\nabla^i(v - I^h v)|^2 dx \\
 &\leq \sum_{T \in \mathcal{T}^h} \left(c \frac{h_T^{k+1}}{r_T^i} \|\nabla^{k+1} v\|_{L^2(T)} \right)^2 \\
 &= \sum_{T \in \mathcal{T}} \left(c \cdot \left[\frac{h_T}{r_T}\right]^i h_T^{k+1-i} \|\nabla^{k+1} v\|_{L^2(T)} \right)^2 \\
 \Rightarrow \|\nabla^i(v - I^h v)\|_{L^2(\Omega)}^2 &\leq c \left[h^{k+1-i} \|\nabla^{k+1} v\|_{L^2(\Omega)} \right]^2
 \end{aligned}$$

(absorbed constants)

□

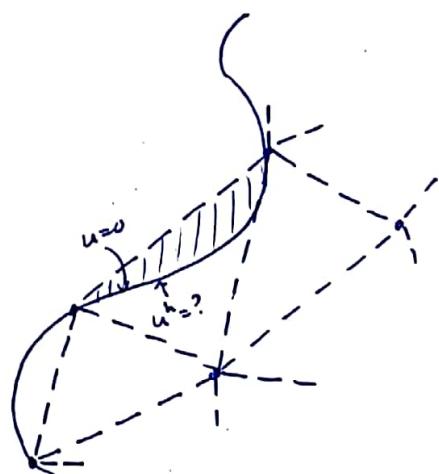
Note the assumption of a polygonal domain in Corollary 2.3.15. This ensures that this domain can be decomposed into triangles and there is no geometric approximation error: $\Omega^h = \Omega$. For curved boundaries, e.g. defined by cubic splines a.k.a. BÉZIER curves as they are typically used in computer aided design (CAD), an additional error arises from the approximation of Ω by Ω^h . This approximation has to be carried out carefully to not deteriorate the overall order of consistency of the discretisation.

In case of $V = H_0^1(\Omega)$ over a convex domain Ω that is approximated by a polygonal domain Ω^h such that all vertices of $\partial\Omega^h$ lie exactly on $\partial\Omega$, the resulting discrete space V^h is a conforming discretisation of V . If Ω is a non-convex domain with curved boundaries, then V^h would usually be non-conforming. This makes the analysis of the non-convex case significantly more difficult (and more suitable to an advanced finite-element course).



44

Convex Domain
(conforming)



Non-Convex Domain
(non-conforming)

2.3.16 Lemma (Estimate on $\Omega \setminus \Omega^h$) For a convex domain $\Omega \subset \mathbb{R}^2$ with C^2 boundary that is approximated with a polygonal domain Ω^h , we have

$$\|w\|_{L^2(\Omega \setminus \Omega^h)} \leq Ch \|w\|_{H^1(\Omega)}$$

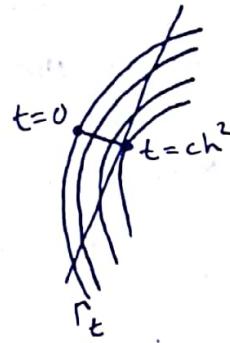
for all functions $w \in H^1(\Omega)$.

Proof. Trace theorem: "a function $w \in H^1(\Omega)$ has traces

(e.g. bdry. values) in $H^{1/2}(\Gamma_t) \subset L^2(\Gamma_t)$.

$$\|w\|_{L^2(\Gamma_t)} \leq C \|w\|_{H^1(\Omega)}$$

$$\text{Now, } \|w\|_{L^2(\Omega \setminus \Omega^h)}^2 = \int_0^{ch^2} \|w\|_{L^2(\Gamma_t)}^2 dt \leq \int_0^{ch^2} C \|w\|_{H^1(\Omega)}^2 dt = ch^2 \|w\|_{H^1(\Omega)}^2$$



With this lemma, we conclude

$$\|I^h v - v\|_{L^2(\Omega \setminus \Omega^h)} = \|v\|_{L^2(\Omega \setminus \Omega^h)} \leq ch^2 \|\nabla v\|_{L^2(\Omega \setminus \Omega^h)} \leq ch^3 \|v\|_{H^2(\Omega)}$$

$$\|\nabla(I^h v - v)\|_{L^2(\Omega \setminus \Omega^h)} = \|\nabla v\|_{L^2(\Omega \setminus \Omega^h)} \leq ch \|v\|_{H^2(\Omega)}$$

Comparing this result with the interpolation error, we see that the polygonal approximation of Ω is 'good enough' for linear finite elements in the sense that the interpolation error doesn't decay at a higher rate than the geometric error: the polygonal boundary approximation does not slow down the rate of convergence for linear finite elements.

For quadratic and higher order elements, though, a boundary approximation that is only piecewise linear will reduce the order of consistency of the approximation. The most common technique to avoid this undesirable effect are isoparametric elements:

2.3.17 Definition (Parametric and Isoparametric Elements) A finite element (T, P, L) is said to be *parametric* if there exists a reference element $(\hat{T}, \hat{P}, \hat{L})$ and an invertible transformation $F_T : \hat{T} \rightarrow T$ such that

- $T = F_T(\hat{T})$
- for any function $p \in P$ there is a corresponding $\hat{p} \in \hat{P}$ such that $p = \hat{p} \circ F_T^{-1}$
- for any degree of freedom $\ell \in L$ there is a corresponding $\hat{\ell} \in \hat{L}$ such that $\ell(p) = \hat{\ell}(p \circ F_T)$ for all $p \in P$.

The element is said to be *isoparametric* if the d components of the transformation F_T are in \hat{P} , i.e. if F_T is a function of the same kind as the shape functions in the reference element.

Stability

We will derive three different stability estimates

- (1) Conforming approximation & exact integration
→ Céa's Lemma
- (2) Conforming approximation & quadrature
→ Strang's 1st Lemma
- (3) Non-conforming approximation
→ Strang's 2nd Lemma

↓ less advanced
↓ more advanced

We recall the abstract weak formulation of a linear PDE

$$\text{Find } u \in V(\gamma) \text{ s.t. } B(u, v) = \langle f, v \rangle_{V^*, V} \quad \forall v \in V$$

its GALERKIN approximation

$$\text{Find } u^h \in V^h(\gamma^h) \text{ s.t. } B(u^h, v^h) = \langle f, v^h \rangle_{V^*, V^h} \quad \forall v^h \in V^h$$

and the property of GALERKIN orthogonality

$$B(e^h = \bar{u} - u^h, v^h) = 0 \quad \forall v^h \in V^h$$

In the sequel we assume the bilinear form $B : V \times V \rightarrow \mathbb{R}$ to be continuous

$$|B(u, v)| \leq C \|u\|_V \|v\|_V \quad \text{"Cauchy-Schwartz" - like}$$

and coercive

$$B(u, u) \geq c \|u\|_V^2$$

with respect to some norm $\|\cdot\|$ norm on the space V .

2.3.18 Lemma (CÉA) *The error $e^h = \bar{u} - u^h$ of the (conforming) GALERKIN approximation satisfies the quasi-best approximation property*

$$\|e^h\| \leq \frac{C}{c} \inf_{v^h \in V^h} \|\bar{u} - v^h\|.$$

Proof. Let $v^h \in V^h$ be arbitrary.

$$\begin{aligned} \|e^h\|^2 &\leq B(e^h, e^h) \leq C \|e^h\|^2 \\ &\sim \\ &= B(e^h, e^h) + B(\underbrace{\bar{u} - \underbrace{u^h}_{e^h}, u^h) \\ &= B(e^h, \bar{u} - v^h) \\ &\leq C \|e^h\| \|\bar{u} - v^h\| \end{aligned}$$

□

$$\Leftrightarrow \|e^h\| \leq \frac{C}{c} \|\bar{u} - v^h\| \quad \forall v^h \in V^h, \text{ thus } \|e^h\| \leq \frac{C}{c} \inf_{v^h \in V^h} \|\bar{u} - v^h\|$$

Note that if we choose the energy norm $\|u\| = \|u\|_B = \sqrt{B(u, u)}$, then $C = c = 1$ and CÉA's lemma gives the best approximation property.

Let us now move on to errors due to inexact numerical integration. Whenever we apply a quadrature formula to assemble e.g. the stiffness matrix, mass matrix or load vector, we obtain a perturbed bilinear form \tilde{B} instead of B and a perturbed right hand side \tilde{f} instead of f and we are solving the perturbed GALERKIN equations

$$\tilde{B}(u^h, v^h) = \langle \tilde{f}, v^h \rangle_{V^h, V^h}, \quad \forall v^h \in V^h$$

instead. To make sure that this problem still possesses a unique solution, we need the additional assumptions that \tilde{B} is continuous and coercive as well

$$\begin{aligned} |\tilde{B}(u^h, v^h)| &\leq \tilde{C} \|u^h\| \|v^h\|, & \forall u^h, v^h \in V^h \\ \tilde{B}(u^h, u^h) &\geq \tilde{c} \|u^h\|^2, & \forall u^h \in V^h \end{aligned}$$

and that \tilde{f} is continuous, just like their unperturbed counterparts B and f .

2.3.19 Lemma (STRANG's First Lemma) *The error $e^h = u - u^h$ of the perturbed (but otherwise conforming) GALERKIN approximation satisfies the estimate*

$$\|e^h\|_h \leq c \left(\inf_{v^h \in V^h} \underbrace{\|\tilde{u} - v^h\|_h + \|B(v^h, \cdot) - \tilde{B}(v^h, \cdot)\|_*}_{\text{approx. error due to } V \rightarrow V^h} + \underbrace{\|\tilde{f} - f\|_*}_{\text{quad. error due to } K^h \rightarrow \tilde{K}^h, M^h \rightarrow \tilde{M}^h, f \rightarrow \tilde{f}^h} \right)$$

with a constant $c > 0$ that is independent of u, u^h and h .

Proof. Let $v^h \in V^h$ be arbitrary. (online notes)

□

Finally, let us add some a few remarks regarding nonconforming approximations. In this case we are solving the discrete problem

$$B^h(u^h, v^h) = (f^h, v^h)_{V^{h*}, V^h}, \quad \forall v^h \in V^h$$

where $V^h \not\subset V$. We need the extra assumptions that B^h can be defined for arguments from V and conversely that B can be defined for arguments from V^h . Additionally we impose continuity and coercivity of B^h

$$\begin{aligned} |B^h(u, v)| &\leq C^h \|u\|_h \|v\|_h, & \forall u, v \in V + V^h \\ B^h(u^h, u^h) &\geq c^h \|u^h\|^2, & \forall u^h \in V^h \end{aligned}$$

where $\|\cdot\|_h$ is some norm on $V + V^h$, the space of all linear combinations $\lambda v + \mu v^h$ with $\lambda, \mu \in \mathbb{R}$, $v \in V$, $v^h \in V^h$. The inhomogeneity f^h is assumed to be continuous in this norm on the space V^h .

We will also need the corresponding operator norm, defined by

$$\|f^h\|_{h*} = \sup_{v^h \in V^h} \frac{|\langle f^h, v^h \rangle_{V^{h*}, V^h}|}{\|v^h\|_h}.$$

2.3.20 Lemma (STRANG's Second Lemma) *The error $e^h = u - u^h$ of the possibly non-conforming finite element approximation satisfies the estimate*

$$\|e^h\|_h \leq c \left(\inf_{v^h \in V^h} \|u - v^h\|_h + \|B^h(u, \cdot) - f^h\|_{h*} \right)$$

with a constant $c > 0$ that is independent of u, u^h and h .

Proof. Let $v^h \in V^h$ be arbitrary. *(online notes)*

□

Convergence

We will now move on to the error analysis of finite element discretisations. To obtain an error estimate in the energy norm $\|e^h\|_B$, we can simply apply the default strategy

consistency of order $m \wedge$ stability \Rightarrow convergence of order m .

The error estimate in the L^2 -norm $\|e^h\|_{L^2}$ can be derived from the energy estimate by applying the so-called AUBIN-NITSCHE trick. Estimates in the L^∞ -norm require very different techniques and therefore we will present these without proof.

We consider a $H_0^1(\Omega)$ -conforming discretisation of POISSON's equation with homogeneous DIRICHLET boundary conditions on a convex polygonal domain $\Omega = \Omega^h \subset \mathbb{R}^2$.

2.3.21 Theorem (Convergence in the Energy Norm) *Let $V^h \subset V = H_0^1(\Omega)$ be a conforming finite element space of piecewise polynomial functions of degree 1. Then the error $e^h = u^h - \bar{u}$ for the approximation of the POISSON-DIRICHLET problem satisfies the a priori estimate*

$$\|e^h\|_B \leq ch \|\nabla^2 \bar{u}\|_{L^2(\Omega)}$$

(lin-fin. elems. of deg. k=1)

and here we even have

$$\|e^h\|_B \leq ch \|f\|_{L^2(\Omega)}$$

with the energy norm $\|e^h\|_B = \|\nabla e^h\|_{L^2}$.

$$\begin{aligned}
 \text{Proof. Céa's Lemma: } \|e^h\|_B &\leq \frac{c}{c} \inf_{v^h} \|\bar{u} - v^h\|_B(\lambda) \\
 &\leq \frac{c}{c} \|\bar{u} - I^h \bar{u}\|_B(\lambda) \\
 &= \frac{c}{c} \|\nabla(\bar{u} - I^h \bar{u})\|_{L^2(\lambda)} \\
 &\stackrel{(2.8)}{\leq} \tilde{c} h \|\nabla^2 \bar{u}\|_{L^2(\lambda)} \\
 &\leq \tilde{c} h \|f\|_{L^2(\lambda)}
 \end{aligned}$$

□

(*) 2.3.22 Theorem (Convergence in the L^2 -Norm) Let $V^h \subset V = H_0^1(\Omega)$ be a conforming finite element space of piecewise polynomial functions of degree $k \geq 1$. Then the error $e^h = u^h - \bar{u}$ for the approximation of the Poisson-Dirichlet problem satisfies the a priori estimate

$$\|e^h\|_{L^2} \leq ch^2 \|\nabla^2 \bar{u}\|_{L^2(\Omega)}$$

and we even have

$$\|e^h\|_{L^2} \leq ch^2 \|f\|_{L^2(\Omega)}.$$

Proof. This proof relies on the AUBIN-NITSCHE trick: consider the so-called dual problem (for $\forall v \in H_0^1(\Omega)$)

$$B(v, z) = \left\langle \frac{e^h}{\|e^h\|_{L^2}}, v \right\rangle_{L^2}, \quad \forall v \in H_0^1(\Omega).$$

test fcn. \rightarrow v

$$\begin{aligned} B(v, z) &= \langle j, v \rangle \\ \bar{u} j &= \frac{e^h}{\|e^h\|_{L^2}} \end{aligned}$$

In our setting, the strong formulation of this problem reads

$$\begin{aligned} -\Delta z &= \frac{e^h}{\|e^h\|_{L^2}} && \text{in } \Omega \\ z &= 0 && \text{on } \partial\Omega \end{aligned}$$

and from (2.8) we conclude that the solution z belongs to $H_0^1(\Omega) \cap H^2(\Omega)$ with

$$\|z\|_{H^2} \leq c \left\| \frac{e^h}{\|e^h\|_{L^2}} \right\|_{L^2} = c.$$

Now, using the test function $v = e^h$ in the dual problem, we obtain

$$\begin{aligned} \|e^h\|_{L^2} &= \left\| \frac{e^h}{\|e^h\|_{L^2}} \right\|_{L^2} \\ \text{dual problem} \rightarrow &= B(e^h, z) \quad \downarrow \text{=0 by Galerkin orthogonality} \\ &= B(e^h, z) - B(e^h, I^h z) \\ &= B(e^h, z - I^h z) \\ &\leq \|e^h\|_3 \|z - I^h z\|_3 \quad \downarrow \tilde{c} h^2 \|\nabla^2 \bar{u}\|_{L^2} \\ &\leq ch \|\nabla^2 \bar{u}\|_{L^2} \quad \downarrow ch \|\nabla^2 z\|_{L^2} \\ &\quad \text{extra } h \text{ from interp.} \quad \text{error on } z \quad \square \end{aligned}$$

In many practical applications, estimates in the L^2 -norm or in the energy norm would not be desirable, as they inherently include some averaging of the error over the entire domain. Even local singularities where the solution blows up may still give a finite error.

In structural engineering, for instance, a local spike in the stress acting on a building or a bridge may result in the failure of the structure. In such settings, pointwise estimates, just like for the finite difference method, are more desirable:

50 (*) Thm. is also true for higher order fin. elem. ($k > 1$), but w/o extra smoothness of true soln \bar{u} , higher order elems. usually do not attain a higher order of convergence.

2.3.23 Theorem (Convergence in the L^∞ -Norm) Let $V^h \subset V = H_0^1(\Omega)$ be the conforming space of linear finite elements on a (sometimes called uniform) triangulation that satisfies the size regularity condition

$$\min_{T \in \mathcal{T}^h} \frac{h_T}{h} \geq c > 0 \quad \text{for all fine grids } h \in]0, h_0].$$

\max edge length on T
 ℓ_{\max} over all h_T 's

Then the error $e^h = u^h - \bar{u}$ for the approximation of the Poisson-Dirichlet problem satisfies the a priori estimate

$$\max_{\bar{\Omega}} |e^h| = \|e^h\|_{L^\infty} \leq ch \|\nabla^2 \bar{u}\|_{L^2}$$

provided that $\bar{u} \in H_0^1(\Omega) \cap H^2(\Omega)$ (e.g. if Ω is a convex polygon). If even $\bar{u} \in H_0^1(\Omega) \cap C^2(\bar{\Omega})$, then

$$\|e^h\|_{L^\infty} \leq ch^2 |\ln h| \max_{\bar{\Omega}} |\nabla^2 \bar{u}|.$$

If $V^h \subset V$ is a conforming finite element space of degree $k \geq 2$ and $\bar{u} \in H_0^1(\Omega) \cap C^{k+1}(\bar{\Omega})$, then the logarithmic term is not needed:

$$\|e^h\|_{L^\infty} \leq ch^{k+1} \max_{\bar{\Omega}} |\nabla^{k+1} \bar{u}|.$$

Proof. Not easy! Please refer to Section 4.4.3 in the book of GROSSMANN, ROOS, STYNES. \square

2.3.24 Remark (Error Estimates for Finite Differences and Finite Elements) We have now encountered a number of different error estimates. This includes the error of the finite difference method

$$\max_{\bar{\Omega}^h} |e^h| \leq ch^2 \max_{\bar{\Omega}} |\nabla^4 \bar{u}| \quad \begin{array}{c} \text{needs assumptions} \\ \text{needs 4 th. derivs} \end{array} \quad (\text{FD})$$

and the L^∞ -error estimates for linear finite elements

$$\max_{\bar{\Omega}} |e^h| \leq ch^2 |\ln h| \max_{\bar{\Omega}} |\nabla^2 \bar{u}| \quad \begin{array}{c} \text{needs assumptions} \\ \text{needs 2 th. derivs} \end{array} \quad (\text{FE1})$$

$$\max_{\bar{\Omega}} |e^h| \leq ch \|\nabla^2 \bar{u}\|_{L^2(\Omega)} \quad \begin{array}{c} \text{needs assumptions} \\ \text{needs 2 weak derivs} \end{array} \quad (\text{FE2})$$

2.3.25 Theorem (Convergence with Numerical Integration) If the integrals of the weak formulation are evaluated with a quadrature formula of order r and if $V^h \subset V$ is a conforming finite element space of piecewise polynomial functions of degree k , then the a priori estimate

$$\|e^h\|_B \leq ch^{\min\{k, r-k+1\}} \|f\|_{L^2(\Omega)}$$

holds. Consequently, the order of convergence is the same as with exact integration provided that $r \geq 2k - 1$.

Proof. i) show that the perturbed bilinear form \tilde{B} is still coercive & continuous
ii) use errors of quadrature formulae in Strang's 1st Lemma

More detail is shown in Section 4.5.3 of the book of GROSSMANN, ROOS, STYNES. \square

On triangles, we have introduced the midpoint rule ($r = 1$), the trapezoidal rule ($r = 2$) and a formula with a cubic rate of convergence ($r = 3$). the latter would be of interest for piecewise quadratic approximation spaces.

A Posteriori Error Estimation

We distinguish between *a priori* and *a posteriori error estimates*. *A priori* and *a posteriori* are Latin expressions which roughly translate into *before* and *after*. In our context, this refers to error estimates which we can already compute before the numerical solution of the problem, based on the theory alone, and error estimates which we compute afterwards using the solution u^h .

A priori error estimates are of the form

$$\|e^h\| = ch^m \quad \text{as } h \rightarrow 0$$

where the norm $\|\cdot\|$ is commonly the energy norm, the H^1 -norm, the L^2 -norm or the L^∞ -norm. *A priori* estimates tell us at what rate the numerical solution u^h converges to the analytical solution \bar{u} as we refine the mesh. E.g. if we have a cubic rate of convergence ($m = 3$), then one mesh refinement $h \rightarrow h/2$ asymptotically results in an error $\|e^{h/2}\| \approx \|e^h\|/8$ eight times smaller on the finer mesh than on the coarser mesh. An *a priori* error estimate does not tell us anything about the actual size of the error e^h , though, as the constant c is generally unknown.

We normally use the approach

$$\text{consistency of order } m \wedge \text{stability} \Rightarrow \text{convergence of order } m$$

to derive an *a priori* error estimate.

In contrast, *a posteriori* error estimates are of the form

$$\|e^h\| \approx \eta(u^h)$$

with a *computable* number (yes, an actual number like $\eta(u^h) = 0.01388$) on the right hand side. Even though the analytical solution u is usually unknown and hence we cannot calculate the error from $e^h = u^h - \bar{u}$ in practice, we still want to ‘postprocess’ our numerical solution u^h in a certain way to at least compute an approximate value $\eta(u^h)$ for the real but unknown error $\|e^h\|$. To qualify as a sensible approximation, such an *a posteriori* error estimator should decay at the same rate as the real error $\|e^h\|$ for finer and finer grids, which can be formalised as follows:

2.3.26 Definition (Efficient and Reliable Error Estimator) An error estimator η is called *efficient* if it does not decay slower than $\|e^h\|$ as $h \rightarrow 0$. That is, there exists a constant $c > 0$ such that

$$c\eta(u^h) \leq \|e^h\|.$$

The error estimator is called *reliable* if it does not decay faster than $\|e^h\|$ as $h \rightarrow 0$. That is, there exists a

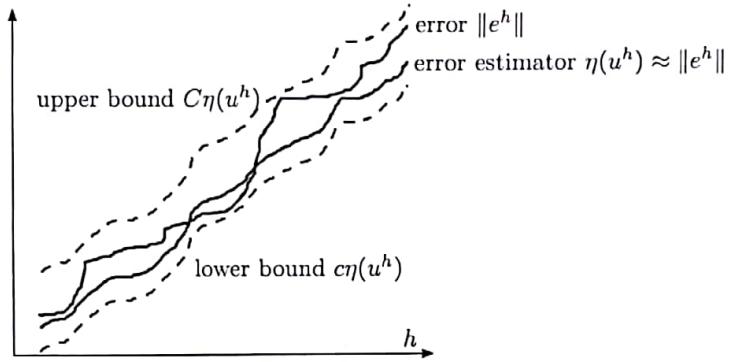
constant $C > 0$ such that

$$\|e^h\| \leq C\eta(u^h).$$

If

$$c\eta(u^h) \leq \|e^h\| \leq C\eta(u^h),$$

then η is said to be both efficient and reliable.



A very powerful framework for a posteriori error estimation is the *dual weighted residual (DWR) method*. The word ‘dual’ already foreshadows that it will be based on the solution to a dual problem, just like the AUBIN-NITSCHE trick for L^2 -error estimates.

In fact, the DWR method is not only applicable to

- (i) different global norms of the error,

but also to the

- (ii) error of a specific ‘quantity of interest’ $J(u)$,

where $J : V \rightarrow \mathbb{R}$ is a functional of the solution u . For example, the solution u of a PDE (the NAVIER-STOKES equations) could describe the flow velocity in a wind tunnel around a car. Then the actual quantity of interest might not always be the complete picture of the overall flow velocity, but more specifically the drag coefficient of the car $c_d = J(u)$. At low speeds, c_d is a linear functional of the velocity u . Due to this linearity, the error between the numerical drag coefficient computed from the numerical solution u^h and the real drag coefficient computed from the exact solution \bar{u} is

$$J(u^h) - J(\bar{u}) = J(u^h - \bar{u}) = J(e^h).$$

Another quantity of interest might be the stress at the tip of a crack, or more generally, the value of the solution u at a given point $P \in \Omega$, $J(u) = u(P)$. In other settings, an average temperature might be of interest, i.e. an average of the solution over a region $R \subset \Omega$, $J(u) = \frac{1}{|R|} \int_R u \, dx$.

For simplicity, we will assume that J is linear in the sequel. With the PDE (‘primal problem’)

$$B(u, v) = \langle f, v \rangle, \quad \forall v \in V$$

we introduce the corresponding dual problem

$$B(v, z) = J(v), \quad \forall v \in V$$

which possesses the dual solution $z \in V$. Testing this equation with $v = e^h$ and exploiting GALERKIN orthogonality results in

$$J(u^h) - J(\bar{u}) = J(e^h) = B(e^h, z) = B(e^h, z - v^h)$$

for arbitrary functions $v^h \in V^h$. In practice, we choose $v^h = I^h z$ as an interpolant or a similar approximation of z .

For the POISSON-DIRICHLET problem we obtain

$$J(u^h) - J(\bar{u}) = \int_{\Omega} \nabla e^h \cdot \nabla(z - v^h) dx.$$

Like for a priori estimates, we assume that the data f is an L^2 -function and that the exact solution is in H^2 . Then we decompose this integral into contributions from each cell and integrate by parts

$$\begin{aligned} J(u^h) - J(\bar{u}) &= \sum_{T \in \mathcal{T}^h} \int_T \nabla e^h \cdot \nabla(z - v^h) dx \\ &= \sum_{T \in \mathcal{T}^h} \int_T (-\Delta e^h)(z - v^h) dx + \int_{\partial T} (\partial_n e^h)(z - v^h) ds. \end{aligned}$$

On each triangle $-\Delta e^h = -\Delta(u^h - \bar{u}) = -\Delta u^h - f$. For triangle edges $\Gamma \subset \partial T$ in the interior of the domain, the boundary integral over Γ appears twice in the sum over all triangles: once for the triangle T on one side of Γ , and once for the neighbouring triangle T' on the other side of Γ . Since the two normal vectors n and n' point in opposite directions, $n = -n'$, we combine the two integrals as follows:

$$\int_{\Gamma} (\partial_n e^h|_T)(z - v^h) ds + \int_{\Gamma} (\partial'_n e^h|_{T'})(z - v^h) ds = \int_{\Gamma} (\partial_n e^h|_T - \partial'_n e^h|_{T'})(z - v^h) ds.$$

Since, due to the regularity assumption $\bar{u} \in H^2(\Omega)$ the analytical solution has no jump in its normal derivative across triangle edges, $\partial_n e^h|_T - \partial'_n e^h|_{T'} = (\partial_n u^h|_T - \partial_n u^h|_{T'}) = [\partial_n u^h]$. We will now use the compact $[]$ -notation for such jumps. In the sum over all triangles, we attribute half of this term to the cell T , the other half to the neighbour T' . On edges $\Gamma \subset \partial\Omega$ which fall on the boundary of the domain $z = v^h = 0$. Overall, we obtain

$$\begin{aligned} J(u^h) - J(\bar{u}) &= \sum_{T \in \mathcal{T}^h} \int_T \underbrace{(-\Delta u^h - f)}_{\text{cell residual}} \underbrace{(z - v^h)}_{\text{dual weight}} dx + \int_{\partial T \setminus \partial\Omega} \underbrace{\frac{1}{2}[\partial_n u^h]}_{\text{jump residual}} \underbrace{(z - v^h)}_{\text{dual weight}} ds \\ &= \sum_{T \in \mathcal{T}^h} \int_T r_T w_T dx + \int_{\partial T \setminus \partial\Omega} r_{\partial T} w_{\partial T} ds \end{aligned} \tag{2.25}$$

In this a posteriori error identity, the *cell residuals* $r_T = -\Delta u^h - f$ measure how accurately the discrete solution satisfies the exact PDE on T . Note that for linear finite elements $-\Delta u^h = 0$ on each triangle so that $r_T = -f$. The *jump residuals* $r_{\partial T} = \frac{1}{2}h_T^{-1/2}[\partial_n u^h]$ provide a measure of the smoothness of the discrete

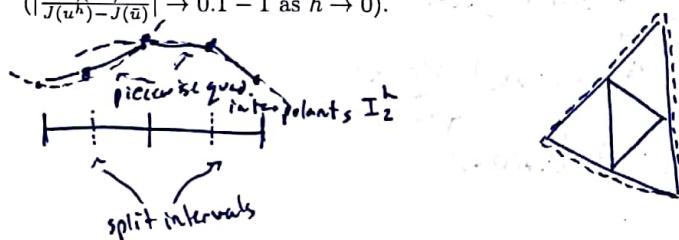
solution. The dual weights $w_T = z - v^h$ and $w_{\partial T} = h_T^{1/2}(z - v^h)$ quantify how strongly the local residual r_T or $r_{\partial T}$, respectively, influence the total error $J(u^h) - J(\bar{u})$:

$$\frac{2[J(\bar{u}) - J(\tilde{u})]}{2r_T} \approx w_T \quad \frac{2[J(u^h) - J(\bar{u})]}{2r_{\partial T}} \approx w_{\partial T}$$

Of course, we can only solve the dual problem numerically, so the exact solution z of the dual problem that arises in the error identity (2.25) is usually unknown.

In a first attempt, we might re-use the same finite elements that we already used for the primal problem (the actual PDE) to compute a numerical solution z^h from the same finite-dimensional space V^h . However, if we set $z \approx z^h$, we obtain $z - I^h z \approx z^h - I^h z^h = 0$ and thus the pointless error estimate $e^h \approx \eta(u^h) = 0$. The following two options are more sensible:

1. Solve the dual problem with finite elements of higher order than the primal problem. This procedure would be asymptotically reliable and efficient with constants $c = C = 1$. However, the computation of the error estimator would be more expensive than the solution of the PDE, at least if it's a linear PDE.
2. Instead, one normally uses the same finite element spaces for both the primal and the dual problem. The dual solution z^h is then interpolated with a higher-order spline over larger patches of cells to approximate z . The interpolant $I^h z$ is then equal to the numerical solution z^h . This simple postprocessing step works surprisingly well, even though the error typically tends to be slightly underestimated ($|\frac{\eta(u^h)}{J(u^h) - J(\bar{u})}| \rightarrow 0.1 - 1$ as $h \rightarrow 0$).



→ fine mesh: for solving original + dual problem
→ rough mesh: for post-processing (error estimates)
only

We can also use this generic approach to obtain a posteriori estimates e.g. for the global energy and L^2 -error norms. In that case, we define the quantity of interest as

$$J(u) = B \left(\frac{e^h}{\|e^h\|_B}, u \right) \Rightarrow \|e^h\|_B = J(e^h)$$

$$J(u) = \left\langle \frac{e^h}{\|e^h\|_{L^2}}, u \right\rangle_{L^2} \Rightarrow \|e^h\|_{L^2} = J(e^h).$$

Since the error e^h is unknown in practice, we cannot actually evaluate the latter two J 's and we cannot solve the dual problem numerically. Instead, we have to use interpolation estimates to estimate the dual weights and to obtain a computable error estimator.

2.3.27 Theorem (A Posteriori Energy and L^2 -Norm Estimates) *For the POISSON-DIRICHLET problem, the*

following a posteriori error estimates hold:

$$\|e^h\|_B \leq \eta_B(u^h) = c \sqrt{\sum_{T \in \mathcal{T}^h} h_T^2 (\|r_T\|_{L^2(T)}^2 + \|r_{\partial T}\|_{L^2(\partial T \setminus \partial \Omega)}^2)}$$

$$\|e^h\|_{L^2} \leq \eta_{L^2}(u^h) = c \sqrt{\sum_{T \in \mathcal{T}^h} h_T^4 (\|r_T\|_{L^2(T)}^2 + \|r_{\partial T}\|_{L^2(\partial T \setminus \partial \Omega)}^2)}$$

These estimators are reliable and, for linear finite elements, at least asymptotically efficient in the sense

$$\eta_B(u^h) \leq c \|e^h\|_B + c \sqrt{\sum_{T \in \mathcal{T}^h} h_T^2 \|f\|_{L^2(T)}^2}$$

$$\eta_{L^2}(u^h) \leq c \|e^h\|_{L^2} + c \sqrt{\sum_{T \in \mathcal{T}^h} h_T^4 \|f\|_{L^2(T)}^2}.$$

Proof. (L^2 -estimate)

$$\|e^h\|_{L^2} = J(e^h) = J(u^h) - J(\bar{u})$$

$$\stackrel{(2.25)}{=} \sum_{T \in \mathcal{T}^h} \int_T (-\Delta u^h - f) (z - I^h z) \, dz + \int_{\partial T \setminus \partial \Omega} \frac{1}{2} [\partial_n u^h] (z - I^h z) \, ds$$

$\leftarrow z = I^h z = 0 \text{ on } \partial \Omega$

$$\text{Cauchy-Schwarz} \leq \sum_{T \in \mathcal{T}^h} \|-\Delta u^h - f\|_{L^2(T)} \|z - I^h z\|_{L^2(T)} + \frac{1}{2} \|[\partial_n u^h]\|_{L^2(\partial T \setminus \partial \Omega)} \|z - I^h z\|_{L^2(\partial T)}$$

$$\text{Interpolation Estimates} \leq \sum_{T \in \mathcal{T}^h} \|-\Delta u^h - f\|_{L^2(T)} \underbrace{h_T^2 \|v^h z\|_{L^2(T)}}_{2.3.14} + \frac{1}{2} \|[\partial_n u^h]\|_{L^2(\partial T \setminus \partial \Omega)} \underbrace{h_T^{3/2} \|\nabla^2 z\|_{L^2(T)}}_{\text{Stability}} \\ = c \sum_{T \in \mathcal{T}^h} h_T^2 (\|r_T\|_{L^2(T)} + \|r_{\partial T}\|_{L^2(\partial T \setminus \partial \Omega)}) \|\nabla^2 z\|_{L^2(T)}$$

$$\text{Cauchy-Schwarz} \leq \underbrace{\sum_{i,j} a_{ij} b_j}_{\sum a_i b_i \leq \sqrt{\sum a_i^2} \cdot \sqrt{\sum b_j^2}} \leq c \sqrt{\sum_{T \in \mathcal{T}^h} h_T^4 (\|r_T\|_{L^2(T)}^2 + \|r_{\partial T}\|_{L^2(\partial T \setminus \partial \Omega)}^2)}$$

$$\text{and } (a+b)^2 \leq 2(a^2 + b^2)$$

□

2.3.28 Remark (Adaptive Finite Element Methods) Here we have used a posteriori error estimates to compute an approximation of the unknown error. Another application is goal-oriented mesh refinement: the goal is to compute $J(u)$ as accurately as possible, while making optimal use of the limited computational resources that are available. E.g. if the memory allows for a computation with N nodes in the mesh, then one would start with computing u^h on a coarse mesh. The error identity (2.25)

$$J(u^h) - J(\bar{u}) = \eta(u^h) = \sum_{T \in \mathcal{T}^h} \eta_T(u^h)$$

gives local error indicators $\eta_T(u^h)$ on every triangle T . One can then refine only those triangles with the largest error indicators, solve the problem again on the locally refined mesh and continue with this procedure until the maximum number N of nodes has been reached.

Discrete Maximum Principle

2.3.29 Definition (Acute Triangular Mesh) A triangular mesh is said to be *weakly acute*, if the interior angles of all its triangles are $\leq \frac{\pi}{2}$ and *strongly acute* if they are $< \frac{\pi}{2}$.



2.3.30 Definition (DELAUNAY Triangulation) A triangulation \mathcal{T}^h is called DELAUNAY triangulation if the three vertices of each triangle $T \in \mathcal{T}^h$ are the only mesh vertices in the circumcircle of T .



One can show that weakly acute triangulations are automatically DELAUNAY triangulations.

2.3.31 Theorem *The stiffness matrix discretising the negative Laplacian with linear finite elements on a DELAUNAY triangulation is an M-matrix.*

Proof. Lemma 21 in TJ BARTH: *Numerical methods for gasdynamic systems on unstructured meshes*. In: *An introduction to recent developments in theory and numerics for conservation laws*. Springer, 1999. \square