

# Numerical Methods for PDEs — TA Summary

Created by [Jonas Bachmann](#), [Paul Fischill](#), [Samuel Russo](#) and [Nico Graf](#)

Ported to Typst by [Valentin Vogt](#)

Last updated on May 09, 2024

## About

### Theorems, definitions and equations

from the lecture notes come in boxes like this one.

### Less important results

that are given for context or completeness look like this.

### Tips and practical advice

from the TAs are highlighted like this.

## Basics

### Theorem 0.3.1.19: Cauchy–Schwarz Inequality

If  $a$  is a symmetric positive semi-definite bilinear form, then

$$|a(u, v)| \leq a(u, u)^{\frac{1}{2}} a(v, v)^{\frac{1}{2}} \quad (1)$$

### Equation 1.3.4.15 (Cauchy–Schwarz for Integrals)

$$\left| \int_{\Omega} u(x)v(x) \, dx \right| \leq \left( \int_{\Omega} |u(x)|^2 \, dx \right)^{\frac{1}{2}} \left( \int_{\Omega} |v(x)|^2 \, dx \right)^{\frac{1}{2}} = \|u\|_{L^2(\Omega)} \|v\|_{L^2(\Omega)} \quad (2)$$

### Norms

- **Supremum norm:**  $\|u\|_{\infty} = \|u\|_{L^{\infty}(\Omega)} := \sup_{x \in \Omega} \|u(x)\|$
- **$L^2$  norm:**  $\|u\|_2 = \|u\|_{L^2(\Omega)} := \left( \int_{\Omega} \|u(x)\|^2 \, dx \right)^{\frac{1}{2}}$

### Theorem 0.3.2.31: Transformation rule for Integration

Given two domains  $\Omega, \widehat{\Omega}$  and a continuous, differentiable mapping  $\Phi : \widehat{\Omega} \rightarrow \Omega$

$$\int_{\Omega} f(x) \, dx = \int_{\widehat{\Omega}} f(\Phi(\hat{x})) |\det D\Phi(\hat{x})| \, d\hat{x} \quad (3)$$

# 1 Second-Order Scalar Elliptic Boundary Value Problems

## 1.2 Quadratic Minimization Problems

In the following, let  $V$  be a vector space over  $\mathbb{R}$ .

### Linear forms

$\ell : V \rightarrow \mathbb{R}$  is a *linear form* / *linear functional*  $\iff$

$$\ell(\alpha u + \beta v) = \alpha \ell(u) + \beta \ell(v) \quad \forall u, v \in V, \forall \alpha, \beta \in \mathbb{R} \quad (4)$$

### Bilinear forms

$a : V \times V \rightarrow \mathbb{R}$  is a *bilinear form*  $\iff$

$$\begin{aligned} a(su_1 + u_2, tv_1 + v_2) \\ = st \cdot a(u_1, v_1) + sa(u_1, v_2) + ta(u_2, v_1) + a(u_2, v_2) \forall u_1, u_2, v_1, v_2 \in V, \forall s, t \in \mathbb{R} \end{aligned} \quad (5)$$

### Positive definiteness

A bilinear form  $a : V \times V \rightarrow \mathbb{R}$  is *positive definite* if

$$u \in V \setminus \{0\} \iff a(u, u) > 0$$

It is *positive semi-definite* if

$$a(u, u) \geq 0 \quad \forall u \in V$$

### Quadratic Functional

A *quadratic functional*  $J : V \rightarrow \mathbb{R}$  is defined by

$$J(u) := \frac{1}{2}a(u, u) - \ell(u) + c, \quad u \in V \quad (6)$$

where  $a : V \times V \rightarrow \mathbb{R}$  a symmetric bilinear form,  $\ell : V \rightarrow \mathbb{R}$  a linear form and  $c \in \mathbb{R}$ .

### Continuity of linear form

A linear form  $\ell : V \rightarrow \mathbb{R}$  is *continuous* / *bounded* on  $V$ , if

$$\exists C > 0 \quad |\ell(v)| \leq C\|v\| \quad \forall v \in V, \quad (7)$$

where  $\|\cdot\|$  is a norm on  $V$ .

## 1.3 Sobolev Spaces

When we solve a minimization problem, we first need to define the space of functions in which we look for the solution. For example, in Physics, we generally want the solution to be continuous. E.g, a function describing the shape of an elastic string should not have jumps.

It turns out that the correct space to describe our minimization problems is the **Sobolev space**. For functions  $u$  in a Sobolev space, the bilinear form  $a$  in the quadratic functional is well defined (i.e.,  $a(u, u) < \infty$ ). Hence the space in which we look for minimizers is determined by the given quadratic functional. To select the space for your problem, follow the guideline ... *Choose the largest space such that the problem is well defined.*

### Sobolev Spaces

$H_0^1(\Omega)$  is a vector space with norm

$$|v|_{H^1} := \left( \int_{\Omega} \|\mathbf{grad} v\|^2 dx \right)^{\frac{1}{2}}$$

$H^1(\Omega)$  is another vector space with norm

$$\|v\|_{H^1}^2 := \|v\|_{L^2}^2 + |v|_{H^1}^2$$

Note that  $|\cdot|_{H^1}$  is not a norm on the space  $H^1(\Omega)$ , but a seminorm.

Both spaces contain all functions for which the norm is finite (and, in the case of  $H_0^1(\Omega)$ , which are 0 on  $\partial\Omega$ ).

**Alternative notation** for norms includes  $\|\cdot\|_0$  for  $\|\cdot\|_{L^2}$  and  $|\cdot|_1$  for  $|\cdot|_{H^1}$ .

If the quadratic minimization problem is well defined, we get the following lemma for existence and uniqueness of minimizers:

#### Theorem 1.3.3.6: Existence of minimizers in Hilbert spaces

On a real Hilbert space  $V$  with norm  $\|\cdot\|_a$  for any  $\|\cdot\|_a$ -bounded linear functional  $\ell : V \rightarrow \mathbb{R}$ , the quadratic minimization problem

$$\begin{aligned} u_* &= \operatorname{argmin}_{v \in V} J(v) \\ J(v) &:= \frac{1}{2} \|v\|_a^2 - \ell(v) \end{aligned} \tag{8}$$

has a unique solution.

Note that here, we use the bilinear form to define the norm  $\|u\|_a = \sqrt{a(u, u)}$ . The main point is that we can see the bilinear form of the quadratic minimization problem as the norm of some Sobolev space. The above theorem guarantees that a solution exists in this space if the linear form is bounded.

For checking boundedness we can often use Cauchy–Schwarz (Equation 1.3.4.15) and Poincaré–Friedrichs (Theorem 1.8.0.20).

## 1.4 Linear Variational Problem

### Definition 1.4.1.6: Linear Variational Problem

Let  $V$  be a vector (function) space,  $\widehat{V} \subset V$  an affine space, and  $V_0 \subset V$  the associated subspace. The equation

$$\text{Find } u \in \widehat{V} \text{ such that } a(u, v) = \ell(v) \quad \forall v \in V_0 \quad (9)$$

is called a (generalized) *linear variational problem*, if

- $a : V \times V_0 \rightarrow \mathbb{R}$  is a bilinear form
- $\ell : V_0 \rightarrow \mathbb{R}$  is a linear form

Theorem 1.3.3.6 tells us that the minimization problem has a solution, but knowing that a solution exists is of course not enough: We want to find it, but an infinite-dimensional minimization problem is hard to solve. To make it easier, we reformulate the problems in a linear variational form (9), which is quite close to something we can solve numerically. To do this transformation, we use the following equivalence:

### Theorem 1.4.1.8: Equivalence of quadratic minimization problem and linear variational problem

For a (generalized) quadratic functional  $J(v) = \frac{1}{2}a(v, v) - \ell(v) + c$  on a vector space  $V$  and with a symmetric positive definite bilinear form  $a : V \times V \rightarrow \mathbb{R}$  the following is equivalent:

- The quadratic minimization problem for  $J(v)$  has the unique minimizer  $u_* \in \widehat{V}$  over the affine subspace  $\widehat{V} = g + V_0, g \in V$ .
- The linear variational problem

$$u \in \widehat{V} \quad a(u, v) = \ell(v) \quad \forall v \in V_0$$

has the unique solution  $u_* \in \widehat{V}$ .

Note that the trial space  $\widehat{V}$ , from which we pick a solution, and the test space  $V_0$  can be different. For an example of different trial and test spaces, see Section 1.7.

## 1.5 Boundary Value Problems

### Lemma 1.5.2.1: General product rule

For all  $\mathbf{j} \in (C^1(\overline{\Omega}))^d$ ,  $v \in C^1(\overline{\Omega})$  holds

$$\operatorname{div}(\mathbf{j}v) = v \operatorname{div} \mathbf{j} + \mathbf{j} \cdot \operatorname{grad} v \quad \text{in } \Omega \quad (10)$$

### Lemma 1.5.2.4: Gauss' Theorem

Let  $\mathbf{n} : \partial\Omega \rightarrow \mathbb{R}^d$  denote the exterior unit normal vector field on  $\partial\Omega$  and  $dS$  denote integration over a surface. We have

$$\int_{\Omega} \operatorname{div} \mathbf{j}(\mathbf{x}) \, d\mathbf{x} = \int_{\partial\Omega} \mathbf{j}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) \, dS(\mathbf{x}) \quad \forall \mathbf{j} \in (C_{\text{pw}}^1(\overline{\Omega}))^d \quad (11)$$

### Lemma 1.5.2.7: Green's first formula

For all vector fields  $\mathbf{j} \in (C_{\text{pw}}^1(\overline{\Omega}))^d$  and functions  $v \in C_{\text{pw}}^1(\overline{\Omega})$  holds

$$\int_{\Omega} \mathbf{j} \cdot \operatorname{grad} v \, d\mathbf{x} = - \int_{\Omega} \operatorname{div} \mathbf{j} v \, d\mathbf{x} + \int_{\partial\Omega} \mathbf{j} \cdot \mathbf{n} v \, dS \quad (12)$$

### Lemma 1.5.3.4: Fundamental lemma of the calculus of variations

Let  $f \in L^2(\Omega)$  satisfy

$$\int_{\Omega} f(\mathbf{x})v(\mathbf{x}) \, d\mathbf{x} = 0 \quad \forall v \in C_0^\infty(\Omega), \quad (13)$$

then  $f \equiv 0$ .

We have seen that minimizing a quadratic functional is equivalent to solving a linear variational problem (9). The variational problem is called the **weak form**. We can transform it (with extra smoothness requirements) into the problem's **strong form**, an elliptic BVP (PDE with boundary conditions).

#### Weak to strong

1. Use Lemma 1.5.2.7 to get rid of derivatives on  $v$  (e.g. turn  $\operatorname{grad} u \cdot \operatorname{grad} v$  into  $-\operatorname{div}(\operatorname{grad} u) + \dots$ )
2. Use properties of the test space (usually that  $v = 0$  on  $\partial\Omega$ ) to get rid of boundary terms
3. Use Lemma 1.5.3.4 to remove the integrals and test functions

## 1.7 Boundary Conditions

For 2nd-order elliptic BVPs we need boundary conditions to get a unique solution. To be more precise, we need **exactly one** of the following boundary conditions on every part of  $\partial\Omega$

### Main boundary conditions for 2nd-order elliptic BVPs

- **Dirichlet:**  $u$  is fixed to be  $g : \partial\Omega \rightarrow \mathbb{R}$

$$u = g \quad \text{on } \partial\Omega$$

- **Neumann:** the flux  $\mathbf{j} = -\kappa(\mathbf{x}) \text{ grad } u$  through  $\partial\Omega$  is fixed with  $h : \partial\Omega \rightarrow \mathbb{R}$

$$\mathbf{j} \cdot \mathbf{n} = -h \quad \text{on } \partial\Omega$$

- **Radiation:** flux depends on  $u$  with an increasing function  $\Psi : \mathbb{R} \rightarrow \mathbb{R}$

$$\mathbf{j} \cdot \mathbf{n} = \Psi(u) \quad \text{on } \partial\Omega$$

In the weak form, Dirichlet conditions have to be imposed directly on the **trial** space. The test space needs to be set to 0 wherever Dirichlet conditions are given (“*Don’t test where the solution is known*”). For example, trial and test spaces for a standard Dirichlet problem are

$$V = \{u \in C^1(\Omega) \mid u = g \text{ on } \partial\Omega\}$$

$$V_0 = \{v \in C^1(\Omega) \mid v = 0 \text{ on } \partial\Omega\}$$

Dirichlet BCs are called **essential boundary conditions**.

Neumann conditions, which are only enforced through some term in the variational equation, are called **natural boundary conditions**.

There are some constraints on the boundary data:

- **Admissible Dirichlet Data:** Dirichlet boundary values need to be continuous.
- **Admissible Neumann Data:**  $h$  needs to be in  $L^2(\Omega)$  (can be discontinuous)

The following theorem is frequently needed when dealing with integrals over the boundary:

### Theorem 1.9.0.19: Multiplicative trace inequality

$$\exists C = C(\Omega) > 0 : \|u\|_{L^2(\partial\Omega)} \leq C \|u\|_{L^2(\Omega)} \cdot \|u\|_{H^1(\Omega)} \quad \forall u \in H^1(\Omega) \quad (14)$$

## 1.8 Second-Order Elliptic Variational Problems

We have seen how we can get from a minimization problem via a variational problem to a BVP. Now we want to move in the opposite direction: from a PDE with boundary conditions to a variational problem.

### Strong to weak

1. Test the PDE with (multiply by  $v$ ) and integrate over  $\Omega$
2. Use Lemma 1.5.2.7 to “shift” one derivative from  $u$  to  $v$  (e.g., from  $-\operatorname{div}(\operatorname{grad} u)$  to  $\operatorname{grad} u \cdot \operatorname{grad} v + \dots$ )
3. Use Neumann BC on boundary terms ( $\operatorname{grad} u \cdot n = h$ )
4. Pick Sobolev trial/test spaces  $V, V_0$  such that

- $a(u, u)$  is finite for  $u \in V, V_0$
- boundary conditions are satisfied ( $u = g$  in  $V \Rightarrow v = 0$  in  $V_0$ )

To fulfill the first condition, we can define the “base” space for both trial and test as  $\{v \mid a(v, v) < \infty\}$ , which is equal to  $H^1$  for the usual  $\Delta u = f$  problem. If there are extra (e.g., boundary) terms in  $a$ , try to bound these with the  $H^1$  norm.

For Neumann problems there is a **compatibility condition**. If we choose test function  $v \equiv 1$  we get the requirement

$$-\int_{\partial\Omega} h \, dS = \int_{\Omega} f \, dx$$

for the existence of solutions. Additionally, the solution of Neumann problems is unique only up to constants. To address this we can use the constrained function space

$$H_*^1(\Omega) := \left\{ v \in H^1(\Omega) : \int_{\Omega} v \, dx = 0 \right\}$$

### Theorem 1.8.0.20: Second Poincaré–Friedrichs inequality

If  $\Omega \subset \mathbb{R}^d$  is bounded and connected, then

$$\exists C = C(\Omega) > 0 : \|u\|_0 \leq C \operatorname{diam}(\Omega) \|\operatorname{grad} u\|_0 \quad \forall u \in H_*^1(\Omega) \quad (15)$$

This theorem tells us that (under some conditions), the  $L^2$  norm of functions from this space is bounded by the  $H^1$ -seminorm.

## 2 Finite Element Method

### 2.2 Galerkin Discretization

The idea of Galerkin discretization is to replace the infinite-dimensional function space  $V_0$  by a finite-dimensional subspace  $V_{0,h} \subset V_0$ .

#### Theorem 2.2.1.5: Unique solution of discrete variational problems

If the bilinear form  $a : V_0 \times V_0 \rightarrow \mathbb{R}$  is symmetric and positive definite and the linear form  $\ell : V_0 \rightarrow \mathbb{R}$  is continuous (7) w.r.t.  $\|\cdot\|_a$ , then the discrete variational problem

$$u_h \in V_{0,h} : a(u_h, v_h) = \ell(v_h), \quad \forall v_h \in V_{0,h} \quad (16)$$

has a unique *Galerkin* solution  $u_h \in V_{0,h}$  satisfying the energy estimate

$$\|u\|_a \leq \sup_{v_h \in V_{0,h}} \frac{|\ell(v_h)|}{\|v_h\|_a}$$

Recall the definition of a basis (here, superscripts are indices and not to be confused with exponents):  $\{b^1, \dots, b^N\} \subset V$  is a basis, if for every  $v \in V$  there are unique coefficients  $\mu_i$  such that  $v = \sum_{i=1}^N \mu_i b^i$  and  $N = \dim V$ . Now we can expand  $u_h = \mu_1 b^1 + \dots + \mu_N b^N$  and our goal is to find the coefficients  $\mu_i$ .

#### Galerkin Discretization

Linear discrete variational problem (16)  $\xrightarrow{\text{Choice of basis } \mathfrak{B}_h}$  Linear system of equations  $\mathbf{A}\vec{\mu} = \vec{\varphi}$

$$\text{Galerkin Matrix : } \mathbf{A} = [a(b^k, b^j)]_{j,k=1}^N \in \mathbb{R}^{N,N}$$

$$\text{RHS vector : } \vec{\varphi} = [\ell(b^j)]_{j=1}^N \in \mathbb{R}^N$$

$$\text{coefficient vector : } \vec{\mu} = [\mu_1, \dots, \mu_N]^\top \in \mathbb{R}^N$$

Note that  $\mathbf{A}_{j,k} = a(b_h^k, b_h^j) \neq a(b_h^j, b_h^k)$  if  $a$  is not symmetric (note the order of arguments  $k, j$ ). Of course the bilinear form  $a$  determines some properties of the Galerkin matrix. If  $a$  is symmetric and/or positive definite, the Galerkin matrix  $\mathbf{A}$  will have the same properties.

The choice of  $V_{0,h}$  alone determines the quality of the solution  $u_h$ . While mathematically the choice of basis  $\mathfrak{B}_h$  does not matter, for solving the equation numerically, the choice is crucial: The basis determines properties of  $\mathbf{A}$  like sparsity, so it influences how stable and efficient the numerical solution is.

#### Computing the energy norm in code

Sometimes, problems ask you to compute  $\|u_h\|_a$ , i.e., the energy norm of a discrete solution.

$$\|u_h\|_a = \sqrt{a(u_h, u_h)} = a\left(\sum_{i=0}^N \mu_i b_i^h, \sum_{j=0}^N \mu_j b_j^h\right) = \sum_{i=0}^N \sum_{j=0}^N \mu_i \mu_j a(b_i^h, b_j^h) = \vec{\mu}^\top \mathbf{A} \vec{\mu}$$

So it can be computed as `sqrt(mu.transpose() * A * mu)`.



## 2.3 Linear FEM in 1D

In FEM, the goal is to approximate  $u$  by piecewise polynomial functions.

### Mesh in one dimension

Let  $\Omega = [a, b]$ , we equip it with  $M + 1$  **nodes** resulting in the set of nodes:

$$\mathcal{V}(\mathcal{M}) = \{a = x_0 < x_1 < \dots < x_M = b\}$$

The nodes define intervals, which build up the mesh:

$$\mathcal{M} = \{]x_{j-1}, x_j[ : 1 \leq j \leq M\}$$

The intervals  $]x_{j-1}, x_j[$  are the **cells** of the mesh. We define local cell size  $h_j = |x_j - x_{j-1}|$  and global mesh width  $h_{\mathcal{M}} = \max_j h_j$ .

A simple space for continuous,  $\mathcal{M}$ -piecewise polynomial functions in  $H_0^1([a, b])$  consists of linear functions on each cell:

$$V_{0,h} = S_{1,0}^0(\mathcal{M}) = \left\{ v \in C^0([a, b]) : v|_{[x_{i-1}, x_i]} \text{ is linear, } i = 1, \dots, M, v(a) = v(b) = 0 \right\} \quad (17)$$

$$\rightarrow N = \dim S_{1,0}^0(\mathcal{M}) = M - 1$$

The 0-superscript stands for  $C^0$ : functions are globally continuous. The 1-subscript denotes local degree 1 polynomial and the 0-subscript denotes value 0 on the boundary. The  $\mathcal{S}$  stands for  $\mathcal{S}$ calar functions.

Common basis functions are the 1D tent functions:

$$b_h^j(x) = \begin{cases} (x - x_{j-1})/h_j & \text{if } x_{j-1} \leq x \leq x_j \\ (x_j - x)/h_{j+1} & \text{if } x_j \leq x \leq x_{j+1} \\ 0 & \text{else} \end{cases} \quad (18)$$

$$\rightarrow b_h^j(x_i) = \delta_{ij} \quad (19)$$

A basis satisfying condition (19) is called a **cardinal** basis. Another key property of tent functions is that their support just comprises two adjacent cells:

$$\text{supp}(b_h^j) = [x_{j-1}, x_{j+1}]$$

The support of a function  $f : \Omega \rightarrow \mathbb{R}$  is defined as  $\overline{\{x \in \Omega : f(x) \neq 0\}}$ , the set of inputs for which the function is nonzero. Since polynomials are easy to differentiate and integrate, computing (bi)linear forms  $a$  and  $\ell$  for them is quite easy.

## 2.4 Linear FEM in 2D

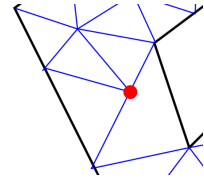
### Mesh in two dimension

Mesheres in 2D rely on **triangulations**. A **triangulation**  $\mathcal{M}$  of  $\Omega$  satisfies:

- $\mathcal{M} = \{K_i\}$ , where  $K_i$  are open triangles
- $i \neq j \rightarrow K_i \cap K_j = \emptyset$
- $\bigcup_{i=1}^M \overline{K_i} = \overline{\Omega}$
- $i \neq j \rightarrow \overline{K_i} \cap \overline{K_j}$  is either  $\emptyset$ , an edge from both triangles or a vertex from both

Again, the vertices are called **nodes** and the triangles are the **cells**.

This definition does not allow for hanging nodes because of point 4. Hanging nodes are those which lie on the edge of a triangle:



We define a space of piecewise-linear functions analogously to (17):

$$V_{0,h} = S_1^0(\mathcal{M}) = \left\{ v \in C^0(\overline{\Omega}) : v_K(x) = \alpha_K + \beta_K \cdot x, \alpha_K \in \mathbb{R}, \beta_K \in \mathbb{R}^2, x \in K \right\} \quad (20)$$

$$\rightarrow \dim S_1^0(\mathcal{M}) = \#\mathcal{V}(\mathcal{M})$$

And  $S_{1,0}^0(\mathcal{M})$  additionally requires functions to be zero on  $\partial\Omega$ , with

$$\dim S_{1,0}^0(\mathcal{M}) = \#\{x \in \mathcal{V}(\mathcal{M}) : x \notin \partial\Omega\}$$

Similarly, the 1D tent functions can be extended to 2D by requiring the cardinal property. This property is already enough since three points fully define a plane — in other words, knowing the values in three points (the vertices of a triangle) fully defines a linear function. Cardinal bases will produce sparse Galerkin matrices, as the support of a basis function only covers the triangles adjacent to its node, which can only interact with neighboring basis functions.

### Computation of Galerkin Matrix

Often, bilinear forms are defined by integration over the whole domain. But we have seen that the support of basis functions is only local. We can exploit this by performing integration only over the cells.

$$A_{ij} = a(b_h^j, b_h^i) = \sum_{K \in \text{supp}(b_h^j) \cap \text{supp}(b_h^i)} a|_K(b_h^j, b_h^i) \quad (21)$$

where  $a|_K$  is the local bilinear form over cell  $K$ .

### Cell oriented assembly

To further take advantage of (21), cell oriented assembly can be performed. Go through all cells and compute  $a|_K(b_h^j, b_h^i)$  for all pairs of basis functions associated with cell  $K$  (element matrix) and add the values to the entry at  $(i, j)$  of  $A$ .

The same procedure can be applied to calculating the right hand side vector  $\varphi$ , just that only one basis function is involved as the RHS comes from a linear functional.

## 2.5 Building Blocks of General Finite Element Methods

The first building block are meshes, see Section 2.3 and Section 2.4. Next we need to choose a space of functions.

### Definition 2.5.2.2: Multivariate Polynomials

Space of d-variate (taking inputs in  $\mathbb{R}^d$ ) polynomials with total degree  $p$ :

$$\mathcal{P}_p(\mathbb{R}^d) = \left\{ \mathbf{x} \in \mathbb{R}^d \rightarrow \sum_{\alpha \in \mathbb{N}_0^d, |\alpha| \leq p} c_\alpha \mathbf{x}^\alpha, c_\alpha \in \mathbb{R} \right\}$$

with  $\alpha = (\alpha_1, \dots, \alpha_d)$ ,  $\mathbf{x}^\alpha = x_1^{\alpha_1} \cdot \dots \cdot x_d^{\alpha_d}$  and  $|\alpha| = \alpha_1 + \dots + \alpha_d$

As an example,  $\mathcal{P}_2(\mathbb{R}^2) = \text{Span}\{1, x_1, x_2, x_1^2, x_2^2, x_1x_2\}$

### Lemma 2.5.2.5: Dimension of spaces of Polynomials

$$\dim \mathcal{P}_p(\mathbb{R}^d) = \binom{d+p}{p}, \quad p \in \mathbb{N}_0, d \in \mathbb{N}$$

which in the limit of  $p \rightarrow \infty$  behaves like  $\mathcal{O}(p^d)$ .

### Definition 2.5.2.7: Tensor Product Polynomials

The space of tensor product polynomials of degree  $p$  in each coordinate is

$$\begin{aligned} \mathcal{Q}_p(\mathbb{R}^d) &= \left\{ \mathbf{x} \in \mathbb{R}^d \rightarrow \sum_{l_1=0}^p \dots \sum_{l_d=0}^p c_{l_1, \dots, l_d} x_1^{l_1} \cdot \dots \cdot x_d^{l_d}, \text{ where } c_{l_1, \dots, l_d} \in \mathbb{R} \right\} \\ &= \text{Span}\{\mathbf{x} \rightarrow p_1(x_1) \cdot \dots \cdot p_d(x_d), p_i \in \mathcal{P}_p(\mathbb{R})\} \end{aligned}$$

As an example,  $\mathcal{Q}_2(\mathbb{R}^2) = \text{Span}\{1, x_1, x_2, x_1x_2, x_1^2, x_1^2x_2, x_1^2x_2^2, x_1x_2^2, x_2^2\}$

### Lemma 2.5.2.8: Dimension of spaces of tensor product polynomials

$$\dim \mathcal{Q}_p(\mathbb{R}^d) = (p+1)^d, \quad p \in \mathbb{N}_0, d \in \mathbb{N}$$

Finally we need **locally** supported basis functions. This basis  $\mathfrak{B}_h = b_h^1, \dots, b_h^N$  should satisfy the following constraints:

1.  $\mathfrak{B}_h$  is a basis of  $V_h$ , hence  $\dim \mathfrak{B}_h = \dim V_h$ .
2. Each  $b_h^i$  is associated with a single mesh entity (cell/edge/face/vertex).
3. Each  $b_h^i$  is only locally supported, i.e., only nonzero on cells adjacent to the associated entity.

## 2.6 Lagrangian Finite Element Methods

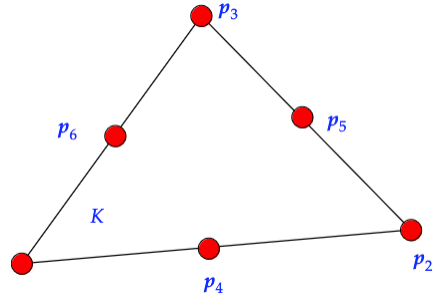
Remember Eqs. (17) and (20) as two examples of finite element spaces. They are examples of the general Lagrangian FE spaces. First we introduce them for *simplicial* meshes, i.e., those consisting of triangles (2D) or tetrahedra (3D).

### Definition 2.6.1.1: Simplicial Lagrangian finite element spaces

$$\mathcal{S}_p^0(\mathcal{M}) = \left\{ v \in C^0(\overline{\Omega}) : v|_K \in \mathcal{P}_p(K), \forall K \in \mathcal{M} \right\} \quad (22)$$

This space is well suited for triangular meshes, as the local dimension  $\binom{d+p}{p} = \binom{2+p}{p}$  (in 2D) is the same as the number of interpolation nodes of a triangle. The local basis functions of  $\mathcal{S}_1^0$  are the barycentric coordinate functions. In  $\mathcal{S}_2^0$ , the local basis functions are combinations of barycentric coordinate functions:

$$\begin{aligned} b_K^1 &= (2\lambda_1 - 1)\lambda_1, & b_K^4 &= 4\lambda_1\lambda_2, \\ b_K^2 &= (2\lambda_2 - 1)\lambda_2, & b_K^5 &= 4\lambda_2\lambda_3, \\ b_K^3 &= (2\lambda_3 - 1)\lambda_3, & b_K^6 &= 4\lambda_1\lambda_3 \end{aligned}$$



where the local basis functions 1-3 are associated with vertices and 4-6 with edges.

Analogously, the following space is well suited for quadrilaterals, as its local dimension  $(p+1)^2$  is the same as the amount of vertices and interpolation points on quads.

### Definition 2.6.2.5: Tensor product Lagrangian finite element spaces

$$\mathcal{S}_p^0(\mathcal{M}) = \left\{ v \in C^0(\overline{\Omega}) : v|_K \in \mathcal{Q}_p(K), \forall K \in \mathcal{M} \right\} \quad (23)$$

Note that the choice of local polynomial space is the only difference,  $\mathcal{Q}_p(K)$  instead of  $\mathcal{P}_p(K)$ . Of course these spaces can be mixed: on mixed (so-called hybrid) meshes, we use (22) on triangles and (23) on quadrilaterals.

## 2.7 Implementation of Finite Element Methods

Remember the principle of cell-oriented assembly. The goal is to rely mostly on local computations. To perform cell-oriented assembly, a map from local to global indices is needed. In LehrFEM++ this is the job of the dofhandler ([lf::assemble::DofHandler](#)). It provides the following main methods:

- `NumDofs()`, returns the total number of global basis functions, the dimension of the FE space.
- `NumLocalDofs(const lf::mesh::Entity &)`, returns the number of global basis functions **covering** any geometric entity, so counts those associated with the entity and with its sub-entities.
- `GlobalDofIndices(const lf::mesh::Entity &)`, returns an array of indices of the global basis function **covering** the given entity. The order is with decreasing co-dimension of the functions's associated entity, so vertex, edge, cell. Within each co-dimension, the order is according to the local numbering of the entity.
- `NumInteriorDofs(const lf::mesh::Entity &)`, returns the number of global basis functions **associated with** the given entity (so for a triangle, the number of basis functions in the interior of the triangle, not those of the edges or vertices).
- `InteriorGlobalDofIndices(const lf::mesh::Entity &)`, similar to `GlobalDofIndices` but returns only the indices of the global basis functions **associated with** the given entity.
- `Entity(unsigned int dofnum)`, returning the entity associated with the global index `dofnum`.

Instead of dimension, in LehrFEM++ the concept of **co-dimension** is used. Instead of going from a point with dimension 0 to a triangle with dimension 2, the co-dimension is the other way around. The highest-dimensional entity has co-dimension 0. This ensures that cells are always of co-dimension 0.

To assemble the Galerkin matrix, [lf::assemble::AssembleMatrixLocally](#) can be used. To use it, we need element matrix providers, constructs which provide the *local* element matrix for given bilinear forms. Some common bilinear forms are already implemented.

- $\int_K \alpha(x) \mathbf{grad} u \cdot \mathbf{grad} v \, dx$  is implemented in [lf::fe::DiffusionElementMatrixProvider](#)
- $\int_K \gamma(x) u v \, dx$  is implemented in [lf::fe::MassElementMatrixProvider](#)
- $\int_e \gamma(x) u v \, dS$  is implemented in [lf::fe::MassEdgeMatrixProvider](#). Note the integration over an edge and not a cell.
- $\int_K \alpha(x) \mathbf{grad} u \cdot \mathbf{grad} v \, dx + \int_K \gamma(x) u v \, dx$  combined is implemented in [lf::uscalfe::ReactionDiffusionElementMatrixProvider](#)
- $\int_K f(x) v \, dx$  is implemented by [lf::fe::ScalarLoadElementVectorProvider](#)
- $\int_e f(x) v \, dS$  is implemented by [lf::fe::ScalarLoadEdgeVectorProvider](#). Note again the integration over an edge.

Note that the last two are actually element vector providers.

The following formula is usually used when computing element matrices by hand:

**Lemma 2.7.5.5: Integration of powers of barycentric coordinate functions**

For a  $d$ -simplex  $K$  (line in 1D, triangle in 2D, tetrahedron in 3D) with barycentric coordinate functions  $\lambda_1, \dots, \lambda_{d+1}$

$$\int_K \lambda_1^{\alpha_1} \dots \lambda_{d+1}^{\alpha_{d+1}} d\mathbf{x} = d!|K| \frac{\alpha_1! \dots \alpha_{d+1}!}{(\alpha_1 + \dots + \alpha_{d+1} + d)!}, \quad \alpha_i \in \mathbb{N}, \quad (24)$$

where  $|K|$  is the volume of the simplex.

### Quadrature rules

$$\int_K f(x) d\mathbf{x} \approx \sum_{l=1}^{P_K} w_l^K f(\zeta_l^K), \quad w_l^K \rightarrow \text{weights}, \zeta_l^K \rightarrow (\text{quadrature}) \text{ nodes}$$

**Order** of a quadrature rule: a quad rule is of order  $q$  if

- for a simplex  $K$ , it is exact for all polynomials  $f \in \mathcal{P}_{p-1}(\mathbb{R}^d)$
- for a tensor product element  $K$ , it is exact for all polynomials  $f \in \mathcal{Q}_{p-1}(\mathbb{R}^d)$

### Essential Boundary Conditions

Remember from Section 1.7 that essential boundary conditions are Dirichlet boundary conditions, i.e.,  $u = g$  on  $\partial\Omega$ , and can be solved with the offset function trick. This trick can also be used in FEM. Assume that the basis functions are sorted such that all interior ones come first, followed by the ones on the boundary (we are free to choose the order of basis functions). Then, it is possible to write  $\mathbf{A}$  as follows:

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_0 & \mathbf{A}_{0\partial} \\ \mathbf{A}_{0\partial}^\top & \mathbf{A}_{\partial\partial} \end{bmatrix}$$

where  $\mathbf{A}_0$  is the Galerkin matrix for  $\mathcal{S}_{p,0}^0(\mathcal{M})$ , containing the interactions among interior basis functions.  $(\mathbf{A}_{0\partial})_{ij} = a(b_h^j, b_h^i)$ , where  $b_h^j$  belongs to the boundary and  $b_h^i$  to the interior, so  $\mathbf{A}_{0\partial}$  contains the interactions of interior with boundary functions. Similarly,  $\mathbf{A}_{\partial\partial}$  consists only of entries calculated from basis functions of the boundary. Then we want to solve:

$$\begin{bmatrix} \mathbf{A}_0 & \mathbf{A}_{0\partial} \\ \mathbf{A}_{0\partial}^\top & \mathbf{A}_{\partial\partial} \end{bmatrix} \begin{bmatrix} \vec{\mu}_0 \\ \vec{\mu}_\partial \end{bmatrix} = \begin{bmatrix} \vec{\varphi} \\ \vec{\varphi}_\partial \end{bmatrix}$$

where  $\vec{\mu}_\partial$  are the coefficients of the basis expansion of  $g$  on the boundary, which are known since  $g$  is given. We only need to solve for  $\vec{\mu}_0$  which results in

$$\mathbf{A}_0 \vec{\mu}_0 = \vec{\varphi} - \mathbf{A}_{0\partial} \vec{\mu}_\partial$$

This can be done in LehrFEM++ with [lf::assemble::FixFlaggedSolutionComponents](#) or [lf::assemble::FixFlaggedSolutionCompAlt](#). Both modify the matrix  $\mathbf{A}$  and RHS vector  $\vec{b}$  such that  $\mathbf{A}\vec{\mu} = \vec{b}$  has the same solution as the above equation, but do it in slightly different ways (see docs).

## Boundary data in LehrFEM

FixFlaggedSolutionComponents (in both versions) requires a lambda function `std::pair<bool, double> selector(unsigned int dof_idx)`. It returns whether the dof is to be fixed and if so, the value it should be fixed to.

We get the boundary flags with

```
auto bd_flags = lf::mesh::utils::flagEntitiesOnBoundary(dofh.Mesh(), 2);
```

The second argument is the codim of entities we are interested in. To fix all dofs on the boundary (including those associated with edges and cells), give no second argument.

### 1. One function $g$ for the whole boundary

To get the function values of  $g$  at the boundary nodes, wrap it in a MeshFunction:

```
auto mf_g = lf::mesh::utils::MeshFunctionGlobal(g);
auto boundary_val =
    lf::fe::InitEssentialConditionFromFunction(*fe_space, bd_flags, mf_g);
```

This gives a `std::vector<std::pair<bool, scalar_t>>`. To get our selector:

```
auto selector = [&](unsigned int dof_idx) -> std::pair<bool, double> {
    return boundary_val[dof_idx];
};
```

If  $g$  has different definitions on different parts of the boundary (e.g.,  $g = 0$  on  $\Gamma_0$  and  $g = 1$  on  $\Gamma_1$ ), try to express  $g$  as a lambda function with an if-else statement.

### 2. Constant value

The selector becomes

```
auto selector = [&](unsigned int dof_idx) -> std::pair<bool, double> {
    if (bd_flags[dof_idx]) {
        return std::make_pair(true, boundary_value);
    } else {
        return std::make_pair(false, 0.0); // value irrelevant
    }
};
```

### 3. BC only on part of the boundary

We need to create our own `bd_flags`. Initialize it with default value false:

```
lf::mesh::utils::AllCodimMeshDataSet<bool> bd_flags(mesh_p, false);
```

Then loop over nodes and edges *separately* and set `bd_flags` to true for the nodes/edges where the BC should be applied.

```
for (const auto& edge : fe_space->Mesh()->Entities(1)) {
    if (...) bd_flags(*edge) = true;
}
for (const auto& node : fe_space->Mesh()->Entities(2)) {...}
```

## 2.8 Parametric Finite Element Methods

### Lemma 2.7.5.14: Affine transformation of triangles

For any triangle  $K$  s.t.  $|K| > 0$ , there is a unique affine transformation  $\Phi_K(\hat{x}) = F_K \hat{x} + \tau_K$ , with  $K = \Phi_K(\hat{K})$  and  $\hat{K}$  the unit triangle.

This is nice since it allows us to transform integrals on an arbitrary triangle to ones on the unit triangle and perform easy integration there. Additionally,  $\Phi_K$  can be computed straightforwardly:

$$\text{Let } K \text{ be a triangle with vertices } a^1, a^2, a^3. \text{ Then } \Phi_K(\hat{x}) = \begin{bmatrix} a_1^2 - a_1^1 & a_1^3 - a_1^2 \\ a_2^2 - a_2^1 & a_2^3 - a_2^2 \end{bmatrix} \hat{x} + \begin{bmatrix} a_1^1 \\ a_2^1 \end{bmatrix}$$

### Definition 2.8.1.2: Pullback

Given domains  $\Omega, \hat{\Omega} \subset \mathbb{R}^d$  and a bijective mapping  $\Phi : \hat{\Omega} \rightarrow \Omega$ , the *pullback*  $\Phi^* u$  of a function  $u : \Omega \rightarrow \mathbb{R}$  is a function on  $\hat{\Omega}$  defined by  $(\Phi^* u)(\hat{x}) := u(\Phi(\hat{x}))$ ,  $\hat{x} \in \hat{\Omega}$ .

An example for this is the affine transformation of triangles (Lemma 2.7.5.14) above. Note that in the following we will use  $\hat{x}$  for an element that "lives" in a reference triangle (or quadrilateral). The reference triangle has corners  $\{(0, 0), (1, 0), (0, 1)\}$ ; the reference quad is the unit square.

When dealing with *parametric* finite element methods, we know that

$$\hat{b}_{\hat{K}}^i = \Phi_K^* b_K^i,$$

where  $b_K^i : K \rightarrow \mathbb{R}$  are the local basis functions on the concrete element  $K$  and  $\hat{b}_{\hat{K}}^i$  are the local basis functions on the reference element  $\hat{K}$  that matches  $K$ .

To make life easier, we call everything on the reference element **local** and everything on the concrete element **global**.

Note that all bilinear forms and linear forms in this course consist of integrals, so we will sooner or later need to use quadrature to approximate the integrals. But in the literature, quadrature rules are defined on reference elements. Hence, we want a change of variables in the integrals of the linear and bilinear forms such that we can apply these quadrature rules. This change of variable is given by the pullback function.

For the simple example of a mass matrix we get (by some multidimensional analysis)

$$(A_K)_{ij} = \int_K b_K^j(x) b_K^i(x) dx = \int_{\hat{K}} (\Phi_K^* b_K^j)(\hat{x}) (\Phi_K^* b_K^i)(\hat{x}) \sqrt{\det(D\Phi_K^\top(\hat{x}) D\Phi_K(\hat{x}))} d\hat{x}, \quad (25)$$

where we have simply applied the pullback to both functions, switched to integration on  $\hat{K}$ , and added an integration element. Similarly, for the diffusion matrix we get

$$\begin{aligned} \int_K \mathbf{grad}_x b_K^j(x) \cdot \mathbf{grad}_x b_K^i(x) dx \\ = \int_{\hat{K}} \Phi_K^* (\mathbf{grad}_x b_K^j)(\hat{x}) \cdot \Phi_K^* (\mathbf{grad}_x b_K^i)(\hat{x}) \sqrt{\det(D\Phi_K^\top(\hat{x}) D\Phi_K(\hat{x}))} d\hat{x}. \end{aligned} \quad (26)$$

by pulling back the global gradients (denoted by  $\mathbf{grad}_x$ ) of the global shape functions.



Why is the integration element  $\det(D\Phi_K^\top(\hat{x})D\Phi_K(\hat{x}))$  more complicated than what we know from Theorem 0.3.2.31? It is needed when  $\Omega$  and  $\hat{\Omega}$  do not live in the same space. If, for example,  $\Omega \subset \mathbb{R}^3$  describes a 2D plane “living” in 3D space and  $\hat{\Omega} \subset \mathbb{R}^2$ , we will have  $D\Phi_K \in \mathbb{R}^{3 \times 2}$  and hence  $\det(\Phi_K)$  is not defined, so the full term is used. In the simpler case of square  $\Phi$ , this simplifies to  $\det(\Phi_K)$ . For both cases, the term is provided by [lf::geometry::IntegrationElement](#).

The next thing which needs clarification is  $\Phi_K^*(\mathbf{grad}_x b_K^i)$ , pullback of the gradient. As the gradient  $\mathbf{grad}_x b_K^i(x)$  depends on the shape of  $K$ , it would be complicated to compute this term directly.

Therefore we use

**Lemma 2.8.3.10: Transformation formula for gradients**

For differentiable  $u : K \rightarrow \mathbb{R}$  and any diffeomorphism  $\Phi : \hat{K} \rightarrow K$  we have

$$(\mathbf{grad}_{\hat{x}} \Phi^* u)(\hat{x}) = (D\Phi(\hat{x}))^\top \underbrace{[(\mathbf{grad}_x u)(\Phi(\hat{x}))]}_{(\Phi^* \mathbf{grad}_x)u}(\hat{x})$$

In words: on the left side we have **local gradient of the pullback** of  $u$ , on the right side we have **pullback of the global gradient** of  $u$ .

This is exactly what we need since the pullback of the global gradient occurs in (26) and the local gradient of local shape functions is easy to compute.

Now we need to get rid of  $(D\Phi(\hat{x}))^\top$  on the right. Let's define  $J = D\Phi(\hat{x})$ . Since  $J$  might not be square, we cannot simply invert it. But we can multiply by  $J(J^\top J)^{-1}$  to get

$$\Phi_K^*(\mathbf{grad}_{\hat{x}} b_K^i) = J(J^\top J)^{-1}(\mathbf{grad}_{\hat{x}} \hat{b}_{\hat{K}}^i)$$

$\hat{b}(\hat{x})$  is a reference basis function on the reference shape — we can compute its gradient easily by hand.

Note that in the script, the case of square  $J$  is assumed, such that this simplifies to  $J^{-\top}(\mathbf{grad}_{\hat{x}} \hat{b}_{\hat{K}}^i)$ . The long term only matters when we have (as above)  $\Omega$  and  $\hat{\Omega}$  not living in the same space. The same remedy applies here, you can in any case use [lf::geometry::JacobianverseGramian](#) which will return  $J^{-\top}$  or  $J(J^\top J)^{-1}$ , respectively.

**Bilinear Transformation for Quadrilaterals** Let  $\{a^1, a^2, a^3, a^4\}$  be the ordered corners of a quadrilateral. Then

$$\begin{aligned} \Phi_K(\hat{x}) &= (1 - \hat{x}_1)(1 - \hat{x}_2)a^1 + \hat{x}_1(1 - \hat{x}_2)a^2 + (1 - \hat{x}_1)\hat{x}_2a^3 + (1 - \hat{x}_1)\hat{x}_2a^4 \\ &= \begin{bmatrix} \alpha_1 + \beta_1\hat{x}_1 + \gamma_1\hat{x}_2 + \delta_1\hat{x}_1\hat{x}_2 \\ \alpha_2 + \beta_2\hat{x}_1 + \gamma_2\hat{x}_2 + \delta_2\hat{x}_1\hat{x}_2 \end{bmatrix} \end{aligned}$$

with

$$\begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = a^1, \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} = a^2 - a^1, \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix} = a^4 - a^1, \begin{bmatrix} \delta_1 \\ \delta_2 \end{bmatrix} = a^4 - a^3 - a^2 + a^1$$

## 3 FEM: Convergence and Accuracy

### 3.1 Abstract Galerkin Error Estimates

The main takeaway here is that the solution given by the Galerkin method is the best one (with respect to the energy norm) in the chosen discrete subspace. This is formalized by the following lemma:

#### Theorem 3.1.3.7: Cea's Lemma

Under some assumptions that guarantee the existence of a unique solution we have

$$\|u - u_h\|_a = \inf_{v_h \in V_{0,h}} \|u - v_h\|_a, \quad (27)$$

where  $u$  is the exact solution and  $u_h$  is the Galerkin solution.

Next we want to discuss the types of refinement, i.e., the steps we can take to increase the accuracy of our method.

#### h-Refinement

Replace the mesh  $\mathcal{M}$  (underlying  $V_{0,h}$ ) with a finer mesh  $\mathcal{M}'$  (underlying larger discrete trial space  $V'_{0,N'}$ ).

#### p-Refinement

Replace  $V_{0,h} := S_p^0(\mathcal{M})$ ,  $p \in \mathbb{N}$ , with  $V'_{0,h} := S_{p+1}^0(\mathcal{M})$ , yielding a larger space:  $V_{0,h} \subset V'_{0,h}$

So h-refinement refines the mesh (smaller and smaller triangles) and p-refinement chooses more powerful basis functions (start with linear, then quadratic, etc.). The h in h-refinement refers to the mesh width:

#### Definition 3.2.1.4: Mesh width

Given a mesh  $\mathcal{M} = \{K\}$ , the **mesh width**  $h_{\mathcal{M}}$  is defined as

$$h_{\mathcal{M}} := \max\{\text{diam } K : K \in \mathcal{M}\}$$
$$\text{diam } K := \max\{|p - q| : p, q \in K\}$$

### 3.2 Empirical (Asymptotic) Convergence of Lagrangian FEM

As in NumCSE, there are two types of convergence, algebraic and exponential. We refer to the number of basis functions (dimension of the trial space) as  $N$  and we study the behavior of errors as  $N \rightarrow \infty$ . Note that both exercises and theorems are often posed in terms of  $h_{\mathcal{M}}$  as  $h \rightarrow 0$ , which is equivalent if we have a fixed polynomial degree  $p$  and only do h-refinement.

**Definition 3.2.2.1: Types of convergence**

$\|u - u_N\|_a = \mathcal{O}(N^{-\alpha})$ ,  $\alpha > 0$  is called **algebraic** convergence with rate  $\alpha$ .

$\|u - u_N\|_a = \mathcal{O}(\exp(-\gamma N^\delta))$ ,  $\gamma, \delta > 0$  is called **exponential** convergence.

**Determining convergence rates**• **algebraic**

$$\alpha \approx \frac{\log \epsilon_{i-1} - \log \epsilon_i}{\log N_i - \log N_{i-1}} = \frac{\log(\epsilon_i / \epsilon_{i-1})}{\log(h_i / h_{i-1})}$$

• **exponential** In general: complicated (see §3.2.2.5 Lecture Notes)

If  $\delta = 1$  (plain exponential convergence) and we let  $N$  increase linearly (e.g.,  $N_i = i$ ),

$$\frac{\epsilon_{i+1}}{\epsilon_i} \approx \exp(-\gamma)$$

Note that in the case of h-refinement we get the relation between  $N$  and  $h$  given by

**Equation 3.2.2.1**

$$N = \dim S_p^0(\mathcal{M}) \approx p^d h_{\mathcal{M}}^{-d} \implies \frac{h_{\mathcal{M}}}{p} \approx N^{-\frac{1}{d}} \quad (28)$$

where  $p$  are the dimensions of the local basis functions, e.g.,  $p = 1$  for linear basis functions, and  $d$  is the dimension of the underlying space  $\Omega$ .

E.g., in the case where we have  $\Omega \subset \mathbb{R}^2$  and piecewise linear basis functions, we get  $h_{\mathcal{M}}^{-2} \approx N$ .

**3.3 A Priori (Asymptotic) Finite Element Error Estimates**

Since FEM is similar to polynomial interpolation, we can use interpolation error estimates to get error bounds. Here are some results for linear interpolation:

**Linear interpolation error 1D**

Using the linear interpolant  $I_1$  we want to study the interpolation error  $u - I_1 u$ . The following interpolation error estimates can be used for sufficiently smooth functions  $u$ :

$$\|u - I_1 u\|_{L^\infty([a,b])} \leq \frac{1}{4} h_{\mathcal{M}}^2 \|u''\|_{L^\infty([a,b])} \quad (29)$$

$$\|u - I_1 u\|_{L^2([a,b])} \leq h_{\mathcal{M}}^2 \|u''\|_{L^2([a,b])} \quad (30)$$

$$|u - I_1 u|_{H^1([a,b])} \leq h_{\mathcal{M}} \|u''\|_{L^2([a,b])} \quad (31)$$

### Linear interpolation error 2D

In 2D, linear interpolation corresponds to using tent functions:  $I_1 u = \sum_{p \in \mathcal{V}(\mathcal{M})} u(p) b^p$ , where  $b^p$  is the tent function associated with point  $p$ .

$$\begin{aligned} \|u - I_1 u\|_{L^2(\Omega)} &\leq C h_{\mathcal{M}}^2 \left\| \|D^2 u\|_F \right\|_{L^2(\Omega)} \\ \|\mathbf{grad}(u - I_1 u)\|_{L^2(\Omega)} &\leq C \rho_{\mathcal{M}} h_{\mathcal{M}} \left\| \|D^2 u\|_F \right\|_{L^2(\Omega)} \end{aligned} \quad (32)$$

Here,  $D^2 u$  is the Hessian of  $u$ ,  $\|\cdot\|_F$  is the Frobenius norm, and  $\rho_{\mathcal{M}}$  is the shape regularity measure of the mesh  $\mathcal{M}$ , defined as  $\rho_{\mathcal{M}} = \max_{K \in \mathcal{M}} \frac{h_{\mathcal{M}}^2}{|K|}$  for a triangular mesh.

To get rid of this cumbersome notation, we can introduce more Sobolev spaces.

#### Definition 3.3.3.1: Higher order Sobolev spaces/norms

The  $m$ -th order Sobolev norm is defined as

$$\|u\|_{H^m(\Omega)}^2 = \sum_{k=0}^m \sum_{\alpha \in \mathbb{N}^d, |\alpha|=k} \int_{\Omega} |D^{\alpha} u|^2 dx, \quad \text{where } D^{\alpha} u = \frac{\partial^{|\alpha|} u}{\partial x_1^{\alpha_1} \dots \partial x_d^{\alpha_d}} \quad (33)$$

Hence we can define the  $m$ -th Sobolev space as

$$H^m(\Omega) = \left\{ v : \Omega \rightarrow \mathbb{R} : \|v\|_{H^m(\Omega)} < \infty \right\} \quad (34)$$

#### Definition 3.3.3.3: Higher order Sobolev semi-norms

The  $m$ -th order Sobolev semi-norm is defined as

$$|u|_{H^m(\Omega)}^2 = \sum_{\alpha \in \mathbb{N}^d, |\alpha|=m} \int_{\Omega} |D^{\alpha} u|^2 dx \quad (35)$$

Remember the multi-index  $\alpha$  already seen in Definition 2.5.2.2. Using this new notation, we can rewrite the error bounds from (32) as

$$\begin{aligned} \|u - I_1 u\|_{L^2(\Omega)} &\leq C h_{\mathcal{M}}^2 |u|_{H^2(\Omega)} \\ \|\mathbf{grad}(u - I_1 u)\|_{L^2(\Omega)} &\leq C \rho_{\mathcal{M}} h_{\mathcal{M}} |u|_{H^2(\Omega)} \end{aligned}$$

It turns out that these bounds are not sharp. There is a very useful result for Lagrangian finite elements:

**Theorem 3.3.5.6: Best approximation error estimates for Lagrangian finite elements**

Given a triangular mesh  $\mathcal{M}$ , if the true solution  $u$  is in  $H^k(\Omega)$ , the best approximation error is bounded by

$$\inf_{v_h \in \mathcal{S}_p^0(\mathcal{M})} \|u - v_h\|_{H^1(\Omega)} \leq C \left( \frac{h_{\mathcal{M}}}{p} \right)^{\min\{p, k-1\}} \|u\|_{H^k(\Omega)} \quad (36)$$

We might not know the constant  $C$  and/or  $\|u\|_{H^k(\Omega)}$ , but we know  $p$  and  $h_{\mathcal{M}}$  as they are imposed by the choice of function space and mesh. Remember the concept of refinement: we can adjust these values. And from Eq. (28) we know that  $h_{\mathcal{M}}/p \approx N^{-\frac{1}{d}}$ . Hence the error displays **algebraic** convergence with rate  $\min\{p, k-1\}/d$ . What still remains a question is  $k$ , the smoothness of the solution  $u$ .

Another useful result hidden in the lecture notes is

**Equation 3.6.3.10 ( $L^2$  estimate)**

Under some assumptions (convex domain, smooth coefficient functions), we have

$$\|u - u_h\|_{L^2(\Omega)} \leq C \frac{h_{\mathcal{M}}}{p} \|u - u_h\|_{H^1(\Omega)} \quad (37)$$

So we gain one order of convergence in the  $L^2$  norm compared to the  $H^1$  norm.

**Rules of thumb for converge**

If we are using  $\mathcal{S}_p^0(\mathcal{M})$  and  $u$  is sufficiently smooth (e.g.,  $u \in C^\infty(\Omega)$ ), we have

$$\|u - u_h\|_{H^1(\Omega)} = \mathcal{O}(h^p)$$

$$|u - u_h|_{H^1(\Omega)} = \mathcal{O}(h^p)$$

$$\|u - u_h\|_{L^2(\Omega)} = \mathcal{O}(h^{p+1})$$

**3.4 Elliptic regularity****Theorem 3.4.0.2: Smooth elliptic lifting theorem**

For domains  $\Omega$  with smooth boundaries  $\partial\Omega$ , i.e. no corners and sufficiently smooth  $\sigma$ , if

$$u \in H_0^1(\Omega) \quad \text{and} \quad -\operatorname{div}(\sigma \operatorname{grad} u) \in H^k(\Omega)$$

or

$$u \in H^1(\Omega), -\operatorname{div}(\sigma \operatorname{grad} u) \in H^k(\Omega) \quad \text{and} \quad \operatorname{grad} u \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega$$

holds, then  $u \in H^{k+2}(\Omega)$  and

$$\|u\|_{H^{k+2}(\Omega)} \leq C \|\operatorname{div}(\sigma \operatorname{grad} u)\|_{H^k(\Omega)}$$

This tells us that when solving the typical PDE  $-\operatorname{div}(\sigma \operatorname{grad} u) = f$  and the source term  $f$  is in  $H^k(\Omega)$ , the solution  $u$  will be in  $H^{k+2}(\Omega)$  (of course under the right assumptions).

The theorem requires smooth domains, but our meshes will have corners, so what can be done there? As long as the domain and all cells are convex, something similar still holds:

**Theorem 3.4.0.10: Elliptic lifting on convex domains**

If  $\Omega \subset \mathbb{R}^d$  is convex,  $u \in H_0^1(\Omega)$  and  $\Delta u \in L^2(\Omega)$ , then  $u \in H^2(\Omega)$ .

If we are solving the Laplace equation on a convex domain, we just need to check if  $f \in L^2(\Omega)$ , since  $-\Delta u = f$ .

**Finding  $k$  in Theorem 3.3.5.6**

- Often, the exercise gives you  $u$  which is in  $C^\infty(\Omega)$ , which means  $u$  is infinitely smooth. In this case,  $k = \infty$ .
- Sometimes, you may use elliptic regularity results like Theorem 3.4.0.10 to find  $k = 2$ , for example.

### 3.5 Variational Crimes

What are variational crimes?

A variational crime is committed when instead of the true variational problem

$$u_h \in V_{0,h} : a(u_h, v_h) = \ell(v_h), \forall v_h \in V_{0,h},$$

we solve a different, “perturbed” variational problem

$$\tilde{u}_h \in V_{0,h} : a_h(\tilde{u}_h, v_h) = \ell_h(v_h), \forall v_h \in V_{0,h}$$

with modified (bi-)linear forms  $a_h, \ell_h$ . With computers, the use of quadrature and approximation of boundaries result in such a crime and are unavoidable. As Hiptmair likes to say, “we are all sinners”.

So the only distinction we can make is between acceptable and unacceptable “crimes”. Crimes which do not affect the type and rate of convergence are acceptable.

So how to not temper with the convergence?

- if  $\|u - u_h\|_1 \in \mathcal{O}(h_{\mathcal{M}}^p)$ , use a quadrature rule of order at least  $2p - 1$
- if  $V_{0,h} = \mathcal{S}_p^0(\mathcal{M})$  then approximate the boundary with polynomials of degree  $p$

### 3.6 FEM: Duality Techniques for Error Estimation

**Theorem 3.6.1.7: Duality estimate for linear functional output**

Given a functional  $F : V_0 \rightarrow \mathbb{R}$  the dual solution  $g_F$  solves

$$g_F \in V_0 : a(g_F, v) = F(v), \forall v \in V_0$$

and we get the estimate

$$|F(u) - F(u_h)| \leq C \|u - u_h\|_a \inf_{v_h \in V_{0,h}} \|g_F - v_h\|_a$$

Why is this useful? If  $g_F$  can be approximated well in  $V_{0,h}$ , then the output error  $|F(u) - F(u_h)|$  can converge to 0 much faster than  $\|u - u_h\|_a$ .

## 5 Non-Linear Elliptic Boundary Value Problems

### 5.1 Elastic String Model

We want to derive the general variational equation for an elastic string. For this, one approximates the string as  $n$  point masses affected by gravity connected with springs, whose energy behaves according to Hooke's law. Then, one takes the limit  $n \rightarrow \infty$  to derive a continuous model. The total energy is then just given by the sum of elastic and gravitational energies – given positions  $(\mu_0, \dots, \mu_n)$  and  $x_i = a + hi$ , assuming the spring constants are 1:

**Total energy of discrete spring system**

$$E(\mu) = \frac{1}{2} \sum_{i=0}^n \left( \sqrt{h^2 + (\mu_{i+1} - \mu_i)^2} \right)^2 + \sum_{i=1}^n m_i \mu_i g$$

Then, the equilibrium position for this model can be found by minimizing this expression over  $\mu$ . A continuous model is derived by replacing the discrete positions  $\mu_i$  by a function  $u(x_i)$ , and the mass by a mass density. Then, performing some manipulations, one obtains

**Total energy for the continuous string model**

$$J_s(u) = \int_a^b \frac{1}{2} \frac{b-a}{L} \sigma(x) \left( \sqrt{1 + |u'(x)|^2} - \frac{L}{b-a} \right)^2 + g\rho(x)u(x) \, dx$$

where  $\sigma$  is the spring stiffness,  $L$  is the total rest length of the string, and  $\rho$  is the mass density. The term  $\sqrt{1 + |u'(x)|^2}$  is the length of the spring at position  $x$ .

In a similar fashion, a membrane model can be derived by assuming a two-dimensional grid of springs containing points masses and taking the limit  $n \rightarrow \infty$  springs. The energy then becomes

**Total energy for the membrane model**

$$J_M(u) = \int_{\Omega} \frac{1}{2L} \sigma(x) \left( \left( \sqrt{1 + \left| \frac{\partial u}{\partial x_1}(x) \right|^2} - \frac{L}{b-a} \right)^2 + \left( \sqrt{1 + \left| \frac{\partial u}{\partial x_2}(x) \right|^2} - \frac{L}{b-a} \right)^2 \right) + g\rho(x)u(x) \, dx$$

In the limit of a taut membrane, i.e. for  $L \ll b - a$ , these equations just reduce to the problem of minimizing the familiar functionals seen in earlier chapters:



### Elastic string/membrane: taut limit

$$\begin{aligned} J_s(u) &= \frac{1}{2} \int_a^b \hat{\sigma}(x) |u'(x)|^2 + g\rho(x)u(x) \, dx \\ J_M(u) &= \frac{1}{2} \int_{\Omega} \hat{\sigma}(x) \|\mathbf{grad} u(x)\|^2 + g\rho(x)u(x) \, dx \end{aligned} \quad (38)$$

## 5.2 Calculus of Variations

The difference between the equations seen in (38), which hold in the limit of a very stretched string/membrane, and the general equations given above, is that the former are **quadratic** minimization problems, while the latter are **nonlinear** minimization problems. The theory used so far mapped quadratic minimization problems to **linear** variational problems, which were then discretized. The new equations, however, yield nonlinear variational equations. Therefore, more general variational problems need to be derived.

The idea employed is that, for a minimizer  $u_*$  of  $J(u)$ , every perturbation  $J(u_* + v)$  would be larger than  $J(u_*)$ . This means that  $f(t) = J(u_* + tv)$  has a minimum at  $t = 0$  for every function  $v$ :

### Theorem 5.2.1.5: Characterization of global minimizers

Assume  $u_* \in V_0$  is a global minimizer of  $J(u)$ , i.e.

$$u_* = \operatorname{argmin}_{u \in V_0} J(u)$$

Then, if  $\varphi_v(t) = J(u_* + tv)$  is differentiable in  $t = 0$ , we have

$$\frac{d\varphi_v}{dt}(0) = 0 \quad \forall v \in V_0$$

This means that nonlinear variational equations can be derived by computing this derivative for an arbitrary  $v$ . As an example, for the elastic string model introduced in the last sub-chapter, this yields

### Variational equations for elastic string model

$$\int_a^b \frac{\sigma(x)}{c} \left( \sqrt{1 + |u'(x)|^2} - c \right) \frac{u'(x)v'(x)}{\sqrt{1 + |u'(x)|^2}} + g\rho(x)v(x) \, dx = 0 \quad \forall v \in H_0^1([a, b])$$

This can be formulated more generally as a **general variational equation**:

### General variational equation

A general, nonlinear variational equation reads

$$u \in \widehat{V} : a(u; v) = 0 \quad \forall v \in V_0$$

Where  $a$  is **linear in the second argument**  $v$  and  $V_0, \widehat{V}$  are function spaces.

## 5.3 Nonlinear Boundary Value Problems

Similarly to the linear case, nonlinear PDEs can be derived from the variational equations by “stripping” of the derivatives of  $v$  by partial integration and employing the fundamental lemma of calculus of variations. As an example, for the string model, this yields for  $u(a) = u_a, u(b) = u_b$

$$\frac{d}{dx} \left( \frac{\sigma(x)}{c} \left( \sqrt{1 + |u'(x)|^2} - c \right) \frac{u'(x)}{\sqrt{1 + |u'(x)|^2}} \right) = g\rho(x) \quad \text{in } ]a, b[$$

## 5.4 Galerkin Discretization of Non-Linear BVPs

The idea of Galerkin discretization for non-linear variational equations is exactly the same as for linear equations, but they yield nonlinear systems of equations instead of linear systems of equations: One restricts  $u$  and  $v$  to a finite function space  $u_h \in \widehat{V}_h$  and  $v_h \in V_{0,h}$ , and expands the functions in some basis of the space. This, then, leads to nonlinear equations for the basis expansion coefficients. These equations could be solved directly by employing some fixed-point iteration seen in NumCSE.

### Nonlinear Galerkin Discretization

Given a variational problem  $a(u; v) = 0 \quad \forall v \in V_{0,h}$ , the Galerkin discretization reads

$$(F(\mu))_i = a \left( u_{0,h} + \sum_{j=1}^N \mu_j b_h^j; b_h^i \right), \quad i = 1, \dots, N$$

where  $b_h^i$  are fixed basis functions,  $\mu_i$  are the basis function coefficients and  $u_{0,h}$  contains Dirichlet boundary conditions.

Another option is to already linearize the continuous problem, and then discretize it to derive linear systems of equations. This is done by employing Newton’s method in function space: The conventional Newton iteration is given as

$$\xi^{(k+1)} = \xi^{(k)} - DF(\xi^{(k)})^{-1} F(\xi^{(k)})$$

Replacing the vector  $\xi$  with a function  $u$  and the derivative by a function derivative now gives

#### Equation 5.3.2.5 (Functional Newton iteration)

$$\begin{aligned} w \in V_0 : a(u^{(k)}; v) + D_u a(u^{(k)}; v)w &= 0 \quad \forall v \in V_0 \\ u^{(k+1)} &= u^{(k)} + w \end{aligned}$$

Here, the directional derivative is defined as

$$D_u a(u^{(k)}; v)w = \lim_{t \rightarrow 0} \frac{a(u + tw; v) - a(u; v)}{t}, \quad u^{(k)} \in \widehat{V}, \quad v, w \in V_0$$

#### Computing the directional derivative

You can simply compute the defining limit by Taylor-expanding  $a(u + tw, v)$ . If you do this once, you notice that only the first-order term in  $t$  remains: the 0th order term cancels out with  $a(u, v)$ , and the higher-order terms vanish in the limit  $t \rightarrow 0$ .

So a more convenient method is to do Taylor expansion and to keep only the first-order term.

Now, the advantage of this equation for  $w$  is that the functional derivative is linear, i.e.  $(v, w) \mapsto D_u a(u^{(k)}; v)w$  is a **bilinear form**. Now, one can employ Galerkin discretization for the linear problem in  $w$ , exactly like it was done in Chapter 2 and 3. The final equations then read

#### Equation 5.3.3.6 (Nonlinear Newton equations for variational problems)

$$\begin{aligned} w_h \in V_{0,h}^{(k)} : D_u a(u_h^{(k-1)}; v_h)w_h &= -a(u_h^{(k-1)}; v_h) \quad \forall v_h \in V_{0,h}^{(k)} \\ u_h^{(k)} &= P_h^{(k)}(u_h^{(k-1)} + w_h) \end{aligned}$$

Here, different function spaces can be used for each iterations, so a projector  $P_h^{(k)}$  needs to be used to project the solution from  $V_h^{(k-1)}$  to  $V_h^{(k)}$ . In all of these equations, the previous iterate  $u_h^{(k-1)}$  is kept fixed, a linear system like derived in Chapter 2 is solved to obtain the intermediate  $w_h$ , and then a new iterate  $u^{(k)}$  is obtained.

## 9 Second-Order Linear Evolution Problems

### 9.2 Parabolic Initial-Boundary Value Problems

#### 9.2.1 Heat Equation

In local form the heat equation is given by

$$\frac{\partial}{\partial t}(\rho u) - \operatorname{div}(\kappa(\mathbf{x}) \mathbf{grad} u) = f \quad \text{in } \tilde{\Omega} = \Omega \times ]0, T[ \quad (39)$$

where  $u$  is the temperature,  $\rho$  the heat capacity,  $\kappa$  the heat conductivity and  $f$  a (time-dependent) heat source/sink. Without the time derivative, this looks very similar to standard PDEs, for which we know the transformation into a nice variational problem.

To solve it we still need boundary conditions. Besides the boundary conditions of the spatial domain — which are now required for all times — one also needs initial conditions over the whole domain at time 0.

$$\begin{aligned} u(\mathbf{x}, t) &= g(\mathbf{x}, t) \quad \text{for } (\mathbf{x}, t) \in \partial\Omega \times ]0, T[ \\ u(\mathbf{x}, 0) &= u_0(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega \end{aligned} \quad (40)$$

Testing with time-independent test functions  $v$  and assuming  $\rho$  to be time independent as well, we get to

$$\begin{aligned} \int_{\Omega} \rho(\mathbf{x}) \dot{u} v \, d\mathbf{x} + \int_{\Omega} \kappa(\mathbf{x}) \mathbf{grad} u \cdot \mathbf{grad} v \, d\mathbf{x} &= \int_{\Omega} f(\mathbf{x}, t) v \, d\mathbf{x} \quad \forall v \in H_0^1(\Omega) \\ u(\mathbf{x}, 0) &= u_0(\mathbf{x}) \in H_0^1(\Omega) \end{aligned} \quad (41)$$

with the shorthand notation

$$\begin{aligned} m(\dot{u}, v) &= \int_{\Omega} \rho(\mathbf{x}) \dot{u} v \, d\mathbf{x} \\ a(u, v) &= \int_{\Omega} \kappa(\mathbf{x}) \mathbf{grad} u \cdot \mathbf{grad} v \, d\mathbf{x} \\ \ell(v) &= \int_{\Omega} f(\mathbf{x}, t) v \, d\mathbf{x} \end{aligned}$$

and the realization that  $m(\dot{u}, v) = \frac{d}{dt}m(u, v)$  as only  $u$  depends on time (note that importantly, the domain  $\Omega$  also stays constant) we can rewrite (41) as

$$\frac{d}{dt}m(u, v) + a(u, v) = \ell(v) \quad (42)$$

which looks like something we know how to solve from NumCSE.

### 9.2.3 Stability

#### Lemma 9.2.3.8.: Decay of solutions of parabolic evolutions

If  $f \equiv 0$ , the solution  $u(t)$  of the heat equation (41) satisfies

$$\|u(t)\|_m \leq e^{-\gamma t} \|u_0\|_m, \quad \|u(t)\|_a \leq e^{-\gamma t} \|u_0\|_a \quad \forall t \in ]0, T[$$

where  $\gamma = \text{diam}(\Omega)^{-2}$ .

Note that this lemma also tells us that if  $f$  is time-independent, the solution  $u(t)$  converges exponentially (in time) to the stationary solution (the solution of (42) without the  $m(\cdot, \cdot)$  part).

### 9.2.4 Method of Lines

Now let's look into how we can solve (42). Let's apply the Galerkin discretization. As  $u$  is now also time-dependent, we let the coefficients of  $u$  (and not the basis) depend on time.

$$u_h(t) = \sum_{i=1}^N \mu_i(t) b_h^i$$

Combining this with (42), we get

$$M \left\{ \frac{d}{dt} \vec{\mu}(t) \right\} + A \vec{\mu}(t) = \vec{\varphi}(t) \quad (43)$$

$$\vec{\mu}(0) = \vec{\mu}_0$$

where  $M_{i,j} = m(b_h^j, b_h^i)$ ,  $A_{i,j} = a(b_h^j, b_h^i)$  and  $[\vec{\varphi}(t)]_i = \ell(b_h^i)$ .

This is now an ODE with respect to time and can be solved by time stepping, learned in NumCSE.

**Recall ODEs** An ODE is given as

$$\dot{\mathbf{u}} = \mathbf{f}(t, \mathbf{u})$$

and is called linear if  $\mathbf{f}(t, \mathbf{u}) = A(t)\mathbf{u}$ . There is an evolution operator associated with the ODE, defined as  $\Phi^{t_0, t} u_0 = u(t)$ . There are some methods to approximate the evaluation operator with a discrete evolution operator  $\Psi$ .

- explicit Euler:  $\Psi^{t, t+\tau} \mathbf{u} = \mathbf{u} + \tau \mathbf{f}(t, \mathbf{u})$
- implicit Euler:  $\Psi^{t, t+\tau} \mathbf{u} = \mathbf{w}, \mathbf{w} = \mathbf{u} + \tau \mathbf{f}(t + \tau, \mathbf{w})$
- implicit midpoint:  $\Psi^{t, t+\tau} \mathbf{u} = \mathbf{w}, \mathbf{w} = \mathbf{u} + \tau \mathbf{f}(t + \frac{1}{2}\tau, \frac{1}{2}(\mathbf{w} + \mathbf{u}))$

Hence we can calculate the time evolution by the sequence

$$\mathbf{u}^{(0)} = \mathbf{u}_0, \quad \mathbf{u}^{(j)} = \Psi^{t_{j-1}, t_j} \mathbf{u}^{(j-1)}, \quad j = 1, \dots, M$$

As  $\Psi$  is the discrete approximation, the question about the error is immediate. One usually considers

- the error at final time:  $\epsilon_M = \|\mathbf{u}(T) - \mathbf{u}^{(M)}\|$
- maximum error in the sequence:  $\epsilon_\infty = \max_j \|\mathbf{u}^{(j)} - \mathbf{u}(t_j)\|$

**Theorem 9.2.6.14: Convergence of single-step methods**

Given the above sequence of solutions, obtained by a single step method of order  $q \in \mathbb{N}$ , then

$$\epsilon_\infty = \max_j \|u^{(j)} - u(t_j)\| \leq C\tau^q$$

with  $\tau = \max_j |t_j - t_{j-1}|$ .

**Runge–Kutta Single-Step Methods****Definition 7.3.3.1: General Runge–Kutta single-step method**

For coefficients  $b_i, a_{i,j} \in \mathbb{R}, c_i = \sum_{j=1}^s a_{i,j}$ , the discrete evolution operator  $\Psi^{s,t}$  of an **s-stage Runge–Kutta single step method** (RK-SSM) for the ODE  $\dot{u} = f(t, u)$  is defined by

$$k_i = f\left(t + c_i\tau, u + \tau \sum_{j=1}^s a_{i,j}k_j\right), \quad i = 1, \dots, s, \quad \Psi^{t,t+\tau}u = u + \tau \sum_{j=1}^s b_jk_j$$

with  $k_j$  the increments.

The RK-SSM methods can be written down in compact form (the butcher scheme) as

$$\begin{array}{c|c} \mathbf{c} & \mathbf{A} \\ \hline & \mathbf{b} \end{array} \quad (44)$$

where  $\mathbf{c}$  is a vector containing the coefficients  $c_i$ ,  $\mathbf{b}$  the coefficients  $b_i$  and  $\mathbf{A}$  a matrix containing the coefficients  $a_{i,j}$ .

So continuing from (43) with different time stepping schemes, we get

- explicit Euler:

$$\vec{\mu}^{(j)} = \vec{\mu}^{(j-1)} + \tau_j \mathbf{M}^{-1}(\vec{\varphi}(t_{j-1}) - \mathbf{A}\vec{\mu}^{(j-1)})$$

- implicit Euler:

$$\vec{\mu}^{(j)} = (\tau_j \mathbf{A} + \mathbf{M})^{-1}(\mathbf{M}\vec{\mu}^{(j-1)} + \tau_j \vec{\varphi}(t_{j-1}))$$

- implicit midpoint (Crank-Nicolson)

$$\vec{\mu}^{(j)} = \left(\mathbf{M} + \frac{1}{2}\mathbf{A}\right)^{-1} \tau_j \left( \left(\mathbf{M} - \frac{1}{2}\mathbf{A}\right) \vec{\mu}^{(j-1)} + \frac{1}{2}(\vec{\varphi}(t_j) + \vec{\varphi}(t_{j-1})) \right)$$

These all involve solving a linear system of equations each time step. However note, that the matrices to invert stay constant with respect to time, so we can calculate the decomposition only once to save a lot of time.

Using a general RK-SSM method as the time step, we get the following system of equations

$$M\vec{\kappa}_i + \sum_{m=1}^s \tau a_{i,m} A\vec{\kappa}_m = \vec{\varphi}(t_j + c_i\tau) - A\vec{\mu}^{(j)}$$

$$\vec{\mu}^{(j+1)} = \vec{\mu}^{(j)} + \tau \sum_{m=1}^s b_m \vec{\kappa}_m$$

With the Kronecker product, this can be rewritten as

$$(I_s \otimes M + \tau \mathbf{A} \otimes A) \begin{bmatrix} \vec{\kappa}_1 \\ \vdots \\ \vec{\kappa}_s \end{bmatrix} = \begin{bmatrix} \vec{\varphi}(t_j + c_1\tau) - A\vec{\mu}^{(j)} \\ \vdots \\ \vec{\varphi}(t_j + c_s\tau) - A\vec{\mu}^{(j)} \end{bmatrix}$$

which can be used to solve for the increments  $\vec{\kappa}_i$ .

Recall stiff initial value problems:

### Stiffness

An initial value problem is called stiff if stability imposes much tighter timestep constraints on explicit single step methods than the accuracy requirements.

To study the stiffness of the method of lines, we first diagonalize it. For (43) let  $\vec{\psi}_1, \dots, \vec{\psi}_N$  denote the  $N$  linearly independent generalized eigenvectors satisfying

$$A\vec{\psi}_i = \lambda_i M\vec{\psi}_i, \quad (\vec{\psi}_j)^\top M\vec{\psi}_i = \delta_{ij}$$

with positive eigenvalues  $\lambda_i$ . With  $T = [\vec{\psi}_1, \dots, \vec{\psi}_N]$  and  $D = \text{diag}(\lambda_1, \dots, \lambda_N)$ , this can be rewritten as

$$AT = MTD, \quad T^\top MT = I$$

The existence of eigenvectors with positive eigenvalues is guaranteed, as  $A, M$  are positive (semi)definite. Thus with a change of basis to the eigenvector basis, one can diagonalize (43).

$$\begin{aligned} \vec{\mu}(t) &= \sum_k \eta_k(t) \vec{\psi}_k \Leftrightarrow \vec{\mu}(t) = T\vec{\eta}(t) \Leftrightarrow T^\top M\vec{\mu}(t) = \vec{\eta}(t) \\ &\rightarrow MT \frac{d}{dt} \vec{\eta}(t) + MTD\vec{\eta}(t) = \vec{\varphi}(t) \\ &\rightarrow \frac{d}{dt} \vec{\eta}(t) + D\vec{\eta}(t) = T^\top \vec{\varphi}(t) \end{aligned}$$

As  $D$  is diagonal, this amounts to  $N$  decoupled scalar ODEs. On those, we can perform our analysis more easily. In NumCSE you have seen that both Euler and Crank-Nicolson can be rewritten as a RK-SSM with appropriate coefficients, so we can study the stability of the general RK-SSM for the scalar case. For  $\dot{u} = -\lambda u$ , with the butcher scheme (44) we obtain  $\Psi_\lambda^{t, t+\tau} u = S(-\lambda\tau)u$ , with the stability function

$$S(z) = 1 + z\mathbf{b}^\top (I - z\mathbf{A})^{-1} \mathbf{1} = \frac{\det(I - z\mathbf{A} + z\mathbf{b}\mathbf{1}^\top)}{\det(I - z\mathbf{A})}$$

**Unconditional stability of single step methods** A necessary condition for unconditional stability of a single step method, is that the discrete evolution operator  $\Psi_\lambda^t$  applied to the scalar ODE  $\dot{u} = -\lambda u$  satisfies

$$\lambda > 0 \rightarrow \lim_{j \rightarrow \infty} (\Psi_\lambda^\tau)^j u_0 = 0, \quad \forall u_0, \forall \tau > 0$$

**Definition 9.2.7.46: L-stability**

An RK-SSM satisfying the above condition, is called L-stable if its stability function satisfies

$$|S(z)| < 1, \forall z < 0 \quad \text{and} \quad S(-\infty) = \lim_{z \rightarrow -\infty} S(z) = 0$$

Plugging  $-\infty$  into  $S$  we obtain  $S(-\infty) = 1 - \mathbf{b}^\top \mathfrak{A}^{-1} \mathbf{1}$ , which is equal to zero if  $\mathbf{b}$  is equal to the last row of  $\mathfrak{A}$ .

**Theorem 9.2.8.5: Meta-theorem — Convergence of fully discrete evolution**

Assume that

- the solution of the parabolic IBVP is “sufficiently smooth”
- its spatial Galerkin finite element discretization relies on degree  $p$  Lagrangian finite elements on uniformly shape-regular families of meshes
- time stepping is based on an L-stable single step method of order  $q$ ,

then we can expect an asymptotic behavior of the total discretization error according to

$$\left( \tau \sum_{j=1}^M \left| u(\tau j) - u_h^{(j)} \right|_{H^1(\Omega)}^2 \right)^{\frac{1}{2}} \leq C(h_{\mathcal{M}}^p + \tau^q)$$

Hence the total error is the spatial error plus the temporal error.

### 9.3 Linear wave equations

In local form, the (linear) wave equation is given by

$$\rho(\mathbf{x}) \frac{\partial^2}{\partial t^2} u - \operatorname{div}(\sigma(\mathbf{x}) \mathbf{grad} u) = f \quad \text{in } \tilde{\Omega} \quad (45)$$

Note the similarity to the heat equation (39). Since the wave equation is a second order ODE  $\ddot{u} = \mathbf{f}(\mathbf{u})$ , two initial conditions are needed. In addition to the initial conditions (40), the initial velocity

$$\frac{\partial}{\partial t} u(\mathbf{x}, 0) = v_0(\mathbf{x}) \quad \text{for all } \mathbf{x} \in \Omega$$

is also needed.

We want to use the time stepping schemes we already know. To apply them, the wave function can be converted into a first order ODE:

$$\begin{aligned} \dot{u} &= v \\ \rho \dot{v} &= \operatorname{div}(\sigma(\mathbf{x}) \mathbf{grad} u) \end{aligned}$$



Remember from Analysis that the particular wave equation  $\frac{\partial^2}{\partial t^2}u - c^2 \frac{\partial^2}{\partial x^2}u = 0$  in 1D results in the d'Alembert solution:

$$u(x, t) = \frac{1}{2}(u_0(x + ct) + u_0(x - ct)) + \frac{1}{2c} \int_{x-ct}^{x+ct} v_0(s) ds$$

with  $u_0$  and  $v_0$  the initial conditions. Hence there is again the concept of domain of dependence and domain of influence. This will be important later. Furthermore, in the absence of a source term, as in the simple case above, the solution will stay undamped. This corresponds to *conservation of total energy*.

We can formulate the variational problem:

$$m(\ddot{u}, v) + a(u, v) = 0 \quad \forall v \in V_0 \quad (46)$$

#### Theorem 9.3.2.16: Energy conservation in wave propagation

If  $u$  solves (46), then its energy is conserved, in the sense that

$$\frac{1}{2}m\left(\frac{\partial}{\partial t}u, \frac{\partial}{\partial t}u\right) + \frac{1}{2}a(u, u) \equiv \text{const}$$

where  $\frac{1}{2}m\left(\frac{\partial}{\partial t}u, \frac{\partial}{\partial t}u\right)$  can be understood as the 'kinetic' energy and  $\frac{1}{2}a(u, u)$  as the 'potential' energy.

### 9.3.3 Method of Lines

The method of lines gives rise to

$$M \left\{ \frac{d^2}{dt^2} \vec{\mu}(t) \right\} + A \vec{\mu}(t) = \vec{\varphi}(t)$$

$$\vec{\mu}(0) = \vec{\mu}_0, \frac{d}{dt} \vec{\mu}(0) = \vec{v}_0$$

Using  $\vec{v} = \dot{\vec{\mu}}$ , we can rewrite it to be a first order ODE.

$$\frac{d}{dt} \vec{\mu} = \vec{v}$$

$$M \frac{d}{dt} \vec{v} = \vec{\varphi}(t) - A \vec{\mu}$$

$$\vec{\mu}(0) = \vec{\mu}_0, \vec{v}(0) = \vec{v}_0$$

Remember that in the case of  $\vec{\varphi} \equiv 0$ , energy is conserved:

$$E_h(t) = \frac{1}{2} \frac{d}{dt} \vec{\mu}^\top M \frac{d}{dt} \vec{\mu} + \frac{1}{2} \vec{\mu}^\top A \vec{\mu} \equiv \text{const}$$

So we would like a time stepping that preserves this. Such time stepping schemes are called *structure preserving*. One such timestepping scheme is the **Crank–Nicolson** one:

$$M \frac{\vec{\mu}^{(j+1)} - 2\vec{\mu}^{(j)} + \vec{\mu}^{(j-1)}}{\tau^2} = -\frac{1}{2} A (\vec{\mu}^{(j-1)} + \vec{\mu}^{(j+1)}) + \frac{1}{2} \left( \vec{\varphi}\left(t_j - \frac{1}{2}\tau\right) + \vec{\varphi}\left(t_j + \frac{1}{2}\tau\right) \right)$$

Another one would be the **Störmer scheme**:

$$M \frac{\vec{\mu}^{(j+1)} - 2\vec{\mu}^{(j)} + \vec{\mu}^{(j-1)}}{\tau^2} = -\mathbf{A}\vec{\mu}^{(j)} + \vec{\varphi}(t_j)$$

For both of these *second order* time stepping schemes, we need  $\vec{\mu}^{(-1)}$  to get  $\vec{\mu}^{(1)}$ . Now the question is, where do we get this from? It can be obtained with a special initial step, using a symmetric (first order) difference quotient:

$$\frac{d}{dt}\vec{\mu}(0) = \vec{v}_0 \rightarrow \frac{\mu^{(1)} - \mu^{(-1)}}{2\tau} = \vec{v}_0$$

And finally there is the **Leapfrog** timestepping. Using the auxiliary variable  $\vec{v}^{(j+1/2)} = \frac{\vec{\mu}^{(j+1)} - \vec{\mu}^{(j)}}{\tau}$  and inserting this into the Störmer scheme results in

$$M \frac{\vec{v}^{(j+1)} - \vec{v}^{(j)}}{\tau} = -\mathbf{A}\vec{\mu}^{(j)} + \vec{\varphi}(t_j)$$

$$\frac{\vec{\mu}^{(j+1)} - \vec{\mu}^{(j)}}{\tau} = \vec{v}^{(j+1/2)}$$

with the initial step  $\vec{v}^{(-1/2)} + \vec{v}^{(1/2)} = 2\vec{v}_0$ .

## 10 Convection-Diffusion Problems

### 10.1 Heat conduction in a Fluid

Consider a flowing fluid. Then there is the key quantity, the *flow field*  $\mathbf{v} : \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}^d$ , where  $d$  is the dimension we consider. The flow field can be understood as  $\mathbf{v}(\mathbf{x}) =$  fluid velocity at point  $\mathbf{x} \in \Omega$ .

Given a flow field  $\mathbf{v}$ , we can consider the autonomous initial value problem

$$\dot{\mathbf{y}} = \mathbf{v}(\mathbf{y}), \quad \mathbf{y}_0 = \mathbf{x}_0$$

The solution  $t \rightarrow \mathbf{y}(t)$  describes how a particle moves, carried by the fluid, also called *streamline*.

As the domain  $\Omega$  is usually bounded, we cannot have fluid leaving the domain. This means the fluid velocity must be zero in the normal direction of the domain boundary. Hence

$$\mathbf{v}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) = 0, \quad \forall \mathbf{x} \in \partial\Omega$$

#### Fourier's law in a moving fluid

$$\mathbf{j}(\mathbf{x}) = -\kappa \mathbf{grad} u(\mathbf{x}) + \mathbf{v}(\mathbf{x})\rho u(\mathbf{x})$$

with  $\kappa$  the heat conductivity and  $\rho$  the volumetric heat capacity. We already know the first part, called diffusive heat flux, from the heat equation. The second part is called convective heat flux. With this new flux, the standard PDE becomes

$$-\operatorname{div}(\kappa \mathbf{grad} u) + \operatorname{div}(\rho \mathbf{v}(\mathbf{x})u) = f \quad \text{in } \Omega \quad (47)$$

#### Incompressible Fluids

A fluid is called incompressible if its associated flow map (evaluation operator)  $\Phi^t$  is volume preserving. This is the case iff.

$$\operatorname{div} \mathbf{v} \equiv 0 \quad \text{in } \Omega$$

Hence the fluid is incompressible if its flow velocity is divergence-free.

In case of incompressibility, equation (47) can be simplified using the general product rule Lemma 1.5.2.1 and  $\operatorname{div} \mathbf{v} = 0$

$$-\operatorname{div}(\kappa \mathbf{grad} u) + \mathbf{v}(\mathbf{x}) \cdot \mathbf{grad} \rho u = f \quad \text{in } \Omega$$

We can also look at the time-dependent heat flow, which is similar to the heat equation (39)

$$\frac{\partial}{\partial t}(\rho u) - \operatorname{div}(\kappa \mathbf{grad} u) + \operatorname{div}(\rho \mathbf{v}(\mathbf{x}, t) u) = f \quad \text{in } \Omega \quad (48)$$

We will see later on how to solve this.

## 10.2 Stationary Convection-Diffusion Problems

Here we will focus on the convection diffusion equation (47) with constant  $\kappa$ ,  $\rho$ , incompressible flow  $\mathbf{v}$  and zero Dirichlet boundary conditions. Hence

$$-\kappa \Delta u + \rho \mathbf{v}(\mathbf{x}) \cdot \mathbf{grad} u = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \partial\Omega$$

Non-dimensionalizing the problem results in

$$-\epsilon \Delta u + \mathbf{v}(\mathbf{x}) \cdot \mathbf{grad} u = f \quad \text{in } \Omega, \quad u = 0 \quad \text{on } \partial\Omega$$

with  $\|\mathbf{v}\|_{L^\infty(\Omega)} = 1$ . This results in the following variational form

$$\epsilon \int_{\Omega} \mathbf{grad} u \cdot \mathbf{grad} w \, d\mathbf{x} + \int_{\Omega} (\mathbf{v} \cdot \mathbf{grad} u) w \, d\mathbf{x} = \int_{\Omega} f(\mathbf{x}) w \, d\mathbf{x} \quad (49)$$

with the left-hand side the bilinear form  $a(u, w)$ . However,  $a$  is not symmetric. This also means that it does not induce an energy norm. However, it is still positive definite (see lecture document).

### 10.2.1 Singular perturbation

A boundary value problem depending on a parameter  $\epsilon$  is called singularly perturbed, if the limit problem for  $\epsilon \rightarrow \epsilon_0$  is not compatible with the boundary conditions.

For  $\epsilon = 0$ , the above PDE is singular perturbed. It cannot satisfy Dirichlet boundary conditions on the *outflow* part of the boundary.

$$\begin{aligned} \Gamma_{\text{out}} &= \mathbf{x} \in \partial\Omega : \mathbf{v}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) > 0 \\ \Gamma_{\text{in}} &= \mathbf{x} \in \partial\Omega : \mathbf{v}(\mathbf{x}) \cdot \mathbf{n}(\mathbf{x}) < 0 \end{aligned}$$

### 10.2.2 Upwinding

When trying to solve Eq. (49) with the Galerkin approach, when  $\epsilon$  is very close to 0, one can observe huge oscillations in the solution, which is not correct. It comes from the fact that the Galerkin matrix becomes close to singular. So our goal is to get a robust method that can solve Eq. (49) no matter the  $\epsilon$ .

Consider again Eq. (49) but in  $d = 1$  and with zero boundary conditions

$$\epsilon \int_0^1 \frac{\partial u}{\partial x} \frac{\partial w}{\partial x} \, dx + \int_0^1 \frac{\partial u}{\partial x} w \, dx = \int_0^1 f(x) w \, dx$$

To calculate the Galerkin matrix for an equidistant mesh with  $M$  cells, we use the global composite trapezoidal rule for the convective term

$$\int_0^1 \psi(x) \, dx \approx h \sum_{j=0}^{M-1} \psi(jh)$$

Hence the convective term of the bilinear form will be approximated by

$$\int_0^1 \frac{\partial u_h}{\partial x} w_h \, dx \approx h \sum_{j=1}^{M-1} \frac{\partial u_h}{\partial x}(jh) w_h(jh)$$

But  $\frac{\partial u_h}{\partial x}(jh)$  is not defined:  $u_h \in \mathcal{S}_{1,0}^0$  is piecewise continuous, so its derivative has jumps at the nodes. However, convection transports the information in the direction of  $\mathbf{v}$  ( $= 1$  in our case). So we use the value of  $\frac{\partial u_h}{\partial x}$  from “upwind”, i.e., in direction  $-\mathbf{v}$ , which here is the value in the previous cell:

$$\frac{\partial u_h}{\partial x}(jh) = \lim_{\delta \rightarrow 0} \frac{\partial u_h}{\partial x}(jh - \delta v) = \frac{\partial u_h}{\partial x} \Big|_{x_{j-1}, x_j[}$$

And generalized in more dimensions

$$\mathbf{v}(\mathbf{p}) \cdot \mathbf{grad} u_h(\mathbf{p}) = \lim_{\delta \rightarrow 0} \mathbf{v}(\mathbf{p}) \cdot \mathbf{grad} u_h(\mathbf{p} - \delta \mathbf{v}(\mathbf{p}))$$

### 10.2.2.2 Streamline Diffusion

A totally different idea for fixing the problem of  $\epsilon \rightarrow 0$  is to add some  $h$ -dependent diffusion. I.e., replace  $\epsilon \leftarrow \epsilon + c(h)$  with  $c(h) > 0$ . However, there is smearing in the internal layers. But the solution is smooth along the direction of  $\mathbf{v}$ , so adding diffusion along the velocity should not do any harm.

The method of *Anisotropic diffusion* is born. On cell  $K$ , replace  $\epsilon \leftarrow \epsilon \mathbf{I} + \delta_K \mathbf{v}_K \mathbf{v}_K^\top$ . Here,  $\mathbf{v}_K$  is the local velocity, obtained by averaging over the vertices, and  $\delta_K > 0$  some controlling parameter.

$$\int_{\Omega} (\epsilon \mathbf{I} + \delta_K \mathbf{v}_K \mathbf{v}_K^\top) \mathbf{grad} u \cdot \mathbf{grad} w \, dx + \int_{\Omega} (\mathbf{v} \cdot \mathbf{grad} u) w \, dx = \int_{\Omega} f(x) w \, dx$$

However this affects the solution  $u$ , such that it will not be the same as the one from Eq. (49). To get rid of this inconsistency, the anisotropic diffusion can be introduced via a residual term

$$\begin{aligned} & \int_{\Omega} \epsilon \mathbf{grad} u \cdot \mathbf{grad} w \, dx + \int_{\Omega} (\mathbf{v} \cdot \mathbf{grad} u) w \, dx \\ & + \sum_{K \in \mathcal{M}} \delta_K \int_K (-\epsilon \Delta u + \mathbf{v} \cdot \mathbf{grad} u - f)(\mathbf{v} \cdot \mathbf{grad} w) = \int_{\Omega} f(x) w \, dx \end{aligned}$$

the added term will be zero for the exact solution (strong PDE) and the anisotropic diffusion is still here. The control parameter is usually chosen according to

$$\delta_K = \begin{cases} \epsilon^{-1} h_K^2 & \text{if } \|\mathbf{v}\|_{K,\infty} h_K \leq 2\epsilon \\ h_K & \text{if } \|\mathbf{v}\|_{K,\infty} h_K > 2\epsilon \end{cases}$$

With this, the  $\mathcal{O}(h_{\mathcal{M}}^2)$  convergence of  $\|u - u_h\|_{L^2(\Omega)}$  for  $h$ -refinement is preserved, while upwind quadrature only achieves  $\mathcal{O}(h_{\mathcal{M}})$  convergence.

### 10.3 Discretization of Time-Dependent Convection-Diffusion IBVPs

Now we will take a look at how time-dependent (also called *transient*) convection-diffusion can be modeled. Assuming the incompressibility condition and non-dimensionalizing, (48) becomes

$$\frac{\partial u}{\partial t} - \epsilon \Delta u + \mathbf{v}(\mathbf{x}, t) \cdot \mathbf{grad} u = f \quad \text{in } \Omega \quad (50)$$

If we solve this with the method of lines (Section 9.3.3) without upwind quadrature, oscillations occur. But with upwinding, we get damping, which is also wrong. Therefore, we need a different method. Again, we need to take care of the limit  $\epsilon \rightarrow 0$ . So let's first look at the pure transport problem ( $\epsilon = 0$ ):

$$\frac{\partial u}{\partial t} + \mathbf{v}(\mathbf{x}, t) \cdot \mathbf{grad} u = f \quad \text{in } \Omega$$

Its solution is given by the *Method of Characteristics*

$$u(\mathbf{x}, t) = \begin{cases} u_0(\mathbf{x}_0) + \int_0^t f(\mathbf{y}(s), s) ds & \text{if } \mathbf{y}(s) \in \Omega \quad \forall 0 < s < t \\ g(\mathbf{y}(s_0), s_0) + \int_{s_0}^t f(\mathbf{y}(s), s) ds & \text{if } \mathbf{y}(s_0) \in \partial\Omega \wedge \mathbf{y}(s) \in \Omega \quad \forall s_0 < s < t \end{cases} \quad (51)$$

where  $\mathbf{y}(t)$  is defined as the solution of  $\dot{\mathbf{y}}(t) = \mathbf{v}(\mathbf{y}(t), t)$  and describes a streamline.

We have the initial condition  $u_0$  on  $\Omega$  and the Dirichlet boundary condition  $g$  on the inflow boundary (see (40) for definitions). This method works as follows: We follow a streamline backwards in time starting at  $\mathbf{x}$ .

Then, the solution consists of two parts: The initial value of  $u$  at  $t = 0$  and the integral of the source function  $f$  from 0 to  $t$  along the streamline. This is described by the first case for a streamline that stays inside the domain. The second case is for a streamline that leaves the domain. Here,  $s_0$  is the time at which the streamline intersects the boundary, and we use the boundary value  $g$  as the initial value at time  $s_0$ .

#### Splitting Methods

Given a general ODE whose right-hand side is the sum of two functions

$$\dot{\mathbf{y}} = \mathbf{g}(t, \mathbf{y}) + \mathbf{r}(t, \mathbf{y})$$

The Strang Splitting single step method provides a method to solve this.

#### Strang Splitting

Compute  $\mathbf{y}^{(j+1)}$  given  $\mathbf{y}^{(j)}$  according to

$$\tilde{\mathbf{y}} = \mathbf{z}\left(t_j + \frac{1}{2}\tau\right), \quad \text{where } \mathbf{z}(t) \text{ solves } \dot{\mathbf{z}} = \mathbf{g}(t, \mathbf{z}), \mathbf{z}(t_j) = \mathbf{y}^{(j)}$$

$$\hat{\mathbf{y}} = \mathbf{w}(t_{j+1}), \quad \text{where } \mathbf{w}(t) \text{ solves } \dot{\mathbf{z}} = \mathbf{r}(t, \mathbf{w}), \mathbf{w}(t_j) = \tilde{\mathbf{y}}$$

$$\mathbf{y}^{(j+1)} = \mathbf{z}(t_{j+1}), \quad \text{where } \mathbf{z}(t) \text{ solves } \dot{\mathbf{z}} = \mathbf{g}(t, \mathbf{z}), \mathbf{z}\left(t_j + \frac{1}{2}\tau\right) = \hat{\mathbf{y}}$$

and  $t_{j+1} = t_j + \tau$

**Theorem 10.3.3.5: Order of Strang splitting single step method**

If the IVP in each sub-step is solved exactly or with a 2nd-order time stepping method, the Strang splitting method is of second order.

We can now apply this to Eq. (50).

$$\begin{array}{ccccc} \frac{\partial}{\partial t} u = & \epsilon \Delta u & f - \mathbf{v} \cdot \mathbf{grad} u & & \\ \updownarrow & \updownarrow & \updownarrow & & \\ \dot{\mathbf{y}} = & \mathbf{g}(\mathbf{y}) & \mathbf{r}(\mathbf{y}) & & \end{array}$$

This amounts to once solving pure diffusion

$$\frac{\partial t}{\partial z} - \epsilon \Delta z = 0$$

and once pure transport

$$\frac{\partial t}{\partial w} + \mathbf{v} \cdot \mathbf{grad} u = f$$

To solve the pure transport problem, we have seen the method of characteristics (51). However, it requires integration along streamlines. We can approximate these integrals by “following” particles along their path defined by the velocity field. This **particle method** works as follows:

1. Pick suitable interpolation nodes  $\mathbf{p}_i$ , the initial particle positions.
2. Solve initial value problems

$$\dot{\mathbf{y}}(t) = \mathbf{y}(\mathbf{v}(t), t) \quad , \quad \mathbf{y}(0) = \mathbf{p}_i$$

with suitable single step methods.

3. Reconstruct the approximation. E.g., with the composite midpoint rule,

$$u_h^{(j)}(\mathbf{p}_i^{(j)}) = u_0(\mathbf{p}_i) + \tau \sum_{l=1}^{j-1} f\left(\frac{1}{2}(\mathbf{p}_i^l + \mathbf{p}_i^{l-1}), \frac{1}{2}(t_l + t_{l-1})\right)$$

Basically, we have an interpolation problem with nodes that change over time. At each step, we need to add particles at the inflow boundary and remove ones that leave the domain. This means we need to re-mesh in each step, i.e., create a new triangular mesh with the advected nodes/particles.

### 10.3.4 Semi-Lagrangian Method

Check out this [video](#) for an intuitive explanation of the method in 1D.

A method which relies on a *fixed* mesh is the Semi-Lagrangian method.

**Definition 10.3.4.2: Material derivative**

Given a velocity field  $\mathbf{v}$ , the material derivative of a function  $f$  is given by

$$\frac{Du}{D\mathbf{v}}(\mathbf{x}, t_0) = \lim_{\tau \rightarrow 0} \frac{u(\mathbf{x}, t_0) - u(\Phi^{t_0, t_0 - \tau} \mathbf{x}, t_0 - \tau)}{\tau}$$

By the chain rule we find

$$\frac{Du}{D\mathbf{v}}(\mathbf{x}, t) = \mathbf{grad} u(\mathbf{x}, t) \cdot \mathbf{v}(\mathbf{x}, t) + \frac{\partial}{\partial t} u(\mathbf{x}, t)$$

Hence the transient convection-diffusion equation (50) can be rewritten as

$$\frac{Du}{D\mathbf{v}} - \epsilon \Delta u = f \quad \text{in } \Omega$$

We will now approximate the material derivative by a backwards difference. It is easy to interpret that expression: the total change of  $u$  in a timestep for a particle is the difference of the current  $u^{(j)}$  at the particle's current position  $\mathbf{x}$  minus the value of the previous function  $u^{(j-1)}$  at the particle's old position  $\Phi^{t_j, t_j - \tau} \mathbf{x}$ . We get a semi-discretization

$$\frac{u^{(j)}(\mathbf{x}) - u^{(j-1)}(\Phi^{t_j, t_j - \tau} \mathbf{x})}{\tau} - \epsilon \Delta u^{(j)} = f(\mathbf{x}, t_j) \quad \text{in } \Omega$$

with additional initial conditions for  $t = t_j$ . We can apply the standard Galerkin method to this semi-discretization.

$$\int_{\Omega} \frac{u^{(j)}(\mathbf{x}) - u^{(j-1)}(\Phi^{t_j, t_j - \tau} \mathbf{x})}{\tau} w(\mathbf{x}) \, d\mathbf{x} + \epsilon \int_{\Omega} \mathbf{grad} u^{(j)} \cdot \mathbf{grad} w \, d\mathbf{x} = \int_{\Omega} f(\mathbf{x}, t_j) w(\mathbf{x}) \, d\mathbf{x}$$

Unfortunately this cannot be implemented as is because the function  $u^{(j-1)}(\Phi^{t_j, t_j - \tau} \mathbf{x})$  is not a finite element function on  $\mathcal{M}$ .

To get around this, we do two things: Replace  $\Phi^{t_j, t_j - \tau} \mathbf{x}$  by  $\mathbf{x} - \tau \mathbf{v}(\mathbf{x}, t_j)$  (an explicit Euler step) and replace  $u^{(j-1)}(\Phi^{t_j, t_j - \tau} \mathbf{x})$  by its linear interpolant  $I_1(u^{(j-1)}(\Phi^{t_j, t_j - \tau} \mathbf{x}))$ .

$$\int_{\Omega} \frac{u^{(j)}(\mathbf{x}) - I_1\{u^{(j-1)}(\mathbf{s} - \tau \mathbf{v}(\mathbf{s}, t_j))\}(\mathbf{x})}{\tau} w(\mathbf{x}) \, d\mathbf{x} + \epsilon \int_{\Omega} \mathbf{grad} u^{(j)} \cdot \mathbf{grad} w \, d\mathbf{x} = \int_{\Omega} f(\mathbf{x}, t_j) w(\mathbf{x}) \, d\mathbf{x}$$

Here,  $\mathbf{s}$  is the variable on which the interpolation operator  $I_1$  acts (in the lecture notes, this variable is written as a dot “ $\cdot$ ”). The above equation still looks quite complicated, but can be implemented.



# 11 Numerical Methods for Conservation Laws

## 11.2 Scalar Conservation Laws in 1D

The goal of this chapter is to solve general conservation laws which are of the form

$$\frac{\partial u}{\partial t}(x, t) + \frac{\partial}{\partial x}(f(u(x, t), x)) = s(u(x, t), x, t). \quad (52)$$

The flux  $f : \mathbb{R} \times \Omega \rightarrow \mathbb{R}$  can be a general function, which can depend non-linearly on the solution  $u$ . Everything in this chapter will be one-dimensional in space and time, so we have  $\Omega \subseteq \mathbb{R}$ .

We usually consider the special case  $s = 0$  and  $f = f(u)$ , known as the **Cauchy problem**:

$$\begin{aligned} \frac{\partial u}{\partial t} + \frac{\partial}{\partial x} f(u) &= 0 \\ u(x, 0) &= u_0(x) \end{aligned} \quad (53)$$

### Example: Particle Model

We model cars as particles with position  $x_i(t)$ . Traffic speed is modeled by

$$\dot{x}_i(t) = v_{\max} \cdot \left(1 - \frac{\Delta_0}{\Delta x_i(t)}\right), \quad \Delta x_i(t) = x_{i+1}(t) - x_i(t),$$

where  $v_{\max}$  is the maximum velocity and  $\Delta_0$  is the minimal distance between cars (i.e., the car length). Using some basic assumptions, we get a PDE:

$$\frac{\partial u}{\partial t}(x, t) = \frac{\partial}{\partial x} u(1 - u).$$

We define the car density  $u(x)$  as the number of cars in an infinitesimal interval around  $x$ .

### 11.2.2 Characteristics

We consider the Cauchy problem (53). Then a characteristic curve is defined as

#### Definition 11.2.2.3: Characteristic curve for 1D scalar conservation law

$\Gamma : [0, T] \rightarrow \mathbb{R} \times [0, T]$  with  $\Gamma(\tau) := (\gamma(\tau), \tau)$ , such that  $\gamma$  satisfies

$$\frac{d}{d\tau} \gamma(\tau) = f'(u(\gamma(\tau), \tau))$$

for  $0 \leq \tau \leq T$ .

Generally, characteristic curves are lines along which information propagates. This means  $u(x, t)$  will only depend on the initial condition at  $x_0$ , namely  $u_0(x_0)$ , if there is a characteristic curve that starts in the point  $x_0$  and travels to the spacetime point  $(x, t)$ . One property of characteristic curves is the following:

**Lemma 11.2.2.6: Classical solution and characteristic curves**

Smooth solutions of (53) are constant along characteristic curves.

For example in the case of linear advection ( $\partial_t u + v \partial_x u = 0$ ) we can use this to solve the equation because  $\gamma(\tau) = (x_0 + \tau v)$ , which implies  $u(x, t) = u_0(x - tv)$ .

But this does not work if the solution is not smooth. For example, in the traffic flow model above, the solution has a jump after a certain time and hence this approach breaks down.

**11.2.4 Jump conditions and Riemann Problem**

The method of characteristics usually only works up to a certain point in time. To get the solution for times after that, we first note that the solution will usually have a discontinuity after the time where the method of characteristics breaks down. So we study how the solution behaves at these jumps (discontinuities).

The setting is as follows: We still study the Cauchy problem (53). We can derive that along jumps, the normal components must be continuous, which leads to the

**Definition 11.2.4.2: Rankine–Hugoniot (jump) condition**

$$\dot{s}(u_l - u_r) = f(u_l) - f(u_r)$$

where  $\dot{s} = \frac{d\gamma}{dt}$  is the time derivative of the discontinuity curve  $\Gamma(t) = (\gamma(t), t) \in \mathbb{R} \times [0, T]$ .

Note that this is useful because it allows us to compute the jump if we know  $u_l$  and  $u_r$ .

The Riemann problem is given as

**Definition 11.2.5.1: Riemann problem**

$$\frac{\partial u}{\partial t} + \frac{\partial f(u)}{\partial x} = 0$$

and

$$u_0(x) = \begin{cases} u_l \in \mathbb{R} & \text{if } x < 0 \\ u_r \in \mathbb{R} & \text{if } x > 0 \end{cases}$$

Note that  $f$  can still be chosen to be any sufficiently smooth flux function.

Using the Rankine–Hugoniot jump condition we then get the following solution for Riemann problems with a shock:

**Lemma 11.2.5.4: Shock solution for Riemann problem**

For any two states  $u_r, u_l \in \mathbb{R}$  the piecewise constant function

$$u(x, t) := \begin{cases} u_l & \text{for } x < \dot{s}t \\ u_r & \text{for } x > \dot{s}t \end{cases} \quad \dot{s} := \frac{f(u_l) - f(u_r)}{u_l - u_r}, x \in \mathbb{R}, 0 < t < T$$

is a weak solution to the Riemann problem.

Note that the solution only holds if the equation implies a shock (jump). This is the case if  $f$  is convex and  $u_l > u_r$  or if  $f$  is concave and  $u_r > u_l$ .

If the jump only exists in the beginning we have a different solution

**Lemma 11.2.5.4: Rarefaction solution for Riemann problem**

If  $f \in C^2(\mathbb{R})$  is strictly  $\begin{cases} \text{convex and } u_l < u_r \\ \text{concave and } u_r < u_l \end{cases}$ , then

$$u(x, t) := \begin{cases} u_l & \text{for } x < \min\{f'(u_l), f'(u_r)\} \cdot t \\ g\left(\frac{x}{t}\right) & \text{for } \min\{f'(u_l), f'(u_r)\} < \frac{x}{t} < \max\{f'(u_l), f'(u_r)\} \\ u_r & \text{for } x > \max\{f'(u_l), f'(u_r)\} \cdot t \end{cases}$$

is a weak solution to the Riemann problem with  $g := (f')^{-1}$ .

The question when to choose which of the two solution is answered by

**Definition 11.2.6.1: Lax entropy condition**

Let  $u$  be a weak solution of (53), and a piecewise classical solution in neighborhood of  $C^2$ -curve  $\Gamma := (\gamma(\tau), \tau), 0 \leq \tau \leq T$ , discontinuous across  $\Gamma$ .

$u$  satisfies the *Lax entropy condition* in  $(x_0, t_0) \in \Gamma$  iff.

$$f'(u_l) > \underbrace{\frac{f(u_l) - f(u_r)}{u_l - u_r}}_{\dot{s}} > f'(u_r)$$

Now if  $u$  satisfies the Lax entropy condition, then we have to pick the shock solution. Otherwise we pick the rarefaction solution.

**11.2.7 Properties of Entropy Solutions**

The essential properties here are that with the propagation speed  $f'(u)$ , we find the **domain of dependence** and the **domain of influence**, which is best illustrated by a picture and hence we encourage the reader to look at the lecture document and the illustrations below Theorem 11.2.7.3.

Moreover the second result is that the number of extrema of the solution is non-increasing in time.

## 11.3 Conservative Finite Volume (FV) Discretization

### 11.3.1 Finite Difference Methods

Finite difference methods are probably the simplest methods for solving PDEs. We just replace the spacial derivatives by some finite difference quotient for example one of the following

Symmetric difference quotient	$\frac{\partial f}{\partial x} \approx \frac{f(x_0 + h) - f(x_0 - h)}{2h}$
Backward difference quotient	$\frac{\partial f}{\partial x} \approx \frac{f(x_0) - f(x_0 - h)}{h}$
Forward difference quotient	$\frac{\partial f}{\partial x} \approx \frac{f(x_0 + h) - f(x_0)}{h}$

Then we construct a solution by time stepping: given  $u(x, t_k)$  we compute  $u(x, t_{k+1})$  by using some Runge–Kutta integrator.

With the example of finite difference methods, we observe that by the nature of the problems we are studying in this chapter, the solutions have to be constructed under consideration of the flux direction. That is, in the spirit of characteristic curves we know that information propagates along curves in space time. And if this curve advances, for example, from left to right in space, then we need to use the forward difference quotient, because the backward difference quotient will not contain the information of the flow direction.

### 11.3.2 Spatially Semi-Discrete Conservation Form

The method we will use in this section of the course is the Finite Volume Method. We use a uniform 1D mesh with  $N$  cells and search a piecewise constant solution. The coefficients are the values of  $u$  at the cell centers:

$$\mu_i \approx u(x_i, t)$$

we build a mesh by taking intervals around the spatial points in which we approximate the solution  $u$ . Now we derive an ODE for  $\mu$  by starting from the problem definition and integrating over cell  $j$ :

$$\frac{\partial u}{\partial t} = -\frac{\partial}{\partial x} f(u)$$

$$\underbrace{\frac{\partial}{\partial t} \int_{x_{j-1/2}}^{x_{j+1/2}} u(x, t) dx}_{\approx h \cdot \mu_j} = -\left(f(u(x_{j+1/2}, t)) - f(u(x_{j-1/2}, t))\right) \quad (54)$$

Because  $u$  is piecewise constant,  $u(x_{j+\frac{1}{2}}, t)$  is undefined.

An approximation for  $f_{j+\frac{1}{2}} := f(u(x_{j+\frac{1}{2}}, t))$  is given by the **numerical flux function**  $F$ :

#### General Numerical Flux Function

$$f_{j+\frac{1}{2}} \approx F(\mu_{j-m_l+1}, \dots, \mu_{j+m_r})$$

$F$  is called a  $(m_l + m_r)$ -point numerical flux function.

A **2-point flux**, which is all we need in this section, is defined as  $f_{j+\frac{1}{2}} \approx F(\mu_j, \mu_{j+1})$ .

With this, (54) becomes (omitting the time argument of  $\mu_j$ )

$$\frac{d\mu_j}{dt}(t) = -\frac{1}{h}(F(\mu_j, \mu_{j+1}) - F(\mu_{j-1}, \mu_j)) \quad (55)$$

We can simply plug the RHS of (55) into a Runge–Kutta method, so the only thing left to do is to find a suitable flux function  $F$ , which should satisfy the following condition:

**Definition 11.3.3.5: Consistent numerical flux function**

A numerical function  $F : \mathbb{R}^{m_l+m_r} \rightarrow \mathbb{R}$  is *consistent* with the flux  $f : \mathbb{R} \rightarrow \mathbb{R}$  if

$$F(u, u, \dots, u) = f(u), \quad \forall u \in \mathbb{R}$$

Theory tells us that a consistent numerical flux will produce a solution with the correct shock speed.

### 11.3.4 Numerical Flux Functions

This section now treats how to find suitable flux functions. There are several options: the simplest one just averages the inputs of  $F(u, w)$ , but suffers from the same problem as finite difference methods.

One remedy for this is the *Lax–Friedrichs / Rusanov* Flux which is useful but flattens the edges of jumps (which is due to its construction with additional diffusion).

**Equation 11.3.4.16 (Lax–Friedrichs / Rusanov Flux)**

$$F_{LF}(v, w) = \frac{1}{2}(f(v) + f(w)) - \frac{1}{2}(w - v) \max_{\min\{v, w\} \leq u \leq \max\{v, w\}} f'(u)$$

The third term is the artificial diffusion term.

As pointed out before, the direction in which the information flows is crucial, so an important idea to choose the right flux is to respect that. Moreover the flux has to reproduce physical solution in the sense explained above when studying two possible solutions for the Riemann problem.

The final flux we look at, which solves these problems, is the Godunov Flux.

**Definition 11.3.4.33: Godunov Flux**

$$F_{GD}(v, w) = \begin{cases} \min_{v \leq u \leq w} f(u) & \text{if } v < w \\ \max_{w \leq u \leq v} f(u) & \text{if } v \geq w \end{cases}$$

### 11.3.5 Monotone Schemes

First, we define what it means for a flux to be monotone:

#### Definition 11.3.5.5: Monotonicity of flux functions

A 2-point numerical flux function  $F = F(v, w)$  is *monotone* if

- $F$  is increasing in its first argument, i.e.  $F(v, w) \leq F(v + \Delta v, w) \quad \forall w \in \mathbb{R}$
- $F$  is decreasing in its second argument, i.e.  $F(v, w) \leq F(v, w + \Delta w) \quad \forall v \in \mathbb{R}$

In Section 11.2.7, we mentioned that the number of extrema of an analytical solution does not increase over time. Here, we show that the two main fluxes (LF and Godunov) we derived in the previous chapter are monotone and therefore both have this property. This is established by the following two lemmas:

#### Lemma 11.3.5.8: Monotonicity of Lax–Friedrichs and Godunov flux

For any continuously differentiable flux function  $f$  the associated Lax–Friedrichs flux and Godunov flux are monotone.

#### Lemma 11.3.5.13: Non-oscillatory monotone semi-discrete evolutions

If  $\mu = \mu(t)$  solves the two-point flux equation (55) with a monotone numerical flux and  $\mu(0)$  has finitely many local extrema, then the number of local extrema of  $\mu(t)$  cannot be larger than that of  $\mu(0)$ .

## 11.4 Timestepping for Finite-Volume Methods

As explained before, once we have chosen the numerical Flux, we just need to apply Runge–Kutta integration. This subsection studies some conditions that have to be considered when applying Runge–Kutta – in particular, constraints on choosing the timestep size  $\tau$ .

#### Definition 11.4.2.5: Numerical domain of dependence

Consider the explicit fully discrete evolution  $\mu^{(k+1)} := \mathcal{H}(\mu^{(k)})$  on a uniform spatio-temporal mesh  $(x_j = hj, j \in \mathbb{Z}, t_k = k\tau, k \in \mathbb{N}_0)$  with

$$\exists m \in \mathbb{N}_0 : (\mathcal{H}(\mu))_j = \mathcal{H}(\mu_{j-m}, \dots, \mu_{j+m}), j \in \mathbb{Z}. \quad (56)$$

Then the **numerical domain of dependence** is given by

$$D_h^-(x_j, t_k) := \{(x_n, t_l) \in \mathbb{R} \times [0, t_k] : j - m(k - l) \leq n \leq j + m(k - l)\}$$

$\mathcal{H}$  is an operator which gives  $\mu$  at the next timestep. It incorporates both the numerical flux and the Runge–Kutta method. If we have a first-order time integrator and the numerical flux function is  $F = F(\mu_{j-m}, \dots, \mu_{j+m})$ , then  $\mathcal{H} = \mathcal{H}(\mu_{j-m}, \dots, \mu_{j+m})$ .

**Example:** Let's say we are solving (55) and call the RHS  $R(\mu_j)$

$$\frac{d\mu_j}{dt}(t) = -\frac{1}{h}(F(\mu_j, \mu_{j+1}) - F(\mu_{j-1}, \mu_j)) = R(\mu_j)$$

and we apply explicit Euler:

$$\begin{aligned} (\mathcal{H}(\mu))_j &= \mu_j + \tau R(\mu_j) \\ &= \mu_j - \frac{\tau}{h}(F(\mu_j, \mu_{j+1}) - F(\mu_{j-1}, \mu_j)) \end{aligned}$$

If we compare this to (56), we see that  $m = 1$  in this case.

The following kind of condition appears over and over in numerical integration and gives an upper bound for the timestep size  $\tau$ : If this upper bound is respected, the numerical solution is stable.

**Definition 11.4.2.11: Courant–Friedrichs–Lewy (CFL) condition**

An explicit local fully discrete evolution  $\mu^{(t+1)} := \mathcal{H}(\mu)$  on uniform spatio-temporal mesh  $(x_j = hj, j \in \mathbb{Z}, t_k = k\tau, k \in \mathbb{N})$  satisfies the CFL condition, if the convex hull of its numerical domain of dependence contains the maximal analytical domain of dependence

$$D^-(x_j, t_k) \subset \text{convex}(D_h^-(x_j, t_k)) \quad \forall j, k.$$

Applied to our problem, this means

**Equation 11.4.2.12 (CFL condition for explicit fully discrete evolution)**

$$\frac{\tau}{h} \leq \frac{m}{\max\{|\dot{s}_{\min}|, |\dot{s}_{\max}|\}}.$$

**11.4.4 Convergence of Fully Discrete FV Method**

The **consistency error** is defined as follows:

$$\epsilon := \max_j \left\{ F(u(x_j, t), u(x_{j+1}, t)) - f\left(u\left(x_{j+\frac{1}{2}}, t\right)\right) \right\},$$

if we assume  $u$  to be an exact solution of If  $\epsilon = \mathcal{O}(h^q)$ , the flux is called *q-th order consistent*.  $q$  is then the order of convergence of the FV method.

We get at most order one convergence in the maximum and the  $L^1$  norm. This can be seen by the following fact:

**Order barrier for monotone numerical fluxes**

Monotone numerical fluxes (Definition 11.3.5.5) are at most first order consistent.

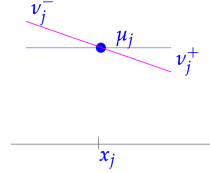
## 11.5 Higher-Order FV Schemes

For the FV Method to achieve a convergence of order  $> 1$ , we need fluxes which are consistent of order  $> 1$ . Instead of coming up with completely new fluxes, we use the ones from Section 11.3.4 and only modify the arguments given to them.

We choose the approach to go from piecewise-constant to **discontinuous piecewise-linear** solutions. Instead of approximating  $f\left(u\left(x_{j+\frac{1}{2}}\right)\right)$  by  $F(\mu_j, \mu_{j+1})$ , we use  $F(\nu_j^+, \nu_{j+1}^-)$ .

$\nu_j^+, \nu_j^-$  are the values of the linear approximation of  $u$  at the boundaries of cell  $j$  (see image). The discrete evolution equation becomes

$$\frac{d\mu_j}{dt}(t) = -\frac{1}{h} \left( F(\nu_j^+(t), \nu_{j+1}^-(t)) - F(\nu_{j-1}^+(t), \nu_j^-(t)) \right) \quad (57)$$



where  $\nu_j^+, \nu_j^-$  are the values of the linear approximation of  $u$  at the cell boundaries. The linear approximation on cell  $j$  is defined by the slope  $\sigma_j$ . So this method boils down to finding formulas for the slopes:

$$\sigma_j = \text{slope}(\mu_{j-1}, \mu_j, \mu_{j+1})$$

### 11.5.2 Slope limiting

Finding a linear approximation is formally expressed as a **reconstruction operator**:

$$R_M \vec{\mu} = \text{discontinuous pw. linear approximation of } u$$

Note that this operator as a mapping between vector spaces is not linear. However, we call it a **linear reconstruction** because the result is a (piecewise) linear function. We get the following formulas which we can plug into (57):

$$v_j^+ = R_M \vec{\mu} \left( x_{j+\frac{1}{2}}^- \right) \quad \text{and} \quad v_j^- = R_M \vec{\mu} \left( x_{j-\frac{1}{2}}^+ \right)$$

where  $x^+$  denotes the limit from the right and  $x^-$  from the left. Note that we always assume  $R_M \vec{\mu}(x_j) = \mu_j$ , i.e., the linear reconstruction has the same values as the piecewise constant function at the cell centers.

Given slopes  $\sigma_j$ , we can now define the linear reconstruction operator:

$$R_M \vec{\mu}(x) = \mu_j + \sigma_j(x - x_j) \quad \text{for } x \in ]x_{j-\frac{1}{2}}, x_{j+\frac{1}{2}}[$$

We introduce an important property of the reconstruction operator:

#### Definition 11.5.2.1: Monotonicity-preserving linear reconstruction

A linear reconstruction operator  $R_M$  is called *monotonicity-preserving* if

$$\mu_j \leq \mu_{j+1} \implies R_M \vec{\mu} \text{ is non-decreasing in } ]x_j, x_{j+1}[$$

$$\mu_j \geq \mu_{j+1} \implies R_M \vec{\mu} \text{ is non-increasing in } ]x_j, x_{j+1}[$$

This ensures that no new extrema are introduced by the numerical scheme.



Finally, we present a reconstruction operator which is monotonicity-preserving:

**Definition 11.5.2.9: Minmod reconstruction**

$$\sigma_j(\mu_{j-1}, \mu_j, \mu_{j+1}) = \text{minmod} \left( \frac{\mu_{j+1} - \mu_j}{x_{j+1} - x_j}, \frac{\mu_j - \mu_{j-1}}{x_j - x_{j-1}} \right)$$

where minmod is defined as

$$\text{minmod}(a, b) = \begin{cases} a, & \text{if } ab > 0, |a| < |b| \\ b, & \text{if } ab < 0, |b| < |a| \\ 0, & \text{otherwise} \end{cases}$$

## 11.6 Systems of Conservation Laws

**Definition 11.6.1.1: Linear system of conservation laws**

Let  $A(x) : \mathbb{R} \rightarrow \mathbb{R}^m \times \mathbb{R}^m$  be a matrix-valued function and  $D \subset \mathbb{R} \times \mathbb{R}^+$ . Then

$$\frac{\partial \mathbf{u}}{\partial t}(x, t) + \frac{\partial}{\partial x} (A(x) \mathbf{u}(x, t)) = 0 \quad \text{on } D \quad (58)$$

is a linear system of conservation laws for  $\mathbf{u} : D \rightarrow \mathbb{R}^m$ .

By diagonalizing  $A$ , we can get a general solution formula. Let  $\mathbf{r}_i$  be an eigenvector of  $A$  with eigenvalue  $\lambda_i$ :  $A\mathbf{r}_i = \lambda_i\mathbf{r}_i$ . We collect the eigenvectors in a matrix  $R$ . Note that both  $\lambda_i$  and  $\mathbf{r}_i$  depend on  $x$ .

**Solution of (58)**

$$\mathbf{u}(x, t) = \sum_{i=1}^m (R^{-1} \mathbf{u}_0)_i (x - \lambda_i t) \mathbf{r}_i$$

We have a superposition of states  $\mathbf{r}_i$  propagating with speed  $\lambda_i$  scaled by  $w_i := (R^{-1} \mathbf{u}_0)_i$

We want to solve the Riemann problem

$$\mathbf{u}(x, 0) = \begin{cases} \mathbf{u}_l & \text{if } x < 0 \\ \mathbf{u}_r & \text{if } x > 0 \end{cases} \quad (59)$$

for which we can find an explicit solution. Let's illustrate it for the case  $m = 4$ :

$$\begin{aligned} \mathbf{u}(x, t) &= \mathbf{u}_l \text{ if } x < \lambda_1 t \\ &\sum_{i=1}^k w_i^r \mathbf{r}_i + \sum_{i=k+1}^m w_i^l \mathbf{r}_i \text{ if } \lambda_k t < x < \lambda_{k+1} t \end{aligned}$$

For the jumps we get a formula similar to the Rankine–Hugoniot condition:

$$\mathbf{A}(x)(\mathbf{u}_k - \mathbf{u}_{k-1}) = \lambda_k(\mathbf{u}_k - \mathbf{u}_{k-1}),$$

where  $\lambda_k$  takes the role of  $\dot{s}$  in the scalar case.

**Definition 11.6.2.2: Non-linear system of conservation laws**

Given  $U \subset \mathbb{R}^m$ ,  $D \subset \mathbb{R} \times \mathbb{R}^+$  and a continuously differentiable flux function  $\mathbf{F} : U \rightarrow \mathbb{R}^m$ ,

$$\frac{\partial \mathbf{u}}{\partial t} + \frac{\partial \mathbf{F}(\mathbf{u})}{\partial x}(x, t) = 0 \quad \text{on } D \quad (60)$$

is a system of conservation laws for  $\mathbf{u} : D \rightarrow U$ .

This system is called *constant-coefficient* because  $\mathbf{F}$  does not depend on time and *translation-invariant* because it does not take  $x$  as an argument.

We apply diagonalization again, this time to the Jacobian  $D\mathbf{F}(\mathbf{u})$ . Note that here, eigenvectors and eigenvalues do not depend on  $x$  anymore, but instead on  $\mathbf{u}$ :

$$D\mathbf{F}(\mathbf{u}) \mathbf{r}_i(\mathbf{u}) = \lambda_i(\mathbf{u}) \mathbf{r}_i$$

Here, we also get a Rankine–Hugoniot condition:

**Theorem 11.6.2.19: Rankine--Hugoniot condition for systems**

If a curve  $\Gamma = (\gamma(t), t)$  separates two domains

$$\Omega_l = \{(x, t) \in D : x < \gamma(t)\} \quad \text{and} \quad \Omega_r = \{(x, t) \in D : x > \gamma(t)\}$$

and  $\mathbf{u}_l, \mathbf{u}_r = \mathbf{u}|_{\Omega_l}, \mathbf{u}|_{\Omega_r} \in C^1$  are strong solutions of (60), then  $\mathbf{u}$  is a weak solution of (60) if and only if

$$\frac{d\gamma}{dt}(t) \cdot (\mathbf{u}_l(\gamma(t), t) - \mathbf{u}_r(\gamma(t), t)) = \mathbf{F}(\mathbf{u}_l(\gamma(t), t)) - \mathbf{F}(\mathbf{u}_r(\gamma(t), t))$$

For the Riemann problem (59), this simplifies to

$$\dot{s}(\mathbf{u}_l - \mathbf{u}_r) = \mathbf{F}(\mathbf{u}_l) - \mathbf{F}(\mathbf{u}_r), \quad \dot{s} = \frac{d\gamma}{dt}$$

There is also an entropy condition for systems:

**Definition 11.6.2.29: Lax entropy condition for systems**

1.  $\exists k \in \{1, \dots, m\} : \lambda_k(\mathbf{u}_l) > \dot{s} > \lambda_k(\mathbf{u}_r)$
2.  $\forall j > k : \lambda_j(\mathbf{u}_l), \lambda_j(\mathbf{u}_r) < \dot{s}$
3.  $\forall j < k : \lambda_j(\mathbf{u}_l), \lambda_j(\mathbf{u}_r) > \dot{s}$