

Scientific Computing WS 2019/2020

Lecture 16

Jürgen Fuhrmann

[juergen.fuhrmann@wias-berlin.de](mailto:juergen.fuhrmann@wias-berlin.de)

**Recap: Weak formulation**

## Weak formulation of homogeneous Dirichlet problem

Find  $u \in H_0^1(\Omega)$  such that

$$\int_{\Omega} \lambda \vec{\nabla} u \cdot \vec{\nabla} v \, d\vec{x} = \int_{\Omega} f v \, d\vec{x} \quad \forall v \in H_0^1(\Omega)$$

- Then,

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

is a self-adjoint bilinear form defined on the Hilbert space  $H_0^1(\Omega)$ .

- It is bounded due to Cauchy-Schwarz:

$$|a(u, v)| = \lambda \cdot \left| \int_{\Omega} \vec{\nabla} u \cdot \vec{\nabla} v \, d\vec{x} \right| \leq \lambda \|u\|_{H_0^1(\Omega)} \cdot \|v\|_{H_0^1(\Omega)}$$

- $f(v) = \int_{\Omega} f v \, d\vec{x}$  is a linear functional on  $H_0^1(\Omega)$ . For Hilbert spaces  $V$  the dual space  $V'$  (the space of linear functionals) can be identified with the space itself.

## The Lax-Milgram lemma

**Theorem:** Let  $V$  be a Hilbert space. Let  $a : V \times V \rightarrow \mathbb{R}$  be a self-adjoint bilinear form, and  $f$  a linear functional on  $V$ . Assume  $a$  is coercive, i.e.

$$\exists \alpha > 0 : \forall u \in V, a(u, u) \geq \alpha \|u\|_V^2.$$

Then the problem: find  $u \in V$  such that

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

admits one and only one solution with an a priori estimate

$$\|u\|_V \leq \frac{1}{\alpha} \|f\|_{V'}$$



## Coercivity of weak formulation

**Theorem:** Assume  $\lambda > 0$ . Then the weak formulation of the heat conduction problem: search  $u \in H_0^1(\Omega)$  such that

$$\int_{\Omega} \lambda \vec{\nabla} u \cdot \vec{\nabla} v \, d\vec{x} = \int_{\Omega} f v \, d\vec{x} \quad \forall v \in H_0^1(\Omega)$$

has an unique solution.

**Proof:**  $a(u, v)$  is cocercive:

$$a(u, u) = \int_{\Omega} \lambda \vec{\nabla} u \cdot \vec{\nabla} u \, d\vec{x} = \lambda \|u\|_{H_0^1(\Omega)}^2$$



## Weak formulation of inhomogeneous Dirichlet problem

$$\begin{aligned} -\nabla \cdot \lambda \vec{\nabla} u &= f \text{ in } \Omega \\ u &= g \text{ on } \partial\Omega \end{aligned}$$

If  $g$  is smooth enough, there exists a *lifting*  $u_g \in H^1(\Omega)$  such that  $u_g|_{\partial\Omega} = g$ . Then, we can re-formulate:

$$\begin{aligned} -\nabla \cdot \lambda \vec{\nabla} (u - u_g) &= f + \nabla \cdot \lambda \vec{\nabla} u_g \text{ in } \Omega \\ u - u_g &= 0 \text{ on } \partial\Omega \end{aligned}$$

Find  $u \in H^1(\Omega)$  such that

$$\begin{aligned} u &= u_g + \phi \\ \int_{\Omega} \lambda \vec{\nabla} \phi \cdot \vec{\nabla} v \, d\vec{x} &= \int_{\Omega} f v \, d\vec{x} + \int_{\Omega} \lambda \vec{\nabla} u_g \cdot \vec{\nabla} v \, d\vec{x} \quad \forall v \in H_0^1(\Omega) \end{aligned}$$

Here, necessarily,  $\phi \in H_0^1(\Omega)$  and we can apply the theory for the homogeneous Dirichlet problem.

## Weak formulation of Robin problem II

Let

$$a^R(u, v) := \int_{\Omega} \lambda \vec{\nabla} u \cdot \vec{\nabla} v \, d\vec{x} + \int_{\partial\Omega} \alpha uv \, ds$$

$$f^R(v) := \int_{\Omega} fv \, d\vec{x} + \int_{\partial\Omega} gv \, ds$$

Find  $u \in H^1(\Omega)$  such that

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

If  $\lambda > 0$  and  $\alpha > 0$  then  $a^R(u, v)$  is coercive, and by Lax-Milgram we establish the existence of a weak solution

# The Galerkin method I

- Weak formulations “live” in Hilbert spaces which essentially are infinite dimensional
- For computer representations we need finite dimensional approximations
- The Galerkin method and its modifications provide a general scheme for the derivation of finite dimensional approximations
- Finite dimensional subspaces of Hilbert spaces are the spans of a set of basis functions, and are Hilbert spaces as well  $\Rightarrow$  e.g. the Lax-Milgram lemma is valid there as well



## The Galerkin method II

- Let  $V$  be a Hilbert space. Let  $a : V \times V \rightarrow \mathbb{R}$  be a self-adjoint bilinear form, and  $f$  a linear functional on  $V$ . Assume  $a$  is coercive with coercivity constant  $\alpha$ , and continuity constant  $\gamma$ .
- Continuous problem: search  $u \in V$  such that

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

- Let  $V_h \subset V$  be a finite dimensional subspace of  $V$
- “Discrete” problem  $\equiv$  Galerkin approximation:  
Search  $u_h \in V_h$  such that

$$a(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h$$

By Lax-Milgram, this problem has a unique solution as well.

## Céa's lemma

- What is the connection between  $u$  and  $u_h$  ?
- Let  $v_h \in V_h$  be arbitrary. Then

$$\begin{aligned}\alpha \|u - u_h\|^2 &\leq a(u - u_h, u - u_h) \quad (\text{Coercivity}) \\ &= a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h) \\ &= a(u - u_h, u - v_h) \quad (\text{Galerkin Orthogonality}) \\ &\leq \gamma \|u - u_h\| \cdot \|u - v_h\| \quad (\text{Boundedness})\end{aligned}$$

- As a result

$$\|u - u_h\| \leq \frac{\gamma}{\alpha} \inf_{v_h \in V_h} \|u - v_h\|$$

- Up to a constant, the error of the Galerkin approximation is the error of the best approximation of the solution in the subspace  $V_h$ .

# From the Galerkin method to the matrix equation

- Let  $\phi_1 \dots \phi_n$  be a set of basis functions of  $V_h$ .
- Then, we have the representation  $u_h = \sum_{j=1}^n u_j \phi_j$
- In order to search  $u_h \in V_h$  such that

$$a(u_h, v_h) = f(v_h) \quad \forall v_h \in V_h$$

it is actually sufficient to require

$$\begin{aligned} a(u_h, \phi_i) &= f(\phi_i) \quad (i = 1 \dots n) \\ a\left(\sum_{j=1}^n u_j \phi_j, \phi_i\right) &= f(\phi_i) \quad (i = 1 \dots n) \\ \sum_{j=1}^n a(\phi_j, \phi_i) u_j &= f(\phi_i) \quad (i = 1 \dots n) \end{aligned}$$

$$AU = F$$

with  $A = (a_{ij})$ ,  $a_{ij} = a(\phi_i, \phi_j)$ ,  $F = (f_i)$ ,  $f_i = F(\phi_i)$ ,  $U = (u_i)$ .

- Matrix dimension is  $n \times n$ . Matrix sparsity ?

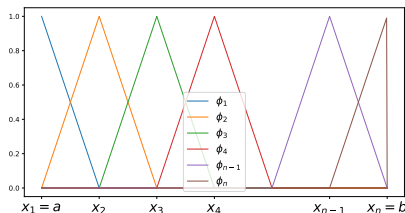
## Obtaining a finite dimensional subspace

- Let  $\Omega = (a, b) \subset \mathbb{R}^1$
- Let  $a(u, v) = \int_a^b \delta \vec{\nabla} u \vec{\nabla} v d\vec{x} + \alpha u(a)v(a) + \alpha u(b)v(b)$
- Calculus 101 provides a finite dimensional subspace: the space of sin/cos functions up to a certain frequency  $\Rightarrow$  *spectral method*
- Ansatz functions have global support  $\Rightarrow$  full  $n \times n$  matrix
- OTOH: rather fast convergence for smooth data
- Generalization to higher dimensions possible
- Big problem in irregular domains: we need the eigenfunction basis of some operator...
- Spectral methods are successful in cases where one has regular geometry structures and smooth/constant coefficients – e.g. “Spectral Einstein Code”

# The finite element idea I

- Choose basis functions with local support.  $\Rightarrow$  only integrals of basis function pairs with overlapping support contribute to matrix.
- Linear finite elements in  $\Omega = (a, b) \subset \mathbb{R}^1$ :
- Partition  $a = x_1 \leq x_2 \leq \dots \leq x_n = b$
- Basis functions (for  $i = 1 \dots n$ )

$$\phi_i(x) = \begin{cases} \frac{x-x_{i-1}}{x_i-x_{i-1}}, & i > 1, x \in (x_{i-1}, x_i) \\ \frac{x_{i+1}-x}{x_{i+1}-x_i}, & i < n, x \in (x_i, x_{i+1}) \\ 0, & \text{else} \end{cases}$$



# FE matrix elements for 1D heat equation I

- Any function  $u_h \in V_h = \text{span}\{\phi_1 \dots \phi_n\}$  is piecewise linear, and the coefficients in the representation  $u_h = \sum_{i=1}^n u_i \phi_i$  are the values  $u_h(x_i)$ .
- Fortunately, we are working with a weak formulation, and weak derivatives are well defined (and coincide with the classical derivatives where the basis functions are smooth)
- Let  $\phi_i, \phi_j$  be two basis functions, regard

$$s_{ij} = \int_a^b \vec{\nabla} \phi_i \cdot \vec{\nabla} \phi_j dx$$

- We have  $\text{supp } \phi_i \cap \text{supp } \phi_j = \emptyset$  unless  $i = j$ ,  $i + 1 = j$  or  $i - 1 = j$ .
- Therefore  $s_{ij} = 0$  unless  $i = j$ ,  $i + 1 = j$  or  $i - 1 = j$ .

## FE matrix elements for 1D heat equation II

- Let  $j = i + 1$ . Then  $\text{supp } \phi_i \cap \text{supp } \phi_j = (x_i, x_{i+1})$ ,  $\phi'_i = -\frac{1}{h}$ ,  $\phi'_j = \frac{1}{h}$  where  $h = x_{i+1} - x_i$

$$\int_a^b \vec{\nabla} \phi_i \vec{\nabla} \phi_j d\vec{x} = \int_{x_i}^{x_{i+1}} \phi'_i \phi'_j d\vec{x} = - \int_{x_i}^{x_{i+1}} \frac{1}{h^2} d\vec{x} = -\frac{1}{h}$$

- Similarly, for  $j = i - 1$ :  $\int_a^b \vec{\nabla} \phi_i \vec{\nabla} \phi_j d\vec{x} = -\frac{1}{h}$
- For  $1 < i < N$ :

$$\int_a^b \vec{\nabla} \phi_i \vec{\nabla} \phi_i d\vec{x} = \int_{x_{i-1}}^{x_{i+1}} (\phi'_i)^2 d\vec{x} = \int_{x_{i-1}}^{x_{i+1}} \frac{1}{h^2} d\vec{x} = \frac{2}{h}$$

- For  $i = 1$  or  $i = N$ :  $\int_a^b \vec{\nabla} \phi_i \vec{\nabla} \phi_i d\vec{x} = \frac{1}{h}$
- For the right hand side, calculate vector elements  $f_i = \int_a^b f(x) \phi_i dx$  using a quadrature rule.

## FE matrix elements for 1D heat equation III

Adding the boundary integrals yields

$$A = \begin{pmatrix} \alpha + \frac{1}{h} & -\frac{1}{h} & & & & & & & \\ -\frac{1}{h} & \frac{2}{h} & -\frac{1}{h} & & & & & & \\ & -\frac{1}{h} & \frac{2}{h} & -\frac{1}{h} & & & & & \\ & & \ddots & \ddots & \ddots & \ddots & & & \\ & & & -\frac{1}{h} & \frac{2}{h} & -\frac{1}{h} & & & \\ & & & & -\frac{1}{h} & \frac{2}{h} & -\frac{1}{h} & & \\ & & & & & -\frac{1}{h} & \frac{2}{h} & -\frac{1}{h} & \\ & & & & & & -\frac{1}{h} & \frac{1}{h} + \alpha & \end{pmatrix}$$

... the same matrix as for the finite difference and finite volume methods



# Simplices

- Let  $\{\vec{a}_1 \dots \vec{a}_{d+1}\} \subset \mathbb{R}^d$  such that the  $d$  vectors  $\vec{a}_2 - \vec{a}_1 \dots \vec{a}_{d+1} - \vec{a}_1$  are linearly independent. Then the convex hull  $K$  of  $\vec{a}_1 \dots \vec{a}_{d+1}$  is called *simplex*, and  $\vec{a}_1 \dots \vec{a}_{d+1}$  are called *vertices* of the simplex.
- *Unit simplex*:  $\vec{a}_1 = (0 \dots 0)$ ,  $\vec{a}_2 = (0, 1 \dots 0) \dots \vec{a}_{d+1} = (0 \dots 0, 1)$ .

$$K = \left\{ \vec{x} \in \mathbb{R}^d : x_i \geq 0 \ (i = 1 \dots d) \text{ and } \sum_{i=1}^d x_i \leq 1 \right\}$$

- A general simplex can be defined as an image of the unit simplex under some affine transformation
- $F_i$ : face of  $K$  opposite to  $\vec{a}_i$
- $\vec{n}_i$ : outward normal to  $F_i$

# Simplex characteristics

- Diameter of  $K$ :  $h_K = \max_{\vec{x}_1, \vec{x}_2 \in K} \|\vec{x}_1 - \vec{x}_2\|$   
 $\equiv$  length of longest edge if  $K$
- $\rho_K$  diameter of largest ball that can be inscribed into  $K$
- $\sigma_K = \frac{h_K}{\rho_K}$ : local shape regularity measure
  - $\sigma_K = 2\sqrt{3}$  for equilateral triangle
  - $\sigma_K \rightarrow \infty$  if largest angle approaches  $\pi$ .

# Barycentric coordinates

**Definition:** Let  $K \subset \mathbb{R}^d$  be a  $d$ -simplex given by the points  $\vec{a}_1 \dots \vec{a}_{d+1}$ . Let  $\Lambda(x) = (\lambda_1(\vec{x}) \dots \lambda_{d+1}(\vec{x}))$  be a vector such that for all  $\vec{x} \in \mathbb{R}^d$

$$\sum_{j=1}^{d+1} \vec{a}_j \lambda_j(\vec{x}) = \vec{x}, \quad \sum_{j=1}^{d+1} \lambda_j(\vec{x}) = 1$$

This vector is called the vector of *barycentric coordinates* of  $\vec{x}$  with respect to  $K$ .

## Barycentric coordinates II

**Lemma** The barycentric coordinates of a given point is well defined and unique. Moreover, for the simplex edges  $\vec{a}_i$ , one has

$$\lambda_j(\vec{a}_i) = \delta_{ij}$$

**Proof:** The definition of  $\Lambda$  given by a  $d + 1 \times d + 1$  system of equations with the matrix

$$M = \begin{pmatrix} a_{1,1} & a_{2,1} & \dots & a_{d+1,1} \\ a_{1,2} & a_{2,2} & \dots & a_{d+1,2} \\ \vdots & \vdots & & \vdots \\ a_{1,d} & a_{2,d} & \dots & a_{d+1,d} \\ 1 & 1 & \dots & 1 \end{pmatrix}$$

Subtracting the first column from the others gives

$$M' = \begin{pmatrix} a_{1,1} & a_{2,1} - a_{1,1} & \dots & a_{d+1,1} - a_{1,1} \\ a_{1,2} & a_{2,2} - a_{1,2} & \dots & a_{d+1,2} - a_{1,2} \\ \vdots & \vdots & & \vdots \\ a_{1,d} & a_{2,d} - a_{1,d} & \dots & a_{d+1,d} - a_{1,d} \\ 1 & 0 & \dots & 0 \end{pmatrix}$$

## Barycentric coordinates III

$\det M = \det M'$  is the determinant of the matrix whose columns are the edge vectors of  $K$  which are linearly independent.

For the simplex edges one has

$$\sum_{j=1}^{d+1} \vec{a}_j \lambda_j(\vec{a}_i) = \vec{a}_i$$

which is fulfilled if  $\lambda_j(\vec{a}_i) = 1$  for  $i = j$  and  $\lambda_j(\vec{a}_i) = 0$  for  $i \neq j$ . And we have uniqueness.  $\square$

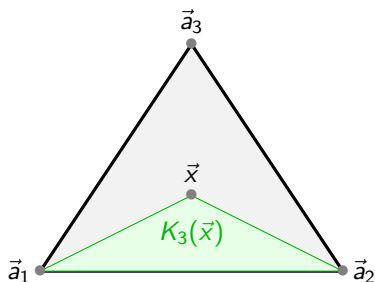
At the same time, the measure (area) is calculated as  $|K| = \frac{1}{d!} |\det M'|$ .

## Barycentric coordinates IV

- Let  $K_j(\vec{x})$  be the subsimplex of  $K$  made of  $\vec{x}$  and  $\vec{a}_1 \dots \vec{a}_{d+1}$  with  $\vec{a}_j$  omitted.
- Its measure  $|K_j(\vec{x})|$  is established from its determinant and a linear function of the coordinates for  $\vec{x}$ .
- One has  $\frac{|K_j(\vec{a}_i)|}{|K|} = \delta_{ij}$  and therefore,

$$\lambda_j(\vec{x}) = \frac{|K_j(\vec{x})|}{|K|}$$

is the ratio of the measures of  $K_j(\vec{x})$  and  $K$ .



# Conformal triangulations

- Let  $\mathcal{T}_h$  be a subdivision of the polygonal domain  $\Omega \subset \mathbb{R}^d$  into non-intersecting compact simplices  $K_m$ ,  $m = 1 \dots n_e$ :

$$\bar{\Omega} = \bigcup_{m=1}^{n_e} K_m$$

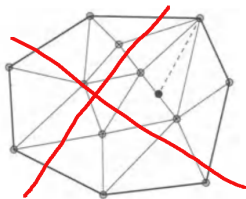
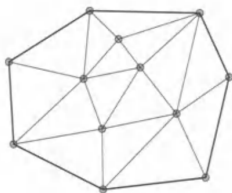
- Each simplex can be seen as the image of a affine transformation of a reference (e.g. unit) simplex  $\hat{K}$ :

$$K_m = T_m(\hat{K})$$

- We assume that it is conformal, i.e. if  $K_m, K_n$  have a  $d - 1$  dimensional intersection  $F = K_m \cap K_n$ , then there is a face  $\hat{F}$  of  $\hat{K}$  and renumberings of the vertices of  $K_n, K_m$  such that  $F = T_m(\hat{F}) = T_n(\hat{F})$  and  $T_m|_{\hat{F}} = T_n|_{\hat{F}}$

# Conformal triangulations II

- $d = 1$  : Each intersection  $F = K_m \cap K_n$  is either empty or a common vertex
- $d = 2$  : Each intersection  $F = K_m \cap K_n$  is either empty or a common vertex or a common edge



- $d = 3$  : Each intersection  $F = K_m \cap K_n$  is either empty or a common vertex or a common edge or a common face
- Delaunay triangulations are conformal



# Shape regularity

- Now we discuss a family of meshes  $\mathcal{T}_h$  for  $h \rightarrow 0$ .
- For given  $\mathcal{T}_h$ , assume that  $h = \max_{K \in \mathcal{T}_h} h_K$
- A family of meshes is called *shape regular* if

$$\forall h, \forall K \in \mathcal{T}_h, \sigma_K = \frac{h_K}{\rho_K} \leq \sigma_0$$

- In 1D,  $\sigma_K = 1$
- In 2D,  $\sigma_K \leq \frac{2}{\sin \theta_K}$  where  $\theta_K$  is the smallest angle

# Polynomial space $\mathbb{P}_k$

- Space of polynomials in  $x_1 \dots x_d$  of total degree  $\leq k$  with real coefficients  $\alpha_{i_1 \dots i_d}$ :

$$\mathbb{P}_k = \left\{ p(x) = \sum_{\substack{0 \leq i_1 \dots i_d \leq k \\ i_1 + \dots + i_d \leq k}} \alpha_{i_1 \dots i_d} x_1^{i_1} \dots x_d^{i_d} \right\}$$

- Dimension:

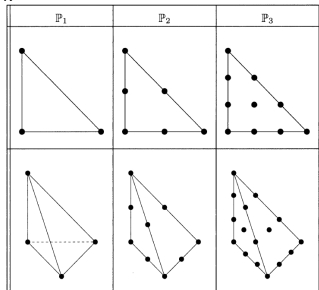
$$\dim \mathbb{P}_k = \binom{d+k}{k} = \begin{cases} k+1, & d=1 \\ \frac{1}{2}(k+1)(k+2), & d=2 \\ \frac{1}{6}(k+1)(k+2)(k+3), & d=3 \end{cases}$$

$$\dim \mathbb{P}_1 = d+1$$

$$\dim \mathbb{P}_2 = \begin{cases} 3, & d=1 \\ 6, & d=2 \\ 10, & d=3 \end{cases}$$

# $\mathbb{P}_k$ simplex finite elements

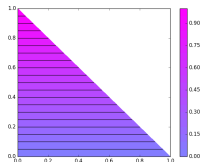
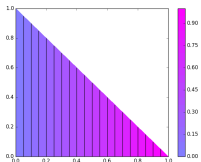
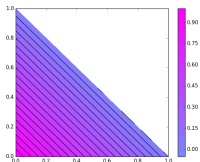
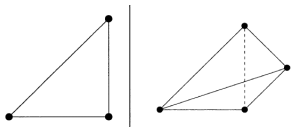
- $K$ : simplex spanned by  $\vec{a}_1 \dots \vec{a}_{d+1}$  in  $\mathbb{R}^d$
- For  $0 \leq i_1 \dots i_{d+1} \leq k$ ,  $i_1 + \dots + i_{d+1} = k$ , let the set of nodes  $\Sigma = \{\vec{\sigma}_1 \dots \vec{\sigma}_s\}$  be defined by the points  $\vec{a}_{i_1 \dots i_{d+1}; k}$  with barycentric coordinates  $(\frac{i_1}{k} \dots \frac{i_{d+1}}{k})$ .



- $s = \text{card } \Sigma = \dim \mathbb{P}_K \Rightarrow$  there exists a basis  $\theta_1 \dots \theta_s$  of  $\mathbb{P}_K$  such that  $\theta_i(\vec{\sigma}_j) = \delta_{ij}$

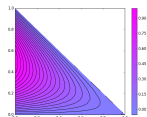
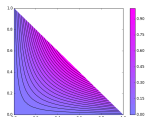
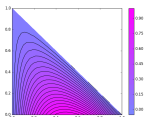
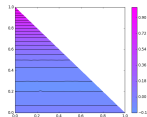
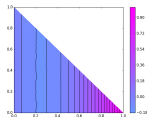
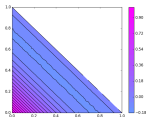
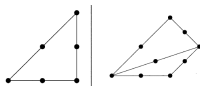
# $\mathbb{P}_1$ simplex finite elements

- $K$ : simplex spanned by  $a_1 \dots a_{d+1}$  in  $\mathbb{R}^d$
- $s = d + 1$
- Nodes  $\equiv$  vertices
- Basis functions  $\theta_1 \dots \theta_{d+1} \equiv$  barycentric coordinates  $\lambda_1 \dots \lambda_{d+1}$



# $\mathbb{P}_2$ simplex finite elements

- $K$ : simplex spanned by  $a_1 \dots a_{d+1}$  in  $\mathbb{R}^d$
- Nodes  $\equiv$  vertices + edge midpoints
- Basis functions:  
 $\lambda_i(2\lambda_i - 1), (0 \leq i \leq d); \quad 4\lambda_i\lambda_j, \quad (0 \leq i < j \leq d)$  ("edge bubbles")



- Finite element: triplet  $\{K, P, \Sigma\}$  where
  - $K \subset \mathbb{R}^d$ : compact, connected Lipschitz domain with non-empty interior
  - $P$ : finite dimensional space of functions  $p : K \rightarrow \mathbb{R}$
  - $\Sigma = \{\sigma_1 \dots \sigma_s\} \subset \mathcal{L}(P, \mathbb{R})$ : set of linear forms defined on  $P$  called *local degrees of freedom* such that the mapping

$$\begin{aligned}\Lambda_\Sigma : P &\rightarrow \mathbb{R}^s \\ p &\mapsto (\sigma_1(p) \dots \sigma_s(p))\end{aligned}$$

is bijective, i.e.  $\Sigma$  is a basis of  $\mathcal{L}(P, \mathbb{R})$ .

- Given a set of points  $\{\vec{\sigma}_1 \dots \vec{\sigma}_s\} \subset K$ , we use the symbols  $\sigma_1 \dots \sigma_s$  to denote the evaluation of a function at the respective points. Elements with this type of degree of freedom functionals are called *Lagrangian*.

# General finite elements

- Simplicial finite elements can be defined on triangulations of polygonal domains. During the course we will stick to this case.
- More general elements are possible: cuboids, but also prismatic elements etc.
- For vector PDEs, one can define finite elements for vector valued functions
- Different types of degrees of freedom (e.g. derivatives) are possible
- A curved domain  $\Omega$  may be approximated by a polygonal domain  $\Omega_h$  which is then triangulated. During the course, we will ignore this difference.
- Curved element geometries are possible. Isoparametric finite elements use the polynomial space to define a mapping of some polyhedral reference element to an element with curved boundary

# Global degrees of freedom

- Given a triangulation  $\mathcal{T}_h$
- Let  $\{\vec{a}_1 \dots \vec{a}_N\} = \bigcup_{K \in \mathcal{T}_h} \{\vec{\sigma}_{K,1} \dots \vec{\sigma}_{K,s}\}$  be the set of global degrees of freedom.
- Degree of freedom map

$$j: \mathcal{T}_h \times \{1 \dots s\} \rightarrow \{1 \dots N\}$$

$(K, m) \mapsto j(K, m)$  the global degree of freedom number



# Lagrange finite element space

- Given a triangulation  $\mathcal{T}_h$  of  $\Omega$ , define the spaces

$$P_h^k = \{v_h \in C^0(\Omega) : v_h|_K \in \mathbb{P}_K \forall K \in \mathcal{T}_h\} \subset H^1(\Omega)$$

$$P_{0,h}^k = \{v_h \in P_h^k : v_h|_{\partial\Omega} = 0\} \subset H_0^1(\Omega)$$

- Global shape functions  $\theta_1, \dots, \theta_N \in P_h^k$  defined by

$$\phi_i|_K(\vec{a}_{K,m}) = \begin{cases} \delta_{mn} & \text{if } \exists n \in \{1 \dots s\} : j(K, n) = i \\ 0 & \text{otherwise} \end{cases}$$

- $\{\phi_1, \dots, \phi_N\}$  is a basis of  $P_h$ , and  $\gamma_1 \dots \gamma_N$  is a basis of  $\mathcal{L}(P_h, \mathbb{R})$ :

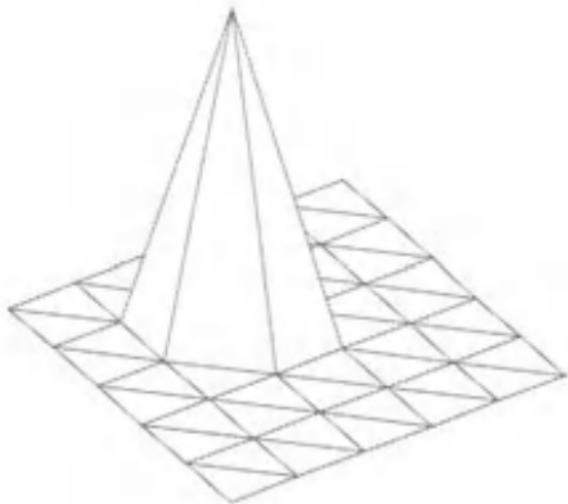
- $\{\phi_1, \dots, \phi_N\}$  are linearly independent: if  $\sum_{j=1}^N \alpha_j \phi_j = 0$  then evaluation at  $\vec{a}_1 \dots \vec{a}_N$  yields that  $\alpha_1 \dots \alpha_N = 0$ .
- Let  $v_h \in P_h$ . Let  $w_h = \sum_{j=1}^N v_h(\vec{a}_j) \phi_j$ . Then for all  $K \in \mathcal{T}_h$ ,  $v_h|_K$  and  $w_h|_K$  coincide in the local nodes  $\vec{a}_{K,1} \dots \vec{a}_{K,2}$ ,  $\Rightarrow v_h|_K = w_h|_K$ .

# Finite element approximation space

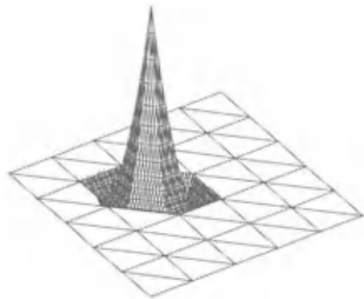
We have

d	k	$N = \dim P_h^k$
1	1	$N_v$
1	2	$N_v + N_{el}$
1	3	$N_v + 2N_{el}$
2	1	$N_v$
2	2	$N_v + N_{ed}$
2	3	$N_v + 2N_{ed} + N_{el}$
3	1	$N_v$
3	2	$N_v + N_{ed}$
3	3	$N_v + 2N_{ed} + N_f$

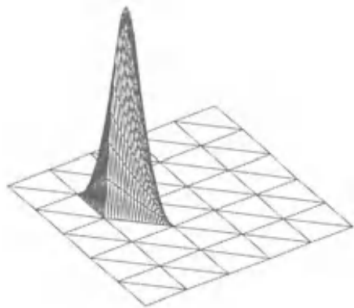
# $P^1$ global shape functions



## $P^2$ global shape functions



Node based



Edge based

# Local Lagrange interpolation operator

- Let  $\{K, P, \Sigma\}$  be a finite element with shape function bases  $\{\theta_1 \dots \theta_s\}$ . Let  $V(K) = \mathbb{C}^0(K)$  and  $P \subset V(K)$
- *local interpolation operator*

$$\mathcal{I}_K : V(K) \rightarrow P$$

$$v \mapsto \sum_{i=1}^s v(\vec{\sigma}_i) \theta_i$$

- $P$  is invariant under the action of  $\mathcal{I}_K$ , i.e.  $\forall p \in P, \mathcal{I}_K(p) = p$ :
  - Let  $p = \sum_{j=1}^s \alpha_j \theta_j$  Then,

$$\begin{aligned} \mathcal{I}_K(p) &= \sum_{i=1}^s p(\vec{\sigma}_i) \theta_i = \sum_{i=1}^s \sum_{j=1}^s \alpha_j \theta_j(\vec{\sigma}_i) \theta_i \\ &= \sum_{i=1}^s \sum_{j=1}^s \alpha_j \delta_{ij} \theta_i = \sum_{j=1}^s \alpha_j \theta_j \end{aligned}$$

# Global Lagrange interpolation operator

Let  $V_h = P_h^k$

$$\mathcal{I}_h : C^0(\bar{\Omega}_h) \rightarrow V_h$$
$$v(\vec{x}) \mapsto v_h(\vec{x}) = \sum_{i=1}^N v(\vec{a}_i) \phi_i(\vec{x})$$

# Reference finite element

- Let  $\{\widehat{P}, \widehat{K}, \widehat{\Sigma}\}$  be a fixed finite element
- Let  $T_K$  be some affine transformation and  $K = T_K(\widehat{K})$
- There is a linear bijective mapping  $\psi_K$  between functions on  $K$  and functions on  $\widehat{K}$ :

$$\begin{aligned}\psi_K : V(K) &\rightarrow V(\widehat{K}) \\ f &\mapsto f \circ T_K\end{aligned}$$

- Let
  - $K = T_K(\widehat{K})$
  - $P_K = \{\psi_K^{-1}(\widehat{p}); \widehat{p} \in \widehat{P}\},$
  - $\Sigma_K = \{\sigma_{K,i}, i = 1 \dots s : \sigma_{K,i}(p) = \widehat{\sigma}_i(\psi_K(p))\}$

Then  $\{K, P_K, \Sigma_K\}$  is a finite element.

- This construction allows to develop theory for a reference element and to lift it later to an arbitrary element.

# Commutativity of interpolation and reference mapping

- $\mathcal{I}_{\widehat{K}} \circ \psi_K = \psi_K \circ \mathcal{I}_K$ ,  
i.e. the following diagram is commutative:

$$\begin{array}{ccc} V(K) & \xrightarrow{\psi_K} & V(\widehat{K}) \\ \downarrow \mathcal{I}_K & & \downarrow \mathcal{I}_{\widehat{K}} \\ P_K & \xrightarrow{\psi_K} & P_{\widehat{K}} \end{array}$$

- $\equiv$  Interpolation and reference mapping are interchangeable



# Affine transformation estimates I

## Lemma T:

- $|\det J_K| = \frac{\text{meas}(K)}{\text{meas}(\widehat{K})}$
- $\|J_K\| \leq \frac{h_K}{\rho_{\widehat{K}}}, \|J_K^{-1}\| \leq \frac{h_{\widehat{K}}}{\rho_K}$
- $\Rightarrow \|J_K\| \cdot \|J_K^{-1}\| \leq c_{\widehat{K}} \sigma_K$

## Proof:

- $|\det J_K| = \frac{\text{meas}(K)}{\text{meas}(\widehat{K})}$ : basic property of affine mappings
- Further:

$$\|J_K\| = \sup_{\hat{x} \neq 0} \frac{\|J_K \hat{x}\|}{\|\hat{x}\|} = \frac{1}{\rho_{\widehat{K}}} \sup_{\|\hat{x}\| = \rho_{\widehat{K}}} \|J_K \hat{x}\|$$

Set  $\hat{x} = \hat{x}_1 - \hat{x}_2$  with  $\hat{x}_1, \hat{x}_2 \in \widehat{K}$ . Then  $J_K \hat{x} = T_K \hat{x}_1 - T_K \hat{x}_2$  and one can estimate  $\|J_K \hat{x}\| \leq h_K$ .

- For  $\|J_K^{-1}\|$  regard the inverse mapping  $\square$

# Estimate of derivatives under affine transformation

- For  $w \in H^s(K)$  recall the  $H^s$  seminorm  $|w|_{s,K}^2 = \sum_{|\beta|=s} \|\partial^\beta w\|_{L^2(K)}^2$
- We have

$$\|v\|_{L^2(K)}^2 = \int_K |v(x)|^2 dx = \int_{\hat{K}} |v(T_K(\hat{x}))|^2 \det J_K d\hat{x}$$

- For the derivative

$$|v|_{H^1(K)}^2 = \int_K \|\nabla v\|^2(x) dx = \int_{\hat{K}} \|J_K^{-T} \vec{\nabla} \hat{v}(\hat{x})\|^2 \det J_K d\hat{x}$$

## Estimate of derivatives under affine transformation II

**Lemma D:** Let  $w \in H^s(K)$  and  $\hat{w} = w \circ T_K$ . There exists a constant  $c$  such that

$$\begin{aligned} |\hat{w}|_{s, \hat{K}} &\leq c \|J_K\|^s |\det J_K|^{-\frac{1}{2}} |w|_{s, K} \\ |w|_{s, K} &\leq c \|J_K^{-1}\|^s |\det J_K|^{\frac{1}{2}} |\hat{w}|_{s, \hat{K}} \end{aligned}$$

**Proof:** Let  $|\alpha| = s$ . By affinity and chain rule one obtains

$$\|\partial^\alpha \hat{w}\|_{L^2(\hat{K})} \leq c \|J_K\|^s \sum_{|\beta|=s} \|\partial^\beta w \circ T_K\|_{L^2(K)}$$

Changing variables in the right hand side yields

$$\|\partial^\alpha \hat{w}\|_{L^2(\hat{K})} \leq c \|J_K\|^s |\det J_K|^{-\frac{1}{2}} |w|_{s, K}$$

Summation over  $\alpha$  yields the first inequality. Regarding the inverse mapping yields the second estimate.  $\square$

# Local interpolation error estimate I

**Theorem:** Let  $\{\widehat{K}, \widehat{P}, \widehat{\Sigma}\}$  be a finite element with associated normed vector space  $V(\widehat{K})$ . Assume that

$$\mathbb{P}_k \subset \widehat{P} \subset H^2(\widehat{K}) \subset V(\widehat{K})$$

Then there exists  $c > 0$  such that for all  $m = 0 \dots 2$ ,  $K \in \mathcal{T}_h$ ,  $v \in H^2(K)$ :

$$|v - \mathcal{I}_K^1 v|_{m,K} \leq ch_K^{2-m} \sigma_K^m |v|_{2,K}.$$

I.e. the the local interpolation error can be estimated through  $h_K$ ,  $\sigma_K$  and the norm of a higher derivative.

# Local interpolation error estimate II

## Draft of Proof

- Estimate on reference element  $\hat{K}$  using deeper results from functional analysis:

$$|\hat{w} - \mathcal{I}_{\hat{K}}^1 \hat{w}|_{m, \hat{K}} \leq c |\hat{w}|_{2, \hat{K}} \quad (*)$$

(From Poincare like inequality, e.g. for  $v \in H_0^1(\Omega)$ ,

$c \|v\|_{L^2} \leq \|\vec{\nabla} v\|_{L^2}$ : under certain circumstances, we can estimate the norms of lower derivatives by those of the higher ones)

- Derive estimate on  $K$  from estimate on  $\hat{K}$ : Let  $v \in H^2(K)$  and set  $\hat{v} = v \circ T_K$ . We know that  $(\mathcal{I}_K^1 v) \circ T_K = \mathcal{I}_{\hat{K}}^1 \hat{v}$ .

$$|v - \mathcal{I}_K^1 v|_{m, K} \leq c \|J_K^{-1}\|^m |\det J_K|^{\frac{1}{2}} |\hat{v} - \mathcal{I}_{\hat{K}}^1 \hat{v}|_{m, \hat{K}} \quad (\text{Lemma E})$$

$$\leq c \|J_K^{-1}\|^m |\det J_K|^{\frac{1}{2}} |\hat{v}|_{2, \hat{K}} \quad (*)$$

$$\leq c \|J_K^{-1}\|^m \|J_K\|^2 |v|_{2, K} \quad (\text{Lemma E})$$

$$= c (\|J_K\| \cdot \|J_K^{-1}\|)^m \|J_K\|^{2-m} |v|_{2, K}$$

$$\leq ch_K^{2-m} \sigma_K^m |v|_{2, K} \quad (\text{Lemma T})$$

## Local interpolation: special cases

- $m = 0$ :  $|v - \mathcal{I}_K^1 v|_{0,K} \leq ch_K^2 |v|_{2,K}$
- $m = 1$ :  $|v - \mathcal{I}_K^1 v|_{1,K} \leq ch_K \sigma_K |v|_{2,K}$

# Global interpolation error estimate

**Theorem** Let  $\Omega$  be polyhedral, and let  $\mathcal{T}_h$  be a shape regular family of affine meshes. Then there exists  $c$  such that for all  $h, v \in H^2(\Omega)$ ,

$$\|v - \mathcal{I}_h^1 v\|_{L^2(\Omega)} + \sum_{m=1}^2 h^m \left( \sum_{K \in \mathcal{T}_h} |v - \mathcal{I}_h^1 v|_{m,K}^2 \right)^{\frac{1}{2}} \leq ch^2 |v|_{2,\Omega}$$

and

$$\lim_{h \rightarrow 0} \left( \inf_{v_h \in V_h^1} \|v - v_h\|_{L^2(\Omega)} \right) = 0$$

# Global interpolation error estimate for Lagrangian finite elements, $k = 1$

- Assume  $v \in H^2(\Omega)$ , e.g. if problem coefficients are smooth and the domain is convex

$$\begin{aligned} \|v - \mathcal{I}_h^1 v\|_{0,\Omega} + h|v - \mathcal{I}_h^1 v|_{1,\Omega} &\leq ch^2|v|_{2,\Omega} \\ |v - \mathcal{I}_h^1 v|_{1,\Omega} &\leq ch|v|_{2,\Omega} \end{aligned}$$

$$\lim_{h \rightarrow 0} \left( \inf_{v_h \in V_h^1} |v - v_h|_{1,\Omega} \right) = 0$$

- If  $v \in H^2(\Omega)$  cannot be guaranteed, estimates become worse.  
Example: L-shaped domain.
- These results immediately can be applied in Cea's lemma.



# Error estimates for homogeneous Dirichlet problem

- Search  $u \in H_0^1(\Omega)$  such that

$$\int_{\Omega} \delta \vec{\nabla} u \vec{\nabla} v \, d\vec{x} = \int_{\Omega} f v \, d\vec{x} \quad \forall v \in H_0^1(\Omega)$$

Then,  $\lim_{h \rightarrow 0} \|u - u_h\|_{1,\Omega} = 0$ . If  $u \in H^2(\Omega)$  (e.g. on convex domains) then

$$\|u - u_h\|_{1,\Omega} \leq ch |u|_{2,\Omega}$$

$$\|u - u_h\|_{0,\Omega} \leq ch^2 |u|_{2,\Omega}$$

Under certain conditions (convex domain, smooth coefficients) one also has

$$\|u - u_h\|_{0,\Omega} \leq ch |u|_{1,\Omega}$$

(“Aubin-Nitsche-Lemma”)

- $u \in H^2(\Omega)$  may be *not* fulfilled e.g.
  - if  $\Omega$  has re-entrant corners
  - if on a smooth part of the domain, the boundary condition type changes
  - if problem coefficients ( $\delta$ ) are discontinuous
- Situations differ as well between two and three space dimensions
- Delicate theory, ongoing research in functional analysis
- Consequence for simulations
  - Deterioration of convergence rate
  - Remedy: local refinement of the discretization mesh
    - using a priori information
    - using a posteriori error estimators + automatic refinement of discretization mesh

## Higher regularity

- If  $u \in H^s(\Omega)$  for  $s > 2$ , convergence order estimates become even better for  $P^k$  finite elements of order  $k > 1$ .
- Depending on the regularity of the solution the combination of grid adaptation and higher order ansatz functions may be successful