

# Chapter 1

## Course Overview

We had a stance that our workflow in numerical analysis goes in the following way.

1. We are given either a mathematical definition, or a theorem stating the existence.
  - (a) Example: The definition of *rank* of a matrix.
  - (b) Example: For an equation  $u'(t) = f(t, u(t))$ ,  $u(0) = \alpha$  with  $f$  smooth, there exists  $T_{max} > 0$  and a solution in  $[0, T_{max}]$ .
2. Not all of them come with *methodology* to capture them.
3. If one comes with a methodology, we specify the algorithm and implement it.
4. If an implementation comes as an *approximation*, we should be able to estimate the size of error.

We will keep the same stance in the course.

Based on what we have learned, we could apply our study to two major fields of mathematics

1. many notions, definitions, and theorems in Linear Algebra,
2. existence theorem of a certain partial differential equation.

1. Each of the two would take one or two semesters.
2. So that students who studies PDE course along with this course can benefit from here,
3. We set our objective to be to implement what are studied in Chapter 6 and 7 in Evans textbook.
4. This means two things:
  - (a) We work with a pde that is either elliptic or parabolic, and the method we employ is the Galerkin method.
  - (b) The important class of pde of hyperbolic type will not be studied. Hyperbolic pde solver has to be implemented differently, where one needs to reflect knowledge of theory of hyperbolic pdes.

The objective of our course is from the following theorem.

We consider the Initial Boundary Value Problem of the following pde:

$$\begin{aligned} - \sum_{i,j=1}^n \partial_{x_i} (a^{ij}(x) \partial_{x_j} u(x)) &= \varphi(x), \quad x \in \Omega, \\ - \sum_{i,j=1}^n (a^{ij}(x) \partial_{x_j} u(x)) \nu_i &= \psi(x), \quad x \in \partial\Omega \end{aligned} \tag{P}$$

with constraint  $\int_{\partial\Omega} \psi = \int_{\Omega} \varphi$ .

We will specify assumptions on  $\Omega$ , coefficients  $a^{ij}(x)$ , and (r-h-s)  $\varphi(x)$  and  $\psi(x)$ .

**Theorem 1.** *There exists a solution of the boundary value problem (P).*

We will make the course as parallel as possible to the previous one.

1. The problem, taken as a root-finding problem for an equation

$$F(x) = y \quad \text{for given } y,$$

is tackled in 2nd half of the course.

2. The first half of the course is to extend our knowledge on piecewise polynomial functions to the multi-dimensional settings.

## The first half: pp functions in $\Omega \subset \mathbb{R}^n$ .

In this course, unless otherwise specified,  $\Omega \subset \mathbb{R}^n$  is a bounded open set with smooth boundary, that is simply connected. Also  $n = 3$  in most cases.

This is to extend our far reaching 1d remainder theorem into multi-dimensional setting. Recall

**Theorem 2** (1d remainder theorem). *Let  $f \in C^{n+1}([a, b])$  and  $x_0, x_1, \dots, x_n \in [a, b]$ . Then, for  $x \in [a, b]$ ,*

$$\begin{aligned} R(x) &= f[x_0, x_1, x_2, \dots, x_n, x](x - x_0)(x - x_1) \cdots (x - x_{n-1})(x - x_n) \\ &= f(x) - \left( f[x_0] + f[x_0, x_1](x - x_0) + f[x_0, x_1, x_2](x - x_0)(x - x_1) + \cdots \right. \\ &\quad \left. + f[x_0, x_1, \dots, x_n](x - x_0)(x - x_1)(x - x_2) \cdots (x - x_{n-1}) \right). \end{aligned}$$

The extension to multi-dimensional setting is not at all trivial.

### Partitioning of $\Omega$

To partition  $\Omega$  into small pieces is not trivial. We will be speaking of the *Simplicial complex*, tetrahedrons in  $\mathbb{R}^3$  for example, borrowing language of *Combinatorial Algebraic Topology*.

### Kinds of functions on $\Omega$

In  $\mathbb{R}^3$ , not only the function  $f : \Omega \subset \mathbb{R}^3 \rightarrow \mathbb{R}$  but also for instance the vector field

$$E : \Omega \subset \mathbb{R}^3 \rightarrow \mathbb{R}^3,$$

is relevant. We will be speaking of  $k$ -covector fields, for  $k = 0, 1, 2, 3$ , borrowing language of *Differential Geometry*.

### Remainder formula or remainder estimates

We present Poincare inequality, Bramble-Hilbert inequality, etc.

## The second half: a rootfinding problem framework

We recall the set up for a rootfinding problem for  $F : X \rightarrow Y$ , solving for a given  $y \in Y$  such that

$$F(x) = y.$$

1. We will specify  $X$  and  $Y$  for the pde.
2. We will recall the design of the solver. This includes
3. Consistency of the parametrized approximation  $(\tilde{F}_h)$ .
4. Consistency of the parametrized approximate solver  $(\tilde{R}_h)$ .
5. Continuity properties of solver.
6. A priori error estimate.
7. A posteriari error estimate.

## Chapter 2

# Partitioning of Domain

- In 1-d, the domain  $[a, b]$  is partitioned simply by specifying points in ascending order

$$a = x_0 < x_1 < x_2 < \cdots < x_n = b,$$

producing  $n$  small intervals

$$I_j = [x_{j-1}, x_j] \quad j = 1, 2, \dots, n.$$

- The role of the interval in 1-d is taken by the oriented triangle in 2-d, by the oriented tetrahedron in 3-d, and by the oriented  $n$ -simplex in  $n$ -d.
- We first introduce a few notions to speak of partitioning a domain.

## ***k*-cells**

- A polyhedral convex set is a finite intersection of closed half spaces of  $\mathbb{R}^n$ .
- A nonempty polyhedral convex set that is  $k$ -dimensional and compact is called a  $k$ -cell.
  1. We make use of the fact that the dimension of any convex set is well-defined.
  2. The  $k$ -dimensional area of  $k$ -cell  $\tau$  is thus

$$0 < \mathcal{H}^k(\tau) < \infty.$$

### **An open set $\Omega \subset \mathbb{R}^n$ equipped with the partition by $n$ -cells**

We say  $(\Omega, \mathcal{P})$  of an open set  $\Omega \subset \mathbb{R}^n$  and a set of  $n$ -cells  $\mathcal{P}$  is an  $n$ -cell partition of  $\Omega$  if

1. If  $K_1, K_2 \in \mathcal{P}$  and  $K_1 \neq K_2$  then  $\text{int}(K_1) \cap \text{int}(K_2) = \emptyset$ .
2.  $\bar{\Omega} = \bigcup_{K \in \mathcal{P}} K$ .

### **Examples of $n$ -cell partition**

CF. For an open set with smooth curved boundary, we in general consider a partition of  $\Omega$  by a homeomorphic image of  $n$ -cell, not  $n$ -cell itself. For simplicity, we assume  $\Omega$  is just a union of  $n$ -cells, omitting the flattening procedure.

1. A partition of  $\Omega$  by  $n$ -cells works fine. But we may want more structured partitionining of  $\Omega$ .
2.  $(\Omega, \mathcal{P})$  may be said to be inconvenient in the following sense.
  - (a)  $n$ -cells in  $\mathcal{P}$  are not conformal to each other. For example in 3-d, some may be tetrahedrons, some may be cubes, octahedrons, and so on.
  - (b) It is not suitable to define the boundary faces.

### Examples of $n$ -cell partition

In our course, we work with  $n$ -simplices. We consider  $(\Omega, \mathcal{S})$  an open set  $\Omega \subset \mathbb{R}^n$  equipped with a *simplicial complex*. We borrow terminology from combinatorial topology, which describes every (topological) detail of  $\Omega$  in precise manner.

## **$k$ -simplex**

1. As a generalization of an interval in 1-d, we define the oriented  $k$ -simplex.
2. (Warning: Let us not bother too much about definitions arising one after another here.)
3. Let  $k \in \{0, 1, \dots, n\}$ . A  $k$ -simplex is a juxtaposition of  $k + 1$  points

$$[x_0 x_1 x_2 \cdots x_k],$$

where  $x_0, x_1, \dots, x_k$  are affine independent.

4.  $x_0, x_1, \dots, x_k$  are affine independent if

$$\lambda_0 + \lambda_1 + \cdots + \lambda_k = 0 \quad \text{and} \quad \lambda_0 x_0 + \lambda_1 x_1 + \cdots + \lambda_k x_k = 0 \implies \lambda_0 = \lambda_1 = \cdots = \lambda_k = 0.$$

This notion is to say that any point of them are not lied in the plane that rest make up.

5. For a  $\tau = [x_0 x_1 x_2 \cdots x_k]$ , we define its closed point set

$$\overline{\text{pts}}(\tau) = \overline{\text{conv}}\{x_0, x_1, \dots, x_k\}.$$

**1d example, 2d example**

We can consider a partition whose  $n$ -cells are all the closed point sets of  $n$ -simplices.

- But we can simply do better by storing all the information by the notion of simplicial complex we soon introduce, and more importantly
- $n$ -simplices partitioning still does not resolve the problem of defining the boundary face.

## A simplicial complex $S$ satisfying further assumptions.

We consider a set  $S$  of closed point sets of  $k$ -simplices for  $k \in \{0, 1, 2, \dots, n\}$  that satisfies the following conditions:

1. If  $A, B \in S$  then either  $A \cap B = \emptyset$  or the intersection is a common face of both  $A$  and  $B$ .
2. If  $A \in S$  then every face of  $A$  is also included in  $S$ .
3. Every  $A \in S$  whose dimension is less than  $n$  is a face of some  $\sigma \in S$  of dimension  $n$ .

We observe from the definition following:

1. Up to item 1, the hanging node problem is resolved.
2. We are able to speak of boundary faces now, which are all stored in  $S$  by the condition item 2. A set  $S$  satisfying item 1 and 2 are called a simplicial complex.

### Example

3.  $S$  satisfying further the item 3 is suitable for our purpose.

### Example

Eventually, we consider  $(\Omega, \mathcal{S})$  where  $\mathcal{S}$  is a simplicial complex satisfying additionally the condition item 3, such that

$$\begin{aligned} \bar{\Omega} &= |\mathcal{S}| \quad \text{that is} \\ &= \bigcup_{A \in \mathcal{S}} A \\ &= \bigcup_{\sigma \in S_n} \sigma, \quad S_k = \{A \in \mathcal{S} \mid \dim(A) = k\} \quad \text{for } k = 0, 1, 2, \dots, n. \end{aligned}$$

Inspite of all efforts borrowing terminology from the combinatorial topology, from now on we assume  $(\Omega, \mathcal{S})$  is given such that  $\Omega$  is simply connected and bounded.

## Implementing $\mathcal{S}$

We can store simplices of  $\mathcal{S}$  in the following manner.

- We store  $\mathcal{S}_0$  and  $\mathcal{S}_n$  legitimately.
- We identify the set  $S_0$  as the set of coordinates  $x = (x^0, x^1, x^2, \dots, x^n)$  and store it.
  1. Let  $n_0$  = number of elements in  $\mathcal{S}_0$ .
  2. Consider enumeration of  $\mathcal{S}_0$  by  $i = 1, 2, \dots, n_0$ .
- We identify the set  $S_n$  as the  $(n + 1)$ -tuples

$$[i_0 i_1 i_2 \dots i_n], \quad i_0, i_1, \dots, i_n \in \{1, 2, \dots, n_0\}.$$

- Now from  $k = n - 1$  to  $k = 1$  we can store  $S_k$  by the following manner.

$$S_{k-1} = \{\text{boundary faces of } a \mid a \in S_k\}.$$

- Along with this, one can store for each  $f \in S_{k-1}$

$$\{a \in S_k \mid f \text{ is a boundary face of } a\}$$

## Chapter 3

# Polynomials on $n$ -simplex

- Started from the preceding chapter, we are in the program of implementing an approximation of a given function  $v$  defined in  $\Omega$ .
- On  $\Omega$ , one thinks of real-valued functions, vector fields, and so on.
- For a while, we first consider a set of smooth real-valued functions defined on  $\Omega$ ,

$$\Lambda_0(\Omega) = C^\infty(\bar{\Omega}).$$

- We recall the thumb rules in making approximation from data:
  1. Under the limited number of available (sampling) data, do the piecewise low order polynomial approximation rather than one high order polynomial approximation.
  2. If we go for the piecewise approximation, in one such a small domain, choose the preferable sampling points and the preferable basis whenever possible.
- Following the thumb rule, we did the partitioning of  $\Omega$  into small nice  $n$ -simplices.
- Now we discuss polynomials on a  $n$ -simplex.
- We consider an  $n$ -simplex  $\sigma = [x_0 x_1 x_2 \cdots x_n]$  and its point set

$$M = \overline{\text{pts}}(\sigma).$$

- We first consider real-valued functions defined on  $M$ ,

$$v \in \Lambda_0(M) = C^\infty(M).$$

- We consider the subspace of  $\Lambda_0$  that consists of polynomials of order at most  $m$ ,

$$\mathbb{P}_m(M) \subset \Lambda_0(M).$$

- We consider a problem of choosing an element  $p \in \mathbb{P}_m(M)$  for an approximation of  $v \in \Lambda_0(M)$ , out of suitable number of sampling data

$$(x_i, v(x_i)), \quad x_i \in M, \quad i = 1, 2, \dots, d.$$

## Polynomials in $M$ and Multi index notation

As an example, let us consider a second order polynomial in  $\mathbb{R}^2$  that is written as

$$p(x, y) = ax^2 + bxy + cy^2 + dx + ey + f, \quad a, b, c, d, e, f \in \mathbb{R}$$

of three quadratic terms, two linear terms, and one constant term.

To study polynomials in  $\mathbb{R}^n$  in a systematic way, we introduce the multi index.

### Multi Index

We introduce a convenient notation for a polynomial in multi dimensions.

- A multi index  $\alpha$  is an  $n$ -tuple of nonnegative integers  $\alpha_1, \alpha_2, \dots, \alpha_n$ ,

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n) \in (\mathbb{N} \cup \{0\})^n.$$

- We let

$$x^\alpha = (x_1, x_2, \dots, x_n)^{(\alpha_1, \alpha_2, \dots, \alpha_n)} = x_1^{\alpha_1} x_2^{\alpha_2} x_3^{\alpha_3} \cdots x_n^{\alpha_n} \in \mathbb{R}.$$

- The degree or the order of  $\alpha$  is

$$|\alpha| = \alpha_1 + \alpha_2 + \cdots + \alpha_n \in \mathbb{N} \cup \{0\}.$$

- A homogeneous  $r$ -th order polynomial is thus a linear combination of

$$\{x^\alpha \mid |\alpha| = r\}.$$

- For a fixed  $r \in \mathbb{N} \cup \{0\}$ , how many distinct multi indices with degree  $r$  are there ?

This is to choose  $r$  numbers out of  $\{1, 2, \dots, n\}$  with repetition allowed,

$$d_{n,r} = \binom{n+r-1}{r}.$$

- An element  $p \in \mathbb{P}_m(M)$  of polynomials of order at most  $m$  is thus written as

$$p(x) = \sum_{0 \leq |\alpha| \leq m} c_\alpha x^\alpha, \quad c_\alpha \in \mathbb{R} \text{ is a coefficient.}$$

- We note that

$$\mathbb{P}_m(M) \simeq \mathbb{R}^d, \quad d = \sum_{r=0}^m d_{n,r}.$$

In fact,  $d$  is to choose  $m$  numbers for the exponents of  $\{1, x_1, x_2, \dots, x_n\}$  of  $n+1$  elements with repetition allowed, and thus

$$d = \sum_{r=0}^m d_{n,r} = \binom{n+m}{m}.$$

- For later purposes, we also define the factorial

$$\alpha! = \alpha_1! \alpha_2! \cdots \alpha_n! \in \mathbb{R}, \quad \text{having in mind that } 0! = 1.$$

**Examples:**  $n = 2$ .

- $\mathbb{P}_0(M)$  is of dimension 1,

$$d_{2,0} = 1.$$

- $\mathbb{P}_1(M)$  is of dimensions

$$1 + d_{2,1} = 1 + 2 = 3.$$

- $\mathbb{P}_2(M)$  is of dimensions

$$1 + 2 + d_{2,2} = 1 + 2 + 3 = 6.$$

**Examples:**  $n = 3$ .

- $\mathbb{P}_0(M)$  is of dimension 1,

$$d_{3,0} = 1.$$

- $\mathbb{P}_1(M)$  is of dimensions

$$1 + d_{3,1} = 1 + 3 = 4.$$

- $\mathbb{P}_2(M)$  is of dimensions

$$1 + 3 + d_{3,2} = 1 + 3 + 6 = 10.$$

Hence, for  $M \subset \mathbb{R}^2$ , to fix an element in  $\mathbb{P}_0(M)$ ,  $\mathbb{P}_1(M)$ , and  $\mathbb{P}_2(M)$ , we need to provide sampling data  $(x_i, v(x_i))$  respectively as many as 1, 3, and 6.

Hence, for  $M \subset \mathbb{R}^3$ , to fix an element in  $\mathbb{P}_0(M)$ ,  $\mathbb{P}_1(M)$ , and  $\mathbb{P}_2(M)$ , we need to provide sampling data  $(x_i, v(x_i))$  respectively as many as 1, 4, and 10.

## Choosing a Basis of $\mathbb{P}_m(M)$

- We do not use *power basis*

$$\{x^\alpha \mid 0 \leq |\alpha| \leq m\}$$

- We use *Lagrange basis* (nodal basis): for each  $i = 1, 2, \dots, d$ , we choose sampling points  $(x_i)_{i=1}^d$  and basis functions  $\theta_i \in \mathbb{P}_m(M)$  so that

$$\theta_i(x_j) = \begin{cases} 0 & \text{if } i \neq j \\ 1 & \text{if } i = j. \end{cases}$$

1. This is because, of course, if we are given sampling data  $(x_i, v(x_i))$ ,  $i = 1, 2, \dots, d$ , we immediately pick up an element in  $\mathbb{P}_m(M)$  that is

$$x \mapsto \sum_{i=1}^d v(x_i) \theta_i(x),$$

compatible with the sampling data.

2. Then it matters that how we choose the sampling points  $x_i$ ,  $i = 1, 2, \dots, d$ . In our course, we will not be bothered too much on the choice of preferable sampling points unlike in 1-d.

### Examples of sampling points in triangle and in tetrahedron

$$\mathbb{P}_0(M), \quad \mathbb{P}_0(M), \quad \mathbb{P}_2(M).$$

- Such nodal basis, as well as other basis in many cases, are better expressed in the *barycentric coordinate system* rather than the given  $\mathbb{R}^n$ -coordinate system. We specify the barycentric coordinate system for a given  $k$ -simplex now.

## Barycentric coordinate system for $k$ -simplex

**Example:** line passing  $x_0$  and  $x_1$

$$\ell : \{(1 - \lambda)x_0 + \lambda x_1 \mid \lambda \in \mathbb{R}\}.$$

For a given  $k$ -simplex  $\tau = [x_0 x_1 x_2 \cdots x_k]$ , there is the unique  $k$ -dimensional plane where  $\tau$  is lied. It is a set of points expressed by a combination

$$P(\tau) = \{\lambda_0 x_0 + \lambda_1 x_1 + \lambda_2 x_2 + \cdots + \lambda_k x_k \in \mathbb{R}^n \mid \lambda_i \in \mathbb{R}, \quad \lambda_0 + \lambda_1 + \lambda_2 + \cdots + \lambda_k = 1\}$$

Now,

1. Consider a hyperplane  $\hat{P}_k$  in  $\mathbb{R}^{k+1}$  constrained by one equation:

$$\hat{P}_k = \{\lambda \in \mathbb{R}^{k+1} \mid \lambda_0 + \lambda_1 + \lambda_2 + \cdots + \lambda_k = 1\} \subset \mathbb{R}^{k+1}$$

2. The barycentric coordinate system is a parametrization from  $\hat{P}_k$  to  $P(\tau)$ :

$$\chi : \hat{P}_k \rightarrow P(\tau), \quad \lambda \mapsto \lambda_0 x_0 + \lambda_1 x_1 + \lambda_2 x_2 + \cdots + \lambda_k x_k.$$

3. In particular, the parametrization of  $\overline{\text{pts}}(\sigma)$  is from the set

$$\chi^{-1}(\overline{\text{pts}}(\sigma)) = \{(\lambda_0, \lambda_1, \dots, \lambda_k) \in \hat{P}_k \mid \forall i \quad \lambda_i \geq 0\} =: L_k \subset \hat{P}_k.$$

Before we specify basis functions in barycentric coordinate system, we get familiar with them by a few observations:

We let  $n = 3$  and consider a tetrahedron  $\sigma = [x_0x_1x_2x_3]$ .

1. Faces parametrized by  $\lambda = (\lambda_0, \lambda_1, \lambda_2, \lambda_3) \in L_3$ .

2. Level sets of  $\lambda_e$ ,  $e = 0, 1, 2, 3$ .

3. The point  $x_c = \frac{1}{4}(x_0 + x_1 + x_2 + x_3)$  corresponds to  $\lambda = \left(\frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4}\right)$ . We will show that  $x_c$  is the center of mass.

## The map from $x$ to $\lambda$ .

1. We should be able to obtain for a given  $x \in \overline{\text{pts}}(\tau)$  the  $(\lambda_0, \lambda_1, \dots, \lambda_k)$ .
2. This can be simply done as below, which is numerically not preferable.

For a given  $(\lambda_0, \lambda_1, \dots, \lambda_k) \in L_k$ ,

$$\begin{aligned} x(\lambda) &= \lambda_0 x_0 + \lambda_1 x_1 + \dots + \lambda_k x_k \\ &= x_0 + \lambda_1(x_1 - x_0) + \lambda_2(x_2 - x_0) + \dots + \lambda_k(x_k - x_0) \\ &= x_0 + \begin{pmatrix} & & & & \\ | & | & \cdots & | & \\ x_1 - x_0 & x_2 - x_0 & \cdots & x_k - x_0 & \\ | & | & \cdots & | & \\ & & & & \end{pmatrix} \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_k \end{pmatrix} \end{aligned}$$

By the affine independence assumption, the matrix must be invertible. This gives that

$$\begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \vdots \\ \lambda_k \end{pmatrix} = \begin{pmatrix} & & & & \\ | & | & \cdots & | & \\ x_1 - x_0 & x_2 - x_0 & \cdots & x_k - x_0 & \\ | & | & \cdots & | & \\ & & & & \end{pmatrix}^{-1} (x(\lambda) - x_0),$$

$$\lambda_0 = 1 - \lambda_1 - \lambda_2 - \dots - \lambda_k.$$

3. Importantly, we record here that  $x \mapsto \lambda$  map is just linear.
4. The inverse map as generalization of internal dividing point in a line segment.

- (a) We denote the  $k$ -volume of  $[x_0 x_1 x_2 \dots x_k]$  by

$$|[x_0 x_1 x_2 \dots x_k]| = \mathcal{H}^k(\overline{\text{pts}}([x_0 x_1 x_2 \dots x_k]))$$

- (b) Then, the  $i$ -th barycentric coordinate is computed by the volume ratio

$$\lambda_i = \frac{|[x_0 x_1 \dots x_{i-1} x x_{i+1} x_{i+2} \dots x_k]|}{|[x_0 x_1 x_2 \dots x_k]|}.$$

We will prove this in the next lecture.

## Nodal basis functions for $\mathbb{P}_0(M)$ , $\mathbb{P}_1(M)$ , $\mathbb{P}_2(M)$

- Here, we make use of  $x \mapsto \lambda_e(x)$  for  $e = 0, 1, 2, \dots, n$ , the barycentric coordinates, to express the nodal basis functions.
- It is nice to know that, for given degree  $m$ , there are exactly  $d$  elements in the set

$$\{\lambda_0(x)^{m_0} \lambda_1(x)^{m_1} \cdots \lambda_n(x)^{m_n} \mid m_0 + m_1 + \cdots + m_n = m\}$$

and that each of  $x \mapsto \lambda_e(x)$  is linear in  $x$ .

- We proceed with  $n = 3$  for  $\mathbb{P}_0(M)$ ,  $\mathbb{P}_1(M)$ , and  $\mathbb{P}_2(M)$ .

### $\mathbb{P}_0(M)$

The only basis function of  $\mathbb{P}_0(M)$  is simply a constant function

$$\theta(x) \equiv 1.$$

### $\mathbb{P}_1(M)$

For  $\mathbb{P}_1(M)$ , we define 4 basis functions to be

$$\theta_e(\lambda) = \lambda_e, \quad e = 0, 1, 2, 3.$$

- Note that

$$\theta_e(\hat{e}') = \begin{cases} 0 & \text{if } e \neq e' \\ 1 & \text{if } e = e'. \end{cases}$$

where  $\hat{e}$  is the  $e$ -th coordinate basis.

### $\mathbb{P}_2(M)$

For  $\mathbb{P}_2(M)$ , we define the 10 basis functions to be

$$\begin{aligned} \lambda_e(2\lambda_e - 1), \quad e = 0, 1, 2, 3, \\ 4\lambda_e\lambda_{e'}, \quad e, e' = 0, 1, 2, 3, \quad e \neq e'. \end{aligned}$$

## Summary up to now

1. We are given

$$\bar{\Omega} = |\mathcal{S}| = \bigcup_{M \in S_n} M, \quad \text{interiors of elements in } S_n \text{ are pairwise disjoint.}$$

2. The objective is to be able to implement an approximation of a function

$$v : \bar{\Omega} \rightarrow \mathbb{R},$$

where  $v \in C^\infty(\bar{\Omega})$ .

3. To ends this, we first consider a local objective to be able to implement an approximation of a function

$$v : M \rightarrow \mathbb{R}, \quad v \in C^\infty(M), \quad M \in S_n.$$

4. For each  $v \in C^\infty(M)$ , the sampling procedure is

$$v \mapsto (x_i, v(x_i))_{i=1}^d, \quad (x_i)_{i=1}^d \text{ of points in } M \text{ are sampling points.}$$

5. A polynomial of order at most  $m$  is a linear combination

$$\sum_{0 \leq |\alpha| \leq m} c_\alpha x^\alpha$$

and  $\text{span} \langle x^\alpha \rangle_{0 \leq |\alpha| \leq m}$  is a vector space of dimensions

$$d = \binom{n+m}{m}.$$

6. For given sample data, we pick up an element in  $\mathbb{P}_m(M)$ . This will be done by selecting basis functions  $\theta_i(x)$  as we want so that the element is

$$x \mapsto \sum_{i=1}^d v(x_i) \theta_i(x).$$

The sampling procedure and picking up procedure are combined:

$$I : C^\infty(M) \rightarrow \mathbb{P}_m(M) \subset C^\infty(M).$$

7. Selecting  $d$  basis of  $\mathbb{P}_m(M)$  we want:

Let  $n = 3$  and fix  $M$ .

Below,  $\lambda$  will be composited with the linear bijective map  $x \mapsto \lambda$ .

- (a) Let  $m = 0$ . Then

$$\mathbb{P}_0(M) = \text{span} \langle \mathbf{1} \rangle.$$

- (b) Let  $m = 1$ . Then  $\mathbb{P}_1(M)$  is spanned by four functions

$$(\lambda_0, \lambda_1, \lambda_2, \lambda_3) \mapsto \lambda_e, \quad e = 0, 1, 2, 3.$$

- (c) Let  $m = 2$ . Then  $\mathbb{P}_2(M)$  is spanned by ten functions

$$\begin{aligned} (\lambda_0, \lambda_1, \lambda_2, \lambda_3) &\mapsto \lambda_e(2\lambda_e - 1), \quad e = 0, 1, 2, 3 \quad \text{and} \\ (\lambda_0, \lambda_1, \lambda_2, \lambda_3) &\mapsto 4\lambda_e \lambda_{e'}, \quad e, e' = 0, 1, 2, 3 \quad e \neq e' \end{aligned}$$



## Chapter 4

# Vector fields and etc.

Let  $n = 3$ .

- Unlike in 1-d case, where the approximation target was a function

$$f : [a, b] \subset \mathbb{R} \rightarrow \mathbb{R},$$

the approximation targets defined on  $\Omega$  include not only  $\mathbb{R}$ -valued functions but include also vector fields, and so on.

- This time, we borrow language from *Differential Geometry*, to identify the target objectives of

$$\Lambda_0(\Omega), \quad \Lambda_1(\Omega), \quad \Lambda_2(\Omega), \quad \Lambda_3(\Omega)$$

that are the vector spaces of  $k$ -covector fields.

- We introduce what are the  $k$ -covectors and what are the  $k$ -covector fields below, which is very crude explanation in the Euclidean space.

## *k*-vectors

We let the vector space  $V$  to be  $\mathbb{R}^3$ .

- A real number  $a \in \mathbb{R}$  is called a 0-vector.
- An element of  $V$  is called a 1-vector.
- The space  $V$  of 1-vectors is equipped with the inner product

$$v_1^T v_2.$$

- The space of 2-vectors is

$$\text{span} \langle e_1 \wedge e_2, e_1 \wedge e_3, e_2 \wedge e_3 \rangle.$$

- For two 1-vectors  $v_1$  and  $v_2$ , the notation

$$v_1 \wedge v_2$$

denotes an element of space of 2-vectors in the following sense.

- $(av_1 \wedge bv_2) = ab v_1 \wedge v_2$ , and distribution law works for two linear combinations of vectors in  $V$ .
- The notation  $v_2 \wedge v_1$  is identified by

$$-v_1 \wedge v_2$$

in the space of 2-vectors.

- Calculus:

If  $v_1 \wedge v_2$  and  $w_1 \wedge w_2$  are given, where  $v_1, v_2, w_1, w_2 \in V$ , we define the inner product

$$\langle v_1 \wedge v_2, w_1 \wedge w_2 \rangle = \det \left( \begin{pmatrix} - & v_1^T & - \\ - & v_2^T & - \\ | & | & | \end{pmatrix} \begin{pmatrix} | & | \\ w_1 & w_2 \\ | & | \end{pmatrix} \right)$$

- Inner product of two linear combinations of such forms are then computed following distribution law.

- The space of 3-vectors is

$$\text{span} \langle e_1 \wedge e_2 \wedge e_3 \rangle.$$

- Calculus works in the same principle.

- If  $v_1 \wedge v_2 \wedge v_3$  and  $w_1 \wedge w_2 \wedge w_3$  are given, where  $v_1, v_2, v_3, w_1, w_2, w_3 \in V$ , we define the inner product

$$\langle v_1 \wedge v_2 \wedge v_3, w_1 \wedge w_2 \wedge w_3 \rangle = \det \left( \begin{pmatrix} - & v_1^T & - \\ - & v_2^T & - \\ - & v_3^T & - \end{pmatrix} \begin{pmatrix} | & | & | \\ w_1 & w_2 & w_3 \\ | & | & | \end{pmatrix} \right)$$

## *k*-covectors

- In the flat Euclidean geometry, a dual element is essentially the same object to the *k*-vector.

(A *k*-vector is identified by the corresponding *k*-covector by the inner product.)

- We will work with *k*-covectors for our purposes.

## *k*-covector fields

- A *k*-covector field on  $\Omega$  is a *k*-covector-valued function defined on  $\Omega$ .
- Here, the role of the vector space  $V$  is taken by  $T_x\Omega$  of tangent space at  $x \in \Omega$ , but every  $T_x\Omega$  is identified by the same space  $\mathbb{R}^3$ .
- Unless otherwise specified, we restrict ourselves in the  $C^\infty$  (up to boundary) fields. The space of such smooth *k*-covector fields are denoted by

$$\Lambda_0(\Omega), \quad \Lambda_1(\Omega), \quad \Lambda_2(\Omega), \quad \Lambda_3(\Omega)$$

- We already discussed about  $\Lambda_0(\Omega)$ .
- Elements of them are our targets of approximation in our course.

## Local polynomial spaces suitable for $(\Lambda_0, \Lambda_1, \Lambda_2, \Lambda_3)$

Let  $n = 3$  and fix  $M$  an  $n$ -simplex.

- The easiest choice for  $(\Lambda_0(M), \Lambda_1(M), \Lambda_2(M), \Lambda_3(M))$  is the following.

$$\begin{aligned}\mathbb{P}_1(M) &\subset \Lambda_0, \\ \mathbb{E}_1(M) &\subset \Lambda_1, \\ \mathbb{J}_1(M) &\subset \Lambda_2, \\ \mathbb{P}_0(M) &\subset \Lambda_3,\end{aligned}$$

- $\mathbb{E}_1(M)$  is the *Nedelec polynomial space* of lowest order, restricted in  $M$ .
- $\mathbb{J}_1(M)$  is the *Raviart-Thomas polynomial space* of lowest order, restricted in  $M$ .

We will delve into them one-by-one: Definition, Basis.

## $k = 3$ : constant approximation

- A function in  $\Lambda_3(M)$  typically represents a *density* of chemical concentration, population, mass, etc.
- In the sense that we require mere integrability of them,  $\mathbb{P}_0(M)$  works fine for  $\Lambda_3(M)$ .

### Sampling Data

- For a given function  $x \mapsto \rho(x) e_1 \wedge e_2 \wedge e_3$ , we store the value

$$L = \frac{1}{|M|} \int_M \rho(x) dx.$$

### Basis functions

- We just recall the basis function is the constant function

$$\theta : x \mapsto 1 e_1 \wedge e_2 \wedge e_3.$$

### The Projector (approximation)

- For a given  $x \mapsto \rho(x) e_1 \wedge e_2 \wedge e_3$  in  $\Lambda_3(M)$ , we let its approximation

$$\left( \frac{1}{|M|} \int_M \rho(x) dx \right) e_1 \wedge e_2 \wedge e_3 \in \mathbb{P}_0(M).$$

This defines the projector

$$I : \Lambda_3(M) \rightarrow \mathbb{P}_0(M).$$

## $k = 2$ : constant outward normals

Let  $M$  be the point set of 3-simplex  $[x_0x_1x_2x_3]$ .

- Here, we may assume that the center of mass of  $M$  is at origin, or
- if  $x_c = \frac{1}{4}(x_0 + x_1 + x_2 + x_3)$ , we let  $z = x - x_c$  and use  $z$  coordinates.

### Writing an element $J \in \Lambda_2(M)$

- We write an element in  $\Lambda_2(M)$  in the following way:

$$\begin{aligned} J : x &\mapsto J_1(x)e_2 \wedge e_3 + J_2(x)e_3 \wedge e_1 + J_3(x)e_1 \wedge e_2, \\ &= J_1(x)\check{e}_1 + J_2(x)\check{e}_2 + J_3(x)\check{e}_3. \end{aligned}$$

- This arrangement is because we want to interprete  $d$  operator for  $J$ -field as div operator:

$$dJ(x) = \left( \frac{\partial J_1}{\partial x_1}(x) + \frac{\partial J_2}{\partial x_2}(x) + \frac{\partial J_3}{\partial x_3}(x) \right) e_1 \wedge e_2 \wedge e_3.$$

- $(\mathbb{P}_1(M))^n$ , the space of first order 1-vector fields ( $(n-1)$ -vector fields as well), has the degrees of freedom  $(n+1)n$ . Indeed,

$$x \mapsto J_0 + A(x - x_c), \quad J_0 \text{ a constant vector, and } A \text{ an } n \times n \text{ matrix.}$$

In 3-d, in total 12 freedoms.

### Definition of $\mathbb{J}_1(M)$

- We define  $\mathbb{RT}(M)$  to be

$$\left\{ J_0 + c_0(x - x_c) \mid J_0 \text{ is a constant 1-vector and } c_0 \in \mathbb{R} \right\}$$

restricted in  $M$ .

- Replacing  $e_i \rightarrow \check{e}_i$  in  $\mathbb{RT}(M)$  let the element be in  $\mathbb{J}_1(M) \subset \Lambda_2(M)$ .
- For notational clarification, we let  $\mathbb{RT}(M)$  to be of 1-vector fields, and let  $\mathbb{J}_1(M)$  of 2-vector fields after the replacement  $e_i \rightarrow \check{e}_i$ .

- Observations:

1.  $\mathbb{J}_1(M)$  is a 4-dimensional subspace  $\mathbb{J}_1(M) \subset (\mathbb{P}_1(M))^3 \subset \Lambda_2(M)$ .
2.  $dJ$  for  $J = J_0 + c_0(x - x_c)$ , (after replacement), in the interior of  $M$ , is simply a constant  $3c_0$ .
3.  $\mathbb{J}_1(M)$ : A minimal requirement so that  $J$  itself and  $dJ$  both are nontrivial and in control.

### Sampling Data

- For a given function  $J : x \mapsto J_1(x)e_2 \wedge e_3 + J_2(x)e_3 \wedge e_1 + J_3(x)e_1 \wedge e_2$ , we store four values.
- We store for  $e = 0, 1, 2, 3$

$$L_e = \frac{1}{|f_{\check{e}}|} \int_{\overline{\text{pts}}(f_{\check{e}})} J^{RT}(x) \cdot \nu_{\check{e}} d\mathcal{H}^2 :$$

where

- For given  $J$ , let  $J^{RT}$  be the 1-vector field where  $\check{e}_i$  is replaced by  $e_i$ .
- $\nu_{\check{e}}$  is the outward unit normal vector seen from  $M$  on  $f_{\check{e}}$ .
- There are four boundaries:

$$[x_1x_2x_3], \quad -[x_0x_2x_3], \quad [x_0x_1x_3], \quad -[x_0x_1x_2].$$

denoted by  $f_{\check{e}}$ , with  $x_e$  missing.

### Basis functions of $\mathbb{RT}(M)$

- Nodal basis: We look for a vector field  $\theta_{\check{e}}$  of the form  $J_0 + c_0(x - x_c)$  such that

$$\theta_{\check{e}}(x) \cdot \nu_{\check{e}'} = \begin{cases} \text{Const.}, & e = e' \quad \text{and} \quad x \in f_{\check{e}} \\ 0, & e \neq e' \quad \text{and} \quad x \in f_{\check{e}'}. \end{cases}$$

- We show that

$$\theta_{\check{e}} : x \mapsto \text{Const.}(x - x_e)$$

do the job.

1. If  $e \neq e'$  and  $x$  is on the face  $f_{\check{e}'}$ , then  $x - x_e$  is a tangent vector of the face  $f_{\check{e}'}$ .

Therefore, the inner product with the normal must be 0.

2. Let  $x$  be on the face  $f_{\check{e}}$ . Let  $e = 3$  for example.

The outward normal

$$\nu_3 \parallel (x_1 - x_0) \times (x_2 - x_0).$$

$x$  is a combination

$$\begin{aligned} x &= \lambda_0 x_0 + \lambda_1 x_1 + \lambda_2 x_2, \quad \lambda_0 + \lambda_1 + \lambda_2 = 1, \\ &= x_0 + \lambda_1(x_1 - x_0) + \lambda_2(x_2 - x_0). \end{aligned}$$

Hence,

$$x - x_3 = x_0 - x_3 + \lambda_1(x_1 - x_0) + \lambda_2(x_2 - x_0)$$

and

$$(x - x_3) \cdot \nu_3 = (x_0 - x_3) \cdot \nu_3$$

and this must be a nonzero constant.

- The normalizing constant can be computed. We present the result:

$$\theta_{\check{e}}^{RT}(x) = \frac{|f_{\check{e}}|}{3|M|} (x - x_e).$$

**The projector (approximation)**

We define an approximation, the projector  $I : \Lambda_2(M) \rightarrow \mathbb{J}_1(M)$  that is

$$J \quad \mapsto \quad \sum_{e=0}^3 L_e \theta_{\check{e}}^{RT} \quad \text{with } \check{e}_i \text{ in place of } e_i.$$

## $k = 1$ : constant tangentials on edges

Let  $M$  be the point set of 3-simplex  $[x_0x_1x_2x_3]$ .

- Here, we may assume that the center of mass of  $M$  is at origin, or
- if  $x_c = \frac{1}{4}(x_0 + x_1 + x_2 + x_3)$ , we let  $z = x - x_c$  and use  $z$  coordinates.
- $(\mathbb{P}_1(M))^n$ , the space of first order 1-vector fields has the degrees of freedom  $(n+1)n$ . Indeed,

$$z \mapsto E_0 + Az, \quad E_0 \text{ a constant vector, and } A \text{ an } n \times n \text{ matrix.}$$

In 3-d, in total 12 freedoms.

- We look for a simpler subspace with 6 degrees of freedom.

### Definition of $\mathbb{E}_1(M)$

- We define  $\mathbb{E}_1(M)$  to be

$$\left\{ E_0 + (c_1, c_2, c_3) \times (z_1, z_2, z_3) \mid E_0 \text{ is a constant 1-vector and } c_1, c_2, c_3 \in \mathbb{R} \right\}$$

restricted in  $M$ .

- Observations:

1.  $\mathbb{E}_1(M)$  is a 6-dimensional subspace  $\mathbb{E}_1(M) \subset (\mathbb{P}_1(M))^3 \subset \Lambda_1(M)$ .
2.  $dE$  for  $E = E_0 + (c_1, c_2, c_3) \times (z_1, z_2, z_3)$ , in the interior of  $M$ , is simply a constant vector  $2(c_1, c_2, c_3)$ .
3.  $\mathbb{E}_1(M)$ : A minimal requirement so that  $E$  itself and  $dE$  both are nontrivial and in control.

### Sampling Data

- For a given function  $E : x \mapsto E_1(x)e_1 + E_2(x)e_2 + E_3(x)e_3$ , we store six values.
- We store for each edges, designated by  $[ee']$  in the set

$$\{[01], [02], [03], [12], [13], [23]\} = Key_1,$$

the line integrals

$$L_{[ee']} = \frac{1}{|[ee']|} \int_{\overline{\text{pts}}([ee'])} E(x) \cdot \frac{x_{e'} - x_e}{|x_{e'} - x_e|} d\mathcal{H}^1.$$

### Basis functions of $\mathbb{E}_1(M)$

- Nodal basis: We look for a vector field  $\theta_{[ee']}$  of the form  $E_0 + (c_1, c_2, c_3) \times (z_1, z_2, z_3)$  such that for  $[\alpha\alpha'] \in Key_1$

$$\theta_{[ee']}(x) \cdot \frac{x_{\alpha'} - x_\alpha}{|x_{\alpha'} - x_\alpha|} = \begin{cases} \text{Const.}, & [ee'] = [\alpha\alpha'] \quad \text{and} \quad x \in [\alpha\alpha'] \\ 0, & [ee'] \neq [\alpha\alpha'] \quad \text{and} \quad x \in [\alpha\alpha'] \end{cases}$$

- We show that

$$\theta_{[ee']} : x \mapsto \text{Const.}(x - x_\gamma) \times (x - x_{\gamma'})$$

do the job, where we let  $(e, e', \gamma, \gamma')$  is an even permutation of  $(0123)$ .

1. First, note that  $\theta_{[ee']}(x)$  is of 1st order in  $x$ , not of 2nd order:

$$(x - x_\gamma) \times (x - x_{\gamma'}) = x \times x - x \times x_{\gamma'} - x_\gamma \times x + x_\gamma \times x_{\gamma'},$$

but the quadratic term is 0.

2. Also, it is in the form of  $E_0 + (c_1, c_2, c_3) \times (x - x_c)$  because

$$x \times (x_\gamma - x_{\gamma'}) + x_\gamma \times x_{\gamma'} = (x - x_c) \times (x_\gamma - x_{\gamma'}) + x_c \times (x_\gamma - x_{\gamma'}) + x_\gamma \times x_{\gamma'}.$$

3. Now, among  $Key_1$  elements, if  $[ee'] \neq [\alpha\alpha']$  and  $x$  is on the face  $[\alpha\alpha']$ , then  $\{\alpha, \alpha'\}$  and  $\{\gamma, \gamma'\}$  has at least one element in common.

In other words, the tangent vector on  $[\alpha\alpha']$  is proportional to either  $x - x_\gamma$  or  $x - x_{\gamma'}$ .

Therefore, the inner product of  $(x - x_\gamma) \times (x - x_{\gamma'})$  with the tangent must be 0.

4. Let  $x$  be on the face  $[ee']$ . Let  $[ee'] = [01]$  for example. Then  $[\gamma\gamma'] = [23]$ .

The tangent

$$x \in [01]$$

$$\begin{aligned} & \frac{x_1 - x_0}{|x_1 - x_0|} \\ & x = \lambda_0 x_0 + \lambda_1 x_1, \quad \lambda_0 + \lambda_1 = 1, \\ & \quad = x_0 + \lambda_1(x_1 - x_0). \\ & x - x_2 = x_0 - x_2 + \lambda_1(x_1 - x_0), \\ & x - x_3 = x_0 - x_3 + \lambda_1(x_1 - x_0), \\ & (x - x_2) \times (x - x_3) = (x_0 - x_2) \times (x_0 - x_3) + \lambda_1(x_1 - x_0) \times (x_2 - x_3), \end{aligned}$$

Hence

$$\frac{x_1 - x_0}{|x_1 - x_0|} \cdot (x - x_2) \times (x - x_3) = \frac{x_1 - x_0}{|x_1 - x_0|} \cdot (x_0 - x_2) \times (x_0 - x_3)$$

and this must be a nonzero constant.

- The normalizing constant can be computed. We present the result:

$$\theta_{[ee']} = \frac{s|[ee']|}{3!|M|} (x - x_\gamma) \times (x - x_{\gamma'}),$$

where  $s$  appears as a sign factor: +1 if  $[x_0x_1x_2x_3]$  is positively oriented, and -1 if  $[x_0x_1x_2x_3]$  is negatively oriented.

### The projector (approximation)

We define an approximation, the projector  $I : \Lambda_1(M) \rightarrow \mathbb{E}_1(M)$  that is

$$E \quad \mapsto \quad \sum_{[ee'] \in Key_1} L_{[ee']} \theta_{[ee']}.$$

$$k = 0$$

- We did this before.

### Sampling Data

- For a given function  $x \mapsto v(x)$ , we store the values

$$v(x_0), v(x_1), v(x_2), v(x_3).$$

### Basis functions

- We just recall the four basis functions in barycentric coordinate are

$$\theta_e(\lambda) = \lambda_e, \quad e = 0, 1, 2, 3.$$

- This time, we prove that

$$\lambda_e(x) = \frac{|[x_0 x_1 \cdots x_{e-1} \ x \ x_{e+1} x_{e+2} \cdots x_n]|}{|[x_0 x_1 x_2 \cdots x_n]|}.$$

Note that this is a Const. multiple of the determinant of matrix with columns  $x_i - x$ , with  $e$ -th column missing.

### The Projector (approximation)

- For a given  $x \mapsto v(x)$  in  $\Lambda_0(M)$ , we let its approximation

$$\sum_{e=0}^3 v(x_e) \lambda_e(x).$$

This defines the projector

$$I : \Lambda_0(M) \rightarrow \mathbb{P}_1(M).$$

**Proposition 1.** Let  $[x_0x_1 \cdots x_n]$  be an  $n$ -simplex. Then the 0-th barycentric coordinate

$$\lambda_0(x) = \frac{|[xx_1x_2 \cdots x_n]|}{|[x_0x_1x_2 \cdots x_n]|}.$$

*Proof.* .

1. For the case  $x$  is on the face  $[x_1x_2 \cdots x_n]$ , then  $\lambda_0(x) = 0$  and the numerator in the volume ratio is also 0. Thus equality holds for this case.
2. Now we assume  $x$  is not on the face  $[x_1x_2 \cdots x_n]$ .

- We know that

$$\begin{pmatrix} \lambda_1(x) \\ \lambda_2(x) \\ \vdots \\ \lambda_n(x) \end{pmatrix} = \begin{pmatrix} | & | & | & | \\ x_1 - x_0 & x_2 - x_0 & \cdots & x_n - x_0 \\ | & | & \cdots & | \end{pmatrix}^{-1} \begin{pmatrix} | \\ x - x_0 \\ | \end{pmatrix}$$

$$= M_0^{-1}(x - x_0).$$

- For each  $e \neq 0$ , we can write

$$\begin{pmatrix} \lambda_1(x) \\ \lambda_2(x) \\ \vdots \\ \lambda_n(x) \end{pmatrix} = M_0^{-1}(x_e - x_0) + M_0^{-1}(x - x_e) \iff M_0^{-1}(x_e - x) = M_0^{-1}(x_e - x_0) - \begin{pmatrix} \lambda_1(x) \\ \lambda_2(x) \\ \vdots \\ \lambda_n(x) \end{pmatrix}.$$

- Listing aboves in the columns of matrix, we write

$$\begin{aligned} & M_0^{-1} \begin{pmatrix} | & | & | & | \\ x_1 - x & x_2 - x & \cdots & x_n - x \\ | & | & \cdots & | \end{pmatrix} \\ &= M_0^{-1} \begin{pmatrix} | & | & | & | \\ x_1 - x_0 & x_2 - x_0 & \cdots & x_n - x_0 \\ | & | & \cdots & | \end{pmatrix} - \begin{pmatrix} \lambda_1(x) & \lambda_1(x) & \cdots & \lambda_1(x) \\ \lambda_2(x) & \lambda_2(x) & \cdots & \lambda_2(x) \\ \vdots & \vdots & \ddots & \vdots \\ \lambda_n(x) & \lambda_n(x) & \cdots & \lambda_n(x) \end{pmatrix} \\ &= I - \begin{pmatrix} \lambda_1(x) \\ \lambda_2(x) \\ \vdots \\ \lambda_n(x) \end{pmatrix} (1 \ 1 \ \cdots 1). \end{aligned}$$

- The determinant of (RHS), which is of identity matrix + rank one matrix, is computed by

$$1 - \lambda_1(x) - \lambda_2(x) - \cdots - \lambda_n(x).$$

- The determinant of (LHS) is the signed volume ratio of  $[xx_1x_2 \cdots x_n]$  and  $[x_0x_1 \cdots x_n]$ . We make use of the fact that  $x$  and  $x_0$  are in the same side of  $\mathbb{R}^n$  divided by the hyperplane  $[x_1x_2 \cdots x_n]$  lies in. Thus, the sign factor must be same. The proof is done.

□

Local Analysis left to do are the remainder calculations.

For example, we will derive equality and inequality such as

For  $x \in M$ ,

$$\begin{aligned} v(x) - \frac{1}{|M|} \int_M v(y) dy &= \frac{1}{|M|} \int_M \nabla v(z) \cdot \frac{(x-z)}{|z-x|^n} (|c(x,z) - x|^n - |z-x|^n) dz, \\ \implies \left| v(x) - \frac{1}{|M|} \int_M v(y) dy \right| &\leq \frac{2 \operatorname{diam}(M)^n}{n |M|} \int_M |D^\alpha v(z)| |z-x|^{1-n} dz. \end{aligned}$$

We do the other easier part first.

## Chapter 5

# Projector to pp functions in $\Omega$

- Now, we discuss about the gluing polynomials defined piecewisely in each 3-simplices.
- We will define the projector

$$I : \Lambda_k(\Omega) \rightarrow \Pi_k(\Omega)$$

from the space of smooth  $k$ -covector fields to the space of certain pp functions.

- We recall that

$$\bar{\Omega} = \bigcup_{M \in \mathcal{S}_n} M, \quad \text{of } M \in \mathcal{S}_n \text{ whose interiors are pairwise disjoint.}$$

Preliminary task:

- To proceed, we choose a convention of orientation on every member of  $\mathcal{S}$ , so that for example  
 $M = \overline{\text{pts}}(\sigma), \quad \text{and } \sigma \text{ is an oriented 3-simplex.}$
- Let us, however, use the same symbol  $\mathcal{S}$  for the set whose members are now oriented simplices.

## $k = 0$ : Gluing Piecewise linear functions

- Gluing local basis to define global basis: For a given vertex point  $a \in S_0$ , we define a function, supported in every 3-simplex the point belongs to, by gluing their local basis.
- For every vertex point  $a \in S_0$ , we define the *global basis pp function*

$$\phi_a : x \mapsto \begin{cases} \theta_{\sigma,e}(x) & \text{if } x \in \text{int } \overline{\text{pts}}(\sigma) \text{ and } \dot{e} = a, \\ 0 & \text{otherwise.} \end{cases}$$

Here, I used the notation  $\dot{e}$  to reference the vertex point in  $S_0$  that is locally the  $e$ -th vertex of the tetrahedron  $\sigma$ .

- The value extensions on faces ?

- Define

$$\Pi_0(\Omega) = \text{span} \langle \phi_a \rangle_{a \in S_0}.$$

- Global projector (approximation)

If  $x \mapsto v(x)$  is given, we define its approximation pp function by

$$\sum_{a \in S_0} v(a) \phi_a.$$

- Dimensions of  $\Pi_0(\Omega)$ .

1. By definition,

$$\dim \Pi_0(\Omega) = n_0, \quad \text{the number of total vertices in the complex } S.$$

2. This number matches to the number obtained by gluing  $\mathbb{P}_1(M)$  with imposed constraints.
3. For every vertex  $a \in S_0$ , we let the number  $m_a$  to be its multiplicity, that is the number of distinct tetrahedrons where  $a$  belongs to as its face. If so,

$$\begin{aligned} 4n_3 &= \sum_{a \in S_0} m_a. \\ \iff \sum_{a \in S_1} (m_a - 1) &= 4n_3 - n_0. \end{aligned}$$

4. Recall that each of local d.o.f of  $\mathbb{P}_1(\overline{\text{pts}}(\sigma))$  was 4, and there are  $n_3$  such tetrahedrons.
5. In gluing function, to impose same value for every vertex, the number of total constraints are

$$\sum_{a \in S_0} m_a - 1 = 4n_3 - n_0.$$

6. Hence,

the total d.o.f - total constraints =  $n_0$ , which is the dimensions of  $\Pi_0(\Omega)$ .

## $k = 1$ : Gluing Piecewise $\mathbb{E}_1(M)$ field

- Gluing local basis to define global basis: For a given edge  $a \in S_0$ , we define a function, supported in every 3-simplex the edge belongs to, by gluing their local basis.
- For every  $a \in S_1$ , we define the *global basis pp function*

$$\phi_a : x \mapsto \begin{cases} \pm \theta_{\sigma, [ee']} (x) & \text{if } x \in \text{int } \overline{\text{pts}}(\sigma) \text{ and } [ee'] = \pm a, \\ 0 & \text{otherwise.} \end{cases}$$

- The value extensions on faces ?

- Define

$$\Pi_1(\Omega) = \text{span} \langle \phi_a \rangle_{a \in S_1}.$$

- Global projector (approximation)

If  $E : x \mapsto E_1(x)e_1 + E_2(x)e_2 + E_3(x)e_3$  is given, we define its approximation pp function by

$$\sum_{a \in S_1} \left( \frac{1}{|a|} \int_{\overline{\text{pts}}(a)} E(x) \cdot \tau_a dx \right) \phi_a.$$

- Dimensions of  $\Pi_1(\Omega)$ .

1. By definition,

$$\dim \Pi_1(\Omega) = n_1, \quad \text{the number of total edges in the complex } \mathcal{S}.$$

2. This number matches to the number obtained by gluing  $\mathbb{E}_1(M)$  with imposed constraints.
3. For every edge  $a \in S_1$ , we let the number  $m_a$  to be its multiplicity, that is the number of distinct tetrahedrons where  $a$  belongs to as its face. If so,

$$\begin{aligned} 6n_3 &= \sum_{a \in S_1} m_a \\ \iff \sum_{a \in S_1} m_a - 1 &= 6n_3 - n_1. \end{aligned}$$

4. Recall that each of local d.o.f. of  $\mathbb{E}_1(\overline{\text{pts}}(\sigma))$  was 6, and there are  $n_3$  such tetrahedrons.
5. In gluing function, to impose continuity of tangential component for every edge, the number of total constraints are

$$\sum_{a \in S_1} m_a - 1 = 6n_3 - n_1.$$

6. Hence,

the total d.o.f - total constraints =  $n_1$ , which is the dimensions of  $\Pi_1(\Omega)$ .

## $k = 2$ : Gluing Piecewise $\mathbb{RT}(M)$ field

- We proceed with the Raviart-Thomas vector field treatment for convenience. This includes also the fact that the choice of convention orientation on a 2-face  $a \in \mathcal{S}_2$  fixes the convention normal  $\nu_a$  on it.
- Gluing local basis to define global basis: For a given 2-face  $a \in \mathcal{S}_2$ , we define a function, supported in every 3-simplex the face belongs to (in fact there can be at most two), by gluing their local basis.
- For every  $a \in \mathcal{S}_2$ , we define the *global basis pp function*

$$\phi_a : x \mapsto \begin{cases} \pm \theta_{\sigma, [\check{e}]}^{RT}(x) & \text{if } x \in \text{int } \overline{\text{pts}}(\sigma) \text{ and } [ee'e''] = \pm a, (\check{e} = ee'e''), \\ 0 & \text{otherwise.} \end{cases}$$

- The value extensions on faces ?

- Define

$$\Pi_2(\Omega) = \text{span} \langle \phi_a \rangle_{a \in \mathcal{S}_2}.$$

- Global projector (approximation)

If  $J^{RT} : x \mapsto J_1^{RT}(x)e_1 + J_2^{RT}(x)e_2 + J_3^{RT}(x)e_3$  is given, we define its approximation pp function by

$$\sum_{a \in \mathcal{S}_2} \left( \frac{1}{|a|} \int_{\overline{\text{pts}}(a)} J^{RT}(x) \cdot \nu_a dx \right) \phi_a.$$

- Dimensions of  $\Pi_2(\Omega)$ .

1. By definition,

$$\dim \Pi_2(\Omega) = n_2, \quad \text{the number of total 2-faces in the complex } \mathcal{S}.$$

2. This number matches to the number obtained by gluing  $\mathbb{RT}(M)$  with imposed constraints.
3. For every edge  $a \in \mathcal{S}_2$ , we let the number  $m_a$  to be its multiplicity, that is the number of distinct tetrahedrons where  $a$  belongs to as its face. If so,

$$\begin{aligned} 4n_3 &= \sum_{a \in \mathcal{S}_2} m_a \\ \iff \sum_{a \in \mathcal{S}_2} m_a - 1 &= 4n_3 - n_2. \end{aligned}$$

4. Recall that each of local d.o.f. of  $\mathbb{RT}(\overline{\text{pts}}(\sigma))$  was 4, and there are  $n_3$  such tetrahedrons.
5. In gluing function, to impose continuity of normal component for every 2-face, the number of total constraints are

$$\sum_{a \in \mathcal{S}_2} m_a - 1 = 4n_3 - n_2.$$

6. Hence,

the total d.o.f - total constraints =  $n_2$ , which is the dimensions of  $\Pi_2(\Omega)$ .

## $k = 3$ : Piecewise constant function

- For a given 3-simplex  $\sigma \in S_3$ , we define a function supported in the simplex, the local basis of the simplex:
- For every  $\sigma$ , we define the *global basis pp function*

$$\phi_\sigma : x \mapsto \chi_{\text{int } \overline{\text{pts}}(\sigma)}(x) e_1 \wedge e_2 \wedge e_3.$$

- The value extensions on faces ?

- Define

$$\Pi_3(\Omega) = \text{span} \langle \phi_\sigma \rangle_{\sigma \in S_n}.$$

- Global projector (approximation)

If  $x \mapsto \rho(x) e_1 \wedge e_2 \wedge e_3$  is given, we define its approximation pp function by

$$\sum_{\sigma \in S_n} \left( \frac{1}{|\sigma|} \int_{\overline{\text{pts}}(\sigma)} \rho(x) dx \right) \phi_\sigma.$$

- Dimensions of  $\Pi_3(\Omega)$ .

1. By definition,

$$\dim \Pi_3(\Omega) = n_3, \quad \text{the number of total tetrahedrons in the complex } \mathcal{S}.$$

### Global remainder estimates

Denote

$$h := \sup_{\sigma \in \mathcal{S}^n} \text{diam}(\sigma).$$

Suppose the following local remainder estimate holds:

$$\begin{aligned} \|R\|_{L^p(M)} &\leq C_0 \text{diam}(M)^q \|Dv\|_{L^p(M)}, \\ \text{i.e., } \left( \int_M |R(x)|^p dx \right)^{\frac{1}{p}} &\leq C_0 \text{diam}(M)^q \left( \int_M |Dv(x)|^p dx \right)^{\frac{1}{p}}. \end{aligned}$$

Then, globally

$$\begin{aligned} \int_{\Omega} |R(x)|^p dx &= \sum_{\sigma \in \mathcal{S}_n} \int_{\overline{\text{pts}}(\sigma)} |R(x)|^p dx \leq \sum_{\sigma \in \mathcal{S}_n} (C_0 \text{diam}(\sigma)^q)^p \int_{\overline{\text{pts}}(\sigma)} |Dv(x)|^p dx \\ &\leq (C_0 h^q)^p \sum_{\sigma \in \mathcal{S}_n} \int_{\overline{\text{pts}}(\sigma)} |Dv(x)|^p dx = (C_0 h^q)^p \int_{\Omega} |Dv(x)|^p dx, \\ \implies \left( \int_{\Omega} |R(x)|^p dx \right)^{\frac{1}{p}} &\leq C_0 h^q \left( \int_{\Omega} |Dv(x)|^p dx \right)^{\frac{1}{p}}, \end{aligned}$$

i.e., we conclude

$$\|R\|_{L^p(\Omega)} \leq C_0 h^q \|Dv\|_{L^p(\Omega)}$$

with the same constant  $C_0$  and exponent  $q$ .

## Chapter 6

# Local Remainder formulas

Here, we let  $M = \overline{\text{pts}}(\sigma)$ , and  $M \subset \mathbb{R}^n$ . We derive formulas with general  $n$ .

We will derive the remainder formulas twice:

1. The simplest case is for  $k = n$ , where we used the constant approximation and the remainder is

$$R(x) = \rho(x) - \frac{1}{|M|} \int_M \rho(y) dy.$$

2. Remaining cases where approximation is higher order.

The first case will be multi-dimensional generalization of the 1d Mean Value Theorem, or the 1d Fundamental Theorem of Calculus:

$$\rho(x) - \rho(y) = \int_y^x \frac{d\rho}{dx}(t) dt.$$

Then the second case will be generalization of it to higher order and multi-dimensional. A generalization of the 1d Taylor's Theorem.

Fix  $x \in M$ , and let  $y \in M$ . Let  $v \in C^\infty(M)$ . We consider the function of one variable

$$t \mapsto f(t) = v(tx + (1-t)y), \quad t \in [0, 1] \quad \text{so that} \quad f(1) = v(x), \quad f(0) = v(y)$$

Then by the Fundamental Theorem of Calculus,

$$\begin{aligned} v(x) &= f(1) = f(0) + \int_0^1 f'(t) dt \\ &= v(y) + \int_0^1 Dv(tx + (1-t)y) \cdot (x - y) dt \end{aligned}$$

**Theorem 1.** *Let  $v \in C^\infty(M)$ . Then*

$$R(x) = \frac{1}{n|M|} \int_M Dv(z) \cdot \frac{(x-z)}{|x-z|^n} \left( |c(x, z) - x|^n - |x-z|^n \right) dz,$$

where  $c(x, z)$  is the boundary point intersecting the half ray from  $x$  to  $z$ .

*Proof.* Integrating in  $y$  variable the both sides of equation

$$v(x) - v(y) = \int_0^1 Dv(tx + (1-t)y) \cdot (x - y) dt,$$

We obtain in average

$$v(x) - \frac{1}{|M|} \int_M v(y) dy = \frac{1}{|M|} \int_M v(x) - v(y) dy = \frac{1}{|M|} \int_M \int_0^1 Dv(tx + (1-t)y) \cdot (x - y) dt dy.$$

The integrating domain is thus for  $(t, y) \in [0, 1] \times M$ . We consider for  $(t, y) \in M$  the function

$$\phi : (t, y_1, y_2, \dots, y_n) \mapsto (t, z_1, z_2, \dots, z_n) = (t, tx + (1-t)y).$$

Then, we compute

$$D\phi(t, y) = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ x_1 - y_1 & 1-t & 0 & \cdots & 0 \\ x_2 - y_2 & 0 & 1-t & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_n - y_n & 0 & 0 & \cdots & 1-t \end{pmatrix}$$

$$|\det D\phi(t, y)| = (1-t)^n.$$

If we were to calculate the integral in  $(t, z)$ -variable in the image set  $\phi([0, 1] \times M)$ ,

$$I = \iint_{\phi([0, 1] \times M)} Dv(z) \cdot (1-t)^{-1}(x-z) (1-t)^{-n} dt dz$$

The integral I coincide with  $R(x)$  by the change of variable formula. Here we also used the fact that

$$(x-z) = x - tx - (1-t)y = (1-t)(x-y).$$

Now, we compute

$$R(x) = \iint_{\phi([0, 1] \times M)} Dv(z) \cdot (x-z) (1-t)^{-n-1} dt dz$$

We make a few observations on the integral domain  $\phi([0, 1] \times M) = A \subset \mathbb{R}^{n+1}$ .

1. For each  $t$ , we let  $A_t$  be the set so that the slice  $\phi(\{t\} \times M) = \{t\} \times A_t$ .
2. For  $\sigma = [x_0 x_1 x_2 \cdots x_n]$ , we notice the slice at  $t$ , or the set  $A_t$ , is the point set of a simplex  $[(1-t)x_0 \ (1-t)x_1 \ \cdots \ (1-t)x_n]$  translated by  $tx$ , namely,

$$A_t = tx + (1-t)M.$$

3. With  $A_0 = M$  and  $A_1 = \{x\}$ .
4. Those amount to say that the integral domain is simply the point set of

$$[(0, x_0) \ (0, x_1) \ (0, x_2) \cdots (0, x_n) \ (1, x)].$$

5. On this  $(n+1)$ -simplex, we consider the iterated integral over the time  $t \in [0, t_*(z)]$  for each  $z \in A_0 = M$  below.

The integral over this  $(n+1)$ -simplex is performed by Fubini theorem

$$R(x) = \frac{1}{|M|} \int_{A_0} \int_0^{t_*(x,z)} (1-t)^{-n-1} Dv(z) \cdot (x-z) dt dz,$$

where for  $z \in A_0 = M$ ,  $t_*(x, z)$  is the intersection of  $A$  and the line  $[0, 1] \times \{z\}$ .

To compute  $t_*(x, z)$  for each  $z \in A_0 = M$ , we do the following:

1. From  $x$ , we draw the half ray passing  $z$ , until the line touches the boundary  $\partial M$ . The unique intersection point is denoted by  $c(x, z)$ .
2. We draw the right-angled-triangle, connecting three points of

$$(0, x), (1, x), (0, c(x, z)).$$

Then  $t_*(x, z)$  is the height on the hypotenuse from the point  $(0, z)$ . This gives rise to that

$$t_*(x, z) = \frac{|c(x, z) - z|}{|c(x, z) - x|} = 1 - \frac{|x - z|}{|c(x, z) - x|}.$$

Finally,

$$\begin{aligned}
R(x) &= \frac{1}{|M|} \int_M Dv(z) \cdot (x - z) \int_0^{t_*(x,z)} (1-t)^{-n-1} dt dz \\
&= \frac{1}{|M|} \int_M Dv(z) \cdot (x - z) \int_{1-t_*(x,z)}^1 t^{-n-1} dt dz \\
&= \frac{1}{n|M|} \int_M Dv(z) \cdot (x - z) \left( (1-t_*(x,z))^{-n} - 1 \right) dz \\
&= \frac{1}{n|M|} \int_M Dv(z) \cdot (x - z) \left( \frac{|c(x,z) - x|^n}{|x-z|^n} - 1 \right) dz \\
&= \frac{1}{n|M|} \int_M Dv(z) \cdot \frac{(x-z)}{|x-z|^n} \left( |c(x,z) - x|^n - |x-z|^n \right) dz.
\end{aligned}$$

□

For each  $M = \overline{\text{pts}}(\sigma)$ , the quantity

$$\frac{\text{diam}(M)^n}{|M|} =: \text{skew}(\sigma)$$

measures the quality of the simplex.

In practice, it is the duty of the mesh generator that for the given (good) domain, the skewness is kept bounded by a constant independent of piece.

**Theorem 2.** Let  $v \in C^\infty(M)$ .

$$|R(x)| \leq \text{skew}(\sigma) \frac{2}{n} \int_M |Dv(z)| |z-x|^{1-n} dz$$

Now, the (RHS) in the Theorem 2 is bounded by

$$\text{skew}(\sigma) \text{diam}(\sigma) \frac{2}{n} \int_M |Dv(z)| |z-x|^{-n} dz.$$

The Hardy-Littlewood-Sobolev inequality (after zero extension of  $Dv$  in  $\mathbb{R}^n \setminus M$ ) then gives the following:

**Theorem 3.** For any  $p \in (1, \infty)$ ,  $\exists C_{n,p} > 0$  such that for any  $v \in C^\infty(M)$

$$\left( \int_M |R(x)|^p dx \right)^{\frac{1}{p}} \leq C_{n,p} \text{skew}(\sigma) \text{diam}(\sigma) \left( \int_M |Dv(z)|^p dz \right)^{\frac{1}{p}}.$$

Theorem 1 ~ 3 will be repeated with the average of higher order taylor remainders.

The case  $p = 1$  can be proved to be true:

$$\begin{aligned} \int_M |R(x)| dx &\leq \text{skew}(\sigma) \frac{2}{n} \int_{x \in M} \int_{z \in M} |Dv(z)| |z - x|^{1-n} dz dx \\ &= \text{skew}(\sigma) \frac{2}{n} \int_{z \in M} |Dv(z)| \int_{x \in M} |z - x|^{1-n} dz dx \\ &\leq \text{skew}(\sigma) C_n \int_M |Dv(z)| \int_{r=0}^{\text{diam}(\sigma)} r^{1-n} r^{n-1} dr dx \\ &\leq C_n \text{skew}(\sigma) \text{diam}(\sigma) \int_M |Dv(z)|. \end{aligned}$$

The case  $p = \infty$  can be proved to be true:

$$\begin{aligned} |R(x)| &\leq \text{skew}(\sigma) \frac{2}{n} \int_M |Dv(z)| |z - x|^{1-n} dz \\ &= \text{skew}(\sigma) \frac{2}{n} \|Dv(z)\|_\infty \int_M |z - x|^{1-n} dz \\ &\leq \text{skew}(\sigma) C_n \|Dv(z)\|_\infty \int_{r=0}^{\text{diam}(\sigma)} r^{1-n} r^{n-1} dr \\ &\leq C_n \text{skew}(\sigma) \text{diam}(\sigma) \|Dv(z)\|_\infty. \end{aligned}$$

## Remark

We recall

**Theorem 1.** *Let  $v \in C^\infty(M)$ . Then*

$$v(x) - \frac{1}{|M|} \int_M v(y) dy = \frac{1}{n|M|} \int_M Dv(z) \cdot \frac{(x-z)}{|x-z|^n} (|c(x,z) - x|^n - |x-z|^n) dz,$$

where  $c(x,z)$  is the boundary point intersecting the half ray from  $x$  to  $z$ .

The remainder equality indeed can be interpreted as the *anti-derivative* formula, available in the  $n$ -simplex  $M$ .

1. If the domain is  $n$ -simplex and  $Dv$  is given,
2. (RHS), a formula involving the gradient  $Dv$  in  $M$ , equals to (LHS), the anti-derivative with the integrating constant  $\frac{1}{|M|} \int_M v(y) dy$ .

Compare this to the FTC in 1d for the domain  $[a, b]$ :

For a differentiable function  $v \in C^\infty([a, b])$ , the 1d fundamental theorem of calculus is the equality

$$\text{for } x \in [a, b] \quad v(x) - v(a) = \int_a^x v'(t) dt$$

and the equality can be read as two different ways:

(i) (LHS), the difference, equals to (RHS), the remainder.

(ii) (RHS), a formula involving  $v'$  in  $[a, b]$ , equals to (LHS), the anti-derivative with the integrating constant  $v(a)$ .

## Plan for the remainder for a projection to $\mathbb{P}_{m-1}$

We may think this case include the approximation for  $k = 0$  case. We consider a projection of  $C^\infty(M)$  to  $\mathbb{P}_{m-1}(M)$  and seek for the remainder formula.

Recall what we did for the  $\mathbb{P}_0(M)$  projection:

- We fixed  $x \in M$ , and let  $y \in M$ .
- Write FTC for each of  $y$ :

$$v(x) - v(y) = \int_0^1 Dv(tx + (1-t)y) \cdot (x - y) dt.$$

- Take average integral in  $y$ :

$$v(x) - \frac{1}{|M|} \int_M v(y) dy = \frac{1}{|M|} \int_M \int_0^1 Dv(tx + (1-t)y) \cdot (x - y) dt dy.$$

- We took the alternative form of (RHS) to conclude.

For  $\mathbb{P}_{m-1}(M)$  projection:

If  $t \mapsto f(t)$  is given by

$$f(t) = v(tx + (1-t)y), \quad t \in [0, 1]$$

we simply replace the second part equality by the 1d Taylor remainder equality:

$$f(1) - \sum_{k=0}^{m-1} \frac{1}{k!} f^{(k)}(0) = \frac{1}{(m-1)!} \int_0^1 (1-t)^{m-1} f^{(m)}(t) dt.$$

and repeat the rest of procedure.

Quite much preparation are needed for what we will do.

## Preparation 1: High order partial derivatives of $v$

### multi index $\alpha$

- Let us recall the multi index  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$  of nonnegative integers.
- Its use in high order derivative is the following. For a given multi index  $\alpha$ , we write

$$D^\alpha v = \partial_{x_1}^{\alpha_1} \partial_{x_2}^{\alpha_2} \cdots \partial_{x_{n-1}}^{\alpha_{n-1}} \partial_{x_n}^{\alpha_n} v.$$

- It is thus one of the  $|\alpha|$ -th order partial derivative of  $v$ .

### measuring the function by $\|\cdot\|_{m,p}$

- We use the  $(m, p)$ -norm for measuring the function  $v \in C^\infty(M)$ .
- We define the norm  $\|v\|_{m,p}$  for  $m \geq 0$  integer and  $p \in [1, \infty]$ :

$$\begin{aligned} \|v\|_{m,p} &= \left( \sum_{0 \leq |\alpha| \leq m} \int_M |D^\alpha v(x)|^p dx \right)^{\frac{1}{p}}, \quad p \in [1, \infty) \\ \|v\|_{m,\infty} &= \sum_{0 \leq |\alpha| \leq m} \sup_{x \in M} |D^\alpha v(x)|. \end{aligned}$$

## Preparation 2: Two projections

Let  $m$  be a positive integer.

Let  $\mathbb{P}_{m-1}(M)$  be the set of all polynomials in  $M$  of order at most  $m-1$ .

In our implementation, we will use an algorithm-specific linear projection operator

$$v \in C^\infty(M) \mapsto \pi(v) \in \mathbb{P}_{m-1}(M)$$

that is convenient in implementation.

To proceed, we define another linear operator

$$v \in C^\infty(M) \mapsto \mathcal{Q}(v) \in \mathbb{P}_{m-1}(M),$$

which is easier in making estimates.

- Let  $(C^\infty(M), \|\cdot\|_{m,p})$  and  $(\mathbb{P}_{m-1}(M), \|\cdot\|_{m,p})$  be equipped with  $(m,p)$ -norm.
- The operator norm of  $\pi$  and  $\mathcal{Q}$  are thus

$$|\pi| = \sup_{\|v\|_{m,p}=1} \|\pi v\|_{m,p}, \quad |\mathcal{Q}| = \sup_{\|v\|_{m,p}=1} \|\mathcal{Q}v\|_{m,p}.$$

- Suppose that I know the remainder estimate for the projector  $\mathcal{Q}$

$$\|v - \mathcal{Q}v\|_{m,p} \leq E.$$

- Then, because the projections are identity maps on  $\mathbb{P}_{m-1}(M)$ ,

$$\begin{aligned} \|v - \pi v\|_{m,p} &\leq \|v - \mathcal{Q}v\|_{m,p} + \|\mathcal{Q}v - \pi v\|_{m,p} = \|v - \mathcal{Q}v\|_{m,p} + \|\pi(v - \mathcal{Q}v)\|_{m,p} \\ &\leq \|v - \mathcal{Q}v\|_{m,p} + |\pi| \|v - \mathcal{Q}v\|_{m,p} = (1 + |\pi|)E. \end{aligned}$$

- Hence, although the inequality is not optimal, the remainder estimate for the projection  $\pi$  is

$$(1 + |\pi|)E$$

with possibly larger factor multiplied to  $E$ .

- Since computations we do here is quite heavy, we do the remainder estimate with calculation-friendly projection  $\mathcal{Q}$ .

## Preparation 3

To compute the expression  $f(1) - \sum_{k=0}^{m-1} \frac{1}{k!} f^{(k)}(0)$ , we compute the  $f^{(k)}(t)$  using the chain rule:

- Chain Rule gives that

$$\begin{aligned} f^{(k)}(t) &= \sum_{i_k=1}^n \sum_{i_{k-1}=1}^n \cdots \sum_{i_2=1}^n \sum_{i_1=1}^n (\partial_{x_{i_k}} \partial_{x_{i_{k-1}}} \cdots \partial_{x_{i_2}} \partial_{x_{i_1}} v)(tx + (1-t)y) \\ &\quad \times (x_{i_k} - y_{i_k})(x_{i_{k-1}} - y_{i_{k-1}}) \cdots (x_{i_2} - y_{i_2})(x_{i_1} - y_{i_1}) \end{aligned}$$

- We want to use the multi index  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n)$  of nonnegative integers.
- We consider the following counting problem: Given a multi index  $\alpha$  with  $|\alpha| = k$ , among all iterations of above  $k$  summations, how many times

$$(\partial_{x_{i_k}} \cdots \partial_{x_{i_2}} \partial_{x_{i_1}} v) (x_{i_k} - y_{i_k}) \cdots (x_{i_2} - y_{i_2})(x_{i_1} - y_{i_1})$$

correponds to

$$D^\alpha v (x - y)^\alpha \quad ?$$

- This number is

$$\begin{aligned} &\binom{k}{\alpha_1} \binom{k-\alpha_1}{\alpha_2} \binom{k-\alpha_1-\alpha_2}{\alpha_3} \cdots \binom{k-\alpha_1-\alpha_2-\cdots-\alpha_{n-1}}{\alpha_n} \\ &= \left( \frac{k!}{(k-\alpha_1)!\alpha_1!} \right) \left( \frac{(k-\alpha_1)!}{(k-\alpha_1-\alpha_2)!\alpha_2!} \right) \left( \frac{(k-\alpha_1-\alpha_2)!}{(k-\alpha_1-\alpha_2-\alpha_3)!\alpha_3!} \right) \cdots \left( \frac{(k-\alpha_1-\alpha_2-\cdots-\alpha_{n-1})!}{0!\alpha_n!} \right) \\ &= \frac{k!}{\alpha_1!\alpha_2!\cdots\alpha_n!} = \frac{k!}{\alpha!}. \end{aligned}$$

- Therefore, summation can be iterated multi index-wisely,

$$f^{(k)}(t) = \sum_{|\alpha|=k} \frac{k!}{\alpha!} D^\alpha v(tx + (1-t)y)(x - y)^\alpha.$$

- Finally, we have

$$f(1) - \sum_{k=0}^{m-1} \frac{1}{k!} f^{(k)}(0) = v(x) - \sum_{k=0}^{m-1} \frac{1}{k!} f^{(k)}(0) = v(x) - \sum_{0 \leq |\alpha| \leq m-1} \frac{1}{\alpha!} D^\alpha v(y)(x - y)^\alpha.$$

## Preparation 4

In the same way the taylor remainder  $\frac{1}{(m-1)!} \int_0^1 (1-t)^{m-1} f^{(m)}(t) dt$  can be expressed in the following way:

$$\begin{aligned} \frac{1}{(m-1)!} \int_0^1 (1-t)^{m-1} f^{(m)}(t) dt &= \frac{1}{(m-1)!} \int_0^1 (1-t)^{m-1} \sum_{|\alpha|=m} \frac{m!}{\alpha!} D^\alpha v(tx + (1-t)y)(x-y)^\alpha dt \\ &= \sum_{|\alpha|=m} \frac{m}{\alpha!} \int_0^1 (1-t)^{m-1} D^\alpha v(tx + (1-t)y)(x-y)^\alpha dt. \end{aligned}$$

Preparation 1 and 2 gives that the equality

$$v(x) - \sum_{0 \leq |\alpha| \leq m-1} \frac{1}{\alpha!} D^\alpha v(y)(x-y)^\alpha = \sum_{|\alpha|=m} \frac{m}{\alpha!} \int_0^1 (1-t)^{m-1} D^\alpha v(tx + (1-t)y)(x-y)^\alpha dt.$$

We do the integration in  $y$ :

$$v(x) - Q(x) = \sum_{|\alpha|=m} \frac{m}{\alpha!} \frac{1}{|M|} \int_M \int_0^1 (1-t)^{m-1} D^\alpha v(tx + (1-t)y)(x-y)^\alpha dt dy = R(x),$$

where

$$Q(x) = \frac{1}{|M|} \int_M \sum_{0 \leq |\alpha| \leq m-1} \frac{1}{\alpha!} D^\alpha v(y)(x-y)^\alpha.$$

**Proposition 2.**  $Q(x)$  is a polynomial of order at most  $m - 1$ .

*Proof.* This is because the expression

$$\int_M \sum_{k=0}^{m-1} h_k(y) x^k dy,$$

where  $h_k(y)$  is an integrable function of  $y$  in the compact set  $M$ , is a polynomial of order at most  $m - 1$ .  $\square$

We define the projection

$$v \mapsto \mathcal{Q}(v)$$

to be the  $Q(x)$  above.

**Theorem 4.** Let  $v \in C^\infty(M)$ . Then

$$v(x) - \mathcal{Q}(v)(x) = \sum_{|\alpha|=m} \frac{m}{n\alpha!|M|} \int_M D^\alpha v(z) \frac{(x-z)^\alpha}{|z-x|^n} \left( |c(x,z) - x|^n - |z-x|^n \right) dz,$$

where  $c(x,z)$  is the boundary point intersecting the half ray from  $x$  to  $z$ .

*Proof.*

$$v(x) - \mathcal{Q}(v)(x) = \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \int_M \int_0^1 (1-t)^{m-1} D^\alpha v(tx + (1-t)y) (x-y)^\alpha dt dy$$

The integrating domain is thus for  $(t,y) \in [0,1] \times M$ . We consider for  $(t,y) \in M$  the function

$$\phi : (t, y_1, y_2, \dots, y_n) \mapsto (t, z_1, z_2, \dots, z_n) = (t, tx + (1-t)y).$$

Then, we compute

$$D\phi(t, y) = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ x_1 - y_1 & 1-t & 0 & \cdots & 0 \\ x_2 - y_2 & 0 & 1-t & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_n - y_n & 0 & 0 & \cdots & 1-t \end{pmatrix}$$

$$|\det D\phi(t, y)| = (1-t)^n.$$

If we were to calculate the integral in  $(t, z)$ -variable in the image set  $\phi([0,1] \times M)$ ,

$$I = \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \iint_{\phi([0,1] \times M)} (1-t)^{m-1} D^\alpha v(z) (1-t)^{-m} (x-z)^\alpha (1-t)^{-n} dt dz$$

The integral I coincide with  $R(x)$  by the change of variable formula. Here we also used the fact that

$$(x-z) = x - tx - (1-t)y = (1-t)(x-y) \quad \text{and} \quad |\alpha| = m.$$

$$v(x) - \mathcal{Q}(v)(x) = \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \iint_{\phi([0,1] \times M)} (1-t)^{-n-1} D^\alpha v(z) (x-z)^\alpha dt dz.$$

As same as in the proof of Theorem 1,

$$t_*(x, z) = \frac{|c(x, z) - z|}{|c(x, z) - x|} = 1 - \frac{|z - x|}{|c(x, z) - x|}$$

and

$$\begin{aligned} v(x) - \mathcal{Q}(v)(x) &= \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \int_{A_0} \int_0^{t_*(x, z)} (1-t)^{-n-1} D^\alpha v(z) (x-z)^\alpha dt dz \\ &= \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \int_M D^\alpha v(z) (x-z)^\alpha \int_0^{t_*(x, z)} (1-t)^{-n-1} dt dz \\ &= \sum_{|\alpha|=m} \frac{m}{\alpha!|M|} \int_M D^\alpha v(z) (x-z)^\alpha \int_{1-t_*(x, z)}^1 t^{-n-1} dt dz \\ &= \sum_{|\alpha|=m} \frac{m}{n\alpha!|M|} \int_M D^\alpha v(z) (x-z)^\alpha \left( (1-t_*(x, z))^{-n} - 1 \right) dz \\ &= \sum_{|\alpha|=m} \frac{m}{n\alpha!|M|} \int_M D^\alpha v(z) (x-z)^\alpha \left( \frac{|c(x, z) - x|^n}{|z-x|^n} - 1 \right) dz \\ &= \sum_{|\alpha|=m} \frac{m}{n\alpha!|M|} \int_M D^\alpha v(z) \frac{(x-z)^\alpha}{|z-x|^n} \left( |c(x, z) - x|^n - |z-x|^n \right) dz \\ &= R(x). \end{aligned}$$

□

By having Theorem 4, we may say that FTC in 1d is generalized to higher order (Taylor, or our 1d remainder Theorem), and also to multi-dimensions, available for the  $n$ -simplex domain.

Now, the (RHS) in the Theorem 4 is bounded by

$$|v(x) - \mathcal{Q}(v)(x)| \leq \frac{\text{diam}(M)^n}{|M|} \sum_{|\alpha|=m} \frac{2m}{n\alpha!} \int_M |D^\alpha v(z)| |z-x|^{m-n} dz$$

**Theorem 5.** *Let  $v \in C^\infty(M)$ .*

$$|v(x) - \mathcal{Q}(v)(x)| \leq \text{skew}(\sigma) \sum_{|\alpha|=m} \frac{2m}{n\alpha!} \int_M |D^\alpha v(z)| |z-x|^{m-n} dz$$

This time, we give the direct proof of the Theorem 6 for  $p \in (1, \infty)$ .

**Theorem 6.** *For any  $p \in (1, \infty)$ ,  $\exists C_{n,m,p} > 0$  such that for any  $v \in C^\infty(M)$*

$$\left( \int_M |v(x) - \mathcal{Q}(v)(x)|^p dx \right)^{\frac{1}{p}} \leq C_{n,m,p} \text{skew}(\sigma) \text{diam}(\sigma)^m \left( \sum_{|\alpha|=m} \int_M |D^\alpha v(z)|^p dz \right)^{\frac{1}{p}}.$$

This is a high order generalization of the Poincare inequality, available in  $M$ .

*Proof.* We estimate  $L^p$  norm of

$$J(x) = \int_M |D^\alpha v(z)| |z - x|^{m-n} dz \quad \text{for a fixed } \alpha$$

For given  $p \in (1, \infty)$ , we write  $\frac{p}{p-1} = p' \in (1, \infty)$ .  $\frac{1}{p} + \frac{p-1}{p} = 1$ .

$$\begin{aligned} \int_M |J(x)|^p dx &= \int_{x \in M} \left[ \int_{z \in M} |D^\alpha v(z)| |z - x|^{m-n} dz \right]^p dx \\ &= \int_{x \in M} \left[ \int_{z \in M} |D^\alpha v(z)| |z - x|^{\frac{m-n}{p}} |z - x|^{\frac{(m-n)(p-1)}{p}} dz \right]^p dx \end{aligned}$$

(Holder inequality in  $z$  integral)

$$\begin{aligned} &\leq \int_{x \in M} \left[ \left( \int_{z \in M} |D^\alpha v(z)|^p |z - x|^{m-n} dz \right)^{\frac{1}{p}} \left( \int_{z \in M} |z - x|^{m-n} dz \right)^{\frac{p-1}{p}} \right]^p dx \\ &= \int_{x \in M} \left( \int_{z \in M} |D^\alpha v(z)|^p |z - x|^{m-n} dz \right) \left( \int_{z \in M} |z - x|^{m-n} dz \right)^{p-1} dx. \end{aligned}$$

Now,

$$\begin{aligned} \left( \int_{z \in M} |z - x|^{m-n} dz \right)^{p-1} &= \left( \int_{w \in M'} |w|^{m-n} dw \right)^{p-1}, \quad (\text{translate } M \text{ so that } x \text{ is moved to origin.}) \\ &\leq \left( \omega_n \int_0^{\text{diam}(\sigma)} r^{m-n} r^{n-1} dr \right)^{p-1} \quad (M' \subset B_{\text{diam}(\sigma)}(0)) \\ &= E_{m,n,p} \text{diam}(\sigma)^{m(p-1)}. \end{aligned}$$

Therefore,

$$\begin{aligned} \int_M |J(x)|^p dx &\leq E_{m,n,p} \text{diam}(\sigma)^{m(p-1)} \int_{x \in M} \int_{z \in M} |D^\alpha v(z)|^p |z - x|^{m-n} dz dx \\ &= E_{m,n,p} \text{diam}(\sigma)^{m(p-1)} \int_{z \in M} |D^\alpha v(z)|^p \int_{x \in M} |z - x|^{m-n} dx dz \\ &= E_{m,n,p} \text{diam}(\sigma)^{m(p-1)} \int_{z \in M} |D^\alpha v(z)|^p \int_{w \in M''} |w|^{m-n} dw dz \\ &\leq D_{m,n,p} \text{diam}(\sigma)^{m(p-1)+m} \int_{z \in M} |D^\alpha v(z)|^p dz \\ &= D_{m,n,p} \text{diam}(\sigma)^{mp} \int_{z \in M} |D^\alpha v(z)|^p dz. \end{aligned}$$

Taking  $\frac{1}{p}$ -th power,

$$\left( \int_M |J(x)|^p dx \right)^{\frac{1}{p}} \leq K_{m,n,p} \text{diam}(\sigma)^m \left( \int_{z \in M} |D^\alpha v(z)|^p dz \right)^{\frac{1}{p}}.$$

□

We leave the proof for  $p = 1$  and  $p = \infty$  cases omitted.

## Derivative estimates of the Remainder

1. We have estimated the  $L^p$  norm of the remainder  $x \mapsto v(x) - \mathcal{Q}v(x)$ .
2. We consider  $k$ -th order ( $0 \leq k \leq m$ ) partial derivative of the remainder

$$D^\beta(v - \mathcal{Q}v)$$

and give an estimate.

**Theorem 7.** For any  $p \in (1, \infty)$ ,  $\exists C_{n,m,k,p} > 0$  such that for any  $v \in C^\infty(M)$  and  $\beta$  with  $|\beta| = k$  and  $0 \leq k \leq m$ ,

$$\left( \int_M |D^\beta(v - \mathcal{Q}v)|^p dx \right)^{\frac{1}{p}} \leq C_{m,n,k,p} \text{skew}(\sigma) \text{diam}(\sigma)^{m-k} \left( \sum_{|\alpha|=m} \int_M |D^\alpha v(x)|^p dx \right)^{\frac{1}{p}}.$$

**Proposition 8.** Let  $x, y \in M$  and  $v \in C^\infty(M)$ . Let the polynomial

$$x \mapsto q(x ; v, y, m-1) = \sum_{0 \leq |\alpha| \leq m-1} \frac{1}{\alpha!} D^\alpha v(y) (x-y)^\alpha.$$

Let  $\beta$  be a multi index with  $|\beta| = k$  and  $0 \leq k \leq m$ . Then

$$D^\beta q(x ; v, y, m-1) = q(x ; D^\beta v, y, m-k-1).$$

*Proof.* .

- Because the partial derivative

$$D^\beta = \partial_{x_1}^{\beta_1} \partial_{x_2}^{\beta_2} \cdots \partial_{x_n}^{\beta_n}$$

is with respect to  $x$  variable, it only hits the  $x$ -dependent polynomial part of

$$\sum_{0 \leq |\alpha| \leq m-1} \frac{1}{\alpha!} D^\alpha v(y) (x-y)^\alpha$$

- Now, for given this  $\beta$  and for some  $\alpha$  in the iteration, suppose it happens

$$\beta_j > \alpha_j \quad \text{for some } j \in \{1, 2, \dots, n\}.$$

Then, the application of  $D^\beta(x-y)^\alpha$  for this  $\alpha$  will be simply zero.

- Let us write  $\alpha \geq \beta$  in case  $\alpha_j \geq \beta_j$  for every  $j = 1, 2, \dots, n$ , a partial order on multi indices. With this notation,

$$D^\beta q(x ; v, y, m-1) = \sum_{0 \leq |\alpha| \leq m-1, \alpha \geq \beta} \frac{1}{\alpha!} D^\alpha v(y) D^\beta(x-y)^\alpha.$$

- The summation is then can be re-written in the way

$$\begin{aligned} D^\beta q(x ; v, y, m-1) &= \sum_{0 \leq |\gamma| \leq m-k-1, \alpha=\beta+\gamma} \frac{1}{\alpha!} D^\alpha v(y) D^\beta(x-y)^\alpha \\ &= \sum_{0 \leq |\gamma| \leq m-k-1} \frac{1}{(\beta+\gamma)!} D^{\beta+\gamma} v(y) D^\beta(x-y)^{\beta+\gamma} \end{aligned}$$

- The computation of  $D^\beta(x - y)^{\beta + \gamma}$  is

$$\begin{aligned} & (\partial_{x_1}^{\beta_1} \partial_{x_2}^{\beta_2} \cdots \partial_{x_n}^{\beta_n}) (x_1 - y_1)^{\beta_1 + \gamma_1} (x_2 - y_2)^{\beta_2 + \gamma_2} \cdots (x_n - y_n)^{\beta_n + \gamma_n} \\ &= \frac{(\beta_1 + \gamma_1)!}{(\gamma_1)!} \frac{(\beta_2 + \gamma_2)!}{(\gamma_2)!} \cdots \frac{(\beta_n + \gamma_n)!}{(\gamma_n)!} (x_1 - y_1)^{\gamma_1} (x_2 - y_2)^{\gamma_2} \cdots (x_n - y_n)^{\gamma_n} \\ &= \frac{(\beta + \gamma)!}{\gamma!} (x - y)^\gamma. \end{aligned}$$

- Substituting this into the  $D^\beta q(x ; v, y, m - 1)$ , we obtain

$$\begin{aligned} D^\beta q(x ; v, y, m - 1) &= \sum_{0 \leq |\gamma| \leq m-k-1} \frac{1}{\gamma!} D^\gamma D^\beta v(y) (x - y)^\gamma \\ &= q(x ; D^\beta v, y, m - k - 1). \end{aligned}$$

□

**Proposition 9.** Let  $x, y \in M$  and suppose  $D^\alpha v$  is integrable for every  $0 \leq |\alpha| \leq m - 1$ . Then

$$D^\beta \int_M q(x ; v, y, m - 1) dy = \int_M D^\beta q(x ; v, y, m - 1) dy.$$

*Proof.* We leave this as an exercise. Use the Lebesgue Dominated Convergence Theorem.

□

Now, we prove the Theorem 7.

*proof of Theorem 7.*

- By Proposition 8 and 9, we have

$$D^\beta (v(x) - \mathcal{Q}(x ; v, m - 1)) = D^\beta v(x) - \mathcal{Q}(x ; D^\beta v, m - k - 1).$$

- Applying Theorem 4 on  $D^\beta v$ , we obtain

$$\begin{aligned} & \left( \int_M |D^\beta v(x) - \mathcal{Q}(x ; D^\beta v, m - k - 1)|^p dx \right)^{\frac{1}{p}} \\ & \leq C_{m-k,n,p} \text{skew}(\sigma) \text{diam}(\sigma)^{m-k} \left( \sum_{|\gamma|=m-k} \int_M |D^{\beta+\gamma} v(x)|^p dx \right)^{\frac{1}{p}} \\ & \leq C_{m-k,n,p} \text{skew}(\sigma) \text{diam}(\sigma)^{m-k} \left( \sum_{|\alpha|=m} \int_M |D^\alpha v(x)|^p dx \right)^{\frac{1}{p}} \end{aligned}$$

□

*Remark 6.1.*

1. In particular, if  $k = m$ , simply

$$D^\beta (v(x) - \mathcal{Q}(x ; v, m - 1)) = D^\beta v$$

because the polynomial is of order at most  $m - 1$ .

2. It is noteworthy that, the inequality we derived is between the semi-norms,

$$|v|_{k,p} = \left( \sum_{|\alpha|=k} |D^\alpha v(x)|^p dx \right)^{\frac{1}{p}},$$

not the norm

$$\|v\|_{k,p} = \left( \sum_{0 \leq |\alpha| \leq k} |D^\alpha v(x)|^p dx \right)^{\frac{1}{p}}.$$

## Chapter 7

### All combined: Example by $(2, p)$ -norm

Finally, we estimate the remainder in global domain by  $(2, p)$ -norm.

We formulate a list of assumptions to get estimate correct.

1. The global domain and triangulation by  $\mathcal{S}$

$$\bar{\Omega} = |\mathcal{S}| = \bigcup_{\sigma \in S_n} \sigma$$

are such that  $\text{skew}(\sigma)$  and  $\text{diam}(\sigma)$  are bounded above uniformly by constants

$$\text{skew}(\sigma) \leq \text{skew}_{\mathcal{S}} < \infty, \quad \text{diam}(\sigma) \leq h \quad \text{for every } \sigma \in \mathcal{S}_n.$$

2. In the example, we assume  $2p > n$ . For instance,  $(m, p) = (2, 2)$  and  $n = 2$  or  $3$ .

## Example of operator norm $|\pi|$ in $M$

- We recall that the projection  $\pi : C^\infty(M) \rightarrow \mathbb{P}_1(M)$  is given by

$$\pi v(x) = \sum_{e=0}^n v(x_e) \theta_e(x).$$

- Thus,

$$\begin{aligned} |\pi| &= \sup_{\|v\|_{2,p}=1} \left\| \sum_{e=0}^n v(x_e) \theta_e(x) \right\|_{2,p} \\ &\leq \sup_{\|v\|_{2,p}=1} \|v\|_\infty \sum_{e=0}^n \|\theta_e(x)\|_{2,p} \end{aligned}$$

- We use that if  $2p > n$ , then in the  $n$ -simplex  $M$ ,

$$\|v\|_\infty \leq C_{sobolev} \|v\|_{2,p}.$$

$$\begin{aligned} |\pi| &\leq \sup_{\|v\|_{2,p}=1} C_{sobolev} \|v\|_{2,p} \sum_{e=0}^n \|\theta_e(x)\|_{2,p} \\ &= C_{sobolev} \sum_{e=0}^n \|\theta_e(x)\|_{2,p}. \end{aligned}$$

- The operator norm  $|\pi|$  also can be made bounded above uniformly under the assumption we listed.

## Global remainder estimates

For given  $v \in C^\infty(\Omega)$ , (or  $v$  is in the Sobolev space  $W^{2,p}(\Omega)$  by density argument), we conduct the local  $\pi$  projection to  $\mathbb{P}_1(M)$ , and assemble to obtain the pp function  $\phi$ , a piecewise linear approximation of  $v$  in  $\Omega$ .

Then, we have

$$\begin{aligned}\|v - \phi\|_{L^p(\Omega)} &\leq C \text{skew}_{\mathcal{S}} h^2 \|v\|_{2,p,\Omega} \\ \|Dv - D\phi\|_{L^p(\Omega)} &\leq C \text{skew}_{\mathcal{S}} h \|v\|_{2,p,\Omega}.\end{aligned}$$



## Chapter 8

# Approximation of solution of PDE in $\Pi_0(\Omega)$

In this chapter, we consider one of the simplest approximation only using  $\Pi_0(\Omega)$ , consisting of continuous and piecewise linear scalar functions.

The (nominal) PDE problem we consider is

$$\begin{aligned}\operatorname{div}(-K(x)\nabla p) &= \varphi, \quad x \in \Omega, \\ -K(x)\nabla p \cdot n &= \psi, \quad x \in \partial\Omega.\end{aligned}$$

To make our statements mathematically correct, we assume the followings.

We state assumptions labelled by (A):

1.  $\Omega \subset \mathbb{R}^3$  is a bounded simply connected polyhedral domain, i.e., a finite union of polyhedral cells.

2.  $\varphi \in L^2(\Omega)$ , i.e.,

$$\left( \int_{\Omega} |\varphi|^2 \right)^{\frac{1}{2}} < \infty.$$

3.  $\psi \in L^2(\partial\Omega)$ , i.e.,

$$\left( \int_{\partial\Omega} |\psi|^2 \right)^{\frac{1}{2}} < \infty.$$

4. The data  $\varphi$  and  $\psi$  are constrained by

$$\int_{\Omega} \varphi = \int_{\partial\Omega} \psi.$$

5. The coefficient  $K$  is a bounded symmetric-positive-definite matrix-valued function and satisfies the uniform ellipticity condition:

$$\exists 0 < \lambda \leq \Lambda, \quad \text{for every } x \in \Omega \text{ and every } w \in \mathbb{R}^3, \quad \lambda|w|^2 \leq w^T K(x) w \leq \Lambda|w|^2.$$

6. For each  $0 < h < 1$ , there is a triangulation of  $\Omega$  by  $\mathcal{S}_h$ , with

$$\sup_{\sigma \in \mathcal{S}_h} \operatorname{diam}(\sigma) \leq h.$$

## Rootfinding problem formulation

$X$ ,  $Y$ , and  $F : X \rightarrow Y$ .

Since we are solving *equation*, i.e., we consider a rootfinding problem for an equation

$$F(x) = y, \quad \text{for given } y,$$

We want to define what is  $F : X \rightarrow Y$ , what is  $X$ , what is  $Y$ .

Although notions here belong to those in functional analysis, we will simply introduce what are needed.

1. We let  $X = W^{1,2}(\Omega)$ , the set obtained by completing the space  $C^\infty(\Omega)$  by the  $(1,p)$ -norm with  $p = 2$ .
2. If you are not comfortable with this, you may be also safe to consider only the case  $K$  is smooth and differentiable functions.
3. For a given vector space  $X$ , we consider the set of its *adjoint elements*,

$$\left\{ \ell : X \rightarrow \mathbb{R} \mid \ell \text{ is linear and continuous} \right\}.$$

4. This set is denoted by  $(W^{1,2}(\Omega))^*$ , and we let  $Y = (W^{1,2}(\Omega))^*$ .

We have let

$$X = W^{1,2}(\Omega), \quad Y = (W^{1,2}(\Omega))^*$$

and now we define what is  $F : X \rightarrow Y$ .

We define  $F : X \rightarrow Y$  to be

$$p \mapsto \ell_p, \quad \ell_p(v) = \int_{\Omega} K(x) \nabla p \cdot \nabla v, \quad \forall v \in X.$$

Although the sets we involved are of infinite dimensional, this is nothing else than a rootfinding problem

Find  $p \in X$  such that  $F(p) = y$ , for given  $y \in Y$ .

**Formulation of  $y \in Y$**

Now, we consider the right-hand-side  $y \in Y$ , so that the rootfinding problem is that of the nominal PDE in a certain sense.

We define the right-hand-side  $\ell \in Y$  for the PDE by

$$\ell(v) = \int_{\Omega} \varphi v - \int_{\partial\Omega} \psi v \in \mathbb{R}, \quad \forall v \in X.$$

**Definition 1.** Let the assumptions in (A) hold. We define the problem  $\mathcal{P}$ : For given  $\varphi$  and  $\psi$ , let  $\ell \in Y$  as defined above. Find  $p \in X$  such that

$$F(p) = \ell.$$

## Modified problem $\mathcal{P}_h$

The problem  $\mathcal{P}$  is difficult to solve. In other words, we do not know what is its inverse

$$F^{-1} : Y \rightarrow X.$$

Hence, we solve other problem,  $F_h(x) = y$ .

For a triangulation by  $\mathcal{S}_h$ ,

1. We let

$$X_h = \Pi_0(\Omega) \subset X.$$

This is  $n_0$  dimensional vector space  $\simeq \mathbb{R}^{n_0}$ .

2. We define  $F_h : X_h \rightarrow Y$ :

$$p_h \mapsto \ell_{p_h}, \quad \ell_{p_h}(v) = \int_{\Omega} K(x) \nabla p_h \cdot \nabla v, \quad \forall v \in X.$$

**Definition 2.** Let the assumptions in (A) hold. We define the problem  $\mathcal{P}_h$ : For given  $\varphi$  and  $\psi$ , let  $\ell \in Y$  as defined above. Find  $p_h \in X_h$  such that

$$F_h(p_h)|_{X_h} = \ell|_{X_h}.$$

In other words, find  $p_h \in X_h$  such that

$$\text{when } v_h \in X_h, \quad \int_{\Omega} K(x) \nabla p_h \cdot \nabla v_h = \int_{\Omega} \varphi v_h - \int_{\partial\Omega} \psi v_h.$$

## Approximate solution operator $\mathcal{R}_h$

The problem  $\mathcal{P}_h$  turns out to be the one solving a system of  $n_0$  linear equations that is going to be written in the form of a matrix equation  $Ax = b$ .

1. Since  $p_h$  we look for is in  $X_h$ , it is in a span of basis of  $X_h = \Pi_0(\Omega)$ :

$$p_h = \sum_{i=1}^{n_0} c_i \phi_i(x).$$

So  $(c_i)_{i=0}^{n_0}$  are unknowns.

2. For the equality to hold for each test function  $v_h \in X_h$ , we test each basis function  $\phi_i$ . The equality in Definition of  $\mathcal{P}_h$  holds if and only if

$$\text{for every } \phi_i, \text{ a basis of } X_h, \quad \int_{\Omega} K(x) \nabla p_h \cdot \nabla \phi_i = \int_{\Omega} \varphi \phi_i - \int_{\partial\Omega} \psi \phi_i.$$

3. Since we have in total  $n_0$  tests, we collect  $b \in \mathbb{R}^{n_0 \times 1}$  column as (RHS) that is

$$b_i = \int_{\Omega} \varphi \phi_i - \int_{\partial\Omega} \psi \phi_i, \quad i = 1, \dots, n_0.$$

(This integration is in fact done by quadrature.)

4. For (LHS), we notice that for each  $i$ ,

$$\begin{aligned} \int_{\Omega} K(x) \nabla p_h \cdot \nabla \phi_i &= \sum_{i'=1}^{n_0} \int_{\Omega} K(x) c_{i'} \nabla \phi_{i'} \cdot \nabla \phi_i \\ &= \sum_{i'=1}^{n_0} \left( \int_{\Omega} K(x) \nabla \phi_{i'} \cdot \nabla \phi_i \right) c_{i'} \\ &= \sum_{i'=1}^{n_0} A_{ii'} c'_{i'}, \quad A_{ii'} = \int_{\Omega} K(x) \nabla \phi_i \cdot \nabla \phi_{i'} \end{aligned}$$

5. Thus, we have a matrix equation

$$Ac = b.$$

We can make the following observations.

- The matrix  $A \in \mathbb{R}^{n_0 \times n_0}$  is symmetric.
- The matrix  $A$  is semi-positive definite: Let  $c \in \mathbb{R}^{n_0}$ .

$$\begin{aligned} c^T A c &= \int_{\Omega} \left( K(x) \left( \sum_{i=1}^{n_0} c_i \nabla \phi_i \right) \right)^T \left( \sum_{i'=1}^{n_0} c_{i'} \nabla \phi_{i'} \right) \\ &\geq \int_{\Omega} \lambda \left( \sum_{i=1}^{n_0} c_i \nabla \phi_i \right)^T \left( \sum_{i'=1}^{n_0} c_{i'} \nabla \phi_{i'} \right) \\ &= \lambda \left\| \sum_{i=1}^{n_0} c_i \nabla \phi_i \right\|_{L^2(\Omega)}^2 \geq 0. \end{aligned}$$

- In fact,  $A$  has precisely rank  $n_0 - 1$ , and this has to be so, naturally.
- $\text{Ker}(A)^\perp$  perspective:

1. We note that  $c^T A c = 0$  if and only if

$$\sum_{i=1}^{n_0} c_i \nabla \phi_i = \nabla \left( \sum_{i=1}^{n_0} c_i \phi_i \right) = 0 \iff \sum_{i=1}^{n_0} c_i \phi_i = \text{const.}$$

2. We look for  $\bar{c} \neq 0$  such that  $\sum_{i=1}^{n_0} \bar{c}_i \phi_i = 1$ . Since  $1 \in \Pi_0(\Omega)$  and  $(\phi_i)_{i=1}^{n_0}$  are basis of  $\Pi_0$ , there is unique such  $\bar{c}$ .

(Of course, if the *const.* is not 1, we have again unique *const.* multiple of  $\bar{c}$ .)

3. In fact, we know that this is the case  $\bar{c} \parallel (1, 1, 1, \dots, 1)^T$ .

4. Hence,  $\text{Ker}(A)$  is precisely one dimensional,

$$\text{Ker}(A)^\perp = \text{span} \langle (1, 1, \dots, 1)^T \rangle^\perp.$$

- $\text{Ran}(A)$  perspective:

1. Since  $\text{rank}(A) = n_0 - 1$ , not every  $b \in \mathbb{R}^{n_0 \times 1}$  will admit the solution.

2. Since  $A$  is symmetric, its  $\text{Ran}(A) = \text{Ker}(A)^\perp$ .

3. Then, is the (RHS)  $b \in \mathbb{R}^{n_0-1}$  we have defined in  $\text{Ran}(A)$ ?

Indeed,

$$\begin{aligned} (1, 1, \dots, 1) \cdot (b_1, b_2, \dots, b_n) &= \sum_i b_i = \sum_i \int_{\Omega} \varphi \phi_i - \int_{\partial\Omega} \psi \phi_i \\ &= \int_{\Omega} \varphi - \int_{\partial\Omega} \psi = 0. \end{aligned}$$

by the constraint in  $(A)$ .

**Definition 3.** *The approximate solution operator  $\mathcal{R}_h$  is defined by*

$$\mathcal{R}_h(\ell) = \sum_{i=1}^{n_0} c_i \phi_i(x), \quad \text{where}$$

*c is given by a linear solver of the system specified by*

$$A_{ii'} = \int_{\Omega} K(x) \nabla \phi_i \cdot \nabla \phi_{i'}, \quad b_i = \int_{\Omega} \varphi \phi_i - \int_{\partial\Omega} \psi \phi_i, \quad i = 1, \dots, n_0.$$

*(More precisely, between the  $\text{Ker}(A)^\perp$ ).*

Considering the discussion in the last class, restrictions of the sets  $X$  and  $Y$  are considered:

$$\begin{aligned}\hat{X} &= \left\{ p \in X \mid \int_{\Omega} p = 0 \right\}, \\ \hat{Y} &= \left\{ \ell \in Y \mid \ell(\mathbf{1}) = 0 \right\}.\end{aligned}$$

Also, we let  $\hat{X}_h = X_h \cap \hat{X}$ .

For today's discussion, we assume in addition (B):

1. We assume

$$\text{skew}(\mathcal{S}_h) \leq M$$

independently for  $h$ .

2. We assume  $\Omega$  is a domain satisfying the following: There exists  $C > 0$  such that

$$\left( \int_{\Omega} |p - (p)|^2 \right)^{\frac{1}{2}} \leq C \left( \int_{\Omega} |\nabla p|^2 \right)^{\frac{1}{2}}, \quad (p) = \frac{1}{|\Omega|} \int_{\Omega} p.$$

In particular, we have the following calculations.

Assume (A) and (B).

1. Let us write

$$\langle \nabla p, \nabla q \rangle_K = \int_{\Omega} K(x) \nabla p \cdot \nabla q, \quad \|\nabla p\|_{L_K^2} = |\langle \nabla p, \nabla p \rangle_K|^{\frac{1}{2}}.$$

2. If  $p \in \hat{X}$ , because its integral vanishes in  $\Omega$ ,

$$\left( \int_{\Omega} |p|^2 \right)^{\frac{1}{2}} \leq C \left( \int_{\Omega} |\nabla p|^2 \right)^{\frac{1}{2}} \leq \frac{C}{\lambda} \left( \int_{\Omega} \lambda |\nabla p|^2 \right)^{\frac{1}{2}} \leq \frac{C}{\lambda} \left( \int_{\Omega} K(x) \nabla p \cdot \nabla p \right)^{\frac{1}{2}} = \frac{C}{\lambda} \|\nabla p\|_{L_K^2}.$$

3. We can use  $\|\nabla p\|_{L_K^2}$  to be a norm on  $\hat{X}$ .



## Chapter 9

# Consistency, Convergence, and Errors

Here, we run our error analysis framework for  $F(x) = y$ . We will re-visit this.

We will speak of

1. Consistency of  $(F_h)_{h>0}$ .
2. Testing error at  $x \in X$ .
3. Consistency of  $(R_h)_{h>0}$ .
4. Testing error at  $y \in Y$ .
5. Solution error and convergence of algorithm.

## Consistency of $(F_h)_{h>0}$

We have defined the modified problem with respect to  $F_h : \hat{X}_h \rightarrow Y$ .

$$p \mapsto \ell_p, \quad \ell_p(v) = \int_{\Omega} K(x) \nabla p \cdot \nabla v, \quad \forall v \in X.$$

Or, alternatively we can make  $F_h$  be defined on all of  $\hat{X}$  by passing to projection  $I^h$

$$p \mapsto I^h p \mapsto \ell_{I^h p}$$

**Definition 1.** We say  $(F_h)_{h>0}$  is consistent at  $p \in \hat{X}$  if

$$\|F(p) - F_h(p)\|_{\hat{Y}} \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

We call  $\|F(p) - F_h(p)\|_{\hat{Y}}$  the testing error at  $p \in X$ .

Consistency of  $(F_h)$  will be used in a priori error analysis, where actually we do not know what the true solution  $p$  is.

To check if  $(F_h)$  is consistent or not at  $p$ , we compute its testing error.

1. First of all, for given  $\ell \in \hat{Y}$ , we use the norm

$$\|\ell\|_{\hat{Y}} = \sup_{v \in \hat{X}, \|v\|_{\hat{X}}=1} |\ell(v)|.$$

2. Testing error at  $p \in \hat{X}$ :

$$\begin{aligned} \|F(p) - F_h(p)\|_Y &= \|\ell_p - \ell_{I^h p}\|_Y, \\ &= \sup_{v \in \hat{X}, \|v\|_{\hat{X}}=1} \left| \int_{\Omega} K(x) \nabla p \cdot \nabla v - K(x) \nabla(I^h p) \cdot \nabla v \right| \\ &\leq \|\nabla(p - I^h p)\|_{L_K^2}. \end{aligned}$$

- As a matter of fact, in case  $p \in W^{1,2}(\Omega)$ ,

$\|\nabla(p - I^h p)\|_{L_K^2}$  does converges to 0 as  $h \rightarrow 0$ , (Lebesgue Differentiation Theorem)

but this lacks the explicit dependency on  $h$ .

- For a while, it will suffice to know that the testing error at  $p$

$$\|F(p) - F_h(p)\|_Y \leq \|\nabla(p - I^h p)\|_{L_K^2} = o(1) \quad \text{as } h \rightarrow 0.$$

## Consistency of approximate solution operator $(R_h)_{h>0}$

**Definition 2.** We say  $(R_h)_{h>0}$  is consistent at  $\ell \in \hat{Y}$  if,

1.  $R_h(\ell) \in \hat{X}$ .

2. Also

$$\|F \circ R_h(\ell) - \ell\|_{\hat{Y}} \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

If 1st item is the case, we call  $\|F \circ R_h(\ell) - \ell\|_{\hat{Y}}$  the testing error at  $\ell \in \hat{Y}$ .

To check if  $(R_h)$  is consistent or not, let us compute its testing error.

1. Let  $p_h = R_h(\ell) \in X_h$ .

In other words, it solves the equation:

$$\int_{\Omega} K(x) \nabla p_h \cdot \nabla v_h = \ell(v_h), \quad \text{when } v_h \in \hat{X}_h.$$

2. We have

$$F \circ \mathcal{R}_h(\ell) : v \mapsto \int_{\Omega} K(x) \nabla p_h \cdot \nabla v.$$

3. Testing error at  $\ell \in \hat{Y}$ :

$$\begin{aligned} \|F \circ \mathcal{R}_h(\ell) - \ell\|_{\hat{Y}} &= \sup_{v \in \hat{X}, \|v\|_{\hat{X}}=1} \left| \int_{\Omega} K(x) \nabla p_h \cdot \nabla (v - I^h v) + \int_{\Omega} K(x) \nabla p_h \cdot \nabla I^h v - \ell(v - I^h v) - \ell(I^h v) \right| \\ &= \sup_{v \in \hat{X}, \|v\|_{\hat{X}}=1} \left| \int_{\Omega} K(x) \nabla p_h \cdot \nabla (v - I^h v) - \ell(v - I^h v) \right| \\ &\leq \|p_h\|_{\hat{X}} \|(v - I^h v)\|_{\hat{X}} + \|\ell\|_{\hat{Y}} \|v - I^h v\|_{\hat{X}} = (\|p_h\|_{\hat{X}} + \|\ell\|_{\hat{Y}}) \|v - I^h v\|_{\hat{X}}. \end{aligned}$$

Consistency of  $(\mathcal{R}_h)$  is used in a posteriori error analysis, where once the algorithm is run, we know  $(\|p_h\|_{\hat{X}} + \|\ell\|_{\hat{Y}})$ .

- Again,

$$\|v - I^h v\|_{\hat{X}} = \|\nabla(v - I^h v)\|_{L_K^2} = o(1) \quad \text{as } h \rightarrow 0.$$



# Chapter 10

## Error Analysis for $F(x) = y$ .

Many problems can be formulated as equation solving problems, and equation solving problem can be analyzed, by the least, in a certain systematic way.

To take an example, we consider the problem of orthogonalization of an  $n \times n$  square matrix:

For given  $A \in R^{n \times n}$ , find the orthogonal matrix  $O \in R^{n \times n}$ , and the upper triangular matrix  $R \in R^{n \times n}$  such that  $A = OR$ .

1. We let  $X = X_1 \times X_2$ ,  $X_1$  the set of all  $n \times n$  orthogonal matrices,  $X_2$  the set of all  $n \times n$  upper triangular matrices.
2. We let  $Y$  be the set of all  $n \times n$  matrices.
3. We let  $F(O, R) = OR$ .

and we solve  $F(O, R) = A$  for given  $A$ .

### $X$ , $Y$ , and $F : X \rightarrow Y$

- $X$  and  $Y$  are metric spaces, and  $F$  is assumed to be continuous and surjective.
- The testing whether  $F(x)$  equals to  $y$  has to be checkable accurate enough.
- In considering the problem to solve  $F(x) = y$ , the existence theory is from upstream, i.e., somebody else other than numerical analyst has done it.

Below, we use notation  $|x - x'|$  for the distance, adapted to the normed space, but you can adjust this to the distance function in metric space.

### (Parametrized) Modified problem $F_h : X \rightarrow Y$

- The objective is to solve the equation  $F(x) = y$ , by taking the limit of a parametrized solvable problems by

$$F_h : X \rightarrow Y.$$

- Adapted to our problem between infinite dimensional spaces, let us be more specific:

- Let  $X_h$  be a subset of  $X$  such that we have a projection  $I_0^h : X \rightarrow X_h$ ,
- Let  $Y_h$  be a subset of  $Y$  such that we have a projection  $I_1^h : Y \rightarrow Y_h$ ,
- The modified problem we consider is between  $X_h$  and  $Y_h$ :

$$\bar{F}_h : X_h \rightarrow Y_h.$$

- Using projections, and  $\bar{F}_h$ , we also write

$$F_h = \bar{F}_h \circ I_0^h : X \rightarrow Y_h.$$

### (Parametrized) Approximate Solver $R_h : Y \rightarrow X$

- The solver of the solvable problem  $\bar{F}_h(x_h) = y_h$  is denoted by  $\bar{R}_h : Y_h \rightarrow X_h$ ,

$$\bar{F}_h(\bar{R}_h(y_h)) = y_h.$$

- For input  $y \in Y_h$ ,  $\bar{R}_h$  produces one output. This means, the step of deciding  $\bar{F}_h$  and  $\bar{R}_h$  includes, if there are many solutions, making the choice among solutions. This may or may not be explicit in the approximate algorighm  $\bar{R}_h$ . From now on, we assume  $\bar{F}_h$  is bijective between  $X_h$  and  $Y_h$  and  $\bar{R}_h$  is its inverse.

- Using projections, and  $\bar{R}_h$ , we also write

$$R_h = \bar{R}_h \circ I_1^h : Y \rightarrow X_h.$$

The consistency is the notion saying we are not doing irrelevant to the original problem.

We may or may not speak about the whole sequence of  $(F_h)$  and  $(R_h)$ , or speak about only one of them. In the below, we give definitions for a fixed  $h$ .

## Consistency of $F_h$

**Definition 1.** We say  $F_h$  is  $\eta$ -consistent at  $x \in X$  if

$$|F(x) - F_h(x)|_Y < \eta.$$

We call  $|F(x) - F_h(x)|_Y$  the testing error at  $x \in X$ .

## Consistency of $R_h$

**Definition 2.** We say  $R_h$  is  $\eta$ -consistent at  $y \in Y$  if

1.  $R_h(y) \in X$ , and
2.  $|F(R_h) - y|_Y < \eta$ .

We call  $|F(R_h(y)) - y|_Y$  the testing error at  $y \in Y$ .

Why do we separate the error analysis into two steps, (i) testing error analysis; (ii) solution error analysis?

1. Consistency is the minimum requirement, and thus if one is a developer for a certain new problem, one should try to build a consistent method. If not, the method must be abandoned.
2. Problemwisely, and purposewisely, we may or may not want to go further than testing error analysis.
  - We may not care if approximate solution  $\tilde{x}$  is close to the solution  $x$  or not, as long as  $F(\tilde{x}) \simeq y$ .
  - There are some problems, this is inherently so, i.e., one cannot reduce solution errors as much as one wants, while the testing error can be reduced sufficiently small.

Orthogonalization algorithms are good examples:

- Gram-Schmidt algorithm  $\tilde{R}_{GS}$  to get  $(O, R)$  out of  $A$  is inconsistent if  $A$  is rank deficient.  
The algorithm gives me  $O$  and  $R$  such that  $OR \simeq A$  but  $O$  is far from being orthogonal matrix.
- Householder algorithm  $\tilde{R}_H$  to get  $(O, R)$  out of  $A$  is consistent. However, it lacks the solution error control. If  $A'$  is a small perturbation of  $A$ , then the algorithm gives me orthogonal  $O'$  and upper triangular  $R'$  and  $O'R' \simeq A'$  but  $O$  and  $O'$  are not at all close to each other.

However, the Householder algorithm is used everywhere for the right purposes.

## Accuracy of $R_h$

**Definition 3.** We say  $R_h$  is  $\epsilon$ -accurate if one of the solution  $x_*$  of  $F(x) = y$  is such that

$$|x_* - R_h(y)|_X < \epsilon.$$

**Definition 4.** We say  $(R_h)$  is a convergent algorithm if one of the solution  $x_*$  of  $F(x) = y$  is such that

$$|x_* - R_h(y)|_X < \epsilon_h, \quad \epsilon_h \rightarrow 0 \quad \text{as } h \rightarrow 0.$$

A Consistent algorithm is upgraded to an Accurate algorithm, in case true solver or approximate solver satisfy a certain continuity.

## Continuity of true local inverse of $F$

Let  $\bar{y} \in Y$  be fixed, and let  $\bar{x} \in X$  be one solution. Suppose that  $F$  is locally continuously invertible, more precisely,

for every  $\epsilon > 0$ , there exists  $\eta > 0$  such that

$$|y - \bar{y}|_Y < \eta \text{ implies there exists } x \in X \text{ with } F(x) = y \text{ and } |x - \bar{x}| < \epsilon.$$

(This is known as the notion for  $F$  being an open map. This is to say  $F(B(\bar{x}, \epsilon))$  is always an open set.)

Assume consistency + continuity:

1. Suppose  $(R_h)$  is consistent:  $(R_h)$  is  $\eta_h$ -consistent and  $\eta_h \rightarrow 0$  as  $h \rightarrow 0$ .
2. Suppose  $F$  is an open map so that the monotone function

$$\eta(\epsilon)$$

is as in the statement at  $\bar{y}$ .

Conclusion:

Let  $\epsilon > 0$  be fixed. Choose  $h$  small so that  $R_h$  is  $\eta(\epsilon)$ -consistent.

1. We let  $\bar{y} = F(R_h(y))$ . By consistency, we know that

$$|F(R_h(y)) - y|_Y = |\bar{y} - y|_Y < \eta.$$

2. By the Local continuous invertibility of  $F$ ,

There exists  $x \in X$  with  $F(x) = y$  and  $|x - R_h(y)|_X < \epsilon$ .

## Continuity of whole approximate inverses ( $R_h$ )

Instead of working with continuity of true local inverse, we work with continuity of approximate inverses ( $R_h$ ) but as a whole.

Let  $\bar{y} \in Y$  be fixed, and let  $\bar{x} \in X$  be a solution. Suppose  $R_h$  is continuous at  $\bar{y}$  in the sense that

for every  $\epsilon > 0$ , there exists  $\eta > 0$  such that

$$|y - \bar{y}|_Y < \eta \text{ implies that for every } h, |R_h(y) - R_h(\bar{y})| < \epsilon.$$

(This is known as the notion of equi-continuity of a sequence of functions.)

Assume consistency + continuity

1. Suppose  $(F_h)$  is consistent:  $(F_h)$  is  $\eta_h$ -consistently and  $\eta_h \rightarrow 0$  as  $h \rightarrow 0$ .
2. Suppose  $(R_h)$  is equi-continuous so that the monotone function

$$\eta(\epsilon)$$

is as in the statement at  $\bar{y}$ .

Conclusion:

Let  $\epsilon > 0$  be fixed,  $\eta = \eta(\epsilon)$ . Choose  $F_h$  for small  $h$  so that it is  $\eta(\epsilon)$ -consistent.

1. We let  $y_h = F_h(x) = \bar{F}_h(I_0^h(x))$ , then  $R_h(y_h) = I_0^h(x)$ .
2. By consistency, we *a priori* know that at unknown  $x$

$$|F_h(x) - F(x)|_Y = |y_h - y|_Y < \eta.$$

3. By the continuity of  $R_h$ , we *a priori* know that

$$|R_h(y_h) - R_h(y)|_X = |I_0^h(x) - R_h(y)|_X < \epsilon$$

4. Then,

$$|x - R_h(y)|_X \leq |x - I_0^h(x)|_X + |I_0^h(x) - R_h(y)|_X$$

and the first term has to be taken care of by the interpolation error estimate.



# Chapter 11

## A priori error analysis

Continuing discussions on solving the *discretized* problem,

1. We look for  $p_h$  in a span of basis of  $\hat{X}_h$ :

$$p_h = \sum_{i=1}^{n_0} c_i \phi_i(x), \quad \int_{\Omega} p_h = \sum_{i=1}^{n_0} c_i = 0.$$

2. It solves a system of equations

$$\text{for every } \phi_i, \text{ a basis of } \hat{X}_h, \quad \int_{\Omega} K(x) \nabla p_h \cdot \nabla \phi_i = \int_{\Omega} \varphi \phi_i - \int_{\partial\Omega} \psi \phi_i.$$

3. This is formulated as a system with  $(n_0 - 1) \times (n_0 - 1)$  matrix

$$Ac = b.$$

4. The solution

$$p_h = R_h(\ell) = \bar{R}_h(b), \quad b_i = \ell(\phi_i) \in Y^h \simeq \mathbb{R}^{n_0-1}.$$

In the next step, we are supposed to check if  $(R_h)$  is continuous, independently of  $h$ .

1. Suppose  $\ell, \tilde{\ell} \in \hat{Y}$  such that

$$\|\ell - \tilde{\ell}\|_{\hat{Y}} = \sup_{v \in \hat{X}, \|v\|_{\hat{X}}=1} |\ell(v) - \tilde{\ell}(v)| < \eta.$$

2. Let  $p_h, q_h \in \hat{X}^h$  be the solutions obtained by solving linear systems,

$$\text{for every } v_h \in \hat{X}_h, \quad \int_{\Omega} K(x) \nabla p_h \cdot \nabla v_h = \ell(v_h), \quad \int_{\Omega} K(x) \nabla q_h \cdot \nabla v_h = \tilde{\ell}(v_h)$$

3. Subtracting the equation,

$$\text{for every } v_h \in \hat{X}_h, \quad \int_{\Omega} K(x) \nabla(p_h - q_h) \cdot \nabla v_h = \ell(v_h) - \tilde{\ell}(v_h).$$

If we take  $v_h = p_h - q_h \in \hat{X}_h$ ,

$$\begin{aligned} |\langle \nabla(p_h - q_h), \nabla(p_h - q_h) \rangle_K| &= |(\ell - \tilde{\ell})(p_h - q_h)| < \eta \|p_h - q_h\|_{\hat{X}} = \eta \|\nabla(p_h - q_h)\|_{L_K^2} \\ \iff \|\nabla(p_h - q_h)\|_{L_K^2}^2 &< \eta \|\nabla(p_h - q_h)\|_{L_K^2} \iff \|\nabla(p_h - q_h)\|_{L_K^2} < \eta. \end{aligned}$$

4. Hence, we can set  $\eta(\epsilon) = \epsilon$ . This is independent of  $h$ .

1. Now, let  $p \in \hat{X}$  be the true solution,  $I^h p \in \hat{X}^h$  be a projection.
2.  $\ell_p$  and  $\ell_{I^h p}$  are defined in  $\hat{Y}$ .
3. Our solution  $p_h \in \hat{X}^h$  satisfies

$$\text{for every } v_h \in \hat{X}_h, \quad \int_{\Omega} K(x) \nabla p_h \cdot \nabla v_h = \ell_p(v_h).$$

4.  $I^h p$  is a solution of

$$\text{for every } v_h \in \hat{X}_h, \quad \int_{\Omega} K(x) \nabla I^h p \cdot \nabla v_h = \ell_{I^h p}(v_h).$$

5. Hence, we can conclude from the continuity before that

$$\|\ell_p - \ell_{I^h p}\|_{\hat{Y}} < \eta \implies \|\nabla(p_h - I^h p)\|_{L_K^2}^2 < \eta.$$

6. By triangle inequality,

$$\|\nabla(p - p_h)\|_{L_K^2}^2 \leq \|\nabla(p - I^h p)\|_{L_K^2}^2 + \|\nabla(p_h - I^h p)\|_{L_K^2}^2 < \|\nabla(p - I^h p)\|_{L_K^2}^2 + \eta.$$

7. From the consistency of  $(F_h)$  calculated before,

$$\|\ell_p - \ell_{I^h p}\|_{\hat{Y}} \leq \|\nabla(p - I^h p)\|_{L_K^2}$$

and thus, we conclude finally

$$\|\nabla(p - p_h)\|_{L_K^2}^2 \leq 2\|\nabla(p - I^h p)\|_{L_K^2}^2.$$

8. Hence, the Galerkin method is a convergent method, converging as fast as  $\|\nabla(p - I^h p)\|_{L_K^2}$  does.

9. In fact, better calculation without using triangle inequality is possible:

This is from the observation that

$$\begin{aligned} \text{for every } v_h \in \hat{X}_h, \quad & \int_{\Omega} K(x) \nabla p_h \cdot \nabla v_h = \ell_p(v_h). \\ \implies \text{for every } v_h \in \hat{X}_h, \quad & \int_{\Omega} K(x) \nabla(p - p_h) \cdot \nabla v_h = 0. \end{aligned}$$

10. Then for any  $v \in \hat{X}$ ,

$$\begin{aligned} \int_{\Omega} K(x) \nabla(p - p_h) \cdot \nabla v &= \int_{\Omega} K(x) \nabla(p - p_h) \cdot \nabla I^h v + \int_{\Omega} K(x) \nabla(p - p_h) \cdot \nabla(v - I^h v) \\ &= \int_{\Omega} K(x) \nabla(p - p_h) \cdot \nabla(v - I^h v). \end{aligned}$$

11. Taking  $v = p - p_h$ , we first observe

$$\begin{aligned} I^h v &= I^h p - I^h p_h = I^h p - p_h, \\ v - I^h v &= p - I^h p. \end{aligned}$$

Plugging-in,

$$\begin{aligned} \|\nabla(p - p_h)\|_{L_K^2}^2 &\leq \|\nabla(p - p_h)\|_{L_K^2}^2 \|\nabla(p - I^h p)\|_{L_K^2} \\ \iff \|\nabla(p - p_h)\|_{L_K^2} &\leq \|\nabla(p - I^h p)\|_{L_K^2}. \end{aligned}$$

with the constant 1 in place of 2 in previous calculation.

*Remark 11.1.* If the coefficient  $K$  has better smoothness, and the domain  $\Omega$  is sufficiently good, so that the true solution  $p$  lies in  $W^{2,2}(\Omega)$ , then there, we know from the Bramble-Hilbert that

$$\|\nabla(p - I^h p)\|_{L_K^2} \leq Ch \|p\|_{W^{2,2}(\Omega)},$$

guaranteeing the first order convergence for the error in gradient and

$$\|p - I^h p\|_{L^2} \leq C'h^2 \|p\|_{W^{2,2}(\Omega)},$$

guaranteeing the 2nd order convergence for the error.



# **Chapter 12**

# **Lab. Project**

1. Mesh generation : use gmsh software and export.
2. Simplicial Complex construction.
3. Matrix system assemble.
4. Linear system solving.
5. Solution plot/table, error plot/table.

## Simplicial Complex construction

In Julia, we use Gmsh Package.

1. getNodes,
2. getElements,
3. getNodesForPhysicalGroup.

The objectives are the followings.

### **Bijective maps between numberings and Keys: domain**

1. Vertices:

$$\{1, 2, \dots, n_0\} \longleftrightarrow \{1, 2, \dots, n_0\}$$

2. Edges:

$$\{1, 2, \dots, n_1\} \longleftrightarrow \{(i_{1,0}, i_{1,1}), (i_{2,0}, i_{2,1}), \dots, (i_{n_1,0}, i_{n_1,1})\}$$

3. Faces:

$$\{1, 2, \dots, n_2\} \longleftrightarrow \{(i_{1,0}, i_{1,1}, i_{1,2}), (i_{2,0}, i_{2,1}, i_{2,2}), \dots, (i_{n_2,0}, i_{n_2,1}, i_{n_2,2})\}$$

### **Bijective maps between numberings and Keys: boundary**

1. Boundary Vertices:

$$\{1, 2, \dots, \ell_0\} \longleftrightarrow \{i_1^b, i_2^b, \dots, i_{\ell_0}^b\}$$

2. Boundary Edges:

$$\{1, 2, \dots, \ell_1\} \longleftrightarrow \{(i_{1,0}^b, i_{1,1}^b), (i_{2,0}^b, i_{2,1}^b), \dots, (i_{\ell_1,0}^b, i_{\ell_1,1}^b)\}$$

### **Simplicial complex $C_0, C_1, C_2$ .**

1. Any information repeatedly used.
2. Face relationships

## Matrix equation assemble

Now, we compose the matrix equation

$$Ac = b,$$

where

$$\text{for } i, j = 1, 2, \dots, n_0 \quad A_{ij} = \int_{\Omega} (\nabla \phi_i)^T K(x) \nabla \phi_j, \quad b_i = \int_{\Omega} \varphi \phi_i - \int_{\partial \Omega} \psi \phi_i.$$

The approximation solution will be  $\tilde{u}(x) = \sum_{j=1}^{n_0} c_j \phi_j(x)$ .

Expressions are global basis based, which are vertex numbering based. It is convenient to write expressions triangle numbering based.

1. The set of triangles  $\Gamma_i$  for  $i = 1, 2, \dots, n_0$ .

$$\Gamma_i = \{\sigma \in Key_2 \mid \text{the vertex } i \text{ is a face of triangle } \sigma\}.$$

2. If  $\sigma \in \Gamma_i$ ,  $\sigma = (i_0, i_1, i_2)$  and  $i_e = i$ , we write

$$\dot{e} = i.$$

$$\phi_i = \sum_{\sigma \in \Gamma_i} \sum_{e, \dot{e}=i} \lambda_{\sigma,e}. \quad \lambda_{\sigma,e} \text{ is the local } e\text{-th basis.}$$

3. We do the similar for the boundary, defining  $\Gamma_i^b$ , the set of boundary edges having  $i$  as its face.

4. For later purpose, we compute the gradient vector of  $\lambda_e$ :

- (a) We recall the coordinate  $z = y - y_{\sigma,c}$ ,  $y_{\sigma,c}$  is the center of mass of  $\sigma$ .
- (b) We also recall that  $\lambda_e(z)$  is linear in  $z$  and thus  $\nabla \lambda_{\sigma,e}$  is a constant vector in interior of  $\sigma$ .
- (c) We compute the gradient, for example let  $e = 0$  and  $f = \dot{e} = (1, 2)$ .
- (d) Then  $\nabla \lambda_0(z) = s \frac{(z_2 - z_1)^\perp}{2|\sigma|}$ .
- (e) The general formula for the constant vector is

$$\nabla \lambda_e(z) = (-1)^e s \frac{(z_{f[2]} - z_{f[1]})^\perp}{2|\sigma|}, \quad \dot{e} = f.$$

- (f)  $\check{0} = (1, 2)$ ,  $\check{1} = (0, 2)$ ,  $\check{2} = (0, 1)$ .

## Assembling right-hand-side $b$

We consider the following triangle-visiting loop.

1. Initialize  $b_i = 0$  for all  $i = 1, 2, \dots, n_0$ .
2. For each triangle  $\sigma \in Key_2$ ,
3. For each  $e \in \{0, 1, 2\}$
4.  $i$  is then the unique vertex number, where  $i = \dot{e}$  and  $\sigma \in \Gamma_i$ .
5. Add the contribution of the integral  $\int_{\overline{\text{pts}}(\sigma)} \varphi \lambda_{\sigma,e}$  to  $b_i$ .

This is followed by the similar boundary edge-visiting loop.

1. For each boundary edge  $a \in BKey_1$ ,
2. For each  $e \in \{0, 1\}$
3.  $i^b$  is then the unique vertex number, where  $i^b = \dot{e}$  and  $a \in \Gamma_i^b$ .
4. Add the contribution of the line integral  $\int_{\overline{\text{pts}}(a)} -\psi \lambda_{a,e}$  to  $b_i$ .

By the moment two loops are done, the right-hand-side  $b_i$  is assembled.

In practice, approximation  $\tilde{b}_i$  of  $b_i$  maybe collected. Namely

- The first contribution is computed simply by  $\frac{|\sigma|}{3} \varphi(y_{\sigma,c})$ .
- The second contribution is computed simply by  $-\frac{|a|}{2} \psi(y_{mid})$ .

## Assembling Matrix $A$

We consider the following triangle-visiting loop.

1. Initialize  $A_{ij} = 0$  for all  $i, j = 1, 2, \dots, n_0$ .
2. For each triangle  $\sigma \in Key_2$ ,
3. For each  $e \in \{0, 1, 2\}$  and  $e' \in \{0, 1, 2\}$
4.  $i$  and  $j$  are vertex numbers such that  $\dot{e} = i$ , and  $\dot{e}' = j$ .
5. Add the contribution of the integral  $\int_{\overline{\text{pts}}(\sigma)} (\nabla \lambda_{\sigma,e})^T K(x) \nabla \lambda_{\sigma,e'} \text{ to } A_{ij}.$

By the moment this loops is done, the matrix  $A$  is assembled.

In practice, approximation  $\tilde{A}_{ij}$  of  $A_{ij}$  maybe collected. Namely

- The contribution is computed simply by  $|\sigma| (\nabla \lambda_{\sigma,e})^T K \nabla \lambda_{\sigma,e'}$ , where the constant matrix  $K = K(y_{\sigma,c})$ .
- In that case, the contribution is (with  $\check{e} = f, \check{e}' = g$ )

$$(-1)^{e+e'} \frac{((z_f[2] - z_f[1])^\perp)^T K (z_g[2] - z_g[1])^\perp}{4|\sigma|}.$$

If  $A$  and  $b$  are legitimately assembled,

1. The rank of  $A$  must be  $n_0 - 1$ ,
2.  $b$  must lie in  $\text{Ran}(A)$ , i.e.,  $b \perp (1, 1, \dots, 1)$ .
3. Since linear solver in Python, Julia, automatically deal with the least square solution, we may let the built-in linear solver solve the problem  $Ax = b$  as it is in the rank-deficient form.



# Chapter 13

# Final Project

We first consider the following local basis products:

- Suppose  $\sigma$  is an oriented triangle  $[z_0, z_1, z_2]$ , with center of mass

$$z_c = \frac{1}{3}(z_0 + z_1 + z_2) = 0.$$

- Let  $\lambda_0(z), \lambda_1(z), \lambda_2(z)$  be the barycentric coordinate of  $z$ .
- We compute for  $e, e' \in \{0, 1, 2\}$

$$K_{e,e'} = \int_{\overline{\text{pts}}(\sigma)} \lambda_e(z) \lambda_{e'}(z) dz.$$

1. We first compute the integral

$$\tilde{K}_{e,e'} = \int_{\Lambda^2} \lambda_e \lambda_{e'} d\mathcal{H}^2(\lambda),$$

over the regular triangle

$$\Lambda^2 = \{(\lambda_0, \lambda_1, \lambda_2) \mid \lambda_0 + \lambda_1 + \lambda_2 = 1, \quad \forall e, \lambda_e \geq 0\}$$

as we did in midterm take home for dimensions 3.

2. We parametrize  $\Lambda^2$  by the set

$$\tau^2 = \{(\tau_1, \tau_2) \mid \tau_1, \tau_2 \geq 0, \quad \tau_1 + \tau_2 \leq 1\},$$

that is  $(\lambda_0(\tau), \lambda_1(\tau), \lambda_2(\tau)) = (1 - \tau_1 - \tau_2, \tau_1, \tau_2)$ .

3. This map has the Jacobian factor

$$\sqrt{\det D\lambda^T D\lambda} = \sqrt{\det \begin{pmatrix} -1 & 1 & 0 \\ -1 & 0 & 1 \end{pmatrix} \begin{pmatrix} -1 & -1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}} = \sqrt{\det \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}} = \sqrt{3}.$$

4. Therefore, the area of the regular triangle  $\Lambda^2$  (with the side length =  $\sqrt{2}$ )

$$I_0 = \int_{\Lambda^2} 1 d\mathcal{H}^2(\lambda) = \int_{\tau^2} 1 \sqrt{3} d\tau = \sqrt{3} \int_{\tau^2} 1 d\tau = \frac{\sqrt{3}}{2}.$$

5. Since  $\Lambda^2$  is a regular triangle,  $\tilde{K}_{00} = \tilde{K}_{11} = \tilde{K}_{22}$ .

$$\begin{aligned}\tilde{K}_{00} &= \int_{\tau_1=0}^1 \int_{\tau_2=0}^{1-\tau_1} (1 - \tau_1 - \tau_2)^2 \sqrt{3} d\tau_2 d\tau_1 \\ &= \sqrt{3} \int_{\tau_1=0}^1 \left[ -\frac{1}{3}(1 - \tau_1 - \tau_2)^3 \right]_{\tau_2=0}^{1-\tau_1} d\tau_1 = \sqrt{3} \int_{\tau_1=0}^1 \frac{(1 - \tau_1)^3}{3} d\tau_1 \\ &= \sqrt{3} \left[ -\frac{1}{12}(1 - \tau_1)^4 \right]_{\tau_1=0}^1 = \frac{\sqrt{3}}{12}.\end{aligned}$$

6. Since  $\Lambda^2$  is a regular triangle,  $\tilde{K}_{01} = \tilde{K}_{02} = \tilde{K}_{12}$ .

$$\begin{aligned}\tilde{K}_{01} &= \int_{\tau_1=0}^1 \int_{\tau_2=0}^{1-\tau_1} (1 - \tau_1 - \tau_2) \tau_1 \sqrt{3} d\tau_2 d\tau_1 \\ &= \sqrt{3} \int_{\tau_1=0}^1 \left[ -\frac{1}{2}(1 - \tau_1 - \tau_2)^2 \tau_1 \right]_{\tau_2=0}^{1-\tau_1} d\tau_1 = \sqrt{3} \int_{\tau_1=0}^1 \frac{\tau_1(1 - \tau_1)^2}{2} d\tau_1 \\ &= \frac{\sqrt{3}}{24}.\end{aligned}$$

7. Hence,

$$\tilde{K} = \frac{\sqrt{3}}{2} \frac{1}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}.$$

8. One can show

$$\frac{1}{|\sigma|} \int_{\overline{\text{pts}}(\sigma)} \lambda_e(z) \lambda_{e'}(z) dz = \frac{1}{I_0} \int_{\Lambda^2} \lambda_e \lambda_{e'} d\mathcal{H}^2(\lambda).$$

Using this,

$$K = \frac{|\sigma|}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix} = |\sigma| \hat{K}, \quad \hat{K} = \frac{1}{12} \begin{pmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{pmatrix}.$$

**Objective 1** : Discretization, and Computation of  $L^2$  norm of  $\sum_j c_j \phi_j(x)$

**Objective 1-1** : You should be able to import mesh data of your domain, using your programming language .

**Objective 1-2** : Given the mesh, you should be able to compute the  $L^2$  norm of a function  $\sum_j c_j \phi_j(x)$ .  $(\phi_j)_{j=1}^{n_0}$  are the global basis.

- Let  $v^h(x) = \sum_{j=1}^{n_0} c_j \phi_j(x)$ .
- The square of its  $L^2$  norm is then

$$\|v^h\|_{L^2(\Omega)}^2 = \int_{\Omega} \sum_{i,j=1}^{n_0} c_i c_j \phi_i(x) \phi_j(x) = \sum_{i,j=1}^{n_0} c_i M_{ij} c_j, \quad M_{ij} = \int_{\Omega} \phi_i(x) \phi_j(x).$$

- Again, to compute the number  $c^T M c$ , we rather want to work with triangle-based expression.

We consider the following triangle-visiting loop.

#### Assembling of matrix $M$

1. Initialize  $M_{i,j} = 0$  for every  $i, j = 1, 2, \dots, n_0$ .
2. For each  $\sigma \in Key_2$ ,
3. For each  $e, e' \in \{0, 1, 2\}$
4.  $i$  and  $j$  are such that  $i = \dot{e}, j = \dot{e}'$ .
5. Add the contribution  $\int_{\overline{\text{pts}}(\sigma)} \lambda_e(z) \lambda_{e'}(z) dz = |\sigma| \hat{K}_{e,e'}$  to  $M_{ij}$ .

Once the loop is done,  $M_{ij}$  is assembled, and

$$\left\| \sum_{j=1}^{n_0} c_j \phi_j(x) \right\|_{L^2}^2 = c^T M c.$$

To check if the computation is correct, compute  $c^T M c$  with  $c = (1, 1, \dots, 1)^T$ . Because  $\sum_{j=1}^{n_0} \phi_j(x) \equiv 1$ , this is to compute

$$\int_{\Omega} 1 = |\Omega| = \frac{3\pi}{4},$$

if the domain  $\Omega$  is the  $\frac{3\pi}{2}$  sector of the disk with radius 1.

## Objective 2 : Solve the Neumann Boundary Value Problem.

You may compose your Neumann Boundary Value Problem, and solve the problem. Below is one instance.

**Preparation of exact solution  $u$ , and data of the problem  $\varphi \equiv 0$  and  $\psi$ .**

$u(r, \theta) = r^{\frac{2}{3}} \sin\left(\frac{2}{3}\theta\right)$  is a Harmonic function in the domain  $\Omega$ :

- The gradient

$$\nabla u(r, \theta) = (u_r, \frac{1}{r}u_\theta) = \left(\frac{2}{3}r^{-\frac{1}{3}} \sin\left(\frac{2}{3}\theta\right), \frac{2}{3}r^{-\frac{1}{3}} \cos\left(\frac{2}{3}\theta\right)\right)$$

- The laplacian

$$\Delta u = \frac{1}{r}(ru_r)_r + \frac{1}{r^2}u_{\theta\theta} = 0 (= -\varphi).$$

- You compute the outward normal component over the  $\partial\Omega$  from the exact solution

$$-\nabla u \cdot n = \psi$$

### The Neumann Boundary Value Problem

$$\begin{aligned} -\Delta u &= 0 && \text{in } \Omega \\ -\nabla u \cdot n &= \psi && \text{on } \partial\Omega. \end{aligned}$$

The solution solved from PDE must be the exact solution above.  
(To fix non-uniqueness, we set  $u(0) = 0$ .)

### Assembling Matrix $A$ and r-h-h $b$

This is done following procedure in page 90-91.

## Result Presentation

**Plot of the approximation solution  $u^h$ .**

Prepare the graph so that the approximation solution  $u^h$  looks similar to the exact solution.

### Error Table

For the project, let us first compare the numerical solution  $u^h$  to the projection

$$I^h u = \sum_{i=1}^{n_0} u(x_i) \phi(x),$$

as an alternative of the exact solution.

We consider the difference

$$u^h - I^h u = \sum_{i=1}^{n_0} d_i \phi(x), \quad \text{for some } (d_i)_{i=1}^{n_0}.$$

We can compute the  $L^2$  norm and gradient  $L^2$  norm of the difference:

$$\begin{aligned} \|u^h - I^h u\|_{L^2}^2 &= d^T M d, \\ \|\nabla(u^h - I^h u)\|_{L^2}^2 &= d^T A d. \end{aligned}$$

Fix  $h_0$  and reduce the mesh size by factor of  $\frac{1}{2}$  and collect errors in table.

Mesh size	$L^2$ Error	Gradient $L^2$ error
$(\frac{1}{2})^0 h_0$		
$(\frac{1}{2})^1 h_0$		
$(\frac{1}{2})^2 h_0$		
$(\frac{1}{2})^3 h_0$		

Table 13.1: Numerical results