

Emory University  
**MATH 572 Numerical PDEs**  
Learning Notes

Jiuru Lyu

April 27, 2026

## Contents

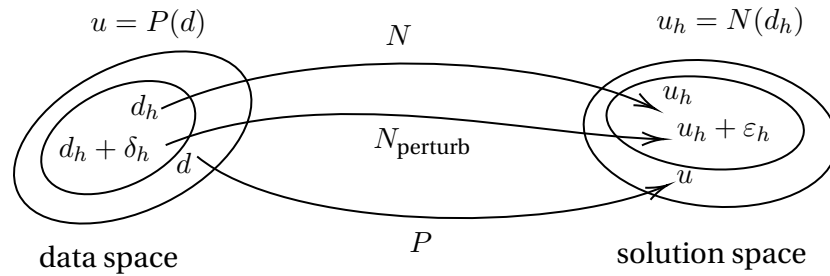
<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Basic Concept . . . . .	3
1.2	Functional Analysis . . . . .	3
1.2.1	Lebesgue Integration . . . . .	6
1.2.2	Distributional Derivatives . . . . .	7
1.2.3	Special Spaces . . . . .	9
1.3	Solving $Ax = b$ . . . . .	11
1.4	General Statements of PDEs . . . . .	13
<b>2</b>	<b>Poisson Equation</b>	<b>14</b>
2.1	Poisson Equation as a Minimization Problem . . . . .	14
2.2	Finite Differences . . . . .	15
2.2.1	Poisson Equation in 1D . . . . .	16
2.3	Galerkin Methods . . . . .	20
2.4	Finite Element: One Possible Choice of $X_N$ . . . . .	25
2.4.1	Finite Element in 1D . . . . .	25
2.4.2	Finite Elements in Multiple Dimension . . . . .	32
2.5	Mixed Problems . . . . .	39
2.6	Remarks on FE . . . . .	41
2.7	Spectral Method . . . . .	42
<b>3</b>	<b>Advection-Diffusion Problem</b>	<b>45</b>
3.1	Finite Difference of 1D Problem . . . . .	45
3.2	Finite Difference in 2D+ . . . . .	47
3.3	Finite Elements in 1D . . . . .	48
3.4	Convection-Dominated problem . . . . .	50

---

3.5	Strongly Consistent Stabilization . . . . .	51
3.6	Reaction-Dominated Problem . . . . .	53
<b>4</b>	<b>Parabolic Equations</b>	<b>55</b>
4.1	Weak Formulation and Well-Posedness . . . . .	55
4.2	Finite Element and Semi-Discretization . . . . .	58
4.3	Space-Time Finite Element (Space-Time FEM) . . . . .	62
<b>5</b>	<b>Least-Square FEM and PINNs</b>	<b>64</b>
5.1	Least-Square FEM (LS-FEM) . . . . .	64
5.2	Physics-Informed Neural Networks (PINNs) . . . . .	67
<b>6</b>	<b>Hyperbolic Problems</b>	<b>70</b>
6.1	Finite Differences . . . . .	70
6.1.1	Conservation Laws . . . . .	70
6.1.2	Methods for Linear System . . . . .	71
6.1.3	Wave Equation . . . . .	72
6.2	Analysis of FD Methods . . . . .	72
6.2.1	Consistency . . . . .	72
6.2.2	Convergence . . . . .	73
6.2.3	Stability . . . . .	73
6.3	Von Neumann Analysis of FD Methods . . . . .	77
6.4	Finite Elements . . . . .	77

# 1 Introduction

## 1.1 Basic Concept



- Physical Problem: Solution  $u$  and data  $d$ .
- Numerical Problem: Numerical solution  $u_h$  and approximated data  $d_h$ , where  $h$  measures how fine the approximation is.

Some key concepts are

- *Convergence*: Does  $u_h \rightarrow u$  when  $h \rightarrow 0$ ?
- *Consistency*: Does  $N \rightarrow P$  when  $h \rightarrow 0$ ?
- *Stability*: Does  $\epsilon_h \rightarrow 0$  when  $\delta_h \rightarrow 0$ ?

Another perspective:

$$P(u, d) = 0 \quad \text{true problem}, \quad P_N(u_N, d_N) = 0 \quad \text{numerical/approximated problem}$$

- Consistency: when  $N \rightarrow \infty$  ( $h \rightarrow 0$ ),  $P_N \rightarrow P$ .

$$P_N(u, d_N) = \tau \quad (\text{consistency error}).$$

When  $N \rightarrow \infty$ ,  $\tau \rightarrow 0$ .

- Stability:  $P_N(u_N + \delta u, d_N + \delta) = 0$ .

$$|\delta u| \leq c|\delta| \implies \text{when } \delta \rightarrow 0, \delta u \rightarrow 0.$$

- Lax-Richtmyer Theorem: For linear PDE problem, consistency and stability  $\implies$  convergence.

## 1.2 Functional Analysis

**Definition 1.2.1 (Functional).** A functional is  $\mathcal{F}(\cdot) : X \rightarrow \mathbb{R}$ . Let  $f(x), g(x) \in X$ . A functional is *linear* if  $\mathcal{F}(\lambda f(x) + \mu g(x)) = \lambda \mathcal{F}(f(x)) + \mu \mathcal{F}(g(x))$ .

**Example 1.2.2**

For example,  $\int_a^b f(x) dx$  is a functional.

**Definition 1.2.3 (Form).** A *form* is  $a(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$ , where  $f, g \in X$ . A form is *bilinear* if:

- $a(\lambda f_1 + \mu f_2, g) = \lambda a(f_1, g) + \mu a(f_2, g)$ , and
- $a(f, \lambda g_1 + \mu g_2) = \lambda a(f, g_1) + \mu a(f, g_2)$ .

**Definition 1.2.4 (Space).** A *space* is set with operations defined. If  $X$  is a *space*, then

- $f \in X \implies \lambda f \in X$ , and
- $f, g \in X \implies \lambda f + \mu g \in X$ .

**Definition 1.2.5 (Distance).** A *distance* is  $d(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$  such that  $\forall f, g, w \in X$ :

- $d(f, g) = d(g, f)$ , [symmetry]
- $d(f, g) \geq 0$  and  $d(f, g) = 0 \iff f = g$ , and
- $d(f, g) \leq d(f, w) + d(w, g)$  [triangle inequality]

A space with a distance is a *metric space*.

**Definition 1.2.6 (Completeness).** Let  $(X, d)$  be a metric space. Let  $f_n$  be a sequence in  $X$  and  $f$  be a candidate limit.

1. The sequence is *convergent* if  $d(f_n, f) \rightarrow 0$  as  $n \rightarrow +\infty$ .
2. The sequence is *Cauchy* if  $d(f_n, f_m) \rightarrow 0$  as  $n, m \rightarrow +\infty$ .

In general,  $1 \implies 2$ , but  $2 \not\implies 1$  [ $d(f_n, f_m) \leq d(f_n, f) + d(f, f_m)$  by triangle inequality].

If in a metric space,  $1 \iff 2$ , then the metric space is *complete*.

**Definition 1.2.7 (Norm).** A *norm* is  $\|\cdot\|_X : X \rightarrow \mathbb{R}$  such that

- $\|f\|_X \geq 0$  and  $\|f\|_X = 0 \iff f = 0$ ,
- $\|\lambda f\|_X = \lambda \|f\|_X$ , and
- $\|f + g\|_X \leq \|f\|_X + \|g\|_X$ .

A space is *normed* if a norm is defined.

A normed space is also a metric space:  $d(f, g) = \|f - g\|_X$ .

**Definition 1.2.8 (Banach Space).** A *Banach space* is a complete normed space.

**Definition 1.2.9 (Scalar Product).**  $(\cdot, \cdot)_X : X \times X \rightarrow \mathbb{R}$  is a symmetric bilinear form such that

- $(f, g)_X = (g, f)_X$
- $(f, f)_X \geq 0$  and  $(f, f)_X = 0 \iff f = 0$ .
- $|(f_1, f_2)_X| \leq (f_1, f_1)_X^{1/2} (f_2, f_2)_X^{1/2}$  [*Cauchy-Swcharz Inequality*]

A scalar product induces a norm:

$$\|f\|_X^2 = (f, f)_X.$$

So, a space with a scalar product is automatically normed.

**Example 1.2.10**

$$\begin{aligned} (f_1 + f_2, f_1 + f_2)_X &= (f_1, f_1)_X + 2(f_1, f_2)_X + (f_2, f_2)_X \\ &\leq \|f_1\|_X^2 + 2\|f_1\|_X\|f_2\|_X + \|f_2\|_X^2 \\ &= (\|f_1\|_X + \|f_2\|_X)^2. \end{aligned}$$

**Definition 1.2.11 (Hilbert Space).** A *Hilbert space* is a Banach space equipped with a scalar product.

**Definition 1.2.12 (Orthogonality).** The *angle* between two functions is defined as

$$\angle(f_1, f_2) : \cos(\varphi) = \frac{(f_1, f_2)}{\|f_1\|\|f_2\|} = \frac{(f_1, f_2)}{(f_1, f_1)^{1/2}(f_2, f_2)^{1/2}}$$

Two functions are *orthogonal* to each other  $\iff (f, g) = 0$ . We denote  $f \perp g$ .

**Find the Best Approximation** Let  $X_L \subset X$ , where  $X_L$  is “small.”  $X_L$  and  $X$  are both Hilbert spaces. We want to find the *best* approximation of  $f_L$  of  $f$  in  $X_L$ . So, we need to find  $f_L$  s.t.

$$(f - f_L, v_L) = 0 \quad \forall v_L \in X_L.$$

This is an extension of the Pythagorean Theorem. If  $g_L \in X_L$ :

$$\begin{aligned}\|f - g_L\|^2 &= (f - g_L, f - g_L) = (f - f_L + f_L - g_L, f - f_L + f_L - g_L) \\ &= \|f - f_L\|^2 + 2 \underbrace{(f - f_L, f_L - g_L)}_{=0, \text{ orthogonality}} + \|f_L - g_L\|^2 \\ &= \|f - f_L\|^2 + \|f_L - g_L\|^2 \\ &\geq \|f - f_L\|^2.\end{aligned}$$

### 1.2.1 Lebesgue Integration

**Definition 1.2.13 ( $L^p$  Spaces).**

$$L^p(a, b) = \left\{ f : (L) \int_a^b f^p dx < +\infty \right\}, \quad \text{with } p = 1, 2, \dots, \in \mathbb{N}.$$

Here,  $(L)$  indicates that the integration is done in Lebesgue sense.  $L^p(a, b)$  is the collection of functions that can be integrated in interval  $(a, b)$  in the Lebesgue sense without blowing up. These are Banach spaces with the norm

$$\|f\|_p = \left( \int_a^b |f|^p dx \right)^{1/p}.$$

Similarly,

$$L^\infty(a, b) = \{f : \text{bounded (up to all-measurable intervals) in } (a, b)\}$$

is also a Banach space with norm

$$\|f\|_\infty = \max_{(a,b)} |f|.$$

Specially,  $L^2(a, b)$  or  $L^2(\Omega)$  with  $\Omega \subseteq \mathbb{R}$  is a Hilbert space:

$$(f, g) = \int_a^b fg dx \quad \text{and} \quad \|f\|_2 = \left( \int_a^b f^2 dx \right)^{1/2} = (f, f)^{1/2}.$$

**Definition 1.2.14 (Conjugate Spaces).**  $L^p$  and  $L^q$  are called *conjugated spaces* if we have the function  $f \in L^p(a, b)$  and  $g \in L^q(a, b)$  with

$$\frac{1}{p} + \frac{1}{q} = 1.$$

We then have

$$(L) \int_a^b fg dx < +\infty \quad \text{or} \quad fg \in L^1(a, b)$$

**Remark.**  $L^2$  is conjugate with itself. We have

$$|(f, g)_{L^2}| \leq \|f\|_2 \|g\|_2.$$

### 1.2.2 Distributional Derivatives

**Definition 1.2.15 (Dual Spaces).** The set of linear bounded functionals obtained on a space  $X$  is per se a functional space, that we call *dual space* and denote as  $X'$ :

$$X' = \{\mathcal{F}(f) \text{ for } f \in X, \text{ linear and s.t. } \|\mathcal{F}(f)\| \leq c\|f\|_X\}.$$

#### Theorem 1.2.16 Riesz Representation Theorem

In  $L^2$ , the dual of  $L^2$  is isometric to  $L^2$ . [This means that any linear and continuous functional in  $L^2$  corresponds to an element within  $L^2$  itself.]

**Definition 1.2.17 (Support).** The *support* of a function is where the function  $\neq 0$  on the interval or region of the domain of definition. In particular, if the support is closed and bounded region, we say the function has *compact support*.

#### Definition 1.2.18 (The Space $\mathcal{D}$ ).

$$\mathcal{D} = \{\text{functions in } \mathcal{C}^\infty(\mathbb{R}^n) \text{ (} n \geq 1 \text{) with compact support}\}.$$

The dual space of  $\mathcal{D}$ ,  $\mathcal{D}'$  is defined in the following notion:

$$\mathcal{D}' : \langle T, f \rangle_{\mathcal{D}} \in \mathbb{R}, \text{ where } T \in \mathcal{D}' \text{ is linear and continuous and } f \in \mathcal{D}.$$

$T$  to be continuous is defined in the following sense: Let  $\text{supp}(f_n) \subseteq K$  and  $f_n \rightarrow f$  and  $n \rightarrow \infty$ . Then,  $T$  is *continuous* if  $\langle T, f_n \rangle \rightarrow \langle T, f \rangle$ .

#### Example 1.2.19

Let  $H(x)$  be the step function:

$$H(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0. \end{cases}$$

Then,

$$\langle T, f \rangle = \int_{-\infty}^{+\infty} H f \, dx = \int_0^{+\infty} f \, dx.$$

The functional  $\langle T_0, f \rangle = f(0)$  belongs to  $\mathcal{D}'$ :

$$|\langle T_0, f \rangle| \leq \max_{x \in \mathbb{R}} |f(x)|.$$

**Definition 1.2.20 (Distributions).** The objects of  $\mathcal{D}'$  are called *distributions* or generalized functions.

**Theorem 1.2.21**

$$\langle T', f \rangle = -\langle T, f' \rangle.$$

**Proof 1.** Let  $g \in L^2$ . Then,

$$\langle g', f \rangle = \int_{\mathbb{R}} g' f = [gf]_{-\infty}^{+\infty} - \int_{\mathbb{R}} g f' = -\langle g, f' \rangle$$

Note that  $[gf]_{-\infty}^{+\infty} = 0$  because  $\lim_{x \rightarrow \infty} f = 0$ , and  $f$  has compact support. Q.E.D. ■

**Example 1.2.22**

Consider

$$H(x) = \begin{cases} 0 & \text{for } x \leq 0 \\ 1 & \text{for } x > 0. \end{cases}$$

Then, we can *differentiate*  $H$ :

$$\int_{-\infty}^{+\infty} H' f = - \int_{-\infty}^{+\infty} H f' = - \int_0^{+\infty} f = -[f]_0^{+\infty} = f(0) = \langle T_0, f \rangle.$$

$T_0$  is called the *Dirac- $\delta$*  and it is the mathematical description of an *impulse*. It is *not* a function in the traditional sense.

**Definition 1.2.23 (Dirac- $\delta$ ).** *Dirac- $\delta$*  is the element in  $\mathcal{D}'$  such that

$$\langle \delta, f \rangle = f(0)$$

We know that  $H' = \delta$ .

**Theorem 1.2.24**

$$L^2 \subset \mathcal{D}'.$$

## 1.2.3 Special Spaces

**Definition 1.2.25 (Sobolev Spaces).** Space  $H^k(a, b)$  (or  $H^k(\Omega)$ ) is the space of functions  $f$  s.t. it belongs to  $L^2$  together with all the derivatives up to the order  $k$ .

$$H^k(a, b) = \left\{ f \text{ s.t. } f, f', \dots, f^{(k)} \in L^2 = H^0 \right\}.$$

In multiple dimensions, the derivatives are all the derivatives of order  $k$ :

$$\frac{\partial^{\alpha_1 + \alpha_2 + \dots + \alpha_m} f}{\partial^{\alpha_1} x_1 \partial^{\alpha_2} x_2 \dots \partial^{\alpha_m} x_m} \in L^2 \quad \text{with} \quad \sum_i \alpha_i = k.$$

These spaces are called *Sobolev spaces*. A more generalized definition:

$$W^{k,p} = \left\{ f \text{ s.t. } f, f', \dots, f^{(k)} \in L^p \right\}.$$

So,  $H^k = W^{k,2}$ .

**Theorem 1.2.26 Properties of  $H^k$** 

- $H^k(a, b)$  is a Hilbert space with the scalar product

$$(f, g)_{H^k} = \sum_{i=0}^k (f^{(i)}, g^{(i)})_{L^2}$$

- $H^1(a, b)$  has the norm

$$\|f\|_{H^1}^2 = \|f\|_{L^2}^2 + \|f'\|_{L^2}^2$$

- $H^1 \subset L^2$ . So,

$$\|f\|_{H^1}^2 \geq \|f\|_{L^2}^2.$$

In general,

$$H^{k+1} \subset H^k \subset \dots \subset L^2.$$

**Definition 1.2.27 (Trace Operator).** *Trace* is an operator that maps a function in  $H^k$  to its value of the boundary, and, in general, reduces regularity.

$$\gamma : H^k(\Omega) \rightarrow H^{k-\frac{1}{2}}(\partial\Omega)$$

is the (continuous) trace operator.

**Definition 1.2.28 (The Space  $H_0^1$ ).** The space  $H_0^1(\Omega)$  is the space of  $H^1$  functions with null trace on  $\partial\Omega$ .

**Theorem 1.2.29 Poincare Inequality**

$\exists c > 0$  s.t. for  $f \in H_0^1(\Omega)$ ,

$$\|f\|_{L^2} \leq c \|\nabla f\|_{L^2}.$$

Therefore,

$$\|\nabla f\|_{L^2}^2 \leq \underbrace{\|f\|_{L^2}^2 + \|\nabla f\|_{L^2}^2}_{\|f\|_{H^1}^2} \leq (1 + c^2) \|\nabla f\|_{L^2}^2.$$

Hence, in  $H_0^1(\Omega)$ , the norm of  $\nabla f$  in  $L^2$  and of  $f$  in  $H^1$  are equivalent. This is also true if the function has only a realization of the boundary (with measure  $> 0$ ) where it vanishes.

**Remark.**

$$\alpha \|u\|_F \leq \|u\|_S \leq \beta \|u\|_F$$

implies that  $\|\cdot\|_F$  and  $\|\cdot\|_S$  are equivalent.

**Definition 1.2.30 (Dual of  $H_0^1$ :  $H^{-1}$ ).**

$$H_0^1 \subset L^2 \subset H^{-1}.$$

**Space-Time Functional Spaces** For time-dependent problems, we have spaces that describe the two dependencies. For instance:

$L^p(0, T; H^k(\Omega))$  is the space of functions s.t.  $\|f\|_{H^k}(t) \in L^p(0, T)$ .

In short, we denote it as  $L^p(H^k)$ . We will use  $L^2(H_0^1)$  and  $L^\infty(L^2)$ .

**Theorem 1.2.31 Embedding Theorems**

If  $u \in H^s$  and  $\Omega \subset \mathbb{R}^n$ ,

- $0 < 2s < n$ :  $H^s \subset L^q$ , where  $1 \leq q \leq q^* = \frac{2n}{n-2s}$  [continuous embedding]
- $2s = n$ :  $H^s \subset L^q \quad \forall q : 1 \leq q < +\infty$
- $2s > n$ :  $H^s \subset C^0(\overline{\Omega})$ .

Also,  $H^s \subset C^n(\overline{\Omega})$  if  $s > n + \frac{n}{2}$ .

**Remark.**  $H_1^1(\Omega) = \{f \in H^1, f(\partial\Omega) = 1\}$  is not a space. In fact:

$$f_1, f_2 \in H_1^1 \implies f_1 + f_2 \notin H_1^1.$$

### 1.3 Solving $Ax = b$

#### Direction Method

- Gaussian Elimination:  $A = LU$ .

$$Ax = b \implies LUx = b \implies \begin{cases} Ly = b \\ Ux = y \end{cases}$$

- With pivoting:  $PAQ = LU$ .

$$Ax = b \implies PAx = Pb \implies PAQ \underbrace{Q^{-1}x}_z = Pb \implies \begin{cases} LUz = Pb = c \\ x = Qz \end{cases}$$

- Problem:
  1. Complexity  $\sim \mathcal{O}(n^3)$
  2. Rounding errors

#### Theorem 1.3.1 Perturbation Theorem

Suppose a perturbation  $(A + E)(x + \delta x) = b + \delta$ . Then, the error

$$\frac{\|\delta x\|}{\|x\|} \leq \chi(A) \left( \frac{\|\delta\|}{\|b\|} + \frac{\|E\|}{\|A\|} \right),$$

where  $\chi(A) = \|A\| \|A^{-1}\| \geq 1$  is the condition number of  $A$ .

3. Storage

#### Iterative Methods

$$x^{(k)} \longrightarrow x_{\text{exact}}$$

#### Example 1.3.2

Given  $Ax = b$  and initial guess  $x^{(0)}$ :

$$x^{(k+1)} = \underbrace{B_k x^{(k)}}_{\sim \mathcal{O}(n^2)} + c_k$$

- Total cost:  $\sim \mathcal{O}(kn^2)$

- If  $k \ll n$ , iterative methods will be faster.
- Such methods converge  $\iff \rho(B) < 1$ , where  $\rho(B)$  is the spectral radius of  $B$ .

- Richardson method:

$$\begin{aligned}\mathbf{0} &= \mathbf{x} - \mathbf{x} = \mathbf{b} - A\mathbf{x} \\ \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} &= \mathbf{b} - A\mathbf{x}^{(k)} \\ \implies \mathbf{x}^{(k+1)} &= (I - A)\mathbf{x}^{(k)} + \mathbf{b}.\end{aligned}$$

Convergent when  $\rho(I - A) < 1$ .

- What if  $\rho(I - A) \geq 1$ ? If  $S$  is SPD, we can introduce a parameter  $\sigma$ :

$$\begin{aligned}\mathbf{0} &= \mathbf{x} - \mathbf{x} = (\mathbf{b} - A\mathbf{x})\sigma \\ \mathbf{x}^{(k+1)} &= (I - \sigma A)\mathbf{x}^{(k)} + \sigma\mathbf{b}.\end{aligned}$$

In fact, choosing  $\sigma = \frac{1}{\lambda_{\min} + \lambda_{\max}}$  yields the best convergence. More generally, choosing  $0 < \sigma < \frac{2}{\lambda_{\max}}$  gives convergence.

- Preconditioning:

$$\begin{aligned}\mathbf{0} - \mathbf{x} - \mathbf{x} &= P^{-1}(\mathbf{b} - A\mathbf{x}) \\ \mathbf{x}^{(k+1)} &= (I - P^{-1}A)\mathbf{x}^{(k)} + P^{-1}\mathbf{b}\end{aligned}$$

We require  $\rho(I - P^{-1}A) < 1$ . Ideally, we should choose  $P = A$ , but it is not practical. In practice, we want some  $P \approx A$ . For example,

1.  $P = \text{diag}(\text{diag}(A))$  yields the *Jacobi* method.
2.  $P = \text{lower}(A)$  yields the *Gauss-Seidel* method.

- Error: Richardson with Optimal  $\sigma$  for SPD  $A$ :

$$\|\mathbf{e}^{(k+1)}\| \leq \frac{\chi_2(A) - 1}{\chi_2(A) + 1} \|\mathbf{e}^{(k)}\|,$$

where  $\mathbf{e}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}_{\text{exact}}$ .

- If  $A$  is SPD,  $A\mathbf{x} = \mathbf{b}$  can be rewritten as a LS problem:

$$\min J \equiv \frac{1}{2} \mathbf{x}^\top A \mathbf{x} - \mathbf{x}^\top \mathbf{b}.$$

By FOC:  $\nabla J = 0$ . So, we can use methods from optimization:

1. Gradient descent, and

## 2. Conjugate gradient.

- If  $A$  is not SPD,  $A^\top A$  is SPD, but  $\chi(A^\top A) = (\chi(A))^2$ . If we want to use CG, use BiCG instead.
- Krylov subspace methods: GMRES.

## 1.4 General Statements of PDEs

**Definition 1.4.1 ( $n$ -th Order PDE).**

$$\mathcal{F}\left(u, \frac{\partial u}{\partial x_i}, \frac{\partial^2 u}{\partial x_i \partial x_j}, \dots, \frac{\partial^n u}{\partial^{i_1} x_1 \partial^{i_2} x_2 \dots \partial^{i_d} x_d}, f\right) = 0,$$

where  $i_1 + i_2 + \dots + i_d = n$ , and  $d = \text{number of independent variables}$ .

**Definition 1.4.2 (Linearity of PDEs).**

- Linear: all derivatives have linear occurrences.
- Quasi-Linear: linear in derivatives of order  $n$ . For example,

$$\frac{\partial u}{\partial t} + (1 - 2u) \frac{\partial u}{\partial x} = 0$$

- Semi-Linear: derivatives are linear, and coefficients of derivatives do not depend on lower-dimensional derivatives. For instance,

$$\frac{\partial u}{\partial t} + \sin(x) \frac{\partial u}{\partial x} + u^3 = 0$$

**Definition 1.4.3 (Boundary Value Problems/BVP).**  $B(u)|_{\partial\Omega} = \text{data}$  is the boundary conditions. Then,

$$\begin{cases} \mathcal{F}(u, \dots, f) \equiv 0 \\ B(u)|_{\partial\Omega} = \text{data} \end{cases}$$

defines a BVP. An *initial boundary value problem (IBVP)* is a BVP with one dimension being the time dimension.

**Definition 1.4.4 (Well-Posedness).** Our problem is *well-posed (WP)* when the solution

- exists,
- is unique, and
- depends continuously on the data [*Suppose  $\mathcal{F}(u, \dots, f) = 0$ . Perturb the problem:  $\mathcal{F}(u^*, \dots, f + \delta f) = 0$ . Does  $u^* \rightarrow u$  as  $\delta f \rightarrow 0$ ?*]

## 2 Poisson Equation

$$\begin{cases} -\Delta u = f & \text{and some BCs} \\ \text{OR: } -\frac{\partial^2 u}{\partial x_1^2} - \frac{\partial^2 u}{\partial x_2^2} - \frac{\partial^2 u}{\partial x_3^2} = f. \end{cases}$$

Boundary Conditions:

- Dirichlet  $u(\partial\Omega) = g$ .

With lifting function, we can turn Dirichlet into a homogeneous condition: Let  $\mathcal{L}_g \in \Omega$  s.t.  $\mathcal{L}_g(\partial\Omega) = g$ . Then,

$$u = \mathring{u} + \mathcal{L}_g \implies \begin{cases} -\Delta \mathring{u} = f + \Delta \mathcal{L}_g \\ \mathring{u}(\partial\Omega) = 0 \end{cases}$$

So, only homogeneous Dirichlet BC will be discussed.

- Neumann:  $\nabla u \cdot \mathbf{n} = g$ .

Poisson equations with Neumann BCs are not WP. The solutions are not unique by adding constants.

- Robin:  $\nabla u \cdot \mathbf{n} + \alpha u(\partial\Omega) = g$ .

### 2.1 Poisson Equation as a Minimization Problem

Consider the energy:

$$J(u) = \frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} f u$$

with boundary condition  $u(\partial\Omega) = 0$ .

**Goal:** Find the configuration ( $u$ ) that minimizes the energy. i.e.,

$$\min J(u).$$

Solve  $\nabla J = 0$ . Let  $v$  be a function such that  $v(\partial\Omega) = 0$ . Then,

$$\begin{aligned} \nabla J &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} (J(u + \varepsilon v) - J(u)) \\ &= \lim_{\varepsilon \rightarrow 0} \frac{1}{\varepsilon} \cdot \frac{1}{2} \int_{\Omega} |\nabla u|^2 + \varepsilon^2 |\nabla v|^2 + 2\varepsilon \nabla u \cdot \nabla v - \int_{\Omega} f u - \varepsilon \int_{\Omega} f v + \frac{1}{2} |\nabla u|^2 + \int_{\Omega} f u \\ &= \lim_{\varepsilon \rightarrow 0} \frac{1}{2} \int_{\Omega} \varepsilon |\nabla v|^2 + 2 \nabla u \cdot \nabla v - \int_{\Omega} f v \\ &= \int_{\Omega} \nabla u \cdot \nabla v - \int_{\Omega} f v = 0 \quad \forall v \text{ s.t. } v(\partial\Omega) = 0. \end{aligned}$$

By Greens' Theorem and Gauss Theorem,

$$\int_{\partial\Omega} \mathbf{v} \cdot \mathbf{n} = \int_{\Omega} \nabla \cdot (\mathbf{v}) \quad \text{and} \quad \int_{\Omega} \nabla \cdot (\beta \rho) = \int_{\Omega} \nabla \rho \cdot \beta + \int_{\Omega} \nabla \cdot \beta \cdot \rho.$$

So,

$$\int_{\Omega} \nabla \rho \cdot \beta = \int_{\partial\Omega} \rho \beta \cdot \mathbf{n} - \int (\nabla \cdot \beta) \cdot \rho.$$

Substitute  $\rho = v$  and  $\beta = \nabla u$ :

$$\int_{\Omega} \nabla u \cdot \nabla v = \underbrace{\int_{\partial\Omega} v \nabla u \cdot \mathbf{n}}_{=0 \text{ b/c } v(\partial\Omega)=0} - \int_{\Omega} (\Delta u)v.$$

Hence, we have

$$\begin{aligned} - \int_{\Omega} (\Delta u)v - \int_{\Omega} f v &= 0 \\ \int_{\Omega} (-\Delta u - f)v &= 0 \\ \boxed{-\Delta u = f}. \end{aligned}$$

### Remark.

- Strong formulation:

$$-\Delta u = f$$

requires second derivatives (Laplace operator)

- Weak formulation:

$$\int_{\Omega} \nabla u \nabla v - \int_{\Omega} f v = 0 \quad \forall v \text{ s.t. } v(\partial\Omega) = 0$$

only needs first-order derivatives (Gradient).

If  $u$  solves  $\int_{\Omega} \nabla u \nabla v - \int_{\Omega} f v = 0 \quad \forall v \text{ s.t. } v(\partial\Omega) = 0$ , then  $J(u) \leq J(u + w)$ .

## 2.2 Finite Differences

Two mathematically equivalent strong formulations may be different numerically:

$$-\Delta u = f.$$

On the other hand, we can write

$$-\nabla \cdot (\nabla u) = f \implies \begin{cases} \nabla u = \sigma \\ \nabla \cdot \sigma = f. \end{cases}$$

## 2.2.1 Poisson Equation in 1D

$$\begin{cases} -u'' = f, & x \in [0, 1] \\ u(0) = u(1) = 0 \end{cases}$$

- Max-Min Principle
- Discretization:

$$-u''(x_j) = f(x_j).$$

Collocation/Mesh.  $|h|$  is the step size.

- **Claim**

$$(u')' \approx \frac{u(x_{j+1}) - 2u(x_j) + u(x_{j-1}))}{h^2}$$

**Proof 1.** Denote  $u(x_{j+1}) = u_{j+1}$ . By Taylor expansion:

$$u_{j+1} = u_j + u'(x_j)h + \frac{1}{2}u''(x_j)h^2 + \frac{1}{3!}u'''(x_j)h^3 + \frac{1}{4!}u^{(4)}(x_j)h^4 + \mathcal{O}(h^5)$$

$$u_{j-1} = u_j - u'(x_j)h + \frac{1}{2}u''(x_j)h^2 - \frac{1}{3!}u'''(x_j)h^3 + \frac{1}{4!}u^{(4)}(x_j)h^4 + \mathcal{O}(h^5)$$

$$u_{j+1} + u_{j-1} = 2u_j + u''(x_j)h^2 + \frac{2}{4!}u^{(4)}h^4 + \mathcal{O}(h^6)$$

$$u''(x_j) = \frac{u_{j+1} + u_{j-1} - 2u_j}{h^2} + \underbrace{\frac{1}{12}u^{(4)}(x_j)h^2}_{\tau(h), \text{ consistency error}}$$

Hence, when we make the mesh finer ( $h \rightarrow \frac{1}{2}h$ ), the consistency error decrease by 4 times (order of 2). Q.E.D. ■

- Form the numerical problem: Except for the endpoints, we can write

$$-\frac{u_{j+1} - 2u_j + u_{j-1}}{h^2} = f_j.$$

Denote

$$\mathbf{u} = \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_{N-1} \end{bmatrix}, \quad \mathbf{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_{N-1} \end{bmatrix}.$$

Define

$$A = \frac{1}{h^2} \begin{bmatrix} 2 & -1 & & 0 \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & -1 \\ 0 & & -1 & 2 \end{bmatrix}.$$

Then, we get  $\boxed{A\mathbf{u} = \mathbf{f}}$ .

- **Proposition**  $A$  is SPD (symmetric positive definite)  $\implies$  all eigenvalues of  $A$  are positive.

**Proof 2.** We will show  $A$  is symmetric, and  $A$  is positive definite.

1. Symmetry is trivial.
2. Positivity: [WTS:  $\forall \mathbf{x} \neq 0, \mathbf{x}^\top A \mathbf{x} > 0.$ ]

$$\begin{aligned} \mathbf{x}^\top A \mathbf{x} &= \frac{1}{h^2} (2x_1^2 - 2x_1x_2 + x_2^2 + \dots) \\ &= \frac{1}{h^2} (x_1^2 + (x_1 - x_2)^2 + \dots + (x_j - x_{j+1})^2 + \dots + x_{N-1}^2) \\ &\geq 0. \end{aligned}$$

If we have a constant vector:  $\mathbf{x} = [x_* \dots x_*]^\top$ . Then,

$$\mathbf{x}^\top A \mathbf{x} = \frac{1}{2} (\mathbf{x}_*^2 + \mathbf{x}_*^2) = 0 \iff \mathbf{x}_* = 0.$$

Hence,  $\forall \mathbf{x} \neq 0, \mathbf{x}^\top A \mathbf{x} \geq 0$ .

So,  $A$  is SPD.

Q.E.D. ■

- **Corollary**  $A$  is non-singular, and the numerical problem is well-posed.
- Error Analysis:

$$\begin{aligned} \mathbf{u}_{\text{num}} &= A^{-1} \mathbf{f} \implies \|\mathbf{u}_{\text{num}}\| \leq \|A^{-1}\| \|\mathbf{f}\| \\ A \mathbf{u}_{\text{ex}} &= \mathbf{f} + \boldsymbol{\tau}, \quad [\tau_j] = [cf''(x_j)h^2] \\ \mathbf{e} &= \mathbf{u}_{\text{ex}} - \mathbf{u}_{\text{num}} = A^{-1} \boldsymbol{\tau} \\ \|\mathbf{e}\| &\leq \|A^{-1}\| \underbrace{\|\boldsymbol{\tau}\|}_{\mathcal{O}(h^2)} \end{aligned}$$

But this is not enough:  $A^{-1}$  also depends on  $h$ .

**Proposition**  $\|A^{-1}\|_2 \propto h$ .

**Proof 3.** Recall:  $\|A\|_2 = \rho(A) = \max_i |\lambda_i|$ , the spectral radius. Then,

$$\|A^{-1}\|_2 = \frac{1}{\min_i |\lambda_i|} = \rho(A^{-1}).$$

[WTS: eigenvalues of  $A$  (or  $A^{-1}$ ) are independent of  $h$ ].

$$A = \frac{1}{h^2} \text{trid}(-1, 2, -1), \quad T = \text{trid}(-1, 2, -1).$$

Solve  $T\mathbf{x} = \lambda\mathbf{x}$ :

$$(T - \lambda I)\mathbf{x} = 0$$

$$x_{j-1} - (2 - \lambda)x_j + x_{j+1} = 0.$$

Characteristic equation:

$$1 - (2 - \lambda)\rho^2 + \rho^2 = 0$$

$$x_j = a\rho_1^j + b\rho_2^j \quad \text{or} \quad x_j = a\rho^j + jb\rho^j.$$

Imposing  $\Delta \geq 0$ , we get  $\lambda \geq 4$ . Imposing BC:  $x_0 = x_N = 0$ , we get  $a = b = 0$ , and thus  $\mathbf{x} = 0$ . Trivial.

Hence, we don't have eigenvalues  $\geq 4$ , and we should only consider complex solutions.

$$\rho_{1,2} = \frac{2 - \lambda \pm i\sqrt{2}\sqrt{4 - \lambda}}{2} = \cos(\theta) \pm i\sin(\theta) \quad [\text{Enforcing } |\rho| = 1]$$

So,

$$x_j = a\rho_1^j + b\rho_2^j.$$

Since  $x_0 = x_N = 0$ :

$$\begin{cases} a + b = 0 \\ a\rho_1^N + B\rho_2^N = 0 \end{cases} \implies \rho_1^N - \bar{\rho}_1^N = 0.$$

$$\cos(N\theta) + i\sin(N\theta) - \cos(N\theta) + i\sin(N\theta) = 0 \implies \theta = \frac{j\pi}{N}, \quad j = 1, \dots, N - 1.$$

Going back, we have

$$\lambda_j = -\frac{1}{\rho} + \rho + 2 = 2\left(1 - \cos\left(\frac{j\pi}{N}\right)\right)$$

$$\cos(2\alpha) = 1 - 2\sin^2(\alpha)$$

$$2\sin^2(\alpha) = 1 - \cos(2\alpha) \implies \lambda_j = 4\sin^2\left(\frac{j\pi}{2N}\right)$$

Hence, eigenvalues of  $A$ :

$$\lambda = \frac{4}{h^2} \sin^2\left(\frac{\pi}{2N}\right) = \frac{4}{h^2} \sin^2\left(\frac{\pi}{2}h\right) \quad [h = \frac{1}{N}]$$

When  $h \rightarrow 0$ ,  $\sin^2\left(\frac{\pi}{2}h\right) \rightarrow \left(\frac{\pi}{2}h\right)^2$ . So,

$$\lambda = \frac{4}{h^2} \left(\frac{\pi}{2}\right)^2 h^2 = \pi^2 \ll h.$$

Q.E.D. ■

**Corollary** *Our numerical method is convergent.*

$$\mathbf{e} = A^{-1}\boldsymbol{\tau}.$$

When  $h \rightarrow 0$ ,  $\boldsymbol{\tau} \rightarrow 0$ ,  $\|A^{-1}\|_2 \rightarrow \pi^2$ . Then,

$$\begin{aligned} \|\mathbf{e}\|_2 &\leq \|A^{-1}\|_2 \|\boldsymbol{\tau}\|_2 = \pi^2 \|\boldsymbol{\tau}\|_2. \\ \|\boldsymbol{\tau}\|_2^2 &= \sum_{i=1}^{N-1} \tau_i^2 = \sum_{i=1}^{N-1} c^2 (f''(x_j))^2 h^4. \end{aligned}$$

If  $f''$  is bounded,  $\|f''\| \leq \hat{c}$ . Then, combining the constants and the sum, we have

$$\|\boldsymbol{\tau}\|_2^2 \leq \tilde{c}(N-1)h^4.$$

Recall that  $N = \frac{1}{h}$ . Then,  $\|\boldsymbol{\tau}\|_2^2 \leq \tilde{c}h^3 \implies \|\boldsymbol{\tau}\|_2 \sim \mathcal{O}(h^{3/2})$ . So,

$$\|\mathbf{e}\|_2 \sim \mathcal{O}(h^{3/2}).$$

- **Practical Notes:**

When constructing  $A$ , we first assemble the full matrix, and then replace the boundaries with actual BCs.

- **Another perspective of FD Convergence.**

**Definition 2.2.1 ( $h$ -Norm).** Let  $v_h$  be a grid function.

$$\|v_h\|_h^2 = \left( \sum_{j=1}^{N-1} v_j^2 + \frac{1}{2}v_0^2 + \frac{1}{2}v_N^2 \right) h,$$

and

$$L_h v_h = \frac{v_{j+1} - 2v_j + v_{j-1}}{h^2}.$$

Some properties include:

1.  $(L_h v_h, w_h) = (v_h, L_h w_h)$
2.  $(L_h v_h, w_h) = \sum_{j=1}^{N-1} \frac{v_{j+1} - v_j}{h} \cdot \frac{w_{j+1} - w_j}{h}$ .
3.  $h(L_h v_h, v_h)$  is a norm. Denote this norm as  $\|\cdot\|_h$ .

**Proof 4.**

$$h(L_h v_h, v_h) = h \sum \left( \frac{v_{j+1} - v_j}{h} \right)^2.$$

Non-negativity and other properties can be easily shown.

Q.E.D. ■

$$4. c\|v_h\|_h^2 \leq \|v_h\|_h^2.$$

$$5. \text{ Problem: } L_h u_h = f_h.$$

$$c\|u_h\|_h^2 \leq h(L_h u_h, u_h) \leq \|f_h\|_h \|u_h\|_h \implies \|u_h\|_h \leq c\|f_h\|_h$$

$$c\|\mathbf{e}_h\|_h^2 \leq h(L_h \mathbf{e}_h, \mathbf{e}_h) \leq \|\boldsymbol{\tau}\|_h \|\mathbf{e}_h\|_h \implies \|\mathbf{e}_h\|_h \leq \frac{1}{c}\|\boldsymbol{\tau}\|_h$$

$$\|\boldsymbol{\tau}\|_h^2 = h \left| \sum_{j=1}^N c_j h^4 \right| \leq h \cdot N \cdot \max_{j=1, \dots, N} c_j \cdot h^4$$

$$= h \cdot \frac{1}{h} \max_{j=1, \dots, N} c_j \cdot h^4$$

$$= h^4 \max_{j=1, \dots, N} c_j$$

$$\|\boldsymbol{\tau}\|_h \leq \sqrt{h^4 \max_{j=1, \dots, N} c_j} = h^2 \sqrt{\max_{j=1, \dots, N} c_j} \sim \mathcal{O}(h^2).$$

### Remark 5. (What's Missing Here?).

$$\text{Differential Problem} \xrightarrow{\textcircled{1}} \underbrace{\mathbf{A}\mathbf{u} = \mathbf{b}}_{\text{Discretized linear system}} \xrightarrow{\textcircled{2}} \mathbf{u}_{\text{num}} \approx \mathbf{A}^{-1}\mathbf{b}.$$

1. In  $\textcircled{1}$ : Convergence analysis

2. In  $\textcircled{2}$ : numerical linear algebra errors and rounding errors. (We were not discussing these previously).

## 2.3 Galerkin Methods

### Set-Up

$$\begin{cases} -\Delta u = f \\ u(\partial\Omega) = 0 \end{cases}$$

Weak formulation: Find  $u$  s.t.

$$\int_{\Omega} \nabla u \nabla v \, dx = \int_{\Omega} f v \, dx \quad \forall v \text{ s.t. } v(\partial\Omega) = 0.$$

**Requirements**  $\nabla u \in L^2, \nabla v \in L^2 \implies u, v \in H_0^1$ , and  $f \in H^{-1}$  (dual of  $H_0^1$ , functional applied to  $H_0^1$ .)

**Lemma 2.1 Max-Milgram Lemma (Sufficient Condition)** Consider the problem  $a(u, v) = \mathcal{F}(v)$ , where  $u, v \in X$ . If

- $a(\cdot, \cdot)$  is a bilinear form,
- **bounded:**  $\exists M > 0$  s.t.  $|a(u, v)| \leq M\|u\|_X\|v\|_X$ ,

- coercivity:  $a(u, u) \geq \alpha \|u\|_X^2 \quad \forall u \in X$ , and
- $\mathcal{F}(\cdot)$  is a continuous linear functional:  $\|\mathcal{F}(v)\|_X \leq \gamma \|v\|_X \quad \forall v \in X$ ,

then, the problem is well-posed. i.e., ①②③④  $\implies$  well-posedness (sufficient but not necessary).

**Theorem 2.3.2**

Poisson equation is well-posed.

**Proof 1.**

$$\begin{cases} -\Delta u = f \\ u(\partial\Omega) = 0 \end{cases}$$

$\implies \forall v \in H_0^1$ , find  $u \in H_0^1$  s.t.

$$\underbrace{\int_{\Omega} \nabla u \nabla v}_{a(u,v)} - \underbrace{\int_{\Omega} f v}_{\mathcal{F}(v)} = 0.$$

- $\|u\|_{H^1}^2 = \|u\|_{L^2}^2 + \|\nabla u\|_{L^2}^2$
- $\left| \int_{\Omega} f v \right| \leq \|f\|_{L^2} \|v\|_{L^2} \leq \|f\|_{H^1} \|v\|_{H^1}$
- $\forall u \in X = H_0^1$ ,

$$\int_{\Omega} |\nabla u|^2 = a(u, u) \geq \alpha \|u\|_X^2.$$

By Poincaré Inequality,

$$\begin{aligned} C_p \int_{\Omega} u^2 &\leq \int_{\Omega} |\nabla u|^2 \\ C_p \|u\|_{L^2}^2 &\leq \|\nabla u\|_{L^2}^2 = \underbrace{\|u\|_{H^1}^2}_{\text{semi-norm}} \\ \|u\|_{H^1}^2 &= \int_{\Omega} u^2 + \int_{\Omega} \nabla u^2 \leq \left( \frac{1}{C_p} + 1 \right) \int_{\Omega} |\nabla u|^2 \end{aligned}$$

So, we have coercivity.

By Lax-Milgram Lemma, Poisson equation is WP.

Q.E.D. ■

**Remark.**  $a(u, v) = \int_{\Omega} \nabla u \nabla v$  is a scalar product in  $H_0^1$ :

- symmetric, and
- non-negativity (by coercivity).

The norm induced by  $a(u, u) = \|u\|_a^2$  is equivalent to  $\|u\|_{H^1}^2$ :

$$C \|u\|_{H^1}^2 \leq a(u, u) \leq \|u\|_{H^1}^2.$$

If we change the problem and the space:

$$\begin{cases} -\Delta u = f \\ \nabla u \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \quad (\text{Neumann}) \end{cases}$$

$\Rightarrow$  Find  $u \in H^1$  s.t.  $\forall v \in H^1$ :

$$\int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v.$$

Now, we don't have Poincaré Inequality anymore, and  $a(u, v)$  is not coercive. We cannot apply LM Lemma. In fact, the problem is not WP. Any  $(u + \text{constant})$  will be a solution.

How to make the problem WP?

$$\begin{cases} -\Delta u + \sigma u = f, \quad \sigma > 0 \\ \nabla u \cdot \mathbf{n} = 0 \quad \text{on } \partial\Omega \end{cases}$$

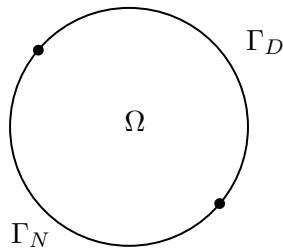
$$\Rightarrow \underbrace{\int_{\Omega} \nabla u \nabla v + \sigma \int_{\Omega} uv}_{a(u, v)} = \int_{\Omega} f v$$

$$|a(u, v)| \leq \max(\sigma, 1) \|u\|_{H^1} \|v\|_{H^1}$$

$$\min(1, \sigma) \|u\|_{H^1}^2 \leq a(u, v).$$

$\Rightarrow a(u, v)$  is coercive, and the problem is WP.

Now, let's consider mixed boundary conditions (i.e., Robin conditions).



$$\begin{cases} -\Delta u = f \\ u(\Gamma_D) = 0 \\ \nabla u \cdot \mathbf{n}|_{\Gamma_N} = 0, \end{cases}$$

where  $\overline{\Gamma_D \cup \Gamma_N} = \partial\Omega$ , and  $\text{measure}(\Gamma_D \cap \Gamma_N) = 0$ .

$$-\int_{\Omega} \Delta u \cdot v = -\int_{\Gamma_N} \underbrace{\nabla u \cdot \mathbf{n}}_{=0} v + \int_{\Gamma_D} \nabla u \cdot \mathbf{n} \underbrace{v}_{=0} + \int_{\Omega} \nabla u \nabla v.$$

**Problem:** Find  $u \in H_{\Gamma_D}^1$  s.t.  $\forall v \in H_{\Gamma_D}^1$  s.t.

$$\int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v.$$

This problem is WP because we have Poincaré on  $H_{\Gamma_D}^1$ .

Now, consider the true problem: Find  $u \in X = H_0^1$  s.t.

$$\int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v \quad \forall v \in X \implies a(u, v) = \mathcal{F}(v).$$

**Assumption**  $a(u, v) = a(v, u)$  (symmetry) and all LM assumptions.

The Galerkin Methods solve the numerical problem in  $X_N \subset X$ : Find  $u_N \in X_N$  s.t.

$$a(u_N, v_N) = \mathcal{F}(v_N), \quad \forall v_N \in X_N.$$

This problem is WP, which can be trivially proven by LM assumptions.

**Proposition 2.3** *This problem is strongly consistent (consistency error is 0).*

**Proof 2.** Note that  $a(u, v) = \mathcal{F}(v)$  and  $a(u_N, v_N) = \mathcal{F}(v_N)$ . Then,

$$a(u, v_N) = \mathcal{F}(v_N) \quad \text{because } X_N \subset X.$$

The consistence error is 0. So, the problem is strongly consistent. We could also show

$$a(u - u_N, v_N) = 0 \quad \forall v_N \in X_N.$$

Q.E.D. ■

**Proposition 2.4** *This problem is stable.*

**Proof 3.**

$$\alpha \|u_N\|_{H^1}^2 \leq a(u_N, u_N) \leq \|f\| \cdot \|u_N\|_{H^1}.$$

So, it is clearly stable.

Q.E.D. ■

**Lemma 2.5 C a Lemma**

$$\|u - u_N\|_{H^1} \leq K \inf_{w_N \in X_N} \|u - w_N\|_{H^1}.$$

**Proof 4.**

$$\alpha \|u - u_N\|_{H^1}^2 \leq a(u - u_N, u - u_N) = a(u - u_N, u - w_N) + a(u - u_N, w_N - u_N)$$

[Bilinearity,  $w_N \in X_N$ ]

$$\alpha \|u - u_N\|_{H^1}^2 \leq a(u - u_N, u - w_N) \leq M \|u - u_N\|_{H^1} \cdot \|u - w_N\|_{H^1}$$

[Continuity]

$$\|u - u_N\|_{H^1} \leq \frac{M}{\alpha} \|u - w_N\|_{H^1}$$

$$\|u - u_N\|_{H^1} \leq K \|u - w_N\|_{H^1} \quad \forall w_N \in X_N$$

$$\|u - u_N\|_{H^1} \leq K \inf_{w_N \in X_N} \|u - w_N\|_{H^1}.$$

Q.E.D. ■

This Lemma implies that Galerkin solution is not the best one, but its  $H_1$  norm error goes to 0 as quick as the best one. i.e., it has the same rate of convergence as the best approximation in  $X_N$ . When

$N \rightarrow \infty, X_N \rightarrow X$ , and  $\inf_{w_N \in X_n} \|u - w_N\|_{H^1} \rightarrow 0$ .

**Claim 2.6**  $P_N$  : Find  $u_N \in X_N$ ,  $a(u_N, v_N) = \mathcal{F}(v_N) \quad \forall v_N \in X_N$  is WP.

**Proof 5.** We can write  $u_N = \sum_{j=0}^{N-1} c_j \varphi_j(\mathbf{x})$ , where  $\varphi_j(\cdot)$  are basis functions. Then,  $X_n = \text{span} \{\varphi_j\}_{j=0}^{N-1}$ .

So,

$$\begin{aligned} a(u_N, v_N) &= a\left(\sum_{j=0}^{N-1} c_j \varphi_j(\mathbf{x}), v_N\right) = \mathcal{F}(v_N) \\ &= \sum_{j=0}^{N-1} c_j a(\varphi_j, v_N) = \mathcal{F}(v_N) \end{aligned} \quad [\text{Bilinearity}]$$

Instead of testing on  $v_N \in X_N$ , as  $X_N = \text{span} \{\varphi_j\}_{j=0}^{N-1}$ , we just need to prove it works for all  $\varphi_i$ 's.

$$\sum_{j=0}^{N-1} c_j a(\varphi_j, \varphi_i) = \mathcal{F}(\varphi_i) \quad [\text{Take } v_N = \varphi_i]$$

Define

$$A := \left[ a(\varphi_j, \varphi_i) \right]_{i,j}, \quad \mathbf{c} = \left[ c_j \right]_j, \quad \mathbf{b} = \left[ \mathcal{F}(\varphi_i) \right]_i.$$

So,  $P_N$  becomes

$$A\mathbf{c} = \mathbf{b}.$$

**Claim 2.7**  $A$  is non-singular

- $A$  is symmetric:  $a(\varphi_j, \varphi_i) = a(\varphi_i, \varphi_j)$  [symmetry]
- $A$  is positive-definite: Assume  $\mathbf{c} \neq 0$ , then

$$\begin{aligned} \mathbf{c}^\top A \mathbf{c} &= \sum_{i,j} c_j A_{ij} c_i = \sum_{i,j} c_j a(\varphi_j, \varphi_i) c_i \\ &= \sum_{i,j} a(c_j \varphi_j, c_i \varphi_i) \quad [\text{bilinearity}] \\ &= a\left(\sum_{j=0}^{N-1} c_j \varphi_j, \sum_{i=0}^{N-1} c_i \varphi_i\right) \quad [\text{bilinearity}] \\ &= a(u_N, u_N) \quad [u_N = \sum_{i=0}^{N-1} c_i \varphi_i] \\ &\geq \alpha \|u_N\|_{H^1}^2 \geq 0 \quad [\text{coercivity}] \end{aligned}$$

So,  $A$  is SPD, and thus  $A$  is non-singular.  $\implies P_N$  is WP.

Q.E.D. ■

**Remark.** What if the problem is non-linear?

$$-\Delta u + u^3 = f \implies \int_{\Omega} \nabla u \nabla v + \int_{\Omega} u^3 v = f.$$

- This problem is WP (although we don't have symmetry).
- Solving: linearize, then discretize:

$$\int_{\Omega} \nabla u^{(k+1)} \nabla v + \int_{\Omega} (u^{(k)})^2 u^{(k+1)} v = f$$

We can use root finding. For example, Newton's method.

- Solving: discretize, then linearize

$$u_N = \sum_{j=0}^{N-1} c_j \varphi_j, \quad X_N = \text{span} \{ \varphi_j \}_{j=0}^{N-1}.$$

Then, we form a non-linear equation  $F(c) = 0$ . We can also use Newton's method.

## 2.4 Finite Element: One Possible Choice of $X_N$

Recall: we need to compute

$$a(\varphi_j, \varphi_i) = \int_{\Omega} \nabla \varphi_j \nabla \varphi_i.$$

So, we want functions that are easy to differentiate and integrate.  $\implies$  **Polynomials!**

### Example 2.4.1 Runge Counterexample

Consider the function  $f = \frac{1}{x^2 + 1}$ ,  $x \in [-5, 5]$ .

If we interpolate  $f(x)$  with equispaced points and polynomials,

$$\inf \|u - w_N\| \xrightarrow{N \rightarrow \infty} 0$$

will not be held.

Solutions:

- Optimize the position of nodes: Gaussian interpolation, Chebyshev nodes.
- Piecewise interpolation: locally fit piecewise polynomial.

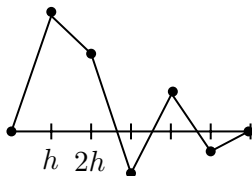
### 2.4.1 Finite Element in 1D

$$-u'' = f \implies \int_0^1 u' v' = \int_0^1 f v$$

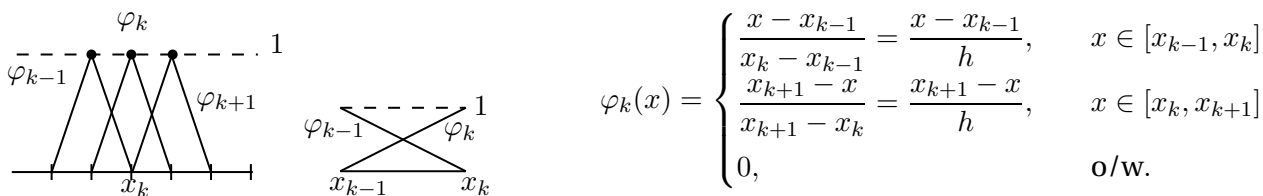
with  $u(0) = u(1) = v(0) = v(1) = 0$ , where

$$X_N = \{u_h \mid u_h \in C([0, 1]) \text{ and } u_h \in \mathbb{P}^1(I_h)\},$$

where  $\mathbb{P}^1$  denotes polynomial of order 1. So,  $X_N^1 \subset H_0^1(0, 1)$ , and  $\dim(X_N) = N - 1 = \frac{1}{h} - 1$ . So, when  $N \rightarrow +\infty, h \rightarrow 0$ , and  $\dim(X_N) \rightarrow +\infty$ .



When doing interpolation, let  $y = f(x)$  be the true function, and  $y_i = f(x_i)$  be the interpolant. There are many ways to construct the interpolant. Finite Element will use the *Lagrange Polynomial* (Hat functions):



Note that

$$\varphi_i(x_k) = \delta_{ik} = \begin{cases} 1, & i = k \\ 0, & i \neq k \end{cases}$$

is the Kronecker- $\delta$  function. So, if  $v_N = \sum_{i=0}^{N-1} y_i \varphi_i$ , we have

$$v_N(x_k) = \sum_{i=0}^{N-1} y_i \varphi_i(x_k) = y_k.$$

**FE In Practice** We will set everything up on a reference interval  $[0, 1]$  with  $\widehat{\varphi}_0 = 1 - \widehat{x}$  and  $\widehat{\varphi}_1 = \widehat{x}$ . We will use the bijection  $x = \text{map}(\widehat{x}) = a\widehat{x} + b$  such that  $\text{map}(0) = x_{j-1}$  and  $\text{map}(1) = x_j$ . Then,



Also,

$$\widehat{x} = \frac{x - b}{a} = \frac{x - x_{j-1}}{h}.$$

**Convergence Theorem and Its Application** Let  $u(x) \in C^0([0, 1])$ . Denote  $\Pi_u^1(x) \in X_h^1$  as

$$\Pi_u^1(x) = \{ \text{continuous function, linear s.t. on } I_k = [x_k, x_{k+1}], \Pi_u^1(I_k) \in \mathbb{P}^1 \},$$

where  $\mathbb{P}^1$  denotes polynomial of order 1. Then,

$$\Pi_u^1(x) = \sum_{j=0}^N u(x_j) \varphi_j(x_j).$$

To measure error, consider the following norms:

$$\|u - \Pi_u^1\|_{L^2} \quad \text{and} \quad \|\nabla(u - \Pi_u^1)\|_{L^2} = \underbrace{|u - \Pi_u^1|_{H^1}}_{\text{semi-norm}}$$

**Theorem 2.4.2**

If  $u \in H^2([0, 1])$ , then

$$\begin{aligned} \|u - \Pi_u^1\|_{L^2} &\leq C_0 |u|_{H^2} h^2 = C \|u''\|_{L^2} h^2 \\ |u - \Pi_u^1|_{H^1} &\leq C_1 \|u''\|_{L^2} h \\ \implies \|u - \Pi_u^1\|_{H^1} &\leq C \|u''\|_{L^2} h. \end{aligned}$$

**Proof 1.**  $u \in H^2([0, 1]) \implies u \in C^1([0, 1])$ .

- Let's work locally on the interval  $[x_j, x_{j+1}]$ . Denote that

$$e = u - \Pi_u^1, \quad e(x_j) = 0, \quad \text{and} \quad e(x_{j+1}) = 0.$$

By Rolle's Theorem,  $\exists \xi_j$  s.t.  $e'(\xi_j) = 0$ . By Fundamental Theorem of Calculus,

$$\begin{aligned} e'(x) &= \int_{\xi_j}^x e''(x) dx = \int_{\xi_j}^x u''(x) dx && [(\Pi_u^1)'' = 0] \\ |e'(x)| &\leq \left| \int_{x_j}^{x_{j+1}} u''(x) dx \right| \\ &\leq \left( \int_{x_j}^{x_{j+1}} 1^2 \right)^{1/2} \left( \int_{x_j}^{x_{j+1}} (u'')^2 \right)^{1/2} = h^{1/2} \|u''\|_{L^2(I_j)}. && [\text{Cauchy-Schwarz}] \\ \implies \|e'(x)\|_{L^2(I_j)}^2 &= \int_{x_j}^{x_{j+1}} |e'(x)|^2 \leq h \cdot h \|u''\|_{L^2(I_j)}^2 = h^2 \|u''\|_{L^2(I_j)}^2. \end{aligned}$$

Now, let's move to global:

$$\|e'\|_{L^2(I)}^2 = \sum_j \|e'\|_{L^2(I_j)}^2 \leq h^2 \|u''\|_{L^2(I)}^2.$$

$$\boxed{|e|_{H^1} = \|e'\|_{L^2} \leq Ch \|u''\|_{L^2} = Ch |u|_{H^2} \sim \mathcal{O}(h)}.$$

- Let's repeat the same argument:

$$|e(x)| = \left| \int_{x_j}^{x_{j+1}} e' dx \right| \leq \left( \int_{x_j}^{x_{j+1}} 1^2 \right)^{1/2} \left( \int_{x_j}^{x_{j+1}} (e')^2 \right)^{1/2}$$

$$e^2(x) = (e(x))^2 \leq h \|e'\|_{L^2(I_j)}^2$$

$$\int_{x_j}^{x_{j+1}} e^2(x) \leq h^2 \|e'\|_{L^2(I_j)}^2 \leq h^4 \|u''\|_{L^2(I_j)}^2 \quad [\text{Use conclusion from ①}]$$

Globally,

$$\|e\|_{L^2}^2 = \sum_j \|e\|_{L^2(I_j)}^2 \leq h^4 \|u''\|_{L^2(I)}^2$$

$$\Rightarrow \boxed{\|e\|_{L^2} \leq Ch^2 \|u''\|_{L^2} \sim \mathcal{O}(h^2)}.$$

Q.E.D. ■

Applications of Theorem 2.4.2:

- Consider FE problem  $a(u_h, v_h) = \mathcal{F}(v_h)$ . For  $u_h \in X_h^1$ , we have

$$\|u - u_h\|_{H^1} \leq \frac{M}{\alpha} \inf_{w_h \in X_h^1} \|u - w_h\|_{H^1} \quad [\text{Cea Lemma}]$$

$$\leq \frac{M}{\alpha} \|u - \Pi_{u,h}^1\|_{H^1} \quad [\text{Theorem 2.4.2}]$$

$$\leq \frac{M}{\alpha} h \|u''\|_{H^2}$$

$$\|u - u_h\|_{L^2} \leq \mathcal{O}(h^2).$$

- If we know the error is bad locally, we can locally refine  $h$  to get better error.

**Advection-Reaction Problem, Local Assembly, and Global Assembly**

$$\begin{cases} -u'' + \sigma u = f, & \sigma > 0 \\ u(0) = u(1) = 0. \end{cases}$$

Weak formulation:

$$\int_0^1 u'v' + \int_0^1 \sigma uv = \int_0^1 fv \quad \forall v \in H_0^1.$$

- Matris  $A$ :

$$A_{ij} = a(\varphi_i, \varphi_j) = \int_0^1 \varphi_j' \varphi_i' + \underbrace{\int_0^1 \sigma \varphi_j \varphi_i}_{\text{mass matrix}}$$

$$b_i = \int_0^1 f \varphi_i$$

How to integrate numerically? We do it interval by interval. Suppose we have first order Lagrange polynomials. Then,

$$\int_0^1 \sum_j \int_{I_j},$$

where

$$\int_{I_k} = \underbrace{\int_{I_k} \varphi_i' \varphi_j'}_{\text{constant}} + \int_{I_k} \sigma \varphi_i \varphi_j.$$

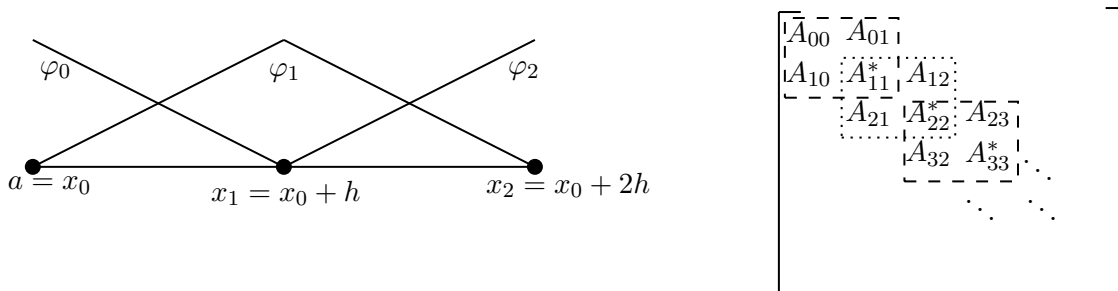
Note that

$$\int_{I_k} \varphi_\ell' \varphi_r' = 0 \quad \text{if } \ell, r \neq k, k+1.$$

So, the matrix will be tridiagonal, and

$$A_{ij} = \int_0^1 \varphi_j' \varphi_i' + \int_0^1 \sigma \varphi_j \varphi_i = \int_{x_{j-1}}^{x_{j+1}} \varphi_j' \varphi_i' + \int_{x_{j-1}}^{x_{j+1}} \varphi_j \varphi_i.$$

Let's only consider the first interval:



$$A_{00} = \int_{x_0}^{x_1} \varphi_0' \varphi_0' + \sigma \int_{x_0}^{x_1} \varphi_0 \varphi_0$$

$$A_{01} = \int_{x_0}^{x_1} \varphi_0' \varphi_1' + \sigma \int_{x_0}^{x_1} \varphi_0 \varphi_1$$

$$A_{11} = \int_{x_0}^{x_2} \varphi_1' \varphi_1' + \sigma \int_{x_0}^{x_2} \varphi_1 \varphi_1 = \underbrace{\int_{x_0}^{x_1} \varphi_1' \varphi_1' + \sigma \int_{x_0}^{x_1} \varphi_1 \varphi_1}_{=: A_{11}^*} + \underbrace{\int_{x_1}^{x_2} \varphi_1' \varphi_1' + \sigma \int_{x_1}^{x_2} \varphi_1 \varphi_1}_{\text{will be calculated in the next interval}}$$

$$A_{10} = \int_{x_0}^{x_1} \varphi_1' \varphi_0' + \sigma \int_{x_0}^{x_1} \varphi_1 \varphi_0 = A_{01}.$$

We will get the rest of  $A_{11}$  in the second interval.

- Local assembly: we will form the  $2 \times 2$  matrix locally:

$$a_{00} = a(\widehat{\varphi}_0, \widehat{\varphi}_0)_{\text{local}}$$

$$a_{10} = (\widehat{\varphi}_0, \widehat{\varphi}_1)_{\text{local}} = a_{01}$$

$$a_{11} = a(\widehat{\varphi}_1, \widehat{\varphi}_1)_{\text{local}}$$

We can compute the integrals on a reference interval  $[0, 1]$ :

$$a_{ii} = \int_{x_i}^{x_{i+1}} \varphi_i' \varphi_i' + \sigma \int_{x_i}^{x_{i+1}} \varphi_i \varphi_i$$

Consider the map:  $x = h\widehat{x} + x_i \in [0, 1]$ .

$$\frac{\partial \varphi_i}{\partial x} = \frac{\partial \widehat{\varphi}_i}{\partial \widehat{x}} \cdot \frac{\partial \widehat{x}}{\partial x}, \quad \frac{dx}{d\widehat{x}} = h.$$

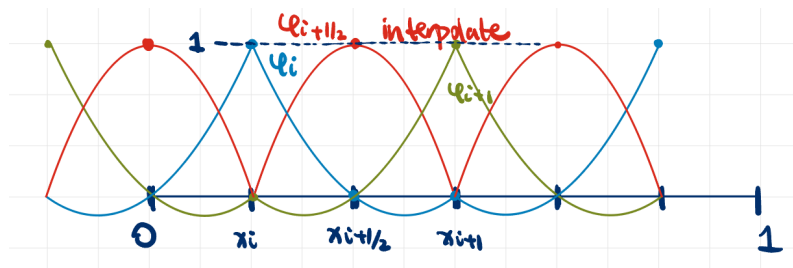
Then,

$$a_{ii} = \int_0^1 \frac{\partial \widehat{\varphi}_i}{\partial \widehat{x}} \cdot \frac{\partial \widehat{x}}{\partial x} \cdot \frac{\partial \widehat{\varphi}_i}{\partial \widehat{x}} \cdot \frac{\partial \widehat{x}}{\partial x} h d\widehat{x} + \sigma \int_0^1 \widehat{\varphi}_i \widehat{\varphi}_i h d\widehat{x}.$$

Now, everything is pre-computable. To compute these integrals efficiently, we will use *Gaussian Quadrature*.

- Global assembly: additivity of integrals.
- Treatment of the boundary conditions: after global assembly, we will impose the BCs.

### Use Quadratic Functions Instead



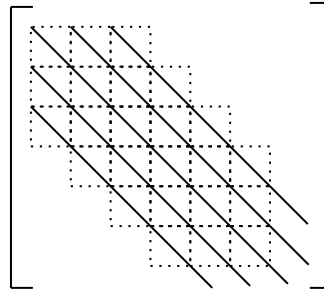
$$\varphi_i = \frac{(x - x_{i+1})(x - x_{i+1/2})}{(x_i - x_{i+1})(x_i - x_{i+1/2})}$$

$$\varphi_{i+1/2} = \frac{(x - x_i)(x - x_{i+1})}{(x_{i+1/2} - x_i)(x_{i+1/2} - x_{i+1})}$$

$$\varphi_{i+1} = \frac{(x - x_i)(x - x_{i+1/2})}{(x_{i+1} - x_i)(x_{i+1} - x_{i+1/2})}$$

- Local assembly:  $3 \times 3$  matrices.

- Global assembly: pena-diagonal (five diagonal entries):



### Convergence and Error

$$\mathbb{P}^1 : e_{H^1} \sim \mathcal{O}(h) \quad e_{L^2} \sim \mathcal{O}(h^2)$$

$$\mathbb{P}^2 : e_{H^1} \sim \mathcal{O}(h^2) \quad e_{L^2} \sim \mathcal{O}(h^3)$$

$$\mathbb{P}^3 : e_{H^1} \sim \mathcal{O}(h^3) \quad e_{L^2} \sim \mathcal{O}(h^4)$$

#### Theorem 2.4.3 Convergence of FE in 1D

Suppose  $u \in H^{s+1}(0, 1)$  with  $s \geq 1$ . For FE, we work in  $X_h^q$ . Then,

$$\|e\|_{H^1} \leq C \|u\|_{H^{\min(q,s)}} h^{\min(q,s)}$$

$$\|e\|_{L^2} \leq C \|u\|_{H^{\min(q,s)}} h^{\min(q,s)+1}.$$

**Remark.** This implies that increasing the degree of polynomials in FE makes sense only when the solution is regular enough.

Table 1: Summary of Rate of Convergence in  $H^1$ , Given Different  $q$  and  $s$

	$s = 1$	2	3	4	5	...
$q = 1$	1	1	1	1	1	
2	1	2	2	2	2	
3	1	2	3	3	3	
4	1	2	3	4	4	
$\vdots$						

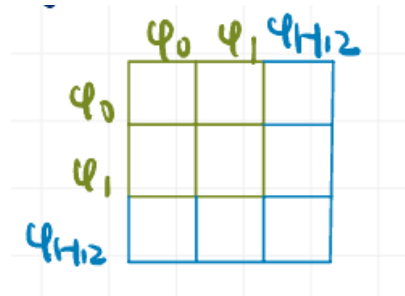
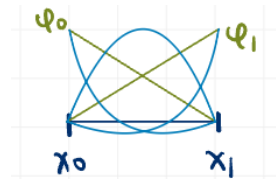
*Optimal choices:* with the regularity we have ( $s$ ), take the minimal optimal polynomial degree ( $q$ ). There are two ways to reduce error: (1) reduce  $h$ , and (2) increase  $q$ .

### Reusing Matrices

**Set-up** We have built linear FE matrices

$$A_L \mathbf{u}_L = \mathbf{b}_L.$$

**Goal** Set up quadratic FE ( $A_Q u_Q = b_Q$ ) from reusing  $A_L$

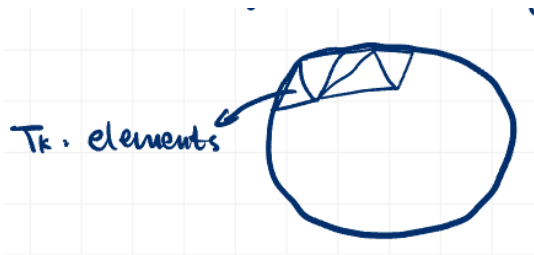


$$\begin{aligned} \varphi_0 &= 1 - \hat{x} \\ \varphi_1 &= \hat{x} \end{aligned} \implies \varphi_{H,2} = \varphi_0 \varphi_1 = (1 - \hat{x})\hat{x}$$

**Conclusion: What is a Finite Element?** We have an interval  $I$ . A finite element is built on  $K \in \mathcal{P}(I)$ , some partition of  $I$ . On each  $K$ , we build basis polynomials of degree  $q$ . i.e.,  $u(K) \in \mathbb{P}^q$ .

**2.4.2 Finite Elements in Multiple Dimension**

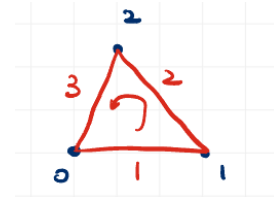
Let's divide the region into triangles:



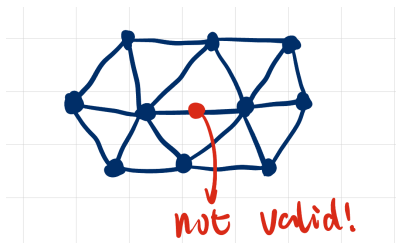
$$\Omega = \bigcup_{k=1}^N T_k,$$

though this relationship could be false in reality.

Triangle = Vertex + Edge



**Requirement** Vertices are not on edges of another triangle

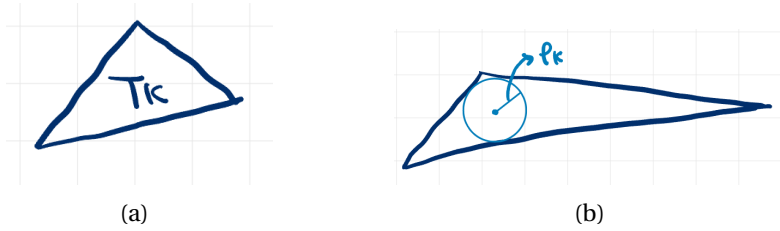


**Property** If  $v \in C^0(\bar{\Omega})$  and  $\in H^k(T_k)$ , then  $v \in H^k(\Omega)$

**How large is the mesh?** We have two measurements:

- $h = \max |x_1 - x_2|$  for  $x_1, x_2 \in T_k$

- $\frac{h_k}{\rho_k} < \delta$ , independent of  $k$ .



Recall the problem:

$$\begin{cases} -\Delta u = f \\ u(\partial\Omega) = 0 \end{cases} \implies \int \nabla u \nabla \varphi = \int f \varphi,$$

we aim to build

$$\Omega = \bigcup_{k=1}^N T_k \quad \text{and} \quad X_k^h = \left\{ v_n \in H^1 : v_h(T_k) \in \mathbb{P}^k \right\}.$$

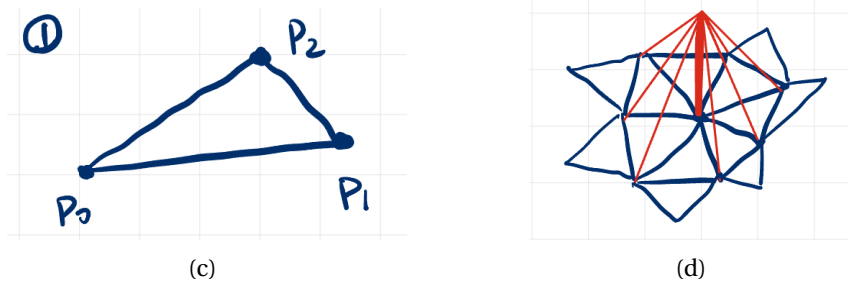
**Linear Polynomial**

$$ax + by + c \quad \left( \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} \right).$$

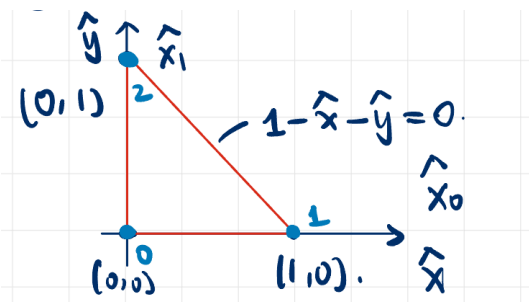
We could rewrite

$$P_k(\mathbf{x}) - P_k(P_j) = c_j,$$

which will give us a “tent” function.

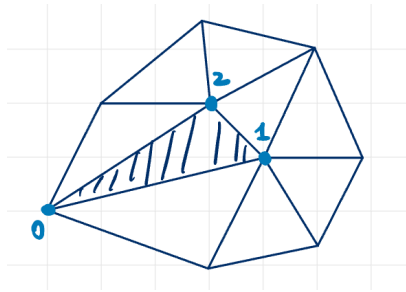


We will work on a reference interval:



$$\begin{aligned} \hat{\varphi}_0 &= 1 - \hat{x} - \hat{y} \\ \hat{\varphi}_1 &= \hat{x} \\ \hat{\varphi}_2 &= \hat{y}. \end{aligned}$$

$$\mathbf{x} = B_k \hat{\mathbf{x}} + \mathbf{c}_k$$



$$x(\hat{x}, \hat{y}) = \begin{cases} x(0, 0) = x_0 \\ x(1, 0) = x_1 \\ x(0, 1) = x_2. \end{cases}$$

$$x = x_0 \hat{\varphi}_0(\hat{x}, \hat{y}) + x_1 \hat{\varphi}_1(\hat{x}, \hat{y}) + x_2 \hat{\varphi}_2(\hat{x}, \hat{y})$$

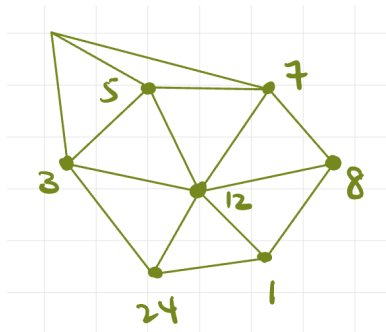
$$y = y_0 \hat{\varphi}_0(\hat{x}, \hat{y}) + y_1 \hat{\varphi}_1(\hat{x}, \hat{y}) + y_2 \hat{\varphi}_2(\hat{x}, \hat{y})$$

Then, we can build the matrix by

$$A_{ij} = \int_{\Omega} \nabla \varphi_j \nabla \varphi_i = \sum_k \int_{T_k} \nabla \varphi_j \nabla \varphi_i.$$

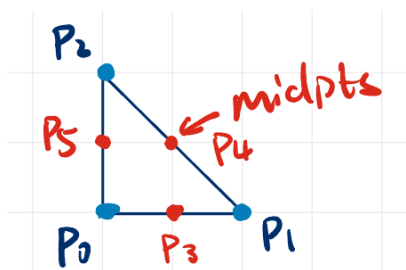
This matrix is symmetric, but we don't have clear pattern.

**Example 2.4.4**



Node 12 sees nodes 1, 8, 7, 5, 3, and 24. With a different mesh generation and numbering, we have a different matrix structure.

$\mathbb{P}^2$  Functions



$$ax^2 + bxy + cy^2 + dx + ey + f$$

$$\hat{\varphi}_0 = 2 \underbrace{(1 - \hat{x} - \hat{y})}_{\text{line b/w } P_1 P_2} \underbrace{\left(\frac{1}{2} - \hat{x} - \hat{y}\right)}_{\text{line b/w } P_3 P_5}$$

$$\hat{\varphi}_1 = 2\hat{x} \left(\hat{x} - \frac{1}{2}\right)$$

$$\hat{\varphi}_2 = 2\hat{y} \left(\hat{y} - \frac{1}{2}\right)$$

$$\hat{\varphi}_3 = 4(1 - \hat{x} - \hat{y})\hat{x}$$

$$\hat{\varphi}_4 = 4\hat{x}\hat{y}$$

$$\hat{\varphi}_5 = 4(1 - \hat{x} - \hat{y})\hat{y}$$

On a general element,  $\mathbf{x} = B_k \hat{\mathbf{x}} + \mathbf{c}_k$ . Transforming the Integral using Chain rule:

$$\left. \begin{aligned} \frac{\partial \varphi}{\partial x} &= \frac{\partial \varphi}{\partial \hat{x}} \cdot \frac{\partial \hat{x}}{\partial x} + \frac{\partial \varphi}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial x} \\ \frac{\partial \varphi}{\partial y} &= \frac{\partial \varphi}{\partial \hat{x}} \cdot \frac{\partial \hat{x}}{\partial y} + \frac{\partial \varphi}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial y} \end{aligned} \right\} \nabla_{\hat{x}, \hat{y}} \varphi = B_k^{-1} \nabla_{\hat{x}, \hat{y}} \hat{\varphi}.$$

Since  $\mathbf{x} = B_k \hat{\mathbf{x}} + \mathbf{c}$ , we have  $\hat{\mathbf{x}} = B_k^{-1} \hat{\mathbf{x}} + \mathbf{d}$ . So,

$$\frac{\partial \hat{x}}{\partial x} = (B_k^{-1})_{1,1}, \quad \frac{\partial \hat{x}}{\partial y} = (B_k^{-1})_{1,2}$$

$$\begin{aligned} \int_K \nabla \varphi_j \cdot \nabla \varphi_i \, d\mathbf{x} &= \int_{\hat{K}_\Delta} \nabla_{\hat{x}, \hat{y}} \hat{\varphi}_j J_k^{-1} \nabla_{\hat{x}, \hat{y}} \hat{\varphi}_i J_k^{-1} |J_k| \, d\hat{\mathbf{x}} \\ &= \int_{\hat{K}_\Delta} (B_k^{-1} \nabla_{\hat{x}, \hat{y}} \hat{\varphi}_i) (B_k^{-1} \nabla_{\hat{x}, \hat{y}} \hat{\varphi}_j) |\det(B_k)| \, d\hat{\mathbf{x}} \end{aligned}$$

We can also transform the semi-norm:

$$\begin{aligned} \int_K (\nabla u)^2 \, dx &= \int_{\hat{K}} (B_k^{-1} \hat{\nabla} u)^2 |\det B_k| \, d\hat{x} \\ |u|_{H^1(K)}^2 &\leq \|B_k^{-1}\|^2 |\det B_k| |u|_{H^1(\hat{K})}^2 \\ |u|_{H^1(K)} &\leq \|B_k^{-1}\| |\det B_k|^{\frac{1}{2}} |u|_{H^1(\hat{K})} \\ |u|_{H^1(\hat{K})} &\leq \|B_k\| |\det B_k|^{-\frac{1}{2}} |u|_{H^1(K)} \\ &\quad [\text{semi-norms in } H^1(K) \text{ and in } H^1(\hat{K}) \text{ are equivalent}] \\ |u|_{H^m(K)} &\leq \|B_k^{-1}\|^m |\det B_k|^{\frac{1}{2}} |u|_{H^m(\hat{K})}. \end{aligned}$$

Bounding  $\|B_k\|: \hat{x} \rightarrow x$ :

$$\begin{aligned} \|B_k\| &= \sup_{\|\xi\|=\hat{\rho}} \frac{\|B\xi\|}{\hat{\rho}} \leq \frac{h_k}{\hat{\rho}} \\ \|B_k^{-1}\| &\leq \dots \leq \frac{\hat{h}}{\rho_k} = \frac{\sqrt{2}}{\delta_k}. \end{aligned}$$

### Theorem 2.4.5 Convergence of FE in Multiple Dimension

Suppose  $u \in H^{s+1}$  and  $u_h \in \mathbb{P}^p = X_h^p$ . The error is given by

$$\|u - u_h\|_{H^1} \leq Ch^q |u|_{H^{q+1}},$$

where  $q = \min(p, s)$ .

**Proof2.** Strategies:

1. Global  $\rightarrow$  Local:  $K$

2.  $K \rightarrow \widehat{K}$
3. Error in  $\widehat{K}$
4. Go back to  $\widehat{K}$  use previous inequalities
5.  $K \rightarrow \Omega = \bigcup_{k=1}^{\infty} T_k$

Recall some results: Ceà Lemma:

$$\|u - u_h\| \leq C \inf_{w_h \in V_h} \|u - w_h\| \leq C \|u - \Pi_h^p u\|,$$

where  $\Pi_h^p$  is the polynomial interpolation of order  $p$ . To move from ②  $\rightarrow$  ③ in our strategy:

$$\int_K (\nabla u - \nabla \Pi_h^p u)^2 dx \rightarrow \int_{\widehat{K}} (\nabla \widehat{u} - \nabla \Pi_h^p \widehat{u})^2 d\widehat{x}$$

**Lemma 2.6 Bramble-Hilbert**  $u \in H^{r+1}$ ,  $\mathcal{L}(u) : H^{r+1} \rightarrow H^m$ ,  $\mathcal{L}(p) = 0 \quad \forall p \in \mathbb{P}^r$ . Then,

$$\|\mathcal{L}(v)\|_{H^m} \leq C \inf_{p \in \mathbb{P}^r} \|u + p\|_{H^{r+1}}$$

*Proof.*

$$\|\mathcal{L}(v)\| \leq \|\mathcal{L}\| \cdot \|v\|$$

Since  $\mathcal{L}(p) = 0$ , we know that

$$\|\mathcal{L}(v + p)\| \leq \|\mathcal{L}\| \cdot \|v + p\| \leq \|\mathcal{L}\| \inf_{p \in \mathbb{P}^r} \|v + p\|.$$

□

**Lemma 2.7 Deny-Lions**

$$\inf \|u + p\|_{H^{r+1}} \leq C |u|_{H^{r+1}}.$$

So, on  $\widehat{K}$ , we have

$$|u - u_{\widehat{h}}|_{H^m(\widehat{K})} \leq C |u|_{H^{r+1}(\widehat{K})}.$$

Let's go to  $K$ : Recall that

$$\begin{aligned} |u|_{H^{r+1}(\widehat{K})} &\leq \|B_k\| \cdot |\det B_k|^{-\frac{1}{2}} |u|_{H^{r+1}(K)} \\ |u|_{H^{r+1}(K)} &\leq \|B_k^{-1}\| \cdot |\det B_k|^{\frac{1}{2}} |u|_{H^{r+1}(\widehat{K})}. \end{aligned}$$

So,

$$|e|_{H^m(K)} \leq C \|B_k^{-1}\|^m |\det B_k|^{\frac{1}{2}} \|B_k\|^{r+1} |\det B_k|^{-\frac{1}{2}} |u|_{H^{r+1}(K)}.$$

Note that

$$\|B_k^{-1}\| \leq \frac{h_k}{\rho_k} \quad \text{and} \quad \|B_k\| \leq \frac{h_k}{\rho},$$

we have that

$$\begin{aligned}
 |e|_{H^m(K)} &\leq C \left( \frac{h_k}{\rho_k} \right)^m h_k^{r+1-m} |u|_{H^{r+1}(K)} \\
 &\quad [\text{Assumption: mesh is regular, so } \frac{h_k}{\rho_k} < \delta \quad (\perp k)] \\
 &\leq C \delta^m h_k^{r+1-m} |u|_{H^{r+1}(K)} \\
 &= \tilde{C} h_k^{r+1-m} |u|_{H^{r+1}(K)}.
 \end{aligned}$$

- If  $m = 1$ :

$$|e|_{H^1(K)} \leq C h^r |u|_{H^{r+1}(K)}.$$

- If  $m = 0$ :

$$|e|_{L^2(K)} \leq C h^{r+1} |u|_{H^{r+1}(K)}.$$

Q.E.D. ■

### Proposition 2.8

$$\|e_{FEM}\|_{L^2} \leq C h^{r+1} |u|_{H^{r+1}}.$$

**Proof 3.** Introduce the auxiliary problem (Aubin-Nitsche Trick):

$$\begin{cases} -\Delta u = f \\ u(\partial\Omega) = 0 \end{cases} \implies \begin{cases} -\Delta \varphi = e \\ \varphi(\partial\Omega) = 0. \end{cases}$$

Elliptic regularity:

$$\begin{cases} -\Delta u = f \in L^2 \\ u(\partial\Omega) = 0 \end{cases} \implies \|u\|_{H^2 \cap H_0^1} \leq C \|f\|_{L^2} \quad (\|u\|_{H^2} \leq C \|\Delta u\|_{L^2}).$$

In weak formulation,  $a(\varphi, v) = (e_h, v)$ . So,

$$\begin{aligned}
 \|e_h\|^2 &= a(\varphi, e_h) = a(e_h, \varphi) = a(e_h, \varphi - \varphi_h) \\
 \|e_h\|_{L^2}^2 &\leq C h |\varphi|_{H^2} \|e_h\|_{H^1} && [\text{Interpolation error, } \|\varphi - \varphi_h\|_{H^1} \leq C h |\varphi|_{H^2}] \\
 &\leq \|e_h\|_{L^2} h \|e_h\|_{H^1} && [-\Delta \varphi = e_h \implies |\varphi|_{H^2} \leq C \|e_h\|_{L^2}] \\
 \|e_h\|_{L^2} &\leq h \|e_h\|_{H^1}.
 \end{aligned}$$

Q.E.D. ■

**Proposition 2.9** The matrix  $A$  is SPD, sparse, and  $\text{cond}_2(A) \sim \mathcal{O}(h^{-2})$ .

**Remark.** Since  $A$  is SPD, solving  $Ax = b$  is equivalent to

$$\min \frac{1}{2} \mathbf{x}^\top A \mathbf{x} - \mathbf{b}^\top \mathbf{x}.$$

We could solve using CG:

$$\left\| \mathbf{e}_{\text{CG}}^{(k+1)} \right\| \leq \frac{\sqrt{\text{cond}_2(A)} - 1}{\sqrt{\text{cond}_2(A)} + 1} \left\| \mathbf{e}_{\text{CG}}^{(k)} \right\|.$$

To accelerate convergence, use pre-conditioner  $P^{-1}$  s.t.

$$\text{cond}_2(P^{-1}A) \ll \text{cond}_2(A).$$

We solve  $P^{-1}A\mathbf{x} = P^{-1}\mathbf{b}$  instead.

**Proof 4.** (of Proposition 2.9)

$$\text{cond}_2(A) = \frac{\lambda_{\max}}{\lambda_{\min}}$$

In Rayleigh coefficient form, eigenvalues are

$$\frac{\mathbf{x}^\top A \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = \frac{\mathbf{x}^\top (\lambda \mathbf{x})}{\mathbf{x}^\top \mathbf{x}} = \lambda,$$

where  $\mathbf{x}^\top A \mathbf{x} = a(v_h, v_h)$ . Note that  $v_h = \sum_i x_i \varphi_i$ . Let  $d$  denote the dimension of the problem. Then,

$$\alpha_1 h^d \mathbf{x}^\top \mathbf{x} \leq \|v_h\|_{L^2}^2 \leq \alpha_2 h^2 \mathbf{x}^\top \mathbf{x} \quad (1)$$

Also, inverse inequality gives

$$c_1 \|v_h\| \leq \|\nabla v_h\|_{L^2} \leq ch^{-1} \|v_h\|_{L^2} \quad (2)$$

Combine the two inequalities,

$$\begin{aligned} c_1 \frac{h^d \cdot \|v_h\|_{L^2}^2}{\|v_h\|_{L^2}^2} &\leq \frac{a(v_h, v_h)}{\mathbf{x}^\top \cdot \mathbf{x}} \leq c_2 \frac{h^d \|v_h\|_{H^1}^2}{\|v_h\|_{L^2}^2} && \begin{array}{l} \text{[1] Divid (1) by } \mathbf{x}^\top \mathbf{x} \\ \text{[2] By continuity} \end{array} \\ &\leq c_2 \frac{h^d h^{-2} \|v_h\|_{L^2}^2}{\|v_h\|_{L^2}^2} && \begin{array}{l} \text{[1] Poincaré: } \|v_h\|_{H^1}^2 \leq C \|\nabla v_h\|_{L^2}^2 \\ \text{[2] By (2)} \end{array} \end{aligned}$$

So,

$$c_1 h^d \leq \underbrace{\frac{a(v_h, v_h)}{\mathbf{x}^\top \mathbf{x}}}_{\text{eigenvalue}} \leq c_2 h^{d-2}.$$

So,  $\lambda_{\min} = c_1 h^d$  and  $\lambda_{\max} = c_2 h^{d-2}$ . Then,

$$\text{cond}_2(A) = \frac{\lambda_{\max}}{\lambda_{\min}} = \frac{c_2 h^{d-2}}{c_1 h^d} = ch^{-2}.$$

Q.E.D. ■

## 2.5 Mixed Problems

$$\begin{cases} -\Delta u = f \\ u(\partial u) = 0. \end{cases}$$

Post-processing: we are also interested in  $\nabla u$ .

$$\begin{aligned} u_h &= \sum_i v_i \varphi_i(x, y, z) \quad \text{from FE} \\ \nabla u_h &= \sum_i v_i \nabla \varphi_i. \end{aligned}$$

Problem with this method:

- No control over accuracy:

$$u_h \sim \mathcal{O}(h^2) \implies \nabla u_h \sim \mathcal{O}(h).$$

- $\nabla u_h$  is usually not continuous on the domain.

What to do? Auxiliary problem

$$\begin{cases} \nabla u - \sigma = 0 & (\nabla u = \sigma) \\ -\nabla \cdot \sigma = f. \end{cases}$$

- This problem is also a minimization problem:

$$\begin{aligned} J^* &= \int_{\Omega} |\mathbf{p}|^2, & W_f &= \underbrace{\{\nabla \cdot \mathbf{w} + f = 0\}}_{\text{constraint}} \\ & & \min_{\mathbf{p} \in W_f} J^*. \end{aligned}$$

Lagrangian multiplier:  $\mathcal{L} = \int_{\Omega} |\mathbf{p}|^2 + \int v(\nabla \cdot \mathbf{p} + f)$ . Then, the problem becomes

$$\min_{\mathbf{p}} \max_v \mathcal{L} \quad \leftarrow \text{saddle point problem}$$

max: pushing  $v$  to be large, we force  $\nabla \cdot \mathbf{p} + f = 0$ .

- FOC: Find  $\sigma \in H(\text{div}; \Omega) \equiv \left\{ \sigma_i \in L^2, \nabla \cdot \sigma = \sum \frac{\partial \sigma_i}{\partial x_i} \in L^2 \right\} = Y, u \in L^2(\Omega) = X$  s.t.

$$\begin{cases} \int_{\Omega} \sigma \cdot \mathbf{p} + \int \nabla \cdot \mathbf{p} u = 0 \\ \int_{\Omega} v(\nabla \cdot \sigma + f) = 0 \end{cases} \quad \forall v \in X, \mathbf{p} \in Y.$$

One can verify that

$$\begin{cases} -\nabla \cdot \boldsymbol{\sigma} = f \\ \boldsymbol{\sigma} - \nabla \mathbf{u} = 0 \end{cases} \implies \begin{cases} \int v(\nabla \cdot \boldsymbol{\sigma} + f) = 0 \\ \int (\boldsymbol{\sigma} - \nabla \mathbf{u}) \cdot \mathbf{p} = 0. \end{cases} \quad [\text{Multiply by test functions and integrate}]$$

For the second integral, consider integration by parts:  $\int \boldsymbol{\sigma} \cdot \mathbf{p} + \int \nabla \cdot \mathbf{p} u = 0$ .

- The Lagrangian multiplier problem is well-posed.

**Proposition 2.1**  $\inf - \sup$  **Condition/LBB**  $\forall v \in X, \exists \mathbf{p} \in Y$  s.t.  $\exists \beta > 0$  s.t.

$$\left| \int_{\Omega} v \nabla \cdot \mathbf{p} \right| \geq \beta \|v\|_X \|\mathbf{p}\|_Y,$$

where  $\|\mathbf{p}\|_{H(\text{div})} = \|\mathbf{p}\|_{L^2} + \|\nabla \cdot \mathbf{p}\|_{L^2}$

$$\iff \inf_{\mathbf{p} \in Y} \sup_{v \in X} \frac{\left| \int_{\Omega} v \nabla \cdot \mathbf{p} \right|}{\|v\|_X \|\mathbf{p}\|_Y} \geq \beta.$$

**Proof 1.** (Sketch) Consider the auxiliary problem:

$$\begin{cases} -\Delta \psi = v \\ \psi(\partial\Omega) = 0, \end{cases} \quad \mathbf{p} = \nabla \psi.$$

Q.E.D. ■

With the inf-sup condition, the Lagrangian multiplier problem is well-posed.

- Is the numerical problem well-posed?

Consider  $X_h \subset X$  and  $Y_h \subset Y$ . Then,

$$\begin{cases} \int_{\Omega} v_h(\nabla \cdot \mathbf{p}_h + f) = 0 \\ \int_{\Omega} \mathbf{p}_h \cdot \boldsymbol{\sigma}_h + (\nabla \cdot \mathbf{p}_h)v_h = 0. \end{cases}$$

We cannot select  $X_h$  and  $Y_h$  as we wish. We have to choose  $X_h$  and  $Y_h$  s.t. the inf-sup condition holds. i.e.,  $\forall v_h \in X_h, \exists \mathbf{p}_h \in Y_h$  s.t.

$$\left| \int_{\Omega} \nabla \cdot \mathbf{p}_h v_h \right| \geq \beta \|v_h\|_X \|\mathbf{p}_h\|_Y.$$

- Raviart-Thomas:

$$\begin{cases} v_h \in \mathbb{P}^{k-1} \\ \mathbf{p}_h \in \mathbb{D}^k \equiv [\mathbb{P}^{k-1}]^d + \mathbf{x}\mathbb{P}^{k-1}. \end{cases}$$

This couple is inf-sup compatible.

$$\|v - v_h\|_{L^2} - \|\mathbf{p} - \mathbf{p}_h\|_{H^d} \sim \mathcal{O}(h^k).$$

There are other compatible couples, such as DBM.

**Remark 2. (Messages from this Section).**

- We can reformulate a second order problem as a system of first order problems, but the requirement of regularity is different.
- In first order formulation, we lose coercivity, and we have an inf-sup (min-max/saddle point) problem. So, we need to find the right pair for discretization.

**2.6 Remarks on FE**

- Why triangles?

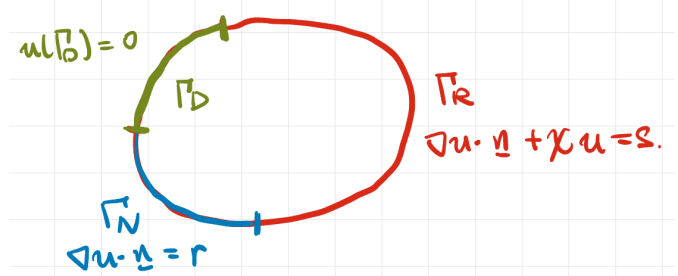
We can also use squares, and we have to modify the polynomials accordingly:

$$\mathbb{Q}^1 = axy + bx + cy + d \quad \leftarrow \text{Bilinear polynomials}$$

If we use pentagon, we use *Mimetic finite difference/Viral finite elements*.

- Weak formulation with General BCs:

$$\begin{cases} u(\Gamma_D) = 0 \\ \nabla u \cdot \mathbf{n} = r & \text{on } \Gamma_N \\ \nabla u \cdot \mathbf{n} + \chi u = s & \text{on } \Gamma_R \end{cases}$$



$$-\Delta u = f \quad \text{where } u \in H_0^1 + \mathcal{L},$$

where  $\mathcal{L}$  is a lifting function.  $\mathcal{L}_g$  is a  $H^1(\Omega)$  s.t.  $\mathcal{L}_g(\Gamma_D) = g$ . So,

$$\begin{aligned} \int_{\Omega} \nabla u \nabla v &= \int_{\Omega} f v + \underbrace{\int_{\partial\Omega} \nabla u \cdot \mathbf{n} v}_{=0 \text{ on } \Gamma_D} \quad \forall v \in H_0^1 \\ &= \int_{\Omega} f v + \int_{\Gamma_N \cup \Gamma_R} \nabla u \cdot \mathbf{n} v \end{aligned}$$

$$\int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v + \int_{\Gamma_N} r v + \int_{\Gamma_R} s v - \int_{\Gamma_R} \chi u v.$$

So, the problem becomes: Find  $u \in \mathcal{L} + H_0^1$  s.t.

$$\int_{\Omega} \nabla u \nabla v + \int_{\Gamma_R} \chi u v = \int_{\Omega} f v + \int_{\Gamma_N} r v + \int_{\Gamma_R} s v \quad \forall v \in H_0^1.$$

## 2.7 Spectral Method

We are still in Galerkin framework.

$$a(u_h, v_h) = \mathcal{F}(v_h) \quad \text{on some } X_h.$$

$$\inf_{w_h \in X_h} \|u - w_h\| \xrightarrow{h \rightarrow 0} 0.$$

**Idea** Use global polynomials (Gaussian interpolation) on special nodes.

- Gaussian Interpolation:

$$\int_a^b f \approx \int_a^b \pi_f.$$

1.  $\int_a^b f(x) = f(\bar{x})(b-a)$ . Do we have an optimal choice for  $\bar{x}$ ?

$$\bar{x} = \frac{a+b}{2}.$$

2. We want two nodes:

$$\int_a^b f(x) = f(x_1) \underbrace{\frac{x-x_1}{x_2-x_1}}_{w_1} + f(x_2) \underbrace{\frac{x-x_2}{x_2-x_1}}_{w_2}$$

We can simply take  $x_1 = a$  and  $x_2 = b$  to recover the Trapezoidal rule, but can we find the optimal  $x_1$  and  $x_2$ ?

3. In general,

$$\int_a^b f(x) \approx \sum_i w_i f(x_i).$$

This problem is not easy to solve... But we can recover it from orthogonal polynomials.

*[Why Orthogonal polynomials? Let's see some Fourier series]*

$$u(x) = \sum_{j=1}^{\infty} c_j \sin(j\pi x)$$

$$c_j = \frac{\int u(x) \sin(j\pi x) dx}{\int \sin^2(j\pi x) dx}$$

We can truncate:

$$u = \underbrace{\sum_{j=1}^N c_j \sin(j\pi x)}_{u_N} + \underbrace{\sum_{j=N+1}^{\infty} c_j \sin(j\pi x)}_{\text{Remainder}}$$

Because of orthogonality, Remainder is orthogonal to  $u_N$ . Up to  $N$  terms,  $u_N$  is the best approximation of  $u$ .

**Bottleneck:** good for periodic problems, but technical difficulties when imposing BCs.

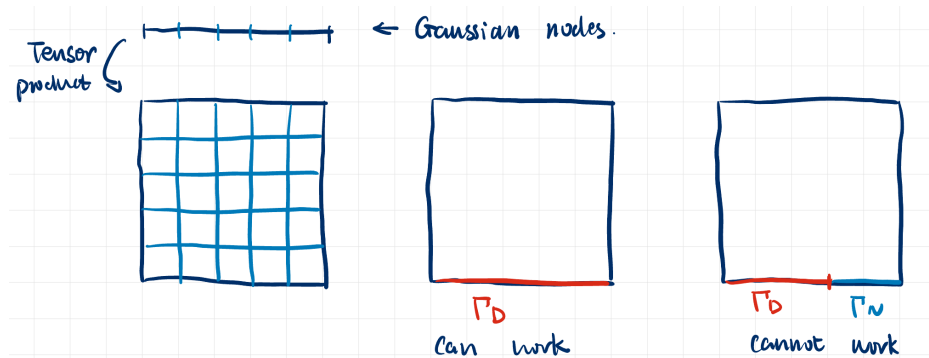
- Legendre polynomial on  $L^2(-1, 1)$ :

$$\int_{-1}^{+1} \pi_L^K \pi_L^R = 0 \quad \text{if } K \neq R.$$

- Gauss Theory: interpolating nodes are always inside the domain. Nothing on the boundary is provided.

Gauss-Lobatto: Correction to Gauss Theory that allows nodes to be on the boundary.

However, we still require the BC to be the same on one edge.



Ceà Lemma:

$$\inf_{p_N \in \mathbb{Q}_N} \|u - p_N\|_{H^1} \leq C \underbrace{N^{-(s-1)}}_{=N^{s-1}} |u|_{H^{s+1}} \quad u \in H^{s+1}$$

$$u \in C^\infty \implies Ce^{-N}.$$

This is *exponential convergence*: If the solution is regular enough, spectral method would work well.

- $\int_{-1}^{+1} u'_N v'_N dx = \int_{-1}^{+1} f v_N dx$ , where  $u_N, v_N \in \mathbb{P}^{N-1}$ , and  $\int_{-1}^{+1} f v_N dx$  gives the source of error: we don't know degree of  $f$ . We usually never compute this exactly.

$$\int_{-1}^{+1} u'_N v'_N dx = \sum w_N u'_N v'_N \quad \leftarrow \text{no quadrature error here}$$

To approximate  $\int_{-1}^{+1} f v_N dx$ , we use GNI (Galerkin Numerical Integration).

**Lemma 2.1 Strang Lemma**

$$\|u - u_N\|_H \leq \underbrace{C \inf_{w_N} \|u - w_N\|}_{\text{Discretization}} + \underbrace{\text{Quadrature Error.}}_{\substack{\text{Could be controlled} \\ \text{by using Gaussian quadrature}}}$$

### 3 Advection-Diffusion Problem

$$\underbrace{-\mu\Delta u}_{\text{diffusion}} + \underbrace{\beta\nabla u}_{\text{advection}} + \underbrace{\sigma u}_{\text{reaction}} = f.$$

- Diffusion problem (Laplace) was derived from

$$\min J(u) \equiv \frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} f v.$$

- $\beta$  breaks the symmetry.

#### 3.1 Finite Difference of 1D Problem

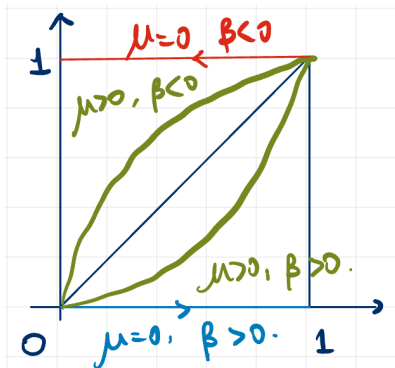
Consider

$$\begin{cases} -\mu u'' + \beta u' = 0, & \mu \geq 0 \\ u(0) = 0 \quad \text{and} \quad u(1) = 1. \end{cases}$$

Then,

$$u = \frac{e^{\beta/\mu x} - 1}{e^{\beta/\mu} - 1}$$

is the analytical solution.



- FD Discretization:

$$-\mu \underbrace{\frac{u_{i+1} - 2u_i + u_{i-1}}{h^2}}_{\sim \mathcal{O}(h^2)} + \beta \underbrace{\frac{u_{i+1} - u_{i-1}}{2h}}_{\substack{\text{centered difference} \\ \sim \mathcal{O}(h^2)}} = 0.$$

To show order of convergence, use Taylor's expansion.

- If  $\beta$  is large, then we have oscillation, and we need finer mesh.

How small  $h$  needs to be to prevent oscillatory solutions?

Use second order difference equation to solve:

$$\begin{cases} \underbrace{\left(\frac{\mu}{h^2} + \frac{\beta}{2h}\right)}_a u_{i+1} - \underbrace{\frac{2\mu}{h^2}}_b u_i + \underbrace{\left(\frac{\mu}{h^2} - \frac{\beta}{2h}\right)}_c u_{i-1} = 0 \\ u_0 = 0, \quad u_N = 1 \end{cases}$$

So, the second order difference equation is

$$a\rho^2 + b\rho + c = 0$$

$$u_i = \alpha\rho_1^i + \beta\rho_2^i \implies u_0 = \alpha + \beta = 0 \implies \alpha = -\beta.$$

Assume  $\rho_1 = 0$ :

$$\text{If } \frac{c}{a} = \rho_1\rho_2 < 0 \implies \rho_2 \text{ is negative} \implies \text{oscillations}$$

1.  $\mathbb{P}_e = \frac{\beta h}{2\mu} > 1$  (Péclet): Oscillations
2.  $\mathbb{P}_e < 1 \implies c > 0 \implies$  no oscillations.

So, we require step size:

$$\frac{\beta h}{2\mu} < 1 \implies h < \frac{2\mu}{\beta}.$$

**Problem:** If  $\frac{\mu}{\beta} = 10^{-6}$ , we need over 1 million points...

- Can we have a better scheme? A upwind scheme.

This scheme breaks the centered symmetry.

$$\beta u' \approx \beta \frac{u_i - u_{i-1}}{h}$$

1. Accuracy: first-order  $\sim \mathcal{O}(h)$

We sacrifice the accuracy to allow larger step size.

2. This method is *physically consistent*, although we lose something mathematically.
3. Prove this scheme avoids oscillations:

Method I Repeat the second order difference equation analysis.

Method II

### Theorem 3.1.1

Upwind never oscillates.

**Proof 1.**

$$\begin{aligned} \beta \frac{u_i - u_{i-1}}{h} &= \underbrace{\frac{\beta}{2} \frac{u_{i+1} - u_{i-1}}{h}}_{\text{centered approx.}} - \frac{\beta}{2h} u_{i+1} + \frac{\beta u_i}{h} - \frac{\beta}{2h} u_{i-1} \\ &= \beta \frac{u_{i+1} - u_{i-1}}{2h} - \frac{\beta h}{2} \underbrace{\frac{u_{i+1} - 2u_i + u_{i-1}}{h^2}}_{\text{approx. of } u''} \\ &= \text{centered approx.} + \frac{\beta h}{2} (-u'') \end{aligned}$$

So, the original problem:  $-\mu u'' + \beta u' = 0$  becomes

$$\begin{aligned}
 \text{UPW} \quad & -\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \frac{\beta}{h}(u_i - u_{i-1}) = 0 \\
 & -\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \frac{\beta}{2h}(u_{i+1} - u_{i-1}) - \frac{\beta h}{2} \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} = 0 \\
 & -\left(\mu + \frac{\beta h}{2}\right) \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \frac{\beta}{2h}(u_{i+1} - u_{i-1}) = 0 \quad [\mathbb{P}_e = \frac{\beta h}{2\mu}] \\
 & \quad \quad \quad -\underbrace{\mu(1 + \mathbb{P}_e)}_{\mu^*} u'' + \beta u' = 0
 \end{aligned}$$

Péclet number of this new problem:

$$\mathbb{P}_e^* = \mathbb{P}_e^{\text{UPW}} = \frac{\beta h}{2\mu^*} = \frac{\beta h}{2\mu(1 + \mathbb{P}_e)} = \frac{\mathbb{P}_e}{1 + \mathbb{P}_e} < 1.$$

Q.E.D. ■

4. Is this method stable?

We have a perturbed problem:

$$-\left(\mu + \frac{\beta h}{2}\right) u'' + \beta u' = 0$$

Note that  $\frac{\beta h}{2} \rightarrow 0$  as  $h \rightarrow 0$ . So, we are stable.

$$\text{upwind} = \text{centered}(\mu(1 + \mathbb{P}_e)).$$

5. Solution is always consistent, regardless of choice of  $h$ .

6. Can we have second order upwind?

Use Taylor's expansion to design a second order approximation of  $u'_i$ , using only  $u_{i-2}$  and  $u_{i-1}$ .

7. If  $\beta < 0$ , the wind is backward. Use

$$\frac{u_{i+1} - u_i}{h}.$$

### 3.2 Finite Difference in 2D+

Wind is a vector.

- $\beta = [\beta \ 0]$ , horizontal wind.

$$-\mu(1 + \mathbb{P}_e)\Delta u + \beta \cdot \nabla u = 0.$$

In this case, only apply upwind to  $\frac{\partial u}{\partial x}$ , as the wind is horizontal. No change on  $y$  direction.

- $\beta = [\beta_1 \ \beta_2]$ .

$$\beta_1 \frac{\partial u}{\partial x} + \beta_2 \frac{\partial u}{\partial y}.$$

Streamline Diffusion:

$$-\mu \nabla \cdot (\nabla u) + \beta \cdot \nabla u - \frac{\frac{1}{2}h}{\|\beta\|} \nabla \cdot ((\beta \cdot \nabla u)\beta) = f.$$

Weak formulation:

$$-\frac{\frac{1}{2}h}{\|\beta\|} \int_{\Omega} \nabla \cdot ((\beta \cdot \nabla u)\beta) v \, dw = \frac{\frac{1}{2}h}{\|\beta\|} \int_{\Omega} (\beta \cdot \nabla u)(\beta \cdot \nabla v) \, dw.$$

### 3.3 Finite Elements in 1D

Consider

$$\begin{cases} -\mu u'' + \beta u' + \sigma u = 0, & x \in [0, 1] \\ u(0) = 0, & u(1) = 1. \end{cases}$$

- What is a weak formulation?

$$v(-\mu u'' + \beta u' + \sigma u) = 0, \quad \text{for } v(0) = v(1) = 0.$$

Define  $\ell(x) = x$ . Then,  $u = \hat{u} + \ell(x)$ , then  $\hat{u}(0) = \hat{u}(1) = 0$ .

$$\begin{aligned} (-\mu \hat{u}'' + \beta \hat{u}' + \sigma \hat{u})v &= v(-\beta \cdot 1 - \sigma \ell) \\ \int_0^1 -\mu \hat{u}'' v + \beta \hat{u}' v + \sigma \hat{u} v &= \int_0^1 v(-\beta \cdot 1 - \sigma \ell) \\ \underbrace{\int_0^1 \mu \hat{u}' v' + \beta \hat{u}' v + \sigma \hat{u} v}_{a(\hat{u}, v)} &= - \underbrace{\int_0^1 (\beta + \sigma \ell) v}_{\mathcal{F}(v)} \end{aligned}$$

So, the variational form:  $a(\hat{u}, v) = \mathcal{F}(v)$ , where

$$a(\hat{u}, v) = \int_0^1 \mu \hat{u}' v' + \beta \hat{u}' v + \sigma \hat{u} v.$$

Still,  $\beta \hat{u}' v$  is breaking the symmetry here.

1.  $a(\hat{u}, v)$  is not symmetric:  $a(\hat{u}, v) \neq a(v, \hat{u})$ .
2.  $a(\hat{u}, v)$  is coercive:  $a(u, u) \geq \alpha \|u\|_{H^1}^2$ .

**Proof 1.**

$$a(u, u) = \int_0^1 \mu (u')^2 + \sigma u^2 + \beta u' u.$$

Note that  $u' u = \frac{1}{2}(u^2)'$ , so  $\int_0^1 \beta u' u = \frac{1}{2} \int_0^1 \beta (u^2)'$ . Using integration by parts,

$$\frac{1}{2} \int_0^1 \beta (u^2)' = \frac{1}{2} \underbrace{[\beta u^2]_0^1}_{=0} - \frac{1}{2} \int_0^1 \beta' u^2$$

So,

$$\begin{aligned} a(u, u) &= \int_0^1 \mu(u')^2 + \sigma u^2 - \frac{1}{2}\beta' u^2 \\ &= \int_0^1 \mu(u')^2 + \left(\sigma - \frac{1}{2}\beta'\right) u^2 \end{aligned}$$

**Assumption:**  $\sigma - \frac{1}{2}\beta' \geq \rho$ .

Then,

$$\begin{aligned} a(u, u) &\geq \int_0^1 \mu(u')^2 + \rho u^2 \\ &\geq \mu \|u'\|_{L^2}^2 + \rho \|u\|_{L^2}^2 \\ &\geq \alpha \|u\|_{H^1}^2, \quad \alpha = \min\{\rho, \mu\}. \end{aligned}$$

Q.E.D. ■

3.  $a(\hat{u}, v)$  is continuous:  $|a(u, v)| \leq \gamma \|u\|_{H^1} \|v\|_{H^1}$ .

**Proof 2.**

$$\mu \int u'v' + \beta \int u'v + \sigma \int uv \leq \max\{\mu, |\beta|, \sigma\} \|u\|_{H^1} \|v\|_{H^1}$$

Q.E.D. ■

**Remark.**

- Usually,  $\mu < \rho$
- For a convection-dominated problem,  $|\beta|$  will be the largest
- For a reaction dominated problem,  $\sigma$  will be the largest.

4. Strong consistent:

$$a(u - u_h, w_h) = 0.$$

**Lemma 3.1 Cèa**

$$\|u - u_h\| \leq \left( \frac{C_{\text{continuous}}}{C_{\text{coercivity}}} \right)^{\frac{\gamma}{\alpha}} \inf_{w_h \in X_h} \|u - w_h\|.$$

Usually,  $\frac{\gamma}{\alpha}$  can be huge, and so Cèa lemma is not very helpful.

- Linear FE:

$$A = [a(\varphi_i, \varphi_j)].$$

$$\mu \int_0^1 \varphi'_i \varphi'_{i\pm 1}, \quad \beta \int_0^1 \varphi'_{i\pm 1} \varphi_i, \quad \text{and} \quad \sigma \int_0^1 \varphi_i \varphi_{i\pm 1}.$$

The FE system is

$$\frac{\mu}{h}(u_{i+1} - 2u_i + u_{i-1}) + \frac{\beta}{2h}(u_{i+1} - u_{i-1}) + \frac{\sigma h}{6}(u_{i+1} + 4u_i + u_{i-1}) = 0.$$

1. Suppose  $\sigma = 0$ :

$$\frac{1}{h}(\text{FE}) = \text{FD}$$

So, FE suffer in the same way as FD. To prevent oscillation, we also require

$$\mathbb{P}_e = \frac{|\beta|h}{2\mu} < 1.$$

2. As lesson from FD: add a diffusive term  $-\frac{\beta h}{2}u''$ . In weak formulation:

$$a(u_h, v_h) + \frac{\beta h}{2} \int_0^1 u'_h v'_h = \mathcal{F}(v_h).$$

This is exactly the upwind FE method.

3. What we lose with UPW-FE: No strong consistency, so we don't have Cèa Lemma anymore.  
4. What we still have: Strang Lemma.

**Lemma 3.2 Strang**

$$\|u - u_h\| \leq C \inf_{w_h} \|u - w_h\| + \sup_{w_h} \underbrace{(|a(u, w_h) - a_h(u, w_h)|)}_{\frac{\beta h}{2} \int_0^1 u'_h v'_h}$$

**Implication:** Regardless of the order of polynomial, due to the additional term, the convergence rate is always  $\sim \mathcal{O}(h)$ . So, it makes no sense to use higher order than linear.

- Scharfter-Gammal: Design a better  $\mu$  such that

$$\mu^{\text{smart}} = \mu(1 + \Phi(\mathbb{P}_e)).$$

**Problem:** Only work for 1D problems.

### 3.4 Convection-Dominated problem

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} + \beta \frac{\partial u}{\partial x} = 0.$$

Analytical approach: change variable  $u \mapsto w$  s.t.

$$\frac{\partial w}{\partial t} - \chi \frac{\partial^2 w}{\partial x^2} = 0.$$

Can we don the same in our numerical method? *Slotboom Variable*

**Proof 1.** In 1D,

$$-\mu u'' + \beta u' = 0, \quad \mu \geq 0.$$

Suppose  $u = \rho e^{\sigma x}$ , where  $\sigma$  is to be determined (Slotboom variable). Then,

$$\begin{aligned} u' &= \rho' e^{\sigma x} + \sigma \rho e^{\sigma x} \\ u'' &= \rho'' e^{\sigma x} + 2\rho' \sigma e^{\sigma x} + \sigma^2 \rho e^{\sigma x}. \end{aligned}$$

Then,

$$\begin{aligned} e^{\sigma x} [-\mu \rho'' - 2\mu \rho' \sigma - \mu \sigma^2 \rho + \beta \rho' + \beta \sigma \rho] &= 0 \\ -2\mu \sigma + \beta &= 0 \implies \sigma = \frac{\beta}{2\mu} \\ \beta \sigma - \mu \sigma^2 &= \beta \frac{\beta}{2\mu} - \mu \frac{\beta^2}{4\mu^2} = \frac{\beta^2}{4\mu}. \end{aligned}$$

So,

$$-\mu \rho'' + \frac{\beta^2}{4\mu} \rho = 0.$$

Weak formulation:

$$\int \mu \rho' v + \frac{\beta^2}{4\mu} \int \rho v = 0.$$

Q.E.D. ■

This method is better than traditional UPW-FE, but still may suffer from instabilities. Moreover, the weak formulation solves for  $\rho$ . To get  $u$  back, we need exponential functions: more errors.

### 3.5 Strongly Consistent Stabilization

$$a(u, v) + b_h(u, v) = \mathcal{F}(v)$$

Generalized Galerkin:

$$a_h(u, v) = \mathcal{F}_h(v).$$

- Streamline Diffusion:

$$b_h = \frac{ch}{\|\beta\|} \int (\beta \cdot \nabla u)(\beta \cdot \nabla v) \, d\Omega.$$

- By Strang's Lemma, as long as  $b_h \rightarrow 0$  as  $h \rightarrow 0$ , we still have convergence:

$$\|u_h - u\| \leq \underbrace{1}_{\text{C\^e a}} + \underbrace{2}_{|a - a_h|} + \underbrace{3}_{|\mathcal{F} - \mathcal{F}_h| = 0}$$

$$\inf_{w_h \in V_h} \|u - w_h\| = \frac{ch^2}{\|\beta\|} (\beta \cdot \nabla u)(\beta \cdot \nabla v)$$

- We don't have strong consistency:

$$a_h(u, v_h) = a(u, v_h) + \underbrace{b_h(u, v_h)}_{\neq 0} \neq \mathcal{F}(v_h)$$

- Strongly Consistent Stabilization (Hughes and Brooks):

$$a(u, v) = \mathcal{F}(v).$$

1. Let's add something that will vanish if the solution is regular enough: residual in the strong form

$$a(u, v) + (a(u, v) - \mathcal{F}(v)) = \mathcal{F}(v).$$

2. In strong form:  $-\mu\Delta u + \boldsymbol{\beta} \cdot \nabla u - f$ .

Go back to weak:

$$\int_{\Omega} (-\mu\Delta u + \boldsymbol{\beta} \cdot \nabla u - f) \mathcal{L}v,$$

where  $\mathcal{L}$  is an operator on  $v$ .

We cannot do this at continuous level, but it is perfectly fine at the discrete level:

$$\sum_K \delta_K \int_K (-\mu\Delta u_h + \boldsymbol{\beta} \cdot \nabla u_h - f) \mathcal{L}v_h.$$

*[If we take  $\mathcal{L}v_h = \boldsymbol{\beta} \cdot \nabla v_h$ , we have exactly the streamline diffusion.]*

3. What is  $\mathcal{L}$ ?

$$A = \frac{1}{2} \underbrace{(A + A^T)}_{\text{symmetric: } B=B^T} + \frac{1}{2} \underbrace{(A - A^T)}_{\text{skew-symmetric } B=-B^T}$$

**Goal:** Write  $\mathcal{L} = -\mu\Delta u + \boldsymbol{\beta} \cdot \nabla u = \mathcal{L}_s + \mathcal{L}_{ss}$ .

$$\begin{aligned} \mathcal{L}_s &= \int_{\Omega} \mu \nabla u \nabla v \\ \int_{\Omega} (\boldsymbol{\beta} \cdot \nabla u) v &= - \int_{\Omega} u \nabla \cdot (\boldsymbol{\beta} v) = - \int_{\Omega} u \boldsymbol{\beta} \cdot \nabla v - \underbrace{\int_{\Omega} uv \nabla \cdot (\boldsymbol{\beta})}_{\text{symmetric}} \\ \underbrace{\int_{\Omega} \boldsymbol{\beta} \cdot \nabla uv + \frac{1}{2} \int_{\Omega} uv \nabla \cdot (\boldsymbol{\beta})}_{\mathcal{G}(u,v) = -\mathcal{G}(u,v), \text{ skew symmetric}} &= - \int_{\Omega} u \boldsymbol{\beta} \cdot \nabla v - \frac{1}{2} \int_{\Omega} uv \nabla \cdot (\boldsymbol{\beta}). \end{aligned}$$

So,

$$\mathcal{L} = \underbrace{\mu \int_{\Omega} \nabla u \nabla v - \frac{1}{2} \int_{\Omega} uv \nabla \cdot (\boldsymbol{\beta})}_{\mathcal{L}_s} + \underbrace{\int_{\Omega} (\boldsymbol{\beta} \cdot \nabla u) v + \frac{1}{2} \int_{\Omega} uv \nabla \cdot (\boldsymbol{\beta})}_{\mathcal{L}_{ss}}.$$

## 4. SUPG: Streamline Upwind Petrov Galerkin:

$$\sum \delta_k \int_k (\mathcal{L}u - f)(\mathcal{L}_{ss}v) dw,$$

$$* \mathcal{L}u = -\mu\Delta u + \beta \cdot \nabla u$$

$$* \mathcal{L}_s = -\mu\Delta u - \frac{1}{2}(\nabla \cdot \beta)u$$

$$* \mathcal{L}_{ss} = \beta \cdot \nabla u + \frac{1}{2}(\nabla \cdot \beta)u$$

## 5. GaLS: Galerkin Least squares

$$\mathcal{L}_\rho = \rho\mathcal{L}_s + \mathcal{L}_{ss}$$

$$* \rho = 0: \text{SUPG}$$

$$* \rho = 1: \text{GaLS}$$

$$* \rho = -1: \text{Douglas-Wang (DW)}$$

## 6. For all the methods:

$$\|u_h - u_{\text{ex}}\|_{*\rho} \leq Ch^{r+\frac{1}{2}}|u|_{H^{s+1}}, \quad r = \min\{p, s\},$$

where  $\|\cdot\|_{*\rho}$  is not a standard norm, and it depends on  $\rho$ .

**Implication:** If the solution is regular enough, we can still benefit from using high order polynomials.

$$\|u\|_{\rho=1}^2 \equiv \mu \|\nabla u_h\|_{L^2}^2 + \min\{\nabla \cdot \beta, \sigma\} \|u_h\|_{L^2}^2 + \sum_K \delta_k \|\mathcal{L}u_h\|_{L^2(T_u)}$$

## 3.6 Reaction-Dominated Problem

$$-\mu\Delta u + \sigma u = f.$$

- FD:

$$-\mu \frac{u_{i+1} - 2u_i + u_{i-1}}{h^2} + \sigma u_i = f_i,$$

and the matrix will be a good matrix:

$$A = \mu \begin{bmatrix} 2 & 1 & & \\ -1 & \ddots & \ddots & \\ & \ddots & \ddots & \end{bmatrix} + \sigma I$$

- FE:

$$\sigma \int_0^1 \varphi_i \varphi_j \rightarrow \frac{h}{6} \begin{bmatrix} 4 & 1 & & \\ 1 & \ddots & \ddots & \\ & \ddots & \ddots & \end{bmatrix} \leftarrow \text{Mass Matrix.}$$

We need  $\mathbb{P}_e = \frac{|\sigma|h^2}{6\mu} < 1$  to avoid oscillation.

We are in trouble, but not as much as in convection-domination,

$$h < \sqrt{6} \sqrt{\frac{\mu}{|\sigma|}},$$

where the square root is helping a lot.

**Notice:** Instead of solving  $\int_0^1 \varphi_i \varphi_j$  exactly, we use a trapezoidal rule:

$$(T) \int_i^{i+1} \varphi_i \varphi_j = \frac{\varphi_i \varphi_j(x_{i+1}) + \varphi_i \varphi_j(x_i)}{2} h.$$

This scheme produces the identity matrix  $I$ .

- **Mass Lumping:** Make the mass matrix  $M$  a diagonal matrix by summing the off-diagonal terms to the diagonal.

1.

$$\begin{aligned} \sigma \int_0^1 \varphi_i \varphi_j &\approx \sigma(T) \int_0^1 \varphi_i \varphi_j = \sigma(T) \int_{x_{i-1}}^{x_{i+1}} \varphi_i \varphi_j \\ &= \sigma \left( (T) \int_{x_{i-1}}^{x_i} \varphi_i \varphi_{i-1} + \underbrace{\int_{x_i}^{x_{i+1}} \varphi_i \varphi_{i-1}}_{=0} + \int_{x_{i-1}}^{x_i} \varphi_i \varphi_{i+1} + \int_{x_i}^{x_{i+1}} \varphi_i \varphi_{i+1} \right) \\ &= \sigma \frac{h}{2} (\varphi_i \varphi_{i-1}(x_i) + \varphi_i \varphi_{i-1}(x_{i-1})) + 0 \\ &= 0. \end{aligned}$$

So, the off-diagonal terms are 0. For the diagonal term,

$$\begin{aligned} \sigma(T) \int_{x_{i-1}}^{x_i} \varphi_i^2 + \sigma(T) \int_{x_i}^{x_{i+1}} \varphi_i^2 &= \frac{\sigma h}{2} (\varphi_i^2(x_{i-1}) + \varphi_i^2(x_i)) + \frac{\sigma h}{2} (\varphi_i^2(x_i) + \varphi_i^2(x_{i+1})) \\ &= \sigma h. \end{aligned}$$

Therefore,  $FD = FE + ML$  (mass lumping)

## 2. Lemma 3.1 Strang's Lemma

$$-\mu \underbrace{\int_0^1 \varphi_i' \varphi_j'}_{\mathbb{P}^0} + \beta \underbrace{\int_0^1 \varphi_0' \varphi_j}_{\mathbb{P}^1} + \sigma \underbrace{\int_0^1 \varphi_i \varphi_j}_{\mathbb{P}^2}.$$

Mass lumping (ML) can be viewed as a **Generalized Galerkin Method**:

$a(u, v) = \mathcal{F}(v) \implies a_h(u, v) = \mathcal{F}_h(v)$ , and  $\sigma \int uv - \sigma(T) \int uv \sim \mathcal{O}(h^2)$ . So, we can go up to quadratic FE and achieve  $\mathcal{O}(h^2)$  with mass lumping.

## 4 Parabolic Equations

$$\begin{aligned} \frac{\partial u}{\partial t} - \mu \Delta u + \boldsymbol{\beta} \cdot \nabla u + \sigma u &= f, \quad \mathbf{x} \in \Omega, \quad \mu \geq 0 \\ + \text{B.C. : } u(\partial\Omega) &= g \quad (\text{Can be Dirichlet, Neumann, Robin}) \\ + \text{I.C.} \end{aligned}$$

### 4.1 Weak Formulation and Well-Posedness

$$\begin{aligned} \left( \frac{\partial u}{\partial t} - \mu \Delta u + \boldsymbol{\beta} \cdot \nabla u + \sigma u \right) v &= f v \\ \int_0^T \int_{\Omega} \left( \frac{\partial u}{\partial t} - \mu \Delta u + \boldsymbol{\beta} \cdot \nabla u + \sigma u \right) v &= \int_0^T \int_{\Omega} f v \\ \iint \frac{\partial u}{\partial t} v - \mu \iint \Delta u v + \iint \boldsymbol{\beta} \cdot \nabla u v + \iint \sigma u v &= \iint f v \end{aligned}$$

- By Green's formula:

$$- \iint \Delta u v = \iint \nabla u \nabla v$$

- If  $\Omega \perp\!\!\!\perp t$ , then

$$\int_0^T \int_{\Omega} \frac{\partial u}{\partial t} v = \int_0^T \frac{d}{dt} \int_{\Omega} u v$$

#### Theorem 4.1.1

$$\frac{d}{dt} \int_{\Omega(t)} u v = \int_{\Omega} \frac{\partial(uv)}{\partial t} + \int_{\partial\Omega} (\mathbf{w} \cdot \mathbf{n}) u v,$$

$$\text{where } \mathbf{w} = \frac{\partial\Omega(t)}{\partial t}.$$

- The bilinear form is

$$\int_0^T \left( \frac{d}{dt} \int_{\Omega} u v + \underbrace{\mu \int_{\Omega} \nabla u \nabla v + \int_{\Omega} \boldsymbol{\beta} \cdot \nabla u v + \int_{\Omega} \sigma u v}_{a(u,v)} \right) = \int_0^T \underbrace{\int_{\Omega} f v}_{\mathcal{F}(v)},$$

where  $u, v \in H_0^1$ .

- **Lemma 4.2 Lax-Milgram Lemma Extension** If  $a(u, v)$  is bilinear, continuous in  $H_0^1(\Omega)$ , and  $\exists \alpha > 0$ ,  $\lambda \geq 0$  such that

$$a(u, u) + \lambda \|u\|_{L^2}^2 \geq \alpha \|u\|_{H^1}^2, \quad (\text{Weak Coercivity})$$

and  $\mathcal{F}(v)$  is continuous, then the weak formulation is WP. *[When  $\lambda = 0$ , we have strong coercivity.]*

**Example 4.1.3**

For elliptic problem,

$$-\mu u'' + \sigma u = f \quad \text{is WP requires } \sigma \geq 0.$$

For parabolic problem,

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} + \sigma u = f \quad \text{is WP } \forall \sigma.$$

So, time-dependent problem is more likely to be WP.

- What space we want to use for time?

1.  $\int_0^T \int_{\Omega} uv$ . If  $u = v$ ,

$$\int_0^T \frac{d}{dt} \|u\|_{L^2}^2 = \|u\|_{L^2}^2(T) - \|u\|_{L^2}^2(0)$$

Space:  $L^\infty(0, T; L^2(\Omega)) \equiv L^\infty(L^2)$ .

2.  $\int_0^T \int_{\Omega} \nabla u \nabla v$ . If  $u = v$ ,

$$\int_0^T \|u\|_{H^1}^2(t)$$

Space:  $L^2(0, T; H_0^1(\Omega)) \equiv L^2(H_0^1)$ .

**Theorem 4.1.4**

Suppose  $f \in L^2(L^2)$  and  $u_0 \in L^2(\Omega)$ . If  $a(u, v)$  is weakly coercive, then there exists a unique  $u \in L^2(H_0^1) \cap L^\infty(L^2)$ , solution to the problem.

**Theorem 4.1.5**

$$\left( \frac{\partial u}{\partial t}, v \right) + a(u, v) = \mathcal{F}(v) \quad (\text{P})$$

Suppose  $u(\mathbf{x}, 0) = u_0(\mathbf{x}) \in L^2(\Omega)$ ,  $f \in L^2(L^2)$ . If

- $a(\cdot, \cdot)$  bilinear, continuous, weakly coercive,

$$\exists \alpha, \lambda > 0, \quad a(u, u) + \lambda \|u\|_{L^2}^2 \geq \alpha \|u\|^2.$$

- $\mathcal{F}(v)$  is continuous linear functional.

Then, (P) is WP.

**Example 4.1.6**

$$\begin{aligned}
w &= e^{-\lambda t} u \\
\frac{\partial w}{\partial t} &= e^{-\lambda t} \frac{\partial u}{\partial t} - \lambda e^{-\lambda t} u \\
e^{-\lambda t} \frac{\partial u}{\partial t} &= \frac{\partial w}{\partial t} - \lambda e^{-\lambda t} u \\
e^{-\lambda t} \left( \frac{\partial u}{\partial t}, v \right) + \left( e^{-\lambda t} u, v \right) &= e^{-\lambda t} \mathcal{F}(v) \\
\left( \frac{\partial w}{\partial t} + \lambda w, v \right) + a(w, v) &= \mathcal{G}(v) \\
\left( \frac{\partial w}{\partial t}, v \right) + \underbrace{\lambda(w, v) + a(w, v)}_{\tilde{a}(w, v) \geq \alpha \|w\|_{H^1}^2} &= \mathcal{G}(v)
\end{aligned}$$

**A Priori Bounds**

$$\left( \frac{\partial u}{\partial t}, v \right) + a(u, v) = \mathcal{F}(v).$$

Take  $v = u$ .

- If  $\Omega$  is time independent,

$$\int_{\Omega} \frac{\partial u}{\partial t} u = \frac{1}{2} \int_{\Omega} \frac{\partial}{\partial t} = \frac{1}{2} \frac{d}{dt} \|u\|_{L^2}^2.$$

- By weak coercivity,

$$a(u, u) \geq \alpha \|u\|_{H^1}^2 + \lambda \|u\|_{L^2}^2.$$

- $\mathcal{F}(u) = \int f u \leq \|f\|_{L^2} \|u\|_{L^2}$

So,

$$\begin{aligned}
\frac{1}{2} \int_0^T \frac{d}{dt} \|u\|_{L^2}^2 + \alpha \int_0^T \|u\|_{H^1}^2 + \lambda \int_0^T \|u\|_{L^2}^2 &\leq \int_0^T \|f\|_{L^2} \|u\|_{L^2} \\
ah &\leq \varepsilon a^2 + \frac{1}{4\varepsilon} b^2 \quad [\text{Young Inequality, } 0 \leq \left( \sqrt{\varepsilon} a - \frac{1}{2\sqrt{\varepsilon}} b \right)^2]
\end{aligned}$$

$$\begin{aligned}
\frac{1}{2} \frac{d}{dt} \|u\|^2 + \alpha \|u\|_{H^2}^2 + \frac{\lambda}{2} \|u\|_{L^2}^2 &\leq \frac{1}{2\lambda} \|f\|_{L^2}^2 \\
\frac{1}{2} \int_0^T \|u\|^2 + \alpha \int_0^T \|u\|_{H^2}^2 + \frac{\lambda}{2} \int_0^T \|u\|_{L^2}^2 &\leq \frac{1}{2\lambda} \int_0^T \|f\|_{L^2}^2 \\
\|u\|_{L^2}^2(T) + 2\alpha \int_0^T \|u\|_{H^2}^2 + \lambda \int_0^T \|u\|_{L^2}^2 &\leq \underbrace{\frac{1}{\lambda} \int_0^T \|f\|_{L^2}^2}_{\text{data}} + \|u_0\|_{L^2}^2
\end{aligned}$$

**Lemma 4.7 Granwall Lemma**

$$g(t) \leq \alpha + \int_0^t \beta g(\tau) d\tau, \quad \beta \geq 0$$

$$g(t) \implies \alpha + \int_0^t \alpha \beta e^{\int_\tau^t \beta(\xi) d\xi} d\tau.$$

**4.2 Finite Element and Semi-Discretization****Solving Systems of ODEs**

$$\frac{d\mathbf{u}}{dt} = M^{-1}A\mathbf{u} + M^{-1}\mathbf{b}, \quad \mathbf{u} = [u(x_i)]$$

$$\frac{d\mathbf{u}}{dt} = A\mathbf{u} + \mathbf{b} \quad (\text{more generally})$$

- Absolute stability in ODE:

A general ODE can be written as

$$\frac{dy}{dt} = f(t, y).$$

Consider only the ODE

$$\frac{dy}{dt} = \lambda y, \quad \lambda < 0.$$

Then,  $y(t) \rightarrow 0$  as  $t \rightarrow \infty$ . So, the numerical solution  $y_{\Delta t}$  does the same.

1. Backward Euler (BE):

$$\frac{y^{n+1} - y^n}{\Delta t} \approx \lambda y^{n+1} \implies y^{n+1}(1 + |\lambda|\Delta t) = y^n$$

$$y^{n+1} = \frac{1}{1 + |\lambda|\Delta t} y^n$$

$$y^{n+1} = \left( \frac{1}{1 + |\lambda|\Delta t} \right)^{n+1} y^0 \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

2. Forward Euler (FE):

$$\frac{y^{n+1} - y^n}{\Delta t} \approx \lambda y^n \implies y^{n+1} = (1 + \lambda\Delta t)y^n.$$

For absolute stability, we require

$$|1 + \lambda\Delta t| < 1$$

$$|1 - |\lambda|\Delta t| < 1$$

$$-1 < 1 - |\lambda|\Delta t < 1$$

$$\Delta t < \frac{2}{|\lambda|}$$

- Moving to systems:

$$\frac{dy}{dt} = Ay(+b), \quad \lambda = \text{eig}(A), \quad \text{Re}(\lambda) < 0.$$

Diagonalization:  $TAT^{-1} = D$ . Then,  $D = TAT^{-1}$ . Denote  $\mathbf{w} = Ty$ . We have

$$\frac{d\mathbf{w}}{dt} = D\mathbf{w}.$$

So, we require

$$\Delta t < \frac{2}{\max |\lambda_i|}.$$

- Absolute stability has nothing to do with accuracy.

### Semi-Discretization

$$\left( \frac{\partial u}{\partial t}, v \right) + a(u, v) = \mathcal{F}(v).$$

Let  $V_h \subset H_0^1$ ,  $u_h = \sum u_j(t)\varphi_j(\mathbf{x})$  (separation of variable), and  $v_h = \varphi_j(\mathbf{x})$ .

The numerical problem is

$$\begin{aligned} & \left( \frac{\partial u_h}{\partial t}, v_h \right) + a(u_h, v_h) = \mathcal{F}(v_h) \\ & \left( \sum \frac{du_j}{dt} \varphi_j, \varphi_i \right) + \underbrace{\sum u_j(t) a(\varphi_j, \varphi_i)}_{\mathbf{u}(t) A} = \underbrace{\mathcal{F}(\varphi_i)}_{\mathbf{f}} \\ & \underbrace{\frac{d}{dt} u_j(t)}_{\frac{d\mathbf{u}}{dt}} \underbrace{(\varphi_j, \varphi_i)}_{\substack{M \\ \text{mass matrix}}} + A\mathbf{u}(t) = \mathbf{f} \\ & M \frac{d\mathbf{u}}{dt} + A\mathbf{u} = \mathbf{f} \\ & \frac{d\mathbf{u}}{dt} = \underbrace{-M^{-1}A}_{\substack{\text{negative} \\ \text{eigenvalue}}} \mathbf{u} + M^{-1}\mathbf{f} \end{aligned}$$

**Convergence** Define  $e_h = u - u_h$ .

$$\left( \frac{\partial u_h}{\partial t}, v_h \right) + a(u_h, v_h) = 0.$$

Assume strong coercivity,

$$\begin{aligned} \alpha \|e_h\|_{H^1}^2 & \leq a(e_h, e_h) = a(e_h, u - u_h) \\ & = a(e_h, u - w_h) + a(e_h, \underbrace{w_h - u_h}_{=0, \in V_h}) \\ & = \underbrace{a(e_h, u - w_h)}_{\textcircled{1}} - \underbrace{\left( \frac{\partial e_h}{\partial t}, w_h - u_h \right)}_{\textcircled{2}} \end{aligned}$$

- Term ①:

$$\begin{aligned} a(e_h, u - w_h) &\leq \gamma \|e_h\| \|u - w_h\| && \text{[Continuity]} \\ &\leq \frac{\alpha}{2} \|e_h\|_{H^1}^2 + C \|u - w_h\|_{H^1}^2 && \text{[Young Inequality]} \end{aligned}$$

- Term ②:

$$\begin{aligned} -\left(\frac{\partial e_h}{\partial t}, w_h - u_h\right) &= -\left(\frac{\partial e_h}{\partial t}, w_h - u + u - u_h\right) \\ &= -\left(\frac{\partial e_h}{\partial t}, w_h - u\right) - \left(\frac{\partial e_h}{\partial t}, e_h\right) \\ &= -\left(\frac{\partial e_h}{\partial t}, w_h - u\right) - \frac{1}{2} \frac{d}{dt} \|e_h\|_{L^2}^2 \end{aligned}$$

Combined, we have

$$\begin{aligned} \alpha \|e_h\|_{H^1}^2 &\leq \frac{\alpha}{2} \|e_h\|_{H^1}^2 + C \|u - w_h\|_{H^1}^2 - \left(\frac{\partial e_h}{\partial t}, w_h - u\right) - \frac{1}{2} \frac{d}{dt} \|e_h\|_{L^2}^2 \\ \frac{1}{2} \frac{d}{dt} \|e_h\|_{L^2}^2 + \frac{\alpha}{2} \|e_h\|_{H^1}^2 &\leq C \|u - w_h\|_{H^1}^2 - \left(\frac{\partial e_h}{\partial t}, w_h - u\right) \\ \int_0^T \frac{1}{2} \frac{d}{dt} \|e_h\|_{L^2}^2 + \int_0^T \frac{\alpha}{2} \|e_h\|_{H^1}^2 &\leq C \int_0^T \|u - w_h\|_{H^1}^2 - \int_0^T \left(\frac{\partial e_h}{\partial t}, w_h - u\right). \end{aligned}$$

So, the rate of convergence is of order  $\sim \mathcal{O}(h^2)$ .

### ⊖-Method

$$\int_{t^n}^{t^{n+1}} \text{RHS} = \Delta t \theta \text{RHS}(t^{n+1}) + \Delta t (1 - \theta) \text{RHS}(t^n), \quad \theta \in (0, 1].$$

- Common choices of  $\theta$ :

1.  $\theta = 1$ : Backward Euler
2.  $\theta = \frac{1}{2}$ : Crank-Nicolson/Trapezoidal Rule

$$\int_{t^n}^{t^{n+1}} f(t) = \frac{f(t^{n+1}) + f(t^n)}{2} \Delta t$$

3.  $\theta = 0$ : Forward Euler

- Stability Analysis:

$$\mathbf{u}^{n+1} = \mathbf{u}^n + \theta \Delta t (-M^{-1}A) \mathbf{u}^{n+1} + (1 - \theta) \Delta t (-M^{-1}A) \mathbf{u}^n$$

Use diagonalization and substitution, we get

$$(1 - \theta) \lambda_i w_i^{n+1} = (1 + (1 - \theta) \Delta t \lambda_i) w_i^n$$

1.  $\theta \geq \frac{1}{2}$ : Unconditionally stable.

2.  $\theta < \frac{1}{2}$ : Conditionally stable. Condition is given by

$$\Delta t < \frac{2}{|\lambda_i|(1-2\theta)}$$

**Theorem 4.2.1**

$u_n$  of FEM and  $\Theta$  method is such that

$$\|u - u_h\|_{L^\infty(L^2)} + \|u - u_h\|_{L^2(H^1)} \leq C(\Delta t^{p(\theta)} + h^r),$$

where

$$p(\theta) = \begin{cases} 2, & \theta = \frac{1}{2} \\ 1, & \text{o/w} \end{cases}$$

and  $r = \min \{\text{deg}, s\}$ .

**Summary**

$$\frac{\partial u}{\partial t} - \mu \Delta u = f \quad \text{on } \Omega.$$

$$\int_{\Omega} \frac{\partial u}{\partial t} v + \mu \int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v$$

$$\frac{d}{dt} \int_{\Omega} u v + \mu \int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v$$

Define  $\mathbf{u}(\mathbf{x}, t) = \sum u_j(t) \rho_j(\mathbf{x})$ ,  $M = \int_{\Omega} \varphi_i \varphi_j$ , and  $K = \mu \int_{\Omega} \nabla \varphi_i \nabla \varphi_j$ . We have

$$M \frac{d\mathbf{u}}{dt} + K\mathbf{u} = \mathbf{b}.$$

With  $\Theta$ -method:

$$M \frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{\Delta t} + \theta K \mathbf{u}^{n+1} + (1-\theta) K \mathbf{u}^n = \theta \mathbf{b}^{n+1} + (1-\theta) \mathbf{u}^n$$

$$\underbrace{(M + \Delta \theta K)}_{\text{left matrix}} \mathbf{u}^{n+1} = \underbrace{(M - \Delta t(1-\theta)K)}_{\text{right matrix}} \mathbf{u}^n + \Delta t \theta \mathbf{b} + (1-\theta) \Delta t \mathbf{u}^n$$

- When  $\theta < \frac{1}{2}$ ,  $\Theta$ -method is conditionally stable. The condition is given by

$$\Delta t < \frac{2}{\rho(1-2\theta)},$$

where  $\rho = \max(|\text{eig}(M^{-1}K)|) \sim \mathcal{O}(h^{-2})$ . So,  $\Delta t \sim \mathcal{O}(h^2)$ . [So, if we want to refine  $h$ , we need

to refine  $\Delta$  at the same time. For example,

$$(\Delta, h) \rightarrow \left( \frac{\Delta t}{4}, \frac{h}{2} \right).$$

] This is annoying! So, we don't use explicit method a lot for parabolic problems.

- On the other hand, running full implicit (FI) solver requires three loops. This is inefficient. What we can do is to run a semi-implicit (SI) solver that only requires two loops.
- However, semi-implicit is not unconditionally stable. The condition is given by  $\Delta t < \mathcal{O}(h)$ .  
[Generally, this is easy to fulfill.]

### 4.3 Space-Time Finite Element (Space-Time FEM)

$$\frac{\partial u}{\partial t} - \mu \frac{\partial^2 u}{\partial x^2} = f.$$

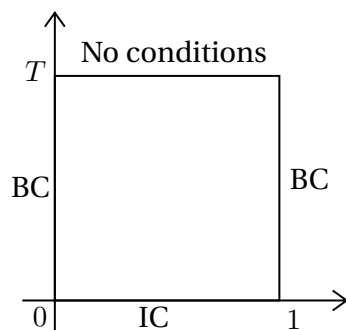
**Attempt:** think of this as a 2D problem:

$$-\mu \frac{\partial^2 u}{\partial x^2} - 0 \cdot \frac{\partial^2 u}{\partial y^2} + 0 \cdot \frac{\partial u}{\partial x} + 1 \cdot \frac{\partial u}{\partial y} = f.$$

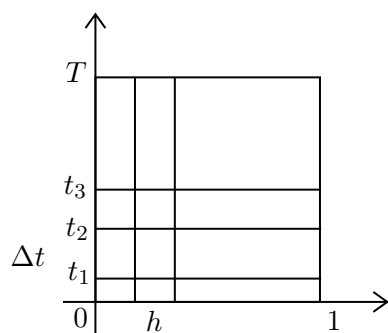
We are back in advection-diffusion problem.

**Problem with this attempt:**

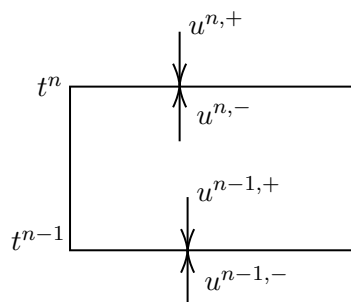
- Time might be huge, so it will be a large (and imbalanced) mesh.
- The physics is different. For time, we have an initial condition, instead of a boundary condition.
- The problem is not WP (time is not reversible).



**Decouple the time slab:** For each time, we can have different space discretization level



- Problem: we will have “jumping” solutions.



See exercise book for reference. Eventually, we will get backward Euler.

- Nowadays, people do not use space-time FEM anymore because it is too complicated.

## 5 Least-Square FEM and PINNs

### 5.1 Least-Square FEM (LS-FEM)

$$\begin{cases} -\Delta u = f & \text{in } \Omega \\ u(\partial\Omega) = g. \end{cases}$$

#### Our Approach:

- Pointwise FD:

$$-\underbrace{\frac{u_{i+1,j} - 2u_{i,j} + u_{i-1,j}}{\Delta x^2}}_{\partial^2 u / \partial x^2} - \underbrace{\frac{u_{i,j+1} - 2u_{i,j} + u_{i,j-1}}{\Delta y^2}}_{\partial^2 u / \partial y^2} = f_{i,j}$$

- Weak formulation (Variational): FEM

$$\frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} f u \implies \text{Find } u \in V \text{ s.t. } \int_{\Omega} \nabla u \nabla v = \int_{\Omega} f v \quad \forall v \in V.$$

- Minimization problem:

$$\min_u \frac{1}{2} \int_{\Omega} |\nabla u|^2 - \int_{\Omega} f u,$$

where  $u \in \text{Deep Neural Network} \rightarrow \text{Deep Ritz}$

- Least Square:

$$J(u; f, g) = w_1 \| -\Delta u - f \|_{L^2(\Omega)}^2 + w_2 \| u - g \|_{\partial\Omega},$$

where  $f \in L^2(\Omega)$  and  $u \in H^2$ . The problem is find  $u$  such that  $J(u; f, g) \leq J(u + \delta v; f, g)$ .

1. Let  $u \in \mathcal{L}_g \oplus H_0^1 \cap H^2$  (when  $g = 0$ ,  $u \in H_0^1 \cap H^2$ ). Consider  $J = \frac{1}{2} \| \Delta u + f \|^2$ .

**Goal:**  $J(u; f) \leq J(u + \delta v; f)$ .

Note that

$$\begin{aligned} J(u + \delta v; f) - J(u; f) &= \frac{1}{2} \int_{\Omega} (\Delta u + \delta \Delta v + f)(\Delta u + \delta \Delta v + f) - \frac{1}{2} \int_{\Omega} (\Delta u + f)(\Delta u + f) \\ &= \int_{\Omega} \Delta v (\Delta u + f) + \delta (\Delta v)^2 \end{aligned}$$

So,

$$\begin{aligned} \lim_{\delta \rightarrow 0} \frac{1}{\delta} (J(u + \delta v; f) - J(u; f)) &= \lim_{\delta \rightarrow 0} \int_{\Omega} \Delta v (\Delta u + f) + \delta (\Delta v)^2 \\ &= \int_{\Omega} \Delta v (\Delta u + f) \\ &= 0. \end{aligned}$$

So, Euler-Lagrangian condition is

$$\int_{\Omega} (\Delta u + f) \Delta v = 0 \implies \int_{\Omega} \Delta u \Delta v = \int_{\Omega} f \Delta v.$$

Numerically, with FEM,  $u_h \in V_h$ , and

$$\int_{\Omega} \Delta u_h \Delta v_h = \int_{\Omega} f \Delta v_h.$$

But, the problem is  $V_h \not\subset V$ . So, this method is not practical.

2. Let  $X = H^2 \cap H_0^1$  and  $Y = L^2$ . Norm equivalence gives

$$\alpha_1 \|u\|_X^2 \leq \underbrace{J(u; 0)}_{\|\Delta u + f\|_Y} \leq \alpha_2 \|u\|_X^2$$

Then,

$$\alpha_1 \|u - u_h\|_X^2 \leq \|\Delta u + \Delta u_h\|_Y^2 = \|\Delta u + f\|_Y^2 = J(u; f) \leq \alpha_2 \|u - u_h\|_X^2.$$

So,  $\exists, \alpha_1, \alpha_2 > 0$  such that

$$\alpha_1 \|u\|_{H^2 \cap H_0^1}^2 \leq \|\Delta u\|_{L^2}^2 \leq \alpha_2 \|u\|_{H^2 \cap H_0^1}^2.$$

**Proof 1.** For the second inequality,

$$\|u\|_{H^2 \cap H_0^1}^2 = \|u\|_{L^2}^2 + \|\nabla u\|_{L^2}^2 + \|D^2 u\|^2,$$

where  $\|D^2 u\|^2 = \|\Delta u\|^2 + \left\| \frac{\partial^2 u}{\partial x \partial y} \right\|^2$ .

For the first inequality, it depicts “coercivity.”

$$\begin{aligned} \int_{\Omega} u \Delta u &= \int_{\partial \Omega} \underbrace{(\nabla u \cdot \mathbf{n})u}_{=0} - \int \nabla u \nabla v \\ \alpha \|u\|_{L^2}^2 &\leq \|\nabla u\|_{L^2}^2 = \left| \int \Omega u \Delta u \right| \leq \|u\|_{L^2} \|\Delta u\|_{L^2} \leq \frac{\alpha}{2} \|u\|_{L^2}^2 + C \|\Delta u\|_{L^2}^2 \\ \beta \|u\|_{L^2}^2 &\leq \|\Delta u\|_{L^2}^2 \implies \|\nabla u\|_{L^2}^2 \leq \gamma \|\Delta u\|_{L^2}^2 \end{aligned}$$

Q.E.D. ■

**Pavel Bacher and Max Gurnzburger** Consider the first order version of the same problem:

$$\begin{cases} \sigma - \nabla u = 0 \\ \nabla \cdot \sigma + f = 0, \end{cases}$$

where  $u \in H_0^1$ ,  $\sigma \in L^2$ , and  $\nabla \cdot \sigma \in L^2$ , and

$$\frac{\partial \sigma_1}{\partial x_1} + \frac{\partial \sigma_2}{\partial x_2} + \frac{\partial \sigma_3}{\partial x_3} \in L^2 \implies \sigma \in H(\Omega; \text{div}).$$

- Variational formulation:

$$J = \frac{1}{2} \|\boldsymbol{\sigma} - \nabla u\|_{L^2}^2 + \|\nabla \cdot \boldsymbol{\sigma} + f\|_{L^2}^2.$$

Find  $\boldsymbol{\sigma}, u$  such that  $J(\boldsymbol{\sigma}, u; f) \leq J(\boldsymbol{\rho}, v; f) \quad \forall \boldsymbol{\rho} \in H(\text{div}), v \in H^1$ . The first order condition is given by

$$\begin{cases} \int (\boldsymbol{\sigma} - \nabla u) \boldsymbol{\rho} + \int (\nabla \cdot \boldsymbol{\rho})(\nabla \cdot \boldsymbol{\sigma} + f) = 0 \\ \int (\boldsymbol{\sigma} - \nabla u) \cdot \nabla v = 0. \end{cases}$$

We can use linear FEM:  $V = H^1(\Omega)$ ,  $Q = H(\text{div})$ . Choose  $Q_h \subset Q$  and  $V_h \subset V$ .

- Do we have norm equivalence? Yes.

$$\begin{aligned} \|\boldsymbol{\sigma}\|_{H(\text{div})}^2 &= \|\boldsymbol{\sigma}\|_{L^2}^2 + \|\nabla \cdot \boldsymbol{\sigma}\|_{L^2}^2 \\ \alpha_1 \left( \|u\|_{H^1}^2 + \|\boldsymbol{\sigma}\|_{H(\text{div})}^2 \right) &\leq J(\boldsymbol{\sigma}, u; 0) \leq \alpha_2 \left( \|u\|_{H^1}^2 + \|\boldsymbol{\sigma}\|_{H(\text{div})}^2 \right). \end{aligned}$$

- Benefits of LS formulation, and compare it with the mixed method.

Recall that we had mixed method

$$\begin{cases} \int_{\Omega} \boldsymbol{\sigma} \cdot \boldsymbol{\rho} + u \nabla \cdot \boldsymbol{\rho} = 0 \\ \int_{\Omega} (\nabla \cdot \boldsymbol{\sigma} + f)v = 0. \end{cases}$$

1. Why we don't go with this formulation?

We require an inf-sup condition for mixed method, but we don't need this condition anymore with the LS formulation.

2. Why?

The solution to the LS problem is the normal equations

$$A^T A \mathbf{u} = A^T \mathbf{b},$$

where  $A^T A$  is SPD. So, we don't need inf-sup condition.

3. Some practical questions:

- We include some additional quadrature error:

$$\int \approx \sum.$$

- With boundary conditions,

$$\|\text{Residual}\|_{L^2}^2 + \|\text{Residual of BCs}\|_{L^2}^2$$

We have no Poincaré Inequality, and norm equivalence is much more harder to obtain. So, we don't have control in  $J$ , which means no control in error.

- What if we show advection-diffusion-reaction problem with LS?

$$\mathcal{L}u = -\mu\Delta u + \beta \cdot \nabla u + \gamma u$$

$$\int_{\Omega} (\mathcal{L}u - f)\mathcal{L}v = 0$$

- (a) We've seen this in GaLS, in which we consider it elementwise and it served as a stabilization term.
- (b) With LS, assume  $u_h \in H^2 \cap H_0^1$ . The normal equations are SPD. So, we have no oscillations.

The first order formulation is

$$\begin{cases} \sigma - (-\mu\nabla u + \beta u) = 0 \\ \nabla \cdot \sigma + f = 0 \end{cases}$$

We can also prove norm equivalence.

- A curl free formulation for Poisson equation:

$$\begin{cases} \sigma - \nabla u = 0 \\ \nabla \cdot \sigma + f = 0 \\ \nabla \times \sigma = 0 \end{cases}$$

1.  $(u, \sigma) \in H_0^1 \times [H^1]^d$
2.  $\nabla \times \sigma = 0$  is always true in theory, but it may not be true numerically.

## 5.2 Physics-Informed Neural Networks (PINNs)

**Set-up:**

$$-\Delta u = f \implies -\Delta u - f = 0$$

$$\sum w_i |\Delta u(x_i) + f(x_i)|^2 = 0$$

Consider  $u$  a neural network  $\text{NN}(\theta)$ .

- Adding BC:  $u = g$  on  $\partial\Omega$ .

$$\min w_1 \sum w_i |\Delta u(x_i) + f(x_i)|^2 + w_2 \sum \tilde{w}_j |u(x_j) - g(x_j)|^2$$

- Strength: We don't need geometry information. We don't need a mesh! (But it is always better if we could have some geometry information...)
- It is natural to use PINNs for Data Assimilation problems. We just add an extra term to measure data mesh mismatch.

**Numerical Evidence:**

- The architecture of the NN has a limited impact on accuracy and efficiency
- Strong enforcement of BC works better than the weak one (including the BC term in the minimization problem)
- First-order formulations work better and with ReLU networks.
- No need for inf-sup for the mixed formulation.
- No need of special formulation for the stabilization.
- Great potential for the Data Assimilation.

**LS Formulation of PINNs:**

- $\mathcal{L} : X \rightarrow Y$ ,  $\mathcal{L}u = f$  in  $\Omega$ ,  $u = 0$  on  $\partial\Omega$ .

Assume BC is strongly encoded in  $X$  (we manually set the boundary entries to 0).

- **Hypothesis: (Norm Equivalence)**  $\exists$  two positive constants  $\alpha_1, \alpha_2$  such that  $\forall v \in X$ ,

$$\alpha_1 \|v\|_X^2 \leq \|\mathcal{L}v\|_Y^2 \leq \alpha_2 \|v\|_X^2.$$

- Define  $\mathcal{J}(u, f) = \|\mathcal{L}u - f\|_Y^2$ . The problem is reformulated as finding  $u \in X$  such that

$$u = \arg \min_{v \in X} \mathcal{J}(v, f).$$

- In PINNs, replace  $\mathcal{J}(v, f)$  with

$$\mathcal{J}_N(v, f) = \frac{|\Omega|}{N} \sum_{i=1}^N |\mathcal{L}v(x_i) - f(x_i)|^2.$$

For  $Y = L^2(\Omega)$ ,  $\mathcal{J}_N(v, f)$  is a quadrature approximation of  $\mathcal{J}(v, f)$ .

- Ceà Lemma in PINNs:

If the adopted NN generates functions in  $X_{\text{NN}} \subset X$ , then,

$$\|u_{\text{PINN}} - u\|^2 \leq C_1 + C_2 N^{-s/2} + \frac{\alpha_2}{\alpha_1} \inf_{u_{\text{NN}} \in X_{\text{NN}}} \|u - u_{\text{NN}}\|_X,$$

where  $s$  depends on the regularity of the residual  $\mathcal{L}u - f$  in  $Y$  and the distribution of the quadrature nodes, and  $C_1$  depends on the accuracy of the optimization algorithm.

**Consequences for the Poisson Problem:**

- The second order (traditional) formulation works ( $\alpha_1 \|v\|_{H^2}^2 \leq \|\Delta v\|_{L^2}^2 \leq \|v\|_{H^2}^2$ ) only if the activation function are regular enough. [*ReLU does not work, but tanh works.*]
- The Poincaré inequality is critical: A weak enforcement of the BC suffers because the coercivity control works in weaker spaces ( $H^{1/2}(\Omega)$ ).

- The first order traditional div-grad formulation works for NN generating functions in  $H^1$  and  $H(\text{div})$ . So, ReLU works here!

$$u \in H_0^1(\Omega), \quad \sigma \in H(\Omega; \text{div}) \text{ s.t. } \begin{cases} \nabla u - \sigma = 0 \\ -\nabla \cdot \sigma = f \end{cases}$$

$$\|v\|_{H^1} + \|\tau\|_{H(\text{div})} \leq \|\nabla v - \tau\|_{L^2} + \|\nabla \cdot \tau\|_{L^2}$$

- The (redundant) curl-free formulation works with a stronger control ( $X = H_0^1 \cap [H^1]^d$ ):

$$u \in H_0^1(\Omega), \quad \sigma \in H(\Omega; \text{div}) \text{ s.t. } \begin{cases} \nabla u - \sigma = 0 \\ \nabla \times \sigma = 0 \\ -\nabla \cdot \sigma = f \end{cases}$$

$$\|v\|_{H^1} + \|\tau\|_{H^1} \leq \|\nabla v - \tau\|_{L^2} + \|\nabla \cdot \tau\|_{L^2} + \|\nabla \times \tau\|_{L^2}.$$

### Remark 1. (Comments on LS-FEM and PINNs).

- For LS-FEM, we have the norm equivalence

$$\alpha_1 \|u\|_X^2 \leq J(u; 0, 0) \leq \alpha_2 \|u\|_X^2 \quad \forall u \in X$$

1. With this norm equivalence, we have control of error.
2. If we don't have Poincaré Inequalities, it is hard to prove norm equivalence. We should impose strong enforcement of BC to have this norm equivalence result.

- Source of error:

1. (LS)-FEM: Discretization Error, Consistency, Quadrature Error, Linear Algebra Error.  
dominating
2. PINNs: Discretization Error, Quadrature Error, Optimization Error.  
dominating

- PINNs are excellent in solving inverse problems but questionable in solving forward problems.

## 6 Hyperbolic Problems

### 6.1 Finite Differences

- We discretize the semi-plane  $\{(x, t) : 0 < t, -\infty < x < \infty\}$  with time step  $\Delta t$  and mesh size  $h$ .
- The grid nodes  $(x_j, t^n)$  denotes where  $x_j = jh$ ,  $j \in \mathbb{Z}$  and  $t^n = n\Delta t$ ,  $n \in \mathbb{N}$ .
- Denote the numerical solution  $u_j^n$  that approximates  $u(x_j, t^n)$  for any  $j$  and  $n$ .
- For simplicity, set  $\lambda = \frac{\Delta t}{h}$ .

#### 6.1.1 Conservation Laws

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0.$$

- Backward Euler/Centered (BE/C): Implicit method.

$$\left. \frac{\partial u}{\partial t} \right|_{t^{n+1}} \approx \frac{u_j^{n+1} - u_j^n}{\Delta t} \quad \text{and} \quad a \frac{\partial u^{n+1}}{\partial x} \approx \frac{a}{2} \frac{u_{j+1}^{n+1} - u_{j-1}^{n+1}}{h} \quad \leftarrow \text{centered difference}$$

$$\boxed{u_j^{n+1} + \frac{\lambda}{2} a (u_{j+1}^{n+1} - u_{j-1}^{n+1}) = u_j^n.}$$

At each time step, the solution is obtained by solving a linear system.

- Upwind (UPW): Explicit method.

$$\begin{aligned} a \frac{\partial u}{\partial x} &= \frac{a}{2} \frac{u_{j+1}^n - u_{j-1}^n}{h} - \frac{|a|}{2} \cdot h \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2} \\ u_j^{n+1} &= u_j^n - \frac{\lambda}{2} a (u_{j+1}^n - u_{j-1}^n) + \frac{\lambda}{2} |a| (u_{j+1}^n - 2u_j^n + u_{j-1}^n) \end{aligned}$$

- Lax-Wendroff (LW): Explicit method.

This method is based on Taylor expansion.

$$\begin{aligned} u_j^{n+1} &= u_j^n + \Delta t \left. \frac{\partial u}{\partial t} \right|_{t^n} + \frac{\Delta t^2}{2} \left. \frac{\partial^2 u}{\partial t^2} \right|_{t^n} + \dots \\ \frac{\partial u}{\partial t} &= -a \frac{\partial u}{\partial x} \implies \frac{\partial^2 u}{\partial x \partial t} = -a \frac{\partial^2 u}{\partial x^2} \implies \frac{\partial^2 u}{\partial t^2} = -a \frac{\partial^2 u}{\partial t \partial x} = -a \left( -a \frac{\partial^2 u}{\partial x^2} \right) = a^2 \frac{\partial^2 u}{\partial x^2} \end{aligned}$$

So,

$$\begin{aligned} u_j^{n+1} &= u_j^n + \Delta t \left( -a \left. \frac{\partial u}{\partial x} \right|_{t^n} \right) + \frac{\Delta t^2}{2} a^2 \left( \left. \frac{\partial^2 u}{\partial x^2} \right|_{t^n} \right) \\ &= u_j^n - \frac{\Delta t}{2} a \frac{u_{j+1}^n - u_{j-1}^n}{h} + \frac{\Delta t^2}{2} a^2 \frac{u_{j+1}^n - 2u_j^n + u_{j-1}^n}{h^2} \\ &= u_j^n - \frac{\lambda}{2} a (u_{j+1}^n - u_{j-1}^n) + \frac{\lambda^2}{2} a^2 (u_{j+1}^n - 2u_j^n + u_{j-1}^n). \end{aligned}$$

- For UPW and LW, they are explicit methods because solution in  $t^{n+1}$  depends on the solutions at the previous time steps, so no linear systems need to be solved.
- The additional term in UPW and LW can be considered as some diffusive term (which are purely numerical). They can also be viewed as stabilizing terms.
- Generalization of UPW/LW to more complicated PDEs: Consider the generalized conservation law:

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

UPW and LW can be written in the form

$$\boxed{u_j^{n+1} = u_j^n - \lambda \left( H_{j+1/2}^n - H_{j-1/2}^n \right)},$$

with a suitable choice of  $H_{j+1/2} = H(u_j, u_{j+1})$ :

1. UPW:

$$H_{j+1/2} = \frac{1}{2} [a(u_{j+1} + u_j) - |a|(u_{j+1} - u_j)]$$

2. LW:

$$H_{j+1/2} = \frac{1}{2} [a(u_{j+1} + u_j) - \lambda a^2(u_{j+1} - u_j)].$$

$H(\cdot, \cdot)$  is called *numerical flux*.

### 6.1.2 Methods for Linear System

$$\frac{\partial \mathbf{u}}{\partial t} + A \frac{\partial \mathbf{u}}{\partial x} = \mathbf{0}.$$

- Generalization of the scalar case:

$$\mathbf{u}_j^{n+1} = \mathbf{u}_j^n - \lambda \left( \mathbf{H}_{j+1/2}^n - \mathbf{H}_{j-1/2}^n \right),$$

where  $\mathbf{u}_j^n$  is the approximation vector of  $\mathbf{u}(x_j, t^n)$  and  $\mathbf{H}_{j+1/2}$  is the *vector numerical flux*. The numerical flux is obtained by generalization of the scalar case.

- Upwind for Systems (UPW-S):

$$\mathbf{u}_j^{n+1} = \mathbf{u}_j^n - \frac{\lambda}{2} A (\mathbf{u}_{j+1}^n - \mathbf{u}_{j-1}^n) + \frac{\lambda}{2} |A| (\mathbf{u}_{j+1}^n - 2\mathbf{u}_j^n + \mathbf{u}_{j-1}^n).$$

The numerical flux is

$$\mathbf{H}_{j+1/2} = \frac{1}{2} [A(\mathbf{u}_{j+1} + \mathbf{u}_j) - |A|(\mathbf{u}_{j+1} - \mathbf{u}_j)].$$

- Lax-Wendroff for Systems (LW-S):

$$\mathbf{u}_j^{n+1} = \mathbf{u}_j^n - \frac{\lambda}{2} A (\mathbf{u}_{j+1}^n - \mathbf{u}_{j-1}^n) + \frac{\lambda^2}{2} A^2 (\mathbf{u}_{j+1}^n - 2\mathbf{u}_j^n + \mathbf{u}_{j-1}^n),$$

and the associated numerical flux is

$$\mathbf{H}_{j+1/2} = \frac{1}{2} [A(\mathbf{u}_{j+1} - \mathbf{u}_j) - \lambda A^2(\mathbf{u}_{j+1} - \mathbf{u}_j)].$$

### 6.1.3 Wave Equation

$$\frac{\partial^2 u}{\partial t^2} - \gamma^2 \frac{\partial^2 u}{\partial x^2} = 0.$$

- Wave equation can be written as system of conservation laws. So, we can use UPW or LW for systems.

- Leap-Frog:

$$\begin{aligned} \frac{u_j^{n+1} - 2u_j^n + u_j^{n-1}}{\Delta t^2} - \gamma^2 \frac{u_{j-1}^n - 2u_j^n + u_{j+1}^n}{h^2} &= 0 \\ u_j^{n+1} - 2u_j^n + u_j^{n-1} &= (\gamma\lambda)^2 (u_{j+1}^n - 2u_j^n + u_{j-1}^n). \end{aligned}$$

- (Optimal Newmark):

$$u_j^{n+1} - 2u_j^n + u_j^{n-1} = \frac{(\gamma\lambda)^2}{4} (w_j^{n-1} + 2w_j^n + w_j^{n+1}),$$

$$\text{with } w_j^n = u_{j+1}^n - 2u_j^n + u_{j-1}^n.$$

## 6.2 Analysis of FD Methods

### 6.2.1 Consistency

**Definition 6.2.1 (Local Truncation Error/LTE).** *Local Truncation Error (LTE)* of a numerical scheme is the residual obtained when the exact solution is plugged into the scheme.

From the classical Taylor expansions, we have

$$\begin{aligned} \left. \frac{\partial u}{\partial t} \right|_{(x_j, t^{n+1})} &= \frac{u(x_j, t^{n+1}) - u(x_j, t^n)}{\Delta t} + \mathcal{O}(\Delta t) \\ \left. \frac{\partial u}{\partial x} \right|_{(x_j, t^{n+1})} &= \frac{u(x_{j+1}, t^{n+1}) - u(x_{j-1}, t^{n+1})}{2h} + \mathcal{O}(h^2) \end{aligned}$$

The scalar problem  $\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0$  gives

$$0 = \left. \frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} \right|_{(x_j, t^{n+1})} = \underbrace{\frac{u(x_j, t^{n+1}) - u(x_j, t^n)}{\Delta t} + a \frac{u(x_{j+1}, t^{n+1}) - u(x_{j-1}, t^{n+1})}{2h}}_{\text{BE/C scheme}} + \underbrace{\tau_j^{n+1}(\Delta t, h)}_{\text{Truncation error}}$$

**Definition 6.2.2 (Global Truncation Error).** The global truncation error is given by

$$\tau(\Delta t, h) = \max_{j,n} |\tau_j^n|.$$

**Definition 6.2.3 (Consistency).** If  $\tau(\Delta t, h)$  tends to 0 for  $\Delta t$  and  $h$  independently going to 0, we say that the numerical scheme is *consistent*. Suppose the solution of the exact problem is smooth enough. The method is said to be of *order  $p$  in time and of order  $q$  in space*, where  $p, q$  are integers, if we have

$$\tau(\Delta t, h) = \mathcal{O}(\Delta t^p + h^q).$$

### Theorem 6.2.4 Order of Convergence of FD Schemes

- BE/C:  $\mathcal{O}(\Delta t + h^2)$  [ $h^2$  because we used centered difference in space]
- UPW:  $\mathcal{O}(\Delta t + h)$
- LW:  $\mathcal{O}(\Delta t^2 + h^2 + \Delta t h)$  [ $\Delta t^2 + h^2$  is the dominating terms. Usually,  $\Delta t$  and  $h$  scales in the same way. So,  $\Delta t h^2$  is a third order term.]

## 6.2.2 Convergence

**Definition 6.2.5 (Convergence).** The scheme is said to be *convergent* if

$$\lim_{\Delta t, h \rightarrow 0} \max_{j,n} |u(x_j, t^n) - u_j^n| = 0.$$

### Theorem 6.2.6 Lax-Richtmyer's Equivalence Theorem

For linear problems, a consistent method is convergent  $\iff$  it is stable.

Hence, convergence is guaranteed provided that numerical errors are kept under control by acting on  $h$  and  $\Delta t$ .

## 6.2.3 Stability

**Definition 6.2.7 (Stability).** A numerical method is *stable* if for any time  $T$ , there are two constants  $C_T > 0$  (possibly depending on  $T$ , but not on  $\Delta t$  or  $h$ ) and  $\delta_0 > 0$  such that

$$\|\mathbf{u}^n\|_{\Delta} \leq C_T \|\mathbf{u}^0\|_{\Delta},$$

for all  $n$  such that  $n\Delta t \leq T$ , for all  $\Delta t, h$  such that  $0 < \Delta t, h \leq \delta_0$ , and for all initial data  $\mathbf{u}^0$ .

Note that in the above definition,  $\|\cdot\|_{\Delta}$  is a suitable *discrete norm*. For example,

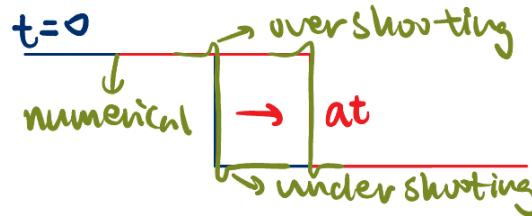
$$\|\mathbf{v}\|_{\Delta,p} = \left( h \sum_{j=-\infty}^{\infty} |v_j|^p \right)^{1/p} \quad \text{for } p = 1, 2 \quad \text{and} \quad \|\mathbf{v}\|_{\Delta,\infty} = \sup_j |v_j|.$$

**Definition 6.2.8 (Strongly Stable).** If the scheme is stable and

$$\|\mathbf{u}^n\|_{\Delta} \leq \|\mathbf{u}^{n-1}\|_{\Delta} \quad \forall n \geq 1,$$

(i.e.,  $C_T = 1$ ), then the scheme is *strongly stable* with respect to  $\|\cdot\|_{\Delta}$ .

Strong stability is useful for asymptotic in time computations ( $T \gg 1$ ). Actually, it ensures that the numerical solution is bounded  $\forall T$ . If we don't have strong stability, we will suffer from overshooting and undershooting problems:



**Definition 6.2.9 ((Un)Conditional Stability).**

- The numerical method is *conditionally stable* if its stability result can hold under restrictions on  $h$  and  $\Delta t$ .
- It is *unconditionally stable* otherwise.

**BE/C** The method BE/C is *unconditionally* strongly stable in norm  $\|\cdot\|_{\Delta,2}$ .

**CFL Condition** An explicit scheme is never unconditionally stable.

A necessary condition is that  $\Delta t$  and  $h$  fulfill the *CFL condition*.

**Definition 6.2.10 (CFL (Courant-Friedrichs-Levy) Condition).** For the problem

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0,$$

the CFL condition is given by

$$|a\lambda| \leq 1 \implies \Delta t \leq \frac{h}{|a|},$$

where  $a$  is the velocity, and  $\lambda = \frac{\Delta t}{h}$ . The adimensional number  $a\lambda$  is the *CFL number*.

### Remark.

- Recall that for parabolic problems, we require  $\Delta t \leq ch^2$ , so the condition is harder to fulfill. But here, we only require  $\Delta t \leq ch$ , as it is easier to fulfill, and the results of explicit methods will not be very bad.
- If  $a$  is not constant, CFL condition reads

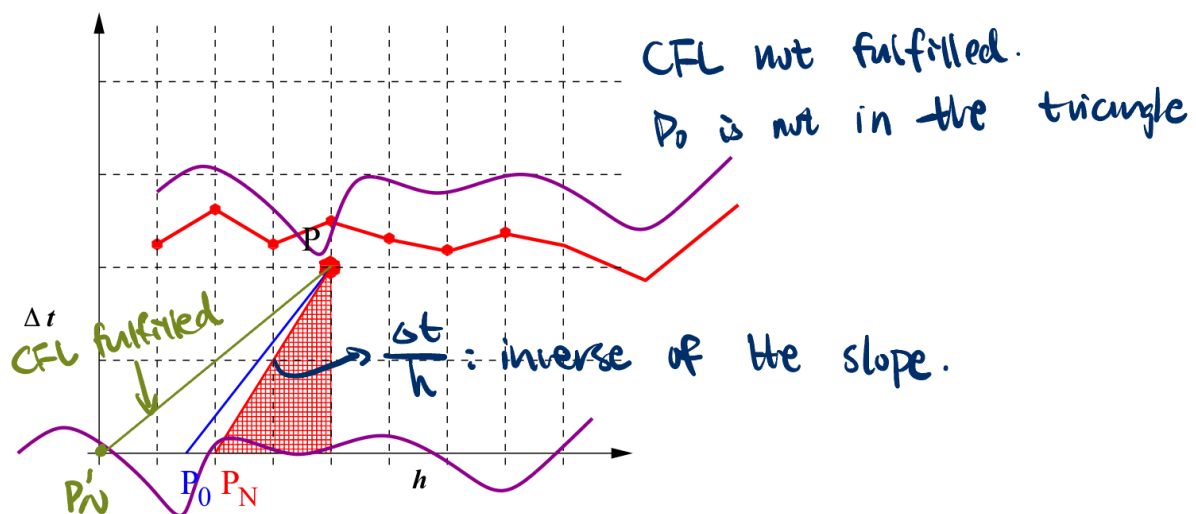
$$\Delta t \leq \frac{\Delta x}{\sup_{x \in \mathbb{R}, t > 0} |a(x, t)|}.$$

- For a system, CFL is given by

$$\left| \lambda_k \frac{\Delta t}{\Delta x} \right| \leq 1, \quad k = 1, \dots, p,$$

where  $\{\lambda_k \mid k = 1, \dots, p\}$  are the eigenvalues of  $A$ .

Graphically, CFL requires that propagation velocity of the numerical solution ( $h/\Delta t$ ) is greater than the velocity of the analytical solution ( $|a|$ ). If CFL condition is satisfied, all the data affecting the exact analytical solution in a point also affect the numerical one.



We prevent initial data to be perturbed so to affect the exact solution and not the numerical one, mandatory for the convergence.

[If CFL is satisfied: numerical domain of convergence contains the exact domain of convergence].

## UPW and LW

**Theorem 6.2.11**

If CFL is satisfied, UPW and LW are *strongly stable* in the norm  $\|\cdot\|_{\Delta,1}$ .

**Theorem 6.2.12 Discrete Maximum Principle**

Under CFL condition, UPW satisfies also the inequality

$$\|\mathbf{u}^n\|_{\Delta,\infty} \leq \|\mathbf{u}^0\|_{\Delta,\infty} \quad \forall n \geq 0.$$

So, UPW is strongly stable in the norm  $\|\cdot\|_{\Delta,\infty}$ .

**Remark.**

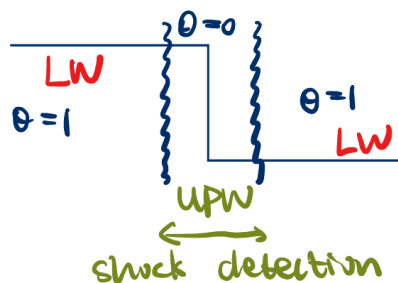
- So, we have no undershooting or overshooting for UPW with CFL.
- For UPW, we trade accuracy with stability: it is less accurate but more stable.
- Improvements: shock capturing methods: combine LW and UPW

1. Numerical flux will be

$$\theta H_{LW} + (1 - \theta) H_{UPW},$$

where  $\theta = 1$  for regular solution, and  $\theta = 0$  for poor regularity.

2. Depending on how we detect shock/irregularity, we form different methods.

**For Wave Equation**

- Leap-Frog is stable under CFL condition

$$\Delta t \leq \frac{\Delta x}{|\gamma|}.$$

- Optimal Newmark method is unconditionally stable.

### 6.3 Von Neumann Analysis of FD Methods

Analyze the stability in  $\|\cdot\|_{\Delta,2}$

$$u_0(x) = \sum_{k=-\infty}^{\infty} a_k e^{ikx}$$

$$a_k = \frac{1}{2\pi} \int_0^{2\pi} u_0(x) e^{-ikx} dx$$

$$u_j^n = \sum_{k=-\infty}^{\infty} a_k e^{ikjh} \gamma_k^n, \quad j = 0, \pm 1, \pm 2, \dots, n \geq 2.$$

Exact solution to

$$\frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0$$

is given by

$$u(x, t^n) = u_0(x - an\Delta t), \quad \forall n \geq 0, \forall x \in \mathbb{R}.$$

With Taylor expansion, we have

$$u(x_j, t^n) = \sum a_k e^{ikjh} g_k^n, \quad g_k = e^{iak\Delta t},$$

where  $g_k^n$  is called the *time travel coefficient*.

$\gamma_j$  corresponds to  $g_k$ , depends on the numerical schemes.

- $g_k$  and  $\gamma_k$  are complex numbers, so we can analyze its magnitude and phase angle.
- Analyzing its magnitude, phase angle, and frequency tells us how solution evolve over time.
- Under CFL,  $|\gamma_k| < 1$  for all methods.
  1. **Dissipation:** how many damping/amplitude.
  2. **Dispersion:** how much delay/velocity.

### 6.4 Finite Elements

FEMs for hyperbolic problems are less stable. We also need to do “semi-discretization.” Depending on how we do time discretization, we form different methods.

- UPW-FEM:

$$\text{CFL} < \frac{1}{3}$$

- LW-FEM:

$$\text{CFL} < \frac{1}{\sqrt{3}}.$$

Under CFL, the two methods are conditionally strongly stable.