

# Numerical Analysis for Partial Differential Equations

Andrea Bonifacio

June 9, 2023

# 1 Boundary Value Problems

## 1.1 Weak Formulation

Let's consider a problem

$$\begin{cases} \mathcal{L}u = f & \text{in } \Omega \\ +\text{B.C.} & \text{on } \partial\Omega \end{cases} \quad (1.1)$$

- $\Omega$ : open bounded domain in  $\mathbb{R}^d$ , with  $d = 2, 3$
- $\partial\Omega$ : boundary of  $\Omega$
- $f$ : given
- B.C. accordingly to  $\mathcal{L}$
- $\mathcal{L}$ : 2<sup>nd</sup> order operator, like:

$$(1) \quad \mathcal{L}u = -\text{div}(\mu \nabla u) + \mathbf{b} \cdot \nabla u + \sigma u \quad (\text{non-conservative form})$$

$$(2) \quad \mathcal{L}u = -\text{div}(\mu \nabla u) + \text{div}(\mathbf{b}u) + \sigma u \quad (\text{conservative form})$$

- $\mu \in L^\infty(\Omega)$ ,  $\mu(\mathbf{x}) \geq \mu_0 > 0$  uniformly bounded from below
- $\mathbf{b} \in (L^\infty(\Omega))^d$  transport term
- $\sigma \in L^2(\Omega)$  reaction term
- $f \in L^2(\Omega)$  can be less regular

## General elliptic problems

Consider

$$\begin{cases} -\text{div}(\mu \nabla u) + \mathbf{b} \cdot \nabla u + \sigma u = f & \text{in } \Omega & g \in L^2(\Gamma_N) \\ u = 0 & \text{on } \Gamma_D & \partial\Omega = \Gamma_D \cup \Gamma_N \\ \mu \nabla u \cdot \mathbf{n} = g & \text{on } \Gamma_N & \Gamma_D^\circ \cap \Gamma_N^\circ = \emptyset \end{cases} \quad (1.2)$$

Suppose that  $f \in L^2(\Omega)$  and  $\mu, \sigma \in L^\infty(\Omega)$ . Also suppose that  $\exists \mu_0 > 0$  s.t.  $\mu(\mathbf{x}) \geq \mu_0$ , and  $\sigma(\mathbf{x}) \geq 0$  a.e. on  $\Omega$ . Then, given a test function  $v$ , we multiply the equation by  $v$ , and integrate on the domain  $\Omega$

$$\int_{\Omega} [-\text{div}(\mu \nabla u) + \mathbf{b} \cdot \nabla u + \sigma u] v = \int_{\Omega} f v$$

By applying Green's formula

$$\underbrace{\int_{\Omega} \mu \nabla u \cdot \nabla v + \int_{\Omega} \mathbf{b} \cdot \nabla u v + \int_{\Omega} \sigma u v}_{=: a(u, v)} = \int_{\Omega} f v + \underbrace{\int_{\Gamma_D} \mu \nabla u \cdot \mathbf{n} v}_{=0 \text{ if } v|_{\Gamma_D}=0} + \int_{\Gamma_N} \underbrace{\mu \nabla u \cdot \mathbf{n} v}_{=g}$$

So the weak formulation of the problem is

$$\begin{cases} \text{find } u \in V & V = \{v \in H^1(\Omega), v|_{\Gamma_D} = 0\} =: H_{\Gamma_D}^1(\Omega) \\ a(u, v) = \langle F, v \rangle & \forall v \in V \end{cases} \quad (1.3)$$

where  $a : V \times V \rightarrow \mathbb{R}$  is a bilinear form and  $F : V \rightarrow \mathbb{R}$  is a linear form s.t.  $\langle F, v \rangle \equiv F(v) = \int_{\Omega} f v + \int_{\Gamma_N} g v$ .

### Theorem 1.1 (Lax-Milgram)

Assume that

- $V$  Hilbert space with  $\|\cdot\|$  and inner product  $(\cdot, \cdot)$

- $F \in V^* : |F(v)| \leq \|F\|_{V^*} \|v\| \quad \forall v \in V$
- $a$  continuous:  $\exists M > 0 : |a(u, v)| \leq M \|u\| \|v\| \quad \forall u, v \in V$
- $a$  coercive:  $\exists \alpha > 0 : a(v, v) \geq \alpha \|v\|^2 \quad \forall v \in V$

Then, there exists a unique solution  $u$  of 1.3

Moreover

$$\alpha \|u\|^2 \leq a(u, u) = F(u) \leq \|F\|_{V^*} \|u\|$$

where  $\alpha$  is the coercivity constant. Hence

$$\|u\| \leq \frac{\|F\|_{V^*}}{\alpha} \rightarrow \text{stability/continuous dependence on data}$$

But what if some of the assumptions of Lax-Milgram (in particular coercivity) are not satisfied? We need a slightly more general problem to formulate Nečas theorem:

$$\begin{cases} \text{find } u \in V \\ a(u, w) = \langle F, w \rangle \quad \forall w \in W \end{cases} \quad (1.4)$$

They belong to different spaces:  $W$  for the test function,  $V$  the solutions

**Theorem 1.2** (Nečas)

Assume that  $F \in W^*$ . Consider the following conditions:

- $a$  continuous:  $\exists M > 0 : |a(u, w)| \leq M \|u\|_V \|w\|_W \quad \forall u \in V, w \in W$
- inf – sup condition:  $\exists \alpha > 0 : \forall v \in V \quad \sup_{w \in W \setminus \{0\}} \frac{a(v, w)}{\|w\|_W} \geq \alpha \|v\|_V$
- $\forall w \in W, w \neq 0, \exists v \in V : a(v, w) \neq 0$

These conditions are necessary and sufficient for the existence and uniqueness of a solution of 1.4, for any  $F \in W^*$ . Moreover

$$\|u\|_V \leq \frac{1}{\alpha} \|F\|_{W^*}$$

When  $W = V$  Lax-Milgram provides necessary and sufficient conditions for existence and uniqueness of solutions.

Going back to

$$\begin{cases} \mathcal{L}u = f & \text{in } \Omega \\ +\text{B.C.} & \text{on } \partial\Omega \end{cases}$$

What could be our choice of  $V$ ? Given that

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

and

$$a(u, v) = \int_{\Omega} \mu \underbrace{\nabla u \nabla v}_{\nabla u, \nabla v \in L^2} + \int_{\Omega} b \underbrace{\nabla uv}_{\in L^1} + \int_{\Omega} \sigma \underbrace{uv}_{\in L^1}$$

We want to choose  $v$  in order to have all of these integrable

$$\Rightarrow V = \left\{ v \in L^2(\Omega), \nabla u \in [L^2(\Omega)]^d, v|_{\Gamma_D} = 0 \right\} = V_{\Gamma_D}$$

.

Knowing that a Sobolev space

$$H^1 = \left\{ v \in L^2(\Omega), \nabla u \in [L^2(\Omega)]^d \right\}$$

we can say  $V_{\Gamma_D} = \{v \in H^1(\Omega) : v|_{\Gamma_D} = 0\}$ , and if  $\Gamma_D = \partial\Omega$ , then  $V_{\Gamma_D} = H_0^1$

## 1.2 Approximation

Recall for a moment the weak formulation of a generic elliptic problem

$$\begin{cases} \text{find } u \in V \\ a(u, v) = \langle F, v \rangle \quad \forall v \in V \end{cases} \quad (1.5)$$

with  $V$  being an appropriate Hilbert space, subset of  $H^1()$ ,  $a(\cdot, \cdot)$  being a continuous and coercive bilinear form from  $V \times V \rightarrow \mathbb{R}$ ,  $F(\cdot)$  being a continuous linear functional from  $V \rightarrow \mathbb{R}$ .

Let  $V_h \subset V$  be a family of spaces that depends on a parameter  $h > 0$ , such that  $\dim V_h = N_h < \infty$ . We can rewrite the weak formulation

$$\begin{cases} \text{find } u_h \in V_h \\ a(u_h, v_h) = \langle F, v_h \rangle \quad \forall v_h \in V_h \end{cases} \quad (1.6)$$

and is called a **Galerkin problem**. Denoting with  $\{\varphi_j, j = 1, 2, \dots, N_h\}$  a basis of  $V_h$ , it is sufficient that the (1.6) is verified for each function of the basis. Also we need that

$$a(u_h, \varphi_i) = F(\varphi_i) \quad i = 1, 2, \dots, N_h$$

Since  $u_h \in V_h$

$$u_h(\mathbf{x}) = \sum_{j=1}^{N_h} u_j \varphi_j(\mathbf{x})$$

where  $u_j$  are unknown coefficients. Then

$$\sum_{j=1}^{N_h} u_j a(\varphi_j, \varphi_i) = F(\varphi_i)$$

We denote by  $A$  the matrix made by  $a_{ij} = a(\varphi_j, \varphi_i)$  and  $\mathbf{f}$  the vector of  $F(\varphi_i) = f_i$  components. If we denote the vector  $\mathbf{u}$  made by the unknown coefficients  $u_h$ .

$$A\mathbf{u} = \mathbf{f} \quad (1.7)$$

### Theorem 1.3

The stiffness matrix  $A$  associated to the Galerkin discretization of an elliptic problem, whose bilinear form is coercive is positive definite.

**Proof.** Recall that a matrix  $B \in \mathbb{R}^{n \times n}$  is said to be positive definite if

$$\mathbf{v}^T B \mathbf{v} \geq 0 \quad \forall \mathbf{v} \in \mathbb{R}^n$$

and

$$\mathbf{v}^T B \mathbf{v} = 0 \Leftrightarrow \mathbf{v} = \mathbf{0}$$

The correspondence

$$\mathbf{v} = (v_i) \in \mathbb{R}^{N_h} \longrightarrow v_h(x) = \sum_{j=1}^{N_h} v_j \varphi_j \in V_h$$

defines a bijection between  $V_h$  and  $\mathbb{R}^{N_h}$ . Given a generic vector  $\mathbf{v} = (v_i)$  of  $\mathbb{R}^{N_h}$ , thanks to the bilinearity and coercivity of  $a$  we obtain

$$\begin{aligned}
\mathbf{v}^T A \mathbf{v} &= \sum_{j=1}^{N_h} \sum_{i=1}^{N_h} v_i a_{ij} v_j \\
&= \sum_{j=1}^{N_h} \sum_{i=1}^{N_h} v_i a(\varphi_j, \varphi_i) v_j \\
&= \sum_{j=1}^{N_h} \sum_{i=1}^{N_h} a(v_j \varphi_j, v_i \varphi_i) \\
&= a \left( \sum_{j=1}^{N_h} v_j \varphi_j, \sum_{i=1}^{N_h} v_i \varphi_i \right) \\
&= a(v_h, v_h) \geq \alpha \|v_h\|_V^2 \geq 0
\end{aligned}$$

Moreover, if  $\mathbf{v}^T A \mathbf{v} = 0$ , then  $\|v_h\|_V^2 = 0$ . ★

### Existence and uniqueness

#### Corollary 1.1

The solution of the Galerkin problem (1.6) exists and is unique.

To prove this we can prove that the solution to (1.7) exists and is unique. The matrix  $A$  is invertible as the unique solution of  $A\mathbf{u} = \mathbf{0}$  is the null solution, meaning that  $A$  is definite positive.

### Stability

#### Corollary 1.2

The Galerkin method is stable, uniformly with respect to  $h$ , by virtue of the following upper bound for the solution

$$\|u_h\|_V \leq \frac{1}{\alpha} \|F\|_{V^*}$$

The stability of the method guarantees that the norm  $\|u_h\|_V$  of the discrete solution remains bounded for  $h \rightarrow 0$ . Equivalently it guarantees that  $\|u_h - w_h\|_V \leq \frac{1}{\alpha} \|F - G\|_{V^*}$  with  $u_h$  and  $w_h$  being numerical solution corresponding to different data  $F$  and  $G$ .

### Convergence

#### Lemma 1.1 (Galerkin orthogonality)

The solution  $u_h$  of the Galerkin method satisfies

$$a(u - u_h, v_h) = 0 \quad \forall v_h \in V_h \tag{1.8}$$

**Proof.** Since  $V_h \subset V$ , the exact solution  $u$  satisfies the weak problem (1.5) for each element  $v = v_h \in V_h$ , hence we have

$$a(u, v_h) = F(v_h) \quad \forall v_h \in V_h \tag{1.9}$$

By subtracting side by side (1.6) from (1.9), we obtain

$$a(u, v_h) - a(u_h, v_h) = 0 \quad \forall v_h \in V_h$$

from which the claim follows. ★

Also this can be generalized in the cases in which  $a(\cdot, \cdot)$  is not symmetric. Consider the value taken by the bilinear form when both its arguments are  $u - u_h$ . If  $v_h$  is an arbitrary element of  $V_h$  we obtain

$$a(u - u_h, u - u_h) = a(u - u_h, u - v_h) + a(u - u_h, v_h - u_h)$$

The last term is null by (1.8). Moreover

$$|a(u - u_h, u - v_h)| \leq M \|u - u_h\|_V \|u - v_h\|_V$$

having exploited the continuity of the bilinear form. Also by the coercivity

$$a(u - u_h, u - u_h) \geq \alpha \|u - u_h\|_V^2$$

hence

$$\|u - u_h\|_V \leq \frac{M}{\alpha} \|u - v_h\|_V \quad \forall v_h \in V_h$$

Such inequality holds for all functions  $v_h \in V_h$  and therefore we find

$$\underbrace{\|u - u_h\|_V}_{\text{Galerkin error}} \leq \frac{M}{\alpha} \underbrace{\inf_{w_h \in V_h} \|u - w_h\|_V}_{\text{Best Approximation Error}} \quad (1.10)$$

In order for the method to converge, it is sufficient that, for  $h \rightarrow 0$  the space  $V_h$  tends to saturate the entire space  $V$ .

$$\lim_{h \rightarrow 0} \inf_{v_h \in V_h} \|v - v_h\|_V = 0 \quad \forall v \in V \quad (1.11)$$

In that case the Galerkin method is convergent and it can be written that

$$\lim_{h \rightarrow 0} \|u - u_h\|_V = 0 \Leftrightarrow \text{convergence}$$

This space  $V_h$  must be chosen carefully to satisfy the saturation property (1.11).

### 1.3 Finite Element Method

#### Partitions

**1D** Let us suppose that  $\Omega$  is an interval  $(a, b)$ . How to create an approximation of the space  $H^1(a, b)$  that depend on a parameter  $h$ . Consider a partition  $\mathcal{T}_h$  in  $N + 1$  subintervals  $K_j = [x_{j-1}, x_j]$ , having width  $h_j = x_j - x_{j-1}$  with

$$a = x_0 < x_1 < \dots < x_{N-1} < x_N = b \quad (1.12)$$

and set  $h = \max_j h_j$ .

**2D** Now we can extend the FEM for multi-dimensional problems. For simplicity we will consider  $\Omega \subset \mathbb{R}^2$  with polygonal shapes  $\mathcal{T}_h$ . In this case the partition is called a triangulation. We can define the discretized domain

$$\Omega_h = \text{int} \left( \bigcup_{K \in \mathcal{T}_h} K \right)$$

in a way that the internal part of the union of the triangles  $\mathcal{T}_h$ . Having set  $\text{diam}(K) = \max_{x, y \in K} |x - y| = h_k$ . Also, given  $\rho_K$  the measure of the diameter of the circle inscribed in the triangle  $K$ , must be satisfied the condition that, for a suitable  $\delta > 0$

$$\frac{h_k}{\rho_k} \leq \delta \quad \forall K \in \mathcal{T}_h \quad (1.13)$$

The condition (1.13) excludes very deformed triangles.

**Definition 1.1** (Seminorms)

A seminorm is defined as

$$|f|_k = |f|_H^k(\Omega) = \sqrt{\sum_{|\alpha|=k} \int_{\Omega} (D^{\alpha} f)^2 d\Omega}$$

In particular

$$\begin{aligned} \text{1D: } |u|_{H^1(a,b)} &= \left( \|u_x\|_{L^2(a,b)}^2 \right)^{\frac{1}{2}} = \|u_x\|_{L^2(a,b)} \\ |u|_{H^2(a,b)} &= \|u_{xx}\|_{L^2(a,b)} \\ \text{2D: } |u|_{H^1(a,b)} &= \left( \|u_x\|_{L^2(a,b)}^2 + \|u_y\|_{L^2(a,b)}^2 \right)^{\frac{1}{2}} \\ |u|_{H^1(a,b)} &= \left( \|u_{xx}\|_{L^2(a,b)}^2 + \|u_{xy}\|_{L^2(a,b)}^2 + \|u_{yx}\|_{L^2(a,b)}^2 + \|u_{yy}\|_{L^2(a,b)}^2 \right)^{\frac{1}{2}} \end{aligned}$$

Always true that  $|u|_{H^q} \leq \|u\|_{H^q}$

The problem is always:

$$\begin{aligned} \text{find } u_h \in V_h : a(u_h, v_h) &= F(v_h) \quad \forall v_h \in V_h \\ \downarrow \\ V_h &= \{v_h \in X_h^r : v_h|_{\Gamma_D} = 0\} \quad r \geq 1 \end{aligned} \tag{1.14}$$

Since the functions of  $H^1(a, b)$  are continuous on  $[a, b]$ , it is possible to create the family of spaces

$$X_h^r = \{v_h \in \mathcal{C}^0(\bar{\Omega}) : v_h|_{K_j} \in \mathbb{P}_r \quad \forall K_j \in \mathcal{T}_h\}, \quad r = 1, 2, \dots \tag{1.15}$$

having denoted by  $\mathbb{P}_r$  the space of polynomials with degree lower or equal to  $r$  in the variable  $x$ . All these spaces are subspaces of  $H^1(a, b)$  as they are constituted by differentiable functions except for at most a finite number of points (the vertices of the partition). It is convenient to select a basis for the  $X_h^r$  space that is *Lagrangian*.

$\mathbb{P}^r :$	1D	$p(x) = \sum_{k=0}^r a_k x^k$	intervals
	2D	$p(x_1, x_2) = \sum_{\substack{k,m=0 \\ k+m \leq r}}^r a_{km} x_1^k x_2^m$	triangles
	3D	$p(x_1, x_2, x_3) = \sum_{\substack{k,m,n=0 \\ k+m+n \leq r}}^r a_{kmn} x_1^k x_2^m x_3^n$	tetrahedra

**The space  $X_h^1$** 

The space is constituted by the functions of the partition (1.12). Since only a straight line can pass through different points, the degrees of freedom (DOF, the number of values we need to assign to the basis to define the functions) of the functions will be equal to the number  $N + 2$  of vertices of the partition. It follows naturally that  $\{\varphi_i\}, i = 0, 1, \dots, N, N + 1$ . In this case the basis functions are characterized by the following properties

$$\varphi_i \in X_h^1 \text{ s.t. } \varphi_i(x_j) = \delta_{ij}, \quad i, j = 0, 1, \dots, N, N + 1$$

where  $\delta_{ij}$  is the Kronecker delta. So we have our basis function that have value 1 in the node  $x_j$  and 0 elsewhere.

The formula for the basis function is then given by

$$\varphi_i(x) = \begin{cases} \frac{x-x_{i-1}}{x_i-x_{i-1}} & \text{for } x_{i-1} \leq x \leq x_i \\ \frac{x-x_{i+1}}{x_{i+1}-x_i} & \text{for } x_i \leq x \leq x_{i+1} \\ 0 & \text{otherwise} \end{cases} \quad (1.16)$$

### The space $X_h^2$

In this case polynomials are of degree 2, so the points necessary to evaluate them are 3. The chosen points for every element of the partition  $\mathcal{T}_h$ . The nodes from the interval goes from  $a = x_0$  to  $b = x_{2N+2}$ , so that midpoints are the nodes with odd indices. As the previous case the basis is Lagrangian

$$\varphi_i \in X_h^2 \text{ s.t } \varphi_i(x_j) = \delta_{ij}, \quad i, j = 0, 1, \dots, 2N+2$$

### The space $V_h$

This space is generated by

$$V_h = \{v_h \in X_h^r : v_h(a) = v_h(b) = 0\}$$

Having defined a basis  $\{\varphi_j(\mathbf{x})\}_{j=1}^{N_h}$  for the space  $V_h$ , each  $v_h$  can be expanded as a linear combination of elements of the basis, suitably weighted by coefficients  $\{v_j\}_{j=1}^{N_h}$

$$v_h(\mathbf{x}) = \sum_{j=1}^{N_h} v_j \varphi_j(\mathbf{x})$$

A basis is called Lagrangian if it satisfies the following properties

$$\varphi_i(\mathbf{x}_j) = \delta_{ij} \quad \forall 1 \leq i, j \leq N_h$$

and then the following property holds:

$$v_h(\mathbf{x}_j) = v_j \quad \forall 1 \leq i, j \leq N_h$$

The solution of the Finite Element Method,  $u_h$  can be written as

$$u_h(\mathbf{x}) = \sum_{j=1}^{N_h} u_j \varphi_j(\mathbf{x}) \quad (1.17)$$

In (1.14) take  $v_h = \varphi_j \quad \forall j = 1, \dots, N_h$  such that  $a(u_h, \varphi_i) = F(\varphi_i) \quad \forall i = 1, \dots, N_h$ . Then use (1.17) to obtain

$$\begin{aligned} a \left( \sum_{j=1}^{N_h} u_j \varphi_j(\mathbf{x}), \varphi_i \right) &= \underbrace{F(\varphi_i)}_{F_i} \\ \Rightarrow \sum_{j=1}^{N_h} \underbrace{a(\varphi_j, \varphi_i)}_{\substack{\text{elements} \\ \text{of } A}} u_j(\mathbf{x}) &= F_i \quad i = 1, \dots, N_h \\ \Rightarrow A \mathbf{u} &= \mathbf{F} \end{aligned}$$

Which is a linear system of dimension  $N_h \times N_h$  with  $\mathbf{F}$  the right hand side (RHS),  $A$  the stiffness matrix and  $\mathbf{u}$  a vector of unknown nodal values of the solution  $u_h$ .



## 1.4 Advection Diffusion Reaction Problem

$$\begin{cases} Lu = \underbrace{-\operatorname{div}(\mu \nabla u)}_{\text{diffusion}} + \underbrace{\mathbf{b} \cdot \nabla u}_{\text{advection}} + \underbrace{\sigma u}_{\text{reaction}} = f & \text{in } \Omega \\ \text{BC} & \text{on } \partial\Omega \end{cases}$$

Lax-Milgram tells us that if  $\sigma - \frac{1}{2}\operatorname{div}\mathbf{b} \geq \gamma > 0$  then  $\exists!$  a solution to the problem. But what if these conditions are not satisfied? We can use Nečas theorem ((1.2)) with equivalent assumptions:

- Weak coercivity (Gårding inequality):

$$\exists \alpha, \lambda : a(v, v) \geq \alpha \|v\|^2 - \|v\|_{L^2(\Omega)}^2 \quad \forall v \in V$$

- Uniqueness condition (typically proven by maximum principle):

$$(a(u, v) = 0 \quad \forall v \in V) \Rightarrow u = 0$$

If  $A$  is spd (symmetric positive defined) then  $K_2(A) = \frac{\lambda_{\max}(A)}{\lambda_{\min}(A)}$

### Proposition 1.1

If  $a(\cdot, \cdot)$  is symmetric and coercive, then  $A$  is spd.

**Proof.** Symmetry:  $A_{ij} = a(\varphi_j, \varphi_i) = a(\varphi_i, \varphi_j) = A_{ji}$   
 $\forall \mathbf{v} \in \mathbb{R}^{N_h}$ :

$$\begin{aligned} \mathbf{v}^T A \mathbf{v} &= \sum_{i,j} A_{ij} v_i v_j = \sum_{i,j} a(\varphi_j, \varphi_i) v_i v_j \\ &= a\left(\sum_j v_j \varphi_j, \sum_i v_i \varphi_i\right) = a(v_h, v_h) \geq \alpha \|v_h\|^2 > 0 \end{aligned}$$

if  $(v_h \neq 0 \Leftrightarrow \mathbf{v} \neq \mathbf{0})$ . Hence  $A$  is positive defined. ★

### Definition 1.2

If  $A$  is spd, we define the  $A$ -norm of  $\mathbf{v}$  as

$$\begin{aligned} \|v\|_A &:= (A\mathbf{v}, \mathbf{v})^{\frac{1}{2}} \\ &= \left( \sum_{i,j} a_{ij} v_i v_j \right)^{\frac{1}{2}} \end{aligned}$$

Since  $A$  is positive defined  $\Rightarrow \operatorname{Re}(\lambda_k(A)) \Rightarrow \lambda_k(A) \neq 0$ . Then, by symmetry of  $A \Rightarrow \lambda_k(A) \in \mathbb{R}$ . Combining the two we have that  $A$  sdp  $\Rightarrow \lambda_k(A) > 0 \Rightarrow \exists!$  solution of  $A\mathbf{u} = \mathbf{f}$

### Definition 1.3

If  $A$  is sdp, then  $K_2(A) = \frac{\lambda_{\max}}{\lambda_{\min}}$  is called **spectral condition number**

If  $K_2(A) \gg 1 \Rightarrow A$  is ill-conditioned  $\Rightarrow$  solving  $A\mathbf{u} = \mathbf{f}$  is hard.

We can also prove that  $\exists C_1, C_2 > 0 : \forall \lambda_h$  eigenvalue of  $A$ :

$$\alpha C_1 h^d \leq \lambda_h \leq M C_2 h^{d-2} \quad d = 1, 2, 3$$

whence

$$\frac{\lambda_{\max}(A)}{\lambda_{\min}(A)} \leq \frac{M C_2}{\alpha C_1} h^{-2}$$

Then

$$K_2(A) = \mathcal{O}(h^{-2})$$

If we use the conjugate gradient method to solve  $A\mathbf{u} = \mathbf{f}$ , then:

$$\|\mathbf{u}^{(k)} - \mathbf{u}\|_A \leq 2 \left( \frac{\sqrt{K_2(A)} + 1}{\sqrt{K_2(A)} - 1} \right)^k \|\mathbf{u}^{(k)} - \mathbf{u}\|_A$$

Same with gradient method, with  $K_2(A)$  instead of  $\sqrt{K_2(A)} \Rightarrow$  need for preconditioners.

## 1.5 Interpolant estimates

$$\|u - u_h\|_V \leq \frac{M}{\alpha} \inf_{v_h \in V_h} \|u - v_h\| \xrightarrow{\text{saturation}} 0 \Leftrightarrow \text{convergence} \quad (1.18)$$

But how fast it saturates?

Note:  $\inf_{v_h \in V_h} \|u - v_h\|_V \leq \|u - \bar{u}_h\|_V \quad \forall \bar{u}_h$  suitable chosen in  $V_h$  and  $\bar{u}_h$  is a smart guy chosen in a smart way (close enough to  $u$ ).

In 1D the finite element interpolant can be defined as  $\prod_h^r u(x_k) = u(x_k) \quad \forall x_k$  node. Then  $\bar{u}_h = \prod_h^r u \in V_h$ .

How good is  $\bar{u}_h$ ?

$$\prod_h^r u(x) = \sum_{j=1}^{N_h} u(x_j) \varphi_j(x)$$

which is a good approximation.

### Interpolant error estimates

Then, for  $m = 0, 1 \exists C = C(r, m, \hat{k})$  s.t.

$$\left| v - \prod_h^r v \right|_{H^m(\Omega)} \leq C \left( \sum_{K \in \mathcal{T}_h} h_K^{2(r+1-m)} |v|_{H^{r+1}(K)}^2 \right)^{\frac{1}{2}} \quad (1.19)$$

where  $h_K = \text{diam}(K)$  and  $h_K \leq h \quad \forall K$  this yields:

$$\left| v - \prod_h^r v \right|_{H^m(\Omega)} \leq C h^{r+1-m} |v|_{H^{r+1}(K)} \quad \forall v \in H^{r+1}(\Omega), m = 0, 1 \quad (1.20)$$

Recall also that

$$\begin{aligned} \|u - u_h\| &= \|u - u_h\|_{H^1(\Omega)} \\ &\leq \frac{M}{\alpha} \inf_{v_h \in V_h} \|u - v_h\| \\ &\leq \frac{M}{\alpha} \left\| u - \prod_h^r u \right\|_{H^1(\Omega)} \end{aligned}$$

Using (1.19) we obtain

$$\|u - u_h\| \leq C \frac{M}{\alpha} \left( \sum_{K \in \mathcal{T}_h} h_K^{2r} |v|_{H^{r+1}(K)}^2 \right)^{\frac{1}{2}} \quad (1.21)$$

Then, by using (1.20):

$$\|u - u_h\| \leq C \frac{M}{\alpha} h^r |u|_{H^{r+1}(\Omega)} \quad (1.22)$$

#### Definition 1.4

Consider a bilinear form  $a : V \times V \rightarrow \mathbb{R}$ . The *adjoint* form  $a^*$  is defined as  $a^* : V \times V \rightarrow \mathbb{R}$

$$a^*(v, w) = a(w, v) \quad \forall v, w \in V$$

Now let's consider the adjoint problem

$$\begin{cases} \text{find } \varphi = \varphi(g) \in V & \forall g \in L^2(\Omega) \\ a^*(\varphi, v) = (g, v) = \int_{\Omega} g v & \forall v \in V \end{cases} \quad (1.23)$$

Assuming that  $\varphi \in H^2(\Omega) \cap V$  (elliptic regularity). Consider now, for example,  $\mathcal{L} = -\Delta$ . Then the solution of

$$\begin{cases} -\Delta\varphi = g & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases}$$

satisfies  $\varphi \in H^2(\Omega)$ . Moreover

$$\exists C_1 > 0 : \|\varphi(g)\|_{H^2(\Omega)} \leq C_1 \|g\|_{L^2(\Omega)} \quad (1.24)$$

Take now  $g = e_h = u - u_h$  in (1.23). Then

$$\begin{aligned} \|e_h\|_{L^2(\Omega)}^2 &= a^*(\varphi, e_h) = a(e_h, \varphi) \\ &= a(e_h, \varphi - \varphi_h) && \text{(Galerkin orthogonality)} \\ &\leq M \|e_h\|_{H^1(\Omega)} \|\varphi - \varphi_h\|_{H^1(\Omega)} \end{aligned}$$

Take then  $\varphi_h = \prod_h^1 \varphi$ :

$$\begin{aligned} \|e_h\|_{L^2(\Omega)}^2 &\leq M \|e_h\|_{H^1(\Omega)} \left\| \varphi - \prod_h^1 \varphi \right\|_{H^1(\Omega)} \\ &\leq M \|e_h\|_{H^1(\Omega)} C_2 h \|\varphi\|_{H^2(\Omega)} && \text{(for (1.20) with m=r=1)} \\ &\leq M \|e_h\|_{H^1(\Omega)} C_2 h C_1 \|e_h\|_{L^2(\Omega)} && \text{(for (1.24))} \end{aligned}$$

Whence:

$$\begin{aligned} \|e_h\|_{L^2(\Omega)} &\leq C_1 C_2 h \|e_h\|_{H^1(\Omega)} \\ &\leq M C_1 C_2 h C_3 h^r |u|_{H^{r+1}(\Omega)} && \text{(for (1.22))} \end{aligned}$$

So

$$\|e_h\|_{L^2(\Omega)} \leq \bar{C} h^{r+1} |u|_{H^{r+1}(\Omega)} \quad (1.25)$$

## 2 Spectral Element Method

### 2.1 Introduction

The problem with the Finite Element Method is that the rate of convergence is limited by the degree of the polynomials used. An alternative can be the Spectral Element Method, for which the convergence rate is limited by the regularity of the solution.

### 2.2 Legendre polynomials

The Legendre polynomials  $\{L_k(x) \in \mathbb{P}_k, k = 0, 1, \dots\}$  are the eigenfunctions of the singular Sturm-Liouville problem:

$$((1-x^2)L'_k(x))' + k(k+1)L_k(x) = 0 \quad -1 < x < 1$$

So they satisfy the recurrence relation

$$\begin{aligned} L_0(x) &= 1, \quad L_1(x) = x, \quad \text{and for } k \geq 1 \\ L_{k+1}(x) &= \frac{2k+1}{k+1}xL_k(x) - \frac{k}{k+1}L_{k-1}(x) \end{aligned} \quad (2.1)$$

Given a weight function  $w(x) \equiv 1$ , they are mutually orthogonal with respect to it on the interval  $(-1, 1)$

$$\int_{-1}^1 L_k(x)L_m(x) dx = \begin{cases} \frac{2}{2k+1} & \text{if } k = m \\ 0 & \text{if } k \neq m \end{cases}$$

The expansion of  $u \in L^2(-1, 1)$  in terms of  $L_k$  is

$$u(x) = \sum_{k=0}^{\infty} \hat{u}_k L_k(x)$$

Given that  $(f, g) = \int_{-1}^1 fg dx$  we know that:

$$(u, L_m) = \sum_{k=0}^{\infty} \hat{u}_k (L_k, L_m) \underset{\text{orth.}}{=} \hat{u}_m \frac{2}{2m+1} \Rightarrow \hat{u}_k = \frac{2k+1}{2} \int_{-1}^1 u L_k dx$$

The truncated Legendre series of  $u$  is the  $L^2$  - projection of  $u$  over  $\mathbb{P}_N$  is

$$P_N u = \sum_{k=0}^N \hat{u}_k L_k \quad (2.2)$$

Given any  $u \in H^s(-1, 1)$  with  $s \in N$ , the projection error  $(u - P_N u)$  satisfies the estimates

$$\begin{aligned} \|u - P_N u\|_{L^2(-1,1)} &\leq CN^{-s} \|u\|_{H^s(-1,1)} & \forall s \geq 0 \\ \|u - P_N u\|_{L^2(-1,1)} &\leq CN^{-s} |u|_{H^s(-1,1)} & \forall s \leq N+1 \end{aligned}$$

There is also a “modified” Legendre basis for function that vanish at  $\pm 1$ . This is because the Legendre basis is not suited to impose Dirichlet B.C.

$$\begin{aligned} \psi_0(x) &= \frac{1}{2}(L_0(x) - L_1(x)) = \frac{1-x}{2} \\ \psi_N(x) &= \frac{1}{2}(L_0(x) + L_1(x)) = \frac{1+x}{2} \\ \psi_{k-1}(x) &= \frac{1}{\sqrt{2(2k-1)}}(L_{k-2}(x) - L_k(x)) \\ &\quad \text{for } k = 2, \dots, N \quad -1 < x < 1 \end{aligned}$$

### 2.3 Spectral Galerkin formulation

Given  $\Omega = (-1, 1)$ ,  $\mu, b, \sigma > 0$  const.,  $f : \Omega \rightarrow \mathbb{R}$ . Look for  $u : \Omega \rightarrow \mathbb{R}$  s.t.

$$\begin{cases} -(\mu u')' + (bu)' + \sigma u = f & \text{in } \Omega \\ u(-1) = 0 \\ u(1) = 0 \end{cases}$$

Set  $V = H_0^1(\Omega)$ , then the weak form of the differential problem reads:

$$\text{find } u \in V \text{ s.t. } a(u, v) = (f, v)_{L^2(\Omega)} \quad \forall v \in V, f \in L^2(\Omega)$$

where

$$\begin{aligned} a(u, v) &= \int_{\Omega} (\mu u' - bu)v' dx + \int_{\Omega} \sigma uv dx \\ (f, v)_{L^2(\Omega)} &= \int_{\Omega} f v dx \end{aligned}$$

Now set  $V_N = \mathbb{P}_N^0$

$$\text{find } u_N \in V_N : a(u_N, v_N) = (f, v_N)_{L^2(\Omega)} \quad (2.3)$$

Now expand  $u_N(x) = \sum_{k=1}^{N-1} \tilde{u}_k \psi_k(x)$  and chose  $v_N = \psi_i(x)$  for any  $i = 1, \dots, N-1$ . The discretization of the problem reads:

$$\text{find } u = [\tilde{u}]_{k=1}^{N-1} : \sum_{k=1}^{N-1} a(\varphi_k, \psi_i) \tilde{u}_k = (f, \psi_i)_{L^2(\Omega)} \quad \text{for any } i = 1, \dots, N-1$$

Given  $u_N \in V_N$  the solution of the problem, then if  $u \in H^{s+1}(\Omega)$  with  $s \geq 0$ , thanks to Ceà Lemma, holds that:

$$\|u - u_N\|_{H^1(\Omega)} \leq C(s) \left( \frac{1}{N} \right)^s \|u\|_{H^{s+1}(\Omega)}$$

So  $u_N$  converges with spectral accuracy with respect to  $N$ . But doing so we would have two full matrices, the stiffness one and the mass one  $M_{ij} = (\psi_j, \psi_i)_{L^2(-1,1)}$  are quite expensive to compute or invert.

To solve this we can use a Lagrange nodal basis instead of a modal one, by using the Legendre-Gauss-Lobatto quadrature formulas. In this case we need a Legendre polynomial  $L_N(x)$ .

Given a  $L_N(x)$  polynomial, we can put one node at each end of the domain, so  $x_0 = -1, x_N = 1$  and  $x_j = \text{zeros of } L'_N$  with  $j = 1, \dots, N-1$ . We also need a set of weights  $w_j = \frac{2}{N(N+1)} \frac{1}{[L_N(x_j)]^2}$  with  $j = 0, \dots, N$ .

With this set of nodes and weights it's possible to obtain the following interpolatory quadrature formula

$$\int_{-1}^1 f(x) dx \approx \sum_{j=0}^N f(x_j) w_j$$

The degree of exactness of this method is  $2N-1$ , meaning that

$$\int_{-1}^1 f(x) dx = \sum_{j=0}^N f(x_j) w_j \quad \forall f \in \mathbb{P}_{2N-1}$$

Some useful operation with LGL nodes

- Discrete inner product in  $L^2(-1, 1)$ :

$$(u, v)_N = \sum_{j=0}^N u(x_j) v(x_j) w_j$$

with degree of exactness  $2N-1$

$$(u, v)_{L^2(\Omega)} = (u, v)_N \quad \text{only if } u, v \in \mathbb{P}_{2N-1}$$

- Discrete norm in  $L^2(-1, 1)$

$$\|u\|_N = (u, u)_N^{\frac{1}{2}}$$

with the following norm equivalence:  $\exists c_1, c_2 > 0$  s.t.

$$c_1 \|v_N\|_{L^2(-1,1)} \leq \|v_N\|_N \leq c_2 \|v_N\|_{L^2(-1,1)} \quad \forall v_N \in \mathbb{P}_N$$

Given  $\{\varphi_0, \dots, \varphi_N\}$  characteristics Lagrange polynomials in  $\mathbb{P}_N$  w.r.t the LGL nodes. then

$$\varphi_j = \frac{1}{n(n+1)} \frac{(1-x^2)}{(x_j-x)} \frac{L'_N(x)}{L_N(x_j)} \quad \text{for } j = 0, \dots, N$$

Also true that  $\varphi_j(x_k) = \delta_{kj}$  and  $\{\varphi_j\}$  are orthogonal w.r.t. the discrete inner product  $(\cdot, \cdot)_N$ , meaning that the mass matrix  $M$  is diagonal. Given  $\{w_i\}$  the set of weights, then

$$M_{ij} = (\varphi_j, \varphi_i)_N = \delta_{ij} w_i \quad i, j = 0, \dots, N$$

## 2.4 Galerkin with Numerical Integration

We can now define the spectral Galerkin method with numerical integration (GNI), by setting our bilinear discrete form as  $a_N(u_N, v_N) = (\mu u'_N - b u_n, v'_N)_N + (\sigma u_n, v_n)_N$ , and the problem as

$$\text{find } u_N^{\text{GNI}} \in V_N : a_N(u_N^{\text{GNI}}, v_N) = (f, v_N)_N \quad \forall v_N \in V_N$$

Then, by the same expansion w.r.t. the Lagrange basis:  $u_N^{\text{GNI}}(x) = \sum_{i=0}^N u_N^{\text{GNI}}(x_i) \varphi_i(x)$  and choose  $v_N(x) = \varphi_i(x)$  for any  $i = 1, \dots, N-1$ .

The GNI discretization of the weak problem reads:

$$\text{look for } u^{\text{GNI}} = [u_N^{\text{GNI}}(x_j)]_{j=0}^N : \begin{cases} u_N^{\text{GNI}}(x_0) = u_N^{\text{GNI}}(x_N) \\ \sum_{j=0}^N a_N(\varphi_j, \varphi_i) u_N^{\text{GNI}}(x_j) = (f, \varphi_i)_N \quad \forall i = 1, \dots, N-1 \end{cases}$$

Now let's have a closer look to the  $\{\varphi_j\}$ :

$$\varphi_j \in \mathbb{P}_N : \varphi_j(x_i) = \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}$$

Given the discrete inner product  $(u, v)_N = \sum_{j=0}^N u(x_j) v(x_j) w_j$  we can write:

$$\begin{aligned} (\varphi_k, \varphi_m)_N &= \sum_{j=0}^N \underbrace{\varphi_k(x_j)}_{\delta_{kj}} \underbrace{\varphi_m(x_j)}_{\delta_{mj}} w_j \quad 0 \leq k, m \leq N \\ &= \sum_{k=0}^N \begin{cases} w_m & \text{if } k = m \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

so  $\{\varphi_k\}$  is orthogonal under the discrete inner product.

The GNI solution is

$$u_N(x) = \sum_{i=0}^N \alpha_i \varphi_i(x) \quad \{\alpha_i\} \text{ unknown coefficients}$$

Set now  $x = x_j$  with LGL nodes:

$$u_N(x_j) = \sum_{i=0}^N \alpha_i \underbrace{\varphi_i(x_j)}_{\delta_{ij}} = \alpha_j$$

So, given  $u_n^{\text{GNI}}(x_j)$  the nodal values, we obtain the nodal expansion:

$$u_N^{\text{GNI}}(x) = \sum_{j=0}^N u_N^{\text{GNI}}(x_j) \varphi_j(x)$$

## Algebraic form of Spectral GNI

Now it's about solving the following linear system

$$A^{\text{GNI}} \mathbf{u}^{\text{GNI}} = \mathbf{f}^{\text{GNI}}$$

with  $A_{ij}^{\text{GNI}} = a_N(\varphi_j, \varphi_i)$  for  $i = 1, \dots, N-1, j = 0, \dots, N$  and  $\mathbf{f}^{\text{GNI}} = (f, \varphi_i)_N$  for  $i = 1, \dots, N-1$ :

$$A^{\text{GNI}} = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ \vdots & \ddots & & & \vdots \\ \vdots & & a_N(\varphi_j, \varphi_i) & & \vdots \\ \vdots & & & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix}, \quad \mathbf{f}^{\text{GNI}} = \begin{bmatrix} 0 \\ \vdots \\ f_i^{\text{GNI}} \\ \vdots \\ 0 \end{bmatrix}$$

Given that  $a(u, v) = \int_{-1}^1 \mu u' v' - \int_{-1}^1 b u v' + \int_{-1}^1 \sigma u v$  and  $(f, v) = \int_{-1}^1 f v$ . We established that  $a_n(u, v) = (\mu u', v')_N - (b u, v')_N + (\sigma u, v)_N$  and that  $(f, v)_N = (f, v)_N$ , so we obtain

$$A_{ij}^{\text{GNI}} = a_N(\varphi_j, \varphi_i) = \underbrace{(\mu \varphi_j', \varphi_i')_N}_{\text{A}} - \underbrace{(b \varphi_j, \varphi_i')_N}_{\text{B}} + \underbrace{(\sigma \varphi_j, \varphi_i)_N}_{\text{C}}$$

Assuming  $\mu, b, \sigma \in \mathbb{R}$  we have that

$$C : \sigma(\varphi_j, \varphi_i)_N = \sigma \delta_{ij} w_i = \begin{cases} \sigma w_i & i = j \\ 0 & i \neq j \end{cases} \rightarrow M = \sigma \underbrace{\begin{bmatrix} w_0 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_N \end{bmatrix}}_{\text{diagonal weight matrix}}$$

$$B : -b(\varphi_j, \varphi_i')_N = -b \sum_{k=0}^N \underbrace{\varphi_j(x_k)}_{\delta_{jk}} \underbrace{\varphi_i'(x_k)}_{D_{ki} \neq 0} w_k \rightarrow \text{full matrix}$$

$$A : \mu(\varphi_j', \varphi_i')_N = \mu \sum_{k=0}^N \underbrace{\varphi_j'(x_k)}_{D_{kj}} \underbrace{\varphi_i'(x_k)}_{D_{ki}} w_k \rightarrow \text{full matrix}$$

where  $D = (D_{ki}) = \varphi_k'(x_i)$  is the differentiation matrix that can be computed only once. The computation of  $(f, \varphi_i)_N$  can be made this way

$$(f, \varphi_i)_N = \sum_m w_m f(x_m) \underbrace{\varphi_i(x_m)}_{\delta_{im}} = w_i f(x_i)$$

In conclusion the GNI method is still as full as the spectral one, but much easier to compute thanks to the nodal expansion.

## Accuracy

We can define the Global Lagrange polynomial of degree  $N$  that interpolates  $u$  at LGL nodes as:

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

And the interpolation error, for any  $u \in H^{s+1}(-1, 1)$  with  $s \geq 0$ , the interpolation error  $u - I_N u$  satisfies the estimate:

$$\|u - I_N u\|_{H^k(-1, 1)} \leq C(s) \left( \frac{1}{N} \right)^{s+1-k} \|u\|_{H^{s+1}(-1, 1)} \quad \text{for } k = 0, 1$$

One important feature of LGL nodes is that they are not uniformly spaced (otherwise there could be problems), so that

$$I_n u(x_k) = u(x_k) \quad 0 \leq k \leq N$$

It's also possible to estimate the  $L^2$  norm of the error as:

$$\|u - I_N u\|_{L^2(-1,1)} \leq C(s) \left(\frac{1}{N}\right)^{s+1} \|u\|_{H^{s+1}(-1,1)} \quad s \geq 1$$

**Theorem 2.1** (Quadrature error)

$\exists c > 0 : \forall f \in H^q(-1,1)$ , with  $q \geq 1$ ,  $\forall v_N \in \mathbb{P}_N$  it holds

$$\left| \int_{-1}^1 f v_N dx - (f, v_N)_N \right| \geq c \left(\frac{1}{N}\right)^q \|f\|_{H^q(-1,1)} \|v_N\|_{L^2(-1,1)}$$

Let now  $u_N^{\text{GNI}} \in V_N$  be the solution of

$$a_N(u_N^{\text{GNI}}, v_N) = (f, v_N)_N \quad \forall v_N \in V_N$$

If  $u \in H^{s+1}(\Omega)$  and  $f \in H^s(\Omega)$  with  $s \geq 0$ , then:

$$\|u - u_N^{\text{GNI}}\|_{H^1(\Omega)} \leq C(s) \left(\frac{1}{N}\right)^s \left( \|u\|_{H^{s+1}(\Omega)} + \|f\|_{H^s(\Omega)} \right)$$

So  $u_N^{\text{GNI}}$  converges with spectral accuracy w.r.t. to  $N$  to the exact solution when the latter is smooth.

## General ideas

The idea proposed until now are the following:

(WP)	$V$ Hilbert	$a$ bilinear form	$F$ functional
(SG)	$V_h$ instead of $V$	same $a$	same $F$
(GNI)	$V_N$	$a_N$	$F_N$

- For the Galerkin method one can use Ceà Lemma

$$\begin{aligned} \|u - u_N\|_{H^1(\Omega)} &\leq \underbrace{\inf_{v_N \in V_N} \|u - v_N\|_{H^1(\Omega)}}_{\text{distance of } V \text{ from } V_N} \\ &\leq \|u - I_n u\|_{H^1(\Omega)} \end{aligned}$$

- For the Galerkin with Numerical Integration we need something more:

$$\begin{aligned} \|u - u_N\|_{H^1(\Omega)} &\leq \text{“distance” of } V \text{ from } V_N \\ &\quad + \text{“distance” of } a(\cdot, \cdot) \text{ from } a_N(\cdot, \cdot) \\ &\quad + \text{“distance” of } F(\cdot) \text{ from } F_N(\cdot) \end{aligned}$$

## 2.5 Strang Lemma

**Lemma 2.1** (Strang lemma)

Consider the problem

$$\text{find } u \in V : a(u, v) = F(v) \quad \forall v \in V \quad (2.4)$$

and its approximation

$$\text{find } u_h \in V_h : a_h(u_h, v_h) = F_h(v_h) \quad \forall v_h \in V_h \quad (2.5)$$

with  $\{V_h\}$  being a family of subspaces of  $V$ . Suppose that  $a_h(\cdot, \cdot)$  is continuous on  $V_h \times V_h$  and uniformly



coercive on  $V_h$  meaning that:

$$\exists \alpha^* > 0 \text{ independent of } h : a_h(v_h, v_h) \geq \alpha^* \|v_h\|_V^2 \quad \forall v_h \in V_h$$

Also suppose that  $F_h$  is linear and bounded on  $V_h$ . Then:

- exist a unique solution  $u_h$  to the problem.
- such solution depends continuously on the data, i.e. we have

$$\|u_h\|_V \leq \frac{1}{\alpha^*} \sup_{v_h \in V_h \setminus \{0\}} \frac{F_h(v_h)}{\|v_h\|_V}$$

- finally, the following a priori error estimate holds

$$\begin{aligned} \|u - u_h\|_V &\leq \inf_{w_h \in V_h} \left\{ \left(1 + \frac{M}{\alpha^*}\right) \|u - w_h\|_V \right. \\ &\quad \left. + \frac{1}{\alpha^*} \sup_{v_h \in V_h \setminus \{0\}} \frac{|a(w_h, v_h) - a_h(w_h, v_h)|}{\|v_h\|_V} \right\} \\ &\quad + \frac{1}{\alpha^*} \sup_{v_h \in V_h \setminus \{0\}} \frac{|F(v_h) - F_h(v_h)|}{\|v_h\|_V} \end{aligned}$$

with  $M$  being the continuity constant of  $a(\cdot, \cdot)$

**Proof.** The assumption of Lax-Milgram are satisfied for (2.5), so the solution exists and is unique. Moreover

$$\|u_h\|_V \leq \frac{1}{\alpha^*} \|F_h\|_{V_h'}$$

with  $\|F_h\|_{V_h'} = \sup_{v_h \in V_h \setminus \{0\}} \frac{F_h(v_h)}{\|v_h\|_V}$  being the norm of the dual space  $V_h'$ .

Now the only thing missing is the error inequality. Let  $w_h$  be any function of the subspace  $V_h$ . Setting  $\sigma_h = u_h - w_h \in V_h$ , we have:

$$\begin{aligned} \alpha^* \|\sigma_h\|_V^2 &\leq a_h(\sigma_h, \sigma_h) && \text{(by coercivity of } a_h) \\ &= a_h(u_h, \sigma_h) - a_h(w_h, \sigma_h) \\ &= F_h(\sigma_h) - a_h(w_h, \sigma_h) && \text{(by (2.5))} \\ &= F_h(\sigma_h) - F(\sigma_h) + F(\sigma_h) - a_h(w_h, \sigma_h) \\ &= [F_h(\sigma_h) - F(\sigma_h)] + a(u, \sigma_h) - a_h(w_h, \sigma_h) && \text{(by (2.4))} \\ &= [F_h(\sigma_h) - F(\sigma_h)] + a(u - w_h, \sigma_h) + [a(w_h, \sigma_h) - a_h(w_h, \sigma_h)] \end{aligned}$$

If  $\sigma_h \neq 0$ , we can divide everything by  $\alpha^* \|\sigma_h\|_V$

$$\begin{aligned} \|\sigma_h\|_V &\leq \frac{1}{\alpha^*} \left\{ \frac{|F_h(\sigma_h) - F(\sigma_h)|}{\|\sigma_h\|_V} + \frac{|a(u - w_h, \sigma_h)|}{\|\sigma_h\|_V} + \frac{|a(w_h, \sigma_h) - a_h(w_h, \sigma_h)|}{\|\sigma_h\|_V} \right\} \\ &\leq \frac{1}{\alpha^*} \left\{ M \|u - w_h\|_V + \sup_{v_h \in V_h \setminus \{0\}} \frac{|a(w_h, v_h) - a_h(w_h, v_h)|}{\|v_h\|_V} + \sup_{v_h \in V_h \setminus \{0\}} \frac{|F_h(v_h) - F(v_h)|}{\|v_h\|_V} \right\} \end{aligned}$$

Clearly, if  $\sigma_h = 0$ , the inequality still holds.

We can now estimate the error between  $u$  and  $u_h$ . Since  $u - u_h = (u - w_h) - \sigma_h$  we obtain

$$\begin{aligned}
\|u - u_h\| &\leq \|u - w_h\|_V + \|\sigma_h\|_V \\
&\leq \|u - w_h\|_V + \frac{1}{\alpha^*} \left\{ M \|u - w_h\|_V + \sup_{v_h \in V_h \setminus \{0\}} \frac{|a(w_h, v_h) - a_h(w_h, v_h)|}{\|v_h\|_V} \right. \\
&\quad \left. + \sup_{v_h \in V_h \setminus \{0\}} \frac{|F_h(\sigma_h) - F(\sigma_h)|}{\|v_h\|_V} \right\} \\
&= \left(1 + \frac{M}{\alpha^*}\right) \|u - w_h\|_V + \frac{1}{\alpha^*} \sup_{v_h \in V_h \setminus \{0\}} \frac{|a(w_h, v_h) - a_h(w_h, v_h)|}{\|v_h\|_V} \\
&\quad + \sup_{v_h \in V_h \setminus \{0\}} \frac{|F_h(\sigma_h) - F(\sigma_h)|}{\|v_h\|_V}
\end{aligned}$$

If this inequality holds  $\forall w_h \in V_h$ , then it holds when taking the infimum. ★

Now we should try to apply Strang's lemma to GNI method in one dimension, to verify its convergence. Obviously, we will have  $V_N$  instead of  $V_h$  and everything that follows from there.

First of all, the error of the LGL numerical integration formula

$$E(g, v_N) = (g, v_N) - (g, v_N)_N$$

with  $g$  and  $v_N$  being a generic continuous function and a generic polynomial of  $\mathbb{Q}_N$  respectively. Introducing the interpolation polynomial  $I_N g$ , we obtain:

$$\begin{aligned}
E(g, v_N) &= (g, v_N) - (I_N g, v_N) \\
&= (g, v_N) - (I_{N-1} g, v_N) + \underbrace{(I_{N-1} g, v_N)}_{\in \mathbb{Q}_{2N-1}} \\
&= (g, v_N) - (I_{N-1} g, v_N) + (I_{N-1} g, v_N)_N - (I_N g, v_N)_N \\
&= (g - I_{N-1} g, v_N) + (I_{N-1} g - I_N g, v_N)_N
\end{aligned}$$

The first summand of the right-hand side can be bounded from above using Cauchy-Schwartz:

$$|(g - I_{N-1} g, v_N)| \leq \|g - I_{N-1} g\|_{L^2(-1,1)} \|v_N\|_{L^2(-1,1)}$$

For the second term, it's a bit more difficult, we need to introduce two new lemmas

### Lemma 2.2

The discrete scalar product  $(\cdot, \cdot)_N$  is a scalar product on  $\mathbb{Q}_N$  and, as such, it satisfies the Cauchy-Schwartz inequality

$$|(\varphi, \psi)_N| \leq \|\varphi\|_N \|\psi\|_N$$

where the discrete norm is defined as

$$\|\varphi\|_N = \sqrt{(\varphi, \varphi)_N} \quad \forall \varphi \in \mathbb{Q}_N$$

### Lemma 2.3

The “continuous” norm of  $L^2(-1, 1)$  and the “discrete” norm  $\|\cdot\|_N$  verify the inequalities

$$\|v_N\|_{L^2(-1,1)} \leq \|v_N\|_N \leq \sqrt{3} \|v_N\|_{L^2(-1,1)}$$

hence they are uniformly equivalent on  $\mathbb{Q}_N$

By using these two lemmas we are able to obtain

$$\begin{aligned}
|(I_{N-1} g - I_N g, v_N)_N| &\leq \|I_{N-1} g - I_N g\|_N \|v_N\|_N \\
&\leq 3 \left[ \|I_{N-1} g - g\|_{L^2(-1,1)} + \|I_N g - g\|_{L^2(-1,1)} \right] \|v_N\|_{L^2(-1,1)}
\end{aligned}$$

Putting all together we obtain the upper bound

$$|E(g, v_N)| \leq \left[ 4\|I_{N-1}g - g\|_{L^2(-1,1)} + 3\|I_{N-1}g - g\|_{L^2(-1,1)} \right] \|v_N\|_{L^2(-1,1)}$$

Using then the interpolation estimate

$$\|f - I_N f\|_{H^k(-1,1)} \leq C(s) \left( \frac{1}{N} \right)^{s-k} \|f\|_{H^s(-1,1)} \quad s \geq 1, k = 0, 1$$

we can bound  $|E(g, v_N)|$  even more

$$|E(g, v_N)| \leq C(s) \left[ \left( \frac{1}{N-1} \right)^s + \left( \frac{1}{N} \right)^s \right] \|g\|_{H^s(-1,1)} \|v_N\|_{L^2(-1,1)}$$

assuming that  $g \in H^s(-1, 1)$ .

Then, since for each  $N \geq 2$  we have that  $\frac{1}{N-1} \leq \frac{2}{N}$ , the error for the LGL integration can be written as

$$|E(g, v_N)| \leq C(s) \left( \frac{1}{N} \right)^s \|g\|_{H^s(-1,1)} \|v_N\|_{L^2(-1,1)}$$

## 2.6 GNI as Collocation method

Let us introduce a problem

$$\begin{cases} Lu = -(\mu u')' + (bu)' + \sigma u = f & -1 < x < 1 \\ u(-1) = u(1) = 0 \end{cases}$$

that has the usual weak formulation

$$\text{find } u \in V = H_0^1(-1, 1) : a(u, v) = F(v), \forall v \in V$$

The GNI formulation follows

$$\begin{cases} \text{find } u_N \in V_N = \mathbb{P}_N^0 = \{v_N \in \mathbb{P}_N : v_N(\pm 1) = 0\} \\ a_N(u_N, v_N) = F_N(v_N) \quad \forall v_N \in V_N \end{cases}$$

Note that, thanks to the exactness of LGL quadrature formula

$$\begin{aligned} a_N(u_N, v_N) &\stackrel{(\text{def. of } I_N)}{=} (I_N(\underbrace{\mu u'_N - bu_N}_{\in \mathbb{P}_N}), \underbrace{v'_N}_{\in \mathbb{P}_{N-1}})_N + (\sigma u_N, v_N)_N \\ &\stackrel{(\text{exactness})}{=} (I_N(\mu u' - bu), v_N)_N + (\sigma u_N, v_N)_N \\ &\stackrel{(\text{int. by parts})}{=} -(\underbrace{I_N(\mu u' - bu)'}_{\in \mathbb{P}_{N-1}}, \underbrace{v_N}_{\in \mathbb{P}_N})_N + (\sigma u_N, v_N)_N \\ &\stackrel{(\text{exactness})}{=} (\underbrace{-(I_N(\mu u' - bu))' + \sigma u_N}_{= L_N u_N}, v_N)_N \end{aligned}$$

So it's obvious that  $(\text{GNI}) \iff (L_N u_N, v_N)_N = F_N(v_N) \quad \forall v_N \in V_N$ , so it's a collocation method.

## 2.7 1D Spectral Elements

Let  $p \geq 1$  integer and  $\mathbb{P}_p$  the space of polynomials of degree  $\leq p$ . We can divide the domain  $\Omega = \bigcup_{n=1}^{N_e} I_k$  with  $I_k$  disjoint elements s.t.  $I_k = F_k((-1, 1))$  and

$$F_k : \xi \mapsto x = \frac{b_k - a_k}{2} \xi + \frac{b_k + a_k}{2}$$

with  $N_p = p \cdot N_e + 1$  the total number of nodes in  $\Omega$ . Then we use the Lagrange basis functions  $\{\varphi_i\}_{i=1}^{N_p}$  w.r.t. the LGL nodes.

Now set  $X_\delta = \{v \in \mathcal{C}^0 : v|_{I_k} \in \mathbb{P}_p, \forall I_k\}$  with  $h_k = \text{meas}(I_k)$ , mesh size  $h = \max_k h_k$  and polynomial degree  $p$  we can define  $\delta = (h, p)$  and

$$v_\delta(x) = \sum_{i=1}^{N_p} v_\delta(x_i) \varphi_i(x) \quad \forall v_\delta \in X_\delta$$

Let now  $(\hat{\xi}_j, \hat{w}_j)$  for  $j = 0, \dots, p$  be the LGL nodes and respective weights in  $\hat{\Omega} = (-1, 1)$ . We can define the local LGL quadrature as

$$\int_{I_k} u(x)v(x) dx \approx (u, v)_{\delta, I_k} = \sum_{j=0}^p u(\xi_j)v(\xi_j)w_j$$

with  $\xi_j = \frac{b_k - a_k}{2} \hat{\xi}_j + \frac{b_k + a_k}{2}$  and  $w_j = \frac{b_k - a_k}{2} \hat{w}_j$ . Meanwhile we can pass this quadrature to the whole domain, obtaining the composite LGL quadrature:

$$\int_{\Omega} u(x)v(x) dx \approx (u, v)_{\delta, \Omega} = \sum_{k=1}^{N_e} (u, v)_{\delta, I_k}$$

with its relative error  $\exists c > 0 : \forall f \in H^r(\Omega), r \geq 1, p \geq 1 : \forall v_\delta \in X_\delta$ :

$$\left| \int_{\Omega} f v_\delta dx - (f, v_\delta)_{\delta, \Omega} \right| \leq c h^{\min(p, r)} \left( \frac{1}{p} \right)^r \|f\|_{H^r(\Omega)} \|v_\delta\|_{L^2(\Omega)}$$

and its interpolation error as:  $\exists c > 0 : \forall v \in H^{s+1}(\Omega), s \geq 1$

$$\left\| v - \Pi_\delta^{LGL} v \right\|_{H^k(\Omega)} \leq C h^{\min(p+1, s+1)-k} \left( \frac{1}{p} \right)^{s+1-k} \|v\|_{H^{s+1}(\Omega)}$$

## 2.8 Spectral Element Method with Numerical Integration

Let's go back to the problem

$$\begin{cases} -(\mu u')' - (bu)' + \sigma u = f & \text{in } \Omega \\ u(a) = u(b) = 0 \end{cases}$$

Given  $V = H_0^1(\Omega)$ , the weak formulation reads

$$\text{find } u \in V : a(u, v) = (f, v)_{L^2(\Omega)} \quad \forall v \in V, f \in L^2(\Omega)$$

with

$$\begin{aligned} a(u, v) &= \int_{\Omega} (\mu u' - bu)v' dx + \int_{\Omega} \sigma u v dx \\ (f, v)_{L^2(\Omega)} &= \int_{\Omega} f v dx \end{aligned}$$

Now set  $a_\delta(\varphi_j, \varphi_i) = (\mu \varphi_j' - b \varphi_j, \varphi_i)_{\delta, \Omega} + (\sigma \varphi_j, \varphi_i)_{\delta, \Omega}$  to get the SEM-GNI formulation:

$$\text{find } u_\delta^{\text{GNI}} \in V_\delta : a_\delta(u_\delta^{\text{GNI}}, v_\delta) = (f, v_\delta)_{\delta, \Omega} \quad \forall v_\delta \in V_\delta \quad (2.6)$$

Now expand  $u_\delta^{\text{GNI}}$  w.r.t the Lagrange basis  $u_\delta^{\text{GNI}}(x) = \sum_{i=1}^{N_p} u_\delta^{\text{GNI}}(x_i) \varphi_i(x)$  and choose  $v_\delta(x) = \varphi_i(x)$  for any  $i = 1, \dots, N_p$ . We can now write the SEM-GNI discretization of the weak formulation:

$$\begin{cases} \text{find } u^{\text{GNI}} = [u_\delta^{\text{GNI}}(x_j)]_{j=1}^{N_p} \\ u_\delta^{\text{GNI}}(x_1) = u_\delta^{\text{GNI}}(x_{N_p}) = 0 \\ \sum_{j=1}^{N_p} a_\delta(\varphi_j, \varphi_i) u_\delta^{\text{GNI}}(x_j) = (f, \varphi_i)_{\delta, \Omega} \quad \forall i = 1, \dots, N_p \end{cases}$$

or, in algebraic form  $A^{\text{GNI}} u^{\text{GNI}} = f^{\text{GNI}}$  with  $A_{ij}^{\text{GNI}} = a_\delta(\varphi_j, \varphi_i)$  and  $f_i^{\text{GNI}} = (f, \varphi_i)_{\delta, \Omega}$ . Now, for the error analysis, we will apply the Strang lemma, so:

$$\begin{aligned} \|u - u_\delta^{\text{GNI}}\|_V &\leq \|u - u_\delta\|_V \\ &\quad + \frac{1}{\mu^*} \sup_{v_\delta \in V_\delta \setminus \{0\}} \frac{|a(u_\delta, v_\delta) - a_\delta(u_\delta, v_\delta)|}{\|v_\delta\|_V} \\ &\quad + \frac{1}{\mu^*} \sup_{v_\delta \in V_\delta \setminus \{0\}} \frac{|f, v_\delta|_{L^2(\Omega) - (f, v_\delta)_{\delta, \Omega}}}{\|v_\delta\|_V} \end{aligned}$$

where  $\mu^*$  is the coercivity constant of  $a_\delta$ :  $a_\delta(v_\delta, v_\delta) \geq \mu^* \|v_\delta\|_V^2$  and  $u_\delta$  the SEM-GNI solution. Thus for any  $u \in H^{s+1}(\Omega)$  and  $f \in H^r(\Omega)$

$$\|u - u_\delta^{\text{GNI}}\|_{H^1(\Omega)} \leq C \left[ h^{\min(p, s)} \left(\frac{1}{p}\right)^s \|u\|_{H^{s+1}(\Omega)} + h^{\min(p, r)} \left(\frac{1}{p}\right)^r \|f\|_{H^r(\Omega)} \right]$$

So  $u_\delta^{\text{GNI}}$  converges with spectral accuracy w.r.t.  $p$  and algebraic accuracy w.r.t.  $h$  to the exact solution.

## 2.9 Convergence rate of SEM-GNI

When  $s, r$  are large ( $s, r > p$ ):

$$\|u - u_\delta\|_{H^1(\Omega)} \leq C \left[ h^p \left(\frac{1}{p}\right)^s \|u\|_{H^{s+1}(\Omega)} + h^p \left(\frac{1}{p}\right)^r \|f\|_{H^r(\Omega)} \right]$$

when  $s$  is small ( $s \leq p$ ):

$$\|u - u_\delta\|_{H^1(\Omega)} \leq C \left(\frac{h}{p}\right)^s \|u\|_{H^{s+1}(\Omega)}$$

### 3 Discontinuous Galerkin methods

The idea behind DG methods is to seek the solution in a discrete space made of polynomials that are completely discontinuous across the elements of the mesh.

$$V_h \subsetneq V$$

#### 3.1 1D case

Let us consider a Poisson problem

$$\begin{cases} -u'' = f & a < x < b \\ u(a) = u(b) = 0 \end{cases}$$

The aim is to use discontinuous piecewise polynomials, so that between every interval  $I_k$  from one node to another we obtain

$$\int_a^b -u''v = \int_a^b fv \Rightarrow - \sum_{k=0}^{N-1} \int_{I_k} u''v = \sum_{k=0}^{N-1} \int_{I_k} fv$$

We must know integrate by parts, but our test functions are discontinuous at the nodes, so we must acknowledge it. Let's call  $x_k^-$  and  $x_k^+$  the left and right side of the  $x_k$  node. Then we can:

$$- \sum_{k=0}^{N-1} \int_{I_k} u''v = \sum_{k=0}^{N-1} \left[ \int_{I_k} u'v' - \left( u'v|_{x_{k+1}^-} - u'v|_{x_k^+} \right) \right] \quad (3.1)$$

$$\begin{aligned} \sum_{k=0}^{N-1} (u'v|_{x_{k+1}^-} - u'v|_{x_k^+}) &= u'(x_1^-)v(x_1^-) - u'(x_0^+)v(x_0^+) \\ &\quad + u'(x_2^-)v(x_2^-) - u'(x_1^+)v(x_1^+) \\ &\quad + \dots \\ &\quad + u'(x_N^-)v(x_N^-) - u'(x_{N-1}^+)v(x_{N-1}^+) \\ &= \sum_{k=0}^N \llbracket u'(x_k)v(x_k) \rrbracket \end{aligned} \quad (3.2)$$

where we have defined the jump function

$$\begin{aligned} \llbracket \varphi(x_0) \rrbracket &:= -\varphi(x_0^+) \\ \llbracket \varphi(x_k) \rrbracket &:= \varphi(x_k^-) - \varphi(x_k^+) & x_k : \text{interior node} \\ \llbracket \varphi(x_N) \rrbracket &:= \varphi(x_N^-) \end{aligned} \quad (3.3)$$

By using (3.1) and (3.3) we obtain

$$\sum_{k=0}^{N-1} \int_{I_k} u'v' - \sum_{k=0}^N \llbracket u'(x_k)v(x_k) \rrbracket = \sum_{k=0}^{N-1} \int_{I_k} fv \quad (3.4)$$

Now define the average operator

$$\begin{aligned} \{\!\!\{ \varphi(x_0) \}\!\!\} &:= \varphi(x_0^+) \\ \{\!\!\{ \varphi(x_k) \}\!\!\} &:= \frac{1}{2} \varphi(x_k^-) + \varphi(x_k^+) & x_k : \text{interior node} \\ \{\!\!\{ \varphi(x_N) \}\!\!\} &:= \varphi(x_N^-) \end{aligned} \quad (3.5)$$

This way we obtain this formula

$$\sum_{k=0}^N \llbracket u'(x_k) v(x_k) \rrbracket = \sum_{k=0}^N \{\!\!\{ u'(x_k) \}\!\!\} \llbracket v(x_k) \rrbracket + \sum_{k=1}^{N-1} \llbracket u'(x_k) \rrbracket \{\!\!\{ v(x_k) \}\!\!\} \quad (3.6)$$

If  $u$  is the exact solution and  $u \in \mathcal{C}^1([a, b])$ , then  $\llbracket u'(x_k) \rrbracket = 0$  for every interior node, and the second sum in (3.6) drops.

We end up with the formulation (by collecting (3.4) and (3.6))

$$\underbrace{\sum_{k=0}^{N-1} \int_{I_k} u' v' - \sum_{k=0}^N \{\!\!\{ u'(x_k) \}\!\!\} \llbracket v(x_k) \rrbracket - \sum_{k=1}^{N-1} \llbracket u'(x_k) \rrbracket \{\!\!\{ v(x_k) \}\!\!\}}_{\mathcal{A}(u, v)} = \sum_{k=0}^{N-1} \int_{I_k} f v \quad \forall v \in V \quad (3.7)$$

where

$$V = H_{\text{broken}}^1(\Omega) := \{v \in L^2(\Omega) : v|_{I_k} \in H^1 I_k \ \forall k = 0, \dots, N-1\}$$

with the broken norm

$$\|v\|_{H_{\text{broken}}^1(\Omega)} = \left( \sum_{k=0}^N \|v|_{I_k}\|_{H^1(\Omega)}^2 \right)^{\frac{1}{2}}$$

Let now  $V_h \subset V$

$$\text{find } u_h \in V_h : \mathcal{A}(u_h, v_h) = \sum_{k=0}^{N-1} \int_{I_k} f v_h \quad \forall v_h \in V_h \quad (3.8)$$

**Remark 3.1**

$V_h$  is not a subspace of  $H^1(\Omega)$

But (3.8) is not well posed, so the (3.7) must be modified such that:

- drop  $3^{rd}$  term because  $\llbracket u'(x_k) \rrbracket = 0$
- add symmetrization term ( $= 0$  if  $u$  is the exact solution)

$$- \sum_{k=0}^N \theta \{\!\!\{ v'(x_k) \}\!\!\} \llbracket u(x_k) \rrbracket$$

with

- $\theta = 1$  SIP (Symmetric Interior Penalty)
- $\theta = -1$  NIP (Non-symmetric Interior Penalty)
- $\theta = 0$  IIP (Incomplete Interior Penalty)

- add the stabilization term ( $= 0$  if  $u$  is the exact solution)

$$+ \sum_{k=0}^N \gamma \llbracket u(x_k) \rrbracket \llbracket v(x_k) \rrbracket$$

We can now obtain a new bilinear form

$$\begin{aligned} \mathcal{A}^*(u_h, v_h) = & \underbrace{\sum_{k=0}^{N-1} \int_{I_k} u_h' v_h'}_{(1)} - \underbrace{\sum_{k=0}^N \{\!\!\{ u_h'(x_k) \}\!\!\} \llbracket v_h(x_k) \rrbracket}_{(2)} \\ & - \underbrace{\sum_{k=0}^N \theta \{\!\!\{ v_h'(x_k) \}\!\!\} \llbracket u_h(x_k) \rrbracket}_{(3)} + \underbrace{\sum_{k=0}^N \gamma \llbracket u_h(x_k) \rrbracket \llbracket v_h(x_k) \rrbracket}_{(4)} \end{aligned} \quad (3.9)$$

### Neumann BC

Impose Neumann BC through  $\{u'(x_k)\}$  in (2). In this case we have  $\sum_{k=1}^{N-1}$  in (2) and, consequently, we write  $\sum_{k=1}^{N-1}$  in (3) for symmetry.

### Non-homogeneous Dirichlet BC

Impose Dirichlet BC as follows. In (3) and (4) replace  $\llbracket u_h(x_0) \rrbracket$  and  $\llbracket u_h(x_N) \rrbracket$  with the following definition:

$$\begin{aligned}\llbracket u_h(x_0) \rrbracket &:= \alpha - u_h(x_0^+) && \text{if } u(a) = \alpha \\ \llbracket u_h(x_N) \rrbracket &:= u_h(x_N^-) - \beta && \text{if } u(b) = \beta\end{aligned}$$

In case  $\alpha = \beta = 0$  we have homogeneous Dirichlet.

Now in (3.9) split sums as follows:

$$\begin{aligned}\mathcal{A}^*(u_h, v_h) &= \sum_{k=0}^{N-1} \int_{I_k} u'_h v'_h \\ &\quad - \sum_{k=1}^{N-1} \{ \{u'_h(x_k)\} \llbracket v_h(x_k) \rrbracket + u'_h(x_0^+) v_h(x_0^+) - u'_h(x_N^-) v_h(x_N^-) \\ &\quad - \sum_{k=1}^{N-1} \theta \{ \{v'_h(x_k)\} \llbracket u_h(x_k) \rrbracket - [\theta v'_h(x_0^+) (\alpha - u_h(x_0^+)) + \theta v'_h(x_N^-) (u_h(x_N^-) - \beta)] \\ &\quad - \sum_{k=1}^{N-1} \gamma \llbracket u_h(x_k) \rrbracket \llbracket v_h(x_k) \rrbracket + \gamma (\alpha - u_h(x_0^+)) (-v_h(x_0^+)) + \gamma (u_h(x_N^-) - \beta) v_h(x_N^-) \end{aligned} \quad (3.10)$$

Now move terms, including  $\alpha$  and  $\beta$  to the right hand side of the formulation.

On the left hand side it remains

$$\begin{aligned}\tilde{\mathcal{A}}(u_h, v_h) &= \sum_{k=0}^{N-1} \int_{I_k} u'_h v'_h \\ &\quad - \sum_{k=1}^{N-1} \{ \{u'_h(x_k)\} \llbracket v_h(x_k) \rrbracket + u'_h(x_0^+) v_h(x_0^+) - u'_h(x_N^-) v_h(x_N^-) \\ &\quad - \sum_{k=1}^{N-1} \theta \{ \{v'_h(x_k)\} \llbracket u_h(x_k) \rrbracket - [\theta v'_h(x_0^+) u_h(x_0^+) + \theta v'_h(x_N^-) u_h(x_N^-)] \\ &\quad - \sum_{k=1}^{N-1} \gamma \llbracket u_h(x_k) \rrbracket \llbracket v_h(x_k) \rrbracket + \gamma u_h(x_0^+) v_h(x_0^+) + \gamma u_h(x_N^-) v_h(x_N^-) \end{aligned} \quad (3.11)$$

On the right hand side instead

$$\mathcal{F}(v_h) = \sum_{k=0}^{N-1} \int_{I_k} f v_h + \theta (\alpha v'_h(x_0^+) - \beta v'_h(x_N^-)) + \gamma (\alpha v_h(x_0^+) + \beta v_h(x_N^-)) \quad (3.12)$$

### Remark 3.2

Note that for  $\theta = 1$ ,  $\tilde{\mathcal{A}}(u_h, v_h) = \tilde{\mathcal{A}}(v_h, u_h)$ , so it's symmetric.



### Non-homogeneous Dirichlet conditions

$$\text{find } u_h \in V_h : \tilde{\mathcal{A}}(u_h, v_h) = \mathcal{F}(v_h) \quad \forall v_h \in V_h \quad (3.13)$$

with  $\mathcal{F}$  depending on  $f$ ,  $\alpha$  and  $\beta$ .

Note that in (3.11), if we define  $\llbracket u_h(x_0) \rrbracket$  and  $\llbracket u_h(x_N) \rrbracket$  as  $\llbracket v_h(x_0) \rrbracket$  and  $\llbracket v_h(x_N) \rrbracket$

$$\begin{aligned} & - \sum_{k=1}^{N-1} \{ \{ u'_h(x_k) \} \} \llbracket v_h(x_k) \rrbracket + u'_h(x_0^+) v_h(x_0^+) - u'_h(x_N^-) v_h(x_N^-) \\ & = - \sum_{k=0}^N \{ \{ u'_h(x_k) \} \} \llbracket v_h(x_k) \rrbracket \\ & - \sum_{k=1}^{N-1} \theta \{ \{ v'_h(x_k) \} \} \llbracket u_h(x_k) \rrbracket + (\theta u_h(x_0^+) v'_h(x_0^+) - \theta u_h(x_N^-) v'_h(x_N^-)) \\ & = - \sum_{k=0}^N \theta \{ \{ v'_h(x_k) \} \} \llbracket u_h(x_k) \rrbracket \\ & + \sum_{k=1}^{N-1} \gamma \llbracket u_h(x_k) \rrbracket \llbracket v_h(x_k) \rrbracket + \gamma u_h(x_0^+) v_h(x_0^+) + \gamma u_h(x_N^-) v_h(x_N^-) \\ & + \sum_{k=0}^N \gamma \llbracket u_h(x_k) \rrbracket \llbracket v_h(x_k) \rrbracket \end{aligned}$$

### 3.2 Multidimensional case

We can take our Poisson problem in multidimension

$$\begin{cases} -\Delta u = f & \text{in } \Omega \\ u = 0 & \text{on } \partial\Omega \end{cases} \quad (3.14)$$

with the triangulation  $\mathcal{T}_h$ , but this time we cannot assume that the conformity constraint is present. So we need to take a test function  $v$  (element-wise smooth), and integrate over an element  $\mathcal{K} \in \mathcal{T}_h$

$$\int_{\mathcal{K}} -\Delta u v = \int_{\mathcal{K}} f v$$

As usual, integrate by parts, and sum over all the elements  $\mathcal{K} \in \mathcal{T}_h$

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\partial\mathcal{K}} \nabla u \cdot \mathbf{n}_{\mathcal{K}} v = \int_{\Omega} f v$$

since for any  $F \in \mathcal{F}'_h$  which is the set of interior faces shared by two elements  $\mathcal{K}^+$  and  $\mathcal{K}^-$

$$\begin{aligned} \{ \{ v \} \} &= \frac{(v^+ + v^-)}{2} & \llbracket v \rrbracket &= v^+ \mathbf{n}^+ + v^- \mathbf{n}^- \\ \{ \{ \boldsymbol{\tau} \} \} &= \frac{(\boldsymbol{\tau}^+ + \boldsymbol{\tau}^-)}{2} & \llbracket \boldsymbol{\tau} \rrbracket &= \boldsymbol{\tau}^+ \mathbf{n}^+ + \boldsymbol{\tau}^- \mathbf{n}^- \end{aligned}$$

while, for the set of boundary faces  $F \in \mathcal{F}_h^B$

$$\begin{aligned} \{ \{ v \} \} &= v & \llbracket v \rrbracket &= v \mathbf{n} \\ \{ \{ \boldsymbol{\tau} \} \} &= \boldsymbol{\tau} & \llbracket \boldsymbol{\tau} \rrbracket &= \boldsymbol{\tau} \cdot \mathbf{n} \end{aligned}$$

in this way we can obtain the following formula  $\forall \boldsymbol{\tau}$  vector function:

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\partial\mathcal{K}} \boldsymbol{\tau} \cdot \mathbf{n}_{\mathcal{K}} v = \sum_{F \in \mathcal{F}_h} \int_F \{ \{ \boldsymbol{\tau} \} \} \cdot \llbracket v \rrbracket + \sum_{F \in \mathcal{F}'_h} \int_F \llbracket \boldsymbol{\tau} \rrbracket \{ \{ v \} \} \quad (3.15)$$

and thanks to that we obtain

$$-\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\partial \mathcal{K}} \boldsymbol{\tau} \cdot \mathbf{n}_{\mathcal{K}} v = -\sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}'_h} \int_F \llbracket \nabla u \rrbracket \{\{v\}\}$$

then

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\partial \mathcal{K}} \nabla u \cdot \mathbf{n}_{\mathcal{K}} v = \int_{\Omega} f v$$

so it becomes

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}'_h} \int_F \llbracket \nabla u \rrbracket \{\{v\}\} = \int_{\Omega} f v$$

but, if we assume  $u \in H^2(\Omega)$ , then  $\llbracket \nabla u \rrbracket = 0 \ \forall F \in \mathcal{F}'_h$ . This regularity assumption is fulfilled if  $f \in L^2$  and the domain is a convex polygon, thanks to the property of elliptic regularity.

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}'_h} \int_F \llbracket \nabla u \rrbracket \{\{v\}\} = \int_{\Omega} f v$$

Now we can assume that  $\llbracket u \rrbracket = 0 \ \forall F \in \mathcal{F}_h$  (since  $u \in H^2(\Omega) \cap H_0^1(\Omega)$ ) to add a symmetry term

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}'_h} \int_F \{\{\nabla_h v\}\} \llbracket u \rrbracket = \int_{\Omega} f v$$

where  $\nabla_h$  is the elementwise gradient ( $v$  is only piecewise smooth). We also add a stabilization term that controls the jumps

$$\sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}_h} \int_F \llbracket u \rrbracket \cdot \{\{\nabla_h v\}\} + \sum_{F \in \mathcal{F}_h} \int_F \gamma \llbracket u \rrbracket \cdot \llbracket v \rrbracket = \int_{\Omega} f v$$

where  $\gamma$  is a stabilization function.

Now we can define the DG discrete space

$$V_h^p = \{v_h \in L^2(\Omega) : v_h|_{\mathcal{K}} \in \mathcal{P}^{p_{\mathcal{K}}}(\mathcal{K}) \ \forall \mathcal{K} \in \mathcal{T}_h\} \not\subset H_0^1(\Omega)$$

Discretize  $u \rightsquigarrow u_h, v \rightsquigarrow v_h$  and obtain the following weak formulation

$$\text{find } u_h \in V_h^p \text{ s.t. } \mathcal{A}(u_h, v_h) = \int_{\Omega} f v \quad \forall v_h \in V_h^p$$

where

$$\begin{aligned} \mathcal{A}(u, v) = & \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \sum_{F \in \mathcal{F}_h} \int_F \llbracket u \rrbracket \cdot \{\{\nabla_h v\}\} \\ & + \sum_{F \in \mathcal{F}_h} \int_F \gamma \llbracket u \rrbracket \cdot \llbracket v \rrbracket \end{aligned}$$

### Interior Penalty DG methods

$$\text{find } u_h \in V_h^p \text{ s.t. } \mathcal{A}(u_h, v_h) = \int_{\Omega} f v \quad \forall v_h \in V_h^p$$

Note that  $\mathcal{A}$  depends on the triangulation and it differs from the original weak formulation in the infinite dimension problem.

$$\begin{aligned} \mathcal{A}(u, v) = & \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{F \in \mathcal{F}_h} \int_F \{\{\nabla u\}\} \cdot \llbracket v \rrbracket - \theta \sum_{F \in \mathcal{F}_h} \int_F \llbracket u \rrbracket \cdot \{\{\nabla_h v\}\} \\ & + \sum_{F \in \mathcal{F}_h} \int_F \gamma \llbracket u \rrbracket \cdot \llbracket v \rrbracket \end{aligned}$$

where

- $\theta = 1$  Symmetric Interior Penalty (SIP)
- $\theta = -1$  Non-symmetric Interior Penalty (NIP)
- $\theta = 0$  Incomplete Interior Penalty (IIP)

### Dirichlet BC

The above formulation holds when applying homogeneous Dirichlet BC, but in the case of non-homogeneous BC, such as

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

the right hand side must be modified as

$$\int_{\Omega} f v - \theta \sum_{F \in \mathcal{F}_h^B} \int_F g_D \nabla_h v \cdot \mathbf{n} + \sum_{F \in \mathcal{F}_h^B} \int_F \gamma g_D v$$

### Neumann BC

In the case of Neumann BC, like

$$\nabla u \cdot \mathbf{n} = g_N \quad \text{on } \partial\Omega$$

the bilinear form has to be modified such as

$$\begin{aligned} \mathcal{A}(u, v) = \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \nabla u \cdot \nabla v - \sum_{F \in \mathcal{F}_h'} \int_F \{\nabla u\} \cdot \llbracket v \rrbracket - \theta \sum_{F \in \mathcal{F}_h'} \int_F \llbracket u \rrbracket \cdot \{\nabla_h v\} \\ + \sum_{F \in \mathcal{F}_h'} \int_F \gamma \llbracket u \rrbracket \cdot \llbracket v \rrbracket \end{aligned}$$

and the right hand side

$$\int_{\Omega} f v - \sum_{F \in \mathcal{F}_h^B} \int_F g_N v$$

### The stabilization function $\gamma$

$$\sum_{F \in \mathcal{F}_h} \int_F \gamma \llbracket u \rrbracket \cdot \llbracket v \rrbracket \quad \gamma = \alpha \frac{p^2}{h}$$

where

$$p = \begin{cases} \max\{p_{\mathcal{K}^+}, p_{\mathcal{K}^-}\} & \text{if } F \in \mathcal{F}_h' \\ p_{\mathcal{K}} & \text{if } F \in \mathcal{F}_h^B \end{cases}$$

and

$$h = \begin{cases} \min\{h_{\mathcal{K}^+}, h_{\mathcal{K}^-}\} & \text{if } F \in \mathcal{F}_h' \\ h_{\mathcal{K}} & \text{if } F \in \mathcal{F}_h^B \end{cases}$$

Since we can make some assumptions

$$h_F \approx h_{\mathcal{K}^+} \approx h_{\mathcal{K}^-}, \quad p_{\mathcal{K}^+} \approx p_{\mathcal{K}^-} \Rightarrow \gamma = \mathcal{O}\left(\frac{p^2}{h}\right)$$

### 3.3 Theoretical reminders

For an integer  $s \geq 1$  define the broken Sobolev space

$$H^s(\mathcal{T}_h) = \{v \in L^2(\Omega) : v|_{\mathcal{K}} \in H^s(\mathcal{K}) \ \forall \mathcal{K} \in \mathcal{T}_h\}$$

$$\|v\|_{H^s(\mathcal{K})}^2 = \sum_{\mathcal{K} \in \mathcal{T}_h} \|v\|_{H^s(\mathcal{K})}^2$$

Define also

$$\|v\|_{L^2(\mathcal{F}_h)}^2 = \sum_{F \in \mathcal{F}_h} \|v\|_{L^2(F)}^2$$

We define then the following norms

$$\|v\|_{DG}^2 = \|\nabla_h v\|_{L^2(\Omega)}^2 + \left\| \gamma^{\frac{1}{2}} \llbracket v \rrbracket \right\|_{L^2(\mathcal{F}_h)}^2 \quad \forall v \in H^2(\mathcal{T}_h)$$

$$\|v\|_{DG} = \|v\|_{DG}^2 + \left\| \gamma^{\frac{1}{2}} \{\nabla_h v\} \right\|_{L^2(\mathcal{F}_h)}^2 \quad \forall v \in H^2(\mathcal{T}_h)$$

where  $\nabla_h v$  is the elementwise gradient:

$$(\nabla_h v)|_{\mathcal{K}} = \nabla(v|_{\mathcal{K}}) \quad \forall \mathcal{K} \in \mathcal{T}_h$$

Notice that  $V_h^p \subset H^2(\mathcal{T}_h)$ . It can be shown that

$$\|v\|_{DG} \underset{(trivial)}{\leq} \|v\|_{DG} \lesssim \|v\|_{DG} \quad \forall v \in H^2(\mathcal{T}_h)$$

$$\|v_h\|_{DG} \underset{(trivial)}{\leq} \|v_h\|_{DG} \underset{(on \ slides)}{\lesssim} \|v_h\|_{DG} \quad \forall v_h \in V_h^p$$

Some key properties:

- Continuity on  $H^2(\mathcal{T}_h) \times V_h^p$ :

$$|\mathcal{A}(v, w_h)| \lesssim \|v\|_{DG} \|w_h\|_{DG} \quad \forall v \in H^2(\mathcal{T}_h), \forall w_h \in V_h^p$$

Also remind that  $|\mathcal{A}(v, w_h)| \lesssim \|v\|_{DG} \|w_h\|_{DG}$

- Coercivity on  $V_h^p \times V_h^p$ :

$$\mathcal{A}(v_h, v_h) \gtrsim \|v_h\|_{DG} \quad \forall v_h \in V_h^p$$

For SIP and IIP, the penalty parameter  $\alpha$  should be large enough.

- Strong-consistency (Galerkin orthogonality):

$$\mathcal{A}(u, v_h) = \int_{\Omega} f v_h \quad \forall v_h \in V_h^p \Rightarrow \mathcal{A}(u - u_h, v_h) = 0 \quad \forall v_h \in V_h^p$$

- Approximation. Let  $\Pi_h^p u \in V_h^p$  be a suitable approximation of  $u$ , then

$$\|u - \Pi_h^p u\|_{DG} \lesssim \frac{h^{\min(p,s)}}{p^{s-\frac{1}{2}}} \|u\|_{H^{s+1}(\mathcal{T}_h)}$$

If  $p \geq s$

$$\|u - \Pi_h^p u\|_{DG} \lesssim \left(\frac{h}{p}\right)^s p^{\frac{1}{2}} \|u\|_{H^{s+1}(\mathcal{T}_h)}$$

### 3.4 Error estimates

Recall the abstract error estimate  $\|u - u_h\|_{DG} \lesssim \|u - \Pi_h^p\|_{DG}$ .

If  $u$  is sufficiently regular then

$$\|u - u_h\|_{DG} \lesssim \frac{h^{\min(p,s)}}{p^{s-\frac{1}{2}}} \|u\|_{H^{s+1}(\mathcal{T}_h)}$$

Then, by using a duality argument, one can obtain an estimate for the  $L^2$  norm.

Assuming that  $\Omega$  is such that the following elliptic regularity result holds: for any  $g \in L^2(\Omega)$ , the solution  $z$  of the problem

$$\begin{cases} -\Delta z = g & \text{in } \Omega \\ z = 0 & \text{on } \partial\Omega \end{cases}$$

satisfies  $z \in H^2(\Omega)$  and

$$\|z\|_{H^2(\Omega)} \lesssim \|g\|_{L^2(\Omega)}$$

If the exact solution  $u \in H^s(\Omega)$ ,  $s \geq 2$  and, if  $u_h$  is obtained with the SIP method, it holds

$$\|u - u_h\|_{L^2(\Omega)} \lesssim \frac{h^{\min(p,s)+1}}{p^{s+\frac{1}{2}}} \|u\|_{H^{s+1}(\Omega)}$$

while for NIP and IIP holds

$$\|u - u_h\|_{L^2(\Omega)} \lesssim \frac{h^{\min(p,s)}}{p^{s-\frac{1}{2}}} \|u\|_{H^{s+1}(\Omega)}$$

## 4 Advection-Diffusion-Reaction equations

### 4.1 Formulation of the problem

Considering the problem  $\mathcal{L}u = f$  in  $\Omega$ ,  $u = 0$  on  $\partial\Omega$  where

$$\begin{aligned}\mathcal{L}u &= -\operatorname{div}(\mu\nabla u + \mathbf{b}u) + \sigma u && \text{(conservative form)} \\ \mathcal{L}u &= -\operatorname{div}(\mu\nabla u) + \mathbf{b} \cdot \nabla u + \sigma u && \text{(non-conservative form)}\end{aligned}$$

with the same assumptions as (1.1).

The weak formulation is written as

$$\text{find } u \in V = H_0^1(\Omega) : a(u, v) = F(v) \quad \forall v \in V \quad (4.1)$$

with

$$F(v) = \int_{\Omega} f v$$

and

$$a(u, v) = \begin{cases} \int_{\Omega} (\mu\nabla u + \mathbf{b}u) \cdot \nabla v + \int_{\Omega} \sigma uv & \text{conservative form} \\ \int_{\Omega} \mu\nabla u \cdot \nabla v + \int_{\Omega} \mathbf{b} \cdot \nabla uv + \int_{\Omega} \sigma uv & \text{non-conservative form} \end{cases}$$

Let's verify the uniqueness of the solution:

#### Coercivity

Sufficient conditions for coercivity:

$$\begin{aligned}\sigma - \frac{1}{2}\operatorname{div}\mathbf{b} &\geq 0 \text{ in } \Omega && \text{non-conservative case} \\ \sigma + \frac{1}{2}\operatorname{div}\mathbf{b} &\geq 0 \text{ in } \Omega && \text{conservative case}\end{aligned}$$

In both cases:  $a(u, v) \geq \mu_0 \|\nabla u\|^2 \rightarrow$  coercivity constant  $\alpha \simeq \mu_0$

#### Continuity

In both cases, continuity constant:  $M \simeq \|\mu\|_{L^\infty} + \|\mathbf{b}\|_{L^\infty} + \|\sigma\|_{L^2}$

Given that the hypotheses of Lax-Milgram holds, the solution exists and is unique. We can now bring in the Galerkin formulation

$$\text{find } u_h \in V_h : a(u_h, v_h) = (f, v_h) \quad \forall v_h \in V_h$$

and move to the error estimate

$$\|u - u_h\| \underset{(\text{Ceà})}{\leq} \frac{M}{\alpha} \inf_{v_h \in V_h} \|u - v_h\| \underset{\substack{\text{(interpolation} \\ \text{error estimate)}}}{\leq} C \frac{M}{\alpha} h^r |u|_{H^{r+1}(\Omega)}$$

If it is a convection dominated flow (or reaction dominated), then  $\frac{M}{\alpha} \gg 1$ , then we need to find a tradeoff between  $\frac{M}{\alpha}$  and  $h^r$ . Also it is numerically prohibitive.

The Péclet number tells us if the flow is dominated by advection or diffusion if its greater or smaller than 1. We can define it as

$$\mathbb{P}e = h \frac{M}{\alpha}$$

Should be less than 1 for stability issues.

## 4.2 Stabilization methods

The idea now is to stabilize the Galerkin method.

- 1D case: Upwind method  $\iff$  Artificial diffusion
- 2D case: Streamline diffusion:

$$+c(h) \int_{\Omega} \frac{1}{\|\mathbf{b}\|} (\mathbf{b} \cdot \nabla u_h) (\mathbf{b} \cdot \nabla v_h)$$

Artificial diffusion:

$$+c(h) \int_{\Omega} \nabla u_h \cdot \nabla v_h$$

Now the solution is stabilized, but is not fully consistent. So the solution is to find a way to obtain a fully consistent solution

$$\text{find } u_h \in V_h : a(u_h, v_h) + \mathcal{L}_h(u_h, f; v_h) = F(v_h) \quad \forall v_h \in V_h$$

with  $\mathcal{L}_h$  suitably chosen such that

$$\mathcal{L}(u_h, f; v_h) = 0 \quad \forall v_h \in V_h$$

so we obtain a strongly consistent approximation of the original problem.

One possibility could be to use an operator proportional to the residual:

$$\mathcal{L}_h(u_h, f; v_h) = \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} (\mathcal{L}u - f) \tau_{\mathcal{K}} \varphi(v_h) \quad \forall v_h \in V_h$$

with  $\tau_{\mathcal{K}}$  as a scaling factor. Typically is chosen, given  $h_{\mathcal{K}} = \text{diam}(\mathcal{K})$ :

$$\tau_{\mathcal{K}}(\mathbf{x}) = \delta \frac{h_{\mathcal{K}}}{|\mathbf{b}(\mathbf{x})|} \quad \forall \mathbf{x} \in \mathcal{K}, \mathcal{K} \in \mathcal{T}_h$$

while, for  $\varphi(v_h)$  there are many possibilities. Two of them are

- $\varphi(v_h) = \mathcal{L}v_h \rightarrow$  GLS - Galerkin Least Squares method
- $\varphi(v_h) = \mathcal{L}_{ss}v_h \rightarrow$  SUPG - Streamline Upwind Petrov-Galerkin method

Brief notation remark:  $\mathcal{L} = \mathcal{L}_s + \mathcal{L}_{ss}$  (symmetric + skew-symmetric part) Which we define as

$$\begin{aligned} {}_{V'}\langle \mathcal{L}_s u, v \rangle_V &= {}_V\langle v, \mathcal{L}_s u \rangle_{V'} \quad \forall u, v \in V \\ {}_{V'}\langle \mathcal{L}_{ss} u, v \rangle_V &= -{}_V\langle v, \mathcal{L}_{ss} u \rangle_{V'} \quad \forall u, v \in V \end{aligned}$$

For matrices it is  $A = A_S + A_{SS}$

$$A_S = \frac{1}{2}(A + A^T) \quad A_{SS} = \frac{1}{2}(A - A^T)$$

Let us see an example in the non conservative form

$$\begin{aligned} \mathcal{L}^1 &= -\mu \Delta u + \mathbf{b} \cdot \nabla u + \sigma u \\ &= \underbrace{\left[ -\mu \Delta u + \left( \sigma - \frac{1}{2} \text{div} \mathbf{b} \right) u \right]}_{\mathcal{L}_s^1 u} + \underbrace{\left[ \frac{1}{2} (\text{div}(\mathbf{b}u) + \mathbf{b} \cdot \nabla u) \right]}_{\mathcal{L}_{ss}^1 u} \end{aligned}$$

Indeed we can see

$$\begin{aligned} {}_{V'}\langle \mathcal{L}_s^1, v \rangle_V &= \int_{\Omega} \mu \nabla u \cdot \nabla v + \left( \sigma - \frac{1}{2} \text{div} \mathbf{b} \right) uv \\ &= \int_{\Omega} \left[ -\mu \Delta v + \left( \sigma - \frac{1}{2} \text{div} \mathbf{b} \right) v \right] u = {}_V\langle v, \mathcal{L}_s^1 \rangle_{V'} \end{aligned}$$

$$\begin{aligned}
{}_{V'}\langle \mathcal{L}_{ss}^1, v \rangle_V &= \frac{1}{2} \int_{\Omega} (\operatorname{div}(\mathbf{b}u)v + (\mathbf{b} \cdot \nabla u)v) \\
&= \frac{1}{2} \int_{\Omega} (-(\mathbf{b}u)\nabla v + (\mathbf{b}v) \cdot \nabla u) \\
&= \frac{1}{2} \int_{\Omega} (-(\mathbf{b} \cdot \nabla v)u - \operatorname{div}(\mathbf{b}v)u) = -_V\langle u, \mathcal{L}_{ss}^1 \rangle_{V'}
\end{aligned}$$

**Remark 4.1**

If  $\operatorname{div} \mathbf{b} = 0$ , which happens if  $\mathbf{b}$  is constant, then the conservative and non conservative forms coincide.

### 4.3 GLS method (conservative form)

$$\text{find } u_h \in V_h : a(u_h, v_h) + \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\Omega} \mathcal{L}u_h \tau_{\mathcal{K}} \mathcal{L}v_h = \int_{\Omega} f v_h + \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\Omega} f \tau_{\mathcal{K}} \mathcal{L}v_h \quad \forall v_h \in V_h$$

**Theorem 4.1**

Consider the conservative case. Suppose that

$$\exists \gamma_0, \gamma_1 > 0 : 0 < \gamma_0 \leq \gamma(\mathbf{x}) \leq \gamma_1$$

then, for a suitable constant  $C$ , independent of  $h$ , we have:

$$\|u_h\|_{GLS}^2 \leq C \|f\|_{L^2(\Omega)}^2$$

where  $\|\cdot\|_{GLS}$  will be defined later

**Proof.** Take  $u_h = v_h$ . We have

$$\begin{aligned}
a_h(u_h, u_h) &= \int_{\Omega} \mu |\nabla u_h|^2 + \underbrace{\int_{\Omega} \operatorname{div}(\mathbf{b} u_h) u_h}_{\substack{= - \int_{\Omega} \mathbf{b} \cdot (u_h \nabla u_h) \\ = - \frac{1}{2} \int_{\Omega} \mathbf{b} \cdot \nabla (u_h^2) \\ = \frac{1}{2} \int_{\Omega} \operatorname{div} \mathbf{b} u_h^2}} + \int_{\Omega} \sigma u_h^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \tau_{\mathcal{K}} (\mathcal{L}u_h)^2 \\
&= \int_{\Omega} \mu |\nabla u_h|^2 + \underbrace{\int_{\Omega} \left( \sigma + \frac{1}{2} \operatorname{div} \mathbf{b} \right) u_h^2}_{=: \gamma(\mathbf{x})} + \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \tau_{\mathcal{K}} (\mathcal{L}u_h)^2 \\
&=: \|u_h\|_{GLS}^2
\end{aligned}$$

On the other hand

$$|F_h(u_h)| \leq \left| \int_{\Omega} f u_h \right| + \left| \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} f \tau_{\mathcal{K}} \mathcal{L}u_h \right|$$

where

$$\begin{aligned}
\left| \int_{\Omega} f u_h \right| &= \left| \int_{\Omega} \frac{1}{\sqrt{\gamma}} f \sqrt{\gamma} u_h \right| \stackrel{\text{Cauchy-Schwartz}}{\leq} \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)} \|\sqrt{\gamma} u_h\|_{L^2(\Omega)} \\
&\stackrel{\text{Young}}{\leq} \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)}^2 + \frac{1}{4} \|\sqrt{\gamma} u_h\|_{L^2(\Omega)}^2
\end{aligned}$$

and where

$$\begin{aligned}
\left| \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} f \tau_{\mathcal{K}} \mathcal{L}u_h \right| &= \left| \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \sqrt{\tau_{\mathcal{K}}} f \sqrt{\tau_{\mathcal{K}}} \mathcal{L}u_h \right| \\
&\stackrel{\text{Cauchy-Schwartz}}{\leq} \sum_{\mathcal{K} \in \mathcal{T}_h} \|\sqrt{\tau_{\mathcal{K}}} f\|_{L^2(\mathcal{K})} \|\sqrt{\tau_{\mathcal{K}}} \mathcal{L}u_h\|_{L^2(\mathcal{K})} \\
&\stackrel{\text{Young}}{\leq} \sum_{\mathcal{K} \in \mathcal{T}_h} \|\sqrt{\tau_{\mathcal{K}}} f\|_{L^2(\mathcal{K})}^2 + \frac{1}{4} \|\sqrt{\tau_{\mathcal{K}}} \mathcal{L}u_h\|_{L^2(\mathcal{K})}^2
\end{aligned}$$



So,  $a_h(u_h, u_h) = F_h(u_h)$  implies:

$$\begin{aligned}
\|u_h\|_{GLS}^2 &= \int_{\Omega} \mu |\nabla u_h|^2 + \int_{\Omega} \gamma u_h^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \int_{\mathcal{K}} \tau_{\mathcal{K}} (\mathcal{L}u_h)^2 \\
&\leq \left[ \left\| \frac{1}{\sqrt{\gamma}} f \right\|_{L^2(\Omega)} + \sum_{\mathcal{K} \in \mathcal{T}_h} \|\sqrt{\tau_{\mathcal{K}}} f\|_{L^2(\Omega)}^2 \right] \\
&\quad + \frac{1}{4} \left[ \int_{\Omega} \gamma u_h^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \tau_{\mathcal{K}} (\mathcal{L}u_h)^2 \right] \\
&\leq \underbrace{\left( \frac{1}{\gamma_0} + \max_{\mathcal{K} \in \mathcal{T}_h} \tau_{\mathcal{K}} \right)}_{=C(\text{if } \tau_{\mathcal{K}} \text{ uniformly bounded w.r.t. } h)} \|f\|_{L^2(\Omega)}^2 + \frac{1}{4} \|u_h\|_{GLS}^2
\end{aligned}$$

In the end

$$\|u_h\|_{GLS}^2 \leq \frac{4}{3} C \|f\|_{L^2(\Omega)}^2$$

★

As we already said, a smart choice for  $\tau_{\mathcal{K}}$  is  $\delta \frac{h_{\mathcal{K}}}{|\mathbf{b}(\mathbf{x})|}$ . But another possibility may be

$$\tau_{\mathcal{K}}(\mathbf{x}) = \frac{h_{\mathcal{K}}}{2|\mathbf{b}(\mathbf{x})|} \xi(\mathbb{P}e_{\mathcal{K}})$$

with  $\xi(\theta) = \coth(\theta) - \frac{1}{\theta}$ . and  $\mathbb{P}e_{\mathcal{K}}(\mathbf{x}) = \frac{|\mathbf{b}(\mathbf{x})|}{2\mu(\mathbf{x})} h_{\mathcal{K}}$  is the local Péclet number. Moreover, if  $\theta \rightarrow 0$ , then  $\xi(\theta) = \frac{\theta}{3} + o(\theta)$ , therefore when  $\mathbb{P}e_{\mathcal{K}}(\mathbf{x}) \ll 1$ , we have  $\tau_{\mathcal{K}}(\mathbf{x}) \rightarrow 0$  and no stabilization is needed.

#### 4.4 Convergence of GLS

To state the convergence of GLS we need the inverse inequality, defined as

$$\sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \int_{\mathcal{K}} |\Delta v_h|^2 d\mathcal{K} \leq C_0 \|\nabla v_h\|_{L^2(\Omega)}^2 \quad \forall v_h \in X_h^r \quad (4.2)$$

##### **Theorem 4.2** (Convergence of GLS)

Assume that the space  $V_h$  satisfies the following local approximation property: for each  $v \in V \cap H^{r+1}(\Omega)$ , there exists a function  $\hat{v}_h \in V_h$  s.t.

$$\|v - v_h\|_{L^2(\mathcal{K})} + h_{\mathcal{K}} \|v - \hat{v}_h\|_{H^1(\mathcal{K})} + h_{\mathcal{K}}^2 \|v - \hat{v}_h\|_{H^2(\mathcal{K})} \leq C h_{\mathcal{K}}^{r+1} |v|_{H^{r+1}} \quad (4.3)$$

for each  $\mathcal{K} \in \mathcal{T}_h$ . Moreover, we suppose that for each  $\mathcal{K} \in \mathcal{T}_h$  the local Péclet number of  $K$  satisfies

$$\mathbb{P}e_{\mathcal{K}}(\mathbf{x}) = \frac{|\mathbf{b}(\mathbf{x})| h_{\mathcal{K}}}{2\mu} > 1 \quad \forall \mathbf{x} \in \mathcal{K} \quad (4.4)$$

that is, we are in the pre-asymptotic regime. Finally, we suppose that the inverse inequality holds and that the stabilization parameters satisfies the relation  $0 < \delta \leq 2C_0^{-1}$ .

Then, as long as  $u \in H^{r+1}(\Omega)$ , the following super-optimal estimate holds:

$$\|u - u_h\|_{GLS} \leq C h^{r+\frac{1}{2}} |u|_{H^{r+1}(\Omega)} \quad (4.5)$$

**Proof.** First of all, rewrite the error as

$$e_h = u_h - u = \sigma_h - \eta \quad (4.6)$$

with  $\sigma_h = u_h - \hat{u}_h$ ,  $\eta = u - \hat{u}_h$ , where  $\hat{u}_h$  is a function that depends on  $u$  and that satisfies property (4.3). If, for instance,  $V_h = X_h^r \cap H_0^1(\Omega)$ , we can choose  $\hat{u}_h = \prod_h^r u$  that is the finite element interpolant of  $u$ .

We start by estimating the norm  $\|\sigma_h\|_{GLS}$ . By exploiting the strong consistency of the GLS scheme we obtain

$$\|\sigma_h\|_{GLS}^2 = a_h(\sigma_h, \sigma_h) = a_h(u_h - u + \eta, \sigma_h) = a_h(\eta, \sigma_h)$$

Now, thanks to the homogeneous Dirichlet boundary conditions it follows that, by adding and subtracting  $\sum_{\mathcal{K} \in \mathcal{T}_h} (\eta, \mathcal{L}\sigma_h)_{\mathcal{K}}$ , suitable computations lead to:

$$\begin{aligned} a_h(\eta, \sigma_h) &= \mu \sigma_\Omega \nabla \eta \cdot \nabla \sigma_h \, d\Omega - \int_\Omega \eta \mathbf{b} \cdot \nabla \sigma_h \, d\Omega + \int_\Omega \sigma \eta \sigma_h \, d\Omega \\ &\quad + \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \mathcal{L}\eta, \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h \right)_{L^2(\mathcal{K})} \\ &= \underbrace{\mu (\nabla \eta, \nabla \sigma_h)_{L^2(\Omega)}}_{(I)} - \underbrace{\sum_{\mathcal{K} \in \mathcal{T}_h} (\eta, \mathcal{L}\sigma_h)_{L^2(\Omega)}}_{(II)} + \underbrace{2 (\gamma \eta, \sigma_h)_{L^2(\mathcal{K})}}_{(III)} \\ &\quad + \underbrace{\sum_{\mathcal{K} \in \mathcal{T}_h} (\eta, -\mu \Delta \sigma_h)_{L^2(\mathcal{K})}}_{(IV)} + \underbrace{\sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \mathcal{L}\eta, \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h \right)_{L^2(\mathcal{K})}}_{(V)} \end{aligned}$$

Now, we bound each of these terms. By using Cauchy-Schwartz and Young's inequalities we obtain

$$\begin{aligned} |(I)| &= \left| \mu (\nabla \eta, \nabla \sigma_h)_{L^2(\mathcal{K})} \right| \leq \frac{\mu}{4} \|\nabla \sigma_h\|_{L^2(\Omega)}^2 + \mu \|\nabla \eta\|_{L^2(\Omega)}^2 \\ |(II)| &= \left| \sum_e it (\eta, \mathcal{L}\sigma_h)_{L^2(\mathcal{K})} \right| \\ &= \left| \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \sqrt{\frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}}} \eta, \sqrt{\frac{\delta h_{\mathcal{K}}}{|\mathbf{b}|}} \mathcal{L}\sigma_h \right)_{L^2(\Omega)} \right| \\ &\leq \frac{1}{4} \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h, \mathcal{L}\sigma_h \right) \\ |(III)| &= 2 \left| (\gamma \eta, \sigma_h)_{L^2(\Omega)} \right| = 2 \left| \left( \sqrt{\gamma \eta}, \sqrt{\gamma} \sigma_h \right)_{L^2(\Omega)} \right| \\ &\leq \frac{1}{2} \|\sqrt{\gamma} \sigma_h\|_{L^2(\Omega)}^2 + 2 \|\sqrt{\gamma} \eta\|_{L^2(\Omega)}^2 \end{aligned}$$

Then, thanks to CS and Young, but also hypotheses (4.4) and (4.2), we obtain

$$\begin{aligned} |(IV)| &= \left| \sum_{\mathcal{K} \in \mathcal{T}_h} (\eta, -\mu \Delta \sigma_h)_{L^2(\mathcal{K})} \right| \\ &\leq \frac{1}{4} \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \mu^2 \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \Delta \sigma_h, \Delta \sigma_h \right)_{L^2(\mathcal{K})} + \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}} \eta, \eta \right)_{L^2(\mathcal{K})} \\ &\leq \frac{1}{8} \delta \mu \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 (\nabla \sigma_h, \nabla \sigma_h)_{L^2(\mathcal{K})} + \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}} \eta, \eta \right)_{L^2(\mathcal{K})} \\ &\leq \frac{\sigma C_0 \mu}{8} \|\nabla \sigma_h\|_{L^2(\Omega)}^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}} \eta, \eta \right)_{L^2(\mathcal{K})} \end{aligned}$$

The last one can be bounded once again thanks to CS and Young inequalities as follows

$$\begin{aligned} |(V)| &= \left| \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \mathcal{L}\eta, \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h \right)_{L^2(\mathcal{K})} \right| \\ &\leq \frac{1}{4} \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h, \mathcal{L}\sigma_h \right)_{L^2(\mathcal{K})} + \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\eta, \mathcal{L}\eta \right)_{L^2(\mathcal{K})} \end{aligned}$$

So we can rewrite everything bounded as

$$\begin{aligned} \|\sigma_h\|_{GLS}^2 &= a_h(\eta, \sigma_h) \leq \frac{1}{4} \|\sigma_h\|_{GLS}^2 \\ &+ \frac{1}{4} \left( \|\sqrt{\gamma}\sigma_h\|_{L^2(\Omega)}^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\sigma_h, \mathcal{L}\sigma_h \right)_{L^2(\mathcal{K})} \right) + \frac{\delta C_0 \mu}{8} \|\nabla \sigma_h\|_{L^2(\Omega)}^2 \\ &+ \underbrace{\mu \|\nabla \eta\|_{L^2(\Omega)}^2 + 2 \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}} \eta, \eta \right)_{L^2(\mathcal{K})} + 2 \|\sqrt{\gamma}\eta\|_{L^2(\Omega)}^2 + \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\eta, \mathcal{L}\eta \right)_{L^2(\mathcal{K})}}_{\mathcal{E}(\eta)} \\ &\leq \frac{1}{2} \|\sigma_h\|_{GLS}^2 + \mathcal{E}(\eta) \end{aligned}$$

Having exploited the assumption that  $\delta \leq 2C_0^{-1}$ . We can state then

$$\|\sigma_h\|_{GLS}^2 \leq 2\mathcal{E}(\eta)$$

It's time to estimate  $\mathcal{E}(\eta)$ , by bounding each of its summands separately. To do this, we will use the local approximation property (4.3) and the local Péclet (4.4). Moreover, we observe that the constant  $C$ , introduced in the remainder, depends neither on  $h$  nor on  $\mathbb{P}e_{\mathcal{K}}$ , but can depend on other quantities such as the constant  $\gamma_1$ , the reaction constant  $\sigma$  or the norm  $\|\mathbf{b}\|_{L^\infty(\Omega)}$ , the stabilization parameter  $\delta$ . Then we have

$$\begin{aligned} \mu \|\nabla \eta\|_{L^2(\Omega)}^2 &\leq C \mu h^{2r} |u|_{H^{r+1}(\Omega)}^2 \\ &\leq C \frac{\|\mathbf{b}\|_{L^\infty(\Omega)} h}{2} h^{2r} |u|_{H^{r+1}(\Omega)}^2 \leq C h^{2r+1} |u|_{H^{r+1}(\Omega)}^2 \\ 2 \sum_{\mathcal{K} \in \mathcal{T}_h} \left( \frac{|\mathbf{b}|}{\delta h_{\mathcal{K}}} \eta, \eta \right)_{L^2(\mathcal{K})} &\leq C \frac{\|\mathbf{b}\|_{L^\infty(\Omega)} h}{2} \sum_{\mathcal{K} \in \mathcal{T}_h} \frac{1}{h_{\mathcal{K}}} h_{\mathcal{K}}^{2r+1} |u|_{H^{r+1}(\Omega)}^2 \\ &\leq c h^{2r+1} |u|_{H^{r+1}(\Omega)}^2 \\ 2 \|\sqrt{\gamma}\eta\|_{L^2(\Omega)}^2 &\leq 2\gamma_1 \|\eta\|_{L^2(\Omega)}^2 \leq C h^{2r+1} |u|_{H^{r+1}(\Omega)}^2 \end{aligned}$$

For the fourth term we have

$$\begin{aligned} \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \frac{h_{\mathcal{K}}}{|\mathbf{b}|} \mathcal{L}\eta, \mathcal{L}\eta \right)_{L^2(\mathcal{K})} &= \sum_{\mathcal{K} \in \mathcal{T}_h} \left\| \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \mathcal{L}\eta \right\|_{L^2(\mathcal{K})}^2 \\ &= \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left\| -\mu \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \Delta \eta + \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \operatorname{div}(\mathbf{b}\eta) + \sigma \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \eta \right\|_{L^2(\mathcal{K})}^2 \\ &\leq C \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left( \left\| \mu \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \Delta \eta \right\|_{L^2(\mathcal{K})}^2 + \left\| \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \operatorname{div}(\mathbf{b}\eta) \right\|_{L^2(\mathcal{K})}^2 \right. \\ &\quad \left. + \left\| \sigma \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \eta \right\|_{L^2(\mathcal{K})}^2 \right) \end{aligned} \tag{4.7}$$

Now it is easy to prove that the second and third term of the summands can be bounded using a term or the form  $Ch^{2r+1}|u|_{H^{r+1}(\Omega)}^2$ , for a suitable choice of the constant  $C$ . For the first term we have

$$\begin{aligned} \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \left\| \mu \sqrt{\frac{h_{\mathcal{K}}}{|\mathbf{b}|}} \Delta \eta \right\|_{L^2(\mathcal{K})}^2 &\leq \sum_{\mathcal{K} \in \mathcal{T}_h} \delta \frac{h_{\mathcal{K}}^2 \mu}{2} \|\Delta \eta\|_{L^2(\mathcal{K})}^2 \\ &\leq C \delta \|\mathbf{b}\|_{L^\infty(\Omega)} \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^3 \|\Delta\|_{L^2(\mathcal{K})}^2 \leq |u|_{H^{r+1}(\Omega)}^2 \end{aligned}$$

having used again (4.3) and (4.4). Now we can conclude that

$$\mathcal{E}(\eta) \leq Ch^{2r+1}|u|_{H^{r+1}(\Omega)}^2$$

that is

$$\|\sigma_h\|_{GLS} \leq Ch^{r+\frac{1}{2}}|u|_{H^{r+1}(\Omega)} \quad (4.8)$$

Reverting to (4.6), to obtain the desired estimate for the norm  $\|u_h - u\|_{GLS}$  we need to estimate  $\|\eta\|_{GLS}$ . But thanks to (4.7) we obtain

$$\|\eta\|_{GLS} \leq Ch^{r+\frac{1}{2}}|u|_{H^{r+1}(\Omega)}$$

Combining this with (4.8) we obtain (4.5). ★

## 5 Parabolic equations

### 5.1 Introduction

Now we consider parabolic equations of the form

$$\frac{\partial u}{\partial t} + \mathcal{L}u = f \quad \mathbf{x} \in \Omega, t > 0 \quad (5.1)$$

where:

- $\Omega$  is a domain of  $\mathbb{R}^d$  with  $d = 1, 2, 3$
- $f = f(\mathbf{x}, t)$  is a given function
- $\mathcal{L} = \mathcal{L}(\mathbf{x})$  is a generic elliptic operator acting on  $u = u(\mathbf{x}, t)$

When solved for a bounded time interval, for example  $0 < t < T$ , the region  $Q_T = \Omega \times (0, T)$  is called cylinder in the space  $\mathbb{R}^d \times \mathbb{R}^+$ . In (5.1) must be assigned an initial condition

$$u(\mathbf{x}, 0) = u_0(\mathbf{x}), \quad \mathbf{x} \in \Omega \quad (5.2)$$

also we'll need some BC, like

$$\begin{aligned} u(\mathbf{x}, t) &= \varphi(\mathbf{x}, t) \quad \mathbf{x} \in \Gamma_D \text{ and } t > 0 \\ \frac{\partial u(\mathbf{x}, t)}{\partial n} &= \psi(\mathbf{x}, t) \quad \mathbf{x} \in \Gamma_N \text{ and } t > 0 \end{aligned} \quad (5.3)$$

where  $u_0, \varphi$  and  $\psi$  are given function and  $\{\Gamma_D, \Gamma_N\}$  provides a boundary partition that is  $\Gamma_D \cup \Gamma_N = \partial\Omega, \Gamma_D^\circ \cap \Gamma_N^\circ = \emptyset$ . For obvious reasons  $\Gamma_D$  is the Dirichlet boundary, while  $\Gamma_N$  is the Neumann one. In the one dimensional case the problem becomes

$$\begin{aligned} \frac{\partial u}{\partial t} - \nu \frac{\partial^2 u}{\partial x^2} &= f & 0 < x < d, t > 0 \\ u(x, 0) &= u_0(x) & 0 < x < d \\ u(0, t) = u(d, t) &= 0 & t > 0 \end{aligned} \quad (5.4)$$

which describes the evolution of the temperature  $u(x, t)$  at point  $x$  and time  $t$  of a metal bar of length  $d$  occupying the interval  $[0, d]$ , whose thermal conductivity is  $\nu$  and whose endpoints are kept at a constant temperature of zero degrees. The function  $u_0$  describes the temperature in the initial state, while  $f$  represents the heat generated per unit of length by the bar. This is called the heat equation.

### 5.2 Weak formulation and approximation

We proceed as usual by multiplying for each  $t > 0$  the differential equation by a test function  $v = v(\mathbf{x})$  and integrating. We set  $V = H_{\Gamma_D}^1(\Omega)$  and, for each  $t > 0$ , we seek  $u(t) \in V$  s.t.

$$\int_{\Omega} \frac{\partial u(t)}{\partial t} v \, d\Omega + a(u(t), v) = \int_{\Omega} f(t) v \, d\Omega \quad \forall v \in V \quad (5.5)$$

where

- $u(0) = u_0$
- $a(\cdot, \cdot)$  is the bilinear form associated to the operator  $\mathcal{L}$
- we have supposed for simplicity  $\varphi = 0$  and  $\psi = 0$

**Definition 5.1**

A bilinear form  $a(\cdot, \cdot)$  is said weakly coercive if

$$\exists \lambda \geq 0, \exists \alpha > 0 : a(v, v) + \lambda \|v\|_{L^2(\Omega)}^2 \geq \alpha \|v\|_V^2 \quad \forall v \in V$$

Rationale for weak coercivity

$$\frac{\partial u}{\partial t} + \mathcal{L}u = f$$

now we perform a change of variable ( $u = e^{\lambda t}$ )

$$\frac{\partial w}{\partial t} + \mathcal{L}w + \lambda w = e^{-\lambda t} f$$

so that the bilinear form

$$\tilde{a}(w, v) := a(w, v) + \lambda(w, v)_{L^2(\Omega)} \Rightarrow \tilde{a}(w, w) := a(w, w) + \lambda \|w\|_{L^2(\Omega)}^2$$

**Theorem 5.1**

Suppose that the bilinear form  $a(\cdot, \cdot)$  is continuous and weakly coercive. Moreover, we require  $u_0 \in L^2(\Omega)$  and  $f \in L^2(Q_T)$ . Then, (5.5) admits a unique solution  $u \in \mathcal{C}^0(\mathbb{R}^+; L^2(\Omega))$ . Also,  $u \in L^2(\mathbb{R}^+; V)$  and  $\frac{\partial u}{\partial t} \in L^2(\mathbb{R}^+; V')$ .

**5.3 Algebraic formulation**

Now we can use Galerkin to approximate, for each  $t > 0$ , we need to find  $u_h(t) \in V_h$  s.t.

$$\int_{\Omega} \frac{\partial u_h(t)}{\partial t} v_h d\Omega + a(u_h(t), v_h) = \int_{\Omega} f(t) v_h d\Omega \quad \forall v_h \in V_h \quad (5.6)$$

with  $u_h(0) = u_{0h}$ , where  $V_h \subset V$  is a suitable space of finite dimension and  $u_{0h}$  is a convenient approximation of  $u_0$  in the space  $V_h$ .

Now we need to discretize the temporal variable, because, as of now, we obtained a semi-discretization of the problem.

We introduce a basis  $\{\varphi_j\}$  for  $V_h$  and we observe that it suffices that (5.6) is verified for the basis function in order to be satisfied by all the functions in the subspace.

Moreover, since for each  $t > 0$ , the solution to the Galerkin problem belongs to the subspace as well, we will have

$$u_h(\mathbf{x}, t) = \sum_{j=1}^{N_h} u_j(t) \varphi_j(\mathbf{x})$$

where the coefficients  $\{u_j(t)\}$  represent the unknown of the problem.

Denoting by  $\dot{u}_j(t)$  the temporal derivatives of  $u_j(t)$ , (5.6) becomes

$$\int_{\Omega} \sum_{j=1}^{N_h} \dot{u}_j(t) \varphi_j \varphi_i d\Omega + a \left( \sum_{j=1}^{N_h} u_j(t) \varphi_j, \varphi_i \right) = \int_{\Omega} f(t) \varphi_i d\Omega, \quad i = 1, 2, \dots, N_h$$

that is

$$\sum_{j=1}^{N_h} \dot{u}_j(t) \underbrace{\int_{\Omega} \varphi_j \varphi_i d\Omega}_{m_{ij}} + \sum_{j=1}^{N_h} u_j(t) \underbrace{a(\varphi_j, \varphi_i)}_{a_{ij}} = \underbrace{\int_{\Omega} f(t) \varphi_i d\Omega}_{f_i(t)} \quad i = 1, 2, \dots, N_h \quad (5.7)$$

If we define the vector of unknowns  $\mathbf{u} = (u_1(t), u_2(t), \dots, u_{N_h}(t))^T$ , the mass matrix  $M = [m_{ij}]$ , the stiffness matrix  $A = [a_{ij}]$  and the right hand side vector  $\mathbf{f} = (f_1(t), f_2(t), \dots, f_{N_h}(t))^T$ , the system (5.7) can be rewritten as

$$M \dot{\mathbf{u}}(t) + A \mathbf{u}(t) = \mathbf{f}(t)$$

## 5.4 Time discretization

For the numerical solution of this ODE system we will use the  $\theta$ -method, which discretizes the temporal difference quotient and replaces the other terms with a linear combination of the value at time  $t^k$  and of the value at time  $t^{k+1}$ , depending on the real parameter  $0 \leq \theta \leq 1$ .

$$M \frac{\mathbf{u}^{k+1} - \mathbf{u}^k}{\Delta t} + A[\theta \mathbf{u}^{k+1} + (1 - \theta) \mathbf{u}^k] = \theta \mathbf{f}^{k+1} + (1 - \theta) \mathbf{f}^k \quad (5.8)$$

The real positive parameter  $\Delta t = t^{k+1} - t^k$ ,  $k = 0, 1, \dots$  denotes the discretization step. Some particular cases of (5.7)

- $\theta = 0$  The forward Euler method

$$M \frac{\mathbf{u}^{k+1} - \mathbf{u}^k}{\Delta t} + A \mathbf{u}^k = \mathbf{f}^k$$

- $\theta = 1$  The backward Euler method

$$M \frac{\mathbf{u}^{k+1} - \mathbf{u}^k}{\Delta t} + A \mathbf{u}^{k+1} = \mathbf{f}^{k+1}$$

- $\theta = \frac{1}{2}$  The Crank-Nicolson (or trapezoidal) method

$$M \frac{\mathbf{u}^{k+1} - \mathbf{u}^k}{\Delta t} + \frac{1}{2} A (\mathbf{u}^{k+1} + \mathbf{u}^k) = \frac{1}{2} (\mathbf{f}^{k+1} + \mathbf{f}^k)$$

which is of the second order in  $\Delta t$ .

Let us consider the two extremal cases  $\theta = 0$  and  $\theta = 1$ . In the first case the system to solve is only the mass matrix  $\frac{M}{\Delta t}$ , while in the other case  $\frac{M}{\Delta t} + A$ .  $M$  is invertible, being positively defined. In the case  $\theta = 0$  the scheme is not unconditionally stable, and in the case where  $V_h$  is a subspace of finite elements, there is the following stability condition

$$\exists c > 0 : \Delta t \leq ch^2 \quad \forall h > 0$$

so  $\Delta t$  cannot be chosen irrespective of  $h$ .

In this case, if we make  $M$  diagonal, we actually decouple the system. This operation is called lumping of the mass matrix.

When  $\theta > 0$ , the system will have the form  $K \mathbf{u}^{k+1} = \mathbf{g}$ , where  $\mathbf{g}$  is the source term and  $K = \frac{M}{\Delta t}$ . The latter is invariant in time (the operator  $\mathcal{L}$  being independent of time), so, if the spatial mesh doesn't change, it can be factorized only once at the beginning of the process.

Then, since  $M$  is symmetric, if  $A$  is symmetric, also  $K$  will be symmetric too. Hence, we can use, for example, the Cholesky factorization  $K = HH^T$ , with  $H$  being lower triangular. At each timestep, then, will be solved two triangular systems

$$\begin{aligned} H \mathbf{y} &= \mathbf{g} \\ H^T \mathbf{u}^{k+1} &= \mathbf{y} \end{aligned}$$

## 5.5 A priori estimate

Let us now consider (5.5). Since the corresponding equation must hold for each  $v \in V$ , it will be legitimate to set  $v = u(t)$  with  $t$  being given. yielding

$$\int_{\Omega} \frac{\partial u(t)}{\partial t} u(t) d\Omega = \frac{1}{2} \frac{\partial}{\partial t} \int_{\Omega} |u(t)|^2 d\Omega = \frac{1}{2} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}^2 \quad (5.9)$$

Considering the individual terms, the first one is

$$\int_{\Omega} \frac{\partial u(t)}{\partial t} u(t) d\Omega = \frac{1}{2} \frac{\partial}{\partial t} \int_{\Omega} |u(t)|^2 d\Omega = \frac{1}{2} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}^2 \quad (5.10)$$

Assuming that the bilinear form is coercive, with  $\alpha$  coercivity constant, we obtain

$$a(u(t), u(t)) \geq \alpha \|u(t)\|_V^2$$

while, thanks to the CS inequality we find:

$$(f(t), u(t)) \leq \|f(t)\|_{L^2(\Omega)} \|u(t)\|_{L^2(\Omega)} \quad (5.11)$$

So (5.9) becomes

$$\frac{1}{2} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}^2 + \alpha \|u(t)\|_V^2 \leq \|f(t)\|_{L^2(\Omega)} \|u(t)\|_{L^2(\Omega)}$$

Now let's define two important inequalities

**Definition 5.2** (Young's inequality)

$\forall a, b \in \mathbb{R}$

$$ab \leq \varepsilon \alpha^2 + \frac{1}{4\varepsilon} b^2 \quad \forall \varepsilon > 0 \quad (5.12)$$

**Definition 5.3** (Poincaré inequality)

If  $\Gamma_D$  is a set of positive measure, then:

$$\exists C_\Omega > 0 : \|v\|_{L^2(\Omega)} \leq C_\Omega \|\nabla v\|_{L^2(\Omega)} \quad \forall v \in H_{\Gamma_D}^1 \quad (5.13)$$

By using these two we obtain

$$\begin{aligned} \frac{1}{2} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}^2 + \alpha \|\nabla u\|_{L^2(\Omega)}^2 &\leq \|f(t)\|_{L^2(\Omega)} \|u(t)\|_{L^2(\Omega)} \\ &\leq \frac{C_\Omega^2}{2\alpha} \|f(t)\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|\nabla u(t)\|_{L^2(\Omega)}^2 \end{aligned} \quad (5.14)$$

Then, by integrating in time, we obtain, for all  $t > 0$

$$\|u(t)\|_{L^2(\Omega)}^2 + \alpha \int_0^t \|\nabla u(s)\|_{L^2(\Omega)}^2 ds \leq \|u_0\|_{L^2(\Omega)}^2 + \frac{C_\Omega^2}{\alpha} \int_0^t \|f(s)\|_{L^2(\Omega)}^2 ds \quad (5.15)$$

This is an a priori energy estimate. Note that

$$\frac{1}{2} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}^2 = \|u(t)\|_{L^2(\Omega)} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)}$$

then, from (5.9), using (5.10), (5.11) and (5.13)

$$\begin{aligned} \|u(t)\|_{L^2(\Omega)} \frac{\partial}{\partial t} \|u(t)\|_{L^2(\Omega)} + \frac{\alpha}{C_\Omega} \|u(t)\|_{L^2(\Omega)} \|\nabla u(t)\|_{L^2(\Omega)} \\ \leq \|f(t)\|_{L^2(\Omega)} \|u(t)\|_{L^2(\Omega)}, \quad t > 0 \end{aligned} \quad (5.16)$$

If  $\|u(t)\|_{L^2(\Omega)} \neq 0$ , we can divide by  $\|u(t)\|_{L^2(\Omega)}$  and integrate in time to obtain another estimate.

$$\|u(t)\|_{L^2(\Omega)} \leq \|u_0\|_{L^2(\Omega)} + \int_0^t \|f(s)\|_{L^2(\Omega)} ds \quad t > 0 \quad (5.17)$$

Let us now use the first inequality in (5.14) and integrate in time to yield

$$\begin{aligned} \|u(t)\|_{L^2(\Omega)}^2 + 2\alpha \int_0^t \|\nabla u(s)\|_{L^2(\Omega)}^2 ds &\stackrel{(5.14)}{\leq} \|u(t)\|_{L^2(\Omega)}^2 + 2 \int_0^t \|f(s)\|_{L^2(\Omega)} \|u(s)\|_{L^2(\Omega)} ds \\ &\stackrel{(5.17)}{\leq} \|u_0\|_{L^2(\Omega)}^2 + 2 \int_0^t \|f(s)\|_{L^2(\Omega)} \cdot \left( \|u_0\|_{L^2(\Omega)} \right. \\ &\quad \left. + \int_0^s \|f(\tau)\|_{L^2(\Omega)} d\tau \right) ds \\ &= \|u_0\|_{L^2(\Omega)}^2 \int_0^t \|f(s)\|_{L^2(\Omega)} \|u_0\|_{L^2(\Omega)} ds \\ &\quad + 2 \int_0^t \|f(s)\|_{L^2(\Omega)} \int_0^s \|f(\tau)\|_{L^2(\Omega)} d\tau ds \end{aligned} \quad (5.18)$$



Now, noticing that

$$\|f(s)\|_{L^2(\Omega)} \int_0^s \|f(\tau)\|_{L^2(\Omega)} d\tau = \frac{\partial}{\partial s} \left( \int_0^s \|f\|_{L^2(\Omega)} d\tau \right)^s$$

we can rewrite the right hand side of (5.18) as:

$$\begin{aligned} \|u_0\|_{L^2(\Omega)}^2 \int_0^t \|f(s)\|_{L^2(\Omega)} \|u_0\|_{L^2(\Omega)} ds + 2 \int_0^t \|f(s)\|_{L^2(\Omega)} \int_0^s \|f(\tau)\|_{L^2(\Omega)} d\tau ds \\ = \left( \|u_0\|_{L^2(\Omega)} + \int_0^t \|f(s)\|_{L^2(\Omega)} ds \right)^2 \end{aligned}$$

Therefore we can obtain the following estimate

$$\begin{aligned} \left( \|u(t)\|_{L^2(\Omega)}^2 + 2\alpha \int_0^t \|\nabla u(s)\|_{L^2(\Omega)}^2 ds \right)^{\frac{1}{2}} \\ \leq \|u_0\|_{L^2(\Omega)} + \int_0^t \|f(s)\|_{L^2(\Omega)} ds \quad t > 0 \end{aligned} \quad (5.19)$$

Now that we have found an estimate for (5.5), we can now estimate its discretization (5.6), like (5.15)

$$\|u_h(t)\|_{L^2(\Omega)}^2 + \alpha \int_0^t \|\nabla u_h(s)\|_{L^2(\Omega)}^2 ds \leq \|u_{0h}\|_{L^2(\Omega)}^2 + \frac{C_\Omega^2}{\alpha} \int_0^t \|f(s)\|_{L^2(\Omega)}^2 ds \quad (5.20)$$

We can proceed as before, simply taking for every  $t > 0$ ,  $v_h = u_h(t)$ . Since  $u_h(0) = u_{0h}$  we obtain the discrete counterparts of (5.17) and (5.19)

$$\|u_h(t)\|_{L^2(\Omega)} \leq \|u_{0h}(t)\|_{L^2(\Omega)} + \int_0^t \|f(s)\|_{L^2(\Omega)} ds \quad t > 0 \quad (5.21)$$

and

$$\begin{aligned} \left( \|u_h(t)\|_{L^2(\Omega)}^2 + 2\alpha \int_0^t \|\nabla u_h(s)\|_{L^2(\Omega)}^2 ds \right)^{\frac{1}{2}} \\ \leq \|u_{0h}(t)\|_{L^2(\Omega)} + \int_0^t \|f(s)\|_{L^2(\Omega)} ds \quad t > 0 \end{aligned} \quad (5.22)$$

## 5.6 Convergence analysis

### Theorem 5.2

There exists a constant  $C > 0$  independent of both  $t$  and  $h$  s.t.

$$\begin{aligned} \left\{ \|u(t) - u_h(t)\|_{L^2(\Omega)}^2 + \alpha \int_0^t \|\nabla u(s) - \nabla u_h(s)\|_{L^2(\Omega)}^2 ds \right\}^{\frac{1}{2}} \\ \leq Ch^r \left\{ |u_0|_{H^r(\Omega)}^2 + \int_0^t |u(s)|_{H^{r+1}(\Omega)}^2 ds + \int_0^t \left| \frac{\partial u(s)}{\partial s} \right|_{H^{r+1}(\Omega)}^2 ds \right\}^{\frac{1}{2}} \end{aligned}$$

## 5.7 Stability analysis of the $\theta$ -method

We can now analyze the stability of discretized problem. By applying the  $\theta$ -method to the Galerkin problem (5.6) we obtain

$$\begin{aligned} \left( \frac{u_h^{k+1} - u_h^k}{\Delta t}, v_h \right) + a(\theta u_h^{k+1} + (1 - \theta) u_h^k, v_h) \\ = \theta F^{k+1}(v_h) + (1 - \theta) F^k(v_h) \quad \forall v_h \in V_h \end{aligned} \quad (5.23)$$

for each  $k \geq 0$ , with  $F^k$  indicating the functional evaluated at the time  $k$ . For the analysis we will consider  $F = 0$ , and the case of implicit Euler  $\theta = 1$ .

$$\left( \frac{u_h^{k+1} - u_h^k}{\Delta t}, v_h \right) + a(u_h^{k+1}, v_h) = 0 \quad \forall v_h \in V_h$$

By choosing  $v_h = u_h^{k+1}$ , we obtain

$$(u_h^{k+1}, u_h^{k+1}) + \Delta t a(u_h^{k+1}, u_h^{k+1}) = (u_h^k, u_h^{k+1})$$

and then, thanks to these inequalities

$$\begin{aligned} a(u_h^{k+1}, u_h^{k+1}) &\geq \alpha \|u_h^{k+1}\|_V^2 \\ (u_h^k, u_h^{k+1}) &\leq \frac{1}{2} \|u_h^k\|_{L^2(\Omega)}^2 + \frac{1}{2} \|u_h^{k+1}\|_{L^2(\Omega)}^2 \end{aligned}$$

which are derived from the coercivity of  $a(\cdot, \cdot)$  and the CS inequality we obtain

$$\|u_h^{k+1}\|_{L^2(\Omega)}^2 + 2\alpha\Delta t \|u_h^{k+1}\|_V^2 \leq \|u_h^k\|_{L^2(\Omega)}^2 \quad (5.24)$$

Observing that  $\|u_h^{k+1}\|_V \geq \|u_h^{k+1}\|_{L^2(\Omega)}$ , we deduce from (5.24) that

$$(1 + 2\alpha\Delta t) \|u_h^{k+1}\|_{L^2(\Omega)}^2 \leq \|u_h^k\|_{L^2(\Omega)}^2$$

hence

$$\|u_h^{k+1}\|_{L^2(\Omega)} \leq \frac{1}{\sqrt{1 + 2\alpha\Delta t}} \|u_h^k\|_{L^2(\Omega)}$$

which entails

$$\|u_h^k\|_{L^2(\Omega)} \leq \left( \frac{1}{\sqrt{1 + 2\alpha\Delta t}} \right)^k \|u_{0h}\|_{L^2(\Omega)}$$

and therefore

$$\lim_{k \rightarrow \infty} \|u_h^k\|_{L^2(\Omega)} = 0$$

so the backward Euler method is absolutely stable without any restriction on  $\Delta t$ .

Now we assume  $f \neq 0$

$$\underbrace{\left( \frac{u_h^{k+1} - u_h^k}{\Delta t}, u_h^{k+1} \right)}_{(1)} + \underbrace{a(u_h^{k+1}, u_h^{k+1})}_{(2)} = \underbrace{\int_{\Omega} f^{k+1} u_h^{k+1}}_{(3)}$$

where

$$\begin{aligned} (1) &\geq \frac{1}{2\Delta t} \left( \|u_h^{k+1}\|_{L^2(\Omega)}^2 - \|u_h^k\|_{L^2(\Omega)}^2 \right) \\ (2) &\geq \alpha \|u_h^{k+1}\|_{L^2(\Omega)}^2 \\ (3) &\stackrel{(CS)}{\leq} \|f^{k+1}\|_{L^2(\Omega)} \|u_h^{k+1}\|_V \stackrel{(Young)}{\leq} \frac{1}{2\alpha} \|f^{k+1}\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|u_h^{k+1}\|_V^2 \end{aligned}$$

Then, after summation on  $k$ , for  $k = 0, \dots, n-1$ :

$$\begin{aligned} \|u_h^n\|_{L^2(\Omega)}^2 + \underbrace{\alpha \sum_{k=1}^n \Delta t \|u_h^k\|_V^2}_{\simeq \alpha \int_0^{t_n} \|u_h(t)\|_V^2 dt} &\leq \|u_{0h}\|_{L^2(\Omega)}^2 + \underbrace{\frac{1}{\alpha} \sum_{k=1}^n \Delta t \|f^k\|_{L^2(\Omega)}^2}_{\simeq \frac{1}{\alpha} \int_0^{t_n} \|f^k\|_{L^2(\Omega)}^2 dt} \end{aligned}$$

meaning that we have unconditional stability.

Now, before analyzing the general case, where  $0 \leq \theta \leq 1$ , we need the following definition

**Definition 5.4**

We say that the scalar  $\lambda$  is an eigenvalue of the bilinear form  $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$  and that  $w \in V$  is its corresponding eigenfunction if

$$a(w, v) = \lambda(w, v) \quad \forall v \in V$$

If the bilinear form  $a(\cdot, \cdot)$  is symmetric and coercive, it has positive and real eigenvalues forming an infinite sequence, moreover, its eigenfunctions form a basis of the space  $V$ . The eigenvalues and eigenfunctions of  $a(\cdot, \cdot)$  can be approximated by finding the pairs  $\lambda_h \in \mathbb{R}$  and  $w_h \in V_h$ , which satisfy

$$a(w_h, w_h) \lambda_h(w_h, v_h) \quad \forall v_h \in V_h \quad (5.25)$$

From an algebraic viewpoint, problem (5.25) can be formulated as

$$A\mathbf{w} = \lambda_h M\mathbf{w},$$

where  $A$  is the stiffness matrix and  $M$  the mass matrix. This is a generalized eigenvalue problem. Such eigenvalue are all positive and  $N_h$  (the usual dimension of  $V_h$ ). After ordering them in ascending order we have

$$\lambda_h^{N_h} \quad \text{for } N_h \rightarrow \infty.$$

Moreover, the corresponding eigenfunctions form a basis for the subspace  $V_h$  and can be chosen to be orthonormal w.r.t. the scalar product of  $L^2(\Omega)$ . This means that, by denoting with  $w_h^i$  the eigenfunction corresponding to the eigenvalue  $\lambda_h^i$ , we have  $(w_h^i, w_h^j) = \delta_{ij}$ ,  $\forall i, j = 1, \dots, N_h$ . Thus, each function  $v_h \in V_h$  can be represented as follows

$$v_h(\mathbf{x}) = \sum_{j=1}^{N_h} v_j w_h^j(\mathbf{x})$$

and, thanks to the eigenfunctions orthonormality

$$\|v_h\|_{L^2(\Omega)}^2 = \sum_{j=1}^{N_h} v_j^2 \quad (5.26)$$

Let us now consider an arbitrary  $\theta \in [0, 1]$  and assume that  $a(\cdot, \cdot)$  is symmetric. Since  $u_h^k \in V_h$  we can write

$$u_h^k(\mathbf{x}) = \sum_{j=1}^{N_h} u_j^k w_h^j(\mathbf{x})$$

In this case, however,  $u_j^k$  no longer represents the nodal values of  $u_h^k$ .

If we now set  $F = 0$  in (5.23) and take  $v_h = w_h^i$ , we find

$$\frac{1}{\Delta t} \sum_{k=1}^{N_h} [u_j^{k+1} - u_j^k] (w_h^i, w_h^j) + \sum_{j=1}^{N_h} [\theta u_j^{k+1} + (1 - \theta) u_j^k] a(w_h^j, w_h^i) = 0,$$

for each  $i = 1, \dots, N_h$ . For each pair  $i, j = 1, \dots, N_h$  we have

$$a(w_h^j, w_h^i) = \lambda_h^j (w_h^j, w_h^i) = \lambda_h^i,$$

and thus, for each  $i = 1, \dots, N_h$ ,

$$\frac{u_i^{k+1} - u_i^k}{\Delta t} + [\theta u_i^{k+1} + (1 - \theta) u_i^k] \lambda_h^i = 0.$$

Solving for  $u_i^{k+1}$ , we find:

$$u_i^{k+1} = u_i^k \frac{1 - (1 - \theta)\lambda_h^i \Delta t}{1 + \theta\lambda_h^i \Delta t}.$$

Now, recalling (5.26), we can conclude that absolute stability comes from the following inequality

$$\left| \frac{1 - (1 - \theta)\lambda_h^i \Delta t}{1 + \theta\lambda_h^i \Delta t} \right| < 1$$

that is

$$-1 - \theta\lambda_h^i \Delta t < 1 - (1 - \theta)\lambda_h^i \theta < 1 + \theta\lambda_h^i \theta.$$

Hence,

$$-\frac{2}{\lambda_h^i \Delta t} - \theta < \theta - 1 < \theta.$$

While the second inequality is always verified, we can rewrite the first one as:

$$2\theta - 1 > -\frac{2}{\lambda_h^i \Delta t}.$$

If  $\theta \geq \frac{1}{2}$ , the left hand side is nonnegative, while the right hand side is negative, so the relation always holds. In the case  $\theta < \frac{1}{2}$ , the inequality is satisfied if

$$\Delta t < \frac{2}{(1 - 2\theta)\lambda_h^i}. \quad (5.27)$$

So, we have that

- if  $\theta \geq \frac{1}{2}$ , the  $\theta$ -method is unconditionally absolutely stable, for every value of  $\Delta t$ .
- if  $\theta < \frac{1}{2}$  the  $\theta$ -method is absolutely stable only if (5.27) is satisfied.

Thanks to the definition of eigenvalue (5.25) and the continuity property of  $a(\cdot, \cdot)$ , we deduce that

$$\lambda_h^{N_h} = \frac{a(w_{N_h}, w_{N_h})}{\|w_{N_h}\|_{L^2(\Omega)}^2} \leq \frac{M\|w_{N_h}\|_V^2}{\|w_{N_h}\|_{L^2(\Omega)}^2} \leq M(1 + C^2 h^{-2}).$$

The constant  $C > 0$ , which appears in the latter step derives from the following inverse inequality

$$\exists C > 0 : \|\nabla v_h\|_{L^2(\Omega)} \leq Ch^{-1}\|v_h\|_{L^2(\Omega)} \quad \forall v_h \in V_h$$

Hence, for  $h$  small enough,  $\lambda_h^{N_h} \leq Ch^{-2}$ . In fact, we can prove that  $\lambda_h^{N_h}$  is indeed of the order  $h^{-2}$ , that is

$$\lambda_h^{N_h} = \max_i \lambda_h^i \simeq ch^{-2}.$$

Knowing that, we obtain the stability of the  $\theta$ -method for  $\theta < \frac{1}{2}$ , which is

$$\Delta t \leq C(\theta)h^2 \quad (5.28)$$

where  $C(\theta)$  denotes a positive constant depending on  $\theta$ . We cannot choose  $\Delta t$  without keeping in mind  $h$ .

## 5.8 Convergence analysis

### Theorem 5.3

Under the hypothesis that  $u_0, f$  and the exact solution are sufficiently regular, the following a priori

estimate holds  $\forall n \geq 1$ :

$$\|u(t^n) - u_h^n\|_{L^2(\Omega)}^2 + 2\alpha\Delta t \sum_{k=1}^n \|u(t^k) - u_h^k\|_V^2 \leq C(u_0, f, u)(\Delta t^{p(\theta)} + h^{2r})$$

where  $p(\theta) = 2$  if  $\theta \neq \frac{1}{2}$ ,  $p(\frac{1}{2}) = 4$  and  $C$  depends on its arguments, but not on  $\Delta t$  or  $h$ .

## 5.9 Parabolic ADR equation

Consider the parabolic PDE, where  $\Omega \subset \mathbb{R}^2$  is an open bounded domain

$$\begin{cases} \frac{\partial u}{\partial t} - \mu \Delta u + \beta \cdot \nabla u + \sigma u = f & \text{in } \Omega \times (0, T) \\ u = 0 & \text{on } \partial\Omega \times (0, T) \\ u(0) = u_0 & \text{in } \Omega \end{cases} \quad (5.29)$$

where  $\mu, \beta, \sigma$  and  $f$  are regular functions, satisfying:

$$\begin{aligned} 0 < \mu_0 \leq \mu \leq \mu_1 & \quad \text{a.e. in } \Omega \\ |\beta| \leq b_1 \text{ a.e. in } \Omega & \\ 0 < \sigma_0 \leq \sigma \leq \sigma_1 \text{ a.e. in } \Omega & \end{aligned}$$

Then, introducing a finite element space  $V_h \subset H_0^1(\Omega)$ , the semi-discrete Galerkin formulation reads

$$\begin{aligned} \text{for all } t \in (0, T] \text{ find } u_h(t) \in V_h : \int_{\Omega} \frac{\partial u_h(t)}{\partial t} v_h dx + \int_{\Omega} \mu \nabla u_h(t) \cdot \nabla v_h + \int_{\Omega} \beta \cdot \nabla u_h(t) v_h \\ + \int_{\Omega} \sigma u_h(t) v_h = \int_{\Omega} f v_h \quad \forall v_h \in V_h \end{aligned} \quad (5.30)$$

And such that  $u_h(0) = u_{0h}$ , where  $u_{0h}$  is the projection of the initial condition into  $V_h$ .

### 5.9.1 A semimplicit scheme

Consider now a time-advancin scheme, where the diffusion and reaction term are treated implicitly, while the advection term is treated explicitly. Let us denote  $t_k = k\Delta t$ , for  $k = 0, \dots, N$ , where  $\Delta t = \frac{T}{N}$ . Let now  $u_h^k$  be the approximation of  $u(t_k)$ . A fully discretized version of (5.30) reads:

$$\begin{cases} \left( \frac{u_h^{k+1} - u_h^k}{\Delta t}, v_h \right) + \left( \mu \nabla u_h^{k+1}, v_h \right) + \left( \beta \cdot \nabla u_h^k, v_h \right) \\ \quad + \left( \sigma u_h^{k+1}, v_h \right) = (f, v_h) \quad \forall v_h \in V_h, \quad k = 0, \dots, N-1 \\ u_h(0) = u_{0h} \end{cases} \quad (5.31)$$

#### Theorem 5.4

If the coefficients of the problem staisfy

$$b_1^2 < 4\mu_0\sigma_0 \quad (5.32)$$

then the semimplicit scheme (5.31) is absolutely stable for any chance of  $\Delta t$ . Consider now the case  $\sigma = 0$ . If the coefficients of the problem satisfy

$$b_1 < \frac{\mu_0}{C_p} \quad (5.33)$$

with  $C_p$  being the Poincaré constant, then the scheme is absolutely stable, for any choice of  $\Delta t$ .

**Proof.** Let us choose  $v_h = u_h^{k+1}$ . We have

$$\begin{aligned} \left( \mu \nabla u_h^{k+1}, \nabla u_h^{k+1} \right) &\geq \mu_0 \left\| \nabla u_h^{k+1} \right\|_{L^2(\Omega)}^2 \\ \left( \sigma u_h^{k+1}, u_h^{k+1} \right) &\geq \sigma_0 \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 \end{aligned}$$

which entails, for every  $k$

$$\left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 + \Delta t \mu_0 \left\| \nabla u_h^{k+1} \right\|_{L^2(\Omega)}^2 + \Delta t \sigma_0 \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 \leq \left| \left( u_h^k, u_h^{k+1} \right) \right| + \Delta t \left| \left( \beta \cdot \nabla u_h^k, u_h^{k+1} \right) \right|$$

The two right-hand side terms can be bounded by combining (5.12) and (5.11):

$$\begin{aligned} \left| \left( u_h^k, u_h^{k+1} \right) \right| &\leq \frac{1}{2\eta_1} \left\| u_h^k \right\|_{L^2(\Omega)}^2 + \frac{\eta_1}{2} \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 \\ \left| \left( \beta \cdot \nabla u_h^k, u_h^{k+1} \right) \right| &\leq \frac{b_1}{2\eta_2} \left\| \nabla u_h^k \right\|_{L^2(\Omega)}^2 + \frac{\eta_2 b_1}{2} \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 \end{aligned}$$

where the positive constant  $\eta_1$  and  $\eta_2$  will be later fixed accordingly. Now we end up with the following inequality

$$\begin{aligned} \underbrace{\left[ 1 + \Delta t \sigma_0 - \frac{\eta_1}{2} - \frac{\Delta t \eta_2 b_1}{2} \right]}_A \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 + \underbrace{\Delta t \mu_0}_B \left\| \nabla u_h^{k+1} \right\|_{L^2(\Omega)}^2 \\ \leq \underbrace{\frac{1}{2\eta_1}}_{A'} \left\| u_h^k \right\|_{L^2(\Omega)}^2 + \underbrace{\frac{\Delta t b_1}{2\eta_2}}_{B'} \left\| \nabla u_h^k \right\|_{L^2(\Omega)}^2 \end{aligned}$$

In order to prove stability we need  $A > A'$  and  $B > B'$ . In fact, if this were true, we would end up with

$$A \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 + B \left\| \nabla u_h^{k+1} \right\|_{L^2(\Omega)}^2 \leq \max \left( \frac{A'}{A}, \frac{B'}{B} \right) \left[ A \left\| u_h^{k+1} \right\|_{L^2(\Omega)}^2 + B \left\| \nabla u_h^{k+1} \right\|_{L^2(\Omega)}^2 \right],$$

which is the sought stability in the norm  $\|\cdot\|_{A,B} := \left( A \|\cdot\|_{L^2(\Omega)}^2 + B \|\nabla \cdot\|_{L^2(\Omega)}^2 \right)^{\frac{1}{2}}$ , equivalent to the standard  $V_h$  norm.

Therefore, we look for a suitable choice of  $\eta_1$  and  $\eta_2$  that ensures  $A > A'$  and  $B > B'$ . The second inequality is satisfied if and only if

$$\eta_2 = \frac{b_1 + \varepsilon}{2\mu_0}$$

for some  $\varepsilon > 0$ .

Hence, the first inequality reads

$$1 + \Delta t \sigma_0 - \frac{\Delta t b_1 (b_1 + \varepsilon)}{4\mu_0} > \frac{1}{2\eta_1} + \frac{\eta_1}{2}$$

The right hand side is minimized for  $\eta_1 = 1$ , thus leading to the condition

$$4 \frac{\mu_0 \sigma_0}{b_1 (b_1 + \varepsilon)} > 1 \tag{5.34}$$

Clearly, it is possible to find  $\varepsilon > 0$  such that (5.34) holds if and only if

$$b_1^2 \leq 4\mu_0 \sigma_0 \tag{5.35}$$

In conclusion, whenever the coefficients of the problem satisfy the condition (5.35), the scheme (5.31) is absolutely stable for any choice of  $\Delta t$ .

Let us now consider  $\sigma = 0$ . Proceeding as before, we have:

$$\begin{aligned} \left[1 - \frac{\eta_1}{2} - \frac{\Delta t \eta_2 b_1}{2}\right] \|u_h^{k+1}\|_{L^2(\Omega)}^2 + \Delta t \mu_0 \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 \\ \leq \frac{1}{2\eta_1} \|u_h^k\|_{L^2(\Omega)}^2 + \frac{\Delta t b_1}{2\eta_2} \|\nabla u_h^k\|_{L^2(\Omega)}^2 \end{aligned}$$

Introducing now a constant  $\omega \in (0, 1)$ , we have (thanks to (5.13)):

$$\begin{aligned} \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 &= (1 - \omega) \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 + \omega \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 \\ &\geq \frac{1 - \omega}{C_p} \|u_h^{k+1}\|_{L^2(\Omega)}^2 + \omega \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 \end{aligned}$$

Combining the latter inequalities, we obtain

$$\begin{aligned} \underbrace{\left[1 - \frac{\eta_1}{2} - \frac{\Delta t \eta_2 b_1}{2} - \frac{(1 - \omega) \Delta t \mu_0}{C_p^2}\right]}_A \|u_h^{k+1}\|_{L^2(\Omega)}^2 + \underbrace{\omega \Delta t \mu_0}_B \|\nabla u_h^{k+1}\|_{L^2(\Omega)}^2 \\ \leq \underbrace{\frac{1}{2\eta_1}}_{A'} \|u_h^k\|_{L^2(\Omega)}^2 + \underbrace{\frac{\Delta t b_1}{2\eta_2}}_{B'} \|\nabla u_h^k\|_{L^2(\Omega)}^2 \end{aligned}$$

As before, we look for conditions such that  $A > A'$  and  $B > B'$ . The second inequality is satisfied if and only if

$$\eta_2 = \frac{b_1 + \varepsilon}{2\omega\mu_0}$$

for some  $\varepsilon > 0$ . Then, the first inequality reads

$$1 - \frac{\Delta t b_1 (b_1 + \varepsilon)}{4\omega\mu_0} + \frac{(1 - \omega) \Delta t \mu_0}{C_p^2} > \frac{1}{2\eta_1} + \frac{\eta_1}{2}$$

The right hand side is minimized for  $\eta_1 = 1$ . Rearranging the term, we get

$$-\omega^2 + \omega - \frac{b_1(b_1 + \varepsilon)C_p^2}{4\mu_0^2} > 0 \tag{5.36}$$

Real solutions  $\omega \in (0, 1)$  exists whenever the discriminant is positive, that is

$$b_1(b_1 + \varepsilon)C_p^2 \leq \mu_0^2$$

The latter condition can be satisfied (by suitable choosing of  $\varepsilon$ ) if and only if

$$b_1 < \frac{\mu_0}{C_p} \tag{5.37}$$

In conclusion, if (5.37) is satisfied, the scheme is absolutely stable for any choice of  $\Delta t$ . ★

## 6 Navier-Stokes equations

### 6.1 Introduction

Navier-Stokes equations describe the motion of a fluid with constant density  $\rho$  in a domain  $\Omega \subset \mathbb{R}^d$  with  $d = 2, 3$ . They read as follows

$$\begin{cases} \frac{\partial \mathbf{u}}{\partial t} - \operatorname{div} [\nu (\nabla \mathbf{u} + \nabla \mathbf{u}^T)] + (\mathbf{u} \cdot \nabla) \mathbf{u} + \nabla p = \mathbf{f} & \mathbf{x} \in \Omega, t > 0 \\ \operatorname{div} \mathbf{u} = 0 & \mathbf{x} \in \Omega, t > 0 \end{cases} \quad (6.1)$$

where

- $\mathbf{u}$  is the fluid's velocity
- $p$  is the pressure divided by the density (which will be simply called “pressure”)
- $\nu$  is the kinematic viscosity
- $\mathbf{f} \in L^2(\mathbb{R}^+; [L^2(\Omega)])$  is a forcing term per unit of mass

The first equation represents conservation of linear momentum, while the second one the conservation of mass. In fact, the term  $(\mathbf{u} \cdot \nabla) \mathbf{u}$  describes the process of convective transport, while the term  $-\operatorname{div} [\nu (\nabla \mathbf{u} + \nabla \mathbf{u}^T)]$  describe the process of molecular diffusion.

When  $\nu$  is constant, from the continuity equation we obtain

$$\operatorname{div} [\nu (\nabla \mathbf{u} + \nabla \mathbf{u}^T)] = \nu (\Delta \mathbf{u} + \nabla \operatorname{div} \mathbf{u}) = \nu \Delta \mathbf{u}$$

and system (6.1) can be rewritten as

$$\begin{cases} \frac{\partial \mathbf{u}}{\partial t} - \nu \Delta \mathbf{u} + (\mathbf{u} \cdot \nabla) \mathbf{u} + \nabla p = \mathbf{f} & \mathbf{x} \in \Omega, t > 0 \\ \operatorname{div} \mathbf{u} = 0 & \mathbf{x} \in \Omega, t > 0 \end{cases} \quad (6.2)$$

The (6.2) are often called incompressible Navier-Stokes equations. More in general, when  $\operatorname{div} \mathbf{u} = 0$ , fluids are said to be incompressible.

In order for (6.2) to be well posed, it is necessary to assign the initial condition

$$\mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \quad \forall \mathbf{x} \in \Omega \quad (6.3)$$

where  $\mathbf{u}_0$  is a given divergence-free vector field. Then we would need need suitable BC

$$\begin{cases} \mathbf{u}(\mathbf{x}, t) = \boldsymbol{\varphi}(\mathbf{x}, t) & \forall \mathbf{x} \in \Gamma_D, t > 0 \\ \left( \nu \frac{\partial \mathbf{u}}{\partial \mathbf{n}} - p \mathbf{n} \right)(\mathbf{x}, t) = \boldsymbol{\psi}(\mathbf{x}, t) & \forall \mathbf{x} \in \Gamma_N, t > 0 \end{cases} \quad (6.4)$$

where  $\boldsymbol{\varphi}$  and  $\boldsymbol{\psi}$  are given vector functions, while  $\Gamma_D$  and  $\Gamma_N$  provide a partition of the boundary, such that  $\Gamma_D \cup \Gamma_N = \partial\Omega$ ,  $\Gamma_D^\circ \cap \Gamma_N^\circ = \emptyset$ . Finally, as usual  $\mathbf{n}$  is the outward unit normal vector to  $\partial\Omega$ .

If we denote with  $u_i, i = 1, \dots, d$  the components of  $uvec$  w.r.t. a Cartesian frame, and, likewise, with  $f_i$  components of  $\mathbf{f}$ , we obtain, from (6.2)

$$\begin{cases} \frac{\partial u_i}{\partial t} - \nu \nabla^2 u_i + \sum_{j=1}^d u_j \frac{\partial u_i}{\partial x_j} + \frac{\partial p}{\partial x_i} = f_i & i = 1, \dots, d \\ \sum_{j=1}^d \frac{\partial u_j}{\partial x_j} = 0 \end{cases}$$



**Remark 6.1**

The Navier-Stokes equations have been written in terms of the primitive variable  $\mathbf{u}$  and  $p$ , but other sets of variables may be used too. For instance, in the two dimensional case it is common to see the vorticity  $\omega$  and the streamfunction  $\psi$  that are related to the velocity

$$\omega = \text{rot } \mathbf{u} = \frac{\partial u_2}{\partial x_1} - \frac{\partial u_1}{\partial x_2}, \quad \mathbf{u} = \begin{bmatrix} \frac{\partial \psi}{\partial x_2} \\ -\frac{\partial \psi}{\partial x_1} \end{bmatrix}.$$

The various formulations are equivalent from a mathematical standpoint, but may give rise to different numerical methods.

**6.2 Weak formulation**

A weak formulation of the problem can be obtained as usual, by multiplying a test function  $\mathbf{v}$ , belonging to a suitable space  $V$ , and integrate in  $\Omega$

$$\int_{\Omega} \frac{\partial \mathbf{u}}{\partial t} \cdot \mathbf{v} \, d\Omega - \int_{\Omega} \nu \Delta \mathbf{u} \cdot \mathbf{v} \, d\Omega + \int_{\Omega} [(\mathbf{u} \cdot \nabla) \mathbf{u}] \cdot \mathbf{v} \, d\Omega + \int_{\Omega} \nabla p \cdot \mathbf{v} \, d\Omega = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega \quad (6.5)$$

Using Green's formula we find

$$\begin{aligned} - \int_{\Omega} \nu \Delta \mathbf{u} \cdot \mathbf{v} \, d\Omega &= \int_{\Omega} \nu \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\Omega - \int_{\partial\Omega} \nabla p \cdot \mathbf{v} \, d\Omega \\ \int_{\Omega} \nabla p \cdot \mathbf{v} \, d\Omega &= - \int_{\Omega} p \text{div } \mathbf{v} \, d\Omega + \int_{\partial\Omega} p \mathbf{v} \cdot \mathbf{n} \, d\gamma \end{aligned}$$

Using these relations in the first of (6.2) we obtain

$$\begin{aligned} \int_{\Omega} \frac{\partial \mathbf{u}}{\partial t} \cdot \mathbf{v} \, d\Omega + \int_{\Omega} \nu \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\Omega + \int_{\Omega} [(\mathbf{u} \cdot \nabla) \mathbf{u}] \cdot \mathbf{v} \, d\Omega \\ - \int_{\Omega} p \text{div } \mathbf{v} \, d\Omega = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\partial\Omega} \left( \nu \frac{\partial \mathbf{u}}{\partial \mathbf{n}} - p \mathbf{n} \right) \cdot \mathbf{v} \, d\gamma \quad \forall \mathbf{v} \in V \end{aligned} \quad (6.6)$$

Similarly, by multiplying the second equation of (6.2) by a test function  $q$ , belonging to a suitable space  $Q$  to be specified, then integrating in  $\Omega$  it follows that

$$\int_{\Omega} q \text{div } \mathbf{u} \, d\Omega = 0 \quad \forall q \in Q \quad (6.7)$$

Usually,  $V$  is chosen so that the test functions vanish on the boundary portion where a Dirichlet data is prescribed on  $\mathbf{u}$

$$V = [H_{\Gamma_D}^1]^d = \{ \mathbf{v} \in [H^1(\Omega)] : \mathbf{v}|_{\Gamma_D} = \mathbf{0} \} \quad (6.8)$$

It will coincide with  $[H_0^1(\Omega)]^d$  if  $\Gamma_D = \partial\Omega$ .

If  $\Gamma_N$  has positive measure, we can choose  $Q = L^2(\Omega)$ . If  $\Gamma_D = \partial\Omega$ , then the pressure space should be  $L_0^2$  to ensure uniqueness for the pressure  $p$ .

Moreover, if  $t > 0$ , then  $\mathbf{u}(t) \in [H^1(\Omega)]^d$ , with  $\mathbf{u}(t) = \boldsymbol{\varphi}(t)$  on  $\Gamma_D$ ,  $\mathbf{u}(0) = \mathbf{u}_0$  and  $p(t) \in Q$ . Having chosen these functional spaces, we can note first of all that

$$\int_{\partial\Omega} \left( \nu \frac{\partial \mathbf{u}}{\partial \mathbf{n}} - p \mathbf{n} \right) \cdot \mathbf{v} \, d\gamma = \int_{\Gamma_N} \boldsymbol{\psi} \cdot \mathbf{v} \, d\gamma \quad \forall \mathbf{v} \in V$$

**Notation 6.1**

For every function  $\mathbf{v} \in \mathbf{H}^1(\Omega)$ , we denote by

$$\|\mathbf{v}\|_{\mathbf{H}^1(\Omega)} = \left( \sum_{k=1}^d \|v_k\|_{H^1(\Omega)}^2 \right)^{\frac{1}{2}}$$

its norm, and by

$$|\mathbf{v}|_{\mathbf{H}^1(\Omega)} = \left( \sum_{k=1}^d |v_k|_{H^1(\Omega)}^2 \right)^{\frac{1}{2}}$$

its seminorm. Thanks to Poincaré's inequality,  $|\mathbf{v}|_{\mathbf{H}^1(\Omega)}$  is equivalent to  $\|\mathbf{v}\|_{\mathbf{H}^1(\Omega)}$  for all functions belonging to  $V$ , provided that the Dirichlet boundary has a positive measure.

All the integrals involving bilinear terms are finite. To be more precise, by using the vector notation  $\mathbf{H}^k(\Omega) = [H^k(\Omega)]^d$ ,  $\mathbf{L}^p(\Omega) = [L^p(\Omega)]^d$ ,  $k \geq 1$ ,  $1 \leq p < \infty$ , we find

$$\begin{aligned} \left| \nu \int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\Omega \right| &\leq \nu |\mathbf{u}|_{\mathbf{H}^1(\Omega)} |\mathbf{v}|_{\mathbf{H}^1(\Omega)} \\ \left| \int_{\Omega} p \operatorname{div} \mathbf{v} \, d\Omega \right| &\leq \|p\|_{L^2(\Omega)} |\mathbf{v}|_{\mathbf{H}^1(\Omega)} \\ \left| \int_{\Omega} q \nabla \mathbf{u} \, d\Omega \right| &\leq \|q\|_{L^2(\Omega)} |\mathbf{u}|_{\mathbf{H}^1(\Omega)} \end{aligned}$$

Also the integral involving the trilinear term is finite. Before we see how let's recall the following result: if  $d \leq 3$

$$\forall \mathbf{v} \in \mathbf{H}^1(\Omega), \text{ then } \mathbf{v} \in \mathbf{L}^4(\Omega) \text{ and } \exists C > 0 : \|\mathbf{v}\|_{\mathbf{L}^4(\Omega)} \leq C \|\mathbf{v}\|_{\mathbf{H}^1(\Omega)}.$$

Then, using the following three-term Hölder inequality

$$\left| \int_{\Omega} f g h \, d\Omega \right| \leq \|f\|_{L^p(\Omega)} \|g\|_{L^q(\Omega)} \|h\|_{L^r(\Omega)},$$

valid for all  $p, q, r > 1$  such that  $p^{-1} + q^{-1} + r^{-1} = 1$ , we conclude that

$$\left| \int_{\Omega} [(\mathbf{u} \cdot \nabla) \mathbf{u}] \cdot \mathbf{v} \, d\Omega \right| \leq \|\nabla \mathbf{u}\|_{\mathbf{L}^2(\Omega)} \|\mathbf{u}\|_{\mathbf{L}^4(\Omega)} \|\mathbf{u}\|_{\mathbf{L}^4(\Omega)} \leq C^2 \|\mathbf{u}\|_{\mathbf{H}^1(\Omega)} \|\mathbf{v}\|_{\mathbf{H}^1(\Omega)}.$$

### 6.3 Solution uniqueness

As for the solution's uniqueness, let us consider again (6.2). If  $\Gamma_D = \partial\Omega$ , when only boundary conditions of Dirichlet type are imposed, the pressure merely appears in terms of its gradient. In that case if we call  $(\mathbf{u}, p)$  a solution, then for any constant  $c$  the couple  $(\mathbf{u}, p + c)$  is a solution too since  $\nabla(p + c) = \nabla p$ .

To avoid that, one can fix a priori  $p$  at a given point  $\mathbf{x}_0$  of the domain  $\Omega$  such that  $p(\mathbf{x}_0) = p_0$ , or, alternatively, require the pressure average to be null,  $\int_{\Omega} p \, d\Omega = 0$ . This condition requires to prescribe a pointwise value for the pressure, but this is inconsistent with the hypothesis that  $p \in L^2$ . For this reason, the pressure space will be considered as

$$Q = L_0^2(\Omega) = \left\{ p \in L^2(\Omega) : \int_{\Omega} p \, d\Omega = 0 \right\}.$$

Then, we observe that if  $\Gamma_D = \partial\Omega$ , the prescribed Dirichlet data  $\boldsymbol{\varphi}$  must be compatible with the incompressibility constraint

$$\int_{\Omega} \boldsymbol{\varphi} \cdot \mathbf{n} \, d\gamma = \int_{\Omega} \operatorname{div} \mathbf{u} \, d\Omega = 0.$$

If  $\Gamma_N$  is not empty, we are ok with using  $L^2(\Omega)$  for our pressure space. So:

$$Q = L^2(\Omega) \quad \text{if } \Gamma_N \neq \emptyset, \quad Q = L_0^2(\Omega) \quad \text{if } \Gamma_N = \emptyset \quad (6.9)$$

The weak formulation of (6.2), (6.3), (6.4) is:

$$\begin{aligned} & \text{find } \mathbf{u} \in L^2(\mathbb{R}^+; [H^1(\Omega)]^d) \cap \mathcal{C}^0(\mathbb{R}^+; [L^2(\Omega)]^d), p \in L^2(\mathbb{R}^+; Q) : \\ & \begin{cases} \int_{\Omega} \frac{\partial \mathbf{u}}{\partial t} \cdot \mathbf{v} \, d\Omega + \nu \int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\Omega + \int_{\Omega} [(\mathbf{u} \cdot \nabla) \mathbf{u}] \cdot \mathbf{v} \, d\Omega \\ - \int_{\Omega} p \operatorname{div} \mathbf{v} \, d\Omega = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma_N} \boldsymbol{\psi} \cdot \mathbf{v} \, d\gamma \quad \forall \mathbf{v} \in V \\ \int_{\Omega} q \operatorname{div} \mathbf{u} \, d\Omega = 0 \quad \forall q \in Q \end{cases} \end{aligned} \quad (6.10)$$

with  $\mathbf{u}|_{\Gamma_D} = \boldsymbol{\psi}_D$  and  $\mathbf{u}|_{t=0} = \mathbf{u}_0$ . The space  $V$  is the one in (6.8), while  $Q$  the one in (6.9).

## 6.4 The Reynolds number

Let us define the Reynolds number,

$$Re = \frac{|\mathbf{U}|L}{\nu}$$

where  $\mathbf{U}$  is a representative length of the domain  $\Omega$  (e.g. the length of the channel in which the fluid flows),  $\mathbf{U}$  a representative fluid velocity and  $\nu$  the kinematic viscosity.

This number measures the extent to which convection dominates over diffusion. When  $Re \ll 1$ , the convective term  $(\mathbf{u} \cdot \nabla) \mathbf{u}$  can be omitted, reducing the equations to the so-called Stokes equation. On the other hand, for large values of  $Re$  problems may arise, concerning the uniqueness of the solution, the existence of stationary and stable solutions, the possibility of strange attractors and the transition towards turbulent flows.

## 6.5 Divergence free formulation of Navier-Stokes equations

By eliminating the pressure, the Navier-Stokes equations can be rewritten in a reduced form, with the sole variable  $\mathbf{u}$ . Introducing the following subspaces of  $[H^1(\Omega)]^d$ :

$$\begin{aligned} V_{div} &= \left\{ \mathbf{v} \in [H^1(\Omega)]^d : \operatorname{div} \mathbf{v} = 0 \right\} \\ V_{div}^0 &= \left\{ \mathbf{v} \in V_{div} : \mathbf{v} = \mathbf{0} \text{ on } \Gamma_D \right\} \end{aligned}$$

If the test function  $\mathbf{v}$  belongs to the space  $V_{div}$ , the term associated with the pressure gradient vanishes, whence we find the following reduced problem for the velocity

$$\begin{aligned} & \text{find } \mathbf{u} \in L^2(\mathbb{R}^+; V_{div}) \cap \mathcal{C}^0(\mathbb{R}^+; [L^2(\Omega)]^d) : \\ & \begin{aligned} \int_{\Omega} \frac{\partial \mathbf{u}}{\partial t} \cdot \mathbf{v} \, d\Omega + \nu \int_{\Omega} \nabla \mathbf{u} \cdot \nabla \mathbf{v} \, d\Omega + \int_{\Omega} [(\mathbf{u} \cdot \nabla) \mathbf{u}] \cdot \mathbf{v} \, d\Omega \\ = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega + \int_{\Gamma_N} \boldsymbol{\psi} \cdot \mathbf{v} \, d\gamma \quad \forall \mathbf{v} \in V_{div}, \end{aligned} \end{aligned} \quad (6.11)$$

with  $\mathbf{u}|_{\Gamma_D} = \boldsymbol{\psi}_D$  and  $\mathbf{u}|_{t=0} = \mathbf{u}_0$ .

Since we are dealing with a nonlinear parabolic problem, we can carry an analysis by using techniques similar to those applied in parabolic problems. Clearly a solution of (6.10) will be a suitable solution of (6.11), while, for the converse, we have the following theorem

### Theorem 6.1

Let  $\Omega \subset \mathbb{R}^d$  be a domain with Lipschitz-continuous boundary  $\partial\Omega$ . Let  $\mathbf{u}$  be a solution of the reduced problem (6.11). Then exists a unique function  $p \in L^2(\mathbb{R}^+; Q)$  such that  $(\mathbf{u}, p)$  is a solution of the full problem (6.10)

In practice, however, the results of this theorem, are quite unsuitable from a numerical viewpoint, since it requires the construction of the  $V_{div}$  subspaces, of the divergence-free velocity functions, etc... Moreover, the result of the above theorem is not constructive, as it does not provide a way to build the solution pressure  $p$ .

## 6.6 Stokes equations and their approximation

In this section we will consider the generalized Stokes problem with homogeneous Dirichlet BC

$$\begin{cases} \sigma \mathbf{u} - \nu \Delta \mathbf{u} + \nabla p = \mathbf{f} & \text{in } \Omega \\ \operatorname{div} \mathbf{u} = 0 & \text{in } \Omega \\ \mathbf{u} = \mathbf{0} & \text{on } \partial\Omega \end{cases} \quad (6.12)$$

for  $\sigma \geq 0$ . This is the motion of an incompressible viscous flow in which the convective term has been neglected, since  $Re \ll 1$ . Moreover, one can generate a problem (6.12) also while using an implicit temporal discretization of the Navier-Stokes equations and by neglecting the convective term.

We have indeed the following scheme, where  $k$  denotes the temporal index

$$\begin{cases} \frac{\mathbf{u}^k - \mathbf{u}^{k-1}}{\Delta t} - \nu \Delta \mathbf{u}^k + \nabla p^k = \mathbf{f}(t^k) & \mathbf{x} \in \Omega, \ t > 0 \\ \operatorname{div} \mathbf{u}^k = 0 & \mathbf{x} \in \Omega, \ t > 0 \\ +B.C. \end{cases}$$

Hence, at each time step  $t^k$  we need to solve the following Stokes-like system of equations

$$\begin{cases} \sigma \mathbf{u}^k - \nu \Delta \mathbf{u}^k + \nabla p^k = \tilde{\mathbf{f}}^k & \text{in } \Omega \\ \operatorname{div} \mathbf{u}^k = 0 & \text{in } \Omega \\ +B.C. \end{cases} \quad (6.13)$$

where  $\sigma = (\Delta t)^{-1}$  and  $\tilde{\mathbf{f}}^k = \tilde{\mathbf{f}}(t^k) + \frac{\mathbf{u}^{k-1}}{\Delta t}$ . The weak formulation of problem (6.12) reads:

$$\text{find } \mathbf{u} \in V \text{ and } p \in Q : \begin{cases} \int_{\Omega} (\sigma \mathbf{u} \cdot \mathbf{v} + \nu \nabla \mathbf{u} \cdot \nabla \mathbf{v}) \, d\Omega - \int_{\Omega} p \operatorname{div} \mathbf{v} \, d\Omega = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} \, d\Omega & \forall \mathbf{v} \in V, \\ \int_{\Omega} q \operatorname{div} \mathbf{u} \, d\Omega = 0 & \forall q \in Q, \end{cases} \quad (6.14)$$

where  $V = [H_0^1(\Omega)]^d$  and  $Q = L_0^2(\Omega)$ . Now with the bilinear forms  $a : V \times V \rightarrow \mathbb{R}$  and  $b : V \times Q \rightarrow \mathbb{R}$ :

$$\begin{aligned} a(\mathbf{u}, \mathbf{v}) &= \int_{\Omega} (\sigma \mathbf{u} \cdot \mathbf{v} + \nu \nabla \mathbf{u} \cdot \nabla \mathbf{v}) \, d\Omega, \\ b(\mathbf{u}, q) &= - \int_{\Omega} q \operatorname{div} \mathbf{u} \, d\Omega. \end{aligned} \quad (6.15)$$

Using these notations, problem (6.14) becomes

$$\text{find } (\mathbf{u}, p) \in V \times Q : \begin{cases} a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) = (\mathbf{f}, \mathbf{v}) & \forall \mathbf{v} \in V \\ b(\mathbf{u}, q) = 0 & \forall q \in Q \end{cases} \quad (6.16)$$

where  $(\mathbf{f}, v) = \sum_{i=1}^d \int_{\Omega} f_i v_i \, d\Omega$ .

Considering non homogeneous BC like in (6.4), the weak formulation of the Stokes problem becomes

$$\text{find } (\overset{\circ}{\mathbf{u}}, p) \in V \times Q : \begin{cases} a(\overset{\circ}{\mathbf{u}}, \mathbf{v}) + b(\mathbf{v}, p) = \mathbf{F}(\mathbf{v}) & \forall \mathbf{v} \in V \\ b(\overset{\circ}{\mathbf{u}}, q) = G(q) & \forall q \in Q \end{cases} \quad (6.17)$$

where  $V$  and  $Q$  are the spaces (6.8) and (6.9), respectively. Having denoted with  $\mathbf{R}\boldsymbol{\varphi} \in [H^1(\Omega)]^d$  a lifting of the boundary datum  $\boldsymbol{\varphi}$ , we have set  $\overset{\circ}{\mathbf{u}} = \mathbf{u} - \mathbf{R}\boldsymbol{\varphi}$ , while the new terms on the right hand side have the following expression

$$\begin{aligned} \mathbf{F}(\mathbf{v}) &= (\mathbf{f}, \mathbf{v}) + \int_{\Gamma_N} \psi \mathbf{v} \, d\gamma - a(\mathbf{R}\boldsymbol{\varphi}, \mathbf{v}), \\ G(q) &= -b(\mathbf{R}\boldsymbol{\varphi}, q) \end{aligned} \quad (6.18)$$

**Theorem 6.2**

The couple  $(\mathbf{u}, p)$  solves the problem (6.16) if and only if it is a saddle point of the Lagrangian functional

$$\mathcal{L}(\mathbf{v}, q) = \frac{1}{2}a(\mathbf{v}, \mathbf{v}) + b(\mathbf{v}, q) - (\mathbf{f}, \mathbf{v})$$

or equivalently

$$\mathcal{L}(\mathbf{u}, p) = \min_{\mathbf{v} \in V} \max_{q \in Q} \mathcal{L}(\mathbf{v}, q)$$

The pressure  $q$  hence plays the role of Lagrange multiplier associated to the divergence-free constraint. We need to define the Fréchet derivative.

**Definition 6.1** (Fréchet derivative)

Let  $F : X \rightarrow Y$  with  $X, Y$  two normed vector spaces.  $F$  is differentiable at  $x \in X$  if  $\exists \mathcal{L}_x : X \rightarrow Y$  linear and bounded such that

$$\forall \varepsilon > 0 \exists \delta > 0 \ \|F(x+h) - F(x) - \mathcal{L}_x h\|_Y \leq \varepsilon \|h\|_X \quad \forall h \in X : \|h\|_X < \delta$$

$F'(x) := \mathcal{L}_x$  is called Fréchet derivative of  $F$  at point  $x$ .

By formally taking the Fréchet derivative of the Lagrangian with respect to the two variables (thanks to the symmetry of  $a(\cdot, \cdot)$ ):

$$\begin{aligned} \left\langle \frac{\partial \mathcal{L}(\mathbf{u}, p)}{\partial \mathbf{u}}, \mathbf{v} \right\rangle &= a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) - (\mathbf{f}, \mathbf{v}) = 0 \quad \forall \mathbf{v} \in V \\ \left\langle \frac{\partial \mathcal{L}(\mathbf{u}, p)}{\partial p}, q \right\rangle &= b(\mathbf{u}, q) = 0 \quad \forall q \in Q \end{aligned}$$

**6.7 Galerkin approximation**

The Galerkin approximation of problem (6.16) has the following form:

$$\text{find } (\mathbf{u}_h, p_h) \in V_h \times Q_h : \begin{cases} a(\mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h) = (\mathbf{f}, \mathbf{v}_h) & \forall \mathbf{v}_h \in V_h \\ b(\mathbf{u}_h, q_h) = 0 & \forall q_h \in Q_h \end{cases} \quad (6.19)$$

If, instead we consider problem (6.17)-(6.18), we need to add  $\mathbf{F}(\mathbf{v}_h)$  and  $G(q_h)$ . These new functionals can be obtained from (6.18) by replacing  $\mathbf{R}\boldsymbol{\varphi}$  with the interpolant of  $\boldsymbol{\varphi}$  at the nodes of  $\Gamma_D$ , and replacing  $\boldsymbol{\varphi}$  with its interpolant at the nodes on  $\Gamma_N$ .

The existence and uniqueness is guaranteed by the theorem

**Theorem 6.3**

The Galerkin approximation (6.19) admits one and only one solution if the following conditions hold

- The bilinear form  $a(\cdot, \cdot)$  is:

- (1) coercive, that is  $\exists \alpha > 0$  (possibly depending on  $h$ ) such that:

$$a(\mathbf{v}_h, \mathbf{v}_h) \geq \alpha \|\mathbf{v}_h\|_V^2 \quad \forall \mathbf{v}_h \in V_h^*,$$

where  $V_h^* = \{\mathbf{v}_h \in V : b(\mathbf{v}_h, q_h) = 0 \quad \forall q_h \in Q_h\}$

- (2) continuous, that is  $\exists \gamma > 0$  such that

$$|a(\mathbf{u}_h, \mathbf{v}_h)| \leq \gamma \|\mathbf{u}_h\|_V \|\mathbf{v}_h\|_V \quad \forall \mathbf{u}_h, \mathbf{v}_h \in V_h$$

- The bilinear form  $b(\cdot, \cdot)$  is continuous, that is  $\exists \delta > 0$  such that

$$|b(\mathbf{v}_h, q_h)| \leq \delta \|\mathbf{v}_h\|_V \|q_h\|_Q \quad \forall \mathbf{v}_h \in V_h, q_h \in Q_h$$

- Finally, there exist a positive constant  $\beta$  (possibly depending on  $h$ ) such that:

$$\forall q_h \in Q_h, \exists \mathbf{v}_h \in V_h : b(\mathbf{v}_h, q_h) \geq \beta \|\mathbf{v}_h\|_{\mathbf{H}^1(\Omega)} \|q_h\|_{L^2(\Omega)} \quad (6.20)$$

Under the previous assumptions the discrete solution fullfills the following a priori estimates:

$$\begin{aligned} \|\mathbf{u}_h\|_V &\leq \frac{1}{\alpha} \|\mathbf{f}\|_{V'} \\ \|p_h\|_Q &\leq \frac{1}{\beta} \left(1 + \frac{\gamma}{\alpha}\right) \|\mathbf{f}\|_{V'}, \end{aligned}$$

where  $V'$  is the dual space of  $V$ .

Moreover, the following convergence results hold:

$$\begin{aligned} \|\mathbf{u} - \mathbf{u}_h\|_V &\leq \left(1 + \frac{\delta}{\beta}\right) \left(1 + \frac{\gamma}{\alpha}\right) \inf_{\mathbf{v}_h \in V_h} \|\mathbf{u} - \mathbf{v}_h\|_V \\ &\quad + \frac{\delta}{\alpha} \inf_{q_h \in Q_h} \|p - q_h\|_Q, \\ \|p - p_h\|_Q &\leq \frac{\gamma}{\beta} \left(1 + \frac{\gamma}{\alpha}\right) \left(1 + \frac{\delta}{\beta}\right) \inf_{\mathbf{v}_h \in V_h} \|\mathbf{u} - \mathbf{v}_h\|_V \\ &\quad + \left(1 + \frac{\delta}{\beta} + \frac{\delta\gamma}{\alpha\beta}\right) \inf_{q_h \in Q_h} \|p - q_h\|_Q \end{aligned}$$

It is worth noticing that condition (6.20) is equivalent to the existence of a positive constant  $\beta$  such that

$$\inf_{\substack{q_h \in Q_h \\ q_h \neq 0}} \sup_{\substack{\mathbf{v}_h \in V_h \\ \mathbf{v}_h \neq 0}} \frac{b(\mathbf{v}_h, q_h)}{\|\mathbf{v}_h\|_{\mathbf{H}^1(\Omega)} \|q_h\|_{L^2(\Omega)}} \geq \beta \quad (6.21)$$

Which is often called inf-sup condition.

## 6.8 A general saddle-point problem

Let  $X$  and  $M$  be two Hilbert spaces endowed with norms  $\|\cdot\|_X$  and  $\|\cdot\|_M$ . Denoting their dual spaces as  $X', M'$ , we introduce the bilinear forms

- $a(\cdot, \cdot) : X \times X \rightarrow \mathbb{R}$
- $b(\cdot, \cdot) : M \times M \rightarrow \mathbb{R}$

that we suppose to be continuous, meaning there exist two constants  $\gamma, \delta > 0$  such that for all  $w, v \in X$  and  $\mu \in M$

$$\begin{aligned} |a(w, v)| &\leq \gamma \|w\|_X \|v\|_X, \\ |b(w, \mu)| &\leq \|w\|_X \|\mu\|_M \end{aligned} \quad (6.22)$$

Consider now the following constrained problem

$$\text{find } (u, \eta) \in X \times M : \begin{cases} a(u, v) + b(v, \eta) = \langle I, v \rangle & \forall v \in X, \\ b(u, \mu) = \langle \sigma, \mu \rangle & \forall \mu \in M \end{cases} \quad (6.23)$$

where  $I \in X'$  and  $\sigma \in M'$  are two assigned linear functionals, while  $\langle \cdot, \cdot \rangle$  denotes the pairing between  $X$  and  $X'$  or  $M$  and  $M'$ .

Formulation (6.23) is general enough to include formulation (6.16) of the Stokes problem. In order to analyze (6.23), we introduce the affine manifold

$$X^\sigma = \{v \in X : b(v, \mu) = \langle \sigma, \mu \rangle \forall \mu \in M\} \quad (6.24)$$

The space  $X^0$  denotes the kernel of  $b$ , that is

$$X^0 = \{v \in X : b(v, \mu) = 0 \ \forall \mu \in M\}$$

This is a closed subspace of  $X$ . We can, therefore, associate (6.23) with the following reduced problem:

$$\text{find } u \in X^\sigma : a(u, v) = \langle I, v \rangle \quad \forall v \in X^0 \quad (6.25)$$

If  $(u, \eta)$  is a solution of (6.23), then it is a solution to (6.25). In the following, we will introduce suitable conditions that allow the converse to hold too. Also, if we are able to prove existence and uniqueness for (6.25), then this would allow us to obtain a result for (6.23).

**Theorem 6.4** (Existence, uniqueness and stability)

Let the bilinear form  $a(\cdot, \cdot)$  satisfy the continuity condition (6.22) and be coercive on  $X^0$ , so

$$\exists \alpha > 0 : a(v, v) \geq \alpha \|v\|_X^2 \quad \forall v \in X^0 \quad (6.26)$$

Suppose also that the bilinear form  $b(\cdot, \cdot)$  satisfies the same continuity condition as well as the following: there exists  $\beta^* > 0$  such that

$$\forall \mu \in M \exists v \in X, v \neq 0 : b(v, \mu) \geq \beta^* \|v\|_X \|\mu\|_M. \quad (6.27)$$

Then, for every  $I \in X'$  and  $\sigma \in M'$ , there exists a unique solution  $(u, \eta) \in X \times M$  to the saddle point problem (6.23).

Moreover, the map  $(I, \sigma) \mapsto (u, \eta)$  is an isomorphism from  $X' \times M'$  onto  $X \times M$  and the following a priori estimates hold

$$\begin{aligned} \|u\|_X &\leq \frac{1}{\alpha} \left[ \|I\|_{X'} + \frac{\alpha + \gamma}{\beta^*} \|\sigma\|_{M'} \right], \\ \|\eta\|_M &\leq \frac{1}{\beta} \left[ \left( 1 + \frac{\gamma}{\alpha} \|I\|_{X'} + \frac{\gamma(\alpha + \gamma)}{\alpha \beta^*} \|\sigma\|_{M'} \right) \right]. \end{aligned} \quad (6.28)$$

### Galerkin approximation

To introduce a Galerkin approximation of the saddle-point problem (6.23), we consider two subspaces  $X_h \subset X$  and  $M_h \subset M$ , respectively. They can be either finite element piecewise polynomial subspaces or spectral element subspaces. We look for a solution for the problem

$$\text{given } I \in X' \text{ and } \sigma \in M' \text{ find } (u_h, \eta_h) \in X_h \times M_h : \begin{cases} a(u_h, v_h) + b(v_h, \eta_h) = \langle I, v_h \rangle & \forall v_h \in X_h \\ b(u_h, \mu_h) = \langle \sigma, \mu_h \rangle & \forall \mu_h \in M_h \end{cases} \quad (6.29)$$

We introduce the subspace

$$X_h^\sigma = \{v_h \in X_h : b(v_h, \mu_h) = \langle \sigma, \mu_h \rangle \ \forall \mu_h \in M_h\} \quad (6.30)$$

which allow us to introduce the following reduced formulation

$$\text{find } u_h \in X_h^\sigma : a(u_h, v_h) = \langle I, v_h \rangle \quad \forall v_h \in X_h^0 \quad (6.31)$$

Since, in general,  $M_h$  is different from  $M$  the space (6.30) is not necessarily a subspace of  $X^\sigma$ .

Clearly, every solution of  $(u_h, \eta_h)$  of (6.29) yields a solution  $u_h$  for the reduced problem (6.31). Now we look for a solution that allows us to prove that the converse is true.

**Theorem 6.5** (Existence, uniqueness and stability)

Let the bilinear form  $a(\cdot, \cdot)$  satisfy the continuity condition (6.22) and be coercive on  $X^0$ , so

$$\exists \alpha > 0 : a(v_h, v_h) \geq \alpha_h \|v_h\|_X^2 \quad \forall v_h \in X_h^0 \quad (6.32)$$

Suppose also that the bilinear form  $b(\cdot, \cdot)$  satisfies the same continuity condition as well as the following: there exists  $\beta_h^* > 0$  such that

$$\forall \mu_h \in M_h \exists v_h \in X_h, v_h \neq 0 : b(v_h, \mu_h) \geq \beta_h^* \|v_h\|_X \|\mu_h\|_M. \quad (6.33)$$

Then, for every  $I \in X'$  and  $\sigma \in M'$ , there exists a unique solution  $(u_h, \eta_h) \in X_h \times M_h$  to the saddle point problem (6.29).

Moreover, the solution satisfies the following conditions:

$$\|u_h\|_X \leq \frac{1}{\alpha_h} \left[ \|I\|_{X'} + \frac{\alpha_h + \gamma}{\beta_h^*} \|\sigma\|_{M'} \right], \quad (6.34)$$

$$\|\eta_h\|_M \leq \frac{1}{\beta_h} \left[ \left( 1 + \frac{\gamma}{\alpha_h} \|I\|_{X'} + \frac{\gamma(\alpha_h + \gamma)}{\alpha_h \beta_h} \|\sigma\|_{M'} \right) \right]. \quad (6.35)$$

The coercivity condition (6.26) does not necessarily guarantees (6.32), as  $X_h^0 \not\subset X^0$ , nor does the compatibility condition (6.27) in general imply the discrete compatibility condition (6.33) due to the fact that  $X_h$  is a proper subspace of  $X$ . (6.33) is called inf-sup condition.

### Theorem 6.6

Let the assumptions of existence and uniqueness theorems (6.4) and (6.5) be satisfied. Then, the solutions  $(u, \eta)$  and  $(u_h, \eta_h)$  of problems (6.23) and (6.29), respectively, satisfy the following error estimates:

$$\|u - u_h\|_X \leq \left( 1 + \frac{\gamma}{\alpha_h} \right) \int_{v_h^* \in X_h^\sigma} \|u - v_h\|_X + \frac{\delta}{\alpha_h} \inf_{\mu_h \in M_h} \|\eta - \mu_h\|_M \quad (6.36)$$

$$\|\eta - \eta_h\|_M \leq \frac{\gamma}{\beta_h} \left( 1 + \frac{\gamma}{\alpha_h} \right) \inf_{v_h^* \in X_h^\sigma} \|u - v_h^*\|_X + \left( 1 + \frac{\delta}{\beta_h} + \frac{\gamma\beta}{\alpha_h \beta_h} \right) \inf_{\mu_h \in M_h} \|\eta - \mu_h\|_M \quad (6.37)$$

where  $\gamma, \delta, \alpha_h$  and  $\beta_h$  are respectively defined by (6.22), (6.32) and (6.33). Also, the following estimate holds

$$\inf_{v_h^* \in X_h^\sigma} \|u - v_h\|_X \leq \left( 1 + \frac{\delta}{\beta_h} \right) \inf_{v_h \in X_h} \|u - v_h\|_X \quad (6.38)$$

The inequalities (6.36) and (6.37) yield error estimates with optimal convergence rate, provided that the constants  $\alpha_h$  and  $\beta_h$  in (6.32) and (6.33) are bounded from below by two constants  $\alpha$  and  $\beta$  independent of  $h$ . Let us also remark that (6.36) holds even if (6.27) and (6.33) are not satisfied.

### Remark 6.2 (Spurious pressure modes)

The compatibility condition (6.33) is essential to guarantee the uniqueness of the  $\eta_h$ -component of the solution. Indeed, if (6.33) does not hold, then

$$\exists \mu_h^* \in M_h, \mu_h^* \neq 0 : b(v_h, \mu_h^*) = 0 \quad \forall v_h \in X_h.$$

Consequently, if  $(u_h, \eta_h)$  is a solution to the problem (6.29), then  $(u_h, \eta_h + \tau \mu_h^*)$ , for all  $\tau \in \mathbb{R}$ , is a solution too.

Any such function  $\mu_h^*$  is called spurious mode.

## Algebraic form of Stokes problem

Let us investigate the structure of the algebraic system associated to the Galerkin approximation (6.19) to the Stokes problem (or, more generally, to a discrete saddle-point problem like (6.29)). Denote with

$$\{\varphi_j \in V_h\}, \quad \{\phi_k \in Q_h\},$$



the basis functions of the spaces  $V_h$  and  $Q_h$ , respectively. Let us expand the discrete solutions  $\mathbf{u}_h$  and  $p_h$  w.r.t. such bases

$$\begin{aligned}\mathbf{u}_h(\mathbf{x}) &= \sum_{j=1}^N u_j \boldsymbol{\varphi}_j(\mathbf{x}), \\ p_h(\mathbf{x}) &= \sum_{k=1}^M p_k \phi_k(\mathbf{x}),\end{aligned}\tag{6.39}$$

having set  $N = \dim V_h$  and  $M = \dim Q_h$ .

By choosing as test functions in (6.19) the same basis functions we obtain the following block linear system

$$\begin{cases} A\mathbf{U} + B^T\mathbf{P} = \mathbf{F}, \\ B\mathbf{U} = \mathbf{0} \end{cases}\tag{6.40}$$

where  $A \in \mathbb{R}^{N \times N}$  and  $B \in \mathbb{R}^{M \times N}$  are the matrices related respectively to the bilinear forms  $a(\cdot, \cdot)$  and  $b(\cdot, \cdot)$ , whose elements are given by

$$\begin{aligned}A &= [a_{ij}] = [a(\boldsymbol{\varphi}_j, \boldsymbol{\varphi}_i)], \\ B &= [b_{km}] = [b(\boldsymbol{\varphi}_m, \phi_k)],\end{aligned}$$

while  $\mathbf{U}$  and  $\mathbf{P}$  are the vectors of the unknowns,

$$\begin{aligned}\mathbf{U} &= [u_j], \\ \mathbf{P} &= [p_j].\end{aligned}$$

The  $(N + M) \times (N + M)$  matrix

$$S = \begin{bmatrix} A & B^T \\ B & 0 \end{bmatrix}\tag{6.41}$$

is block-symmetric (as  $A$  is symmetric) and indefinite, featuring real eigenvalues with variable sign.  $S$  is non-singular if and only if is null, a property that follows from the inf-sup condition (6.21).

To prove the latter statement we proceed as follows.

Since  $A$  is non-singular, because it is associated to the coercive bilinear form  $a(\cdot, \cdot)$ , from the first of (6.40), we can formally obtain  $\mathbf{U}$  as

$$\mathbf{U} = A^{-1}(\mathbf{F} - B^T\mathbf{P})\tag{6.42}$$

Using (6.42) in the second equation of (6.40) yields

$$R\mathbf{P} = BA^{-1}\mathbf{F}, \text{ where } R = BA^{-1}B^T\tag{6.43}$$

This corresponds to having carried out a block Gaussian elimination on system (6.41). In this way we obtain a reduced system for the sole unknown  $\mathbf{P}$ , which admits a unique solution in case  $R$  is non-singular and positive definite, we want to prove that the latter condition is satisfied if and only if  $B^T$  has a null kernel, that is

$$\ker(B^T) = \{\mathbf{0}\}\tag{6.44}$$

where  $\ker(B^T) = \{\mathbf{x} \in \mathbb{R}^M : B^T\mathbf{x} = \mathbf{0}\}$ .

We proceed as follows:

$$R\mathbf{p} = \mathbf{0} \Rightarrow \mathbf{p} = \mathbf{0}$$

that is

$$\langle BA^{-1}B^T\mathbf{p}, \mathbf{q} \rangle = \mathbf{0} \ \forall \ \mathbf{q} \Rightarrow \mathbf{p} = \mathbf{0}$$

Let us take  $\mathbf{q} = \mathbf{p}$ . We require:

$$\langle A^{-1}B^T\mathbf{p}, B^T\mathbf{q} \rangle = \mathbf{0} \Rightarrow \mathbf{p} = \mathbf{0}$$

Set  $\mathbf{w} = B^T\mathbf{p}$ . Since  $A$  is spd, we have that  $\langle A^{-1}\mathbf{w}, \mathbf{w} \rangle = \mathbf{0}$ , which implies  $\mathbf{w} = \mathbf{0}$ . Finally

$$(\mathbf{w} = B^T\mathbf{p} = \mathbf{0} \Rightarrow \mathbf{p} = \mathbf{0}) \iff \ker(B^T) = \{\mathbf{0}\}$$

**Remark 6.3**

Condition (6.44) is equivalent to the inf-sup condition (6.21).

On the other hand, since  $A$  is non-singular, from the existence and uniqueness  $\mathbf{P}$  we infer that there exists a unique vector  $\mathbf{U}$  which satisfies (6.42).

In conclusion, system (6.40) admits a unique solution  $(\mathbf{U}, \mathbf{P})$  if and only if condition (6.44).

**Remark 6.4**

We remark that, for an arbitrary matrix,  $B^T(N \times M)$ , we have  $\text{rank}(B^T) + \dim \ker(B^T) = \min(M, N)$ . Then, condition (6.44) is equivalent to asking that  $B^T$  has a full rank, because  $\text{rank}(B^T)$  is the maximum number of linearly independent row vectors of  $B^T$ .

Let us consider the remark (6.2) concerning the general saddle-point problem and suppose that the inf-sup condition (6.21) does not hold. Then,

$$\exists q_h^* \in Q_h : b(\mathbf{v}_h, q_h^*) = 0 \quad \forall \mathbf{v}_h \in V_h. \quad (6.45)$$

Consequently, if  $(\mathbf{u}_h, p_h)$  is a solution to the Stokes problem (6.19), then  $(\mathbf{u}_h, p_h + q_h)$  is a solution too:

$$\begin{aligned} a(\mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h + q_h) &= a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}_h, p_h) + \underbrace{b(\mathbf{v}_h, q_h^*)}_{=0} \\ &= a(\mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h) = (\mathbf{f}, \mathbf{v}_h) \end{aligned} \quad \forall \mathbf{v}_h \in V_h$$

Functions  $q_h^*$  which fail to satisfy the inf-sup condition are invisible to the Galerkin problem (6.19). For this reason, as already observed, they are called spurious pressure modes, or parasitic modes. Their presence yields numerical instabilities.

Two strategies are generally adopted in order to guarantee well-posedness of the numerical problem:

- choose spaces  $V_h$  and  $Q_h$  that satisfy the inf-sup condition.
- Stabilize the finite dimensional problem by eliminating spurious modes.

**6.9 Suitable spaces**

Let us analyze the first type of strategy. To start with, we will consider the case of finite element spaces. To characterize  $Q_h$  and  $V_h$  it suffices to choose on every element their degrees of freedom. Since the weak formulation does not require a continuous pressure, we will consider the case of a discontinuous pressure first.

As Stokes equation are of order one in  $p$  and two in  $\mathbf{u}$  it makes sense to use polynomials of degree  $k \geq 1$  for  $V_h$  and of degree  $k - 1$  for  $Q_h$ . When looking for a compatible couple of spaces, the larger the velocity space  $V_h$ , the more likely the inf-sup condition is satisfied.

Suitable choices of spaces that fulfill the inf sup pressure. are, for example, the couple  $(\mathbb{P}_2, \mathbb{P}_0)$ ,  $((\mathbb{Q}_2, \mathbb{P}_0))$  or the so-called Crouzeix Raviart elements, in which the velocity components are piecewise quadratic functions with a bubble function on each element, while the pressure components are piecewise linear discontinuous element.

In the case of continuous pressure, examples of incompatible spaces could be piecewise linear elements on triangles for both velocity and pressure. More in general, not a good idea to use the same polynomial degree.

The smallest degree possible for stability is the pair  $(\mathbb{P}_2, \mathbb{P}_1)$ , which are called Taylor-Hood elements.

**Spectral methods**

If we use spectral methods using equal-order polynomial spaces for both velocity and pressure yields subspaces that violate the inf-sup condition. Compatible spaces can be the ones with polynomial degree  $N \geq 2$  for the velocity and  $N - 2$  for the pressure, yielding the so called  $(\mathbb{Q}_N, \mathbb{Q}_{N-2})$  approximation. The degrees of freedom for each velocity components are represented by the  $(N + 1)^2$  LGL nodes.

In the case of discontinuous pressure, at least two sets of interpolation nodes can be used: either the subset represented by the  $(N - 1)^2$  internal nodes of the choice above, or the set  $(N - 1)^2$  LGL nodes. This choice stands at the base of the spectral-type approximation, such as collocation, G-NI or SEM-NI.

## 6.10 A stabilized problem

We have seen that finite element or spectral methods that make use of equal-degree polynomial do not fullfill the inf-sup condition and are therefore unstable. However, stabilizing them is possible by using SUPG or GLS techniques like those enconuntered in the approximation of ADR equation. Here we will limit ourselves to the case of piecewise continuous linear finite elements  $(\mathbb{P}_1, \mathbb{P}_1)$ , stabilized using GLS.

$$V_h = [\overset{\circ 1}{X}_h]^2, \quad Q_h = \left\{ q_h \in X_h^1 : \int_{\Omega} d\Omega = 0 \right\}.$$

This choice is urged by the need of keeping the global numer of degrees of freedom as low as possible, especially when dealing with three-dimensional problem. However, as it violates the inf-sup condition, it will be unstable.

So, we set  $W_h = V_h \times Q_h$ , and, instead of (6.19), we consider the problem

$$\text{find } (\mathbf{u}_h, p_h) \in W_h : A_h((\mathbf{u}_h, p_h), (\mathbf{v}_h, q_h)) = F_h(\mathbf{v}_h, q_h) \quad \forall (\mathbf{v}_h, q_h) \in W_h \quad (6.46)$$

We have set

$$\begin{aligned} A_h : W_h \times W_h &\rightarrow \mathbb{R}, \\ A_h((\mathbf{u}_h, p_h), (\mathbf{v}_h, q_h)) &= a(\mathbf{u}_h, \mathbf{v}_h) + b(\mathbf{v}_h, p_h) - b(\mathbf{u}_h, q_h) \\ &\quad + \delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \int_{\mathcal{K}} (-\nu \Delta \mathbf{u}_h + \nabla p_h) (-\nu \Delta \mathbf{v}_h + \nabla q_h) d\mathcal{K}, \\ F_h : W_h &\rightarrow \mathbb{R}, \\ F_h(\mathbf{v}_h, q_h) &= (\mathbf{f}, \mathbf{v}_h) + \delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \int_{\mathcal{K}} \mathbf{f} (-\nu \Delta \mathbf{v}_h + \nabla q_h) d\mathcal{K} \end{aligned}$$

and  $\delta$  is a positive parameter chosen accordingly.

This is a strongly consistent approximation of problem (6.12). As a matter of fact, the additional term is null when calculated on the exact solution thanks to (6.14). Note that, since  $k = 1$ ,  $\Delta \mathbf{u}_h|_{\mathcal{K}} = \Delta \mathbf{v}_h|_{\mathcal{K}} = \mathbf{0} \forall \mathcal{K} \in \mathcal{T}_h$  as we are using piecewise linear finite elements functions.

We obtain the following stability inequality

$$\nu \|\nabla \mathbf{u}_h\|_{\mathbf{L}^2(\Omega)}^2 + \delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \|\nabla p_h\|_{\mathbf{L}^2(\mathcal{K})}^2 \leq C \|\mathbf{f}\|_{\mathbf{L}^2(\Omega)}^2, \quad (6.47)$$

$C$  being a constant that depends on  $\nu$  but not on  $h$ .

By applying Strang's Lemma (2.1) we can now show that the solution to the generalized Galerkin problem (6.46) satisfies the following estimate

$$\|\mathbf{u} - \mathbf{u}_h\|_{\mathbf{H}^1(\Omega)} + \left( \delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \|\nabla p - \nabla p_h\|_{\mathbf{L}^2(\mathcal{K})}^2 \right)^{\frac{1}{2}} \leq Ch.$$

We can show that (6.46) admits the following matrix form

$$\begin{bmatrix} A & B^T \\ B & -C \end{bmatrix} \begin{bmatrix} \mathbf{U} \\ \mathbf{P} \end{bmatrix} = \begin{bmatrix} \mathbf{F} \\ \mathbf{G} \end{bmatrix}. \quad (6.48)$$

This system differs from (6.40) without stabilization because of the presence of the non-null block occupying the position (2, 2) which is associated to the stabilization term:

$$C = [c_{km}], c_{km} = \delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \int_{\mathcal{K}} \nabla \phi_m \cdot \nabla \phi_k d\mathcal{K}, k, m = 1, \dots, M$$

while the components of the right hand side  $\mathbf{G}$  are

$$g_k = -\delta \sum_{\mathcal{K} \in \mathcal{T}_h} h_{\mathcal{K}}^2 \int_{\mathcal{K}} \mathbf{f} \cdot \nabla \phi_k d\mathcal{K}, k = 1, \dots, M$$

Note that num

$$A\mathbf{U} + B^T\mathbf{P} = \mathbf{F} \Rightarrow \mathbf{U} = A^{-1}(\mathbf{F} - B^T\mathbf{P}), \quad (6.49)$$

$$B\mathbf{U} - C\mathbf{P} \Rightarrow BA^{-1}\mathbf{F} - BA^{-1}B^T\mathbf{P} - C\mathbf{P} = \mathbf{G} \quad (6.50)$$

$$\Rightarrow (BA^{-1}B^T + C)\mathbf{P} = BA^{-1}\mathbf{F} - \mathbf{G}. \quad (6.51)$$

Then, in this case, the reduced system reads

$$R\mathbf{P} = BA^{-1}\mathbf{F} - \mathbf{G}.$$

In contrast to (6.43), this time the matrix  $R$  is non-singular as  $C$  is positive defined., but  $B^T$  is not full-rank.

### 6.11 Time discretization

We consider the following semi discretized formulation

$$\begin{cases} M \frac{d\mathbf{u}(t)}{dt} + A\mathbf{u}(t) + C(\mathbf{u}(t))\mathbf{u}(t) + B^T\mathbf{p}(t) = \mathbf{f}(t), \\ B\mathbf{u}(t) = \mathbf{0}, \end{cases} \quad (6.52)$$

with  $\mathbf{u}(0) = \mathbf{u}_0$ .  $C(\mathbf{u}(t))$  is in fact a matrix depending on  $\mathbf{u}(t)$ , whose generic coefficient is  $c_{ij}(t) = c(\mathbf{u}(t), \varphi_j, \varphi_i)$ .

For the temporal discretization of this system let us use the  $\theta$ -method. By setting

$$\begin{aligned} \mathbf{u}_{\theta}^{n+1} &= \theta \mathbf{u}^{n+1} + (1 - \theta) \mathbf{u}^n, \\ \mathbf{p}_{\theta}^{n+1} &= \theta \mathbf{p}^{n+1} + (1 - \theta) \mathbf{p}^n, \\ \mathbf{f}_{\theta}^{n+1} &= \theta \mathbf{f}^{n+1} + (1 - \theta) \mathbf{f}^n, \\ C_{\theta}(\mathbf{u}^{n+1,n}) \mathbf{u}^{n+1,n} &= \theta C(\mathbf{u}^{n+1}) \mathbf{u}^{n+1} + (1 - \theta) C(\mathbf{u}^n) \mathbf{u}^n \end{aligned}$$

we obtain

$$\begin{cases} M \frac{\mathbf{u}^{n+1} - \mathbf{u}^n}{\Delta t} + A\mathbf{u}_{\theta}^{n+1} + C_{\theta}(\mathbf{u}^{n+1,n}) \mathbf{u}^{n+1,n} + B^T \mathbf{p}_{\theta}^{n+1} = \mathbf{f}_{\theta}^{n+1} \\ B\mathbf{u}^{n+1} = \mathbf{0} \end{cases} \quad (6.53)$$

Since, except for the forward Euler method, the solution of this system is quite involved, a possible alternative is a semi-implicit scheme in which the linear part of the equation is advanced implicitly, while the nonlinear ones explicitly. in this way, for  $\theta \geq \frac{1}{2}$ , the resulting scheme is unconditionally stable, whereas, in all other cases, it must obey a stability restriction on the timestep  $\Delta t$ .

### Finite difference methods

We consider an explicit temporal discretization of the first equation in (6.52), corresponding to  $\theta = 0$  in (6.53). If all the quantities are known at  $t^n$ , then we can write

$$\begin{cases} M\mathbf{u}^{n+1} = H(\mathbf{u}^n, \mathbf{p}^n, \mathbf{f}^n) \\ B\mathbf{u}^{n+1} = \mathbf{0} \end{cases}$$

with  $M$  mass matrix

$$m_{ij} = \int_{\Omega} \varphi_i \varphi_j$$

This system does not allow the determination of the pressure  $\mathbf{p}^{n+1}$ . In particular, there is no way to enforce the divergence-free constraints on  $\mathbf{u}^{n+1}$ . However if we replace  $\mathbf{p}^n$  by  $\mathbf{p}^{n+1}$  we obtain

$$\begin{cases} \frac{1}{\Delta t} M \mathbf{u}^{n+1} + B^T \mathbf{p}^{n+1} = \mathbf{G}, \\ B \mathbf{u}^{n+1} = \mathbf{0}, \end{cases} \quad (6.54)$$

with  $\mathbf{G}$  being a suitable vector. This system corresponds to a semi-explicit discretization of (6.52). Since  $M$  is symmetric and positive definite, if condition (6.44) is satisfied, then the reduced system  $BM^{-1}B^T \mathbf{p}^{n+1} = BM^{-1}\mathbf{G}$  is non-singular. Once solved, the velocity vector  $\mathbf{u}^{n+1}$  can be recovered from the first equation of (6.54). This discretization method is temporally stable if

$$\Delta t \leq C \min \left( \frac{h^2}{\nu}, \frac{h}{\max_{\mathbf{x} \in \Omega} |\mathbf{u}^n(\mathbf{x})|} \right).$$

Let us now consider an implicit discretization of (6.52), for instance the backward Euler method. As already seen, this is an unconditionally stable method. Its algebraic system can be regarded as the finite element space approximation of the Navier-Stokes problem

$$\begin{cases} -\nu \Delta \mathbf{u}^{n+1} + (\mathbf{u}^{n+1} \cdot \nabla) \mathbf{u}^{n+1} + \nabla p^{n+1} + \frac{\mathbf{u}^{n+1}}{\Delta t} = \tilde{\mathbf{f}}, \\ \operatorname{div} \mathbf{u}^{n+1} = 0. \end{cases}$$

The solution of such nonlinear algebraic system can be achieved by Newton-Krylov techniques, that is using a Krylov method to solve, at each timestep, a system obtained by a Newton method.

### Fractional step methods

Let us consider an abstract time dependent problem,

$$\frac{\partial w}{\partial t} + \mathcal{L}w = f,$$

where  $\mathcal{L}$  is a differential operator that splits into the sum of two operators  $\mathcal{L}_1$  and  $\mathcal{L}_2$  such that

$$\mathcal{L}v = \mathcal{L}_1 v + \mathcal{L}_2 v$$

We perform an advancement in time only on the first operator, and then correct the solution by performing the advancement on the second operator. The idea is to separate a complex problem into smaller simpler ones, for example separating diffusion from transport.

#### 6.11.1 Chorin Temam method

This strategy also involves two operators:  $\mathcal{L}_1(\mathbf{w}) = -\nu \Delta \mathbf{w} + (\mathbf{w} \cdot \nabla) \mathbf{w}$  whereas  $\mathcal{L}_2$  is associated to the remaining terms of (6.2). Thanks to this solution we can split the main difficulty of the Navier-Stokes equations, the nonlinear part from the incompressibility constraint.

The corresponding scheme reads

- (1) Solve the diffusion-transport equation for the velocity  $\tilde{\mathbf{u}}^{n+1}$

$$\begin{cases} \frac{\tilde{\mathbf{u}}^{n+1} - \mathbf{u}^n}{\Delta t} - \nu \Delta \tilde{\mathbf{u}}^{n+1} + (\mathbf{u}^* \cdot \nabla) \mathbf{u}^{**} = \mathbf{f}^{n+1} & \text{in } \Omega, \\ \tilde{\mathbf{u}}^{n+1} = \mathbf{0} & \text{on } \partial\Omega \end{cases} \quad (6.55)$$

- (2) solve the following problem for  $\mathbf{u}^{n+1}$  and  $p^{n+1}$

$$\begin{cases} \frac{\tilde{\mathbf{u}}^{n+1} - \mathbf{u}^n}{\Delta t} - \nu \Delta \tilde{\mathbf{u}}^{n+1} + (\mathbf{u}^* \cdot \nabla) \mathbf{u}^{**} = \mathbf{f}^{n+1} & \text{in } \Omega, \\ \operatorname{div} \mathbf{u}^{n+1} = 0 & \text{in } \Omega, \\ \mathbf{u}^{n+1} \cdot \mathbf{n} = 0 & \text{on } \partial\Omega, \end{cases} \quad (6.56)$$

where  $\mathbf{u}^*$  and  $\mathbf{u}^{**}$  can be either  $\tilde{\mathbf{u}}^{n+1}$  or  $\mathbf{u}^n$  depending on whether the nonlinear terms are treated. In such a way, in the first step the velocity is calculated as  $\tilde{\mathbf{u}}^{n+1}$  and then it is corrected in the second step to satisfy the constraints.

While the first step is a classic advection-diffusion problem, the second step is a bit more tricky. First we need to apply the divergence operator to the first equation of (6.56)

$$\operatorname{div} \frac{\mathbf{u}^{n+1}}{\Delta t} - \operatorname{div} \frac{\tilde{\mathbf{u}}^{n+1}}{\Delta t} + \Delta p^{n+1} = 0,$$

that is an elliptic boundary value problem with Neumann boundary conditions

$$\begin{cases} -\Delta p^{n+1} = -\operatorname{div} \frac{\tilde{\mathbf{u}}^{n+1}}{\Delta t} & \text{in } \Omega, \\ \frac{\partial p^{n+1}}{\partial n} = 0 & \text{on } \partial\Omega. \end{cases} \quad (6.57)$$

The Neumann condition follows from  $\mathbf{u}^{n+1} \cdot \mathbf{n} = 0$  on  $\partial\Omega$ .

From the solution of (6.57) we obtain  $p^{n+1}$  and  $\mathbf{u}^{n+1}$  by using the first equation on (6.56),

$$\mathbf{u}^{n+1} = \tilde{\mathbf{u}}^{n+1} - \Delta t \nabla p^{n+1} \quad (6.58)$$

This is precisely the correction to operate on the velocity field in order to fullfill the divergence-free constraint.

To conclude, the full algorithm reads:

- (1) Solve the elliptic system (6.55) to obtain  $\tilde{\mathbf{u}}^{n+1}$

$$\begin{cases} \frac{\tilde{\mathbf{u}}^{n+1} - \mathbf{u}^n}{\Delta t} - \nu \Delta \tilde{\mathbf{u}}^{n+1} + (\mathbf{u}^* \cdot \nabla) \mathbf{u}^{**} = \mathbf{f}^{n+1} & \text{in } \Omega, \\ \tilde{\mathbf{u}}^{n+1} = \mathbf{0} & \text{on } \partial\Omega \end{cases} \quad (6.59)$$

- (2) solve the scalar elliptic problem (6.57) to obtain  $p^{n+1}$

$$\begin{cases} -\Delta p^{n+1} = -\operatorname{div} \frac{\tilde{\mathbf{u}}^{n+1}}{\Delta t} & \text{in } \Omega, \\ \frac{\partial p^{n+1}}{\partial n} = 0 & \text{on } \partial\Omega. \end{cases} \quad (6.60)$$

- (3) obtain the corrected velocity  $\mathbf{u}^{n+1}$  thanks to the correction equation (6.58)

$$\mathbf{u}^{n+1} = \tilde{\mathbf{u}}^{n+1} - \Delta t \nabla p^{n+1} \quad (6.61)$$

As with many explicit schemes, if we take  $\mathbf{u}^* = \mathbf{u}^{**} = \mathbf{u}^n$  we obtain a method that has a stability restriction on the timestep like

$$\Delta t \leq C \frac{h}{\max_{\mathbf{x} \in \Omega} |\mathbf{u}^n(\mathbf{x})|}.$$

On the other hand, this system splits into independent systems of smaller size, one for each spatial component of the velocity fields.

If we decided to use  $\mathbf{u}^* = \mathbf{u}^{**} = \mathbf{u}^{n+1}$  we obtain an unconditionally stable method, but the downside is that the spatial components aren't separated due to the nonlinear convective term.

## 7 Hyperbolic equations

### 7.1 Introduction

Let us consider the following scalar hyperbolic problem

$$\begin{cases} \frac{\partial u}{\partial t} + a \frac{\partial u}{\partial x} = 0 & x \in \mathbb{R}, t > 0, \\ u(x, 0) = u_0(x), & x \in \mathbb{R}, \end{cases} \quad (7.1)$$

where  $a \in \mathbb{R} \setminus \{0\}$ . The solution of such problem is a wave travelling at a velocity  $a$ , in the  $(x, t)$  plane, given by

$$u(x, t) = u_0(x - at), \quad t \geq 0.$$

We consider the curves  $x(t)$  in the plane  $(x, t)$ , solutions of the following ordinary differential equation

$$\begin{cases} \frac{dx}{dt} = a, & t > 0, x(0) = x_0, \end{cases}$$

for varying values of  $x_0 \in \mathbb{R}$ . They read  $x(t) = x_0 + at$  and are called characteristic lines