

# Numerical Analysis for Partial Differential Equations

Andrea Bonifacio

May 24, 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

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

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

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

**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

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

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

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

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

The problem is always:

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

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

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

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

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

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

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

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):

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

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

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

Since  $A$  is positive defined  $\Rightarrow \text{Re}(\lambda_k(A)) \Rightarrow \lambda_k(A) \neq 0$ . 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

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

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

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

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

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

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

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

## 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.

$$\psi_0(x) = \frac{1}{2}(L_0(x) - L_1(x)) = \frac{1-x}{2} \quad (2.3)$$

$$\psi_N(x) = \frac{1}{2}(L_0(x) + L_1(x)) = \frac{1+x}{2} \quad (2.4)$$

$$\psi_{k-1}(x) = \frac{1}{\sqrt{2(2k-1)}}(L_{k-2}(x) - L_k(x)) \quad (2.5)$$

$$\text{for } k = 2, \dots, N-1 < x < 1 \quad (2.6)$$

$$(2.7)$$

## 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  $V_N = \mathbb{P}_N^0$

$$\text{find } u_N \in V_N : a(u_N, v_N) = (f, v_N)_{L^2(\Omega)} \quad (2.8)$$

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  $x_j =$  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  $a_N(u_N, v_N) = (\mu u'_N - b u_n, v'_N)_N + (\sigma u_n, v_n)_N$ , 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  $i = 1, \dots, N-1$ :

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

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

$$A_{ij}^{\text{GNI}} = a_N(\varphi_j, \varphi_i) = \underbrace{(\mu \varphi_j', \varphi_i')_N}_A - \underbrace{(b \varphi_j, \varphi_i')_N}_B + \underbrace{(\sigma \varphi_j, \varphi_i)_N}_C$$

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

$$\text{find } u \in V : a(u, v) = F(v) \quad \forall v \in V \quad (2.9)$$

and its approximation

$$\text{find } u_h \in V_h : a_h(u_h, v_h) = F_h(v_h) \quad \forall v_h \in V_h \quad (2.10)$$

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 (2.10), 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:

$$\begin{aligned} \alpha^* \|\sigma_h\|_V^2 &\leq a_h(\sigma_h, \sigma_h) && \text{(by coercivity of } a_h) \\ &= a_h(u_h, \sigma_h) - a_h(w_h, \sigma_h) \\ &= F_h(\sigma_h) - a_h(w_h, \sigma_h) && \text{(by (2.10))} \\ &= F_h(\sigma_h) - F(\sigma_h) + F(\sigma_h) - a_h(w_h, \sigma_h) \\ &= [F_h(\sigma_h) - F(\sigma_h)] + a(u, \sigma_h) - a_h(w_h, \sigma_h) && \text{(by (2.9))} \\ &= [F_h(\sigma_h) - F(\sigma_h)] + a(u - w_h, \sigma_h) + [a(w_h, \sigma_h) - a_h(w_h, \sigma_h)] \end{aligned}$$

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

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

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

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

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

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

So it's obvious that  $(\text{GNI}) \iff (L_N u_N, v_N)_N = F_N(v_N) \forall v_N \in V_N$ , 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:  $\exists c > 0 : \forall v \in H^{s+1}(\Omega), s \geq 1$

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

## 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

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

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:

$$\text{find } u_{\delta}^{\text{GNI}} \in V_{\delta} : a_{\delta}(u_{\delta}^{\text{GNI}}, v_{\delta}) = (f, v_{\delta})_{\delta, \Omega} \quad \forall v_{\delta} \in V_{\delta} \quad (2.11)$$

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:

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

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

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

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

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)

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