

Numerical Analysis for Partial Differential Equations

Andrea Bonifacio

May 5, 2023

1 Boundary Value Problems

1.1 Weak Formulation

Let's consider a problem

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

- Ω : open bounded domain in \mathbb{R}^d , with $d = 2, 3$
- $\partial\Omega$: boundary of Ω
- f : given
- B.C. accordingly to \mathcal{L}
- \mathcal{L} : 2nd 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 Ω . Then, given a test function v , we multiply the equation by v , and integrate on the domain Ω

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

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 α 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 Ω 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 ρ_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:

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

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

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

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 δ_{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 ± 1 . This is because the Legendre basis is not suited to impose Dirichlet B.C.

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

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

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

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

$$(2.7)$$

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:

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