
TABLE of CONTENTS

1. Introduction to PDE	pag 1.
2. Parametric variational problem	pag 3.
3. Full order Numerical model	pag 4.
4. Low fidelity Model	pag 6.
5. Greedy strategy	pag 9.
6. Different approach to compute the coercivity function	pag 11.
7. Intro to neural networks and Deep learning	pag 14.
8. How to compute the weights	pag 17.
9. Convolutional neural network	pag 19.
10. Physics informed neural network	pag 23.
11. POD - neural network	pag 25.
12. Stokes Equations in MOR	pag 27.
13. Non linear problems	pag 31.
14. Navier Stokes Equations	pag 33.
15. Not Affine problem	pag 35.
16. Geometric Parametrization	pag 37.
17. Optimal Control	pag 38.

1 INTRODUCTION TO PDE

Notation. Consider $x \in \mathbb{R}^d$, $\underline{x} := (x_1, \dots, x_d)$, $A \in \mathbb{R}^{m,n}$, V an Hilbert space and $u, v \in V$. The scalar product is $\langle u, v \rangle_V$ and the norm $\|u\|_V^2 = \langle u, u \rangle_V$. The dual space of V is $V' := \{F: V \rightarrow \mathbb{R} \text{ linear and continuous}\}$ and $\langle F, u \rangle_V = F(u) \in \mathbb{R}$.

vectorial space with scalar product & complete

Example. Consider

- $V = L^2(\Omega)$, with $\Omega \subset \mathbb{R}^d \Rightarrow \langle u, v \rangle_V = \int_{\Omega} u \cdot v$ and $\|u\|_V^2 = \int_{\Omega} u^2 = \langle u, u \rangle_V$
- $V = H^1(\Omega) \Rightarrow \langle u, v \rangle_V = \int_{\Omega} u \cdot v + \int_{\Omega} \nabla u \cdot \nabla v$ and $\|u\|_V^2 = \int_{\Omega} u^2 + \int_{\Omega} \nabla u \cdot \nabla u$

The idea is to reduce the order of partial differential equations so that finding an approximate solution (whose distance from the real one is theoretically controlled) is much easier for the reduced problem

Example. Second order partial differential equations in \mathbb{R} - PDE -

$$\begin{cases} -\Delta u(\underline{x}) = f(\underline{x}) & \text{on } \Omega \subset \mathbb{R}^d \Rightarrow \text{required } u \in C^2, \text{ twice differentiable} \\ u(\underline{x}) = 0 & \text{on } \partial\Omega \quad \text{where } \Delta u(\underline{x}) = \frac{\partial^2 u}{\partial x^2} \end{cases}$$

The solution is a function $u: \mathbb{R}^d \rightarrow \mathbb{R}$ and $f: \mathbb{R}^d \rightarrow \mathbb{R}$

We reformulate $-\Delta u = f$ as : search $u \in V$ such that

$$-\int_{\Omega} \Delta u \cdot v = \int_{\Omega} f \cdot v \quad \forall v \in V \quad \text{weak formulation}$$

$$\int_{\Omega} \nabla u \cdot \nabla v - \int_{\Omega} \nabla v \cdot \nabla u = \int_{\Omega} f \cdot v \quad \downarrow \text{by parts} \quad \downarrow \text{since less hypothesis are needed to find } u$$

We don't want to solve it directly
(needed too many hypothesis)
So we change its formulation so that the
solutions requires less hypothesis

Then we want to find $u \in H^1(\Omega)$ s.t. $\forall v \in H^1(\Omega)$, $f \in L^2(\Omega)$, so in general we introduce a hypothesis to integrate f

$$\Rightarrow a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v - \int_{\Omega} \nabla v \cdot \nabla u \quad a: V \times V \rightarrow \mathbb{R} \quad \Rightarrow \text{We need } u \text{ only once differentiable}$$

Hence, the problem becomes: finding $u \in V$ s.t. $a(u, v) = F(v) \quad \forall v \in V$ that is the **variational problem VP**

$$a(u, v) = F(v) \quad V \neq \mathbb{R}$$

Theorem - Lax Milgram

If a is bilinear and continuous ($\exists c \in \mathbb{R} < \infty$ s.t. $\forall u, v \in V$ $|a(u, v)| \leq c \|u\|_V \cdot \|v\|_V$) and a is coercive ($\exists \alpha \in \mathbb{R} > 0$ s.t. $\forall v \in V$ $a(v, v) \geq \alpha \|v\|_V^2$, with $\alpha = \inf_{v \in V \setminus \{0\}} \frac{a(v, v)}{\|v\|_V^2} > 0$)

with $\beta = \sup_{v \in V \setminus \{0\}} \frac{\|F(v)\|}{\|v\|_V}$, $F: V \rightarrow \mathbb{R}$ \Rightarrow (VP) has a unique solution and $\|u\|_V \leq \frac{1}{\alpha} \|F\|_V$

\bullet VP is well-posed if $\forall F \in V'$ $\exists! u$ and u is controlled by F and $\alpha \rightarrow$ controlled by data

Proof. Assume u exists and is unique, we want to prove $\|u\|_V \leq \frac{1}{\alpha} \|F\|_V$

property of dual norm

• We know $a(u, v) = F(v) \quad \forall v \in V$. Take $v = u$: $a(u, u) = F(u) \Rightarrow \|u\|_V^2 \leq a(u, u) = F(u) \leq \|F\|_V \cdot \|u\|_V \Rightarrow \|u\|_V \leq \frac{1}{\alpha} \|F\|_V$

• To prove uniqueness, assume $u_1, u_2 \in V$ solutions of VP: $a(u_1, v) = F(v)$ and $a(u_2, v) = F(v) \quad \forall v \in V$

$\Rightarrow a(u_1 - u_2, v) = 0$ with $w = u_1 - u_2 \quad \forall v \in V$ is VP, with solution w

$$\Rightarrow \|w\|_V \leq \frac{1}{\alpha} \|F\|_V = 0, \quad \|u_1 - u_2\|_V = 0 \Rightarrow u_1 = u_2 \quad \Rightarrow F(v) = F(v) = 0 \Rightarrow a(u_1 - u_2, v) = 0 \quad \square$$

both u_1, u_2 : $a(u_1, v) = F(v) = a(u_2, v) \Rightarrow a(u_1 - u_2, v) = F(v) - F(v) = 0 \Rightarrow a$ is bilinear, then $a(u_1 - u_2, v) = 0$

Observation. Taking $F_1 \neq F_2$ and u_1 such that $a(u_1, v) = F_1(v)$ and u_2 s.t. $a(u_2, v) = F_2(v) \quad \forall v \in V$

$\Rightarrow a(u_1 - u_2, v) = F_1(v) - F_2(v)$ is VP₂ and it has as solution $w = u_1 - u_2$

Moreover, we know $\|u_1 - u_2\|_V = \|w\|_V \leq \frac{1}{\alpha} \|F_1 - F_2\|_V$, in which $\frac{1}{\alpha}$ is the condition number of VP fixed $a(\cdot, \cdot)$

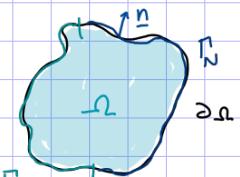
\downarrow related to coercivity of $a(\cdot, \cdot)$, the smaller α , the harder is to solve (ill-conditioned)
a small perturbation on F is reflected on a small perturbation on u if α is not too small

Example. Consider the system

$$\begin{cases} -\nabla \cdot (K(x) \nabla u(x)) = f(x) & \text{on } \Omega \subset \mathbb{R}^d \rightarrow \text{before we have } -\nabla \cdot (\nabla u), \text{ now we add } K(x) \\ u(x) = 0 & \text{on } \Gamma_D \quad \text{boundary conditions} \\ K(x) \nabla u(x) \cdot \underline{n} = \psi & \text{on } \Gamma_N = \partial\Omega \setminus \Gamma_D \end{cases}$$

\downarrow diffusion parameter

- $u(x) = 0$ on Γ_D is the **Dirichlet boundary Condition**
- $K(x) \nabla u(x) \cdot \underline{n} = \psi$ on Γ_N is the **Neumann boundary Condition**



Let's compute the weak formulation: $\forall v \in V$ find $u \in V$ such that $-\int_{\Omega} \nabla K(x) \nabla u(x) v = \int_{\Omega} f v$, $\forall v \in V$

\Rightarrow we must define V such that it satisfies one boundary condition: we assume $v = 0$ on Γ_0

$$\text{by parts } \int_{\Omega} K(x) \nabla u(x) \nabla v(x) - \int_{\Gamma_0} v K(x) \nabla u \cdot n - \int_{\Gamma_0} v K(x) \nabla u \cdot n = \int_{\Omega} f v \Rightarrow \int_{\Omega} K(x) \nabla u \nabla v = \int_{\Omega} f v + \int_{\Gamma_0} v \Psi$$

$$\|v\|_V = \|\nabla v\|_{L^2(\Omega)} + \|\nabla \Psi\|_{[L^2(\Omega)]^d}$$

$$\|v\|_V = \|\nabla v\|_{[L^2(\Omega)]^d}$$

$$= 0 \text{ since } v|_{\Gamma_0} = 0$$

$$= \Psi$$

$$\text{if } v=0 \text{ on } \Gamma_0$$

$$a(u, v)$$

$$F(v)$$

VP

NB. Dirichlet boundary condition is imposed to the whole space (strong boundary condition)

Neumann boundary condition is imposed in the VP equations (weak boundary condition)

$$\hookrightarrow V = \{v \in L^2(\Omega), \nabla v \in L^2(\Omega)\} \cup \{v = 0 \text{ on } \Gamma_0\} = H^1(\Omega) \cup \{v = 0 \text{ on } \Gamma_0\} = H_0^1(\Omega) \text{ and } f \in L^2(\Omega), \Psi \in L^2(\Omega)$$

To have $\int_{\Omega} f v < \infty$, $\int_{\Omega} \nabla f \cdot v < \infty$, $\int_{\Omega} f v < \infty$ we need V defined \downarrow sol is 0 on the Dirichlet boundary

$$\hookrightarrow K \in L^\infty(\Omega) \text{ and } \exists K_*, K^*: 0 < K_* \leq K(x) \leq K^* < \infty \quad \forall x \in \Omega$$

$$\alpha \text{ of } a(u, v) \quad \beta \text{ of } a(u, v)$$

d times

Poincaré Inequality. If $v \in H_0^1(\Omega) \Rightarrow \exists c_1 \in \mathbb{R} < \infty$ s.t. $\|v\|_{L^2(\Omega)} \leq c_1 \|\nabla v\|_{L^2(\Omega)^d}$, where $L^2(\Omega)^d = [L^2(\Omega) \times \dots \times L^2(\Omega)]^d$, then if

$$v \in H_0^1(\Omega) \Rightarrow \|v\|_{H_0^1(\Omega)} = \|\nabla v\|_{L^2(\Omega)} + \|\nabla \Psi\|_{L^2(\Omega)^d} \leq \tilde{c} \|\nabla v\|_{L^2(\Omega)^d} \text{ and } \|\nabla \Psi\|_{L^2(\Omega)^d} \leq \|v\|_{H_0^1(\Omega)} = \|v\|_{L^2(\Omega)} + \|\nabla \Psi\|_{L^2(\Omega)^d}$$

$$\Rightarrow \|v\|_{L^2(\Omega)^d} \text{ and } \|v\|_{H_0^1(\Omega)}$$
 are equivalent and $\|v\|_{H_0^1(\Omega)} \approx \|\nabla v\|_{L^2(\Omega)^d}$ if $v \in H_0^1(\Omega)$

If $v \in H_0^1(\Omega) \Rightarrow \exists c_2 \in \mathbb{R} < \infty$ s.t. $\|v\|_{L^2(\Gamma_N)} \leq c_2 \|v\|_{H_0^1(\Omega)}$. We can define the operator continuous trace operator, then:

$$\begin{aligned} \mathcal{T}: V \rightarrow W = H^{1/2}(\Gamma_N) \\ v \mapsto \mathcal{T}(v) = v|_{\Gamma_N} \end{aligned}$$

the

$$\downarrow \quad V = H_0^1(\Omega), f \in L^2(\Omega), \Psi \in L^2(\Omega) \text{ and } K \in L^\infty(\Omega) \quad (\exists K_*, K^* \text{ s.t. } 0 < K_* \leq K(x) \leq K^* < \infty)$$

The idea is to give a meaning to the value that the function has on the border of the domain

map function from a space to the border of this space

2- PARAMETRIC VARIATIONAL PROBLEMS

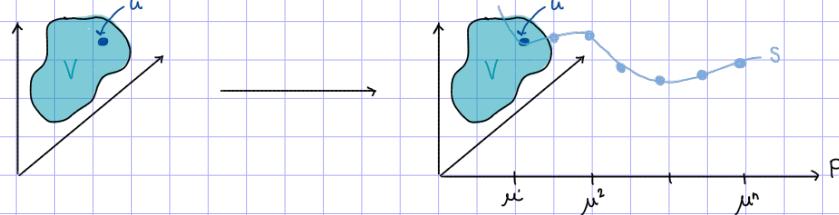
Consider $\mu \in \mathbb{R}^P$, $\mu \in P$ defined as vector of parameters $\mu = [\mu_1, \dots, \mu_P]$. Now a VP(μ) looks for $u(x; \mu) \in V$ s.t. $a(u, v; \mu) = F(v, \mu) \quad \forall v \in V, \forall \mu \in P$ where $a: V \times V \times P \rightarrow \mathbb{R}$, $F: V \times P \rightarrow \mathbb{R}$

Theorem - Lax Milgram for VP(μ) If

- I. $a(\cdot, \cdot; \mu)$ is uniformly continuous: $\exists \gamma \in \mathbb{R}, \gamma < \infty$ s.t. $\forall \mu \in P$ $|a(u, v; \mu)| \leq \gamma \|\mu\|_V \cdot \|u\|_V \cdot \|v\|_V \leq \gamma \|u\|_V \cdot \|v\|_V < \infty$
where $\gamma = \sup_{\mu, v \in V} \frac{|a(u, v; \mu)|}{\|u\|_V \cdot \|v\|_V} \leq \gamma < \infty$
- II. $a(\cdot, \cdot; \mu)$ is uniformly coercive: $\exists \alpha \in \mathbb{R}, \alpha > 0$ s.t. $\forall \mu \in P$ $a(v, v; \mu) \geq \alpha \|\mu\|_V^2 \geq \alpha \|v\|_V^2$, $\alpha \|\mu\| = \inf_{v \in V} \frac{a(v, v; \mu)}{\|v\|_V} \geq \alpha > 0$
- III. F is uniformly continuous: $\exists \rho \in \mathbb{R}, \rho < \infty$ s.t. $|F(v, \mu)| \leq \rho \|\mu\|_V \leq \rho \|v\|_V$, $\rho(\mu) = \sup_{v \in V} \frac{|F(v, \mu)|}{\|v\|_V} = \|F(\cdot, \mu)\|_V \leq \rho < \infty$
 $\Rightarrow VP(\mu)$ has a unique solution $\forall \mu \in P$

we can find a solution for every value of μ

Example. We had one solution in $V \rightarrow$ But, if we add parameters the solution manifold is $M = \{u(\mu) \in V; \mu \in P\} \subset V$



Proposition - If $a(\cdot, \cdot; \mu)$ and $F(\cdot, \mu)$ depend smoothly on $\mu \in P \Rightarrow S: P \rightarrow V$ is smooth

(*) Proof on the notes if you're interested in

Example - Consider $VP(\mu)$, $\int_{\Omega_p} K(\mu) \nabla u \cdot \nabla v = \int_{\Omega} f v + \int_{\Gamma_N} \psi v \quad \forall v \in V$ and $K(x, \mu) = K_0 + \sum_{i=1}^P K_i(x) \mu_i$, $\mu = (\mu_1, \dots, \mu_P) \in [0, 1]^P = P$
with $K_0, K_i \in L^\infty(\Omega)$ s.t. \bullet $K_0(x) + \sum_{i=1}^P |K_i(x)| \leq K_* < \infty \quad \forall x \in \Omega$
 \bullet $K_0(x) - \sum_{i=1}^P |K_i(x)| \geq K_* > 0 \quad \forall x \in \Omega$

} for the uniformity property of a

So we define a curve of solutions
 $S: P \rightarrow V$
 $\downarrow \mu \mapsto u(\mu)$
solution mapping
of $VP(\mu)$
In general non-linear

not reduced

3- FULL ORDER NUMERICAL MODEL

NOT computable

we select a finite subspace of V, V_δ , and we create a VP problem on V_δ

$VP(\mu)$ can not be solved automatically for all parameter \Rightarrow we solve it numerically. To do so, we must discretize the space: $\langle\langle \cdot \rangle\rangle$ is the discretization parameter for the discrete problem

related to the finite dimension of V_δ

Family of subspaces of V

$VP_\delta(\mu)$: find $u_\delta \in V_\delta \subseteq V$ st $\forall \mu \in P, a(u_\delta, v_\delta; \mu) = F(v_\delta, \mu) \quad \forall v_\delta \in V_\delta \rightarrow$ It is NOT guaranteed that there is a solution

Consider $\{V_\delta\}_{\delta>0}$ subspaces of V st the consistency assumption holds: $\forall u \in V \Rightarrow \text{dist}_V(u, V_\delta) = \min_{v_\delta \in V_\delta} \|u - v_\delta\|_V \rightarrow 0$

$\delta \downarrow 0 \Rightarrow V_\delta \rightarrow V$ How fast $\text{dist}_V(u, V_\delta) \rightarrow 0$? we choose V_δ so that \exists an unique solution

We take $\{V_\delta\}_{\delta>0} \subseteq V$ such that $\exists Z \subset V$ normed space such that $\exists \Psi_Z: \mathbb{R} \rightarrow \mathbb{R}, \Psi_Z(\delta) \xrightarrow{\delta \downarrow 0} 0, \quad \text{dist}_V(u, Z) \leq \Psi_Z(\delta) \|u\|_Z \quad \forall u \in V$

controls the distance of the consistency assumptions so that we know how much fast goes $\text{dist}_V(u, V_\delta) \rightarrow 0$

NB tells us that $u_\delta \in V_\delta$, u_δ is the BEST APPROXIMATION we can obtain of $u \in V$ fixing $\{V_\delta\}_{\delta>0}$

Lemma - Cea. If u_δ solution of $VP_\delta(\mu)$ and u is a solution of $VP(\mu)$

$$\Rightarrow \|u(\mu) - u_\delta(\mu)\|_V \leq \frac{\delta(\mu)}{\alpha(\mu)} \inf_{v_\delta \in V_\delta} \|u(\mu) - v_\delta\|_V \quad \forall \mu \in P \quad (2)$$

for consistency assumption, if we chose the right δ , goes to zero
distance controlled by continuity and coercivity

Theorem If $\{V_\delta\}_{\delta>0}$ satisfy 1. and 2. $\Rightarrow \|u(\mu) - u_\delta(\mu)\|_V \leq \frac{\delta(\mu)}{\alpha(\mu)} \Psi_Z(\delta) \|u(\mu)\|_Z \rightarrow$ by hypothesis $u \in Z$

↳ Combined 1+2 \Rightarrow A-priori error estimation

Example Consider $Z = H^1(\Omega), r_1, \forall \mu \in P, V = H^1(\Omega)$ and $Z \subset V$.

1. $H^1 \Rightarrow \Psi_Z(\delta) = \delta^{r-1}: \text{dist}_V(w, V_\delta) \leq \delta^{r-1} \|w\|_{H^1} \quad \forall w \in H^1 \Rightarrow \|u(\mu) - u_\delta(\mu)\|_V \leq \frac{\delta(\mu)}{\alpha(\mu)} \delta^{r-1} \|u(\mu)\|_{H^1}$
2. $\nabla(K \cdot \nabla \mu) = f \quad u \in H^2 = Z, \quad \Psi_Z(\delta) = \delta^1 \Rightarrow \|u(\mu) - u_\delta(\mu)\|_V \leq \frac{\delta(\mu)}{\alpha(\mu)} \delta \|u(\mu)\|_{H^2}$

In general



$VP_\delta(\mu): a(u_\delta, v_\delta) = F(v_\delta) \quad \forall v_\delta \in V_\delta \Rightarrow$ but $V \subset V$, so $\forall v_\delta \in V_\delta: a(u_\delta, v_\delta) - a(u, v_\delta) = F(v_\delta) - F(v) = 0$

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

If $a(\cdot, \cdot)$ is symmetric, since is also bilinear and coercive, it's a scalar product

$$\Rightarrow a(u_\delta - u, v_\delta) = \langle u_\delta - u, v \rangle_a = 0 \Rightarrow u_\delta - u \perp V_\delta$$

$\Rightarrow u_\delta$ is the orthogonal projection of u on V_δ with respect to a-product

Subtracting $Vf(\mu)$ from $VP(\mu)$: $a(u, v_\delta) - a(u, v_\delta) = F(v_\delta) - F(v) = 0 \quad \forall v \in V \Rightarrow a(u - u_\delta, v_\delta) = 0 \quad \forall v \in V_\delta$

Galerkin orthogonalization. Consider $a(u - u_\delta, v_\delta) = 0 \quad \forall v_\delta \in V_\delta$

$$\Rightarrow \|u - u_\delta\|_a = \sqrt{a(u - u_\delta, u - u_\delta)_a} = \sqrt{a(u - u_\delta, u - u_\delta)} \leq \|u - u_\delta\|_V \quad \forall v_\delta \in V_\delta, \forall \mu \in P$$

↳ being it the orthogonal projection, it's the element of V closest to $u \Rightarrow$ all the other v_δ s are at a greater distance

3) How to find $u_\delta(\mu) \quad \forall \mu \in P$? Consider $V \subset V$ finite, $N_\delta = \dim(V_\delta)$. Let's take $\{\varphi_k\}_{k=1}^{N_\delta}$ basis of V_δ ,

$$\varphi_k: S \rightarrow \mathbb{R} \Rightarrow \forall v_\delta \in V_\delta \quad v_\delta(\mu) = \sum_{k=1}^{N_\delta} \varphi_k \cdot v_\delta(\mu) \quad \forall \mu \in P$$

Putting this expression in $VP_\delta(\mu)$, considering $v_\delta = \varphi_j$ G.I.R.

we want to find $u_\delta \in V_\delta$ st $a(u_\delta, v_\delta; \mu) = F(v_\delta; \mu) \quad \forall v_\delta \in V_\delta$

$$\Rightarrow a\left(\sum_{k=1}^{N_\delta} u_{\delta,k} \varphi_k, v_\delta; \mu\right) = F(v_\delta; \mu) \Rightarrow \sum_{k=1}^{N_\delta} a(u_\delta, \varphi_k, v_\delta; \mu) = F(v_\delta; \mu)$$

$$\Rightarrow \sum_{k=1}^{N_\delta} u_{\delta,k} \cdot a(\varphi_k, v_\delta; \mu) = F(v_\delta; \mu) \quad \forall v_\delta \in V_\delta \quad \text{where } N_\delta \text{ and } u_{\delta,k} \text{ unknowns}$$

What v_δ we take? We can take $v_\delta = \varphi_j$ for every $j \in 1, \dots, N_\delta$

one for every basis function φ_j

$$\Rightarrow \text{it holds } \sum_{k=1}^{N_\delta} u_{\delta,k} \cdot a(\varphi_k, \varphi_j; \mu) = F(\varphi_j; \mu) \rightarrow N_\delta \text{ equations and } N_\delta \text{ unknowns: I can solve it for } u_{\delta,k}$$

In matrix form

$$u_\delta = \begin{bmatrix} u_{\delta,1} \\ \vdots \\ u_{\delta,N_\delta} \end{bmatrix} \begin{bmatrix} & & \\ & & \\ & & \end{bmatrix} = \begin{bmatrix} f_1 \\ \vdots \\ f_{N_\delta} \end{bmatrix} = \begin{bmatrix} F(\varphi_1; \mu) \\ \vdots \\ F(\varphi_{N_\delta}; \mu) \end{bmatrix}$$

$$A^\mu \text{ st } A_{jk} = a(\varphi_k, \varphi_j; \mu)$$

$$\Rightarrow \begin{bmatrix} A^\mu & & \\ & I_{N_\delta \times N_\delta} & \\ & & f^\mu \end{bmatrix} \begin{bmatrix} u_\delta \\ & I_{N_\delta \times N_\delta} \\ & & f^\mu \end{bmatrix}$$

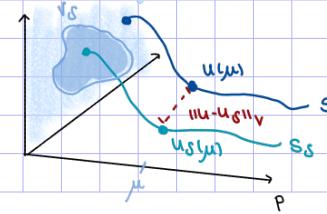
full problem
High fidelity problem
 $A^\mu u_\delta^\mu = f^\mu$

So that we obtain a discrete map and a discrete manifold

$$S_\delta: P \rightarrow V_\delta$$

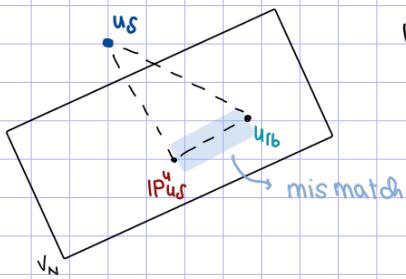
$$\mu \mapsto u_\delta(\mu)$$

$$M_\delta = \{u_\delta(\mu) \in V_\delta : \mu \in P\} \subset V_\delta$$



Step forward...

Let us make a plot of the projection of u_δ in V_N



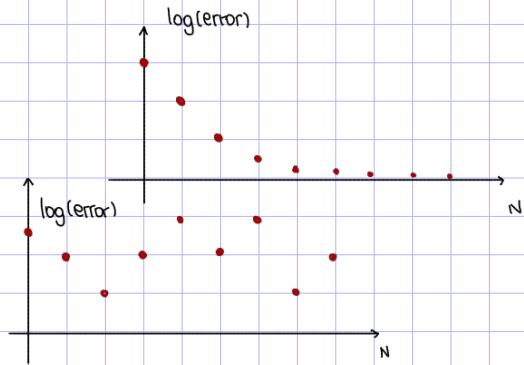
Where $P \in \mathbb{R}^{N_d \times N}$ and $\|P_{u_\delta}^u - u_{rb}\|$ is the mismatch from $P_{u_\delta}^u$ & u_{rb}

Moreover, focusing on the error, we can plot

POD: via theorems we're sure that the error goes to zero

POD-NN: we do NOT have theorems → enlarging the space,
↓ hence using more N , the error does NOT decrease

The issue is the accuracy: the error cannot go under 0.01



4 LOW FIDELITY MODEL

The idea is to reduce S_δ to solve the solution in a cheaper way \Rightarrow we approximate S_δ using a space $V_N \subseteq V_\delta$ for which $N = \dim(V_N)$, $N_\delta = \dim(V_\delta)$ and $N < N_\delta$. Hence, S_N is the approximation of S_δ that is named as **low fidelity model** in contrapposition to S_δ that is **high fidelity model** \rightarrow cheaper to solve it

Given the manifold M_δ we wish to build an approximation in which is cheaper to evaluate $A_{\mu \in P}$. For the approximation of S_δ we imagine to select a subset $\{\mu_1, \dots, \mu_N\} \subset P$ and to compute the **high fidelity** solution we use the so called **snapshots** $\{w_i = u(\mu_i), \dots, w_N = u(\mu_N)\} \in V_\delta$ that produce a "surrogate" of S_δ that we call **low-fidelity model** S_N or **reduced model**.

Different way to approximate S_δ :

1. Interpolation - not taken in account
2. Least Square - not taken in account
3. Reduced basis Method RBM \Rightarrow Galerkin approach
4. Neural Network: supervision of unsupervised learning \rightarrow later

4.1 Reduced basis method RBM

Consider $V \subseteq V_\delta$, $a(u_\delta, v_\delta; \mu) = F(v_\delta, \mu)$, the subspace of the solution is $\{w_1, \dots, w_N\} \in V_\delta$ that could be linear independent:

$$\sum_{i=1}^N a_i w_i = 0 \Leftrightarrow \forall i \neq j: a_i = 0 \quad \forall i \in \{1, \dots, N\}$$

that solutions form a basis, so they build a subspace $V_N \subseteq V_\delta$, $N < N_\delta$

We formulate a variational problem on V_N :

$$\text{Find } u_N \in V_N \text{ st } A_{\mu \in P}, a(u_N, v_N; \mu) = F(v_N; \mu) \quad \forall v_N \in V_N \Rightarrow a(u_\delta - u_N, v_N) = 0 \quad \forall v_N \in V_N$$

\downarrow can be solved because is a bilinear form continuous, coercive, etc

So if the hp holds, it's true $\|u_\delta(\mu) - u_N(\mu)\|_V \leq \inf_{v_N \in V_N} \|u_\delta(\mu) - v_N\|_V \Rightarrow$ it holds the CG lemma

we express $u_N(\mu)$ in terms of $\{w_i\}_{i=1}^N$

Generation of the linear system. We have to select a basis $\{w_1, \dots, w_N\}$ of V_N and we rewrite the solution in that way to compute the degree of freedom \rightarrow points on the mesh

$\Rightarrow S_N: \mu \in P \mapsto u_N(\mu) \in V_N$

$$u_N(\mu) = \sum_{i=1}^N u_{N,i}(\mu) w_i \Leftrightarrow u_N(\mu) := [u_{N,1}(\mu), \dots, u_{N,N}(\mu)] \in \mathbb{R}^N$$

array of coefficient

$\Rightarrow M_N := \begin{bmatrix} u_{N,1}(\mu_1) & \dots & u_{N,N}(\mu_1) \\ \vdots & \ddots & \vdots \\ u_{N,1}(\mu_N) & \dots & u_{N,N}(\mu_N) \end{bmatrix} \in \mathbb{R}^{N \times N}$

Put that in the bilinear form: $a(u_N, v_N; \mu) = F(v_N; \mu) \Rightarrow \sum_{i=1}^N a(w_i, w_i; \mu) = F(v_N; \mu) \quad \forall v_N \in V_N$

\downarrow we have N unkowns, we need N equations
set of linear independent vector from which we can derive the solution

Linear system parametric

dimension very low, easy to solve

$$\text{In matrix form } u_N \in \mathbb{R}^N, A_N \in \mathbb{R}^{N \times N}, f_N \in \mathbb{R}^N \Rightarrow A_N^T u_N = f_N$$

\downarrow NB $A_N^T = [\langle w_i, w_j; \mu \rangle]_{1 \leq i, j \leq N}$, $f_N := [F(w_i; \mu)]_{1 \leq i \leq N}$ what we want to compute

In this way we can compute the solution in a cheaper way, with a not to huge error

priori formula, exponential decay

Exercise - How the matrix A_N is related to A_δ ?

\Rightarrow From A_N^T we can deduce $\{w_i\}_{i=1}^N \in V_N \subseteq V_\delta$ and $\{v_j\}_{j=1}^N \in V_\delta$ so $w_i = \sum_{j=1}^N b_{ij} v_j$

To answer, we can use the definition of the matrix A :

$$(A_N^T)_{ji} = a(w_i, w_j; \mu) = a\left(\sum_{k=1}^N b_{ik} v_k, \sum_{e=1}^N b_{je} v_e; \mu\right) = \sum_{k,e=1}^N b_{ik} b_{je} a(v_k, v_e; \mu) b_{je}$$

\downarrow In matrix form $A_N^T = B^T A_\delta^T B$, where $(B)_{ki} = b_{ki}$, so the two matrices are similar

$B \in \mathbb{R}^{N \times N}$, $A_\delta^T \in \mathbb{R}^{N \times N}$

How to select $\{w_i\}_{i=1}^N$ basis of V_N ? Proper Orthogonal Decomposition POD

NB The computation of A_N^T requires element of size N_δ , that is NOT good if $N_\delta \gg N$!

4.2 How to select the snapshots?

We saw that if we have a set of snapshots $\{w_i\}_{i=1}^N$ linearly independent we are able to reduce and solve the VP $_N(\mu)$

But, how to detect them? $P_M = \{\mu_1, \dots, \mu_M\} \subseteq P$

Suppose we know $H \leq N$ snapshots $w_m = w(\mu_m) \in V$, $m=1, \dots, H$. They may be linearly dependent $\Rightarrow X_H := \text{Span}\{w_m\}_{m=1}^H$ may have dimension $< H$.

The idea is to reduce the unnecessary snapshots



Also known as PRINCIPAL COMPONENT ANALYSIS

4.2 Proper orthogonal decomposition POD \rightarrow can be seen from the Singular-Value-Decomposition's point of view

Consider $P \in \mathbb{R}^{P \times H}$ the parameter space, so the subspace $P_H \subseteq P$ with $\|P_H\|_F = H$ is defined as $P_H = \{u_i, \dots, u_H\} \in P$. Here, we can compute a high fidelity solution $\forall \mu_i \in P_H : \underbrace{w_i = w(\mu_i)}_{\text{solution}}$ of VPS \rightarrow for that problem

So we can build $X_H := \text{Span}\{w_i\}_{i=1}^H$ that are linear independent

The main idea is to reduce X_H so that to have an independent space. To do that we introduce an operator

$$\begin{aligned} C: X_H &\rightarrow X_H \\ v &\mapsto \sum_{m=1}^H \langle v, w_m \rangle_v w_m \end{aligned}$$

In which w_m is called snapshot and represent a solution

Properties of C :

- ① C is self-adjoint $\forall v, z \in X_H : \langle Cv, z \rangle_v = \langle Cz, v \rangle_v \Rightarrow \langle Cv, z \rangle_v = \sum_{m=1}^H \langle v, w_m \rangle_v \langle w_m, z \rangle_v = \langle Cz, v \rangle_v$ eigenfunctions
- ② C is semi-definite positive $\forall v \in X_H : \langle Cv, v \rangle_v = 0$, so C has H eigenpairs $\{\lambda_m, \psi_m\}_{m=1}^H$
- ↳ So $\{\lambda_m \geq 0\}, \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_H > 0$ and $\|\psi_m\|_F = 1 \Rightarrow \langle Cv, v \rangle_v = \sum_{m=1}^H \langle v, w_m \rangle_v^2 \geq 0$ eigenvalue $C\psi_m = \lambda_m \psi_m, \forall m=1, \dots, H$
- ③ C is injective iff $\{w_m\}_{m=1}^H$ are linearly independent

Aim. Find linear independent snapshot: we need eigenpair decomposition



Theorem. We define $(C)_{mn} := \langle w_m, w_n \rangle_v$ and $C \in \mathbb{R}^{H \times H}$ that is the covariance matrix

\Rightarrow The eigenvalues $\{\lambda_m\}_{m=1}^H$ of C are the eigenvalues of the operator C

\Rightarrow The eigenfunctions $\psi_m := \frac{1}{\sqrt{\lambda_m}} \sum_{k=1}^H w_k [a_m]_k \forall m=1, \dots, H$, where a_m are the eigenvectors of C , $[a_m]_k = a_k$

Proof on the lecture notes

More Details on POD on Lecture 4

4.3 Affine case

If the bilinear form $a(\cdot, \cdot; \mu)$ admits an affine decomposition s.t. $a(u, v; \mu) = \sum_{q=1}^{Q_a} \tilde{a}_q(\mu) a_q(u, v)$ where $\tilde{a}_q \in \mathbb{R}, a_q$

$a_q: V \times V \rightarrow \mathbb{R}$ and $\tilde{a}_q: P \rightarrow \mathbb{R}$, then

$$A^H = B^T A^H B = B^T \left(\sum_{q=1}^{Q_a} \tilde{a}_q(\mu) A_q \right) B = \sum_{q=1}^{Q_a} \tilde{a}_q(\mu) B^T A_q B \in \mathbb{R}^{N \times N}$$

can be computed once for all $\mu \in P \Rightarrow$ off-line phases

and A^H is assemble fastly and cheaply each value of $\mu \Rightarrow$ on-line phases

Similarly, the computation of f_N^H can be deduced from f^H with $(f_N^H)_{ij} = F(w_i; \mu) = \sum_{j=1}^{Q_f} b_{ij} F(\varphi_j; \mu) = (B^T f^H)_{ij}$ and, if $F(\cdot; \mu)$ admits affine decomposition i.e

$$F(\cdot; \mu) = \sum_{q=1}^{Q_f} \tilde{f}_q(\mu) f_q(\cdot)$$

with $F_q: V \rightarrow \mathbb{R}$ and $\tilde{f}_q: P \rightarrow \mathbb{R}$ then $f_N^H = B^T f^H = B^T \sum_{q=1}^{Q_f} \tilde{f}_q(\mu) f_q = \sum_{q=1}^{Q_f} \tilde{f}_q(\mu) (B^T f_q) \in \mathbb{R}^N$ and again $B^T f_q$ is independent of μ (off-line)

4.4 Error estimation in RBM

Consider the problem $a(u, v) = f(v)$ in the space V , projecting on lower dimension fixing $\mu \in P$
 $a(u_\mu, v_\mu) = f(v_\mu)$ in the space $V_\mu \subset V$ that represents the high fidelity problem
 $a(u_N, v_N) = f(v_N)$ in the space $V_N \subset V$ that is the low fidelity problem

Before going into details how to find v_N , we need to say something about the decay of the error $\|u_\mu - u_N(\mu)\|_V$.
We want to compute the error $\rightarrow \|u_\mu(\mu) - u_N(\mu)\|_V$

Definition. Consider the finite subset $Y \subset V$ and $w \in V$, the distance of an element from a space is

$$\text{dist}_V(w, Y) = \min_{y \in Y} \|w - y\|_V$$

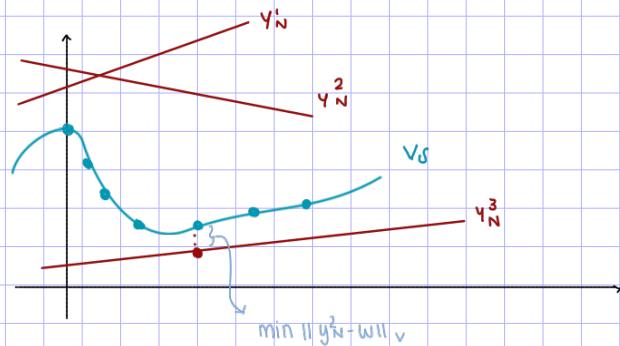
Definition. Given $M \subset V$ finite, the distance of M from Y is

$$\text{dist}_V(M, Y) = \sup_{w \in M} (\text{dist}_V(w, Y)) = \sup_{w \in M} (\min_{y \in Y} \|w - y\|_V)$$

Definition. Given an integer $N \geq 1$, the Kolmogorov N -th width of $M \subset V$ is

$$d_N(M) = \inf_{\substack{Y \subset V \\ \dim(Y) = N}} \{ \text{dist}_V(M, Y) \} = \inf_{\substack{Y \subset V \\ \dim(Y) = N}} \left\{ \sup_{w \in M} (\min_{y \in Y} \|w - y\|_V) \right\}$$

Why we need all this definition? With reduced Space method we want to compute $V_N \subset V$ in the way that the error that I perform is lower as possible, to do that we exploit these definitions



NB \Rightarrow From the definition, we can say that $d_N(M)$ is the maximum error that we can commit approximating an element of M with element contained in a n -dimensional subspace

$$\Rightarrow d_N(M) \leq d_N(H) \quad \forall N \geq 1$$

\Rightarrow If H is compact, $d_N(H) \rightarrow 0, N \rightarrow \infty$ BUT HOW IS IT FAST?

More detail on Lecture 3

4.5 More info NOT required during the exam

In our case V is the variational set of the weak formulation, and $M \subset V$ FEM dimension and $Y_N \subset M \subset V$ RB
Suppose that Ω is compact and $d_{N+1}(\Omega) \leq d_N(\Omega)$ and $d_N(\Omega) \xrightarrow[N \rightarrow \infty]{} 0$: if we enlarge N , the error decrease, but we know that we cannot set $N \rightarrow \infty$ because we are in finite dimension space



How fast this error goes to zero?

Property. Consider $M_\mu := \{u_\mu | u \in V_\mu : \mu \in P\} \subset V_\mu$, where $u_\mu(\mu)$ is a solution of $V_\mu(\mu)$ in which we have $a(\cdot, \mu) = f(\cdot, \mu)$

1. If $a(\cdot, \mu)$ and $f(\cdot, \mu)$ are V -differentiable wrt μ $\Rightarrow d_N(M_\mu) \leq C N^{-r}$

2. If $a(\cdot, \mu)$ and $f(\cdot, \mu)$ are analytic wrt $\mu \in P$ $\Rightarrow d_N(M_\mu) \leq C e^{-\gamma N}, \gamma > 0$

↓
analytic for $f: \Omega \rightarrow \mathbb{R}$ wrt $x \in \Omega$: $\exists n \sum_{n=0}^{\infty} a_n(x-x_0)^n = f(x) \quad \forall x \in B_\varepsilon(x_0)$ with an succession, $a_n := \frac{f^{(n)}(x_0)}{n!}$

The reduced space converges with exponential decay

This is a property of the weak formulation, hence of the problem itself, does not depend on the discretization

5. **GREEDY STRATEGY** \Rightarrow with the SVD we're able to choose a reduced space V_N s.t. the average error $\Phi(y_N)$ is minimum

The POD produce a space V_N in which we can minimize the function

$$\Phi(y_N) := \frac{1}{M} \sum_{m=1}^M \inf_{v \in V_N} \|w_m - v\|_V^2 \Rightarrow \min_{y_N} \Phi(y_N) = \frac{1}{M} \sum_{m=1}^M \inf_{v \in V_N} \|w_m - v\|_V^2$$

Hence, we chose the V_N space s.t. the error on average is the smallest one \rightarrow we can find other way

AH. We want to find this relation: $\|u_s(\mu) - u_N(\mu)\|_V < \epsilon_{tol}$ \rightarrow we have to produce a space in this way GREEDY

\hookrightarrow To do this we need the generation of $M > N$ snapshots and the computation of SVD decomposition (expensive), even if it is off-line

\downarrow for this reason the greedy strategy is implemented \Rightarrow This strategy requires an a-posteriori error estimator $\eta = \eta(u, V_N)$

5.2 Idea of greedy strategy. We start from $P \subset P$, we pick one value of the parameter, compute the HF solution and

\downarrow if it is good, this enters in the set

Starting from $V_N = \{v\}$

We select $\mu \in P$ and we compute $u_s(\mu)$ \rightarrow we add $u_s(\mu)$ to V_N but to do that we have to modify it a little bit

Indeed, we compute $w_1 := \frac{u_s(\mu)}{\|u_s(\mu)\|_V}$, so we have just to normalize it before to put in the basis

\downarrow computed one each iteration \Rightarrow cheap!

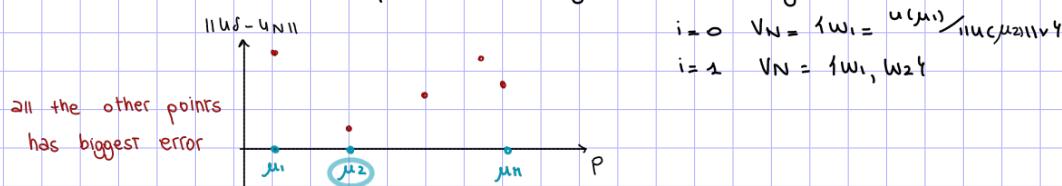
Now we have $V_N := \text{span}\{w_1\}$ and $N = 1$

• Is the space enough? $\max_{\mu \in P} \|u_s(\mu) - u_N(\mu)\|_V < \epsilon_{tol}$ so all the error has to be under a tolerance

• Yes we have finished

• No we need an iterative algorithm to enlarge the space: $\mu_{\text{int}} = \arg \max_{\mu \in P} \|u_s(\mu) - u_N(\mu)\|_V$

So I take the parameter that gives me the largest error



I compute the high-fidelity solution $u_s(\mu_{\text{int}})$

Hence, I normalize $w_{\text{int}} := \frac{u_s(\mu_{\text{int}})}{\|u_s(\mu_{\text{int}})\|_V} + V_N$ NB Has to be orthogonal to all the other functions inside the space V_N to do that we exploit

So that $V_{N+1} := \text{span}\{V_N, w_{\text{int}}\}$

We check again the tolerance

GRAM SCHMIDT

BUT, we have to know all the values of the HF solution? No, we want to compute that in a cheaply way
So we will use a quantity that controls the error η , so that when this quantity $\eta < \epsilon$ we know also that all our values are under the tolerance \Rightarrow we will use an a posteriori estimator that controls the error

$\|u_s(\mu) - u_N(\mu)\|_V \leq \eta(\mu, V_N) < \epsilon_{tol}$ \Rightarrow $\eta(\mu)$ shall be computed for many values of $\mu \in P$, then we can find $\eta(\mu)$ cheaply

\hookrightarrow u is the reduced basis approximation of $u_s(\mu)$ in the space V_N

Why greedy is better than POD?

1. Controls the maximum error and not the average error (as POD)
2. Compute at most N solution
3. It is cheaper

5.3 A posteriori error computation

We recall the $V_N(\mu)$ and we want to find $u_N(\mu) \in V_N$ \Rightarrow $a(u_N(\mu), v_N; \mu) = F(v_N; \mu) \quad \forall v_N \in V_N, \forall \mu \in P$

We introduce the weak residual: $R_s(u_N(\mu); \mu) \in V_N'$ \rightarrow stay in the dual of V_N

$$\langle R_s(u_N(\mu); \mu), v_N \rangle_{V_N'} := F(v_N; \mu) - a(u_N(\mu), v_N; \mu) \quad \forall v_N \in V_N \quad \text{from the definition of } \|\cdot\|_{V_N'}$$

$$\|R_s(u_N(\mu); \mu)\|_{V_N'} := \sup_{v_N \in V_N} \frac{\langle R_s(u_N(\mu); \mu), v_N \rangle_{V_N'}}{\|v_N\|_V} \quad \text{This is a sort of bilinear function}$$

This is the duality pair

Property It is true that . $\frac{1}{\delta(\mu)} \|R_s(u_N(\mu); \mu)\|_V \leq \|u_s(\mu) - u_N(\mu)\|_V \leq \frac{1}{\delta(\mu)} \|R_s(u_N(\mu); \mu)\|_V$

high-fidelity solution

reduce solution

$\delta(\mu)$

coercivity constraint

We are in the discrete formulation of the

problem

$$\begin{aligned} \text{Proof. From coercivity we have } \delta(\mu) \|u_s(\mu) - u_N(\mu)\|_V^2 &\leq a(u_s(\mu) - u_N(\mu); u_s(\mu) - u_N(\mu); \mu) \\ &= F(u_s(\mu) - u_N(\mu); \mu) - a(u_N(\mu), u_s(\mu) - u_N(\mu); \mu) \\ &= \langle R_s(u_N(\mu); \mu), u_s(\mu) - u_N(\mu) \rangle_V \leq \|R_s(u_N(\mu); \mu)\|_V \|u_s(\mu) - u_N(\mu)\|_V \end{aligned}$$

The other direction follow from the continuity hypothesis $\forall \delta \in V'$

$$r_s(u_N(\mu), \delta; \mu) := a(u_s(\mu), \delta; \mu) - a(u_N(\mu), \delta; \mu) = a(u_s(\mu) - u_N(\mu), \delta; \mu) \leq \tau_s(\mu) \|u_s(\mu) - u_N(\mu)\|_V \| \delta \|_V$$

Using the definition of $\|R_s(u_N(\mu), \mu)\|_V$

It is good if it can be computer faster (cheaper!)

Now, using the upper bound of the property, we can derive an a-posteriori error estimator. Let's make it computable using the Riesz' representation theorem

Riesz' Representation Theorem $\exists \hat{r}_s(u_N(\mu); \mu) \in V'$ s.t. $\langle \hat{r}_s(u_N(\mu); \mu), \delta \rangle_V = \langle R_s(u_N(\mu), \mu), \delta \rangle_V \quad \forall \delta \in V'$
and $\|\hat{r}_s(u_N(\mu); \mu)\|_V = \|R_s(u_N(\mu), \mu)\|_V$

in the original space V

Duality pair

The dual space V'

- NB
- ① The computation of $\hat{r}_s(\mu) := \hat{r}_s(u_N(\mu); \mu)$ requires to solve an high-fidelity problem (expensive)
However, in the affine case the computation can be efficient
but we can perform offline-online phase
 - ② The error a-posteriori estimation $\|u_s(\mu) - u_N(\mu)\|_V \leq \frac{1}{\delta(\mu)} \|\hat{r}_s(\mu)\|_V := \gamma_V(\mu)$ requires also the computation of $\delta(\mu)$ → we need to know for all the parameter
This can be accomplished with the lower bound estimation $\alpha_{LB}(\mu) \leq \delta(\mu)$ cheap to be computed
To use in algorithm with relative tolerance

There is the possibility to compute also an a-priori relative error estimator, observing that

$$\|u_N(\mu)\|_V = \|u_s(\mu) + u_N(\mu) - u_s(\mu)\|_V \leq \|u_s(\mu)\|_V + \|u_N(\mu) - u_s(\mu)\|_V \leq \|u_s(\mu)\|_V + \gamma_V(\mu)$$

If we define $\gamma_{rel,V}(\mu) := \frac{\gamma_V(\mu)}{\|u_s(\mu)\|_V}$, we obtain

$$\|u_N(\mu)\|_V \leq \|u_s(\mu)\|_V + \frac{1}{2} \|u_N(\mu)\|_V \stackrel{\gamma_V(\mu)}{\leq} \|u_s(\mu)\|_V + \frac{1}{2} \|u_N(\mu)\|_V \Rightarrow \|u_N(\mu)\|_V \leq 2 \|u_s(\mu)\|_V$$

Property

$$1. \|u_s(\mu) - u_N(\mu)\|_V \leq \gamma_V(\mu) := \frac{\|\hat{r}_s(\mu)\|_V}{\alpha(\mu)}$$

$$\|u_s(\mu) - u_N(\mu)\|_V \leq \frac{\gamma_V(\mu)}{\alpha(\mu)}$$

$$2. \frac{\|u_s(\mu) - u_N(\mu)\|_V}{\|u_s(\mu)\|_V} \leq \gamma_{rel,V}(\mu) := \frac{1}{2} \frac{\gamma_V(\mu)}{\|u_N(\mu)\|_V}$$

$$3. eff_V(\mu) := \frac{\gamma_V(\mu)}{\|u_s(\mu) - u_N(\mu)\|_V}$$

$$eff_{rel,V}(\mu) := \frac{\gamma_{rel,V}(\mu) \|u_s(\mu)\|_V}{\|u_s(\mu) - u_N(\mu)\|_V}$$

Algorithm 1 Greedy Algorithm for Reduced Basis Construction in POD

1: Initialization:

2: Choose a set of training parameters (or a continuous parameter domain) P_H .

3: Select an initial parameter $\mu_1 \in P_H$.

4: Compute the high-fidelity solution $u(\mu_1)$.

5: Define the initial reduced basis $V_1 = \text{span}\{u(\mu_1)\}$ (or its normalization w_1).

6: Set $N = 1$.

7: Define an error tolerance ϵ_{tol} .

8: while $\max_{\mu \in P_H} \|u(\mu) - u_N(\mu)\|_{E_{tol}} > \epsilon_{tol}$ do

9: Find the parameter $\mu_{max} \in P_H$ that maximizes the estimated error: $\mu_{max} = \arg \max_{\mu \in P_H} \|u(\mu) - u_N(\mu)\|_{E_{tol}}$ where $u_N(\mu)$ is the approximated solution in the space V_N .

10: Compute the high-fidelity solution for this parameter: $u(\mu_{max})$.

11: Orthogonalize $u(\mu_{max})$ with respect to the functions already in V_N using Gram-Schmidt to obtain a new basis function w_{N+1} .

12: Extend the reduced basis: $V_{N+1} = \text{span}\{V_N, w_{N+1}\}$.

13: Increment the counter: $N = N + 1$.

14: end while

15: Result:

16: The final reduced basis $V_N = \text{span}\{w_1, w_2, \dots, w_N\}$.

5. DIFFERENT APPROACH TO COMPUTE THE COHERCIVITY FUNCTION $\alpha(\mu)$

→ To compute η we need $\alpha(\mu)$ for each value of the parameters

Recall the error estimator $\eta(\mu) := \|\tilde{v}_S(\mu)\|_V$

compute cheaply with online/offline phase

$\alpha(\mu)$ This part, coercivity function, difficult to compute

↓

The smallest eigenvalue of the spectral problem

Associated to high fidelity problem $\forall \mu \in P$ $\alpha(\mu) = \inf_{\tilde{v} \in V_S} \frac{\alpha(v, \tilde{v}; \mu)}{\|\tilde{v}\|_V^2}$

so this is an eigenvalue problem

$\Rightarrow \lambda_{\min} = \alpha(\mu)$ coercivity $\alpha(w_S, v_S; \mu) = \lambda(w_S, v_S) v$ $\lambda \in \sigma_S$ smallest eigenvalue
 $\lambda_{\max} = \alpha(\mu)$ continuity

How we can compute the eigenvalue? We fix the basis $\{\varphi_i\}_{i=1}^N$ in V_S and we have this linear combination $w_S = \sum_{j=1}^N w_{Sj} \varphi_j$, so that we obtain the matrix form

$A_S^\mu = \lambda X_S w_S$, where $(X_S)_{ij} = \langle \varphi_i, \varphi_j \rangle_V$ is the scalar product matrix, and $w_S = \begin{pmatrix} w_{S1} \\ \vdots \\ w_{SN} \end{pmatrix}$
Generalized eigenvalue problem

But this is very costly, so this has to be modified

Different approach to compute $\alpha(\mu)$ - Most of this approach approximate $\alpha(\mu)$ from the bottom computing lower-bound $\alpha_L(\mu) \leq \alpha(\mu)$

↓ We want now to analize these different approaches.

5.1 Min- η approach → This method works for the affine case and for the so called parametrically coercive problems

We need * The affinity hypothesis (separate the parameter part to the other)

Hypothesis $\alpha(u, v; \mu) = \sum_{q=1}^Q \tilde{\alpha}_q(\mu) a_q(u, v)$

* The positivity hypothesis $\tilde{\alpha}_q(\mu) > 0 \quad \forall \mu \in P$

* A sort of Coercivity $a(u, u) \geq 0 \quad \forall u \in V_S$

So that here we can compute, selecting one parameter $\tilde{\mu} \in P$, $\alpha(\tilde{\mu}) \Rightarrow$ let's assume we know $\alpha(\tilde{\mu})$
By definition, it holds $\alpha(\mu) := \inf_{\tilde{v} \in V_S} \frac{\alpha(\tilde{v}, \tilde{v}; \mu)}{\|\tilde{v}\|_V^2}$ for one value of $\tilde{\mu} \in P$

$$\begin{aligned} &= \inf_{\tilde{v} \in V_S} \sum_{q=1}^Q \tilde{\alpha}_q(\mu) \frac{a_q(\tilde{v}, \tilde{v})}{\|\tilde{v}\|_V^2} \quad \text{using } \frac{\tilde{\alpha}_q(\mu)}{\|\tilde{v}\|_V^2} = \frac{\tilde{\alpha}_q(\tilde{\mu})}{\|\tilde{v}\|_V^2} \\ &= \min_{1 \leq q \leq Q} \frac{\tilde{\alpha}_q(\mu)}{\tilde{\alpha}_q(\tilde{\mu})} \inf_{\tilde{v} \in V_S} \frac{a_q(\tilde{v}, \tilde{v}, \tilde{\mu})}{\|\tilde{v}\|_V^2} \quad \text{min over } \frac{\tilde{\alpha}_q(\mu)}{\tilde{\alpha}_q(\tilde{\mu})} \text{ computed in offline phase} \\ &= \alpha(\tilde{\mu}) \end{aligned}$$

Hence, we obtain $\alpha(\mu) \geq \min_{1 \leq q \leq Q} \frac{\tilde{\alpha}_q(\mu)}{\tilde{\alpha}_q(\tilde{\mu})} \alpha(\tilde{\mu}) := \alpha_{LB}(\mu) \Rightarrow$ lower bound of α

↓ Extending this idea we can build

5.2 Multi-parameter min- η approach → This case requires to know $I \geq 1$ values of $\alpha(\tilde{\mu}_i)$, $i=1, \dots, I$, and following the previous Select $I \geq 1 \Rightarrow \tilde{\mu}_i, i=1, \dots, I$ so that we can compute $\{\alpha(\tilde{\mu}_i)\}_{i=1, \dots, I}$ analize we can obtain
Hence,

$$\alpha(\mu) \geq \max_{1 \leq i \leq I} \left(\min_{1 \leq q \leq Q} \frac{\tilde{\alpha}_q(\mu)}{\tilde{\alpha}_q(\tilde{\mu}_i)} \alpha(\tilde{\mu}_i) \right) := \alpha_{LB}(\mu)$$

5.3 Successive constraint method SCH: complex, but works very properly based on the definition of a functional
Idea: build a functional $S: P \times Y \rightarrow \mathbb{R}$ where $Y \subseteq \mathbb{R}^Q$, $y \in Y := [y_1, \dots, y_Q]$ → affine element

$$S(\mu, y) := \sum_{q=1}^Q \tilde{\alpha}_q(\mu) y_q$$

To better see the relation between the bilinear form and the functional, we can define the space:

$$Y := \left\{ y \in \mathbb{R}^{Q_d} : \exists \forall s \in V_s \text{ s.t. } y_q = \frac{\alpha_q(s, s)}{\|s\|_s^2}, q=1, \dots, Q_d \right\}$$

Hence, the coercivity can be seen as

$$\alpha(\mu) := \inf_{s \in V_s} \frac{\alpha(s, s; \mu)}{\|s\|_s^2} = \min_{y \in Y} S(\mu; y) \quad \forall \mu \in P$$

equivalent problem to avoid to compute the eigenvalue problem

But, instead of solving the minimization problem in Y , we solve it in Y_{LB} so that the problem is easier to solve

In a bigger set Y_{LB} , I can achieve a better solution

How can we find Y_{LB} ? The method suggest to find

$$0 < \alpha_{LB}(\mu) \leq \alpha_{UB}(\mu) \Rightarrow 0 < 1 - \frac{\alpha_{LB}(\mu)}{\alpha_{UB}(\mu)} < 1$$

$$Y_{LB} \subseteq Y \subseteq Y_{UB}$$

$$\alpha_{LB} \leq \alpha \leq \alpha_{UB}$$

NB in a smaller set Y_{LB} , I have a bigger constant α_{UB}

because α is the solution of the

minimization problem, so that in $Y_{LB} \subseteq Y$

could not exist a smaller value of α_{UB}

Moreover, we define $\gamma_\alpha(\mu) := 1 - \frac{\alpha_{LB}(\mu)}{\alpha_{UB}(\mu)}$

$$\text{Define a sequence of sets } \dots \subseteq Y_{LB} \subseteq \frac{\alpha_{UB}(\mu)}{\alpha_{LB}(\mu)} Y_{LB} \subseteq Y_{LB} \subseteq \dots \subseteq Y_{LB}$$

Idea: iterative algorithm to improve the value with goal $Y_{LB} \subseteq Y \subseteq Y_{UB}$

- Algorithm -
1. Selecting $P \subseteq P$, take $\mu_i \in P$ and exactly compute $\alpha(\mu_i)$ with eigenvalue problem
 2. Compute w_i 's eigenfunction associated to $\alpha(\mu_i)$

$$\downarrow \quad \alpha(w_i, s) = \lambda(w_i, s)v \quad \forall s \in V_s$$

$$\downarrow \quad Y_{LB} := \left\{ y^1 \in \mathbb{R}^{Q_d} : y_q^1 = \frac{\alpha_q(w_i, s)}{\|s\|_s^2}, 1 \leq q \leq Q_d \right\} \quad \text{this is an approximation of } Y, \text{ an upper bound approxim}$$

3. $\forall q \in \{1, \dots, Q_d\}$ we have to solve $\alpha_q(w_i, s) = \lambda(w_i, s)v \quad \forall s \in V_s$ obtaining the eigenvalue

$\rightarrow \sigma_q^- := \text{smallest eigenvalue}$

$\rightarrow \sigma_q^+ := \text{largest eigenvalue}$

\downarrow

$$Y_{LB} := \bigcap_{q=1}^{Q_d} [\sigma_q^-, \sigma_q^+]$$

, because we select $\mu_i \in P$

\hookrightarrow Hence, here we build $y^1 \subseteq y \subseteq Y_{LB}$

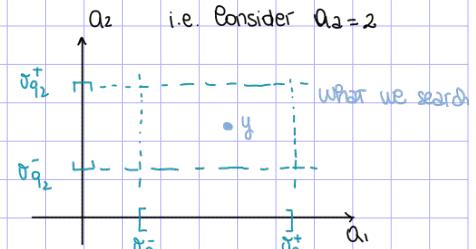
4. We have to improve with the recursion step

\rightarrow we have $\{\mu_1, \dots, \mu_{n+1}\} \subseteq P \subseteq P$

\rightarrow we have $Y_{LB} \subseteq Y \subseteq Y_{UB}$

$$\downarrow \quad \alpha_{UB}(\mu_1) \leq \alpha_{LB}(\mu_1)$$

\rightarrow we compute $\gamma_{\alpha_{LB}}(\mu_1) := 1 - \frac{\alpha_{LB}(\mu_1)}{\alpha_{UB}(\mu_1)} \rightarrow$ Is this lower than a tolerance? $\eta \leq \epsilon$



\rightarrow If so, we have to enlarge the spaces: select $\mu_{n+1} \in P$ and compute $\{\alpha(\mu_{n+1}), w^{n+1}\}$

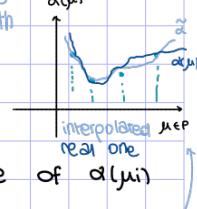
$$\downarrow \quad Y_{LB}^{n+1} := \left\{ y^{n+1} \in \mathbb{R}^{Q_d} : y_q^{n+1} = \frac{\alpha_q(w_i^{n+1}, s)}{\|s\|_s^2} \quad \forall q \in \{1, \dots, Q_d\} \right\}$$

\rightarrow So that $Y_{LB}^{n+1} := Y_{LB} \cup Y^{n+1}$

\rightarrow $\forall q \in \{1, \dots, Q_d\}$ solve the eigenvalue problem $\alpha_q(w_i, s) = \lambda(w_i, s)v \quad \forall s \in V_s$

\rightarrow Compute $B = \bigcap_{q=1}^{Q_d} [\sigma_q^-, \sigma_q^+] \rightarrow$ solving the eigenvalue problem computing min e max eigenvalue

$$\downarrow \quad Y_{LB}^{n+1} := \left\{ y \in B : \sum_{q=1}^{Q_d} \sigma_q^-(\mu) y_q \geq \alpha(\mu) \quad \forall \mu \in \{\mu_1, \dots, \mu_{n+1}\} \right\} \quad \text{The minimization is a linear programming minimization problem with at most } 2Q_d + (n+1) \text{ constraints for } Q_d \text{ variable}$$



5. Interpolation approach \rightarrow The idea is to interpolate the function $\alpha(\mu)$ with $\underline{\alpha}$ over a subset $P \subseteq P$

Consider $\underline{\alpha}: P \rightarrow \mathbb{R}$, we can select $P_i \subseteq P$ and $\{\mu_i \in P_i\}$ we can compute the real value of $\alpha(\mu_i)$

$$\mu \mapsto \alpha(\mu)$$

\downarrow as small as possible

After computing all the exact value, we interpolate the point to compute our $\alpha(\mu)$

But interpolation in huge space is not so good \rightarrow NOT polynomial one

\Rightarrow CURSE OF DIMENSIONALITY: to avoid we use RBF (exponential function)

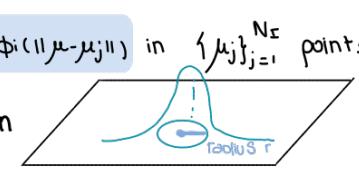
\rightarrow we cannot use the "classic" polynomial interpolation, as it suffers of the curse of dimensionality,

i.e. the number of interpolant points is too large when P increase

\Rightarrow To AVOID this we use radial basis function

In particular, we seek for the logarithm of $\alpha_I(\mu)$, s.t. $\log(\alpha_I(\mu)) = w_0 + w^\top \mu + \sum_{i=0}^{N_I} \gamma_i \phi_i(\|\mu - \mu_j\|)$ in $\{\mu_j\}_{j=1}^{N_I}$ points

belong to P s.t. $\sum_{i=0}^{N_I} \gamma_i = 0$, $\sum_{i=0}^{N_I} \gamma_i \mu_i^p = 0 \forall p \in \{1, \dots, P\}$ where $\phi: P \rightarrow \mathbb{R}$ is the RBF, i.e. the gaussian $r \mapsto e^{-r^2}$



Algorithm 1 Successive Constraint Method (SCM)

- 1: **Idea.** Build a functional $S(\mu, y)$ where $y \in Y : S(\mu, y) = \sum_{q=1}^{Q_a} \Theta_q^a(\mu) y_q$
 - 2: Define the **space**
- $$Y := \{y \in \mathbb{R}^{Q_a} : \exists v_\delta \in V_\delta \text{ s.t. } y_q = \frac{a_q(v_\delta, v_\delta; \mu)}{\|v_\delta\|_V^2}, q = \{1, \dots, Q_a\}\}$$

- 3: The **coercivity** is

$$\alpha(\mu) \geq \inf_{v_\delta \in V_\delta} \frac{a(v_\delta, v_\delta; \mu)}{\|v_\delta\|_V^2}, \forall \mu \in P$$

- 4: **Equivalent problem.** Instead of solving in Y , we solve it in a cheaper way in Y_{UB} and Y_{LB}

- 5: **Goal.** Find $\alpha_{LB} \in Y_{LB}$ and $\alpha_{UB} \in Y_{UB}$ to bound the solution $\alpha \in Y$. using an **iterative algorithm**

- 6: **Quality measure.** $\eta_\alpha = 1 - \frac{\alpha_{LB}(\mu)}{\alpha_{UB}(\mu)}$

- 7: **Input.**
 - Set of candidate functions P_a
 - Tolerance ϵ

- 8: **Initialize.** Take $\mu_1 \in P_a$ and exactly compute $\alpha(\mu_1)$ by solving the eigenvalue problem

- 9: Compute w_δ^1 eigenfunction associated to $\alpha(\mu_1)$

- 10: Define

$$Y_{UB}^1 = \{y^1 \in \mathbb{R}^{Q_a} : y_q^1 = \frac{a_q(w_\delta^1, v_\delta; \mu_1)}{\|v_\delta\|_V}, 1 \leq q \leq Q_a\}$$

- 11: $\forall q \in \{1, \dots, Q_a\}$ solve the eigenvalue problem $\forall v_\delta \in V_\delta$

$$a_q(w_\delta, v_\delta) = \lambda(w_\delta, v_\delta)_V$$

obtaining $\theta_q^- :=$ the smallest eigenvalue and $\theta_q^+ :=$ the biggest eigenvalue

- 12: Build $Y_{LB}^1 = \prod_{q=1}^{Q_a} [\theta_q^-, \theta_q^+]$

- 13: **while** $\eta_\alpha(\mu) > \epsilon$ **do**

- 14: Compute $\eta_\alpha(\mu) = 1 - \frac{\alpha_{LB}(\mu)}{\alpha_{UB}(\mu)}$

- 15: **Enlarge the space.** Select $\mu_{n+1} \in P_a$ and compute $\alpha(\mu_{n+1}; w_\delta^{n+1})$

- 16: Define

$$Y_{UB}^{n+1} = \{y^{n+1} \in \mathbb{R}^{Q_a} : y_q^{n+1} = \frac{a_q(w_\delta^{n+1}, v_\delta; \mu_{n+1})}{\|v_\delta\|_V^2} \forall q \in \{1, \dots, Q_a\}\}$$

$$\rightarrow Y_{UB}^{n+1} = Y_{UB}^n \cup Y^{n+1}$$

- 17: $\forall q \in \{1, \dots, Q_a\}$ solve the eigenvalue problem $\forall v_\delta \in V_\delta$

$$a_q(w_\delta, v_\delta) = \lambda(w_\delta, v_\delta)_V$$

obtaining $\theta_q^- :=$ the smallest eigenvalue and $\theta_q^+ :=$ the biggest eigenvalue

- 18: Build $B = \prod_{q=1}^{Q_a} [\theta_q^-, \theta_q^+]$

- 19:

$$Y_{LB}^{n+1} = \{y^{n+1} \in B : \sum_{q=1}^{Q_a} \Theta_q^a(\mu) y_q \geq \alpha(\mu) \forall \mu \in \{\mu_1, \dots, \mu_{n+1}\}\}$$

- 20: **end while**

- 21: **Output.** Approximate solution α

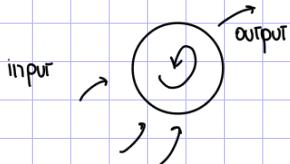
7. NEURAL NETWORKS : DEEP LEARNING

We're going to talk about deep learning based on the **artificial neural network**, a model inspired by the biological network of the human brain

- Terminology**
1. AI: artificial intelligence → simulation of the human intelligence in a machine
 2. ML: machine learning → the machine learn something by not give it directly the instruction
 3. DL: deep learning → the machine learn by using neural networks

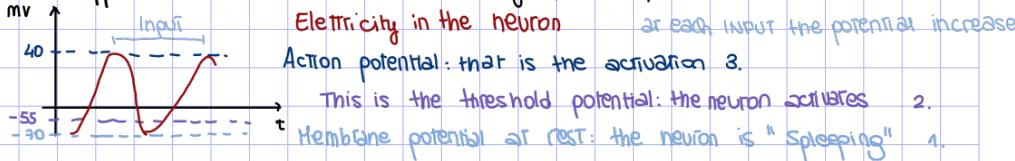


Biological model of NEURON



The IDEA is to simulate the neural inside our head
There are many inputs that are called **dendrites** and there is only one output that is the **axon** → High parallel processed

What happen inside the neuron? The voltage (electrical) increases a lot ↗



7.1 What is the first model of NN? PERCEPTRON

This is the first mathematical surrogate

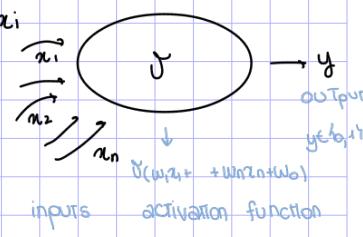
↳ We have different input x_i that represent the neuron (binary: 0-1)

In the middle there is the σ : activation function that takes a linear combination of x_i

$$\sigma(w_1x_1 + \dots + w_nx_n + w_0) \quad \text{weight of the rest potential}$$

The output is a binary one: active (1) or not active (0)

All the weights are collected in an array $w = (w_1, \dots, w_n)$, BUT the w_0 is called bias

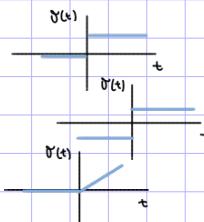


Example of activation function
usually NON-LINEAR

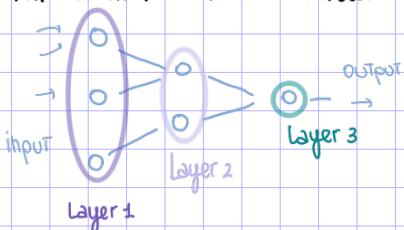
$$\sigma(t) = \begin{cases} 0 & t \leq 0 \\ 1 & t > 0 \end{cases}$$

$$2. \text{ Sign function} \quad \sigma(t) = \begin{cases} -1 & t \leq 0 \\ 1 & t > 0 \end{cases}$$

$$3. \text{ ReLU function} \quad \sigma(t) = \max(0, t) \quad \text{continuous}$$



An artificial neural network is a connection of different perceptions collected in different layers



7.2 A little bit of History

- 1943: Description of the neural model → Pitts & McCulloch propose the first and simplest mathematical model of a neuron
- 1949: Study the «activation» memory → Hebb proposed a mathematical formula to describe the learning process (memory of activation)
- 1957: Presentation of the perceptron model → Rosenblatt
- 1969: Study the limitation of the perceptron → STOP using NN because there are limitation in the mathematical model ↗
- 1986: Proposal - back propagation theory → STUDY START AGAIN: Rumelhart, Hinton, & Willcocks
- 1990: Unsupervised learning → Kohonen
- 1990 - 2000: Study about the universal approximation theorems → mathematical results

=> Our purpose is to use NN to solve PDE in a faster way

Notation - • $y :=$ output, generated by the activation function

• $\sigma :=$ activation function

• $x :=$ input

• $w :=$ weight and $w_0 :=$ bias

$$y = \sigma(\tilde{w} \cdot \tilde{x} + w_0) = \sigma(w \cdot x) \quad \text{where } w = (w_0, \tilde{w}) \in \mathbb{R}^{n+1}$$

$$x = (1, \tilde{x}) \in \mathbb{R}^{n+1}$$

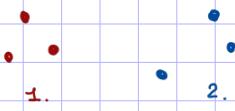
$$y = \sigma(w_1x_1 + \dots + w_nx_n + w_0)$$

1.3 How we can use perception to classify data?

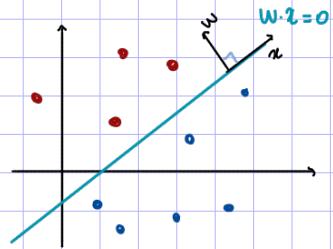
Imagine to have two class and we want to associate to the \neq element of the class the corresponding layer randomly distributed in the space

PURPOSE: I want that a machine, given an input, tells me if it is blue or red, so the belonging class of the element

We want to performe the line that separate the two classes



$w \cdot z = 0$ hyperplane in higher dimension



When I found w , given the input z I can have as answer the label
 \Rightarrow I create the classifier

$$\text{let's use } \sigma(t) = \begin{cases} 1 & t \geq 0 \\ 0 & t < 0 \end{cases}$$

How to build the CLASSIFIER? Hence, we can find the weights? TRAINING PHASE
 Consider to have a Known-sample $(x^n, y^n)_{n=1, \dots, N}$
 The goal is to find w

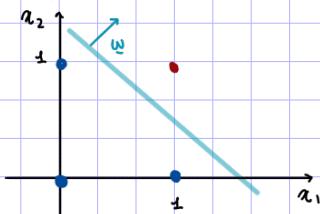
We can build the functional

$$E(w) = \frac{1}{N} \sum_{n=1}^N |y_n - \sigma(w \cdot z^n)|^2 \geq 0$$

Hence, we want $w := \arg \min_{w \in \mathbb{R}^{N+1}} E(w)$ so to find the minimization of the LOSS function

this is why NN is difficult: we have to minimize in higher space

Example. We want to approximate with the NN the AND function



We're looking for a line
 BUT it's evident that we have different solutions

logical operator

x_1	x_2	$x_1 \wedge x_2$	
0	0	0	y_1
0	1	0	y_2
1	0	0	y_3
1	1	1	y_4

ISSUE: no globale minimum in the majority of the situation

We can build the LOSS function: $E(w) = \frac{1}{N} \sum_{n=1}^N |y_n - \sigma(w \cdot z^n)|^2 \rightarrow N=4, w=(w_0, w_1, w_2) \quad z=(1, x_1, x_2)$

$$\downarrow \quad E(w) = \frac{1}{4} (|0 - \sigma(w_0 + 0)|^2 + |0 - \sigma(w_0 + 0 + w_1)|^2 + |0 - \sigma(w_0 + w_1 + 0)|^2 + |1 - \sigma(w_0 + w_1 + w_2)|^2)$$

We decide as activation function the step function: $\sigma(t) = \begin{cases} 0 & t \leq 0 \\ 1 & t > 0 \end{cases}$

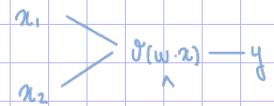
Hence, to solve the problem $w := \arg \min_{w \in \mathbb{R}^N} E(w)$ we have to solve the linear system

$$\begin{cases} |0 - \sigma(w_0)|^2 = 0 \\ |0 - \sigma(w_0 + w_2)|^2 = 0 \\ |0 - \sigma(w_0 + w_1)|^2 = 0 \\ |0 - \sigma(w_0 + w_1 + w_2)|^2 = 0 \end{cases} \Rightarrow \begin{cases} w_0 \leq 0 \\ w_0 + w_2 \leq 0 \\ w_0 + w_1 \leq 0 \\ w_0 + w_1 + w_2 \geq 0 \end{cases}$$

Taking in account $\sigma(t)$

This is the linear system that we have to solve to find the green line

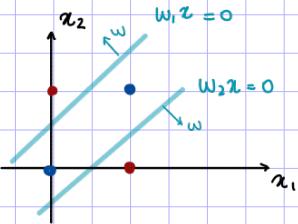
Moreover, there are \neq solutions, i.e. $w = \left(\frac{2}{3}, \frac{1}{2}, \frac{1}{2}\right)$ & $w = (-1, 1, 1)$



7.4 Feed forward neural network → We need at least two perceptrons

Limitation - XOR operator.

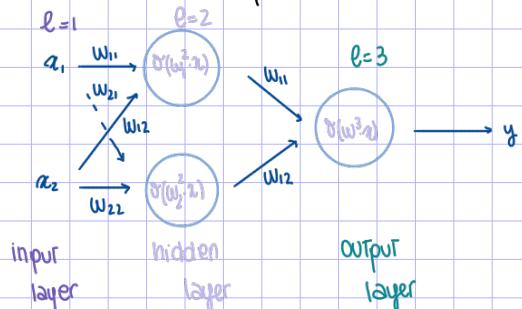
$a_1 \ a_2$	$x_1 \text{ XOR } x_2$
0 0	0 •
0 1	1 •
1 0	1 •
1 1	0 •



BUT using a line we cannot separate the point that XOR generate → one perceptron fails

↓
SOLUTION: If we use two line (two perceptron) we can face this limitation

Now we can represent the NN of the XOR



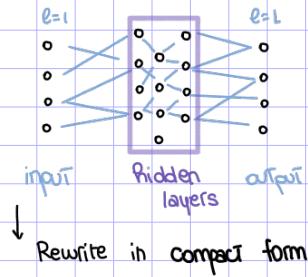
So we have different layer and 3 perceptron to have the solution and we build feed forward NN in which we can use \neq & \wedge in different layer → This is difficult to solve!

BUT in each layer, to solve fastly and parallelize the computation, we usually use the same activation function

For more details on the generalization of FF-NN see the lecture note's lesson 7

8 HOW TO COMPUTE THE WEIGHT?

8.1 FFNN recall



All the layers are fully connected with σ^e the activation function of the layer l

$$y_i^e = \sum_{j=1}^{N_{l-1}} \sigma^e(w_{ij}^e x_j^{e-1}) \quad \forall i \in \{1, \dots, N_l\} \quad \forall e \in \{2, \dots, L\} \Rightarrow \text{The non-linearity is contained in } \sigma^e$$

↓
neuron i
↓
 $y := \text{output}$

↓ Rewrite in compact form $y^e = f_e(w^e, y^{e-1})$ $y^{e-1} = z^{e-1}$ in which we consider the input of layer j is the output of layer $j-1$

Imagine to know the weight w (you already solve the optimization problem), the algorithm to the evaluation process is

ALGORITHM
to compute
the weight
y = $\sigma^L(z)$ → output
for $e=2, \dots, L$
 $z^e = w^e y^{e-1}$
 $y^e = \sigma^e(z^e)$
end
 $y = y^L$ → output

↳ The optimization problem is the offline phase that is NOT fast

↳ This is similar to the online phase that is very fast

How to compute the weight? TRAINING PHASE: we have different strategies

- ① Supervised learning: we have known snapshot $\{(x^n, y^n)\}_{n=1, \dots, N}$ that minimize the loss on input, n-th layer
↳ the network is trained using example of input & output given by a teacher \Rightarrow goal: learn the rule to map $x \mapsto y$
- ② Unsupervised learning: we do NOT have snapshot \rightarrow PINN: Physics Informed Neural Network \rightarrow See next lessons
↳ the network is trained using only inputs, no outputs are given \Rightarrow goal: discover hidden patterns in the data
- ③ Reinforcement learning: You give information and reward and the net learn
↳ the network is trained in a dynamic environment to perform a fixed goal: while training, the environment gives feedback (positive & negative reward) and the network try to maximizes

8.2 Supervised learning: we want to approximate the target function

Consider the function $f: K \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ that is the function in the hidden layer that produce the output y
↳ unknown \downarrow
↳ compact \downarrow
 $x \mapsto y$

Problem - Given FFNN fixed $(L, N_l, \sigma^e)_{e=2, \dots, L}$

Find w^e st.

$y = f(w^e, x)$ \rightarrow where f is an approximation of f

IDEA: we want to find w so the f approximate f (the hidden function) that we DO NOT have, but we have the model (FFNN) so we know \hookrightarrow we know N "samples": $(x^n, y^n = f(x^n))$

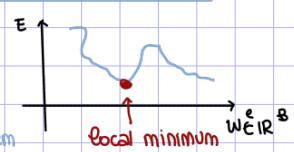
We suppose to know some snapshots $(x^n, y^n)_{n=1, \dots, N}$ and we select the loss function n -th loss error

$$E_n = \frac{1}{2} \|y^n - \bar{y}^n\|_2^2 = \frac{1}{2} \|f(w^e, x^n) - \bar{y}^n\|_2^2$$

↓ predicted value ↓ real value

The total loss error is $E = \sum_{n=1}^N E_n(w)$ and we want to find a local minimum

we try to get a minimization problem



could be NOT unique

We need a minimization method \rightarrow to solve the minimization problem

- ② GRADIENT DESCENT METHOD: based on the gradient of the loss function

Find $\min_{w \in \mathbb{R}^B} E(w)$ where $\Sigma_A := \{w \in \mathbb{R}^B : E(w) = a\}$ $\forall a \in \mathbb{R}$ are the level curve

To go in a curve with lower value we use the gradient \Rightarrow Iterative step: $w_i = w^{i-1} - \alpha \nabla E(w^{i-1})$ $\alpha > 0 \in \mathbb{R}^+$

BUT this works well iff the dimension is NOT so high, because we have $\nabla E(w) = \sum_{n=1}^N \nabla E_n(w)$ so it depends on N = number of snapshots that has to be high to learn properly

↓ Garantees only to reach a local minimum, not a global one

- ③ SOLUTION? STOCHASTIC GRADIENT DESCENT ALGORITHM: compute the gradient for only some snapshots

$$w_i = w_{i-1} - \alpha \frac{1}{N} \sum_{n=1}^N \nabla E_n(w_{i-1}) \quad M \ll N$$

I take $M \neq$ snapshots, chosen in a random way, S is called mini-batch

Example - \rightarrow Real gradient $\nabla E = \sum_{n=1}^N \nabla E_n$

\rightarrow Stochastic gradient set $S_M = \{n_1, \dots, n_M\}$ $n_i \in \{1, \dots, N\}$ yields $\tilde{\nabla} E = \frac{1}{M} \sum_{n \in S_M} \nabla E_n$

mini-batch: selected randomly at each epoch, is a subset of all the snapshots

Each iteration of the stochastic gradient descent is called epoch

↓ Adaptive Moment Estimation

ADAM ALGORITHM: improvement of the stochastic in which exploit the previous computation
all the previous gradient are called momentum of E

This is the standard used

L-BFGS: Second order algorithm \rightarrow exploit the Hessian, so it is more precise \rightarrow following quasi-Newton method

↓ Limited Memory Broyden-Fletcher-Goldfarb-Shanno

c) How to compute $\nabla E(w)$? [BACKPROPAGATION]

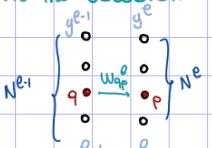
Consider $E(w) = \frac{1}{M} \sum_{m=1}^M \frac{1}{2} \|y^m(w) - \bar{y}^m\|^2$, where i.e. $w = \begin{pmatrix} w_{11} \\ w_{12} \\ \vdots \\ w_{24} \end{pmatrix}$ if $w^1 \in \mathbb{R}^{2,2}$ $w^2 \in \mathbb{R}^{2,4}$ perception

From definition, we have $\nabla E(w) = \frac{1}{2} \cdot 2 \left(y^n(w) - \bar{y}^n \right) \frac{\partial y^n(w)}{\partial w}$, where $y^n(w) = f(w, x^n)$

HOW TO COMPUTE THAT?

We know everything, but not the partial derivative $f: \mathbb{R}^S \rightarrow \mathbb{R}^M$ $\frac{\partial f}{\partial w} = \left\{ \frac{\partial f_j}{\partial w_{pq}} \right\}$ \rightarrow we need to compute the Jacobian

Imagine to fix the layers $l-1, l$ that are fully connected, and select neuron p and q $y^l = \phi_{l-1}(w^l \cdot y^{l-1})$



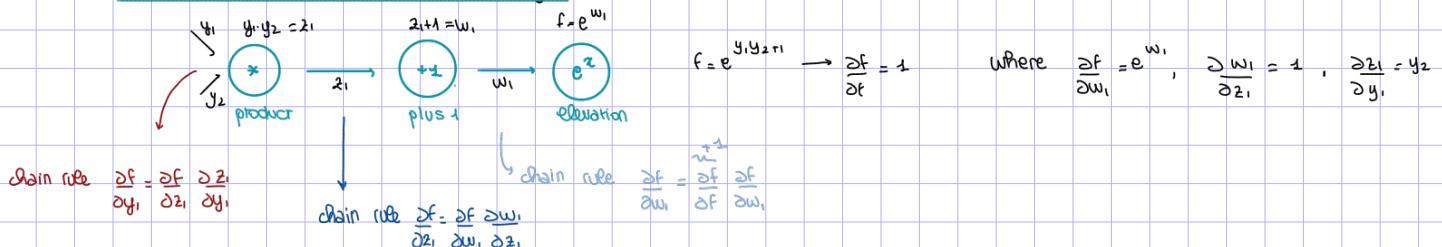
So we can use the chain rule to compute the partial derivative

$$\frac{\partial f_j}{\partial w_{pq}} = \frac{\partial f_j}{\partial y^l_p} \cdot \frac{\partial y^l_p}{\partial w_{pq}}$$

because I derive wrt w_{pq} so all the other element becomes ϕ

we DO NOT KNOW i.e. if $f_j = e^{y_1 y_2 + 1} \Rightarrow \frac{\partial f_j}{\partial y_1} = e^{y_1 y_2 + 1} \cdot y_2 \quad \frac{\partial f_j}{\partial y_2} = e^{y_1 y_2 + 1} \cdot y_1$

BUT How to compute that in NN logic?



Hence, go back we have all the elements and we can compute everything

$$\frac{\partial f}{\partial y_1} = \frac{\partial z_1}{\partial z_1} \cdot \frac{\partial w_1}{\partial z_1} \cdot \frac{\partial f}{\partial w_1} = y_2 e^{y_1 y_2 + 1}$$

Starting to the end, we know everything so we can compute the derivative!

Now, we can apply this IDEA to compute $\frac{\partial f_j}{\partial y^l_p} = \frac{\partial f_j}{\partial y_h^{l+1}} \cdot \frac{\partial y_h^{l+1}}{\partial y^l_p}$, where by definition

$$y_h^{l+1} = \phi_{l+1} \left(\sum_{q=1}^{N_l} w_{hq}^{l+1} y_q^l \right)$$

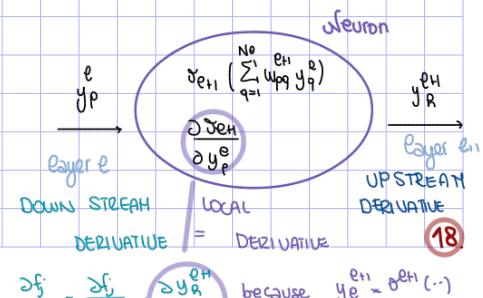
$$f(w, x^n) = y^n$$

BUT! Know we do NOT know this element

We can put in the case $l+1=L \Rightarrow \frac{\partial f}{\partial y^L} = \frac{\partial y^L}{\partial y^L} = I$ identity

So we can do all the steps: $\frac{\partial f}{\partial y^{l-1}} \rightarrow \frac{\partial f}{\partial y^l} \rightarrow \dots \rightarrow \frac{\partial f}{\partial y^l} \rightarrow \frac{\partial y^l}{\partial x} \Rightarrow$

$$\text{And } \frac{\partial y^l}{\partial x} = \frac{\partial}{\partial x} (\phi_l(w^l \cdot x)) = \phi'_l(w^l \cdot x) w^l$$



$$\frac{\partial f}{\partial y^l_p} = \frac{\partial f}{\partial y_h^{l+1}} \cdot \frac{\partial y_h^{l+1}}{\partial y^l_p} \text{ because } y_h^{l+1} = \phi_{l+1}(\dots)$$

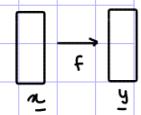
9- CONVOLUTIONAL NEURAL NETWORK

9.1 Increase dimensional input of FFNN

In feed forward neural network we can map $\mathbb{R}^n \rightarrow \mathbb{R}^m$ with the non-linear function

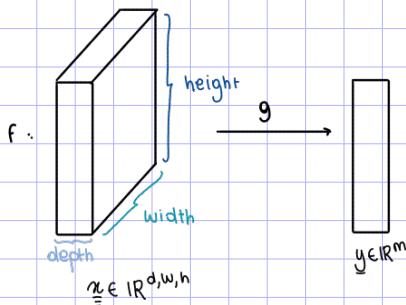
$$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$x \mapsto y = f(x)$$



How we can generalize this kind of NN so that we have an input of bigger dimension, i.e. 3 dimensional input?

Convolutional neural network CNN. This NN is typically used to process image



We want to add a new function so that the 3d input can transform, via g , in a long array so that y can be the input of the FFNN to find the output

the result

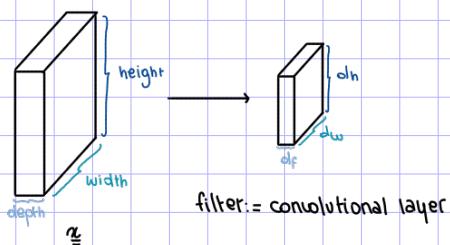
Problem: How to pass from a 3dimensional array to 1dimensional?

Solutions: Ingredients of g

- a. Convolutional layer
- b. Pooling layer
- c. Normalization layer

NB Convolutional NN extend the idea of FFNN but maintaining the structure of the original data

9.2 Convolutional layer → is a filter

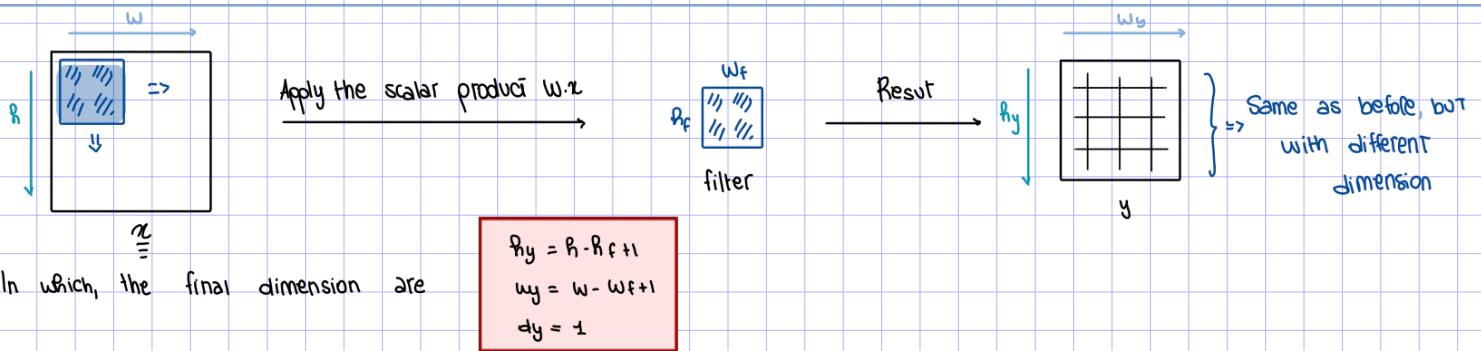


Consider $d_f = d$, $w_f :=$ the weight $\in \mathbb{R}^{df \times wf \times hf}$ and we want to build a filter so that

$$y = w_f \cdot x + w_0$$

, where w_0 is the bias and W the associated matrix

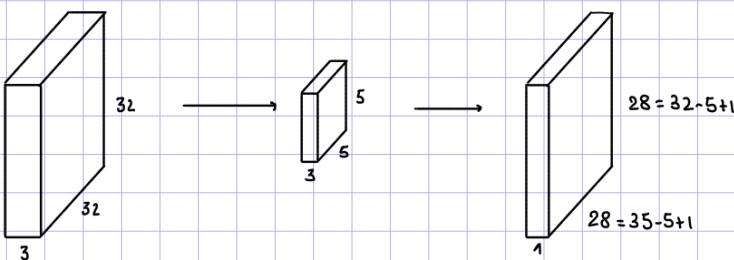
The idea is that the filter can move left-right & up-down on the image
↳ We shift the convolutional filter around x along convolution $y = Wx + w_0$



In which, the final dimension are

$$\begin{aligned} wy &= R_f \cdot w_f + w_0 \\ wy &= W \cdot w_f + w_0 \\ dy &= 1 \end{aligned}$$

Example -



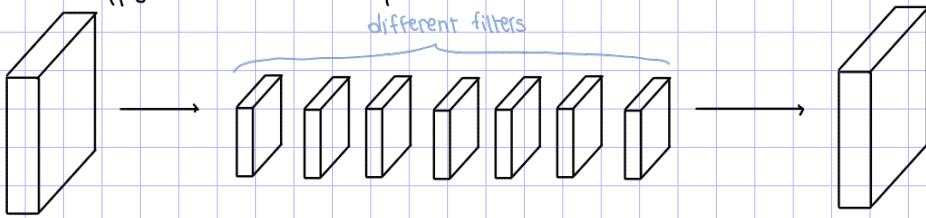
NB In general, we can use more than one convolutional filter, obtaining different filters inputs

Observation - After the convolutional layer you have reduced the dimension of the image, this caused some issue.

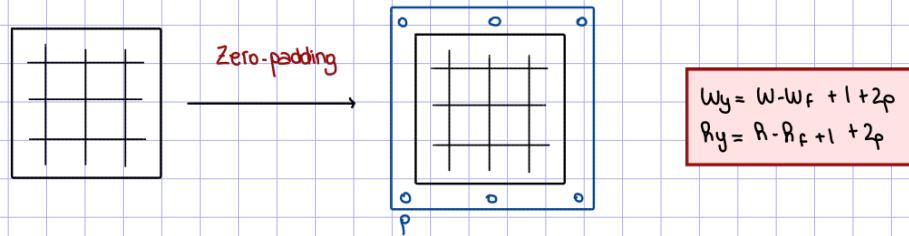
↓ Let's now try to solve them

a. **Channel dimension.** The depth is $d_y = 1$. How we can deal with this?

We can apply different filter in parallel in which I can choose the dimension d called channel to use for all the layer



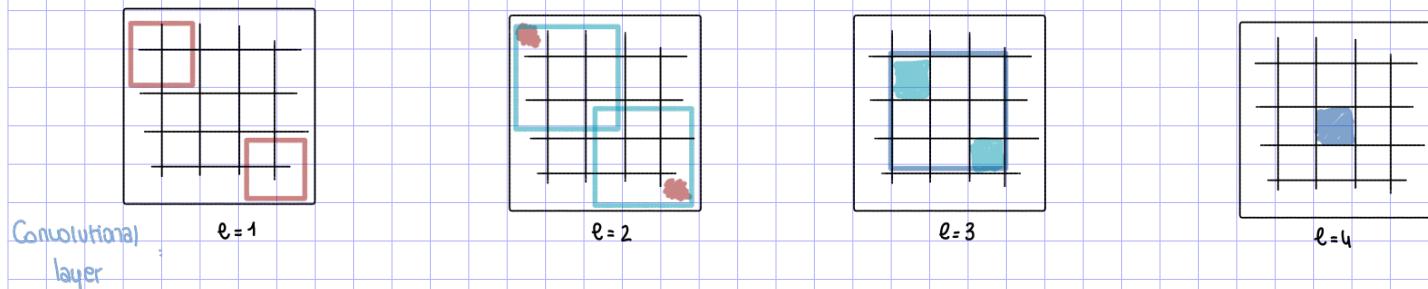
b. **Padding.** If we apply K convolution one after the other it happens that the dimension of x goes to zero rapidly. How we can deal with this? We generalize the convolutional layer with the operation of padding, to avoid constraints on the number of layer in the NN.



c. **Receptive fields.** With convolutional layer the next data is influenced by a particular region of the previous data and by a particular region of the input.

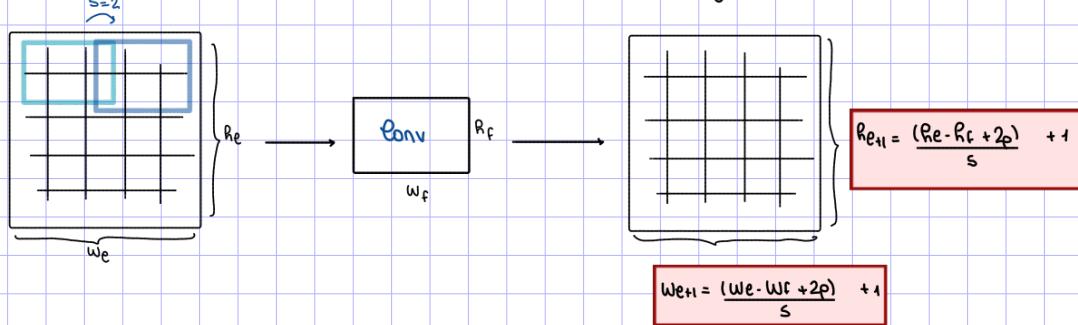
Let's now imagine to apply the filter to my input.

How many layer we need to see all the information of the input?



In this simple example, we need 4 layer, think about bigger one!

We add the last parameter to move it "faster" in the image: **Stride**



Example - Consider

$\Rightarrow x^e: 3 \times 32 \times 32$ is a convolutional filter where $W: 10 \times 3 \times 5 \times 5$ with $s=1, p=2 \rightarrow W = 760$ unknowns to minimize

$\Rightarrow x^{eH}: 10 \times 32 \times 32$

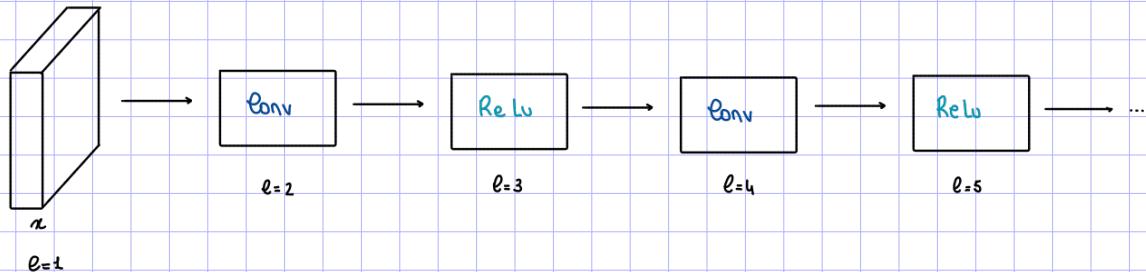
The number of operations to do are ≈ 1 million to perform on each convolutional layer \rightarrow done in parallel via GPU

Definitions.

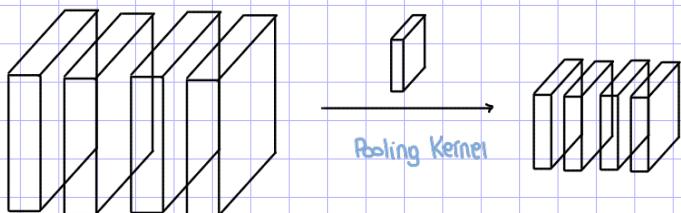
- **Learnable parameters**: weights inside the filter for each channel w_{ij} → learn by the minimization problem
- **Hyperparameters**: not learnable parameters → fixed by the user before the learning phase, i.e. padding, stride

Till now, we have NOT used a non-linearity ...

... Hence, we add the activation function, i.e. **ReLU** to add non-linearity

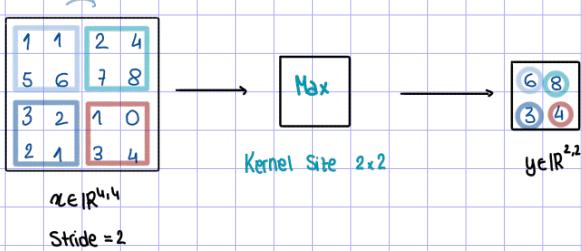


9.3 Pooling layers → the idea is to shrink the input



Let's make an example to better understand

Max pooling



The **pooling Kernels** are just function

i.e.

- Avg pool $K=2, S=2$
- Max pool $K=2, S=2$
- AlexNet pool $K=3, S=2$

are used to downsize the sample without any learnable parameters

The **pooling**, here, just select one element, i.e. the maximum



we apply non linear function use just to reduce the dimension
NB it's NOT mandatory that the pooling layer is NON-linear

With 9.2 I enlarge the number of channel, maintaining w_f and b_f applying convolutional layer, activation & After with 9.3 I modify the w_f and b_f in which I enlarge the number of channel in output

Till we can transform the dimension in an array so that we can apply FFNN

9.4 Normalization → important to find the right solution, used especially to improve the speed of the computation

Consider to minimize $J(z, y) = \|z\|_2^2 + \lambda \|y - \hat{y}\|^2$ → if y is bigger, you focus on it to improve

But we do not want to focus only in one variable ⇒ we normalize so that z & y has the same order of magnitude

$$\hat{z} = \frac{z - E(z)}{\sqrt{\text{Var}(z)}}$$

so that we have $E(\hat{z}) = 0$ & $\text{Var}(\hat{z}) = 1$

$$E(\hat{z}) = \frac{1}{N} \sum z_i$$

$$\text{Var}(\hat{z}) = \frac{1}{N-1} \sum (z_i - E(z))^2$$

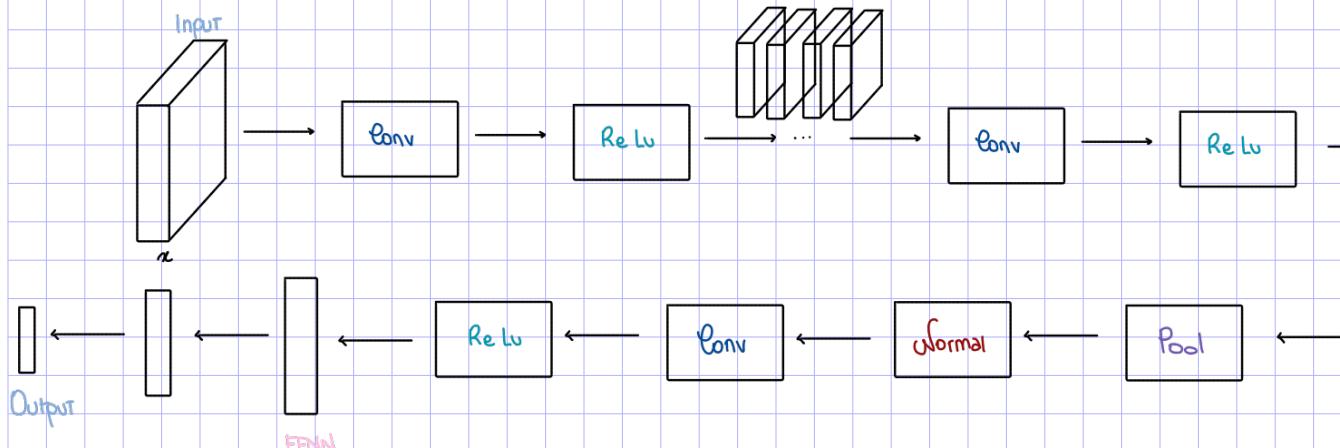
Important: the mean & variance are computed only offline during the training phase, in the validation phase they're fixed

The final output has to belong to the original scale, hence

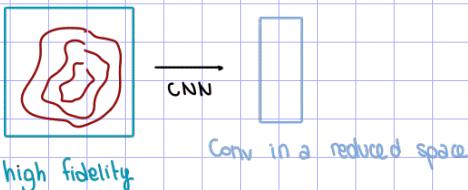
$$y = \sigma \hat{z} + \beta$$
 with σ, β learnable parameters learned in the training phase

↓
capacity of the NN to get back z

9.5 Structure of a CNN



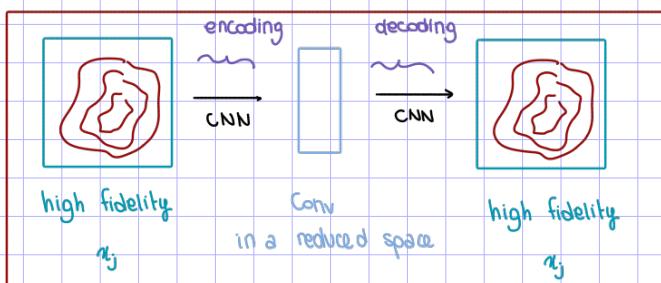
During this process, I do a sort of conservation of the volume during all this step
 we start from the dimension of n^3 and at the end we reach a smaller dimension
 → We can exploit this idea in PDE to reduce the image of the domain of PDE, i.e.



This process is called **encoding**, but we can see also the opposite one

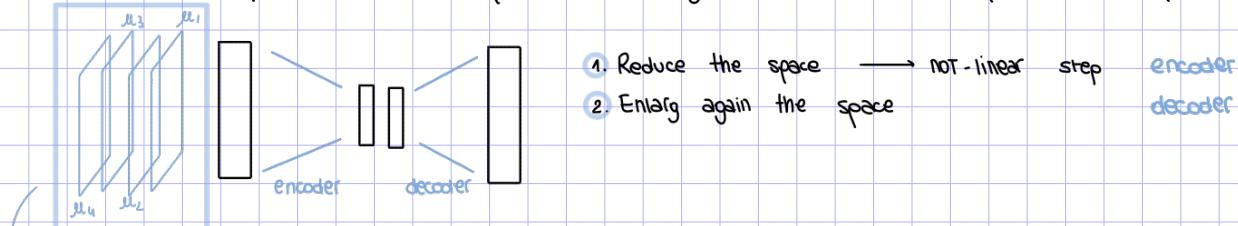
How we can encode the data in a supervised way?

↪ We add another convolutional neural network to enlarg again the dimension



This is an **AUTOENCODER NEURAL NETWORK**
 → After that we can use the decoders to build the final solution of the PDE

The idea is that you want exactly to reproduce the image. The scheme that you have to keep in mind is this one



What happens if I want to work with parametric stuff? We want to exploit the idea of NN, but for solving our problem. Hence we introduce another **NN** that takes the parameters μ and give in result the reduced space

- ↪ Offline you train the net for each different parameter μ and build the encoder
- ↪ Online you exploit the decoder to obtain the result in the huge initial space

↓ The idea is to move from pixel to graph neural network to increase the accuracy of the encoder
 Hence in the μ_i layer there are all the nodes and edges of the graph

NB We do NOT have any kind of idea of the error that we're doing during this step

↪ **BIG PROBLEM:** we have to fix this. but can we?

10 PHYSICS-INFORMED NEURAL NETWORK (PINN) → first way to apply NN to PDE's solver, train in unsupervision way

From the previous lessons, NN are just "magic operator" applies on number to compute classes.

We're moving on unsupervision learning so we do NOT need snapshots to compute the solution
 For more information → 1998: Lagaris first implementation
 → 2018: Karniadakis today best implementation

What is the idea beyond the PINN? Create a NN that is able to achieve \hat{u} → INPUT to NN: point of domain
 Our goal is to solve a PDE → OUTPUT of NN: evaluation of \hat{u}

$$(1) \begin{cases} A(\hat{u}) = f & \text{on } \Omega \subseteq \mathbb{R}^d \text{ open, bounded, Lipschitz} \\ B(\hat{u}) = g & \text{on } \Gamma \subseteq \partial\Omega \end{cases} \rightarrow \text{boundary condition (strong)}$$

with $f: \Omega \rightarrow \mathbb{R}$ and $g: \Gamma \rightarrow \mathbb{R}$ given

solution

- NB → A(.) operator can represents stationary & evolutionary problems, moreover it can be non-linear, time-dependent
 → B(.) is a boundary operator, i.e. trace operator or normal derivate

The idea of PINN is to build a NN that learns how to solve PDE

Hence, we want

$$\begin{aligned} \mathcal{Y}_{\text{NN}}: \Omega &\longrightarrow \mathbb{R} \\ \underline{x} &\longmapsto u(\underline{x}) \text{ approximation of } \hat{u}(\underline{x}) \\ \underline{x} = \begin{pmatrix} \underline{x} \\ y \\ t \\ \mu \end{pmatrix} &\longmapsto u(x, y, t, \mu) \end{aligned}$$

with simulate FOM or using snapshots

it is a network with n inputs and 1 outputs

We need physics to build the NN in an unsupervision way → we train \mathcal{Y}_{NN} unsupervised: we do NOT required snapshots

So, we create loss function using the residuals in the strong formulation

$$\begin{aligned} \bullet L_A^*(\underline{x}, \underline{z}) &= \frac{1}{2} \| f|_{\underline{z}} - A(\underline{x})|_{\underline{z}} \|^2 \\ \bullet L_B^*(\underline{x}, \underline{z}) &= \frac{1}{2} \| g|_{\underline{z}} - B(\underline{x})|_{\underline{z}} \|^2 \end{aligned} \quad \begin{aligned} \text{, Notice that we're NOT using weak formulation, but directly} \\ \text{the strong formulation} \end{aligned}$$

and combined them via convex combination via λ known a-priori

→ the total loss function is $L^*(\underline{x}, \underline{z}) = L_A(\underline{x}, \underline{z}) + \lambda L_B(\underline{x}, \underline{z})$ that does not depends on snapshots

this is an hyperparameter (learn before the training) and balance the order of magnitude of the error

you can find this by hand or automatically somehow

The solution is $\underline{u}_{\text{NN}} := \underset{\underline{z} \in \mathcal{Y}_{\text{NN}}}{\operatorname{argmin}} L^*(\underline{x}, \underline{z})$, hence we have to minimize this function

$\underline{u}_{\text{NN}} \in \mathcal{Y}_{\text{NN}}$

BUT it's NOT easy to find this solution in this way, due to non-linearity of equation of higher dimension (2D, 3D)

↪ The computation of L^* requires to strong hypothesis of f, g, A, B : one way to relax this is using integration

Thus we define a new loss function $L(\underline{x}, \underline{z})$ given by $L(\underline{x}) = L_A(\underline{x}) + \lambda L_B(\underline{x})$ where $L_A(\underline{x}) = \frac{1}{2} \int_{\Omega} \| f - A(\underline{x}) \|^2$ & $L_B(\underline{x}) = \frac{1}{2} \int_{\Gamma} \| g - B(\underline{x}) \|^2$
 and use a quadrature formula to approximate the integrals

How we can minimize the loss function? We have to compute some \underline{z} to evaluate the loss function.

How to get quadrature points? There are different methods.

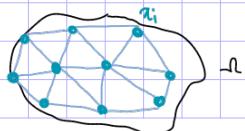
• Meshless method → Similar to monte carlo method

We have the domain Ω , we randomly generate some $\{\underline{z}_i\}$ and evaluate here the loss function



• Mesh generation

We discretize the space / domain and $\{\underline{z}_i\}$ are the vertices of the mesh



1.1 Error PINN —> what we can say about the error? The quality of u_{PINN} depends on many factors

As we have something on the physics, we can measure the error that depends on 3 part:

1. Error capacity of NN: depends on the architecture of NN —> number of layers, neurons

$$\text{ENN} = \inf_{\mathcal{V} \in \mathcal{Y}_{\text{NN}}} \|\hat{u} - v\|_X \quad \rightarrow \text{best approximation of } \hat{u} \text{ in } \mathcal{Y}_{\text{NN}}$$

2. Error associated to the loss function: associated to the minimization process

Eloss

→ related to the minimization problem

3. Error accuracy of the approximation of the residuals: measure how good I am to compute L^* , so

approximate correctly the physics \Rightarrow error integration

$$L^*(v) = \frac{1}{2} \|f - A(v)\|^2 + \frac{1}{2} \|g - B(v)\|^2$$

↓
Here we need strong hypothesis because we use strong formulation

BUT we can relate the loss function ↓

$$\text{weak formulation: } L_w^*(v) = \frac{1}{2\Omega} \int_{\Omega} |f - A(v)|^2 + \frac{1}{2\Omega} \int_{\Omega} |g - B(v)|^2 \quad \rightarrow \text{disadvantages: we have to solve integral}$$

→ use this formulation in the project!

Hence, the final error of the PINN is

$$\|u_{\text{PINN}} - \hat{u}\|_X \leq \text{ENN} + \text{Eloss} + \text{Equad}$$

To relax more the hypothesis, we can study the

1.2 Variational PINN —> can we do something more? Yes, to relax the hypothesis we can use the weak form of (1)

We're introducing the variational (weak) form in the loss function

$$L_{ww}^*(w) = \frac{1}{2\Omega} \int_{\Omega} f(w) - a(v, w) + \dots \quad \forall w \in VR$$

① POD NEURAL NETWORKS → reference paper: Non-intrusive reduced order modeling of non-linear problems using neural networks, Hesthaven & Ubbiali - 2018 JCP

The goal of the paper is to provide a fast tool to deal with non-linear problems
 ↓ How to build this tool?
 They avoid projection online!

Problem formulation for non-linear problems.

Given a parameter $\mu \in \mathcal{P}$, find $u(\mu) \in V$ s.t. $g(u(\mu); \mu) = 0$ (1) strong version of the equation

i.e. Navier-Stokes $-\Delta u + (u \cdot \nabla u) + \nabla p = f$ non-linear problem

when we discretize → $Au + C(u)u + Bp = f$

↓ there is this dependence that represents the non-linearity

assuming $f=0$, put that in the left-part

In weak formulation, we have $g(u(\mu), v; \mu) = 0 \quad \forall v \in V$ (2)

Considering the discretized problem

→ $V_s \subset V$ and given the parameter $\mu \in \mathcal{P}$, find $u_s(\mu) \in V_s$ s.t. $g(u_s(\mu), v; \mu) = 0 \quad \forall v \in V_s$

Algebraically we can build the residual vector $G_f(u_s(\mu), \mu) = 0 \in \mathbb{R}^{N_s}$ where $\dim(V_s) = N_s$

How G is built? $G_f(\bar{u}_s(\mu); \mu) := g(u_s(\mu), \varphi_i; \mu)$ where $\{\varphi_i\}_{i=1}^{N_s}$ is a FOM basis

End of FOM level

What should I do at the ROM level?

- ① Build the space with POD:
 - Collect snapshots
 - Solve the generalized eigenvalue problem over covariance matrix
 - Take eigenvectors related to the N -largest eigenvalues
 - Build the basis matrix
 - Project the matrices & the vectors → Attention! To build

Here, we cannot separate the variable, so we cannot perform this

- ② Solve the reduce system is NOT fast in N_s → NOT good

$$\mathbb{R}^N \ni G_n(u_n(\mu), \mu) = B^T G_f(B u_n(\mu), \mu)$$

ISSUE: for each μ I have to assemble again the matrix?

IDEA: can I avoid 2. using NL? We want to avoid the 2. step. How to do that?

POD-NN is an algorithm that exploits the direct projection of the snapshots in the reduce space to create a map

Π^{NN} (Neural Network) s.t.

$$\mu^* \longmapsto \Pi^{NN}(\mu^*) = u_N(\mu^*) \in \mathbb{R}^N$$

I am learning the reduced coefficients

To build the net we need a geometric interpretation of a reduced problem

→ ROMs build a reduced basis matrix $B \in \mathbb{R}^{N_s \times N}$ where $B = [\varphi_1 \ | \ \dots \ | \ \varphi_N]$ → φ_i : basis function

→ Let us consider a linear problem: the simplest way to define direct projection

$$\text{i.e. } \begin{cases} -\Delta u = f \\ \text{boundary condition} \end{cases} \quad \xrightarrow{\text{weak formulation}} \int_L g \nabla u \nabla v \, dx = \int_L f v \, dx$$

This can be seen as minimization problem ↓ FOM $Au=f$ where $A_{ij} = \int_L \nabla \varphi_i \cdot \nabla \varphi_j \, dx$ for $\{\varphi_i\}_{i=1}^{N_s}$ basis of the FOM
 matrix related to the grad-gradi operator

↓ H_h -semi norm equivalent to H_1 norm

Considering, $A = X_g \in \mathbb{R}^{N_s \times N_s}$ this is the inner product matrix, that during the labs we called stiffness matrix

↓ A small version of my matrix & vector can be written like that: $B^T AB = A_N$
 $B^T f = f_N$

11.1 Algebraic properties of the Galerkin-ROM

The goal here is to find u_N without solving the system

↓

Hence, let us define the vector representation of the error $e_g(\mu) = u_g(\mu) - B^T u_N(\mu)$

The ROM residual, on the contrary, can be written as $r_g(\mu) = f - AB u_N(\mu)$ → take the small vector and project back
depends on u_N and μ

The following hold

1. $Ae_g(\mu) = Au_g(\mu) - AB u_N(\mu) = f - AB u_N(\mu) = r_g(u_N; \mu)$ → The error goes to zero if $r_g(u_N)$ of the ROM solution is zero
2. $B^T Ae_g(\mu) = B^T f = f_N$
3. $B^T r_g(u_N; \mu) = B^T f - B^T AB u_N = f_N - A_N u_N$ → $B^T f(u_N, \mu) \rightarrow 0$ if B is a "good" representation
 $f_N - A_N u_N \rightarrow 0$ is the reduced system (RP)

I do NOT want to solve $f_N - A_N u_N = 0$, but I want to express u_N in terms of $u_g(\mu)$ using the properties 1, 2, 3.

Let us take the ROM problem evaluated in u_N in the ROM space

$$AB u_N(\mu) = f = Au_g(\mu)$$

↓
if B is good by ROM definition

Now we project, everything, in the reduced space

$$B^T AB u_N(\mu) = B^T f = B^T Au_g(\mu)$$

→ Then, $u_N(\mu) = (B^T AB)^{-1} B^T Au_g(\mu)$ → if we have $u_g(\mu)$, we have a direct way to compute u_N
↓
IP

and IP is the direct projection of $u_g(\mu)$ in the ROM space

Finally, we can train our net, because we can define the loss

$$\mu^* \xrightarrow{\pi^{NN}} u_N(\mu^*) \quad \text{BUT How to train it?}$$

- The input is μ
- The output is $u_N(\mu) \in \mathbb{R}^N$
- The loss could be $\mathcal{L} = \sum_{\mu} \|\pi^{NN}(\mu) - u_N(\mu)\|_2^2 = \sum_{\mu} \|\pi^{NN} - \text{IP } u_g(\mu)\|_2^2 = \sum_{\mu} \|\pi^{NN} - (B^T X_g B)^{-1} B^T X_g u_g(\mu)\|_2^2$
↓ number of the snapshots ↓ basis function matrix related to POD
- Hence, now we have all we need!

OFFLINE

- Snapshots collection
- POD including Non linearity
- Train the net → You're losing the theorem
that tells us that the projection
- + I exploit POD is the best thing
- The offline phase is longer

ONLINE

- Evaluation of the net: $\mu^* \xrightarrow{\pi(\mu^*)}$

- + Very fast: give me the parameter μ , I'll give you the representation
- No control on the accuracy

12- STOKES EQUATIONS in H.O.R

→ we're in a linear setting

Consider the domain Ω and we want to study → The fluid velocity $u = (u_1, u_2)$ on Ω
 ↓ find → The pressure

Make some assumption

H1. The fluid is incompressible → volume constant in time, it holds $\nabla \cdot u = 0$

H2. Do NOT consider the time: stationary system

Let us define the system

$$\begin{cases} -\gamma \Delta u + \nabla p = f & \text{in } \Omega \\ \nabla \cdot u = 0 & \text{in } \Omega \\ u = g & \text{in } \Gamma_D \subseteq \partial \Omega, \neq \emptyset \\ -\gamma (\nabla \cdot u)^T \gamma + p \gamma = \psi & \text{in } \Gamma_N := \partial \Omega \setminus \Gamma_D \end{cases}$$

γ : viscosity, that for WE it is a constant

→ strong boundary condition: we need to have a unique solution (Dirichlet)

→ Neuman boundary condition

If this boundary exist, we fix some condition on the pressure

So if this boundary does NOT exist, we solve the system and find the solution of p up to a constant

i.e. If $\Gamma_N = \emptyset$, for example you can state → $p=0$ on a point of Ω to solve the uniqueness of the solution
 $\int_{\Omega} p = 0$ on Ω

↓ Definition of the term

The definition of $\nabla u \in \mathbb{R}^{2,2}$ → $(\nabla u)_{ij} = \frac{\partial u_i}{\partial x_j} = (J_u)_{ij}$

because it's the transpose of the jacobian

The divergence is $\nabla \cdot u = \sum_{i=1}^2 \partial_{x_i} u_i = \partial_x u_1 + \partial_y u_2$

The last definition is $\Delta u := \nabla(\nabla \cdot u) - \nabla \times (\nabla \times u) = -\nabla \times (\nabla \times u)$, where the vectorial product is defined as

$$\nabla \times u := \begin{vmatrix} i & j & k \\ \partial_x & \partial_y & \partial_z \\ u_1 & u_2 & u_3 \end{vmatrix} = i(-\partial_y u_2) - j(-\partial_x u_1) + k(\partial_x u_2 - \partial_y u_1) = (0, 0, \partial_x u_2 - \partial_y u_1)$$

$$\downarrow \text{Hence, doing the computation we have } \nabla \times (\nabla \times u) = \begin{vmatrix} i & j & k \\ \partial_x & \partial_y & \partial_z \\ 0 & 0 & \partial_x u_2 - \partial_y u_1 \end{vmatrix} = i(\partial_y \partial_x u_2 - \partial_y \partial_y u_1) - j(\partial_x \partial_x u_2 - \partial_x \partial_y u_1) + k 0$$

$$\text{Then } \nabla \cdot u = 0 \Leftrightarrow \partial_x u_1 + \partial_y u_2 = 0 \quad (=) \quad \partial_x u_1 = -\partial_y u_2$$

$$\rightarrow (-\partial_y^2 u_1 + \partial_y \partial_x u_2, -\partial_x^2 u_2 + \partial_x \partial_y u_1, 0) = (-\partial_y^2 u_1, -\partial_x^2 u_1, -\partial_x^2 u_2 - \partial_y^2 u_1, 0) = -(\Delta u_1, \Delta u_2)$$

we remove u_2 we remove u_1

Hence, returning to $\Delta u := \nabla(\nabla \cdot u) - \nabla \times (\nabla \times u) = (\Delta u_1, \Delta u_2)$ → we can split and compute the

Laplacian of the 2 variable separately

Can we go to the weak form? Yes, rewrite the system

$$\begin{cases} -\int_{\Omega} \gamma \Delta u \cdot \eta + \int_{\Omega} \nabla p \cdot \eta = \int_{\Omega} f \cdot \eta & \text{in } \Omega \quad \forall \eta \in ? \\ -\int_{\Omega} (\nabla \cdot u) q = 0 & \text{in } \Omega \quad \forall q \in ? \end{cases}$$

in which space? we now do NOT know, it is a test function for velocity

another space that we do NOT know, it is a test function for the pressure

Now, we want to define the Green's identities

$$(1) \int_{\partial \Omega} \Psi \nabla \phi \cdot \eta = \int_{\Omega} \Psi \Delta \phi + \int_{\Omega} \nabla \Psi \cdot \nabla \phi \quad \rightarrow \text{in Stokes equation: } \nabla \phi = u, \Psi = p$$

for square matrix and obtain a scalar

$$(2) \int_{\partial \Omega} \eta \cdot [(\nabla u)^T \eta] = \int_{\Omega} \eta \cdot \Delta u + \int_{\Omega} \nabla u \cdot \nabla \eta \quad \text{where } A:B := \sum_{i,j=1}^d A_{ij} B_{ij} \text{ double dot product}$$

Applying to the system, we obtain

$$\begin{cases} \int_{\Omega} \gamma \nabla u \cdot \nabla \eta - \int_{\Omega} p \nabla \cdot \eta - \int_{\Omega} \gamma \nabla \cdot ((\nabla u)^T \eta) + \int_{\Omega} p \nabla \cdot \eta = \int_{\Omega} f \cdot \eta \quad \forall \eta \\ -\int_{\Omega} (\nabla \cdot u) q = 0 \end{cases}$$

BUT we are in a weak formulation, hence we want that $\mathbf{v}_{|\Gamma_D} = \mathbf{0} = (0, 0)^T$ so we can consider that

$$-\int_{\partial\Omega} \mathbf{r} \cdot ((\nabla \mathbf{u})^T \mathbf{r}) + \int_{\partial\Omega} p \mathbf{r} \cdot \mathbf{n} \xrightarrow{\begin{array}{l} \mathbf{0} \text{ in } \Gamma_0 \\ \neq \mathbf{0} \text{ in } \Gamma_N \end{array}} \mathbf{0}$$

↓ Rewrite the system considering this observations → Weak form Stokes equations

$$\begin{cases} \int_{\Omega} \mathbf{r} \nabla \mathbf{u} : \nabla \mathbf{r} - \int_{\Omega} p \nabla \mathbf{r} = \int_{\Omega} \mathbf{f} \cdot \mathbf{r} - \int_{\Gamma_N} \mathbf{Y} \cdot \mathbf{r} & \forall \mathbf{r} \in H_0^1(\Omega) \\ -\int_{\Omega} (\nabla \cdot \mathbf{u}) q = 0 & \forall q \in L^2(\Omega) \end{cases}$$

velocity
pressure
↓ This terms are "the same"

Let's now deal with the space: 1) pressure space $p \in L^2(\Omega)$, and also the divergence has to belong to the same space
 ↓ we have NOT to derive it, no gradient, very weak hypothesis

2) forcing term $\mathbf{f} \in [L^2(\Omega)]^2$

because they are array

3) velocity space $\mathbf{u} \in [H_0^1(\Omega)]^2 := \{ \mathbf{r} \in [H^1(\Omega)]^2 : \mathbf{r}_{|\Gamma_D} = \mathbf{0} \}$
 $\mathbf{v} \in [H_0^1(\Omega)]^2 := \{ \mathbf{v} \in [H^1(\Omega)]^2 : \nabla \cdot \mathbf{v}_{|\Gamma_N} = 0 \}$

- Remark - R1. If $\Gamma_N = \emptyset$, we have to set conditions to have a unique solution → $p \in L^2(\Omega) := \{ r \in L^2(\Omega) : \int_{\Omega} r = 0 \}$
 R2. Implementing this problem, we define $\mathbf{u} = \mathbf{u}_0 + \mathbf{G}$, where $\mathbf{u}_0 \in [H_0^1(\Omega)]^2$, $\mathbf{G}_{|\Gamma_D} = \mathbf{g}$ is the detection of the Dirichlet condition

Make all these informations together...

$$\begin{cases} \int_{\Omega} \mathbf{r} \nabla \mathbf{u}_0 : \nabla \mathbf{r} - \int_{\Omega} p \nabla \cdot \mathbf{r} = \int_{\Omega} \mathbf{f} \cdot \mathbf{r} - \int_{\Gamma_N} \mathbf{Y} \cdot \mathbf{r} - \int_{\Omega} \mathbf{G} : \nabla \mathbf{r} & \forall \mathbf{r} \in [H_0^1(\Omega)]^2 \equiv V \\ -\int_{\Omega} (\nabla \cdot \mathbf{u}) q = 0 & \forall q \in L^2(\Omega) \equiv Q \end{cases}$$

Introducing the bilinear form, to write everything in a compact way, we have

where $a(\cdot, \cdot) : V \times V \rightarrow \mathbb{R}$, $a(u, v) = \int_{\Omega} \mathbf{r} \nabla u \cdot \nabla v$, $b(\cdot, \cdot) : V \times Q \rightarrow \mathbb{R}$, $b(v, q) := -\int_{\Omega} q \nabla \cdot v$
 $F(\cdot, \cdot) : V \rightarrow \mathbb{R}$, $F(v) = \int_{\Omega} \mathbf{f} \cdot \mathbf{v} - \int_{\Gamma_N} \mathbf{Y} \cdot \mathbf{v} - \int_{\Omega} \mathbf{G} : \nabla \mathbf{v}$

$$\begin{cases} a(\cdot, \cdot) - b(\cdot, \cdot) = F(\cdot, \cdot) & \forall v \in V \\ -b(\cdot, \cdot) = 0 & \forall q \in Q \end{cases}$$

[VP]
↓ Variational problem

↓ Hence we obtain the saddle point problem

Find $\mathbf{u} \in V$, $p \in Q$ s.t.

$$\begin{cases} a(\mathbf{u}, \mathbf{v}) + b(\mathbf{v}, p) = \langle \mathbf{F}, \mathbf{v} \rangle_V & \forall \mathbf{v} \in V \\ b(\mathbf{u}, q) = 0 & \forall q \in Q \end{cases}$$

Theorem - Brezzi's Theorem

Hypothesis - H1. $a(\cdot, \cdot)$ is continuous in V → $\exists \gamma_1 > 0 : |a(\mathbf{u}, \mathbf{v})| \leq \gamma_1 \|\mathbf{u}\|_V \|\mathbf{v}\|_V \quad \forall \mathbf{u}, \mathbf{v} \in V$

H2. $a(\cdot, \cdot)$ is coercive in V_0 → $\forall \mathbf{v}_0 \in V_0 : \mathbf{b}(\mathbf{v}, \mathbf{v}) = 0 \quad \forall q \in Q \subset V$

equivalence

$$\exists \alpha_1 > 0 : a(\mathbf{u}, \mathbf{u}) \geq \alpha_1 \|\mathbf{u}\|_V^2 \quad \forall \mathbf{u} \in V_0 \iff \sup_{\mathbf{w} \in V_0} \frac{a(\mathbf{w}, \mathbf{u})}{\|\mathbf{w}\|_V} \geq \alpha_1 \|\mathbf{u}\|_V \quad \forall \mathbf{u} \in V$$

H3. $b(\cdot, \cdot)$ is continuous on $V \times Q$ → $\exists \gamma_2 > 0 : b(\mathbf{v}, q) \leq \gamma_2 \|\mathbf{v}\|_V \|\mathbf{q}\|_Q \quad \forall \mathbf{v} \in V \quad \forall q \in Q$

H4. $b(\cdot, \cdot)$ is inf-sup stable in $V \times Q$ → $\exists \beta_0 > 0 : \inf_{q \in Q} \sup_{\mathbf{v} \in V} \frac{b(\mathbf{v}, q)}{\|\mathbf{v}\|_V} \geq \beta_0 > 0 \iff \sup_{\mathbf{v} \in V} \frac{b(\mathbf{v}, q)}{\|\mathbf{v}\|_V} \geq \beta_0 \|q\|_Q \quad \forall q \in Q$

extension of coercivity for bilinear form in 2 dimension

Thesis - Then $\exists! (\mathbf{u}, p) \in V \times Q$ solution of [VP]

↳ Unique continuous solution

1. $\|\mathbf{u}\|_V \leq \frac{1}{\alpha_1} \|\mathbf{F}\|_V$

2. $\|p\|_Q \leq \frac{1}{\beta_0} \left(1 + \frac{\gamma_2}{\alpha_1} \right) \|\mathbf{F}\|_V$ ↳

What happen in the Galerkin approximation? See lecture 11 for more details

1.2.1 Supermixer operator

Is the operator define such that $T: Q \longrightarrow V$, where $\forall q \in Q, \forall v \in V: (Tq, v)_V = b(w, q) \quad \forall w \in V$ that is an implicite definition of the supermixer operator T

↓
This is useful only in the discrete space, not in the continuity one

Property - $\forall q \in Q, \exists Tq \in \arg \sup_{w \in V} \frac{b(w, q)}{\|w\|_V}$ when you apply the operator to the pression you obtain the velocity

1.2.2 What happen when we want to study the Stokes equation in a discrete space?

We discretize the domain V_h, Q_h with the same tessellation S to study the discrete space

↓
What happen to the H^1_0 of the Brezzi Theorem?

Consider $\beta_b := \inf_{q \in Q_h} \sup_{v \in V_h} \frac{b(v, q)}{\|v\|_V \|q\|_Q} > 0$, if I restrict Q_h

↓ We have to study the discrete inf-sup stability

- The pressure resulting from the inf-sup belongs to Q_h ✓
- The pressure is NOT in Q_h : it is inf ✓
- The velocity is in V_h : ✓
- The velocity is NOT in V_h : this is really a big problem! ✗
- The value β_b could be 0 → Huge Problem ✗

Hence, NOT ALL THE DISCRETIZATIONS ARE GOOD, then the usual space used is the Taylor-Hood FEM spaces

$$\begin{cases} V_h \text{ FEM } K=2 \\ Q_h \text{ FEM } K=1 \end{cases}$$



basis function

$$\psi_i : \begin{cases} 1 & \text{in the vertex } i \\ 0 & \text{in the other points} \end{cases}$$

basis also in the middle
of the edge

This is a good and stable tessellation, moreover the error of the velocity goes to zero faster wrt the one on the q

$$(p_h = p_{\text{exact}})$$

1.2.3 Reduced Model Order for Stokes equation:

we want to reduce V_h in the classic way, hence we can take the snapshots $\{u_m = u_{S, m}\}$
The same observation as before for the inf-sup value, How can we effort? Add v to the space so that we have
no problem with the value of inf-sup

↓ We have to build the reduced space in order that the inf-sup stability holds. indeed

in high fidelity snapshot the inf-sup stability is given by the choice of Q_h, V_h

Add other velocity to the ones find with POD until the new basis $\{u_m\}$ satisfies the condition of inf-sup

To do that we exploit the supermixer operator. Let's see the discrete version

$$T: Q_h \longrightarrow V_h \text{ s.t. } \forall q_h \in Q_h, \quad Tq_h \in V_h$$

$$\forall w_h \in V_h, \quad \langle Tq_h, w_h \rangle_V = b(w_h, q_h)$$

↓ This now is a linear system

But, it holds the same property

$$Tq_h \in \arg \sup_{w_h \in V_h} \frac{b(w_h, q_h)}{\|w_h\|_V}$$

↓ This is good because in that way it guarantees
the existence of the inf-sup

Corollary - The constants

$$\beta_s := \inf_{q_h \in Q_h} \sup_{w_h \in V_h} \frac{b(w_h, q_h)}{\|w_h\|_V \|q_h\|_Q} := \inf_{q_h \in Q_h} \frac{\|Tq_h\|_V}{\|q_h\|_Q}$$

↓ This holds thanks to the properties
Then $\beta_s \geq \beta_b > 0$ ✓

So that we can build the reduced spaces

Fixing $Q_h \subset Q$, $Q_h := \text{span} \{ p_j \}_{j=1}^m$ find, i.e. with POD, the right choice for the space of the velocity is

$$V_h := \text{span} \{ w_j \}_{j=1}^m \cup \text{span} \{ Tp_j \}_{j=1}^m$$

↓ this is inf-sup stable & $\beta_s \geq \beta_b > 0$ ✓

NB This approach can be done also with other model of PDE to deal with the same problems

12.4 Remark on Stokes High-Fidelity implementation

Lab 10

We have this specific problem, that is the high-fidelity problem

$$\begin{cases} a(u_s, v_s) - b(v_s, p_s) = F(v_s) & \forall v_s \in V_s \\ -b(u_s, q_s) = 0 \end{cases}$$

In FEM, we fix the basis function

$$\Rightarrow Q_s: \{ \varphi_k \}_{k=1}^{N_s^p}$$

$$a(u_s, v_s) := \int_{\Omega} \gamma \nabla u_s : \nabla v_s, \quad b(u_s, q_s) := \int_{\Omega} q_s \nabla \cdot u_s, \quad F(v_s) = \int_{\Omega} f \cdot v_s$$

How to discretize with the FEM?

↪ We need to define a smart basis

$$\text{Standard FEM basis } P_2(\Omega) \quad \psi_j = \{\psi_j^1, \psi_j^2, \psi_j^3 \} \in \{[\varphi_i, 0], [0, \varphi_i]\}$$

$$\text{where } \{\varphi_i\}_{i=1}^{N_s^u} \text{ is the FEM standard basis } N_s^u = \frac{N_s}{2}$$

$$\left\{ \begin{array}{l} [\varphi_1, 0] \\ [\varphi_2, 0] \\ \vdots \\ [0, \varphi_{N-1}] \\ [0, \varphi_N] \end{array} \right\} = \psi_j$$

$$\text{Consider now } A_{ij} = a(\psi_j, \psi_i) = \int_{\Omega} \gamma \nabla \psi_i : \nabla \psi_j \Rightarrow \text{if } k \in \{0, \dots, N_s\} \quad \psi_k = (\varphi_k, 0)$$

(1)

we can split the computation thanks to the basis that we choose

$$\begin{aligned} i \in \{0, \dots, N_s\}, j \in \{0, \dots, N_s\} \quad A_{ij} &= \int_{\Omega} \gamma \nabla \begin{pmatrix} \varphi_i \\ 0 \end{pmatrix} : \nabla \begin{pmatrix} \varphi_j \\ 0 \end{pmatrix} = \int_{\Omega} \gamma \begin{bmatrix} \partial_x \varphi_i & 0 \\ \partial_y \varphi_i & 0 \end{bmatrix} : \begin{bmatrix} \partial_x \varphi_j & 0 \\ \partial_y \varphi_j & 0 \end{bmatrix} = \int_{\Omega} \gamma (\partial_x \varphi_i \partial_x \varphi_j + \partial_y \varphi_i \partial_y \varphi_j) \\ &= \int_{\Omega} \gamma \nabla \varphi_i \cdot \nabla \varphi_j \quad \text{gradient } \varphi_i \cdot \text{dot gradient } \varphi_j \end{aligned}$$

$$i \in \{0, \dots, N_s\}, j \in \{N_s+1, \dots, 2N_s\} \quad A_{ij} = \int_{\Omega} \gamma \begin{bmatrix} \partial_x \varphi_i & 0 \\ \partial_y \varphi_i & 0 \end{bmatrix} : \begin{bmatrix} 0 & \partial_x \varphi_j \\ 0 & \partial_y \varphi_j \end{bmatrix} = 0 \quad \text{by definition dot product}$$

$$i \in \{N_s+1, \dots, 2N_s\}, j \in \{0, \dots, N_s\} \quad A_{ij} = 0$$

$$i \in \{N_s+1, \dots, 2N_s\}, j \in \{N_s+1, \dots, 2N_s\} \quad A_{ij}^2 = \int_{\Omega} \gamma \nabla \varphi_i : \nabla \varphi_j$$

$$(2) \quad B_{kj} = b(\psi_j, \varphi_k) = \int_{\Omega} \varphi_k \nabla \cdot \psi_j = \left(\int_{\Omega} \varphi_k \partial_x \varphi_j^1 \right)_{j \in \{1, \dots, N_s\}} \left(\int_{\Omega} \varphi_k \partial_y \varphi_j^2 \right)_{j \in \{N_s+1, \dots, 2N_s\}}$$

$$= \left(\int_{\Omega} \varphi_k [1, 0] \cdot \nabla \varphi_j^1 \right)_{j \in \{1, \dots, N_s\}} \left(\int_{\Omega} \varphi_k [0, 1] \cdot \nabla \varphi_j^2 \right)_{j \in \{N_s+1, \dots, 2N_s\}}$$

$(B_{FEM}^1)_{kj}$

$(B_{FEM}^2)_{kj}$

$$(3) \quad q_i = F(\psi_i) = \int_{\Omega} f \cdot \psi_i = \left(\int_{\Omega} f_1 \varphi_i^1 \right)_{i \in \{1, \dots, N_s\}} \left(\int_{\Omega} f_2 \varphi_i^2 \right)_{i \in \{N_s+1, \dots, 2N_s\}}$$

q_{FEM}^1

q_{FEM}^2

Hence, the Stokes linear system becomes

$$\begin{array}{c} \begin{array}{ccc} & N_s^u & N_s^p \\ \begin{array}{c} N_s^u \\ N_s \\ N_s^p \end{array} & \xrightarrow{\quad A^1 \quad} & \begin{array}{c} B_{FEM}^1 \\ O \\ B_{FEM}^2 \end{array} \\ \begin{array}{c} A^2 \\ O \\ B_{FEM}^1 \\ B_{FEM}^2 \end{array} & \xrightarrow{\quad N_s^p \quad} & \begin{array}{c} u_1 \\ u_2 \\ p \end{array} \end{array} & = & \begin{array}{c} \begin{array}{c} N_s^u \\ N_s \\ N_s^p \end{array} \\ \begin{array}{c} q_{FEM}^1 \\ q_{FEM}^2 \\ O \end{array} \\ \begin{array}{c} N_s^u \\ N_s \\ N_s^p \end{array} \end{array} \\ \hline & J & \end{array}$$

This is the solution that we want to find

13 - NON LINEAR PROBLEMS

The non-linearity is in the bilinear form $a(u, v) = f(v)$, hence $a(u, v)$ is non linear in u , i.e. $a(u, v) = \max(u, v)$

↓
We have to linearize the problem, i.e. using Newton Method

13.1 Newton Method: linear system resolver

It's an iterative method in which we assume $f(u) = 0$ and the iterative step is $u_{n+1} = u_n + \delta u$

Exploits the Taylor expansion centered in u_n , we can compute, using $\delta u = u_{n+1} - u_n$

$$f(u_{n+1}) = f(u_n) + J_f(\delta u)|_{u_n} \delta u + O(J_u^2)$$

Jacobian first

order difference

negligible second order term

This term has to go to zero

Hence → while $\delta u \approx 0$ do → $\delta u \approx 0$ or $f(u_{n+1}) \approx 0$ are two equivalent stopping criteria

$$J_f(\delta u)|_{u_n} \delta u = -f(u_n)$$

direction v

Definition of the Jacobian: $J_f[v] := \lim_{h \rightarrow 0} \frac{f(x_0 + hv) - f(x_0)}{h}$

if we can compute this, we can solve the system

12

Gâteaux differentiable: if exists the limit of the incremental ratio

In the laboratory we'll study the **Burger equation** $-\Delta u + u \nabla u = g$, where $\nabla u := \partial_x u + \partial_y u$ where u is scalar
sort of with scalar variable

$\nabla \cdot u := \partial_x u + \partial_y u$ where u is scalar
definition of this strange divergence

The first thing to do is to go to the variational form to write the correct form of f

$$f: \int_{\Omega} \nabla u \cdot \nabla v + \int_{\Omega} u \nabla \cdot u v - \int_{\Omega} g v = 0$$

f_1 f_2 f_3

Now we can split f in f_i so that we can compute the Jacobian

→ Both f_1 & f_3 are linear function, so let's try to compute the Jacobian

$$\text{Diffusion } J_{f_1}(\delta u)|_{u_N} = \lim_{h \rightarrow 0} \frac{f_1(u_N + h \delta u) - f_1(u_N)}{h} = \lim_{h \rightarrow 0} \frac{f_1(u_N) + h f'_1(\delta u) - f_1(u_N)}{h} = f'_1(\delta u) = \int_{\Omega} \nabla u \cdot \nabla v \quad \text{This holds for all linear / function}$$

$J_{f_3}(\delta u)|_{u_N} = 0$ because it doesn't appear u in the function

→ For f_2 , that is not linear, we have to compute the Jacobian in a different way

$$\begin{aligned} J_{f_2}(\delta u)|_{u_N} &= \lim_{h \rightarrow 0} \frac{f_2(u_N + h \delta u) - f_2(u_N)}{h} = \lim_{h \rightarrow 0} \frac{\int_{\Omega} (u_N + h \delta u) \nabla (u_N + h \delta u) v - \int_{\Omega} u_N \nabla u_N v}{h} \\ &= \lim_{h \rightarrow 0} \frac{h}{h} \left(\int_{\Omega} u_N \nabla u_N v - \int_{\Omega} u_N \nabla u_N v \right) + \frac{h^2 \dots}{h} + \int_{\Omega} u_N \nabla \cdot (\delta u) v + \int_{\Omega} \delta u \nabla \cdot u_N v \\ &\quad \text{advection term} \qquad \qquad \qquad \text{reaction term} \\ &\quad \text{negligible} \end{aligned}$$

The integral is linear
We can split

$$\int_{\Omega} (u_N \nabla u_N) \nabla v$$

$$\int_{\Omega} (\delta u \nabla \cdot u_N) v$$

Newton Schema

$$J_f[\delta u]|_{u_N} \cdot \delta u = -f(u_N)$$

$$\rightarrow \left(\int_{\Omega} \nabla \delta u \cdot \nabla v + \int_{\Omega} (u_N \nabla u_N) \cdot \nabla \delta u v + \int_{\Omega} (\delta u \nabla \cdot u_N) \delta u v \right) = - \int_{\Omega} \nabla u_N \cdot \nabla v - \int_{\Omega} u_N \nabla \cdot u_N v + \int_{\Omega} g v \quad \forall v \in H^1_0(\Omega)$$

With respect to Lab 12.

$$\int_{\Omega} h(u_N) \ell(\delta u) m(v) \approx \sum_{q=1}^Q h(u_N, x_q) \ell(\delta u, x_q) m(v, x_q) w_q$$

$$\text{burger_non_linear_c} \Rightarrow h(u_N) = \nabla \cdot u_N = \partial_x u_N|_{x_q} + \partial_y u_N|_{y_p}$$

Algorithm 3 Newton's Method for Nonlinear PDEs

- 1: **Input:** Nonlinear PDE operator R , Tolerance ϵ , Maximum iterations m_{max} , Initial guess u_0 .
- 2: **Output:** Approximate solution u .
- 3: Initialize $m = 0$.
- 4: Choose an initial guess u_0 .
- 5: **for** $m = 0, 1, 2, \dots, m_{max} - 1$ **do**
- 6: Evaluate the residual $R(u_m)$. #Often in variational form: $F(u_m)(v)$
- 7: Evaluate the linearized operator $R'(u_m)$. #Often in variational form: find $J_F\cdot$
- 8: Solve the linear equation $R'(u_m)[\delta u_m] = -R(u_m)$ for the correction δu_m . #Often in variational form: $J_F[\delta u_m](v) = -F(u_m)(v)$
- 9: Update the solution: $u_{m+1} = u_m + \delta u_m$.
- 10: Check for convergence.
- 11: **if** $\|R(u_{m+1})\| < \epsilon$ or $\|\delta u_m\| < \epsilon$ **then**
- 12: Set $u = u_{m+1}$.
- 13: **break**
- 14: **end if**
- 15: **end for**
- 16: **if** $m == m_{max} - 1$ and convergence not reached **then**
- 17: **Warning:** Maximum number of iterations reached without convergence.
- 18: **end if**
- 19: Set $u = u_{m+1}$.
- 20: **Return:** u .

13.2 Evaluation of the error of the solution

u_{ex} := exact solution
 u_R := high fidelity approximation
 u_{PIN} := PINN solution } Compute error with different norm, i.e.

- L2
1. Discretize the domain to take mesh
 2. Compute integral on each triangle via quadrature formula
- $$\int_{\Omega} (u_R - u_{PIN})^2 \approx \sum_T \int_T (u_R - u_{PIN})^2 \approx \sum_T \sum_{z_q} \frac{(u_R - u_{PIN})^2}{w_q} w_q$$
- ↑
quadrature points
3. $Sol.x :=$ derivatives of numeric solution in quadrature points
 $Sol.y :=$ derivatives wrt y of numeric solution in quadrature points

14 - NAVIER-STOKES EQUATIONS

the property of the fluid in a fixed point do not change in the time

These equations describes the motion of fluids. We work with **Steady NS equations** for incompressible flows.
The problem in **strong formulation** is:

(*) find $(u, p) \in \mathcal{U} \times Q$ s.t.

$$\begin{cases} -\mu \Delta u + (\nabla \cdot u) u + \nabla p = f & \text{in } \Omega \subset \mathbb{R}^d \\ \operatorname{div} u = 0 & \text{in } \Omega \subset \mathbb{R}^d \\ u = g & \text{in } \Gamma_0 \subset \partial \Omega \\ -p \nabla \cdot u + \mu (\nabla u) \cdot \nabla = h & \text{in } \Gamma_N \subset \partial \Omega \end{cases}$$

non linear term

the density of the fluid is constant
in every point at in every time

with d is the dimension of the problem

f is the forcing term $f \in [L^2(\Omega)]^d$

μ, ρ is the parameter

g, h have proper regularity \rightarrow for simplicity, $g=h=0$ we can study the homogeneous case

We can pass to **weak formulation**

(*) find $(u, p) \in \tilde{\mathcal{U}} \times \tilde{Q}$ s.t.

$$\begin{cases} u \cdot \nabla u + c(u, w, v) + b(u, p) = f(v) & \forall v \in \tilde{\mathcal{U}} \\ b(u, p) = 0 & \forall q \in \tilde{Q} \end{cases}$$

with $\tilde{\mathcal{U}} := [H_{\Gamma_0}^1(\Omega)]^d$

$\tilde{Q} := L^2(\Omega)$ if $\Gamma_N = \emptyset$, $\tilde{Q} := L^2_0(\Omega)$ to fix an added condition to have unique solution $\int_{\Omega} p dx = 0$

$a(u, v) := \int_{\Omega} \nabla u \cdot \nabla v dx$ \rightarrow Laplacian wrt x of u · Laplacian wrt y of v

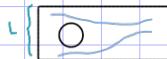
$$b(u, p) := - \int_{\Omega} p \operatorname{div}(v) dx$$

$$c(u, w, v) := \int_{\Omega} (u \cdot \nabla) w \cdot v dx$$

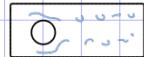
Scalar product between u, v the gradient operator that produces two equations

$$\text{Equation 1: } \int_{\Omega} \left(u_1 \frac{\partial w}{\partial x_1} \right) dx + \int_{\Omega} \left(u_2 \frac{\partial w}{\partial x_2} \right) dx$$

$$\text{Equation 2: } \int_{\Omega} \left(u_1 \frac{\partial v}{\partial x_2} \right) dx + \int_{\Omega} \left(u_2 \frac{\partial v}{\partial x_1} \right) dx$$



- Laminar: the flows goes around the obstacle and nothing happen
- Turbulent: the flows creates a lot of vortices after obstacle



It is determined by the **Reynolds Number** $Re := \frac{LU}{\mu}$

where μ := Kinematic Viscosity

L := characteristic length of Ω

U := characteristic velocity of Ω \rightarrow related to boundary condition

How can I solve this problem numerically?

14.2 General formulation of nonlinear problems

The problem becomes

$$\text{find } u \in \mathcal{U} \text{ s.t. } g(u; \mu) = 0 \quad \forall v \in \mathcal{U}$$

At the FOM level I introduce u_{ref} and the problem now is

$$\text{find } u \in \mathcal{U} \text{ s.t. } g(u_{\text{ref}}, v; \mu) = 0 \quad \forall v \in \mathcal{U}$$

The problem is well posed if the **Frechet derivative** of g at u_{ref} is $dg(u_{\text{ref}})$

applied to 2 variable

continuous
inf-sup stable

We solve the problem with Newton Method: given $u_s^0, \mu_s \in \mathcal{U}_s$ we iterate over K

$$dg[u_s^K(\mu_s)](s^K u_s, v_s; \mu_s) = -g(u_s^K(\mu_s), v_s; \mu_s) \quad \text{Newton step}$$

$$u_s^{K+1}(\mu_s) = u_s^K(\mu_s) + s^K u_s$$

Residual Vector

More precisely, for NSE:

find $(\delta u^k, \delta p^k) \in \tilde{U}_s \times \tilde{A}_s$ s.t.

$$\left\{ \begin{array}{l} (1) \text{ } a(\delta u_s^k, v_s; \mu) + c(u_s^k, \delta u_s^k, v_s) + c(\delta u_s^k, u_s^k, v_s) + b(v_s, \delta p_s^k) = f(v_s) - a(u_s^k, v_s) - c(u_s^k, u_s^k, v_s) - b(v_s, p_s^k) \\ (2) \text{ } b(\delta u_s^k, q_s) = -b(u_s^k, q_s) \end{array} \right. \quad \begin{array}{l} \text{Unknown} \\ \text{Unknown} \end{array}$$

with

$$(u_s^{kh}, p_s^{kh}) = (u_s^k, p_s^k) + (\delta u_s^k, \delta p_s^k) \quad \longrightarrow \quad X^{kh} = X^k + \delta^k x \quad x = (u_s, p_s)$$

Algebraically

$$J_s(X^k, \mu) \delta^k x = -G_s(X^k, \mu), \text{ where the jacobian matrix for NSE is } J_f = \begin{pmatrix} A + C_1(u_s^k) + C_2(u_s^k) & B^T \\ B & 0 \end{pmatrix}$$

where $A :=$ stiffness matrix
 $B_{ij} := - \int_{\Omega} \Psi_i \cdot \operatorname{div}(b_j) \, dx \longrightarrow \{\Psi_i\}_{i=1}^{N_s^q} \text{ span } \tilde{A} \quad u \quad \{\Psi_j\}_{j=1}^{N_s^u} \text{ span } \tilde{U}$
 different from Vionni's forces

$$C_1(u_s^k)_{ij} = c(u_s^k, \Psi_j, \Psi_i)$$

$$C_2(u_s^k)_{ij} = c(\Psi_j, u_s^k, \Psi_i)$$

What about the ROM level? We have to create the basis, i.e. with POD and compute the supermizer

To compute the basis we use POD and

1. We proceed standard for pressure to obtain $\Psi_p \in \mathbb{R}^{N_s^p, N}$

2. We enrich the velocity space with supermizer obtaining $\Psi_u \in \mathbb{R}^{N_s^u, 2N}$ $\Psi_u = [\Psi_{\text{standard}} | \Psi_s]$ a global stabilized basis $\Psi = \begin{bmatrix} \Psi_u & 0 \\ 0 & \Psi_p \end{bmatrix} \in \mathbb{R}^{N_s^u + N_s^p, 3N}$

We want to apply a Reduced Newton

Algo

For every k

$$J_N(X_N^k, \mu) \delta X_N = -G_N(X_N^k, \mu), \text{ where } J_N(X_N; \mu) = \begin{pmatrix} A_N + C_{1N}(u_N^k) + C_{2N}(u_N^k) & B_N^T \\ B_N & 0 \end{pmatrix}$$

$$A_N = \Psi_u^T A \Psi_u$$

$$B_N = \Psi_p^T B \Psi_u$$

$$C_{1N} = \Psi_u^T C_1(\Psi_u^T u_N^k, \mu) \Psi_u \quad \longrightarrow \text{depends on ROM level}$$

$$C_{2N} = \Psi_u^T C_2(\Psi_u^T u_N^k, \mu) \Psi_u$$

PROBLEM! No separation of the variable, at each step at the ROM level you have to compute ROM level

How to go from ROM to FOM during the iteration?

Considering k -th iteration

- I want to solve $J_N(X_N^k, \mu) \delta X_N = -G_N(X_N^k, \mu)$
- Namely at each iteration of the reduced solver
 - Compute the FOM representation of X_N^k .
 - Assemble $J_s(\Psi X_N^k, \mu)$ and $G_s(\Psi X_N^k, \mu)$.
 - Project jacobian and residual vector.
 - Solve reduced system obtaining X_N^{kh} .

$$\Psi^T J_s(\Psi X_N^k, \mu) \Psi \delta X_N^k = -\Psi^T G_s(\Psi X_N^k, \mu)$$

But, all these steps are so consuming! To avoid this you can apply some ML strategy that we already study

15. NOT AFFINE PROBLEM

Affinity property: consider $g: \bar{\Omega} \times P \rightarrow \mathbb{R}$
 \downarrow
 $(\boldsymbol{x}, \boldsymbol{\mu}) \mapsto g(\boldsymbol{x}, \boldsymbol{\mu})$
 be able to do offline vs online phase

that can be split in this formulation $g(\boldsymbol{x}, \boldsymbol{\mu}) = \sum_{q=1}^Q \tilde{g}_q(\boldsymbol{\mu}) g_q(\boldsymbol{x})$

BUT if we do not have this hypothesis?

1. NN

2. Interpolation with Empirical Interpolation Method

15.1 EIM: empirical interpolation method

We want to exploit a NOT-AFFINE approximation

$$g(\boldsymbol{x}, \boldsymbol{\mu}) = \sum_{q=1}^Q \tilde{g}_q(\boldsymbol{\mu}) g_q(\boldsymbol{x})$$

To do that, we have to define the interpolation operator $\mathcal{Y}_{\alpha}^{\boldsymbol{x}}: \mathcal{Y} \rightarrow X_{\alpha}^{\text{EIM}}$ where

$$\mathcal{Y} := \{g(\cdot, \boldsymbol{\mu}): \boldsymbol{\mu} \in P \subset C^0(\bar{\Omega})\}$$

X_{α}^{EIM} of dimension α generated by $\{h_1(\boldsymbol{x}), \dots, h_{\alpha}(\boldsymbol{x})\}$

such that, the operator can be written as

$$\mathcal{Y}_{\alpha}^{\boldsymbol{x}}: \mathcal{Y} \rightarrow X_{\alpha}^{\text{EIM}}$$

$$g(\boldsymbol{x}, \boldsymbol{\mu}) \mapsto (\mathcal{Y}_{\alpha}^{\boldsymbol{x}}(\boldsymbol{\mu}))_{q=1}^{\alpha} := \sum_{q=1}^{\alpha} \tilde{g}_q(\boldsymbol{\mu}) h_q(\boldsymbol{x}) \quad \forall \boldsymbol{\mu} \in P$$

IDEA: via interpolation, find an approximation of g that is affine so that we can do offline/online

$$(\tilde{g}_q(\boldsymbol{\mu}))_{q=1, \dots, \alpha} \text{ are computed } (\mathcal{Y}_{\alpha}^{\boldsymbol{x}}(\boldsymbol{\mu}))_{q=1}^{\alpha} = g(\boldsymbol{x}, \boldsymbol{\mu}) \quad \forall \boldsymbol{\mu} \in P, \quad q=1, \dots, \alpha$$

This interpolation points are called magic points := points in which the rule of interpolation holds

We do that with a greedy strategy

Initialization

\downarrow can be also the forcing term

$$1. \text{ Select some parameter } \boldsymbol{\mu}_{\text{EIM}}^1 := \underset{\boldsymbol{\mu} \in P}{\operatorname{argmax}} \|g(\cdot, \boldsymbol{\mu})\|_{L^{\infty}(\bar{\Omega})}$$

$$2. \text{ Create a set to collect the parameter } S_1 := \{\boldsymbol{\mu}_{\text{EIM}}^1\}$$

$$3. \text{ Define the function } \mathcal{G}_1(\boldsymbol{x}) := g(\boldsymbol{x}, \boldsymbol{\mu}_{\text{EIM}}^1)$$

$$4. \text{ Define the magic point } t^1 := \underset{\boldsymbol{x} \in \bar{\Omega}}{\operatorname{argmax}} |\mathcal{G}_1(\boldsymbol{x})|$$

$$5. \text{ Define a set to collect the magic point } T_1 = \{t^1\}$$

$$6. \text{ Define the basis function } h_1(\boldsymbol{x}) := \frac{\mathcal{G}_1(\boldsymbol{x})}{\mathcal{G}_1(t^1)} \Rightarrow X_1^{\text{EIM}} = \text{span}\{h_1\}$$

Now, this is an iterative method, hence we have to define the q -th step

Iterative orthogonalization steps

$$7. \text{ We have } S_q = \{\boldsymbol{\mu}_{\text{EIM}}^1, \dots, \boldsymbol{\mu}_{\text{EIM}}^q\}, \{\mathcal{G}_1(\boldsymbol{x}), \dots, \mathcal{G}_q(\boldsymbol{x})\}, T_q = \{t^1, \dots, t^q\}, X_q^{\text{EIM}} = \{h_1(\boldsymbol{x}), \dots, h_q(\boldsymbol{x})\}$$

$$8. \text{ Compute new parameter } \boldsymbol{\mu}_{\text{EIM}}^{q+1} := \underset{\boldsymbol{\mu} \in P}{\operatorname{argmax}} \|g(\cdot, \boldsymbol{\mu}) - (\mathcal{Y}_{\alpha}^{\boldsymbol{x}})(\cdot, \boldsymbol{\mu})\|_{L^{\infty}(\bar{\Omega})}$$

$$9. \text{ Enlarge } S \rightarrow \mathcal{G}_{q+1}(\boldsymbol{x}) := g(\boldsymbol{x}, \boldsymbol{\mu}_{\text{EIM}}^{q+1}), S_{q+1} := S_q \cup \{\boldsymbol{\mu}_{\text{EIM}}^{q+1}\}$$

$$10. \text{ Update the magic point } \rightarrow t^{q+1} := \underset{\boldsymbol{x} \in \bar{\Omega}}{\operatorname{argmax}} |\mathcal{G}_{q+1}(\boldsymbol{x}) - (\mathcal{Y}_{\alpha}^{\boldsymbol{x}})(\boldsymbol{x})| \quad \text{the parameter is fixed}$$

$$11. \text{ Update the basis } h_{q+1}(\boldsymbol{x}) := \frac{\mathcal{G}_{q+1}(\boldsymbol{x}) - (\mathcal{Y}_{\alpha}^{\boldsymbol{x}})(\boldsymbol{x})}{\mathcal{G}_{q+1}(t^{q+1}) - (\mathcal{Y}_{\alpha}^{\boldsymbol{x}})(t^{q+1})}$$

12. Then $X_{q+1}^{\text{EIM}} := X_q^{\text{EIM}} \cup \text{span}\{h_{q+1}\}$

13. Hence, the stopping criteria is $\max_{\mu \in P} \|g(x, \mu) - (\mathcal{Y}_q^\alpha g)(x, \mu)\|_{L^\infty(\Omega)} < \epsilon_{\text{tol}}$

How fast this approximation converges to the real function?

This property holds

$$\sup_{\mu \in P} \|g(\cdot, \mu) - (\mathcal{Y}_{q+1}^\alpha g)(\cdot, \mu)\|_{L^\infty(\Omega)} < c e^{-\alpha q}$$

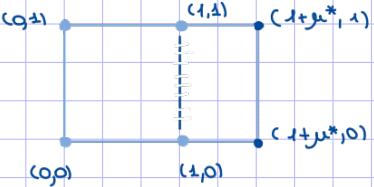
Is this the best method?

If the non-linearity is very complex, this method is NOT a good method via computation error of approximation.
Indeed every time we do

- $\underset{\mu \in P}{\text{argmax}}$ we have to discretize the space P with $P_{\text{train}}^{\text{EIM}} \subset P$
- $\underset{x \in \Omega}{\text{argmax}}$ with $x \in \bar{\Omega} \subset \Omega$

16 - GEOMETRIC PARAMETRIZATION

Let us assume $\mu \in \mathbb{P}$ one input parameter. The parameter can change the physics of the domain, but also the geometry. The μ -dependent domain is $\tilde{\Omega}(\mu)$ that is called **original domain**



To solve the problem we need to define a **reference domain** Ω that is μ -independent! The original domain Ω is obtained by means of a parametric map $\Phi: \Omega \times \mathbb{P} \rightarrow \tilde{\Omega}$ and namely $\tilde{\Omega}(\mu) = \Phi(\Omega; \mu) \quad \forall \mu \in \mathbb{P}$
 $\tilde{x} \mapsto \tilde{x}(\mu)$

Let us specify the nature of μ : $\mu = (\mu_1^{ph}, \dots, \mu_p^{ph}, \mu_g)$ where μ_g is the only parameter interacting with the domain.

Our PDE problem will be:

$$\text{given } \mu \in \mathbb{P}, \text{ find } \tilde{u}(\mu) \in \tilde{V}(\mu, g) : \tilde{a}(\tilde{u}(\mu), \tilde{v}; \mu) = \tilde{f}(v; \mu) \quad \forall v \in \tilde{V}(\mu, g)$$

Thanks to the map Φ , we can pull-back the problem onto the reference domain Ω : given $\mu \in \mathbb{P}$, find $u(\mu) \in V$: $a(u(\mu), v; \mu) = f(v; \mu) \quad \forall v \in V$

⚠ The two problems are different! There are some geometric factors induced by the map Φ

Let us recall that $x \in \Omega$ and $\tilde{x} \in \tilde{\Omega}(\mu)$. Moreover, both the subdomains are subsets of \mathbb{R}^d . We can define the **JACOBIAN MATRIX** of Φ , $J_\Phi(x, \mu) \in \mathbb{R}^{d,d}$ as

$$[J_\Phi(x, \mu)]_{ke} = \frac{\partial \Phi_k(x, \mu)}{\partial x_e}, \quad k, e = 1, \dots, d$$

If for any $\mu \in \mathbb{P}$ the determinant $|J_\Phi(x, \mu)| \neq 0$, then Φ is well-posed! For any $\tilde{\Psi}: \tilde{\Omega} \rightarrow \mathbb{R}$ it holds:

$$\int_{\tilde{\Omega}(\mu)} \tilde{\Psi}(\tilde{x}) d\tilde{x} = \int_{\Omega} \Psi(x) |J_\Phi(x, \mu)| dx, \text{ where } \Psi = \tilde{\Psi} \circ \Phi$$

16.1 Affine transformation

We will always assume Φ being an affine map of the form

$$\Phi(x, \mu) = A(\mu)x + c(\mu) \quad \text{with } A(\mu) \in \mathbb{R}^{d,d} \text{ and } c(\mu) \in \mathbb{R}^d$$

In this case $J_\Phi(x, \mu) = A(\mu)$ and the determinant depends only on the parameter μ . If Φ is well-posed, $\Phi^{-1}: \tilde{\Omega} \rightarrow \Omega$ is defined.

Recall: $\Omega = \Phi^{-1}(\tilde{\Omega}(\mu); \mu)$ and $J_{\Phi^{-1}}(\tilde{x}; \mu) = (J_\Phi(x; \mu))^{-1}$ for $\tilde{x} = \Phi(x; \mu)$

↳ by means of the inverse function theorem:

Φ continuous and non-singular in x than Φ^{-1} is invertible in a neighborhood of \tilde{x} and $J_{\Phi^{-1}} \circ \Phi = (J_\Phi)^{-1}$

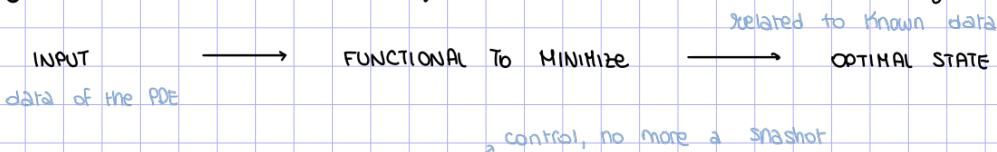
Namely, the determinant verifies: $|J_{\Phi^{-1}}(\tilde{x}; \mu)| = \frac{1}{|J_\Phi(x; \mu)|}$.

More details in **Parametric Geometry** with additional examples

17 - OPTIMAL CONTROL PROBLEMS

In this setting the PDE is a state system needs to be OPTIMIZED with some "external help", i.e. other variables called **CONTROLS**

The goal of OCP is minimize an objective functional to reach beneficial configuration



Then, our problem is: given a control $u \in U$ with U a Hilbert space, the state variable $y(u)$, with Y an Hilbert space, verifies the **strong form of the controlled PDE**

$$\textcircled{1} \quad G(y(u), u; \mu) = 0 \quad \text{in } Y' \quad \text{with } \mu = \text{physic or geometric parameter}$$

The main goal is to minimize over $Y \times U$

$J: Y \times U \times P \rightarrow \mathbb{R}$ defined as

$$J(y(u), u; \mu) = \frac{1}{2} \|y(u) - y_d\|_{Y_d}^2 + \frac{\alpha}{2} \|u\|_u^2 \quad \text{such that the strong form holds}$$

where

• $y_d \in Y_d \subset Y$ is the observation

• $\alpha > 0$ is the penalization parameter

$\alpha :=$ small
 $\alpha :=$ big
big control
small control

The function is quadratic to guarantee the uniqueness

In an abstract way, an OCP can be recast as

$$\text{find } u^* := \underset{u \in U}{\operatorname{argmin}} J(y(u), u; \mu)$$

$$\text{st. } G(y(u), u; \mu) = 0$$

$u^* :=$ the optimal control

$y^*(u^*) :=$ the optimal state

fixing a parameter, if you change μ you have to recompute everything

From now on, we work with linear PDEs:

$$G(y(u), u; \mu) = K(\mu)y - C(\mu)u - f(\mu) = 0$$

where

• $C(\mu): U \rightarrow Y'$, continuous

• $K(\mu): Y \rightarrow Y'$, continuous, positive definite coercive problem

• $f(\mu) \in Y'$

The related bilinear form are

- $a: Y \times Y \rightarrow \mathbb{R}$
- $c: U \times Y \rightarrow \mathbb{R}$
- $f: Y \rightarrow \mathbb{R}$

and the weak formulation of $\textcircled{1}$ is

$$\textcircled{2} \quad a(y, v; \mu) - c(u, v; \mu) - f(v; \mu) = 0 \quad \forall v \in V$$

But also the functional can be written in "weak form", rewriting the norm into integral

$$J(y, u; \mu) = \frac{1}{2} m(y - y_d, y - y_d; \mu) + \frac{\alpha}{2} n(u, u; \mu)$$

where

• $m: Y_d \times Y_d \rightarrow \mathbb{R}$ is Y_d -norm product

• $n: U \times U \rightarrow \mathbb{R}$ is U -inner product

Example. $Y_d = L^2(\Omega_d)$ $\Omega_d \subseteq \Omega$, $U = L^2(\Omega_u)$ $\Omega_u \subseteq \Omega$

$$\text{Then } J \text{ in strong form is } J(y, u; \mu) = \frac{1}{2} \|y - y_d\|_{L^2(\Omega_d)}^2 + \frac{\alpha}{2} \|u\|_{L^2(\Omega_u)}^2$$

$$\begin{aligned} &\downarrow \\ &(y - y_d, y - y_d)_{L^2(\Omega_d)} \\ &\int_{\Omega_d} (y - y_d)(y - y_d) \, dx \end{aligned}$$

What happen if I add another term? Consider $J(y, u; \mu) = \frac{1}{2} m(y - y_d, y - y_d; \mu) + \frac{\alpha}{2} n(u, u; \mu) - \frac{1}{2} m(y_d, y_d)$, but the solution remain the same

To solve I use The Lagrangian Approach

Let us define $p \in Y$ (same regularity of y) and define the lagrangian functional:

$$\mathcal{L}(y, u, p; \mu) : Y \times U \times P \rightarrow \mathbb{R}$$

$$(y, u, p; \mu) \mapsto J(y, u; \mu) + a(y, p) - c(u, p) - f(p)$$

resting the weak formulation on p
↓ in strong formulation
 $\langle K_{\mu} y, y - C_{\mu} u \rangle = F(\mu), p \rangle_Y$

where p is the adjoint variable and the role is the Lagrangian multiplier

We can derive the necessary conditions differentiating the Lagrangian and imposing the quantities to zero.

Explicitly,

$$\begin{aligned}\mathcal{L}(y, u, p; \mu) &= \frac{1}{2} m(y - y_d, y - y_d) + \alpha n(u, u) + a(y, p; \mu) - c(u, p; \mu) - f(p; \mu) - \frac{1}{2} m(y_d, y_d) \\ &= \frac{1}{2} m(y, y) - m(y_d, y) + \frac{1}{2} m(y_d, y_d) + \alpha n(u, u) + a(y, p; \mu) - c(u, p; \mu) - f(p; \mu) - \frac{1}{2} m(y_d, y_d) \\ &= \frac{1}{2} m(y, y) - m(y_d, y) + \frac{\alpha}{2} n(u, u) + a(y, p; \mu) - c(u, p; \mu) - f(p; \mu)\end{aligned}$$

The conditions are related to the derivatives:

$$\begin{aligned}1. \quad \text{Dy } \mathcal{L}(y, u, p; \mu) [q] &= \frac{1}{2} m(q, y) + \frac{1}{2} m(y, q) - m(y_d, q) + a(q, p; \mu) \quad \rightarrow \text{all the other terms does not depend on } y \\ &\quad \text{direction of differentiation} \\ &\quad \text{This is the way to variational differentiation } y \cdot y = 1 \cdot y + y \cdot 1 \xrightarrow{\text{Identity}} \text{but we compute on } q \text{ so that is} \\ &\quad \text{the reason why we put } q \text{ in } m(q, y) \text{ or } m(y, q) \\ &\quad = m(y, q) - m(y_d, q) + a(q, p; \mu) \rightarrow a(q, p; \mu) \text{ the test function is on the left, this is not good, we have to switch} \\ &\quad L^2 \text{ product between } y \otimes q\end{aligned}$$

If the operator is not symmetric, we have to built the adjoint operator of a to put q on the right
Hence

$$\text{Dy } \mathcal{L}(y, u, p; \mu) [q] = m(y - y_d, q) + a^*(p, q; \mu) \text{ with } a^* \text{ the adjoint operator, i.e. } a^*(p, q) = a(q, p)$$

Example. Linear operators $\langle Aq, p \rangle = \langle q, A^T p \rangle \longleftrightarrow a(q, p) = a^*(p, q)$

$$\text{i.e. } A := \text{grad.grad} \longrightarrow \int_{\Omega} \nabla q \nabla p \, dx = \int_{\Omega} \nabla p \nabla q \, dx \quad a(p, q) = a(q, p) \quad \text{symmetric}$$

$$A := \text{convection} \longrightarrow \int_{\Omega} \frac{\partial q}{\partial x} p \, dx = \int_{\Omega} \frac{\partial p}{\partial x} q \, dx = a^*(q, p) \neq a(p, q)$$

= 0 for boundary condition

$$2. \quad \text{Du } \mathcal{L}(y, u, p; \mu) [\nu] = \frac{\alpha}{2} n(\nu, u) + \frac{\alpha}{2} n(u, \nu) - c(\nu, p) = \alpha n(u, \nu) - c(\nu, p) = \alpha n(u, \nu) - c^*(p, \nu) = 0$$

$$3. \quad \text{Dp } \mathcal{L}(y, u, p; \mu) [z] = a(y, z; \mu) - c(u, z; \mu) - f(z) = 0$$

We obtain the constrain equal to zero, hence we obtain the weak formulation

The minimum has reached at the stationary point of the Lagrangian, imposing

$$\left\{ \begin{array}{ll} \text{Dy } \mathcal{L}(y, u, p; \mu) [q] = 0 & \forall q \in Y \\ \text{Du } \mathcal{L}(y, u, p; \mu) [\nu] = 0 & \forall \nu \in U \\ \text{Dp } \mathcal{L}(y, u, p; \mu) [z] = 0 & \forall z \in P \end{array} \right. \begin{array}{l} \text{adjoint equation} \\ \text{optimality equation} \\ \text{state equation} \end{array}$$

Namely

$$\left\{ \begin{array}{ll} m(y, q) + a^*(p, q; \mu) = m(y_d, q) & \forall q \in Y \\ \alpha n(u, \nu) - c^*(p, \nu; \mu) = 0 & \forall \nu \in U \\ a(y, z; \mu) - c(u, z; \mu) = f(z; \mu) & \forall z \in P \end{array} \right.$$

$$\text{Hence, matrix wise (FOM)} \quad \left\{ \begin{array}{l} M_{\mu}(y, y) + K_{\mu}^T(p) p = H_{\mu}(y_d, y_d) \\ a(N_{\mu}(u) u - C_{\mu}(p) p = 0 \\ K_{\mu}(y, z) - C_{\mu}(u, z) = f(z; \mu) \end{array} \right. \quad \text{finite element matrix}$$

Algebraically, in compact form

$$\begin{bmatrix} M_s(\mu) & 0 \\ 0 & \alpha N_s(\mu) \\ K_s(\mu) & -C_s(\mu) \end{bmatrix} \begin{bmatrix} y \\ u \\ p \end{bmatrix} = \begin{bmatrix} M_s(\mu) y_d \\ 0 \\ 0 \end{bmatrix}$$

observation vector

Saddle-point structure. Consider the problem that is in this structure, we can write

$$\begin{bmatrix} A(\mu) & B^T(\mu) \\ B(\mu) & 0 \end{bmatrix} \begin{bmatrix} z \\ p \end{bmatrix} = \begin{bmatrix} g^u \\ f^s \end{bmatrix} \quad \text{where } z = \begin{bmatrix} y \\ u \end{bmatrix}, \quad g^u = \begin{bmatrix} M_s(\mu) y_d \\ 0 \end{bmatrix}$$

- $A(\mu)$ is always well-posedness
- For $B(\mu)$, the well-posedness is related to the inf-sup stability and it is provable if $\mu \in \gamma$



ROM level. Collect three sets of snapshots and build three different basis matrices IB_y, IB_u, IB_p

- y_N spanned by IB_y
- u_N spanned by IB_u
- p_N spanned by $IB_p \neq B_p$ → How can I recover the inf-sup aggregated spaces?

We define $IB_{yp} = [IB_y \mid IB_p]$ the aggregated space and consider y_N, p_N spanned by IB_{yp}

In order to guarantee well-posedness, we need z_N basis wrt N basis without control!

Once we build the basis, assuming affinity, we can pre-compute

$$M_N = IB_{yp}^T M_s IB_{yp}$$

$$N_N = IB_u^T N_s IB_u$$

$$K_N = IB_{yp}^T K_s IB_{yp}$$

$$C_N = IB_{yp}^T C_s IB_u$$

and solve

$$\begin{pmatrix} M_N & 0 & K_N^T \\ 0 & \alpha N_N & C_N^T \\ K_N & C_N & 0 \end{pmatrix} \begin{pmatrix} y_N \\ u_N \\ p_N \end{pmatrix} = \begin{pmatrix} M_N y_d \\ 0 \\ f_N \end{pmatrix} \quad \text{where} \quad \begin{pmatrix} y_d \\ f_N \end{pmatrix} = IB_{yp}^T y_d$$

$$f_N = IB_{yp}^T f_s$$

1.1 A very simple practice example

Consider $\Omega = [0, 1] \times [0, 1]$, $y_d \in L^2(\Omega)$, $\Omega \subseteq \Omega$, $\mu \in P \subset \mathbb{R}$ s.t.

Given a $\mu \in P$, find $(y, u, p) \in H_0^1(\Omega) \times L^2(\Omega)$ s.t. they verify

$$\min_{(y,u)} \frac{1}{2} \|y - y_d\|_{L^2(\Omega)}^2 + \frac{\alpha}{2} \|u\|_{L^2(\Omega)}^2$$

$$\text{s.t. } \begin{cases} -\mu \Delta y = u + f & \text{in } \Omega \\ y = 0 & \text{in } \partial\Omega \end{cases}$$

Let us define $p \in H_0^1(\Omega)$ and the Lagrangian

$$\mathcal{L}(y, u, p; \mu) = \frac{1}{2} \int_{\Omega} (y - y_d)^2 dx + \frac{\alpha}{2} \int_{\Omega} u^2 dx + \mu \int_{\Omega} \nabla y \cdot \nabla p dx - \int_{\Omega} u p dx - \int_{\Omega} f p dx$$

$a(y, p)$ $-c(u, p)$

Now we have to differentiate

$$\frac{\partial \mathcal{L}(y, u, p; \mu)}{\partial y} [q] = \int_{\Omega} (y - y_d) q dx + \mu \int_{\Omega} \nabla q \cdot \nabla p dx$$

self adjoint

$$\frac{\partial \mathcal{L}(y, u, p; \mu)}{\partial u} [v] = \alpha \int_{\Omega} u v dx - \int_{\Omega} p v dx$$

self adjoint

$$\frac{\partial \mathcal{L}(y, u, p; \mu)}{\partial p} [z] = \mu \int_{\Omega} \nabla y \cdot \nabla z dx - \int_{\Omega} u z dx - \int_{\Omega} f z dx \rightarrow \text{This is the variational formulation}$$

⚠ If you have an advection term

$$a(y, z; \mu) = \int_{\Omega} \beta \nabla y \cdot z dx = \int_{\Omega} \operatorname{div}(\beta y) z dx - \int_{\Omega} \beta \nabla y \cdot z dx = \int_{\partial\Omega} \beta y z ds - \int_{\Omega} \beta \cdot \nabla y z dx$$

= 0 b.c. = $a^*(z, y; \mu)$