

A statistical view of physics-informed machine learning

Focus on physics-informed neural networks (PiNNs)

Setting:

$\Omega \subset \mathbb{R}^{d_{in}}$ ^{$d_{in} = d^*$} bounded Lipschitz domain ^{* for simplicity.}

ex: bounded convex domains

ex: manifolds with C^1 boundaries. ^{C^1}

$C^K(\Omega, \mathbb{R}^{d_{out}})$ space of functions from Ω to $\mathbb{R}^{d_{out}}$
K-time continuously differentiable.

$$C^\infty(\Omega, \mathbb{R}^{d_{out}}) = \bigcap_{K \geq 0} C^K(\Omega, \mathbb{R}^{d_{out}})$$

Hölder norm

$$\|f\|_{C^K(\Omega)} = \max_{|\alpha| \leq K} \|\partial^\alpha f\|_\infty, \Omega$$

hybrid modelling

- Goal: to estimate an unknown regression function f^*
such that $f^*: \Omega \rightarrow \mathbb{R}$ ^{$d_{out} = 1$} and

$$Y = f^*(x) + \epsilon$$

random noise

$$E[\epsilon|x] = 0$$

In traditional supervised learning, to do so we use

Training samples $(X_i, Y_i) \quad i = 1 \dots n$ independent copies of (X, Y) distributed as \mathbb{X}

- Additional information: physical modeling / prior knowledge

$$\mathcal{D}(f^*, \cdot) \approx 0 \quad \text{in } \Omega.$$

↑
known differential operator

- Additional information: Initial/Boundary conditions on $E \subseteq \partial\Omega$
↓
partial knowledge allowed.

ex: when $E = \partial\Omega$, Dirichlet boundary condition.

- To estimate f^* , we have 3 different data sets

① The "traditional" training samples (unknown distribution)

$$(X_1, Y_1) \dots (X_n, Y_n)$$

② Calibration points: a collection of iid points (uniformly distributed on Ω)

$$X_1^{(n)}, \dots, X_m^{(n)}$$

These points can be seen as a naive discretization of Ω .

③ Condition points: a collection of iid points $N \sim \mu_E$ on E

$$X_1^{(e)}, \dots, X_n^{(e)}$$

(Known distribution μ_E in E)

- Physics-informed empirical risk

can be an absolute value
for simplicity

$$R_{n, n_d, n_e}(f) = \underbrace{\frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|_2^2}_{\text{Data-fidelity term}} + \underbrace{\frac{\lambda_n}{n_d} \sum_{i=1}^n \|\mathcal{D}(f, X_i^{(n)})\|_2^2}_{\text{PDE term}} + \underbrace{\frac{\lambda_e}{n_e} \sum_{i=1}^{n_e} \|f(X_i^{(e)}) - h(X_i^{(e)})\|_2^2}_{\text{Initial/boundary conditions}}$$

I may drop the initial/boundary conditions for the sake of presentation.

- When this physics-informed risk is minimized over a class of neural networks, we obtain the so-called physics-informed neural networks (PINNs)

To do so, consider the class of fully-connected feed forward NN with H hidden layers of width $(L_1, L_2, \dots, L_H) = (D, D, \dots, D)$
 an input layer of width $L_0 = \text{dim } = d$ ($\Omega \subset \mathbb{R}^{\text{dim } = d}$)
 an output layer of width $L_{H+1} = \text{clout } = 1$ ($y \in \mathbb{R}^{\text{clout } = 1}$)
 activation function tanh

Overall,

$$f_\theta = f_{\text{out}} \circ (\tanh \circ \alpha_H) \circ \dots \circ (\tanh \circ \alpha_1)$$

↑
applied component-wise

with $\alpha_k : \mathbb{R}^{L_{k-1}} \rightarrow \mathbb{R}^{L_k}$

$$x \mapsto W_k x + b_k$$

weight matrix $W_k \in \mathbb{R}^{L_{k-1} \times L_k}$

bias

$$b_k \in \mathbb{R}^{L_k}$$

The parameters of the neural network - to be learned - are encoded

by

$$\theta = (W_1, b_1, \dots, W_{H+1}, b_{H+1}) \in \mathbb{R}^{H+1 \times H+D}$$

$$\mathbb{R}^{\sum_{l=1}^H (L_l+1) \times L_{l+1}}$$

Notation: $NN_H(D)$ = neural networks

of H hidden layers

H hidden layers

of width D .

of width D

Approximation results

Ideally, the NN set $NN_H(D)$ should be chosen large enough to approximate both the solution of the PDE and its derivatives.

Prop: [Density of NN in Hölder space, De Ryck et al (2021), Deumeche et al (2023)]

For $H \geq 2$, NN_H is dense in the space $(C^\infty(\bar{\Omega}, \mathbb{R}^{\text{data}}), \| \cdot \|_{C^K(\bar{\Omega})})$

i.e. for any function $u \in C^\infty(\bar{\Omega}, \mathbb{R}^{\text{data}})$, $\exists (u_p)_{p \in \mathbb{N}} \in \overline{NN_H}$

such that $\lim_{p \rightarrow \infty} \| u_p - u \|_{C^K(\bar{\Omega})} = 0$

for any K

Theorem [Gorovitsky, 2017]

For every $\epsilon > 0$, for every $f \in C^k([0,1]^d)$, there exists a tank neural network \hat{f} with $O(\epsilon^{-d/k})$ neurons such that $\|f - \hat{f}\|_\infty < \epsilon$.

NN structure allows for constructive proofs

Key properties

- Sum or composition of NN is a NN
- For smooth activation function σ with $\sigma'(0) \neq 0$, and $|x| \leq B$

$$\forall h > 0 \quad x = \frac{\sigma(xh) - \sigma(-xh)}{2h\sigma'(0)} + O(B^3 h^2)$$

Approximation of identity with 2 neurons.

Theorem [De Ryck, Lanthaler, Mishra, 2021]

For every $N \in \mathbb{N}$ and every $f \in C^s([0,1]^d)$, \exists a tank neural network \hat{f} with 2 hidden layers of width N^d such that for every $0 \leq k \leq s$,

$$\|f - \hat{f}\|_{C^k} \leq C \left(\ln(cN) \right)^k N^{-s+k}$$

$\uparrow \quad \uparrow$
 $\|f\|_N$

Remark 1

also holds for Sobolev norms -

Remark 2

How to translate this in terms of PDE residuals

Example of the heat equation

when $f^* \in C^1([0,T] \times [0,1])$

$$\partial(f, x) = \partial_t f - \partial_{xx}^2 f$$

$$① \|\partial_t \hat{f} - \partial_{xx}^2 \hat{f}\|_{L^2(\Omega)} \leq \sqrt{|\Omega|} \|\partial_t \hat{f} - \partial_{xx}^2 \hat{f}\|_\infty / c_0$$

$$② \|\partial_t \hat{f} - \partial_{xx}^2 \hat{f}\|_{C^0} = \|\partial_t \hat{f} - \underbrace{\partial_t f^* + \partial_{xx}^2 f^*}_{=0} - \partial_{xx}^2 \hat{f}\|_{C^0}$$

$$\leq \|\partial_t \hat{f} - \partial_t f^*\|_{C^0} + \|\partial_{xx}^2 \hat{f} - \partial_{xx}^2 f^*\|_{C^0}$$

$$\leq \|\hat{f} - f^*\|_{C^1} + \|\hat{f} - f^*\|_{C^2}$$

$$\leq 2 \|\hat{f} - f^*\|_{C^2}$$

Therefore, there is a NN with 2 hidden layers and N^2 neurons s.t.

$$\|\partial_t \hat{f} - \partial_{xx}^2 \hat{f}\|_{L^2} \leq 2\sqrt{|\Omega|} \|\hat{f} - f^*\|_{C^2} \leq C (\ln(cN))^2 N^{-8+2}.$$

Remark 3 : in practice, the PINNs are of size .

NN₄(50) Keerthapriyan et al. "Characterizing failure modes in PINNs"

NN₅(100) Xu et al. "How MN extrapolate : from FFNN to GNN".

NN₁₀(100) Argani et al. "Blood flow".

Back to learning

- In traditional supervised learning, we have the statistical model

$$Y = f^*(x) + \text{noise} \quad (\Delta)$$

such that $E[\text{noise} | X] = 0$.

target function

- Introduce the quadratic rule

$$\text{Risk}(f) = E[(Y - f(x))^2]$$

- This risk is minimal for $f : x \mapsto E[Y | X=x]$.

$$f = f^*. \quad (\text{Bayes predictor})$$

- To estimate f^* , we use n training samples $(X_i; Y_i)$ $i=1 \dots n$ iid as (X, Y) distributed as (Δ) , and construct an estimator \hat{f}_n .

- The first statistical requirement is to hope for **risk consistency**

$$\lim_{n \rightarrow \infty} \text{Risk}(\hat{f}_n) = \text{Risk}(f^*)$$

- How to adapt this notion in the case of PINNs?

Overfitting

- Study of limitations of PINNs
- To this end, we introduce a particular risk function.

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|_2^2 + \frac{\lambda_r}{|\Omega|} \int_{\Omega} \|\mathcal{D}(f, x)\|^2 dx + \lambda_e E \left[\|f(x^a) - h(x^a)\|^2 \right]$$

↑
expectation w.r.t. μ_E .

Regime where

$$\begin{cases} n \text{ is fixed} & (\text{costly physical measurements}) \\ n_e, n_r \rightarrow +\infty & (x_j^{(e)} x_j^{(r)} \text{ can be freely sampled} \\ & \text{up to computational resource}) \end{cases}$$

$$R_{n, n_e, n_r}(f) = \frac{1}{n} \sum_{i=1}^n \|f(x_i) - y_i\|_2^2$$

$$+ \frac{\lambda_r}{n_r} \sum_{i=1}^{n_r} \|\mathcal{D}(f, x_i^{(r)})\|^2$$

$$+ \frac{\lambda_e}{n_e} \sum_{i=1}^{n_e} \|f(x_i^{(e)}) - h(x_i^{(e)})\|_2^2$$

- Consider $\hat{\Theta}(q, n_e, n_r, D)$ a minimizing sequence of R_{n, n_e, n_r}

$$\lim_{q \rightarrow \infty} R_{n, n_e, n_r}(\hat{f}_{\hat{\Theta}(q, n_e, n_r, D)}) = \inf_{\Theta \in \mathbb{H}_{H, D}} R_{n, n_e, n_r}(f_\Theta)$$



- Focus on ERN-type estimators (no implicit bias here)

Definition: Physical risk consistency

$$\lim_{n_e, n_r \rightarrow \infty} \lim_{q \rightarrow \infty} R_n(\hat{f}_{\hat{\Theta}(q, n_e, n_r, D)}) = \inf_{f \in \text{NN}_H(D)} R_n(f)$$

We will show that PINNs can dramatically fail to be risk-consistent.

- Hybrid modelling with friction

$$\Omega =]0, T[$$

$$\forall f \in C^2(\bar{\Omega}, \mathbb{R})$$

$$D(u, x) = mu''(x) + \eta u'(x)$$

Dynamics of an object of mass m subjected to a fluid force of friction coefficient $\eta > 0$.

In this case, the empirical risk reads as

$$R_{n, n_r}(f_\theta) = \frac{1}{n} \sum_{i=1}^n |f_\theta(x_i) - y_i|^2 + \frac{n_r}{n_r} \sum_{l=1}^{n_r} \left(m f_\theta''(x_l^{(r)}) + \eta f_\theta'(x_l^{(r)}) \right)^2$$

Prop: provided that $\exists y_i \neq y_j$, whenever $D \geq n-1$,

for any n_r , for all $(x_l^{(r)})_{l=1..n_r}$, there exists

a minimizing sequence $(\hat{f}_\theta(p, n, n_r, D))$ such that

$$\begin{cases} \lim_{p \rightarrow \infty} R_{n, n_r}(\hat{f}_\theta(p, n, n_r, D)) = 0 & \text{(interpolation)} \\ \lim_{p \rightarrow \infty} R_n(\hat{f}_\theta(p, n, n_r, D)) = +\infty. & \text{(PDE term exploding)} \end{cases}$$

The PINN estimator is not (physical) risk consistent.

This phenomenon is explained by the existence of piecewise constant functions interpolating the observations x_1, \dots, x_n whose derivatives are

null at the points $X_1^{(r)}, \dots, X_{n_r}^{(r)}$ but diverge between these points.

Proof: Consider the following NN

$$\hat{f}_{\theta(p,n,n_r,D)}(x) = Y_{(1)} + \sum_{i=1}^n \frac{Y_{(i+1)} - Y_{(i)}}{2} \left[\tanh_p^{0+} \left(x - X_{(i)} - \frac{\delta(n,n_r)}{2} \right) + 1 \right]$$

where the observations $(X_1, Y_1), \dots, (X_n, Y_n)$ are reordered into $(X_{(1)}, Y_{(1)}), \dots, (X_{(n)}, Y_{(n)})$ such that $X_{(1)} \leq X_{(2)} \leq X_{(3)} \leq \dots \leq X_{(n)}$

$\delta(n,n_r)$ is the minimal distance between 2 distinct points in the input observations (X_i, Y_i) and the collocation points $(X_e^{(r)})_e$

$$\tanh_p(\cdot) = \tanh(p \cdot)$$

\tanh_p^{0+} composition of \tanh_p H times

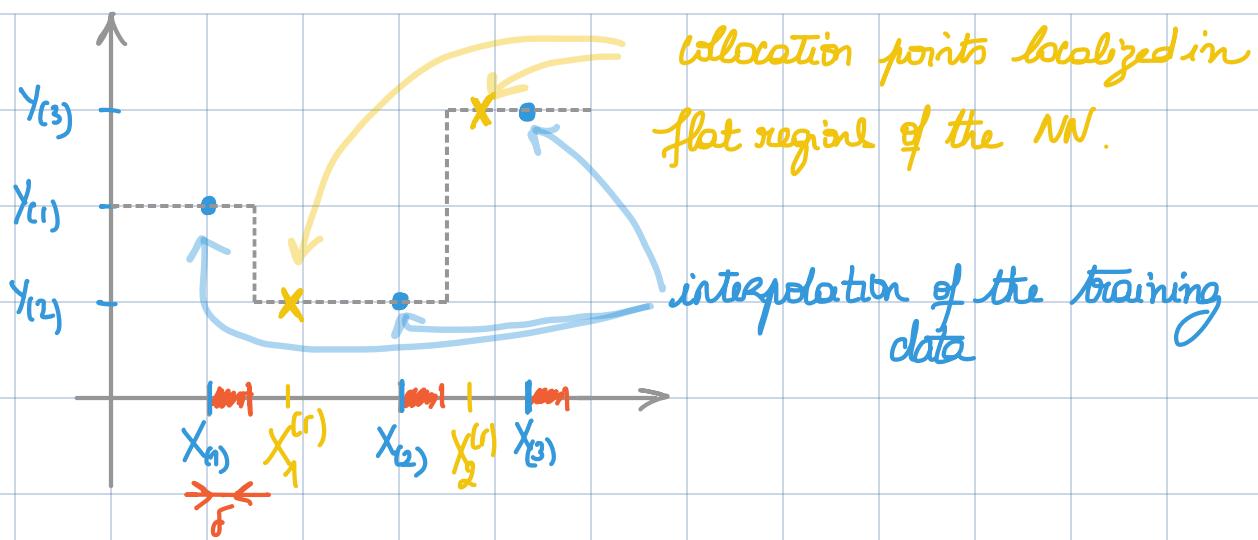
This NN has the following properties

(i) Interpolation of the training points

$$\lim_{p \rightarrow +\infty} \hat{f}_{\theta(p,n,n_r,D)}(X_i) = Y_i$$

$$\text{Indeed for } \delta > 0 \quad \lim_{p \rightarrow +\infty} \|\tanh_p^{0+} - \text{sgn}\|_{C^k(R \setminus [-\delta, \delta])} = 0$$

Therefore $\frac{1}{2} \left[\tanh_p^{0+} \left(x - X_{(i)} - \frac{\delta(n,n_r)}{2} \right) + 1 \right] \xrightarrow[p \rightarrow +\infty]{} \begin{array}{l} \text{heaviside of } x \\ \text{centred at } \left(X_{(i)} + \frac{\delta}{2} \right) \end{array}$



② All derivatives higher than order 1 vanish almost everywhere.

$$\mathcal{D}(\hat{f}_p, x_e^{(r)}) = \left(m \hat{f}_p''(x_e^{(r)}) + m \hat{f}_p'(x_e^{(r)}) \right)^2 = 0 \quad \forall x_e^{(r)}.$$

\parallel
 \parallel

③ But explode locally.

For any $0 < \varepsilon < \delta(n, r)/4$,

$$\begin{aligned}
 R_n(\hat{f}_p(p_{n,r}, \cdot)) &\geq \frac{1}{T} \int_{[0,T]} \mathcal{D}(\hat{f}_p(p_{n,r}, \cdot), x)^2 dx \\
 &\geq \frac{1}{T} \sum_{i=1}^n \int_{x_{ci}+\delta/2-\varepsilon}^{x_{ci}+\delta/2+\varepsilon} \mathcal{D}(\hat{f}_p(p_{n,r}, \cdot), x)^2 dx.
 \end{aligned}$$

Put the square outside the integral
 * fundamental theorem of integration.

$$\begin{aligned}
 &\geq \frac{1}{T} \sum_{i=1}^n \frac{1}{2\varepsilon} \left[m \hat{f}_p''(x_{ci} + \delta/2 + \varepsilon) - \hat{f}_p''(x_{ci} + \delta/2 - \varepsilon) + m \left(\hat{f}_p'(x_{ci} + \delta/2 + \varepsilon) - \hat{f}_p'(x_{ci}) \right) \right]^2 \\
 &\geq m^2 \cdot \frac{1}{2\varepsilon T} \sum_{i=1}^n (y_{ci+1} - y_{ci})^2.
 \end{aligned}$$

Finally, $\lim_{p \rightarrow \infty} R_n(\hat{f}_p(p_{n,r}, \cdot)) = +\infty$

Remarks:

- ① This pathological sequence of Θ is such that $\|\Theta\|_2 \rightarrow \infty$.
- ② The proof does not depend on the geometry of the collocation points $(X_e^{(r)})_e$ or the condition points $(X_e^{(c)})_e$ (holds with grids or quasi-Monte Carlo methods)
- ③ Holds for any PDE with derivatives ≥ 1 involved.

- Similar behaviour in the case of PDE solve

Consider

$$\Omega = [-1; 1] \times [0, \pi]$$

$$\mathcal{D}(f, x) = \partial_t f - \partial_{xx}^2 f \quad \text{that equation}$$

$$\text{Boundary conditions } f(-1, t) = f(1, t) = 0$$

$$\text{Initial condition } f(x, 0) = \tanh^{0.5}(x + 0.5)$$

$$\begin{aligned} & -\tanh^{0.5}(x - 0.5) \\ & + \tanh^{0.5}(0.5) \\ & - \tanh^{0.5}(1.5) \end{aligned}$$

Bell function

Empirical risk function

$$R_{n_r, n_c}(f) = \frac{1}{n_c} \sum_{j=1}^{n_c} |f(x_j^{(c)}) - h(x_j^{(c)})|^2 + \frac{1}{n_r} \sum_{l=1}^{n_r} \mathcal{D}(f, X_e^{(r)})^2$$

Physics theoretical risk function

$$R(f) = \lambda \mathbb{E} \left[|f(x^{(c)}) - h(x^{(c)})|^2 \right] + \frac{1}{|\Omega|} \int_{\Omega} \mathcal{D}(f, x)^2 dx$$

Prop: When $D \geq G$, $\forall (n_e, n_r) \in \mathbb{N}^2$ $\forall (X_1^{(r)}, \dots, X_n^{(r)}) \in \mathbb{R}^{n_r}$ $\forall (X_1^{(e)}, \dots, X_{n_e}^{(e)}) \in \mathbb{R}^{n_e}$

\exists a minimizing sequence $(\hat{f}_{\theta(p_{n_e, n_r}, D)})_p$ s.t

$$\lim_{p \rightarrow \infty} R_{n_e, n_r}(\hat{f}_{\theta(p_{n_e, n_r}, D)}) = 0$$

$$\lim_{p \rightarrow \infty} R(\hat{f}_{\theta(p_{n_e, n_r}, D)}) = +\infty$$

This PINN estimator is not physical risk consistent.

G Similarly, this estimator corresponds to a function that equals zero on Ω (and thus satisfies the linear PDE) while satisfying the initial condition on $\partial\Omega$.

G Again this corresponds to a limit of neural networks f_θ such that $\|\theta\|_2 \rightarrow +\infty$.

Fighting overfitting

- Connection between L^2 -norm of the weights and the NN regularity.

Prop: The following holds: $\forall \theta \in \mathbb{H}_{H,D}$

$$\|f_\theta\|_{C^K(\mathbb{R}^{d_H})} \leq C_{K,H} (D+1)^{HK+1} (1 + \|\theta\|_2)^{HK} \|\theta\|_2.$$

Remark: this bound is tight in the sense that there exists a sequence $(\theta_p) \in \mathbb{H}_{H,D}$ such that

$$(i) \lim_{p \rightarrow \infty} \|\theta_p\|_2 = +\infty$$

$$(ii) \|f_{\theta_p}\|_{C^K(\mathbb{R}^{d_H})} \geq \bar{C}_{K,H} \|\theta_p\|_2^{HK+1}.$$

- Idea: use a ridge-type / Tikhonov-type regularization to avoid pathological minimizing sequences

Def: Ridge risk function

$$R_{n,n_r,n_c}^{(\text{ridge})}(f_\theta) = R_{n,n_r,n_c}(f_\theta) + \lambda_{(\text{ridge})} \|\theta\|_2^2$$

hyperparameter

Denote $(\hat{\theta}_{(n,n_r,n_c,D)}^{(\text{ridge})})_n$ a minimizing sequence of the ridge risk function

$$\lim_{n \rightarrow \infty} R_{n,n_r,n_c}^{(\text{ridge})} \left(f_{\hat{\theta}_{(n,n_r,n_c,D)}^{(\text{ridge})}} \right) = \inf_{\theta \in \mathbb{H}_{H,D}} R_{n,n_r,n_c}^{(\text{ridge})}.$$

Note that ridge regularization is implementable in most deep learning libraries (e.g. keras, pytorch, ...), implemented via "weight decay".

Class of polynomial operators

ex: Navier-Stokes equation on $\Omega = \Omega_1 \times [0, \pi]$, $f = (f_x, f_y, f_z, p)$
 $(3^{\text{rd}} \text{ eqt}) \quad \mathcal{D}(f, z) = \partial_t f_z - f_3 \partial_z f_z - \eta \partial_z^2 f_z + p^{-1} \partial_z p + g(z).$

It can be seen as a polynomial P of $f_z, \partial_z f_z, \partial_z^2 f_z, \partial_z p$, where P is given by $P(z_1, z_2, z_3, z_4, z_5) = z_3 - z_1 z_2 - \eta z_4 + p^{-1} z_5 + g -$

$\deg P = 2$ but $\deg \mathcal{D} = 3$ (counting the degree of the polynomial combined with the order of differentiation).

Encompass linear + nonlinear PDEs.

Theorem: ($H \geq 2$) Assume the condition $f \in H$ is Lipschitz and that

\mathcal{D} is a polynomial operator.

By choosing $\lambda_{\text{ridge}} = \frac{1}{(\min(n_r, n_r))^K}$ with $K = \frac{1}{12+4H(1+2H \deg \mathcal{D})}$

$$\lim_{n_r \rightarrow \infty} \lim_{p \rightarrow \infty} \mathcal{P}_n \left(\frac{\lambda_{\text{ridge}}}{\mathcal{D}^p} \right) = \inf_{\mathcal{N}H(D)} \mathcal{P}_n$$

- The choice of the hyperparameter λ_{ridge} which vanishes when $n_e, n_r \rightarrow \infty$ makes the ridge bias vanish.
- This choice depends on known quantities (NN depth, degree of the differential operator, nb of discretization points).
- Extension of this result to remove any approximation bias introduced by the class $NN_H(D)$:

$$\lim_{D \rightarrow \infty} \lim_{n_e, n_r \rightarrow \infty} \lim_{p \rightarrow \infty} R_n(f_{\theta}^{\text{ridge}}) = \inf_{G^{\infty}(\bar{\Omega}, R^{\text{out}})} R_n$$

One can make D depend on n_e, n_r as well

$$(D = \min(n_e, n_r)^{\xi} \text{ with } \xi \text{ function of } H \text{ and } \deg D)$$

In such a case

$$\lim_{n_e, n_r \rightarrow \infty} \lim_{p \rightarrow \infty} R_n(f_{\theta}^{\text{ridge}}) = \inf_{G^{\infty}(\bar{\Omega}, R^{\text{out}})} R_n$$

Proof: In the case of the 1D heat equation only

$$\partial_t f - \partial_{xx}^2 f = 0 \text{ meaning that } \mathcal{D}(f, \cdot) = \partial_t f - \partial_{xx}^2 f$$

① $\det f_0 = 0 \in NN_H(D)$ the NN with parameters all 0.

Initially $R_{n_e, n_r}^{(\text{ridge})}(f_0) = R_{n_e, n_r}(f_0)$

Therefore,

$$R_{n,n_r,n_r}(f_0) = \frac{1}{n} \sum_{i=1}^n (f_0(x_i) - y_i)^2 + \lambda_r \sum_{i=1}^n \left[\mathcal{D}(f_0, X_i^{(n)}) \right]^2$$

$$\leq \frac{1}{n} \sum_{i=1}^n |y_i|^2 + \lambda_r L =: LB$$

L "coefficients" in the

does not depend on
 n_r (n_r if any) and λ_{ridge}

polynomial operator of degree 1.

$$\mathcal{D}(f, \cdot) = \phi_1 \partial_x f - \phi_2 \partial_{xx} f$$

$$\|\phi_1\|_{\infty, \bar{\Sigma}} = \|\phi_2\|_{\infty, \bar{\Sigma}} = 1.$$

② Let $\hat{\theta}_{p,n_r,n_r,D}^{(\text{ridge})}$ be any minimizing

sequence of the empirical ridge risk : $\lim_{p \rightarrow \infty} R_{n,n_r,n_r}^{(\text{ridge})}(\hat{f}_0^{(\text{ridge})}) = \inf_{\theta \in \mathbb{H}_{HD}} R_{n,n_r,n_r}(\theta)$

call $n_{r,e} = \min(n_r, n_r)$

Define $E_1^{(\text{large})}(n_{r,e}) = \{ \theta \in \mathbb{H}_{HD}, \|\theta\|_2 \geq n_{r,e}^K \}$

$E_2^{(\text{inter})}(n_{r,e}) = \{ \theta \in \mathbb{H}_{HD}, n_{r,e}^{K/4} \leq \|\theta\|_2 \leq n_{r,e}^K \}$

$E_3^{(\text{small})}(n_{r,e}) = \{ \theta \in \mathbb{H}_{HD}, \|\theta\|_2 \leq n_{r,e}^{K/4} \}$

Remark that $\mathbb{H}_{HD} = E_1^{(\text{large})} \cup E_2^{(\text{inter})} \cup E_3^{(\text{small})}$.

Summary of what's next :

θ small = no problem.

in-between = be useful.

θ large = does not happen

Idea of the proof: show that almost surely, given any n_r and

for p large enough, $\hat{\theta}_{(p,n_r,D)}^{(\text{ridge})} \in E_2^{\text{inter}} \cup E_3^{\text{small}}$.

Moreover on $E_2^{\text{inter}} \cup E_3^{\text{small}}$, the empirical risk function $R_{n,n_r}^{(\text{ridge})}$ is close to the theoretical risk R_n .

③ @ $\forall \theta \in E_1^{\text{large}}$

$$R_{n,n_r}^{(\text{ridge})}(\theta) \geq \lambda_{\text{ridge}} \|\theta\|_2^2 \geq n_r^K.$$

$$\text{Once } n_r \geq (UB+1)^{\frac{1}{K}},$$

$$\begin{aligned} \inf_{\theta \in E_3^{\text{small}}} R_{n,n_r}^{(\text{ridge})}(f_\theta) + 1 &\leq R_{n,n_r}^{(\text{ridge})}(f_\theta) + 1 \\ &\leq UB + 1 \\ &\leq n_r^K \\ &\leq \inf_{\theta \in E_1^{\text{large}}} R_{n,n_r}^{(\text{ridge})}(f_\theta) \end{aligned}$$

This shows that for n_r large enough,

$$\hat{\theta}_{(p,n_r,D)}^{(\text{ridge})} \notin E_1^{\text{large}}$$

④ We can show that

$$\sup_{\theta \in E_2^{\text{inter}} \cup E_3^{\text{small}}} \left| \frac{1}{n_r} \sum_{l=1}^{n_r} \mathcal{D}(f_\theta, X_l^{(r)})^2 - \frac{1}{|S|} \int_S \mathcal{D}(f_\theta, z)^2 dz \right| \leq \log(n_r) n_r^{\beta - \frac{1}{2}}$$

$$\text{with } \beta = (\alpha + H(1 + (2 + H) \deg \mathcal{D})) .$$

(Admitted).

Uniform control of stochastic processes (Hoeffding / Dudley / McDiarmid)

Thus almost surely for all n_r large enough and $\theta \in E_2(n_{r,c})$

$$R_{n,n_r}^{(\text{ridge})}(f_\theta) \geq R_n(f_\theta) + \lambda_{(\text{ridge})} \|\theta\|_2^2 - \cancel{\frac{2}{n} \log^2(n_{r,c}) n_r^{-\beta/2}}$$

~~F~~

coming from the absolute value.

≥ 0 ?

Yes! Because

$$\text{for all } \theta \in E_2^{\text{inter}}, \quad \lambda_{(\text{ridge})} \|\theta\|_2^2 \geq n_{r,c}^{-K/2}$$

$n_{r,c}^{-K} \parallel \lambda_{(\text{ridge})} \parallel n^{K/4}$

$$\text{and } -K/2 > \beta - 1/2$$

Then almost surely for $n_{r,c}$ large enough and for all $\theta \in E_2(n_{r,c})$

$$R_{n,n_r}^{(\text{ridge})}(f_\theta) \geq R_n(f_\theta)$$

There, on E_2 , the empirical risk controls the theoretical one

c) $\forall \theta \in E_3^{\text{small}} \quad \lambda_{(\text{ridge})} \|\theta\|_2^2 < n_{r,c}^{-K/2}$

Using \otimes , we deduce that for all n_r large enough

for all $\theta \in E_3^{\text{small}}$

$$|R_{n,n_r}^{(\text{ridge})}(f_\theta) - R_n(f_\theta)| \leq \cancel{\frac{2}{n} \log^2(n_{r,c}) n_r^{-K/2}}$$

one coming from integral approx.
one coming from the ridge reg.

②

Fix $\varepsilon > 0$

Call $(\theta_p)_p$ any minimizing sequence of the theoretical risk R_n

$$\lim_{p \rightarrow \infty} R_n(f_{\theta_p}) = \inf_{\theta \in \mathbb{H}_{n,D}} R_n(f_\theta)$$

By def² $\exists p_\varepsilon$ such that $|R_n(f_{\theta_{p_\varepsilon}}) - \inf_{\theta \in \mathbb{H}_{n,D}} R_n(f_\theta)| \leq \varepsilon$.

For fixed n_r , according to ①, for p large enough

$$\hat{\theta}_{p,n_r,D}^{(\text{ridge})} \in E_2(n_r) \cup E_3(n_r).$$

According to ④ & ⑤

$$R_n(f_{\hat{\theta}_{p,n_r,D}^{(\text{ridge})}}) \leq R_{n,n_r}^{(\text{ridge})}(f_{\hat{\theta}_{p,n_r,D}^{(\text{ridge})}}) + 2 \log^2(n_r) n_r^{-K/2}$$

By def² of the minimizing sequence $\hat{\theta}_{p,n_r,D}^{(\text{ridge})}$, for p large enough,

$$R_{n,n_r}^{(\text{ridge})}(f_{\hat{\theta}_{p,n_r,D}^{(\text{ridge})}}) \leq \inf_{\mathbb{H}_{n,D}} R_{n,n_r}^{(\text{ridge})} + \varepsilon$$

According to ③

$$\begin{aligned} \inf_{\theta \in E_2(n_r) \cup E_3(n_r)} R_{n,n_r}^{(\text{ridge})}(f_\theta) &\leq \inf_{\theta \in E_3(n_r)} R_{n,n_r}^{(\text{ridge})}(f_\theta) \\ &\leq \inf_{\theta \in E_3(n_r)} R_n(f_\theta) + 2 \log^2(n_r) n_r^{-K/2} \end{aligned}$$

Given for all n_r large enough $\theta_p \in E_\delta^{small}(n_r)$

Therefore $\inf_{\theta \in E_\delta^{small}(n_r)} R_n(f_\theta) \leq R_n(f_{\theta_p})$

Combining all, almost surely for n_r and p large enough

$$R_n(f_{\theta^{(ridge)}_{p,n_r,D}}) \leq \inf_{\theta \in H_{H,D}} R_n(f_\theta) + 3\epsilon$$

As ϵ is arbitrary, we get $\lim_{n_r \rightarrow \infty} \lim_{p \rightarrow \infty} R_n(f_{\theta^{(ridge)}_{p,n_r,D}}) = \inf_{\theta \in H_{H,D}} R_n(f_\theta)$.



Open question :

- Implicit regularization in PINN training -
- Spectral bias -

Strong convergence of PINNs for linear PDE

- Good CV properties in terms of $\|u_h\|_{H^k}$ do not imply good CV properties in terms of function, for instance

$$\hat{f}_{\text{PINN}} \xrightarrow{?} f^{\infty} \text{ in } L^2(\text{some}?)$$

- A first example in hybrid modeling setting

$$d_{\text{in}} = 2 \quad d_{\text{out}} = 1$$

$$\Omega = [0,1] \times [0,T]$$

$$\partial_t f + \partial_x f = 0$$

(advection equation)

$$h(x,0) = 1 + x$$

(initial condition)

$$h(0,t) = 1 + t$$

(boundary condition)

Goal: the regression function f^{∞} should be close to 1.

Consider the sequence of NN:

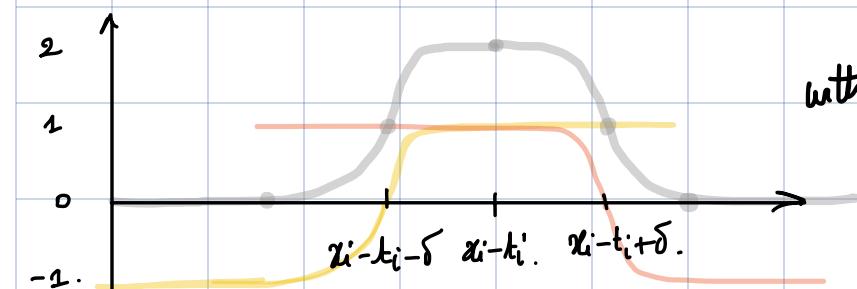
$$\text{sgn}(x-t-x_i+t_i+\delta)$$

$$f_{\text{NN}}(x,t) = 1 + \sum_{i=1}^n \frac{(y_i - 1)}{2} \left[\tanh_p^{\text{OH}}(x-t-x_i+t_i+\delta) - \tanh_p^{\text{OH}}(x-t-x_i+t_i-\delta) \right]$$

$$x_i = (x_i, t_i)$$

$$-\tanh_p^{\text{OH}}(x-t-x_i+t_i-\delta)$$

$$\text{sgn}(x-t-x_i+t_i-\delta)$$



with

$$\delta \leq \frac{1}{2} \min_{i \neq j} |x_i - x_j + t_i - t_j|$$

Indeed, note that by the step shape, we get

$$|f(x+\delta) - f(x)| = 1 = \left| \int_x^{x+\delta} f'(u) du \right| = \left| \langle f', \mathbb{1}_{[x, x+\delta]} \rangle_{L^2} \right|$$

$$\leq \sqrt{\int_x^{x+\delta} (f')^2(u) du} \sqrt{\delta}.$$

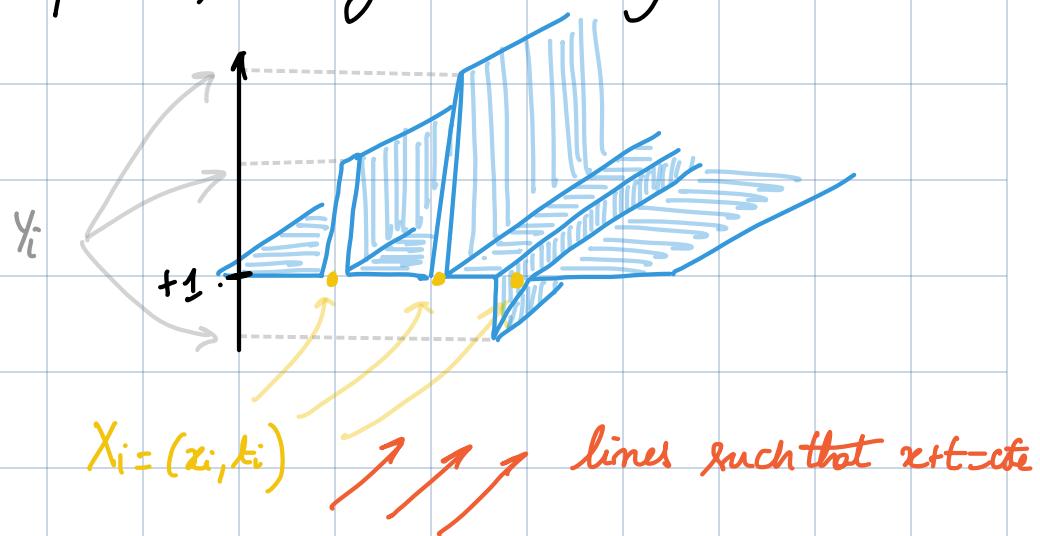
by Cauchy-Schwarz

therefore $\|f'\|_{L^2(S)} \geq 1/\delta$.

□

④ $f_{\delta,p}(x_i, t_i) \xrightarrow[p \rightarrow +\infty]{} y_i$

- ④ $f_{\delta,p}$ satisfies the equation, as being constant along the lines $x+t=c_0t$.



Therefore $\lim_{p \rightarrow +\infty} R(f_{\delta,p}) = 0$ by ④ and ③.

If $D \geq 2n$, $\inf_{f \in \mathcal{H}_n(D)} R_n(f) = 0 \Rightarrow \hat{f}_{\theta} \xrightarrow[\text{per. nrr.}]{L^2(\Omega)} \mathbb{1}_{\Omega}$

- ✓ Fine if the model is accurate independently of n and f*!!!
- ✗ Problematic if the model is not exact.

• A second example in PDE theory

$$\begin{cases} \Omega =]-1; 1[\\ h(1) = 1 \\ \mathcal{D}(f, x) = xf'(x) \end{cases}$$

Clearly $f^*(x) = 1$ is the only strong solⁿ of the PDE

Consider the NN sequence to be $\hat{f}_p = \tanh_p \circ \tanh^{(H-1)}$

$\lim_{p \rightarrow \infty} R(\hat{f}_p) = R(f^*) = 0$ (as \hat{f}_p satisfies the ODE and $\hat{f}_p(1) = 1$)

(\hat{f}_p) is a minimizing sequence of R .

But $\hat{f}_0(x) = \text{sgn}(x) \neq f^*$.

• To retrieve strong convergence properties, one can use Sobolev

regularization

for the sake of
regularization

• Assumption [linear PDEs]

Assume that \mathcal{D} is an affine operator (of order K)
of f and its derivatives

$$\mathcal{D}(f, x) = \sum_{|\alpha| \leq K} \langle A_\alpha(x), \partial^\alpha f(x) \rangle + B(x)$$

of order K .

$$\sum_{|\alpha| \leq K} \langle A_\alpha(x), \partial^\alpha f(x) \rangle$$

$\in \mathbb{B}^\infty(\bar{\Omega}, \mathbb{R}^{\text{out}})$

in the paper, we can handle a 2nd member.

ex: advection, heat, wave, Maxwell equations.

- Regularized theoretical risk

$$R_n^{(\text{reg})}(f) = R_n(f) + \lambda_{(\text{reg})} \|f\|_{H^{m+1}(\Omega)}^2$$

with $\|f\|_{H^{m+1}(\Omega)}^2 = \sum_{|\alpha| \leq m+1} \|\partial_x^\alpha f\|_{L^2(\Omega)}^2$.

- Regularized empirical risk

$$R_{n,n_e,n_r}^{(\text{reg})}(f) = R_{n,n_e,n_r}(f) + \lambda_{(\text{ridge})} \|\partial\|_2^2 + \frac{\lambda_{(\text{reg})}}{m_e} \sum_{l=1}^{m_e} \sum_{1 \leq k \leq m+1} \|\partial_x^\alpha f(x_e^{(k)})\|^2$$

✓ Can be straightforwardly implemented in the usual PINN framework by considering the Sobolev penalty as additional PDE.

In what follows, for the sake of presentation, we discard the boundary condition (cf. gray parts).

Proposition [Characterization of the unique minimizer of $\mathcal{R}_n^{(\text{reg})}$]

Assume \mathfrak{D} to be of order K and to encode a linear PDE.

$$m \geq \max \left(\lfloor \frac{\dim \mathfrak{D}}{2} \rfloor, K \right)$$

then, $\mathcal{R}_n^{(\text{reg})}$ has a unique minimizer \hat{f}_n over $H^{m+1}(\Omega)$, characterized to be the unique element of $H^{m+1}(\Omega)$ such that

$$\forall v \in H^{m+1}(\Omega) \quad \boxed{c_{\hat{f}_n}(\hat{f}_n, v) = \mathcal{B}_n(v)}$$

Weak formulation of PDE on $H^m(\Omega, \mathbb{R}^{\text{dout}})$

with

$$\begin{aligned} c_{\hat{f}_n}(\hat{f}_n, v) &= \frac{1}{n} \sum_{i=1}^n \langle \Pi(\hat{f}_n)(x_i), \Pi(v)(x_i) \rangle \\ &\quad + \lambda \mathbb{E} \langle \Pi(\hat{f}_n)(x^e), \Pi(v)(x^e) \rangle \\ &\quad + \frac{\gamma}{|\Omega|} \int_{\Omega} \mathfrak{D}^{(\text{lin})}(\hat{f}_n, x) \mathfrak{D}^{(\text{lin})}(v, x) dx \\ &\quad + \frac{\lambda \kappa b}{|\Omega|} \sum_{k \leq m+1} \int_{\Omega} \langle \partial^\alpha \hat{f}_n(x), \partial^\alpha v(x) \rangle dx \end{aligned}$$

$$\begin{aligned} \mathcal{B}_n(v) &= \frac{1}{n} \sum_{i=1}^n \langle y_i, \Pi(v)(x_i) \rangle + \\ &\quad + \lambda \mathbb{E} [\langle \Pi(v)(x^e), h(x^e) \rangle] \\ &\quad - \frac{\gamma}{|\Omega|} \int_{\Omega} B(x) \mathfrak{D}^{(\text{lin})}(v, x) dx. \end{aligned}$$

$$\Pi: H^{m+1}(\Omega, \mathbb{R}^{\text{dout}}) \rightarrow C^0(\Omega, \mathbb{R}^{\text{dout}})$$

B-called Sobolev embedding

such that $\Pi(f)$ is the unique continuous function that coincides with f almost everywhere.

↳ Lax-Milgram based result.

↳ Remark that

$$R_n^{(cg)}(f) = c_n(f, f) - 2B_n(f) + \frac{1}{n} \sum_{i=1}^n \|y_i\|^2 + \mu \int_{\Omega} h(x^{(i)}) \|u\|_2^2 + \frac{1}{|\Omega|} \int_{\Omega} B^2(u) dx.$$

Proof: For the sake of simplicity, in the proof, we discard the boundary conditions & 2nd member in the PDE

Remark that

$$\textcircled{1} \quad c_n(f, f) - 2B_n(f) = R_n^{(cg)}(f) - \frac{1}{n} \sum_{i=1}^n \|y_i\|^2$$

$$\textcircled{2} \quad c_n(f, f) \geq \lambda_{(sob)} \|f\|_{H^{m+1}(\Omega)}^2$$

meaning that c_n is coercive on the normed space

$$H^{m+1}(\Omega, \| \cdot \|_{H^{m+1}(\Omega)})$$

Since $m+1 > \max(d_1/2, \kappa)$, one has

$$|c_n(f, g)| \leq \left[C_{\Omega}^2 + \left(\sum_{k \leq K} \|A_{k,\alpha}\|_{\infty, \Omega} \right)^2 + \lambda_{sob} \right] \|f\|_{H^{m+1}(\Omega)}^2 \|g\|_{H^{m+1}(\Omega)}^2$$

and

$$|B_n(f)| \leq C_{\Omega} \frac{1}{n} \sum_{i=1}^n \|y_i\|_2 \|f\|_{H^{m+1}(\Omega)}^2$$

therefore, c_n and B_n are continuous

By Lax-Milgram theorem [Brezis, 2010, Corollary 5.8]

$\exists! \hat{f} \in H^{m+1}(\Omega, R^{dat})$ such that

$$c_n(\hat{f}, \hat{f}) - 2B_n(\hat{f}) = \min_{f \in H^{m+1}(\Omega)} c_n(f, f) - 2B_n(f)$$

\hat{f} is the unique minimizer of $R_n^{(\text{reg})}$ over $H^{m+1}(\Omega, \mathbb{R}^{\text{dat}})$.
 Again Lax-Milgram thus gives that \hat{f} is the unique element of $H^{m+1}(\Omega)$ s.t. $\forall g \in H^m(\Omega)$ $R_n(\hat{f}, g) = R_n(g)$. ◻

✓ The former proposition ensures that the Sobolev regularization makes the PINN problem well-posed, i.e., the theoretical risk function admits a unique minimizer.

❓ Next question: what about CV to this unique minimizer?

Prop [From regularized risk consistency to strong convergence]

Assume that $m \geq \max([d/2], K)$

Let $(f_p)_p \in C^\infty(\bar{\Omega}, \mathbb{R}^{\text{dat}})$
 such that $\lim_{p \rightarrow \infty} R_n^{(\text{reg})}(f_p) = \inf_{f \in C^\infty(\bar{\Omega}, \mathbb{R}^{\text{dat}})} R_n^{(\text{reg})}(f)$

Then theoretical risk

$$\lim_{p \rightarrow \infty} \|f_p - \hat{f}\|_{H^m(\Omega)} = 0$$

unique minimizer of $R_n^{(\text{reg})}$
 over $H^{m+1}(\Omega, \mathbb{R}^{\text{dat}})$

Proof: admitted

Message: minimizing seqs of $R_n^{(\text{reg})}$ converge to the unique minimizer of $H^{m+1}(\Omega)$

Thm: [Strong convergence of regularized PINN]

Assume that

\mathcal{D} is an affine operator

$$m \geq \max \left[\frac{\dim}{2}, \kappa \right]$$

the condition function h is Lipschitz

$\det \left(\hat{\mathcal{D}}^{(\text{reg})} (p, n_e, n_r, \mathcal{D}) \right)_n$ be a min. sequence of the regularized empirical risk function

$$R_{n, n_e, n_r}^{(\text{reg})} (f_\theta) = R_{n, n_e, n_r} (f_\theta) + \lambda_{\text{ridge}} \|\theta\|_2^2 + \frac{\lambda_{\text{stab}}}{n_r} \sum_{l=1}^{n_r} \sum_{k \leq m+1} \|\mathcal{D}^\alpha f(x_k^l)\|^2$$

are the class of $\text{NN}_{\mathcal{H}} (\mathcal{D})$.

Then, choosing $\lambda_{\text{ridge}} = \min(n_e, n_r)^{-\kappa}$ with $\kappa = \frac{1}{12 + 4H(1 + (2 + H)(m + 2))}$ one has

$$\lim_{D \rightarrow \infty} \lim_{n_e, n_r \rightarrow \infty} \lim_{p \rightarrow \infty} \| f_{p, n_e, n_r, \mathcal{D}}^{(\text{reg})} - \hat{f}_n \|_{\mathcal{H}^m(\Omega)} = 0 .$$

unique minimize of $R_n^{(\text{reg})}$
over $\mathcal{H}^{m+1} (\Omega, \mathbb{R}^d)$

Message: Minimizing seqs of the empirical regularized risk converge to the unique minimize of the regularized theoretical risk

✓ The PINN $f_p^{(\text{reg})}$ converges to the unique minimize \hat{f}_n
of the theoretical risk

? What can we say in terms of CV to f^* ?

Then [Strong CV of linear PDE solver].

Same assumptions as before -

+ The PDE admits a unique solution f^* in $H^{m+1}(\Omega, \mathbb{R}^{d_{\text{out}}})$

Then,

$$\lim_{\lambda \xrightarrow{(2)} 0} \lim_{D \rightarrow 0} \lim_{n, r \rightarrow \infty} \lim_{p \rightarrow \infty} \| f_{\text{reg}}^{(p, n, r, D, \lambda, \omega)} - f^* \|_{H^m(\Omega)} = 0$$

For the hybrid modeling framework, we need to measure the gap between f^* and the physical prior.

Def: [Physical inconsistency]

For any $f \in H^{m+1}(\Omega, \mathbb{R}^{d_{\text{out}}})$,

$$\text{PI}(f) = \lambda \mathbb{E} \| \Pi(f)(X^{(1)}) - h(X^{(1)}) \|^2 + \frac{1}{|\Omega|} \int_{\Omega} D(u, x) dx$$

→ PI measures the modeling error.

Better the model, lower $\text{PI}(f)$

→ Remark that $P_n(f) = \frac{1}{n} \sum_{i=1}^n \| \Pi(f)(x_i) - y_i \|_2^2 + \text{PI}(f)$

Prop: Under the same assumptions as before
 Assume that $f^{\infty} \in \mathcal{H}^{m+1}(\Omega, R^{\text{data}})$, let

$$\lambda_{\text{reg}} = \frac{\log n}{\sqrt{n}} \quad \lambda_{\text{obs}} = \frac{1}{\sqrt{n}}$$

Then

STATISTICAL ACCURACY.

$$\lim_{D \rightarrow \infty} \lim_{n_r \rightarrow \infty} \lim_{p \rightarrow \infty} E \int_{\Omega} \| f_{\theta}^{(n)}_{(\text{reg})} - f^{\infty} \|_2^2 d\mu_X \lesssim \frac{\log^2(n)}{\sqrt{n}}$$

PHYSICAL CONSISTENCY.

$$\lim_{D \rightarrow \infty} \lim_{n_r \rightarrow \infty} \lim_{p \rightarrow \infty} E \left[PI \left(f_{\theta}^{(n)}_{(\text{reg})} \right) \right] \leq PI(f^{\infty}) + o(1)_{n \rightarrow \infty}$$

In the lecture : $\lambda_{\lambda} = \lambda_e = \frac{\log n}{\sqrt{n}}$ $\lambda_{\text{obs}} = \frac{1}{\sqrt{n}}$ $\lambda_{\text{data}} = 1$

Choice of λ_{ridge} unaffected by the rescaling since we take the asymptotic in n_r , n_r first.

* Choice in the article $\lambda_{\lambda} = 1$ $\lambda_e = 1$ $\lambda_{\text{obs}} = 1/\log n$ $\lambda_{\text{data}} = \sqrt{n}/\log n$

Choose λ_{ridge} al
 lybre w.r.t n_r (and
 me) such that
 $\lambda_{\text{ridge}} \rightarrow 0$
 $n_r, n_r \rightarrow \infty$