

Deep Learning in Scientific Computing 2023: Lecture 7

Siddhartha Mishra

Seminar for Applied Mathematics (SAM), D-MATH (and),
ETH AI Center,
ETH Zürich, Switzerland.



What you learnt so far

- physics informed NNs .*
- ▶ PINNs to Solve PDEs.
 - ▶ Great for some PDEs, particularly with low amount of training data.
→ *Fairability, spectral bias*
 - ▶ Have several negatives. → *Fairability, spectral bias*
 - ▶ We need alternatives !!
 - ▶ When more data is available: → *Learn the solution*
 - ▶ The next several lectures: Use of Supervised Deep Learning for PDEs
→ *Data.*

What does solving a PDE mean? - finding the sol. Operator

- ▶ Example 1: Consider Darcy PDEs:

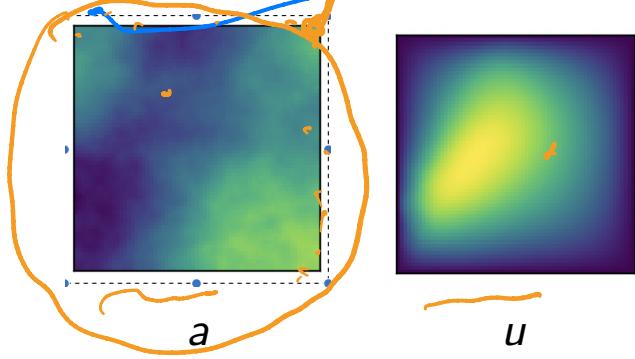
$$-\operatorname{div}(a \nabla u) = f,$$

prototypical
elliptic PDE
 $\operatorname{div}(a) = \sum \partial_{x_i} u_i$

- ▶ Quantities of interest are:

- ▶ u is temperature or pressure.
- ▶ a is conductance or permeability.
- ▶ f is the source.

Fix the source S
Given a
Find u ?



- ▶ Find the solution

Operator $\mathcal{G} : a \mapsto \mathcal{G}a = u$.
function

What does solving a PDE mean ? - Find the Soln Operators

- Example 2: Consider the Compressible Euler equations:

$$E = \frac{P}{\delta-1} + \frac{L}{2} P |u|^2$$

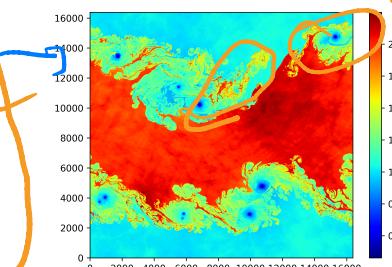
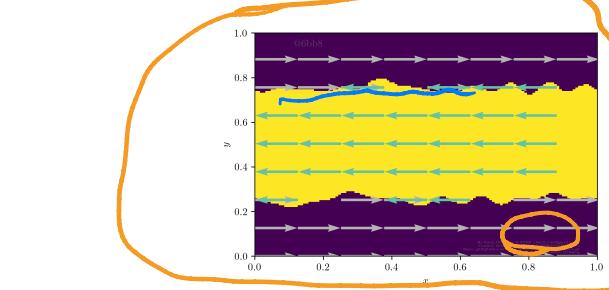
density

$$\rho_t + \operatorname{div}(\rho v) = 0, \quad \text{--- may be true}$$

$$(\rho v)_t + \operatorname{div}(\rho v \otimes v + p I) = 0, \quad \text{---}$$

$$E_t + \operatorname{div}((E + p)v) = 0., \quad \text{---}$$

$$u(x, 0) = (\rho, v, E)(x, 0) = a(x).$$

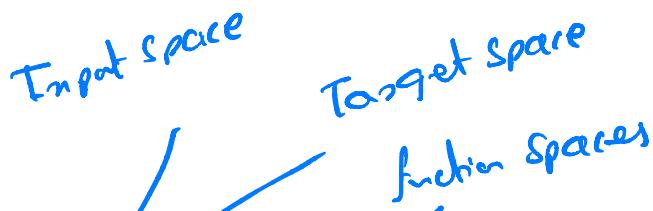


Initial Condition

Solution at time T

- Find the solution Operator $\mathcal{G} : a \mapsto \mathcal{G}a = u(T)$.

Setup



- X, Y are Banach spaces and $\mu \in \text{Prob}(X)$
- Abstract PDE: $\mathcal{D}_a(u) = f$ — Darcy, Euler, Navier-Stokes.
- Solution Operator: $\mathcal{G}: X \mapsto Y$ with $\mathcal{G}(a, f) = u$
- Task: Learn Operators from data
- Core of Operator Learning
- A Problem: DNNs map finite dimensional inputs to outputs !!

$$\mathcal{L}^*: \mathbb{R}^{d_m} \xrightarrow{\quad} \mathbb{R}^{d_{\text{out}}} \quad d_m, d_{\text{out}} \gg ?$$

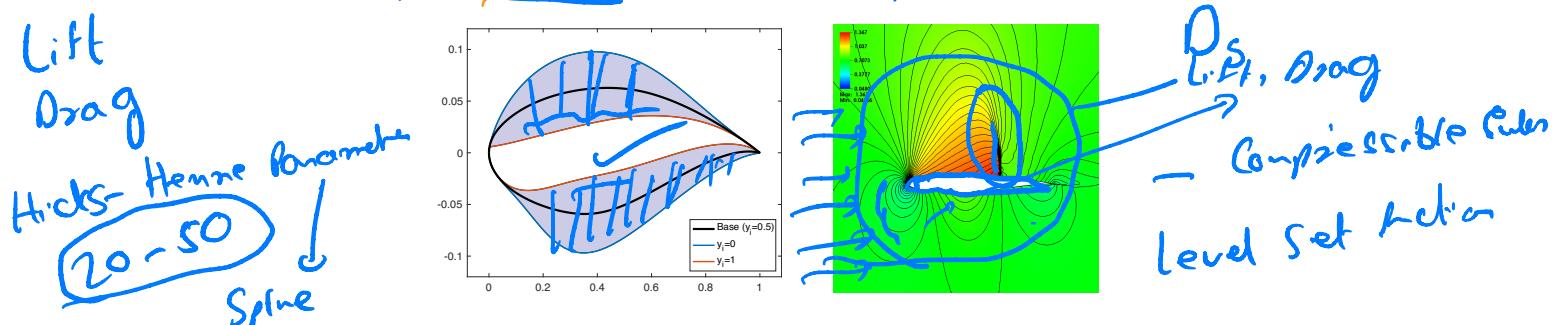
$\downarrow c + \infty$ $\downarrow c + \infty$

Solution I: Use Parametric PDEs instead :-)

"Cheat"

- ▶ X, Y are Banach spaces and $\mu \in \text{Prob}(X)$
 - ▶ Abstract PDE: $\mathcal{D}_a(u) = f$
 - ▶ Solution Operator: $\mathcal{G} : X \mapsto Y$ with $\mathcal{G}(a, f) = u$
 - ▶ Simplified Setting: $\dim(\overline{\text{Supp}(\mu)}) = d_y < \infty$
 - ▶ Corresponds to Parametrized PDEs with finite parameters.
 - ▶ Find Soln $u(t, x, y)$ or Observable $\mathcal{L}(y)$ for $y \in Y \subset \mathbb{R}^{d_y}$.

↓
 Finite dimensions.
 $a = a(y)$
 $y \in Y \subset \mathbb{R}^{d_y}$



- ▶ Approximate Fields or **observables** with **deep neural networks**

$G: a \mapsto u$

$a \in L^2(\Omega)$

φ_k is some basis - Fourier basis

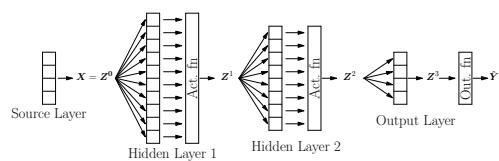
$$a = \sum_{k=1}^{\infty} c_k \varphi_k$$

$$a = \sum_{k=1}^{kd} c_k \varphi_k$$

Approximation:

$$a \mapsto \{c_k\}_{k=1}^{kd}$$

Supervised learning of target \mathcal{L} with Deep Neural networks



- ▶ $\mathcal{L}^*(z) = \sigma_o \odot C_K \odot \sigma \odot C_{K-1} \dots \odot \sigma \odot C_2 \odot \sigma \odot C_1(z)$.
- ▶ At the k -th **Hidden layer**: $z^{k+1} := \sigma(C_k z^k) = \sigma(W_k z^k + B_k)$
- ▶ **Tuning Parameters**: $\theta = \{W_k, B_k\} \in \Theta$, *High-dimensional Inputs*
- ▶ σ : scalar **Activation function**: ReLU, Tanh
- ▶ **Random Training set**: $\mathcal{S} = \{z_i\}_{i=1}^N \in Z$, with i.i.d z_i *Random Training data*
- ▶ Use **SGD (ADAM)** to find $\mathcal{L} \approx \mathcal{L}^* = \mathcal{L}_{\theta^*}^*$

$$\theta^* := \arg \min_{\theta \in \Theta} \sum_{i=1}^N |\mathcal{L}(z_i) - \mathcal{L}_\theta^*(z_i)|^p,$$

↳ obtained from *measur. nets*
/Numerical Simulation

Supervised learning for high-d Parametric PDEs

- d = Ground Truth ($\mathbb{R}^d \rightarrow \mathbb{R}$)*
- ▶ Can we find DNN such that $\|\mathcal{L}^* - \mathcal{L}\| \sim \mathcal{O}(\epsilon)$?
 - ▶ YES: Universal Approximation Property of DNNs \Rightarrow :
Given any Continuous (measurable) \mathcal{L} , exists a $\hat{\mathcal{L}}$:
 - Integrability $\|\mathcal{L} - \hat{\mathcal{L}}\| < \epsilon$ (✓)
 - ▶ If $\mathcal{L} \in W^{s,p}$, \exists DNN $\hat{\mathcal{L}}$ with M parameters (Yarotsky):
Smoothness $\|\mathcal{L} - \hat{\mathcal{L}}\|_p \sim \mathcal{O}(M^{-\frac{s}{d}})$ *smooth man one causes to learn*
 - ▶ But in Scientific Computing (often):
 - ▶ \mathcal{L} is not be very Regular and $\bar{d} \gg 1$ 10^{12}
 - ▶ If $\mathcal{L} \in W^{1,\infty}$, $\bar{d} = 6$ 1% error, need network of size 10^{12} !!
 - ▶ Curse of dimensionality: DNN Size $M \sim \epsilon^{-\frac{\bar{d}}{s}}$ - Exponentially in dimension

Refined Error Estimates

- ▶ Error $\mathcal{E} := \|\mathcal{L} - \mathcal{L}^*\|_p$ Decomposition: $\mathcal{E} \leq \mathcal{E}_{app} + \mathcal{E}_{gen} + \mathcal{E}_{opt}$.
- ▶ Approximation error $\mathcal{E}_{app} = \|\mathcal{L} - \hat{\mathcal{L}}\|_p$,
 - ▶ $\hat{\mathcal{L}}$ is best approximation of \mathcal{L} in $NN(M)$.
 - ▶ One can prove that $\mathcal{E}_{app} \sim \mathcal{O}(\bar{d}^\sigma M^{-\eta})$ for, —
 - ▶ Linear Elliptic PDEs: (Schwab, Kutyniok et al).
 - ▶ Semi-linear Parabolic PDEs: (E, Jentzen et al).
 - ▶ Nonlinear Hyperbolic PDEs: (DeRyck, SM, 2021).
- ▶ Optimization Error $\mathcal{E}_{opt} \sim$ Computable Training error: .
- ▶ Generalization Error — finite sample size error

$$\mathcal{E}_{gen}(\theta) := \|\mathcal{L} - \mathcal{L}_\theta^*\|_p^p - \frac{1}{N} \sum_i |\mathcal{L}_\theta^*(y_i) - \mathcal{L}(y_i)|^p$$

S. ~~the~~ (Quadratic Error)

- ▶ Using Concentration inequalities + Covering number bounds:

$$\mathcal{E}_{gen} \sim \frac{C(M, \log(\|W\|)) \log(\sqrt{N})}{\sqrt{N}}$$

Well-trained Networks

- ▶ Assume we can find DNN such that $\mathcal{E}_{app}, \mathcal{E}_T \ll 1$
- ▶ Still **Overall Error** behaves as

$$\mathcal{E} \sim \frac{C(M)}{N^\alpha}, \quad \alpha \leq \frac{1}{2}.$$

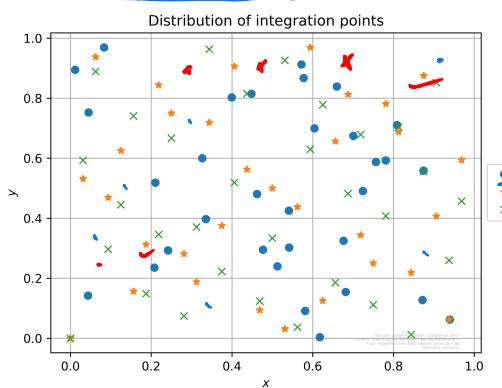
10000
100

- ▶ If $C(M) \sim \mathcal{O}(1)$, error of 1% requires 10^4 training samples !!
- ▶ Challenge: learn maps of low regularity in a data poor regime
- ▶ Contrast with Big Data successes of machine learning.

A Trick: Lye, SM, Ray, 2020

Tricks from Numerical Analysis.

- ▶ Use Low discrepancy sequences $\{y_i\}_{i=1}^N \in Y$ as Training Set



$$\int f(t) dy \approx \frac{1}{m} \sum_{i=1}^m f(t_i)$$

$t_i \rightarrow \text{LOS}$

- ▶ These sequences are **Equidistributed** (better spread out).
- ▶ Examples: **Sobol**, Halton, Owen, Niederreiter ++
- ▶ Basis of **Quasi-Monte Carlo** (QMC) integration.

quasi-random number generators

Training on Low-Discrepancy Sequences

$$\int_0^1 \left| \frac{\partial^d \mathcal{L}}{\partial y_1, \dots, \partial y_d} \right| dy < +\infty$$

- ▶ For \mathcal{L} with Bounded Hardy-Krause variation and smooth σ .
- ▶ Generalization Error for Sobol sequences: (SM, Rusch, 2020),

$$\mathcal{E} \leq \mathcal{E}_T + C(V_{HK}(\mathcal{L}), V_{HK}(\mathcal{L}^*)) \frac{(\log N)^d}{N},$$

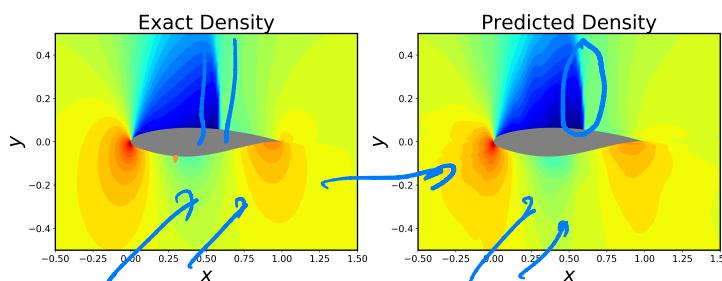
Turing Error

Random: $\mathcal{E} \leq \mathcal{E}_r + \frac{C(\alpha)}{\sqrt{N}} \left[\frac{1}{N} \right]$

Prediction

- 20-dimesn* *Sobol pts*
- Given Hicks-Henne parameter: Predict Drag, Lift, Flow
 - DNN with $10^3 - 10^4$ parameters and 128 training samples :

	Run time (1 sample)	Training	Evaluation	Error
Lift	2400 s <i>-40</i>	700 s	10^{-5} s	0.78%
Drag	2400 s	840 s	10^{-5} s	1.87%
Field	2400 s	<u>1 hr</u>	0.2 s	1.9%



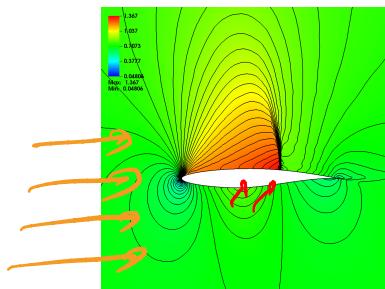
- Errors with **Random** Training pts: Lift 8.2%, Drag: 23.4%

Forward UQ – Uncertainty quantification

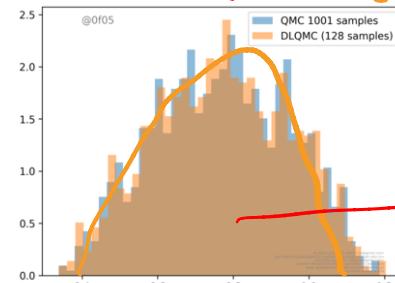
- DL-UQ algorithm of Lye,SM,Ray, 2020 is

$$\mathcal{L} \# \mu \approx \frac{1}{M} \sum_{i=1}^M \delta_{\mathcal{L}(y_i)} \approx \frac{1}{M} \sum_{i=1}^M \delta_{\underline{\mathcal{L}^*(y_i)}} - \text{on N}$$

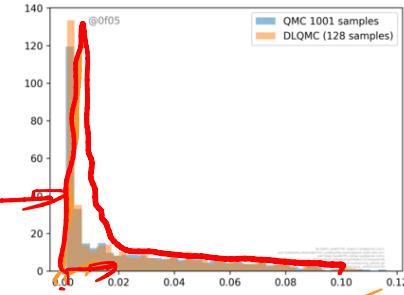
Empirical distrib.



Sample



Lift PDF ✓



Drag PDF

Observable	Speedup (MC)	Speedup (QMC)
Lift	246.02	6.64
Drag	179.54	8.56