

Deep Learning in Scientific Computing

Shizheng Wen

shiwen@student.ethz.ch

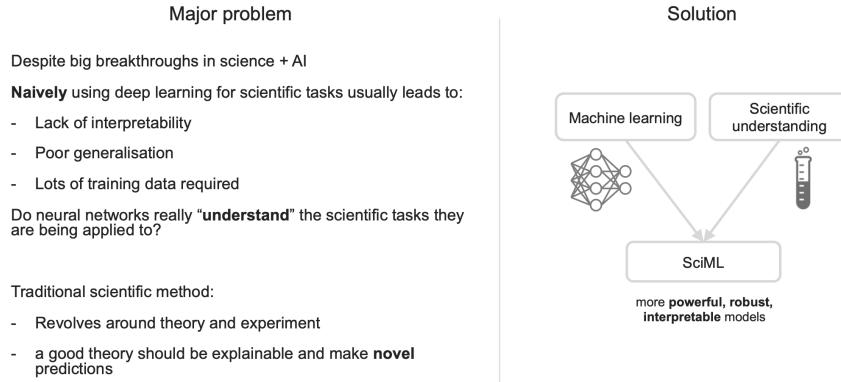
Table of contents

- 1 Introduction
- 2 Introduction to Deep Learning
- 3 Introduction to Deep Learning
- 4 PINN
- 5 PINN - applications
- 6 PINN
- 7 Operator Learning
 - 7.1 Introduction
 - 7.2 Operator Networks and DeepONet
 - 7.3 Lecture 10

1 Introduction

Definition of Scientific machine learning (SciML)

Scientific machine learning (SciML)



Key scientific tasks and how machine learning can help

In science, there are mainly three types of question. We will see how ML can help us.

- Simulation
 - Usually used as a starting point for other tasks (e.g. inverse / control / design problems)
 - $$\mathbf{b} = F(\mathbf{a})$$

\mathbf{a} = set of input conditions

F = physical model of the system

\mathbf{b} = resulting properties given F and \mathbf{a}

- Analytical solutions do not exist, and we must resort to numerical modeling. By the way, most of the questions based on the system modelled by PDEs.
- Inverse, control, and data assimilation problems

–
$$\mathbf{b} = F(\mathbf{a})$$

\mathbf{a} = set of input conditions

F = physical model of the system

\mathbf{b} = resulting properties given F and \mathbf{a}

- Tricks for solving inverse problems

◦

$$b = F(\alpha)$$

- Fundamentally, inverse problems are **search** problem
- It is often useful to frame them as an optimisation problem, for example:

$$\min_{\hat{\alpha}} \|b - F(\hat{\alpha})\|^2$$

- If F is differentiable, one option is to use gradient-based methods (e.g. **gradient descent**)
- Otherwise, we can use gradient-free methods (e.g. evolutionary algorithms, Bayesian optimisation, brute-force search, ...)

- If we want to derive the parameters that gives rise to the solution of b , then we treat it as a optimization problem and use gradient-based or gradient-free methods to solve it. Gradient-based means that the update of α is based on gradient of the function.

- For control, α is a function.

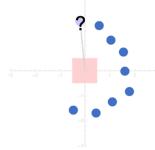
◦

$\alpha = f(t)$, force applied to cart
 b = pendulum stays balanced (i.e. $\theta(t) \rightarrow 0$)
 F = method for solving equations of motion

- For data assimilation

◦

Data assimilation



$\alpha = x(t = 0), \theta(t = 0)$,
 b = Noisy measurements of $x(t_i), \theta(t_i)$
 F = Method for solving equations of motion and noise model

- We have a bunch of noisy measurements of data, we want to know where we start of the system.

- Equation discovery and anomaly detection

—

$$b = F(\alpha)$$

α = set of input conditions
 F = physical model of the system
 b = resulting properties given F and α

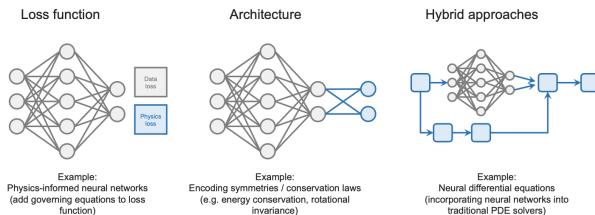
- Historically, F (= laws of physics) has been found through remarkable human **intuition**. From a computational standpoint, discovering physics is like solving an inverse problem (trying to fit a model to observed data). But the model should be **explainable**, **generalisable** and make **novel** predictions.

Ways to incorporate scientific principles into deep learning

- Example - Simulation: Fourier Neural Operator (Application: FourCastNet)

- Example - Inverse problems: use a neural network to update gradients. It's hard to understand such problems. View this paper for more information: Solving ill-posed inverse problems using iterative deep neural networks, Inverse Problems (2017)
- Example - Equation Discovery: Fit a neural network to observed data, then regress over a library of gradients to "discover" underlying equations. Here, we don't need to use complex regression methods, only linear regression is enough for finding the underlying equation.

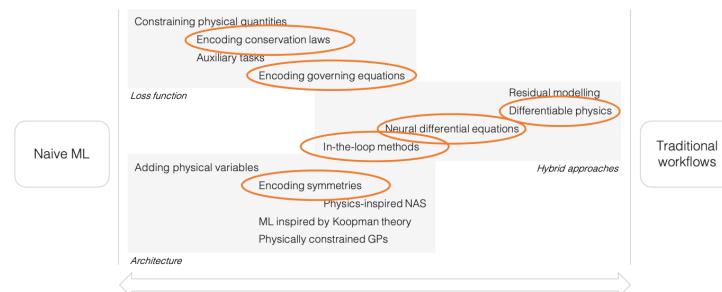
Ways to incorporate scientific principles into machine learning



A plethora of SciML techniques

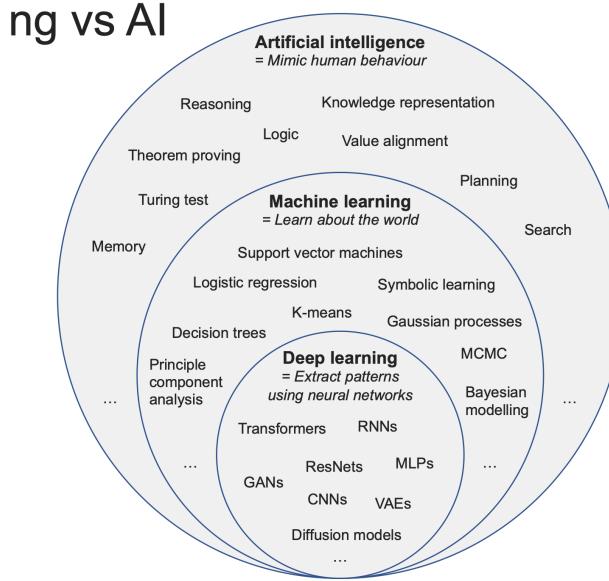
	Scientific task		
	Forward simulation	Inversion	Equation discovery
Architecture			
Adding physical variables	Daw et al.		
Encoding symmetries	Ling et al., Wang et al., Anderson et al., Schmitt et al.	Panju and Ghodsi	Udrescu et al.
Physics-inspired NAS	Ba et al., Panju and Ghodsi		
ML inspired by Koopman theory	Geneva and Zabaras, Lusch et al.		
Physically constrained GPs		Raissi et al., Raissi and Karniadakis	
Other approaches	Jumper et al., Mohan et al.		
Loss function			
Constraining physical quantities	Karpatne et al., Zhang et al., Benjamin et al., Erichson et al., Xie et al., Brehmer et al.		
Encoding conservation laws	Beucler et al., Zeng et al.		Greydanus et al., Toth et al., Cranmer et al.
Auxiliary tasks	de Oliveira et al.		
Encoding governing equations	Raissi et al., Jin et al., Jin et al., Chen et al., Kharazmi et al., Yang et al., Yang et al., et al., Wang et al., Wang et al., Li et al., Zhu et al., Geneva and Zabaras, Gao et al.	Chen et al., Champion et al.	
Hybrid approaches			
Residual modelling	Pawar et al.	Jiang et al., Ren et al., Minkov et al., Wirth et al., Zhang et al.	
Differentiable physics			Chen et al., Rackauckas et al., Long et al.
Neural differential equations			
In-the-loop methods	Um et al., Rasp et al.	Adler and Oktem, Morningstar et al., Hammerl et al., Li et al., Lunz et al., Bora et al., Mosser et al.	

SciML
technique



2 Introduction to Deep Learning

In this chapter, firstly, the lecturer introduce the concept of Machine learning and neural network. Unlike the polynomial regression having closed-form solution, the only way to train neural network is by optimization method. Then the lecturer introduce the CNN, which has advantages for processing spatial correlations input.



Multilayer perceptrons

Simple polynomial regression

$$\hat{y}(x; \theta) = \theta_4 x^3 + \theta_3 x^2 + \theta_2 x + \theta_1$$

To fit, use least-squares:

$$\theta^* = \min_{\theta} \sum_i^N (\hat{y}(x_i; \theta) - y(x_i))^2 \quad (1)$$

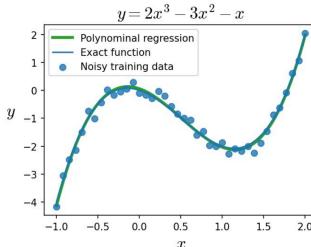
Re-write using linear algebra:

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \dots \end{pmatrix} = \begin{pmatrix} 1 & x_1 & x_1^2 & x_1^3 \\ 1 & x_2 & x_2^2 & x_2^3 \\ \dots & \dots & \dots & \dots \end{pmatrix} \begin{pmatrix} \theta_1 \\ \theta_2 \\ \theta_3 \\ \theta_4 \end{pmatrix} \text{ or } \hat{Y} = \Phi^T \theta$$

$$\theta^* = \min_{\theta} \|\Phi^T \theta - Y\|^2$$

In this case, it can be shown (1) has an analytical solution:

$$\theta^* = (\Phi^T \Phi)^{-1} \Phi^T Y$$



Neural network regression

$$\hat{y}(x; \theta) = NN(x; \theta)$$

To fit, use least-squares:

$$\theta^* = \min_{\theta} \sum_i^N (NN(x_i; \theta) - y(x_i))^2 \quad (2)$$

In general, no analytical solution to (2) exists, so we must use **optimization**

For example, gradient descent:

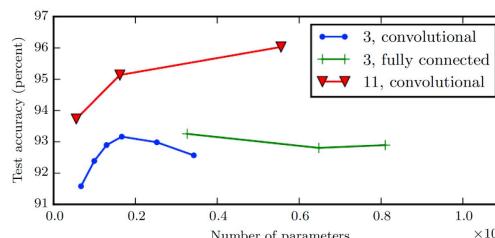
$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial \sum_i^N (NN(x_i; \theta) - y(x_i))^2}{\partial \theta_j}$$

or equally

$$\theta_j \leftarrow \theta_j - \gamma \frac{\partial L(\theta)}{\partial \theta_j}$$

Where γ is the learning rate and $L(\theta)$ is the **loss function**

- Function fitting between polynomial regression and neural network. For polynomial regression, we have the closed form.



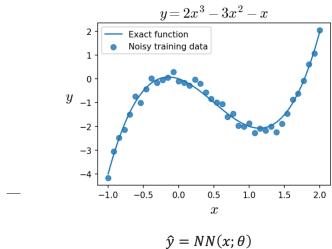
Goodfellow et al, Multi-digit number recognition from street view imagery using deep convolutional neural networks, ICLR (2014)

- 当模型复杂度上升，模型有过拟合的迹象，在深度学习中，model parameters > training dataset的时候，模型的accuracy又会进一步上升，效果又会进一步变好。
- Empirically, deep neural networks perform better than shallow neural networks => encode a very general belief that the true function is **composed** of simpler functions

Popular deep learning tasks

Supervised learning

- Regression

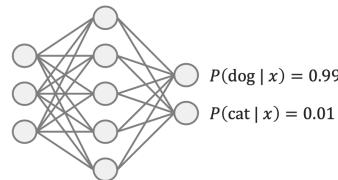


Loss function (mean squared error)

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y(x_i))^2$$

- Here, due to the reason that the observed data will be influenced by noise. It is better to use probabilistic perspective.
- 最大似然估计是一种统计方法，它用来求一个样本集的相关概率密度函数的参数。要在数学上实现最大似然估计法，我们首先要定义可能性： $lik(\theta) = f_D(x_1, x_2, \dots, x_n | \theta)$ 。并且在 θ 的所有取值上，使这个函数最大化。这个使可能性最大的值即被称为 θ 的最大似然估计。
- Regression problem corresponds to Continuous PDF, while classification corresponds to the discrete PDF.

- Classification



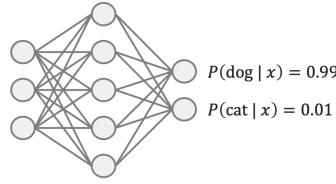
Supervised learning - classification:

Given a set of example inputs and outputs (labels) $\{(x_1, y_1), \dots, (x_N, y_N)\}$ drawn from the discrete probability distribution $P(y|x)$

where $y \in Y$, for example, $Y = \{\text{dog, cat}\}$

Find

$$\hat{P}(y|x, \theta) \approx P(y|x)$$



Then assume

$$\hat{P}(y|x, \theta) = \prod_j^c NN(x; \theta)_j^{y_j}, \quad \sum_j^c NN(x; \theta)_j = 1$$

Then, assume each training datapoint is independently and identically distributed (iid), then the **data likelihood** can be written as:

$$P(D|\theta) = \hat{P}(x_1, y_1, \dots, x_n, y_n | \theta) = \prod_i^N \hat{P}(y_i | x_i, \theta)$$

Then use **maximum likelihood estimation (MLE)** to estimate θ^* :

Let each class be encoded as a one-hot vector of length C , e.g.

$$y = (0, 0, 1, 0)$$

Typically, we use a softmax output layer to assert $\sum_j^c NN(x; \theta)_j = 1$:

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_j^C e^{z_j}}$$

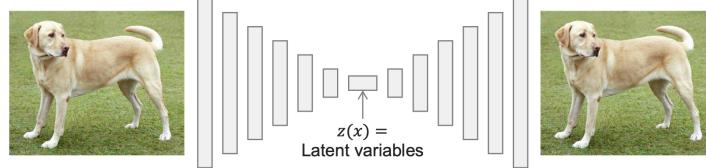
$$\begin{aligned} \theta^* &= \max_{\theta} \hat{P}(D|\theta) \\ &= \max_{\theta} \prod_i^N \prod_j^C NN(x_i; \theta)_j^{y_{ij}} \\ &= \min_{\theta} - \sum_i^N \sum_j^C y_{ij} \log NN(x_i; \theta)_j \end{aligned}$$

Also known as the **cross-entropy loss**

Unsupervised learning

- Feature learning

—



Loss function

Many different possibilities, a simple choice is

$$L(\theta) = \sum_i^N (NN(x_i; \theta) - x_i)^2$$

Unsupervised learning – feature learning

Given a set of examples $\{x_1, \dots, x_N\}$, find some features $z(x)$

Which are salient descriptors of x , where $x \in \mathbb{R}^n, z \in \mathbb{R}^d$

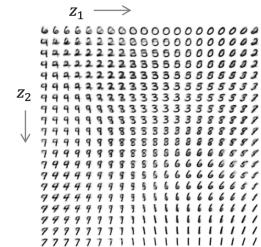
Typically, $d \ll n$ (= compression)

z can be used for downstream tasks, e.g. clustering / classification

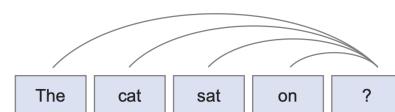
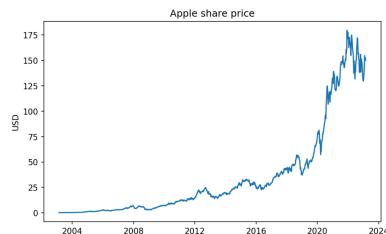
For example:

Variational autoencoders (VAEs)

Kingma et al, 2014

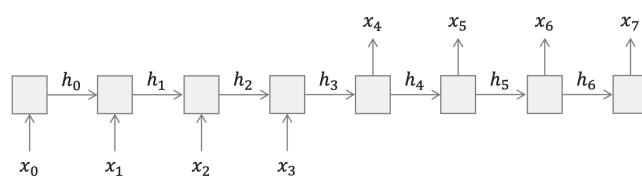


- Autoregression



For example:

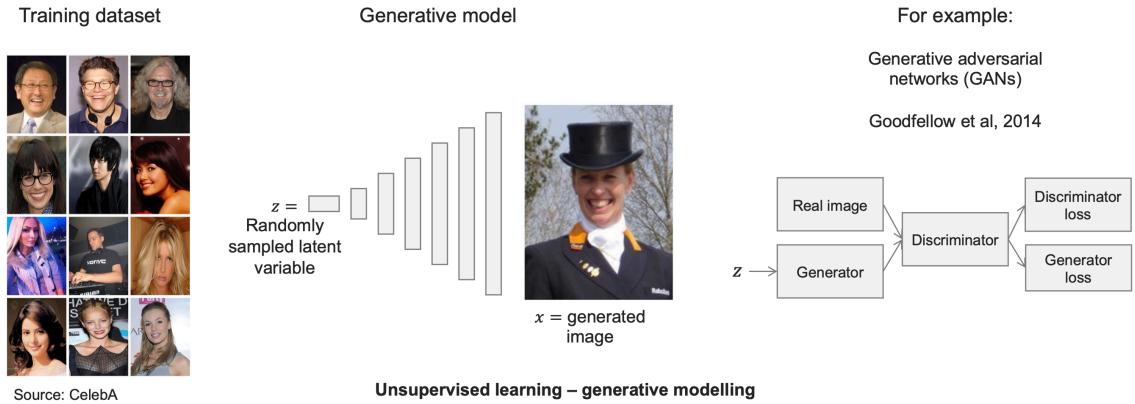
ChatGPT



—

- Given many examples sequences, train a model to predict future values from past values

- Generative models



- - High-dimensional Statistics

Training a deep neural network

Gradient descent

We have discussed the gradient and chain rule, let's think about dimensionality:

$$NN(x; \theta) = W_3(\sigma(W_2\sigma(W_1x + b_1) + b_2) + b_3) = f \circ g \circ h(x; \theta)$$

Let's think about **dimensionality**. Let f have 1 element (scalar output), g and h both have 100 elements (100 hidden units) and x be a flattened 128 x 128 image, then W_1 has shape $(100 \times (128 \times 128)) = 1.6M$ elements and

$$\frac{\partial NN}{\partial W_1} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial W_1}$$

$$\begin{array}{c} \text{---} \\ \text{---} \end{array} = \begin{array}{c} \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

$(1 \times 1.6M) = (1 \times 100)(100 \times 100)(100 \times 1.6M)$

⚠ The last Jacobian is extremely large! = 160M elements (or 0.7 GB!)

Forward- versus reverse-mode differentiation

$$\frac{\partial NN}{\partial W_1} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial h} \frac{\partial h}{\partial W_1}$$

$$\begin{array}{c} \text{---} \\ \text{---} \end{array} = \begin{array}{c} \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ \text{---} \end{array} \quad \begin{array}{c} \text{---} \\ \text{---} \end{array}$$

$(1 \times 1.6M) = (1 \times 100)(100 \times 100)(100 \times 1.6M)$

Consider evaluating the chain rule (RHS):

1) From right to left (**forward-mode**)

$$(100 \times 100)(100 \times 1.6M) \rightarrow (100 \times 1.6M) \\ (1 \times 100)(100 \times 1.6M) \rightarrow (1 \times 1.6M)$$

= lots of computation (large matrix-matrix multiply)

2) From left to right (**reverse-mode**)

$$(1 \times 100)(100 \times 100) \rightarrow (1 \times 100) \\ (1 \times 100)(100 \times 1.6M) \rightarrow (1 \times 1.6M)$$

= much less computation (vector-matrix multiplies)

💡 => Order matters! Typically, reverse mode differentiation is more efficient for scalar loss functions with large numbers of parameters ("wide" Jacobians), whilst forward mode is more efficient for evaluating "tall" Jacobians

- - 这就是自动微分方法的引入背景，我们在前向传播时，会记录下当前参数下各层的数值，则我们可以进行left to right的计算。

- 接着，我们计算发现其实有些层是可以用outer product来进行表示。

—
But note:

$$\frac{\partial \mathbf{h}}{\partial W_1} = \frac{\partial(W_1 \mathbf{x} + \mathbf{b}_1)}{\partial W_1} = \begin{pmatrix} x_1 & x_2 & \cdots & x_n & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 & x_1 & x_2 & \cdots & x_n \end{pmatrix}$$

Then

$$\begin{aligned} \frac{\partial NN}{\partial W_1} &= \frac{df}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial W_1} = \left(\frac{df}{\partial h_1}(x_1, \dots, x_N), \dots, \frac{df}{\partial h_{100}}(x_1, \dots, x_N) \right) \\ &= \frac{\partial f^T}{\partial \mathbf{h}} \otimes \mathbf{x} \end{aligned}$$

This is just the (flattened) **outer product** of two vectors (100×1) \otimes ($16,384 \times 1$)

\Rightarrow We don't have to fully populate the last Jacobian ($\frac{\partial h}{\partial W_1}$) when computing its vector-Jacobian product
 \Rightarrow Computing vector-Jacobian products is usually simpler and **more** efficient than directly computing Jacobians

$$\frac{\partial NN}{\partial W_1} = \frac{\partial f}{\partial \mathbf{g}} \frac{\partial \mathbf{g}}{\partial \mathbf{h}} \frac{\partial \mathbf{h}}{\partial W_1}$$



$$(1 \times 1.6M) = (1 \times 100)(100 \times 100)(100 \times 1.6M)$$

Note:

Training neural networks =
Lots of matrix – matrix
multiplications

These types of parallelisable
operations are perfect for GPUs

10x to 100x speed-ups compared
to CPUs

2) From left to right (**reverse-mode**)

$$(1 \times 100)(100 \times 100) \rightarrow (1 \times 100)
(1 \times 100)(100 \times 1.6M) \rightarrow (1 \times 1.6M)$$

= Efficient training code

Allows us to train neural networks with **billions** of
parameters

- \circ 这告诉我们不需要完全计算出Jacobian之后再去做矩阵乘法，计算的时候就可以算vector-Jacobian。这会更加简单和高效。
- 因此我们策略如下，反向传播算法

Forward pass:

$$\mathbf{x}_i \rightarrow \mathbf{h}_i = W_1 \mathbf{x}_i + \mathbf{b}_1 \rightarrow \mathbf{g}_i = W_2 \sigma(\mathbf{h}_i) + \mathbf{b}_2 \rightarrow f_i = W_3 \sigma(\mathbf{g}_i) + \mathbf{b}_3$$

Save layer outputs in forward pass

Backward pass:

$$\frac{\partial L}{\partial W_1} = \sum_i^N 2(f_i - y(\mathbf{x}_i)) W_3 \text{ diag}(\sigma'(\mathbf{g}_i)) W_2 \text{ diag}(\sigma'(\mathbf{h}_i)) \otimes \mathbf{x}_i$$

Evaluate from left to right (reverse-mode)

Batch matrix multiplies over training examples \mathbf{x}_i

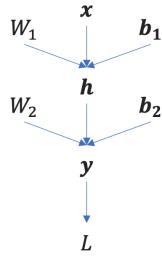
Similar equations for other weight matrices and bias vectors

- \circ In practice, **Autodifferentiation** tracks all your forward computations and their gradients and applies the chain rule automatically for you, so you don't have to worry!

◦

$$\mathbf{y} = W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

Autodifferentiation = gradients of arbitrary programs ☀️



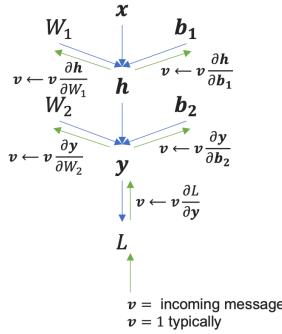
Fundamentally, autodifferentiation builds a directed **graph of primitive operations**

Each primitive operation defines:

- 1) **Forward operation**
- 2) **vector-Jacobian product** (for reverse mode autodiff)
- 3) **Jacobian-vector product** (for forward-mode autodiff)

◦

$$\mathbf{y} = W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$

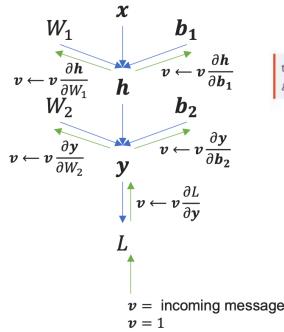


Then, autodifferentiation can be thought of as **message-passing** between the graph's nodes

Each message computes the vjp or jvp

Combined, evaluates the **chain rule**

$$\mathbf{y} = W_2 \sigma(W_1 \mathbf{x} + \mathbf{b}_1) + \mathbf{b}_2$$



PyTorch

```
torch.autograd.backward(tensors, grad_tensors=None, retain_graph=None, create_graph=False, grad_variables=None, inputs=None) [SOURCE]
```

Computes the sum of gradients of given tensors with respect to graph leaves.

The graph is differentiated using the chain rule. If any of `tensors` are non-scalar (i.e. their data has more than one element) and require gradient, then the Jacobian-vector product will be computed, in this case the function additionally requires specifying `grad_tensors`. It should be a sequence of matching length, that contains the “vector” in the Jacobian-vector product, usually the gradient of the differentiated function w.r.t. corresponding tensors (`None` is an acceptable value for all tensors that don't need gradient tensors).

Note autodiff is **not**

- Symbolic differentiation
- Finite differences

It is a way of efficiently computing exact gradients

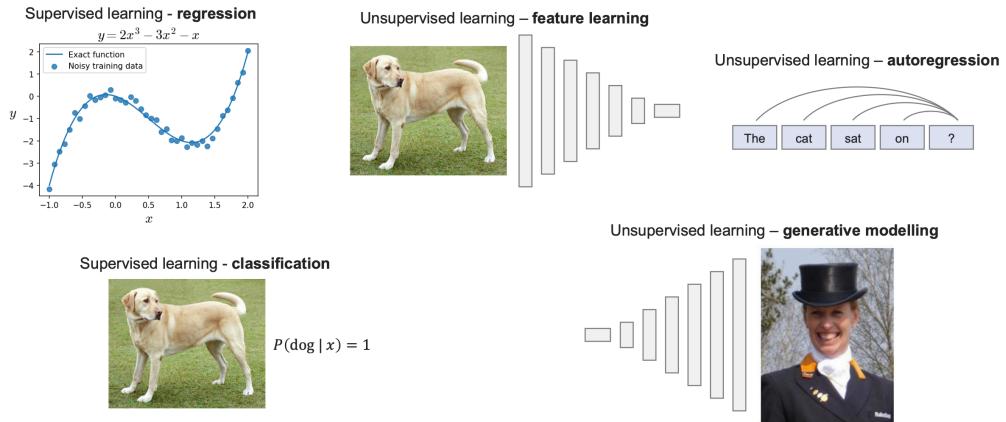
`loss.backward()`

- Autodiff是一种高效计算计算图在某点出的gradients的方法，不是基于符号微分（相当于任意点处的微分）。

3 Introduction to Deep Learning

Recap

Popular deep learning tasks



Automatic differentiation

Chain rule can be efficiently evaluated using vector-Jacobian products (**reverse-mode differentiation**)

Lecture Overview

- Challenges of function fitting
 - Overfitting / underfitting
 - Bias / variance
- Regularising deep neural networks
 - Architecture
 - Training data
 - Loss function
- Optimising deep neural networks
 - Stochastic gradient descent
 - Adam / higher-order
- State-of-the-art models
 - Transformers, ChatGPT

Error Sources

Training loss:

$$L(\theta) = \frac{1}{N} \sum_i^N (NN(x_i; \theta) - y_i)^2$$

$$D = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

This is also known as the **empirical loss**, and can be more generally written as:

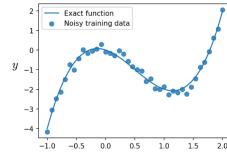
$$L(\theta) = \frac{1}{N} \sum_i^N l(NN(x_i; \theta), y_i), \quad x, y \sim p(x, y)$$

But what we really want to minimise is the **expected loss**:

$$\mathcal{L}(\theta) = \iint l(NN(x; \theta), y) p(x, y) dx dy$$

$$= E_{(x,y) \sim p}[l(x, y; \theta)]$$

But this is impossible in practice!



This means there are two sources of error:

1) **Approximation error** (limited expressivity of neural network)

$$\epsilon_{app} = \mathcal{L}(NN^*) - \mathcal{L}(y^*)$$

NN which minimises empirical loss True function which minimises expected loss

2) **Estimation error** (finite amount of training data) (aka generalisation error)

$$\epsilon_{est} = \mathcal{L}(NN) - \mathcal{L}(NN^*)$$

NN which minimises empirical loss NN which minimises expected loss

Regularization

Regularisation = **restrict** the space of possible functions a neural network can learn / **prefer** certain regions of the function space / impart a **prior** on the learning algorithm

Ways to regularise neural networks:

- Reduce model complexity
- Modify model architecture

4 PINN

curse of dimensionality:

- Conventional numerical methods (FD): error: $\epsilon = O(\delta x^2)$ $\#(\text{pts}) = (\frac{1}{\delta x})^{d+2} \Rightarrow \delta x = (\#\text{point})^{d+2}$
 - which means that if we increase the dimension of the problem. If we don't change the sampling points, the error of the problem will increase exponentially. Conversely, if we want to ensure the resolution/accuracy, the increased work will grow exponentially. **curse of dimensionality**

- PINN

► PINN with Depth 4, Width 20, Interior training points 2^{16} , Boundary points 2^{15} , Total

Dimension	Training Error	Generalization error
1	2.8×10^{-5}	0.0035%
5	0.0002	0.016%
10	0.0003	0.03%
20	0.006	0.79%
50	0.006	1.5%
100	0.004	2.6%

- - sublinearly growth.

5 PINN - applications

6 PINN

7 Operator Learning

7.1 Introduction

What does solving a PDE mean? -> find the solution operator

Two example:

- Darcy PDEs : $-\operatorname{div}(\mathbf{a} \nabla u) = f$ solution operator $\mathcal{G} : \mathbf{a} \rightarrow \mathcal{G}\mathbf{a} = u$
 - prototypical elliptic PDE (time-independent)
- Compressible Euler equations

—

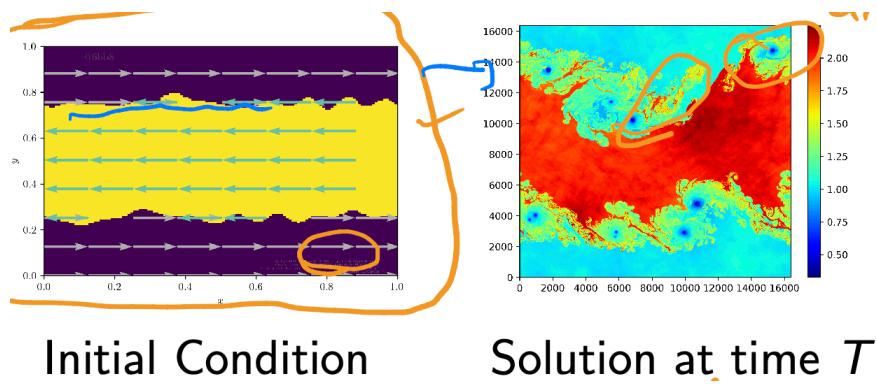
$$E = \frac{\rho}{\delta-1} + \frac{1}{2} \rho |v|^2$$

Handwritten annotations:

- $\rho_t + \operatorname{div}(\rho v) = 0$, — way
- $(\rho v)_t + \operatorname{div}(\rho v \otimes v + p I) = 0$, —
- $E_t + \operatorname{div}((E + p)v) = 0$, —
- $u(x, 0) = (\rho, v, E)(x, 0) = a(x)$.

—

- Solution operator $\mathcal{G} : \mathbf{a} \rightarrow \mathcal{G}\mathbf{a} = u(T)$.



- Learn the evolution operator, input $u(T) \rightarrow \mathcal{G}u(T) = u(2T)$

在time-evolution PDE中，建立的还是figure到figure的一个映射，而不是figure到video的一个mapping。

Setup

这里考虑的一类问题是，当我们有很多的数据时，我们怎么样求解我们的PDE？这里给出的求解PDE的本质是find a solution operator。因此，我们的Setup如下：

- Input space
 Target space
 Function spaces
- 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 !!
- $\mathcal{L}^*: \mathbb{R}^{d_m} \xrightarrow{\quad} \mathbb{R}^{d_{out}}$
 $d_m, d_{out} \gg ?$
 $c < \infty$

但是，我们想用NN学习这个无穷维到无穷维的映射很困难，因此，我们可以考虑如下解决方案：

Solution I: Use Parametric PDEs instead

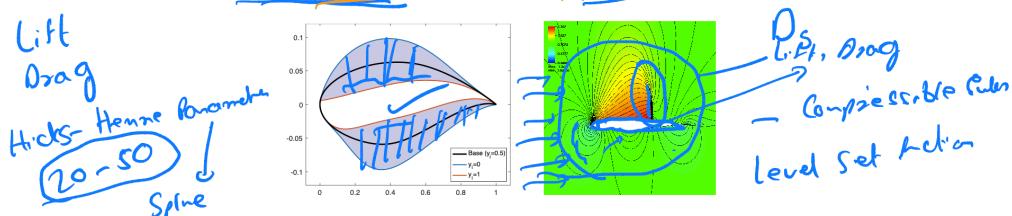
也就是将无穷维的function space参数化（基于一些basis），然后无穷维的function space reduced into a finite space.。接着，学习这些参数到solution的映射（有限维到有限维）。这种方法我们称为“cheat”：

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: $\dim(\text{Supp}(\mu)) = d_y < \infty$
- Corresponds to Parametrized PDEs with finite parameters.
- Find Soln $u(t, x, y)$ or Observable $\mathcal{L}(y)$ for $y \in Y \subset \mathbb{R}^{d_y}$.

$$\begin{aligned} a &\sim \mu \\ &\downarrow \\ a &= a(y) \\ &y \in Y \subset \mathbb{R}^{d_y} \end{aligned}$$

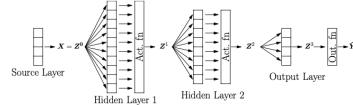


- Approximate Fields or observables with deep neural networks

- 这种方法可行的理论背景是函数的fourier basis展开。此外，这个思路和有限元很像，有限元也是Galerkin Discretization之后，将解reduced到一个finite space上去。
- 但是，即使这种parametric (reduced space)的方法也会存在问题，这面我们会介绍。

当我们使用传统的supervised learning的时候，我们通常有以下setup:

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 meas'n
/Numerical Simulation*

- For high dimensional input, we will often employ random training set.

但是在理论上，这会引起一些问题：

Supervised learning for high-d Parametric PDEs

- ▶ Can we find DNN such that $\|\mathcal{L}^* - \mathcal{L}\| \sim \mathcal{O}(\epsilon)$? *d* - Ground Truth? $d: \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ YES: **Universal Approximation Property** of DNNs $\Rightarrow: d: \mathbb{R}^d \rightarrow \mathbb{R}$
- ▶ Given any **Continuous** (measurable) \mathcal{L} , exists a $\hat{\mathcal{L}}$: *Smooth* \Rightarrow *easy to learn*
- Integrability* $\|\mathcal{L} - \hat{\mathcal{L}}\| < \epsilon$ (✓)
- ▶ If $\mathcal{L} \in W^{1,p}$, \exists DNN $\hat{\mathcal{L}}$ with M parameters (**Yarotsky**): *Smooth* \Rightarrow *easy to learn*
- Smoothness* $\|\mathcal{L} - \hat{\mathcal{L}}\|_p \sim \mathcal{O}(M^{-\frac{s}{d}})$
- ▶ But in Scientific Computing (often):
- ▶ \mathcal{L} is not be very **Regular** and $\bar{d} \gg 1$ 10^{-2}
- ▶ 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*

- 收敛速度由target function的smoothness和integrability来决定。首先就是你如果想控制你的error，你的神经网络得非常非常大。当你的输入维度比较多的时候，所需要的parameters will exponentially increase.

下面，我们考虑的是我们能否refined our error estimates，然后在一些情况下，误差的某些部分能够被控制，我们则专注于提升另一部分。这就是下面这一部分：

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 $\mathcal{E}_{app} \sim \mathcal{O}(d^\sigma M^{-\eta})$ for, —
 - Linear Elliptic PDEs: (Schwab, Kutyniok et al).
 - Semi-linear Parabolic PDEs: (E, Jentzen et al).
 - Nonlinear Hyperbolic PDEs: (DeRyck, SM, 2021). ↑
Lines of quadratic
- Optimization Error $\mathcal{E}_{opt} \sim$ Computable Training error:
- Generalization Error \sim Finite Sample Size Error

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

.. C_{quad} Quadratic Error

- Using Concentration inequalities + Covering number bounds:

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



- 我们将error归结为approximation error, generalization error和optimization error的和。
- 这个地方不是很明白三种error所表达的意思

好像我们的重点在generalizatoin error上面：

Well-trained Networks

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

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


- ▶ 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.

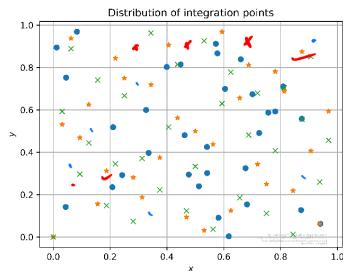
- 而generalization error我们也会遇到一个问题，就是分母是 \sqrt{N} 。意味着我们想要提升1%的error需要多 10^4 个数据。而在scientific computation中，获得一个数据是非常耗时的。
- 因此，有没有什么trick。这里老师提了一种方法，也就是我们知道我们要学习的是pde，而这些pde不是完全随机的，我们可以有一些先验知识在里面，基于这些先验知识，我们可以借助数值分析里面的一些手段，让我们更好地学习模型。
 - 这里采用的是在采样上做变化。这里提到一种方法，就是使用low discrepancy sequences来作为我们的traing dataset，我们可以理论上将我们的error随 N 的变化控制的更小：

—

A Trick: Lye, SM, Ray, 2020

Toads from Numerical Analysis.

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



$$\int f(y) dy \approx \frac{1}{m} \sum_{i=1}^m f(y_i) \rightarrow \text{LDS}$$

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

quasi-random number generator

Training on Low-Discrepancy Sequences

$$\int_D \left| \frac{\partial^d \mathcal{L}}{\partial y_1 \dots \partial y_d} \right| dy < \infty$$

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

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

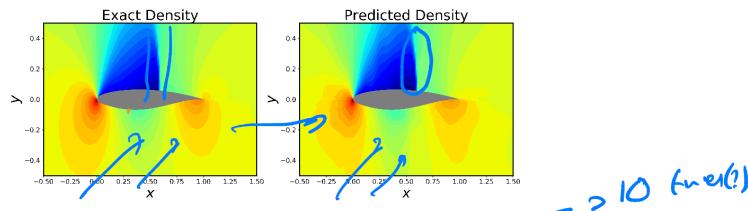
Strong error

Random: $\mathcal{E} \leq \mathcal{E}_r + \frac{C(\alpha)}{\sqrt{N}}$

- 通过下面这个例子我们也可以看到相比于random sampling的优势:

- 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 - 40	700 s	10^{-5} s	0.78%
Drag	2400 s	840 s	10^{-5} s	1.87%
Field	2400 s	1 hr	0.2 s	1.9%



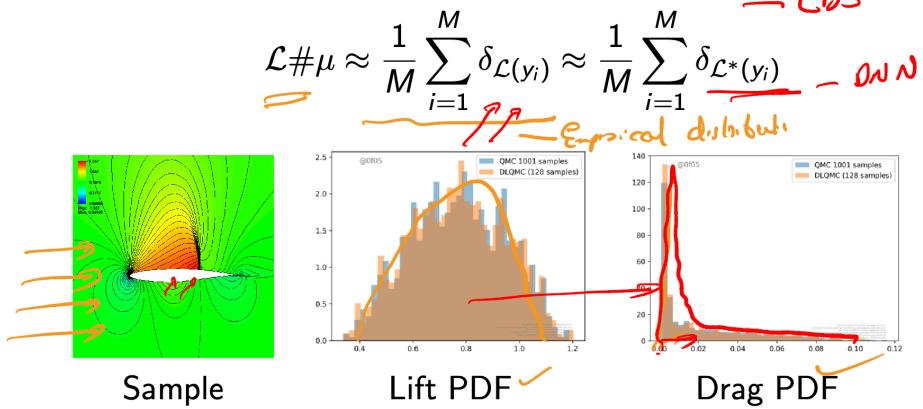
- Errors with Random Training pts: Lift 8.2%, Drag: 23.4%

然后，当我们训练好一个surrogate model之后，我们可以做downstream task，这里列举了两个：

- UQ:

Forward UQ – Uncertainty quantification

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



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

$O(10)$ -speedup

- Inverse design (见tutorial)

7.2 Operator Networks and DeepONet

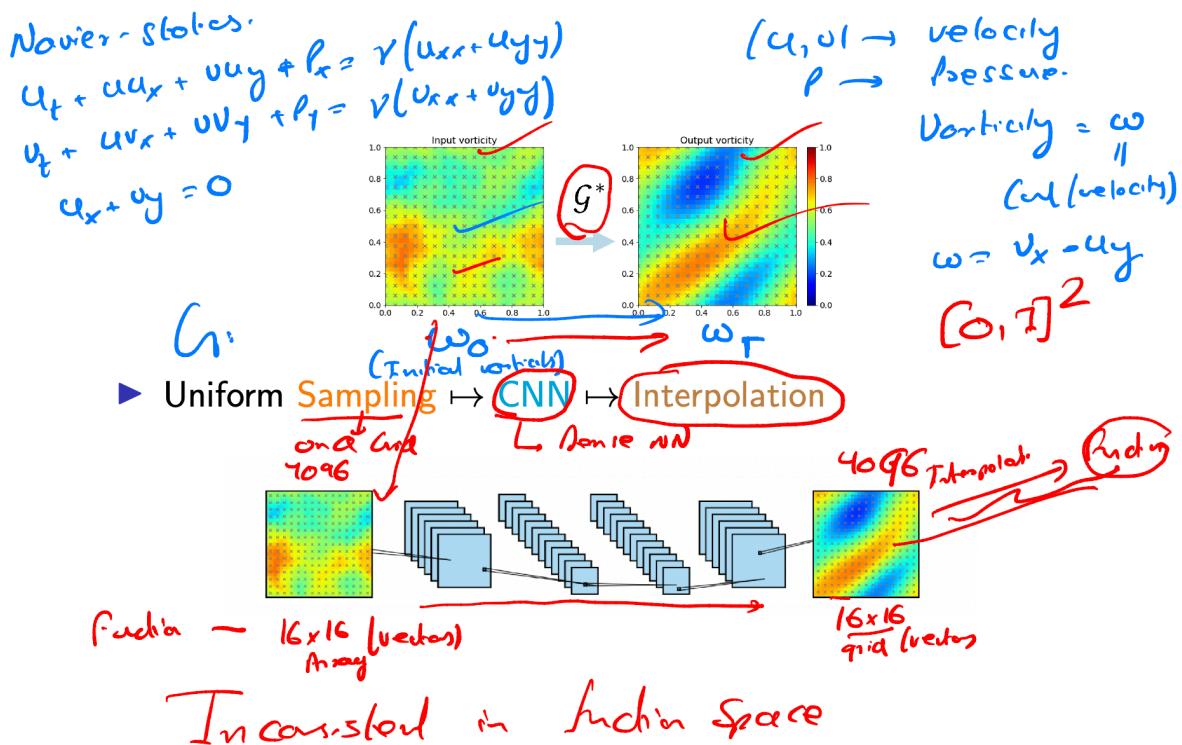
Issue with Parametrization

$\mathcal{G} : x \rightarrow y$ (PDE solution operator) $\dim(\text{supp}(\mu)) = d_y < \infty, \mu \in \text{Prob}(x)$

- Difficult to come up with a suitable one \rightarrow expert domain knowledge
- Not unique
- Does not Generalize well

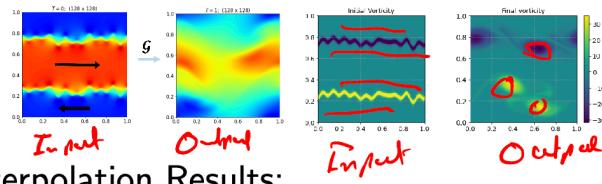
Solution I: Just use DNN + Interpolation

Uniform Sampling \rightarrow CNN \rightarrow Interpolation

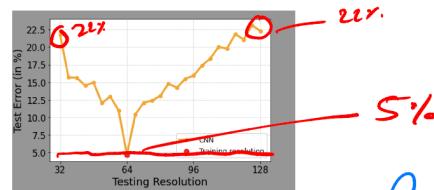


- ▶ Shear flow with Navier-Stokes with $Re \gg 1$ (ϵ_{uler})

64²



- ▶ CNN + Interpolation Results:



Resolution
invariance

- ▶ Consistent with Zhu, Zabaras, 2019.

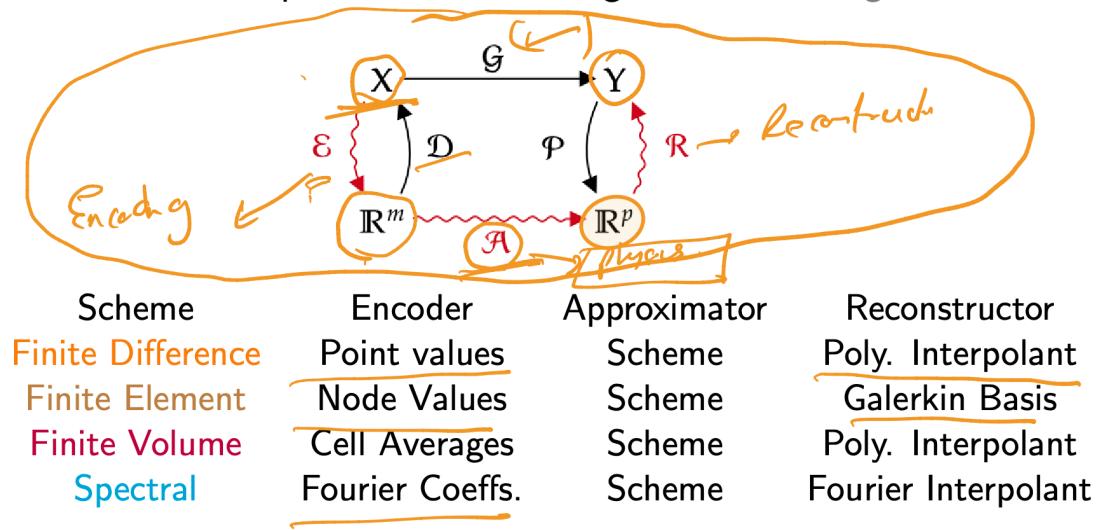
- ▶ Desiderata for Operator Learning:

- ▶ Input + Output are functions.
- ▶ Some notion of Continuous-Discrete Equivalence

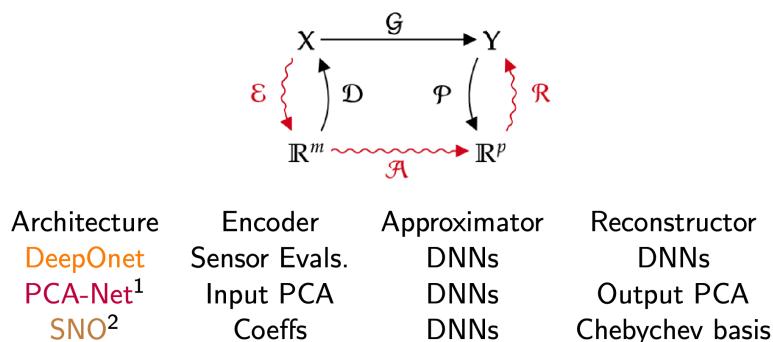
- This method didn't work, just learning the mapping from $64 \times 64 \rightarrow 16 \times 16$, not a operator learning, don't have the resolution invariance.
- Practice - discrete representation. Therefore, we have to satisfy the continuous-discrete equivalence.

Revisiting Numerical Methods

- ▶ Can be reinterpreted in the following abstract Paradigm:



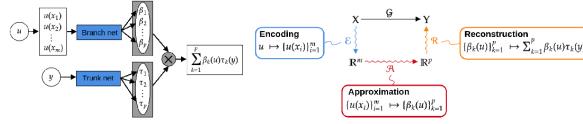
Other OpLearn Architectures



DeepONets

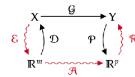
- DeepONets: Chen, Chen 1995, Lu et al, 2020:

$$\mathcal{N}(a)(y) = \tau_0(y) + \sum_{k=1} \beta_k(a) \tau_k(y) \approx \mathcal{G}(a)(y)$$



Why do ONet work?

ONets will face the issue of **curse of dimensionality**.



- **Universal Approximation Thm:** For $\mu \in \text{Prob}(L^2(D))$ and any measurable $\mathcal{G} : H^r \mapsto H^s$ and $\epsilon > 0$, $\exists \mathcal{N}$ (Onet): $\hat{\epsilon} < \epsilon$
- Upper (and lower) bounds via $\mathcal{E}_{\mathcal{R}} \leq \mathcal{E} \leq C(\mathcal{E}_{\mathcal{E}} + \mathcal{E}_{\mathcal{A}} + \mathcal{E}_{\mathcal{R}})$.
- $\mathcal{E}_{\mathcal{E}, \mathcal{R}}$ decay as spectrum of **Covariance operator** of μ and $\mathcal{G} \# \mu$.
- For $G = \mathcal{P} \circ \mathcal{G} \circ \mathcal{D} \in C^k(\mathbb{R}^m, \mathbb{R}^p) \Rightarrow \text{size } (\mathcal{N}) \sim \mathcal{O}\left(\epsilon^{-\frac{m(\epsilon)}{k}}\right)$.
- **Curse of Dimensionality** (CoD) for Onets !!!

- covariance operator可以简单理解为函数更容易通过几个domaint basis来进行展开

However, when the operator is not random and comes from the PDE, DeepONet breaks the **curse of dimensionality**.

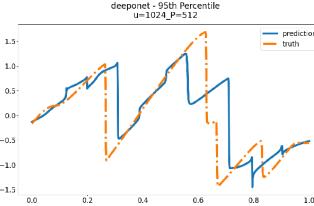
- For operators \mathcal{G} corresponding to many PDEs:

PDE	Operator	Complexity
$-u'' = \sin(u) + f$	$\mathcal{G} : f \rightarrow u(T)$	$M \sim \mathcal{O}(\epsilon^{-\eta}) \quad \eta \approx 0$
$-\operatorname{div}(a \nabla u) = f$	$\mathcal{G} : a \rightarrow u$	$M \sim \mathcal{O}(\epsilon^{-\eta}) \quad \eta \approx 0$
$u_t = \Delta u + f(u)$	$\mathcal{G} : u_0 \rightarrow u(T)$	$M \sim \mathcal{O}(\epsilon^{-2(d+1)})$
$u_t + u \cdot \nabla u + \nabla p = \nu \Delta u$	$\mathcal{G} : u_0 \rightarrow u(T)$	$M \sim \mathcal{O}(\epsilon^{-(d+1)})$
$u_t + \operatorname{div}(f(u)) = 0$	$\mathcal{G} : u_0 \rightarrow u(T)$	$M \sim \mathcal{O}(\epsilon^{-\alpha(d+1)})$

- Case by Case basis: **Operator Holomorphy**, Emulation of Numerical Schemes etc...
- Unified framework in DeRyck, SM, 2022.

Issues

- **Affine Reconstructors** $\Rightarrow \sqrt{\sum_{\ell>p} \lambda_\ell} \leq \mathcal{E}_R \leq \mathcal{E}$
- Large error for Slow decay of eigenvalues of **Covariance operator** of $\mathcal{G}\#_\mu$.
- Holds for **Transport dominated problems**
- Error of 30% for Burgers' equation with **GRF initial data !!**



- Upper bound 的问题是 issue。

↳ Therefore, alternative is neural operator

CNO

Sampling theorems: tell you that if your function has a certain range of frequencies then you need so much of resolution in physical space you need so many great points. So if it has 10 frequencies you need roughly 20 grid points. Anythings more than 23 points can resolve it if not then you will get a aliasing error.

7.3 Lecture 10

Three important points in using DeepOnet for approximation:

- Covariance Operator
- Curse of dimensionality
- The curse of dimensionality can be alleviated by leveraging the structures in PDE

Covariance Operator

The performance of the operator network depends on how fast or how slow the eigenvalues of this object(covariance operator) decay.

After knowing the issues of these Operator Networks, we want to find some alternatives, One possible solution is neural operators.

Because we need to calculate the integral, therefore the computational cost is super large. In order to reduce the complexity, we introduce the fourier operator.

