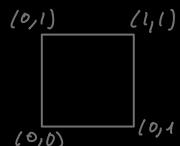


laplacean

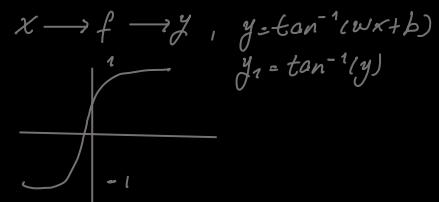
$$-\frac{\partial^2}{\partial x^2} - \frac{\partial^2}{\partial y^2} = f = -\Delta u$$



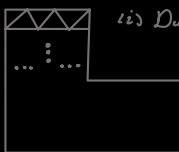
Physics-Informed Neural Networks (PINNs): e.g.  $m \frac{d^2 u}{dx^2} + \mu \frac{du}{dx} + k u = 0$

Forward Difference:  $D^+ f(x) = \frac{f(x+h) - f(x)}{h}$ ,  $x_{\text{norm}} = \frac{x - x_{\text{mean}}}{x_{\text{std}}}$

foreign System:  $\begin{cases} \dot{x} = \sigma(y - z) \\ \dot{y} = x(\rho - z) - y \\ \dot{z} = \beta z \end{cases}$ ,  $x(t=0)=1$ ,  $y(t=0)=1$ ,  $\frac{dq}{dt} = f(q)$ ,  $q(t+\Delta t) - q(t) = \frac{f(q)}{\Delta t}$ ,  $q(t+\Delta t) = q(t) + \Delta t f(q)$



Finite element method:  $\int -\Delta u \cdot v \, dx$  (test function)



iid Discretize domain  $u = \sum a_i \varphi_i(x)$ ;  $\Delta u = \Delta(\sum a_i \varphi_i) = 0$ ,  $\int \Delta u \cdot v \, dx = \int f \cdot v \, dx$ ;

$$m \frac{\partial^2 u}{\partial t^2} + d \frac{\partial u}{\partial t} - \nabla \cdot (c \nabla u) + au = f$$

$$x \in [-4, 4]$$

$$f(x) = x^2, y = x^2$$



$$y_{h_1} = w_1 x + b$$

$$(ii) y_{NL_1} = \phi(y_{h_1})$$

$$y_{h_2} = w_2 y_{NL_1} + b_2, y_{h_2} = \sum (y_{pred} - y)^2$$

Algorithm

Minimizer

Butti

Target

approx. error

\* truth

$\min_{w, b} J(y, \hat{y} | w, b)$  (convex),  $w_{n+1} = w_n - \eta \frac{\partial L}{\partial w_n}$  i.e.  $\theta_{n+1} = \theta_n + \eta \nabla L(\theta_n)$  (Gradient-based optimization)

NN for Regression:  $T^L(x) = W^L x + b^L$ , Activation function ( $\sigma$ ): e.g.  $\tanh(x)$ ,  $\max(0, x)$ , (iii) The  $L-1$  hidden layers of a feedforward N-N.

$N_{HL} = \sigma \circ T^{L-1} \circ \dots \circ T^1(x)$ , NN for Classification:  $f_{SM}(\xi_i) = \frac{\exp(\xi_i)}{\sum_j \exp(\xi_j)}$ ,  $0 \leq f_{SM}(\xi_i) \leq 1$ ,  $\sum_i f_{SM}(\xi_i) = 1$ . Convert any vector  $\xi_i$  to a probability vector  $N(x; \theta) = f_{SM} \circ T^L \circ N_{HL}(x)$

$x \rightarrow [ ] \rightarrow y, y = \varphi(wx + b)$ . Def: We say  $\sigma$  is discriminatory if for measure  $\mu \in \mathcal{M}(I_n)$   $\int_{I_n} \sigma(Wx + b) \, d\mu(x) = 0$

$\forall W^l \in \mathbb{R}^{n \times N}$  and  $b^l \in \mathbb{R}^N$  implies  $\mu = 0$ . (In compact on  $\mathbb{R}^N$ ),  $\mathcal{M}(I_n)$ : Space of finite regular Borel measures on  $I_n$

Def: We say  $\sigma$  is sigmoidal if  $\sigma(x) = \begin{cases} 1, & x \rightarrow +\infty \\ 0, & x \rightarrow -\infty \end{cases}$ . Theorem: Let  $\sigma$  be any continuous discriminatory function, then finite sums of the form  $y = \sum_{j=1}^N W_j^2 \sigma(W_j^l x + b_j)$  are dense in  $C(I_n)$ :  $\forall f \in C(I_n)$  and  $\epsilon > 0$ ,  $\exists y = \sum_{j=1}^N W_j^2 \sigma(W_j^l x + b_j)$  s.t.  $|y - f(x)| < \epsilon \forall x \in I_n$

Theorem: Suppose  $T$  is a compact set in  $C(a, b)$ ,  $f$  continuous functional defined on  $T$ ,  $\sigma(x)$

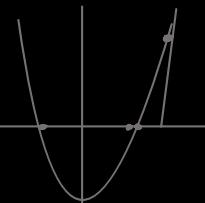
Adaptive Basis Viewpoint:  $N_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$  consisting of  $L-1$  hidden layers of width  $N$  composed with a final linear layer admitting the representation  $N_\theta(x) = \sum_{i=1}^N W_i^L \Phi_i(x; \theta_H)$ . Loss Function: (i) To learn  $u: I \rightarrow \mathbb{R}$ , given a dataset  $\{(x_i, u(x_i))\}_{i=1}^m$

MSE Loss:  $L(\theta) = \|u(x) - N(x, \theta)\|_2^2 \approx \frac{1}{m} \sum_{i=1}^m (u(x_i) - N(x_i, \theta))^2$ . In general let  $\{\tilde{F}_k\}_{k=1}^K$  be a linear/nonlinear operator

(ii)  $x \rightarrow Wx + b \rightarrow \sigma(Wx + b) = y_{NN}$ , (iii)  $J = \frac{1}{N} \sum_{i=1}^N (y_{NN} - y_i)^2 \Rightarrow \underset{[w, b]}{\operatorname{argmin}} J$ . (iii)  $W_{n+1} = W_n + \eta \frac{\partial J}{\partial W_n}$

Suppose  $f(x): y = x^2 - z \equiv 0$

(Use GPU's for Matrix Multiplication)  $b_{n+1} = b_n - \eta \frac{\partial J}{\partial b_n}$



$$f'(x) = f(x_n) + (x - x_n) f'(x) + \dots$$

$$\text{Newton's: } x_{n+1} = x_n - \frac{f(x_n)}{f'(x_n)}$$

$$\ell_p \text{ Norm: } \|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{1/p}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \xrightarrow{f_{INN}} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

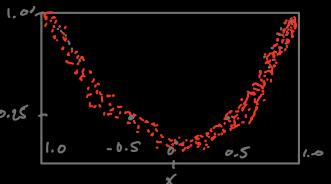
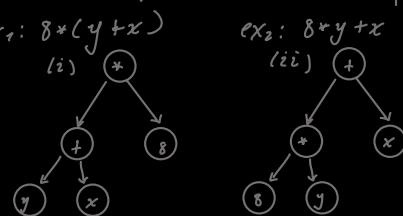
$$f: \mathbb{R}^{3 \times 1} \rightarrow \mathbb{R}^{3 \times 1}, \quad \partial f = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \frac{\partial y_1}{\partial x_3} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \frac{\partial y_2}{\partial x_3} \\ \frac{\partial y_3}{\partial x_1} & \frac{\partial y_3}{\partial x_2} & \frac{\partial y_3}{\partial x_3} \end{bmatrix}_{3 \times 3}$$

NN + Function Approximation:  $y = x^2 + E$ ,  $E \sim \mathcal{U}[0, 1]$ ,  $x \in [-1, 1]$ , Noise = 10%

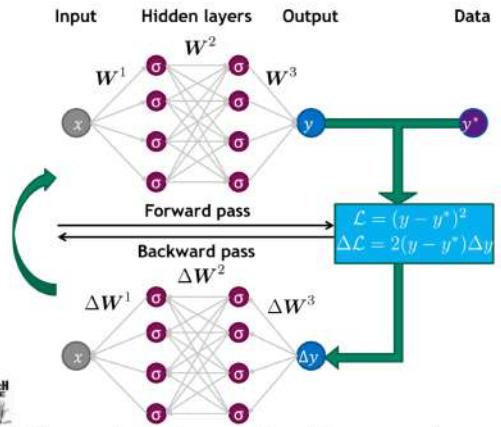
in-dim = 1, out-dim = 1

DL Networks:

$$\text{Finite Difference: } \frac{\partial f}{\partial x_i} \approx \frac{f(x + \Delta x_i) - f(x)}{\Delta x_i}$$



## Workflow in a Neural Network

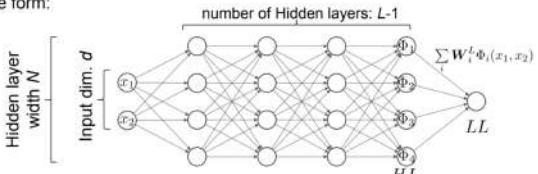


- Input layer (layer 0):  $z^0 = x \in \mathbb{R}^d$
- Hidden layers:
  - Layer 1:  $z^1 = \sigma(W^1 x + b^1) \in \mathbb{R}^{N_1}$
  - Layer 2:  $z^2 = \sigma(W^2 z^1 + b^2) \in \mathbb{R}^{N_2}$
- Output layer (layer 3):
  - $y = z^3 = W^3 z^2 + b^3 \in \mathbb{R}$

## A Neural Network for Regression

- Define the affine transformation in  $i$ -th layer  $T^i(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$
- Activation function  $\sigma$   
Popular choices:  $\tanh(x)$ ,  $\max\{x, 0\}$  (Rectified Linear Unit, ReLU)
- The  $L-1$  hidden layers of a feedforward neural network:  
 $N_{HL}(x) = \sigma \circ T^{L-1} \circ \dots \circ \sigma \circ T^1(x)$   
Where  $\circ$  denotes composition of functions
- For regression, a DNN is typically of the form:

$$N(x; \theta) = T^L \circ N_{HL}(x)$$



- Network parameters:  $\theta = \{W^l, b^l\}_{1 \leq l \leq L}$

## A Neural Network for Classification

- For classification, define the softmax function for  $K$  classes

- $f_{SM}(\xi_i) = \frac{\exp(\xi_i)}{\sum_{j=1}^K \exp(\xi_j)}$
- $0 \leq f_{SM}(\xi_i) \leq 1$
- $\sum_i f_{SM}(\xi_i) = 1$
- Convert any vector  $\xi$  to a probability vector

- The DNN is typically of the form

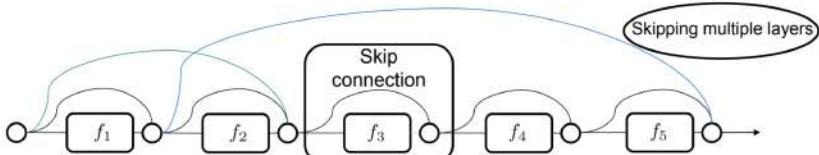
$$N(x; \theta) = f_{SM} \circ T^L \circ N_{HL}(x)$$

- Residual network (ResNet): skip/residual connections

- Replace  $\sigma \circ T^l$  with  $I + \sigma \circ T^l$

- $I$ : Identity function

$$N(x) = T^L \circ (I + \sigma \circ T^{L-1}) \circ \dots \circ (I + \sigma \circ T^2) \circ \sigma \circ T^1(x)$$



## Adaptive Basis Viewpoint

We consider a family of neural networks  $N_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$  consisting of  $L-1$  hidden layers of width  $N$  composed with a final linear layer, admitting the representation

$$N_\theta(x) = \sum_{i=1}^N W_i^L \Phi_i(x; \theta_H)$$

where  $W^L$  and  $\theta_H$  are the parameter corresponding to the final linear layer and the hidden layers respectively. We interpret  $\theta$  as a concatenation of  $W^L$  and  $\theta_H$ .

Gradient-Based Optimization:  $\theta_{k+1} = \theta_k - \eta \nabla L(\theta_k)$

Universal Function Approximation (1-layer):

Dif: We say  $\sigma$  discriminatory if for measure  $\mu \in M(I_N)$

$\exists \lambda: \text{finite-signed-regular Borel measure on } I_N \subset \mathbb{R}^N$

$\int_I \sigma(W^1 x + b^1) d\mu(x) = 0 \quad \forall W^1 \in \mathbb{R}^{N \times N}, b^1 \in \mathbb{R}^N \Rightarrow \mu = 0. (I_N \text{ compact})$

Dif:  $\sigma$  sigmoidal if  $\sigma(x) \rightarrow \begin{cases} 1, & x \rightarrow +\infty \\ 0, & x \rightarrow -\infty \end{cases}$ . Theorem: If  $\sigma$  continuous discriminatory, then  $y = \sum_{j=1}^N W_j^2 \sigma(W_j^1 x + b_j)$  are dense in  $C(I_N)$

$\therefore \forall f \in C(I_N), \epsilon > 0 \quad \exists y = \dots \text{ s.t. } |y - f(x)| < \epsilon \quad \forall x \in I_N$

Theorem: If  $T$  compact in  $C[a, b]$ ,  $f$  continuous functional on  $T$ ,

$\sigma(x)$  bounded sigmoidal function:  $\forall \epsilon > 0, \exists m+1$  points  $a = x_0 < \dots < x_m = b$

$\forall N \in \mathbb{N}$ , constants  $W_i^2, b_i, W_{i,j}^1, j = 1, 2, \dots, N, i = 1, 2, \dots, m$  s.t.

$|f(x) - \sum_{i=1}^N W_i^2 \sigma(\sum_{j=0}^m W_{i,j}^1 u(x_j) + b_i)| < \epsilon \quad \forall u \in T$ . Loss Functions:

To learn  $u: \Omega \rightarrow \mathbb{R}$ , Given dataset  $\{(x_i, u(x_i))\}_{i=1}^m$ .

MSE:  $L(\theta) = \|u(x) - N(x; \theta)\|_2^2 \approx \frac{1}{m} \sum_{i=1}^m (u(x_i) - N(x_i; \theta))^2$

Let  $\{F_k\}_{k=1}^K$  be linear/nonlinear operator  $L(\theta) = \sum_{k=1}^K \lambda_k \|F_k[u] - F_k[N]\|_2^2$

Classical Theory of Universal Approximation Theorem:

Dif: (Tauber-Wiener function) If a function  $\sigma: \mathbb{R}^m \rightarrow \mathbb{R}$  (continuous or discontinuous) satisfies that all the linear combinations  $\sum_{i=1}^N c_i \sigma(W_i^T x + b_i)$ ,  $W_i \in \mathbb{R}^{m \times m}$ ,  $b_i \in \mathbb{R}^m$ ,  $c_i \in \mathbb{R}$  are dense in  $C(K)$  or  $L^p(K)$ ,  $K \subset \mathbb{R}^m$  compact. Then  $\sigma$  is called a Tauber-Wiener (TW) function on  $C(K)$  or  $L^p(K)$ .

Theorem: Let  $K$  compact in  $\mathbb{R}^m$ . Then the Tauber-Wiener function on  $C(K)$  or  $L^p(K)$  exists. Since  $\sigma \in TW$ .

Approximation Theory (discontinuous functions)

Let  $\sigma \in L^p[a, b]$   $\sigma \neq$  a polynomial a.e. Let  $K \subset \mathbb{R}^d$  (compact). Then  $\forall f \in L^p(K) \quad \forall \epsilon > 0 \quad \exists$  real  $c_i$ 's,  $W_i$ 's,  $b_i$ 's  $1 \leq i \leq N$  s.t.  $\left\| f - \sum_{i=1}^N c_i \sigma(W_i^T x + b_i) \right\|_{L^p(K)} < \epsilon$ ,  $f$  (dis)continuous  $\Rightarrow \sigma$  a Tauber-Wiener function.

Maximum Norm: If  $\mu$ -finite on  $(\mathbb{R}^d, \mathcal{B}(\mathbb{R}^d))$   $\sigma$  sigmoidal,  $\forall$  Borel measurable  $f: \mathbb{R}^d \rightarrow \mathbb{R}$   $\epsilon > 0, \exists c_i, W_i, b_i, 1 \leq i \leq N, K$  (compact)  $\subset \mathbb{R}^d$  with  $\mu(K) + \epsilon > \mu(\mathbb{R}^d)$  s.t.  $\sup_{x \in K} |f(x) - \sum_{i=1}^N c_i \sigma(W_i^T x + b_i)| < \epsilon$

Multilayer Neural Networks:

Kolmogorov Superposition Theorem:  $\forall$  continuous  $f$  of  $n$ -variables on  $[0, 1]^n \quad \exists \lambda_1, \dots, \lambda_n, \sum_{j=1}^n \lambda_j \leq 1$  and  $2n+1$  strictly continuous functions  $\phi_1, \dots, \phi_{2n+1}$   $\phi_i: [0, 1] \rightarrow [0, 1]$  s.t.  $f(x_1, \dots, x_n) = \sum_{i=1}^{2n+1} g_i \left( \sum_{j=1}^n \lambda_j \phi_i(x_j) \right)$ ,  $g_i \in C[0, 1]$ .

Maiorov-Pinkus Theorem:  $\forall f \in C[0, 1]^n, \epsilon > 0 \quad \exists d_i, c_{ij}, b_{ij}, r_i, \forall i \in \mathbb{N}$   $\forall j \in \mathbb{N}$   $\sigma$  activation  $f_j \in C^\infty$  increasing-sigmoidal s.t.  $\max_{x \in [0, 1]^n} \left| f(x) - \sum_{i=1}^{d_i} \sum_{j=1}^{r_i} c_{ij} \sigma(\tilde{w}_{ij}^i x + b_{ij}) + \eta_i \right| < \epsilon$

# Regression of a Discontinuous/Oscillatory Function in Physical & Fourier Domains

## 0. Data preparation

```
% Code for Function regression
% Data Generation
% Define f(x) using anonymous definition
f = @(x) (x < 0).* (5.0 + sin(x)) + sin(2.0*x) + ...
sin(3.0*x) + sin(4.0*x)) + (x>=0).*.cos(10*x);

% Training and Testing data
x_in_l = linspace(-pi, -1.0e-3, 201);
x_in_r = linspace(0, pi, 501);
y_in_l = f(x_in_l);
y_in_r = f(x_in_r);
x_in = [x_in_l, x_in_r];
y_in = [y_in_l, y_in_r];
plot(x_in, y_in)
```

## 1. Initialize Neural Network Parameters

```
% Parameters Initialization
parameters = struct;
numLayers = 2;
numNeurons = 48;
sz = [numNeurons 1];
parameters.fcl.Weights = initializeHe(sz,1,"single");
parameters.fcl.Bias = initializeZeros([numNeurons 1],"single");
for layerNumber=2:numLayers-1
    name = "fc"	layerNumber;
    sz = [numNeurons numNeurons];
    numIn = numNeurons;
    parameters.(name).Weights = initializeHe(sz,numIn,"single");
    parameters.(name).Bias = initializeZeros([numNeurons 1],"single");
end
sz = [1 numNeurons];
numIn = numNeurons;
parameters.(["fc" + numLayers].Weights = initializeHe(sz,numIn, "single");
parameters.(["fc" + numLayers].Bias = initializeZeros([1 1], "single");
```

## 2. Data preparation, mini-batches, and GPU placement

```
sz_in = size(X);
numEpochs = 60000;
miniBatchSize = sz_in(2);
executionEnvironment = "auto";
learningRate = 0.001;
ds = arrayDatstore(X);
ds = arrayDatstore(ds, ...
    MiniBatchSize=miniBatchSize, ...
    MiniBatchFormat="CB", ...
    OutputEnvironment=executionEnvironment);
Y = dlarray(Y, "CB");

if (executionEnvironment == "auto" && canUseGPU) || ...
    (executionEnvironment == "gpu")
    X = gpuArray(X);
    Y = gpuArray(Y);
end
```



## 3. Training Loop

```
averageGrad = {};
averageSqrGrad = {};
lossHist = zeros(numEpochs,1);
accfun = diaccelerate(@modelLoss);
start = tic;
iteration = 0;
for epoch = 1:numEpochs
    reset(mbq);
    while hasdata(mbq)
        iteration = iteration + 1;
        X = next(mbq);
        [loss,gradients] = difeval(accfun, parameters,X,Y);
        [parameters,averageGrad,averageSqrGrad] = adamupdate(parameters, ...
            gradients,averageGrad, ...
            averageSqrGrad,iteration,learningRate);
    end
    loss = double(gather(extractdata(loss)));
    lossHist(epoch) = loss;
    D = duration(0,B,toc,start).Format="hh:mm:ss";
    fprintf('On Epoch: %d, and Loss %f\n', epoch, loss)
end
toc;
```

EP FLAND INSTITUTE | BROWN

Approximation Rates:  $B^n \triangleq \{x : \|x\|_2 = (\chi_1^2 + \dots + \chi_n^2)^{1/2} \leq 1\}$ , Sobolev Space  $W_p^m = W_p^m(B^n)$  as the completion of  $C^m(B^n)$   
w.r.t.  $\|f\|_{W_p^m}$   $= \begin{cases} \left( \sum_{0 \leq k \leq m} \|D^k f\|_p^p \right)^{1/p}, & 1 \leq p < \infty \\ \max_{0 \leq k \leq m} \|D^k f\|_\infty, & p = \infty \end{cases}$

Theorem (Lower/upper bounds):  $n \geq 2, m \geq 1$  For  $M_N = \{g = \sum_{i=1}^N c_i \sigma(w_i^T x + b_i)\}$ :  $\sup_{p \in W_p^m} \|f\|_{W_p^m} \leq \inf_{g \in M_N} \|f - g\| \sim N^{-\frac{m}{(n-1)}}$  (Shallow Networks)

Theorem (Deep Networks): Given  $f \in C([0,1]^n)$ ,  $\forall d, N \in \mathbb{N}$   $p \in [1, \infty]$   $\exists$  ReLU forward Network  $f^N$  with width  $C_1: \max\{N^{1/d}, N+2\}$   
depth:  $11L + C_2$  s.t.  $\|f - f^N\|_{L^p([0,1]^n)} \leq 131\sqrt{n} \omega(f, (N^2 L^2 \log_3(N+2))^{-1/n})$ ,  $C_1 = 16$ ,  $C_2 = 18$  if  $p < \infty$ ;  $C_1 = 3^{n+3}$ ,  $C_2 = 18 + 2n$  if

$p = \infty$  and  $\omega(f, \cdot)$  modulus of continuity  $\omega(f, \delta) = \sup_{|z| \leq \delta} \|f(\cdot + z) - f(\cdot)\|_{L^\infty}$ ,  $x, y \in [0,1]^n$  when  $|f(x) - f(y)| \leq \lambda|x-y|^\alpha$ ,  $0 < \alpha \leq 1$   
 $\|f - f^N\|_{L^p([0,1]^n)} \leq 131\sqrt{n} \lambda (N^2 L^2 \log_3(N+2))^{-\alpha/n}$

Theorem: Given any positive numbers and a Hölder continuous function  $f$  on  $[0,1]^n$ ,  $|f(x) - f(y)| \leq \lambda|x-y|^\alpha$ ,  $\exists$   
Floor ReLU-network  $f^N$  width( $f^N$ )  $\leq \max\{n, 5N+13\}$  depth  $64nL+3$   $\max_{x \in [0,1]^n} |f - f^N| \leq 3\lambda n^{\alpha/2} N^{-\alpha/\sqrt{n}}$

Hierarchical Binary Tree: e.g.  $f(x_1, \dots, x_8) = h_3(h_{11}(h_{11}(x_1, x_2), h_{12}(x_3, x_4)), h_{22}(h_{13}(x_5, x_6), h_{14}(x_7, x_8)))$

Theorem: Let  $\sigma$ -infinitely differentiable-univariate which is not a polynomial on any subinterval of  $\mathbb{R}$ .  $\therefore$  for  $f \in \mathcal{N}_p^m([-1, 1]^n)$

$f$ : compositional function conforming to a binary tree  $\Rightarrow$   $\exists$  hierarchical binary tree network  $f \in \mathcal{D}_N = \{\text{deep networks}\}$  s.t.

$\max_{x \in [-1, 1]^n} |f - f^N| \leq CN^{-m/2}$  (No curse of dimensionality), For shallow networks with  $\sigma$  activation:  $\leq N^{-m/n}$

Activation Functions: (i)  $\sigma(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$   $\sigma'(x) = (-\sigma^2(x))$ , (ii) Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$ ,  $\sigma'(x) = \sigma(x)(1-\sigma(x))$

(iii) ReLU:  $\sigma(x) = \begin{cases} 0, & x \leq 0 \\ x, & x > 0 \end{cases}$   $\sigma'(x) = \begin{cases} 0, & x \leq 0 \\ 1, & x > 0 \end{cases}$ , (iv) Leaky ReLU:  $\sigma(x) = \begin{cases} 0.01x, & x \leq 0 \\ x, & x > 0 \end{cases}$

(v) ELU:  $\sigma(x) = \begin{cases} d(e^{x-1}), & x \leq 0 \\ x, & x > 0 \end{cases}$  (vi) Swish:  $\sigma(x) = \frac{x}{1+e^{-x}} \sigma(x)$

$\sigma'(x) = \begin{cases} d(e^{x-1}) / (1+e^{-x}), & x \leq 0 \\ 1, & x > 0 \end{cases}$

Loss Functions: Entropy:  $H(p) = -\sum_i p_i \log_2(p_i)$ , Cross-Entropy:  $H(p, q) = -\sum_i p_i \log_2(q_i)$ , KL-Divergence:  $D(p||q) = H(p, q) - H(p)$

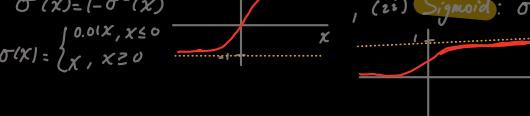
KL-Distribution-Convergence:  $\forall$  distributions  $p(x), q(x)$  over  $X$ :  $D[p(x)||q(x)] = \sum_x p(x) \log \frac{p(x)}{q(x)} \geq 0$

Dual and Jacobian: Dual numbers  $x = v + \tilde{v}\epsilon$   $v, \tilde{v} \in \mathbb{R}$ ,  $\epsilon^2 = 0$ ,  $f(x) = f(v + \tilde{v}\epsilon) = f(v) + f'(v)\tilde{v}\epsilon$   $\tilde{v} = 1$  computes  $f'(x)|_{x=v}$

Chain-Rule:  $f(g(v + \tilde{v}\epsilon)) = f(g(v) + g'(v)\tilde{v}\epsilon) = f(g(v)) + f'(g(v))g'(v)\tilde{v}\epsilon$ ,  $\tilde{v} = 1 \Rightarrow (f(g(x)))'|_{x=v}$

Jacobian:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ ,  $\partial f(x) = \mathbb{R}^n \rightarrow \mathbb{R}^{m \times n}$

$$\partial f(x) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(x) & \frac{\partial f_1}{\partial x_2}(x) & \dots & \frac{\partial f_1}{\partial x_n}(x) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(x) & \frac{\partial f_m}{\partial x_2}(x) & \dots & \frac{\partial f_m}{\partial x_n}(x) \end{bmatrix}_{m \times n} \in \mathbb{R}^{m \times n}$$



JVP: directional derivative of  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  along  $v \in \mathbb{R}^n$  is  $\nabla f(x)^T v = \partial f(x)/\partial v$

$f(x) = \frac{1}{2} \|x\|_2^2 \Rightarrow \nabla f(x)^T v = x^T v$

VJP:  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  the Jacobian  $|_X$  is  $\partial f(x) \in \mathbb{R}^{1 \times n}$ ,  $\nabla f(x) = \partial f(x)^T I$ ,  $f(x) = \frac{1}{2} \|x\|_2^2$   
 $\Rightarrow \nabla f(x) = x \cdot 1$

## Automatic Differentiation

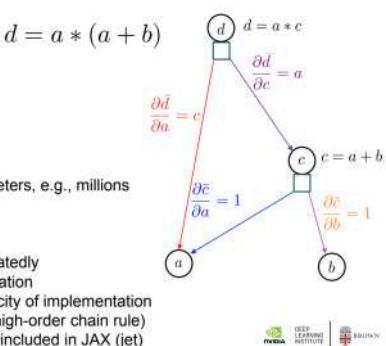
- Exploits the fact that all computations are compositions of a small set of elementary expressions with known derivatives
- Employs the chain rule to combine these elementary derivatives of the constituent expressions.

- Two ways to compute first-order derivative:

- Forward mode AD (details not discussed)
  - Cost scales linearly w.r.t. the input dimension
  - Cost is constant w.r.t. the output dimension
- Reverse mode AD
  - Cost is constant w.r.t. the input dimension
  - Cost scales linearly w.r.t. the output dimension

- In deep learning, backpropagation == Reverse mode AD
  - The input dimension of the loss function is # of parameters, e.g., millions
  - The output dimension is 1: the loss value

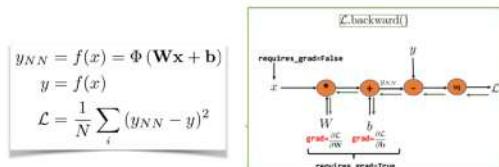
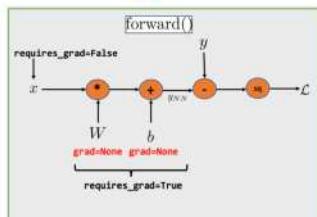
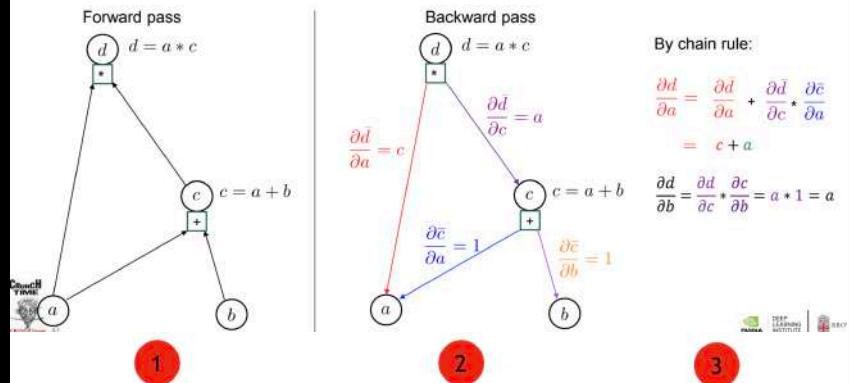
- High-order derivatives:
  - Nested-derivative approach: Apply first-order AD repeatedly
    - Cost scales exponentially in the order of differentiation
    - What we will use in this class, because the simplicity of implementation
  - More efficient approaches, such as Taylor-mode AD (high-order chain rule)
  - Not supported in TensorFlow/PyTorch yet but it is included in JAX (jet)



## Backpropagation

- We apply recursively the chain rule to implement Backprop
- Use computational graphs to accomplish backprop

- Example:  $d = a * (a + b)$



- Neural Networks can approximate functions and functionals with arbitrary accuracy.
- Deep Neural Networks are more expressive and can beat the curse-of-dimensionality.
- Activation functions are important for accuracy and convergence speed.
- Deep Neural Networks can be thought as nonlinear approximations with adaptive basis functions.
- Regression and classification use different loss functions – loss functions can be meta-learned.
- Forward/Backpropagation and Automatic Differentiation have optimal cost and machine precision.
- Generalization and Optimization errors are often greater than Approximation errors.
- Traditional methods, like finite elements, can be related to adaptive Deep Neural Networks.

$$x_{n+1} = x_n - \alpha \nabla f(x_n)$$

$$x_{n+1} = x_n - \alpha H_n^{-1} \nabla f(x_n)$$

### First Order

#### Momentum

Nestrov

Adagrad

RMSProp

Adam

Nadam

Rprop

Lion

### Second Order

#### Newton Method

#### Quasi-Newton Methods: BFGS

#### Quasi-Newton Methods: L-BFGS

Optimization: Stationary Points:  $\nabla L(\theta) = 0$  w.r.t.  $\theta$

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta), \quad \theta_{n+1} = \theta_n - \eta \nabla_{\theta} L(\theta)$$

$$L(\theta - \eta \nabla L(\theta)) = L(\theta) - \eta \|\nabla L(\theta)\|_2^2 + O(\eta^2)$$

Saddle Points: (i) Degenerate: Hessian positive semidefinite and 0-eigenvalue, Hessian: Suppose  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ . If all second order partial derivatives of  $f$  exist.  $\Rightarrow H_{nn}$ :

$$H_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad \forall \vec{x} \neq 0 \in \mathbb{R}^n, H_f \text{ positive def.} \iff \vec{x}^T H_f \vec{x} \geq 0$$

Batch Normalization:

## Algorithm:

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1, \dots, x_m\}$

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Learning rate Scheduling: (i) Train model for a few hundred iterations, exponentially increasing  $\eta(t)$

Dynamic Scheduling:  $\begin{cases} \eta(t) = \eta_0, & t_i \leq t \leq t_{i+1} : \text{Piecewise Decay} \\ \eta(t) = \eta_0 \exp(-\lambda t), & : \text{Exponential decay} \\ \eta(t) = \eta_0 (\beta t + 1)^{-\alpha}, & : \text{Polynomial decay} \end{cases}$

Theoretical Analysis:  $N^L$  will die in probability as  $L \rightarrow \infty$

Theorem: Let  $N^L$  be a ReLU-NN with  $L$  layers each having  $N_1, \dots, N_L$  neurons. Suppose weights are initialized from a symmetric distribution around zero (ii) Biases are 0 or from a symmetric distribution  $\therefore P(N^L(x) \text{ dies}) \leq \prod_{l=1}^L (1 - (1/2)^{N_l})$ . Furthermore, assuming  $N_l = N \ \forall l : \lim_{L \rightarrow \infty} P(N^L \text{ dies}) = 1, \lim_{N \rightarrow \infty} P(N^L \text{ dies}) = 0$

Proof: Lemma: Let  $N^L$  be a ReLU-NN with  $L$  layers. Suppose weights ind. from distribution s.t.  $P(W_j | z=0) = 0 \ \forall z \in \mathbb{R}^{N_L}$  and any  $j$ th row of  $W_L$ . Then  $P(N^L(x) \text{ dies}) = P(\exists l \in \{1, \dots, L-1\} \text{ s.t. } \phi(N^l(x)) = 0 \ \forall x \in D)$ . For a given  $x$

$$P(W_s^j \phi(N^{j-1}(x)) + b_s^j < 0 | \tilde{A}_{j-1}^c, x) = 1/2$$

$$\tilde{A}_{L,x}^c = \{ \forall 1 \leq j \leq L, \phi(N^j(x)) = 0 \}$$

## Loss Regularizers $L^2$

- Regularization strategy is used to reduce test errors for new inputs but may increase the training errors.
- Loss function  $\mathcal{L}$  with  $L^2$  regularizer is expressed as

$$\mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}) \leftarrow \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- After taking the gradient

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}) \leftarrow \alpha \mathbf{w} + \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

- Single gradient step to update the weights is expressed as

$$\mathbf{w} \leftarrow \mathbf{w} - \eta (\alpha \mathbf{w} + \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}))$$

- After rearranging the terms in the above expression

$$\mathbf{w} \leftarrow (1 - \eta \alpha) \mathbf{w} - \eta \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

## Loss Regularizers: $L^1$

- The  $L^1$  regularization model parameter  $\mathbf{w}$  is defined as

$$\Omega(\boldsymbol{\theta}) = \|\mathbf{w}\|_1$$

- Thus, the regularized loss function  $\boldsymbol{\theta}$  is expressed as

$$\mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}) + \alpha \|\mathbf{w}\|_1$$

- The gradient of  $L^1$  regularized loss function is expressed as

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \text{sign}(\mathbf{w}) + \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{X}, \mathbf{y}; \mathbf{w})$$

where  $\text{sign}(\mathbf{w})$  is sign  $\mathbf{w}$  applied element-wise.

## Summary

- The main difficulty in training a NN is due to back propagation and the need to minimize a high-dimensional non-convex function.
- Stochastic gradient descent (SGD) and mini-batch GD with decaying learning rate are effective and practical methods.
- Underfitting occurs for a low-capacity model while overfitting occurs for a high-capacity model.
- The Xavier and He initialization as well as the data normalization are key components of effective NN training.
- Batch normalization can lead to faster convergence, avoid vanishing gradients, allows higher learning rates and act as regularizer.
- There are many optimizers but the combination of Adams first, followed by L-BFGS is usually the winner.
- Decaying learning rate and cyclical learning rate are important in accelerating convergence and avoiding bad minima.
- L2 regularization leads to better generalization while L1 encourages sparsity; which one to use is problem-dependent.
- Dropout regularization is effective but depends strongly on the dropout rate, which should be variable across the layers.
- Information bottleneck theory suggests a phase transition in deep learning, from fitting to early layers to true learning in the higher layers.
- Deep neural networks are expressive but they can collapse during training yielding erroneous results, e.g. a deep and narrow ReLU network.

## Structure of deep convolutional neural networks

- Rectified linear unit (ReLU):  $\sigma(u) = (u)_+ = \max\{u, 0\}$
- sequence of convolutional filter masks  $\mathbf{w} = \{w^{(j)} = (w_k^{(j)})_{k=-\infty}^{\infty}\}_j$
- filter length s:** Assume  $w^{(j)}$  is supported in  $\{0, 1, \dots, s\}$
- convolution**  $w * x = (\sum_k w_k x_{i-k} = \sum_{k=0}^s w_{i-k} x_k)_i$  supported in  $\{0, 1, \dots, N+s\}$ , so  $(w * x)_{i=0}^{N+s} = W^w x$  with  $W^w = (w_{i-k})_{0 \leq i \leq N+s, 0 \leq k \leq N}$

$$W^w = \begin{bmatrix} w_0 & 0 & 0 & 0 & \dots & \dots & \dots & 0 \\ w_1 & w_0 & 0 & 0 & \dots & \dots & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ w_s & w_{s-1} & \dots & w_0 & 0 & \dots & \dots & 0 \\ 0 & w_s & w_1 & w_0 & 0 & \dots & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & w_s & \dots & w_1 & w_0 & \\ 0 & \dots & 0 & 0 & w_s & \dots & w_1 & \\ \vdots & \ddots & \ddots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & 0 & w_s & \end{bmatrix} \in \mathbb{R}^{(N+s) \times N}$$

## Convolutional Neural Network (CNN)

- Input: A 3D tensor (width, height, depth); e.g., if the input is an image, the depth is 3, i.e., Red, Green, Blue channels.
- Output: a 3D tensor (width, height, depth)
- A convolutional layer: Convolution between two functions  $f$  and  $g$  is defined as

$$(f * g)(\mathbf{x}) = \int f(\mathbf{z})g(\mathbf{x} - \mathbf{z})d\mathbf{z}$$

- For 1D discrete case, the integral turns into a sum

$$(f * g)(i) = \sum_a f(a)g(i-a)$$

- For 2D case

$$(f * g)(i, j) = \sum_a \sum_b f(a, b)g(i-a, j-b)$$

- Rewrite above as

$$[\mathbf{H}]_{i,j} = u + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a, j+b}$$

- For multiple channels:

$$[\mathbf{H}]_{i,j,d} = \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} \sum_v [\mathbf{V}]_{a,b,v} [\mathbf{X}]_{i+a, j+b, v}$$



## FNN vs CNN

<b>FNN</b> $T^l(x) = \mathbf{W}^l x + \mathbf{b}^l$	<b>CNN</b> $[\mathbf{H}]_{i,j} = u + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [\mathbf{V}]_{a,b} [\mathbf{X}]_{i+a, j+b}$
$[\mathbf{H}]_{i,j} = [\mathbf{U}]_{i,j} + \sum_k \sum_l [\mathbf{W}]_{i,j,k,l} [\mathbf{X}]_{k,l}$	$= [\mathbf{U}]_{i,j} + \sum_a \sum_b [\mathbf{V}]_{i,j,a,b} [\mathbf{X}]_{i+a, j+b}$

## Properties of CNN

- Efficiency: Much fewer parameters  $(W \times H)^2 \rightarrow K^2$
- Locality:
  - Nearby pixels are typically related to each other
  - We should not have to look very far away from location  $(i, j)$  to compute  $H(i, j)$
- Translation Invariance: A shift in the input  $\mathbf{X}$  should simply lead to a shift in the hidden representation  $\mathbf{H}$
- Two probability distributions  $P$  and  $Q$ .
- Total variation (TV) distance

$$\delta(P, Q) = \sup |P(A) - Q(A)|$$

- Kullback-Leibler (KL) divergence (also called relative entropy)

- Discrete probability distributions

$$D_{KL}(P||Q) = \sum_{x \in \mathcal{X}} P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

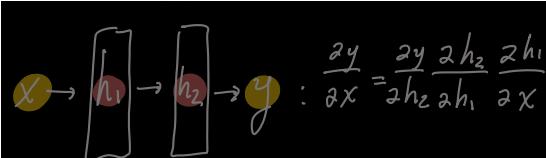
- Continuous random variable

$$D_{KL}(P||Q) = \int_{-\infty}^{\infty} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx$$

- $H(P, Q)$ : cross entropy of  $P$  and  $Q$

$$D_{KL}(P||Q) = H(P, Q) - H(p), \quad \text{Asymmetric!}$$

where  $H(P)$  denotes the entropy of  $P$ .



$$Y \rightarrow \boxed{h_1} \rightarrow \boxed{h_2} \rightarrow Y : \frac{\partial Y}{\partial X} = \frac{\partial Y}{\partial h_2} \frac{\partial h_2}{\partial h_1} \frac{\partial h_1}{\partial X}$$

A deep CNN of depth  $L$  is an  $L$ -layer neural network with widths  $\{N_j = N + js\}$  having a convolutional sparse structure:

$$z^{(j)}(x) = \sigma(W^{(j)})z^{(j-1)}(x) - b^{(j)}, \quad j = 1, 2, \dots, L$$

where  $W^{(j)} = W^{w^{(j)}}$  is a Toeplitz type  $N_j \times N_{j-1}$  matrix generated by the filter mask  $w = w^{(j)}$

The generated hypothesis space of functions:

$$\mathcal{H}_L^{\mathbf{w}, \mathbf{b}} = \left\{ \sum_{k=1}^{N+Ls} c_k z_k^{(L)}(x) : c = (c_i)_{i=1}^{N+Ls} \in \mathbb{R} \right\}$$

Take  $b^{(j)}$  of the form  $b = [b_1 \dots b_{s-1} b_s b_s \dots b_s b_{N_j-s+1} \dots b_{N_j}]^T$

Bias vector sequence  $\mathbf{b} = \{b^{(j)}\}_{j=1}^L$

Total number of free parameters:  $N = 3s(L-1) + 2N + 2L$

## Theorem

Let  $2 \leq s \leq N$ . For any compact subset  $\Omega$  of  $\mathbb{R}^N$  and any  $f \in C(\Omega)$ , there exist sequences  $\mathbf{w}$  of filter masks,  $\mathbf{b}$  of bias vectors, and  $f_L^{\mathbf{w}, \mathbf{b}} \in \mathcal{H}_L^{\mathbf{w}, \mathbf{b}}$  such that  $\lim_{L \rightarrow \infty} \|f - f_L^{\mathbf{w}, \mathbf{b}}\|_{C(\Omega)} = 0$

- Learning an autonomous dynamical system:

$$\dot{z}(t) = F(z), z \in \mathbb{R}^d, z(0) = z_0$$

where  $F(z)$  is the unknown function.

- Trajectory data:  $k = 1, \dots, K$

the  $k$ -th trajectory =  $\{(t_0, z_0^{(k)}), (t_1, z(t_1; z_0^{(k)})), \dots, (t_N, z(t_N; z_0^{(k)}))\}$

- Goal: Approximate  $F(z)$  using data and neural networks  $F^\theta(z)$

1. For given  $\theta$ , generate trajectories from  $\frac{d\tilde{z}}{dt} = F^\theta(\tilde{z}), \tilde{z}(0) = z_0$

the  $k$ -th trajectory =  $\{(t_0, \tilde{z}_0^{(k)}), (t_1, \tilde{z}^{(k)}(t_1)), \dots, (t_N, \tilde{z}^{(k)}(t_N))\}$

2. Define a loss:  $\mathcal{L}(\theta) = \frac{1}{K} \sum_{k=1}^K \frac{1}{N} \sum_{i=1}^N \|\tilde{z}^{(k)}(t_i) - z(t_i)\|^2$

3.  $\min_\theta \mathcal{L}(\theta) \rightarrow \theta^* \rightarrow F^{\theta^*}(z)$

Pure Data Driven

Q: How to incorporate physics?

A: Embed physical laws into  $F^\theta(z)$

## Dynamical Systems:

### (1) Data Driven

### (2) Physics Informed Data Driven:

$$\int \sum_{\text{f.d.}} (U_{NN} - U)^2 + (\nabla \cdot U_{NN})^2, \text{ s.t. } \nabla \cdot U_{NN} = 0$$

$\frac{dq}{dt} = F(q, t)$ : Discovery of closed form governing laws - equations with sparse/noise and noisy data

To solve  $\frac{dq}{dt} = F(q, t)$  we need  $q(t=0) = q_0$

Bounded Domain  $\Omega$

Black Box to Grey Box

Problem: design NN  $\phi$  s.t.

$x(t_i) \approx \phi(x(t_i - dt))$ .

### Learning Autonomous Dynamical System:

$\dot{z}(t) = F(z), z \in \mathbb{R}^d, z(0) = z_0, F(z)$  unknown

$$\begin{aligned} & \therefore -mg\sin\theta = ma, s = \text{de} \\ & \text{mg}\sin\theta \quad \text{mg}\cos\theta \quad a = l \frac{d^2\theta}{dt^2} \\ & \therefore \frac{d^2\theta}{dt^2} + \frac{g}{l} \sin\theta = 0 \quad \Rightarrow -mg\sin\theta = m l \frac{d^2\theta}{dt^2} \end{aligned}$$

Lagrangian:  $L(x, v) := T - U$

Hamiltonian  $H(q, p) := \sum_i p_i q_i - L(q, \dot{q}, t), p_i = \frac{\partial L}{\partial \dot{q}_i}$

Euler-Lagrange  $\ddot{q}: \frac{d}{dt} \left( \frac{\partial L}{\partial \dot{q}} \right) = \frac{\partial L}{\partial q}$

## Approach 1: Neural ODE

Question: How to get  $\tilde{z}$  from  $\frac{d\tilde{z}}{dt} = F^\theta(\tilde{z}), \tilde{z}(0) = z_0$   
Using numerical integrator

- Forward Euler: Residual block in the ResNet, (He et al., 2016)

$$\tilde{z}(t_{i+1}) = \tilde{z}(t_i) + hF^\theta(\tilde{z}(t_i))$$

- Runge-Kutta: (Rico-Martinez et al., 1992)

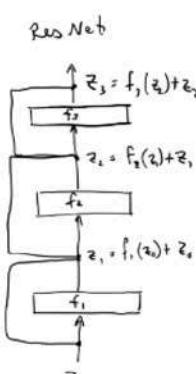
$$k_1 = hF^\theta(\tilde{z}(t_i)), \quad k_2 = hF^\theta(\tilde{z}(t_i) + \frac{k_1}{2})$$

$$k_3 = hF^\theta(\tilde{z}(t_i) + \frac{k_2}{2}), \quad k_4 = hF^\theta(\tilde{z}(t_i) + k_3)$$

$$\tilde{z}(t_{i+1}) = \tilde{z}(t_i) + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4)$$

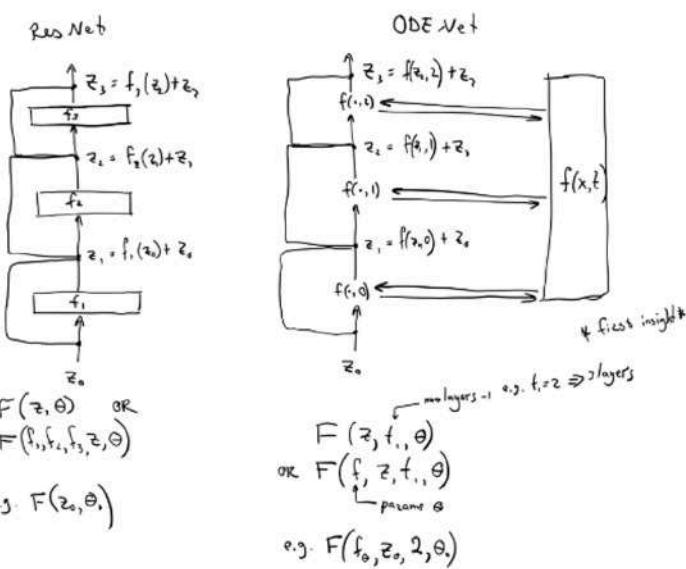
- Black-box/ adaptive integrators: Neural ODE (Chen et al., 2018)

$$\tilde{z}(t_{i+1}) = \text{ODESolve}(\tilde{z}(t_i), F^\theta; h)$$



$$\begin{aligned} F(z, \theta) \quad \text{OR} \\ F(f_1, f_2, f_3, z, \theta) \\ \text{e.g. } F(z_0, \theta) \end{aligned}$$

$$\begin{aligned} \text{e.g. } F(f_0, z_0, 2, \theta) \\ \text{e.g. } F(f_0, z_0, t, \theta) \end{aligned}$$



## Adjoint method

Consider we are dealing with a constrained optimization problem

$$\min_{\mathbf{p}} g(\mathbf{u}(\mathbf{p}), \mathbf{p})$$

where  $\mathbf{u}(\mathbf{p}) \in \mathbb{R}^N$  is the solution to some equation

$$\mathbf{f}(\mathbf{u}(\mathbf{p}), \mathbf{p}) = \mathbf{0}$$

If we wish to optimize  $g$  using some gradient-type algorithm, we need an efficient way to compute the derivative of  $g$  w.r.t.  $\mathbf{p} \in \mathbb{R}^M$ :

$$\frac{dg}{d\mathbf{p}} = \frac{\partial g}{\partial \mathbf{p}} + \frac{\partial g}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} \quad (1)$$

The adjoint method is a general procedure to reduce the computation of  $\frac{dg}{d\mathbf{p}}$  to solving one adjoint equation.

-  $(\mathbf{c}^T \mathbf{A}^{-1}) \frac{\partial \mathbf{b}}{\partial \mathbf{p}}$ . This results in a computational cost of  $O(N(N+M))$ .

-  $\mathbf{c}^T (\mathbf{A}^{-1} \frac{\partial \mathbf{b}}{\partial \mathbf{p}})$ . This results in a computational cost of  $O(NM(1+N))$ .

Clearly, the first method is more efficient as  $M > 1$ .

Summarizing, we can re-formulate the above as

$$\frac{\partial g}{\partial \mathbf{p}} = \boldsymbol{\lambda}^T \frac{\partial \mathbf{b}}{\partial \mathbf{p}}$$

where  $\boldsymbol{\lambda}$  is the solution to

$$\mathbf{A}^T \boldsymbol{\lambda} = \mathbf{c}$$

This is basically the main idea of the adjoint method. Notice that, in fact, we moved to solving a linear system whose matrix is the adjoint of the original linear system (2).

## Adjoint Proof

Let  $\mathbf{z}(t)$  follow the differential equation  $\frac{d\mathbf{z}(t)}{dt} = f(\mathbf{z}(t), t, \theta)$ , where  $\theta$  are the parameters. We will prove that if we define an adjoint state

$$\mathbf{a}(t) = \frac{dL}{d\mathbf{z}(t)}$$

then the adjoint state satisfies the differential equation

$$\frac{da(t)}{dt} = -\mathbf{a}(t) \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} \quad (4)$$

For ease of notation, we denote vectors as row vectors, whereas the main text uses column vectors. The adjoint state is the gradient with respect to the hidden state at a specified time  $t$ . In standard neural networks, the gradient of a hidden layer  $\mathbf{h}_t$  depends on the gradient from the next layer  $\mathbf{h}_{t+1}$  by chain rule

$$\frac{dL}{d\mathbf{h}_t} = \frac{dL}{d\mathbf{h}_{t+1}} \frac{d\mathbf{h}_{t+1}}{d\mathbf{h}_t}$$

With a continuous hidden state, we can write the transformation after an  $\varepsilon$  change in time as

$$\mathbf{z}(t+\varepsilon) = \int_t^{t+\varepsilon} f(\mathbf{z}(t), t, \theta) dt + \mathbf{z}(t) = T_\varepsilon(\mathbf{z}(t), t)$$

and chain rule can also be applied

$$\frac{dL}{d\mathbf{z}(t)} = \frac{dL}{d\mathbf{z}(t+\varepsilon)} \frac{d\mathbf{z}(t+\varepsilon)}{d\mathbf{z}(t)} \quad \text{or} \quad \mathbf{a}(t) = \mathbf{a}(t+\varepsilon) \frac{\partial T_\varepsilon(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)}$$

The proof follows from the definition of derivative:

$$\begin{aligned} \frac{da(t)}{dt} &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}(t+\varepsilon) - \mathbf{a}(t)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}(t+\varepsilon) - \mathbf{a}(t+\varepsilon) \frac{\partial}{\partial \mathbf{z}(t)} T_\varepsilon(\mathbf{z}(t))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}(t+\varepsilon) - \mathbf{a}(t+\varepsilon) \frac{\partial}{\partial \mathbf{z}(t)} (\mathbf{z}(t) + \varepsilon f(\mathbf{z}(t), t, \theta) + \mathcal{O}(\varepsilon^2))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}(t+\varepsilon) - \mathbf{a}(t+\varepsilon) \left( I + \varepsilon \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon^2) \right)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{-\varepsilon \mathbf{a}(t+\varepsilon) \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon^2)}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} -\mathbf{a}(t+\varepsilon) \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon) \\ &= -\mathbf{a}(t) \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} \end{aligned}$$

We point out the similarity between adjoint method and backpropagation. Similarly to backpropagation, ODE for the adjoint  $s$  needs to be solved backwards in time. We specify the constraint on the last time point, which is simply the gradient of the loss at the last time point, and can obtain the gradients with respect to the hidden state at any time, including the initial value,

$$\mathbf{a}(t_N) = \frac{dL}{d\mathbf{z}(t_N)} \quad \mathbf{a}(t_0) = \mathbf{a}(t_N) + \int_{t_N}^{t_0} \frac{da(t)}{dt} dt = \mathbf{a}(t_N) - \int_{t_N}^{t_0} \mathbf{a}(t)^T \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} dt$$

Consider the IVP:  $\frac{dU}{dt} = \mathcal{F}(U, t)$ , where  $U(t=0) = U_0$

Examples:

- 2<sup>nd</sup>-order Adams-Basforth:  $\frac{U^{n+1} - U^n}{\Delta t} = \frac{3}{2} \mathcal{F}^n - \frac{1}{2} \mathcal{F}^{n-1}$
- 3<sup>rd</sup>-order Adams-Basforth:  $\frac{U^{n+1} - U^n}{\Delta t} = \frac{23}{12} \mathcal{F}^n - \frac{16}{12} \mathcal{F}^{n-1} + \frac{1}{12} \mathcal{F}^{n-2}$
- 3<sup>rd</sup>-order Backwards differentiation:  $\frac{U^{n+1} - \frac{1}{11} (18U^n - 9U^{n-1} + 2U^{n-2})}{\Delta t} = \frac{6}{11} \mathcal{F}^{n+1}$ , implicit
- 3<sup>rd</sup>-order Adams-Moulton method:  $\frac{U^{n+1} - U^n}{\Delta t} = \frac{5}{12} \mathcal{F}^{n+1} + \frac{8}{12} \mathcal{F}^n - \frac{1}{12} \mathcal{F}^{n-1}$

General stiffly stable schemes (due to GEAR):				
	$k$	2	3	4
$\alpha_0$	2/3	6/11	12/25	
$\alpha_1$	4/3	18/11	48/25	
$\alpha_2$	-1/3	-9/11	-36/25	
$\alpha_3$	-	2/11	16/25	
$\alpha_4$	-	-	-3/25	

Typically, the coefficients of the Adams-Moulton (AM) methods are smaller than the coefficients of the Adams-Basforth (AB) method, which means lower truncation errors and round-off. Also, for the same accuracy AM uses fewer points but it is implicit.

## Example: A linear problem

Consider the case of a linear function  $g$

$$g(\mathbf{u}) = \mathbf{c}^T \mathbf{u}$$

where  $\mathbf{u}(\mathbf{p})$  is the solution to a (parametric) linear system

$$\mathbf{A}\mathbf{u}(\mathbf{p}) = \mathbf{b}(\mathbf{p}) \quad (2)$$

where  $\mathbf{b} : \mathbb{R}^M \rightarrow \mathbb{R}^N$  and  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is a fixed matrix.  $\mathbf{u}(\mathbf{p})$  is the solution to the system and is thus given by  $\mathbf{u}(\mathbf{p}) = \mathbf{A}^{-1}\mathbf{b}(\mathbf{p})$ . If we derive  $\mathbf{u}$  w.r.t.  $\mathbf{p}$ , we get

$$\frac{\partial \mathbf{u}}{\partial \mathbf{p}} = \mathbf{A}^{-1} \frac{\partial \mathbf{b}}{\partial \mathbf{p}}$$

Plugging this into equation (1), we get

$$\frac{dg}{d\mathbf{p}} = \frac{\partial g}{\partial \mathbf{u}} \frac{\partial \mathbf{u}}{\partial \mathbf{p}} = \mathbf{c}^T \mathbf{A}^{-1} \frac{\partial \mathbf{b}}{\partial \mathbf{p}} \quad (3)$$

Recall that  $\mathbf{c}$  is a  $N$  dimensional vector,  $\mathbf{A}$  is a  $N \times N$  matrix, and  $\frac{\partial \mathbf{b}}{\partial \mathbf{p}}$  is a  $N \times M$  matrix. There are two possible ways to compute the matrix product in the RHS of equation (3):

## Neural ODE and backpropagation

$$L(\mathbf{z}(t_1)) = L\left(\mathbf{z}(t_0) + \int_{t_0}^{t_1} f(\mathbf{z}(t), t, \theta) dt\right) = L(\text{ODESolve}(\mathbf{z}(t_0), f, t_0, t_1, \theta))$$

- To optimize  $L$ , we require gradients with respect to  $\theta$ .
- The first step is to determine how the gradient of the loss depends on the hidden state  $\mathbf{z}(t)$  at each instant. This quantity is called the adjoint and given as
- $\mathbf{a}(t) = \partial L / \partial \mathbf{z}(t)$ .
- Its dynamics are given by another ODE, which can be thought of as the instantaneous analog of the chain rule:

$$\frac{da(t)}{dt} = -\mathbf{a}(t)^\top \frac{\partial f(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}}$$

- We can compute  $\partial L / \partial \mathbf{z}(t_0)$  by another call to an ODE solver. This solver must run backwards, starting from the value of  $\partial L / \partial \mathbf{z}(t_1)$ .
- One complication is that solving this ODE requires knowing value of  $\mathbf{z}(t)$  along its entire trajectory. However, we can simply recompute  $\mathbf{z}(t)$  backwards in time together with the adjoint, starting from its final value  $\mathbf{z}(t_1)$ .

### Algorithm Complete reverse-mode derivative of an ODE initial value problem

```

Input: dynamics parameters  $\theta$ , start time  $t_0$ , stop time  $t_1$ , final state  $\mathbf{z}(t_1)$ , loss gradient  $\partial L / \partial \mathbf{z}(t_1)$ 
def aug_dynamics ([ $\mathbf{z}(t)$ ,  $\mathbf{a}(t)$ ,  $\cdot$ ,  $\cdot$ ,  $\cdot$ ,  $\theta$ ]):
    return  $[f(\mathbf{z}(t), t, \theta), -\mathbf{a}(t)^\top \frac{\partial f}{\partial \mathbf{z}}, -\mathbf{a}(t)^\top \frac{\partial f}{\partial \theta}, -\mathbf{a}(t)^\top \frac{\partial f}{\partial t}]$ 
     $[\mathbf{z}(t_0), \frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_1}] = \text{ODESolve}(s_0, \text{aug\_dynamics}, t_1, t_0, \theta)$ 
    return  $\frac{\partial L}{\partial \mathbf{z}(t_0)}, \frac{\partial L}{\partial \theta}, \frac{\partial L}{\partial t_0}, \frac{\partial L}{\partial t_1}$ 

```

Know  $F$ , solve the following equations for  $z$

$$\sum_{m=0}^M \alpha_m z(t_{n-m}) = h \sum_{m=0}^M \beta_m F(z(t_{n-m}))$$

How to choose  $\alpha$  and  $\beta$ ?

We have 3 families:

- Adams-Basforth(AB)
- Adams-Moulton (AM)
- Backwards Differentiation Formula (BDF)

$$z(t_n) \approx z(t_{n-1}) + \int_{t_{n-1}}^{t_n} \sum_{k \in \Lambda} F(t_{n+k}) l_{k,n}(t; \Lambda) dt$$

where,

$$l_{k,n}(t; \Lambda) = \prod_{i \in \Lambda \setminus \{k\}} \frac{t - t_{n+i}}{t_{n+k} - t_{n+i}}, k \in \Lambda$$

If  $\Lambda_0 = \{-M+1, -M+2, \dots, -1, 0\}$ , we get AB

If  $\Lambda_1 = \{-M+2, -M+3, \dots, 0, 1\}$ , we get AM

BDF:

$$\sum_{i \in \Lambda} z(t_{n+k}) \frac{dl_{k,n}}{dt}(t_n; \Lambda) \approx \frac{d}{dt} z(t_n) = F(z(t_n))$$

## Consistency and stability of AB, AM and BDF

### Definition

A linear multistep method is consistent with the differential equation for dynamics discovery provided  $\|\tau_h\|_\infty \rightarrow 0$  as  $h \rightarrow 0$ , and it is strongly consistent if  $\|\tau_h\|_1 \rightarrow 0$  as  $h \rightarrow 0$ .

### Definition

A linear  $M$ -step method for the dynamics discovery is called stable if there exists a constant  $K < \infty$ , not depending on  $N$ , such that, for any two grid functions  $u, v$ , we have

$$\|u - v\|_1 \leq K \left( \max_{1 \leq i \leq M-1} \|u_i - v_i\|_\infty + \|\hat{R}_h(u - v)\|_1 \right)$$

### Definition

A linear  $M$ -step method for the dynamics discovery is called strongly stable if there exists a constant  $K < \infty$ , not depending on  $N$ , such that, for any two grid functions  $u, v$ , we have

$$\|u - v\|_\infty \leq K \left( \max_{1 \leq i \leq M-1} \|u_i - v_i\|_\infty + \|\hat{R}_h(u - v)\|_\infty \right)$$

### Theorem

If a linear multistep method is strongly consistent and stable, or consistent and strongly stable, then it is convergent.

Pf:

$$\begin{aligned} \|\mathbf{f} - \hat{\mathbf{f}}\|_W &\leq K_W \left( \max_{0 \leq i \leq M-1} \|\mathbf{f}_i - \hat{\mathbf{f}}_i\|_\infty + \|\hat{R}_h \mathbf{f} - R_h \hat{\mathbf{f}}\|_W \right) \\ &\leq K_W (\alpha(1) + \|\tau_h\|_W) \end{aligned}$$

Let's use this theorem to see if AB, AM, BDF are convergent for inverse problem.

## Implementation of Multistep Neural Network: Problem Setup

- Multistep Neural Network: We define the ODE for a non-linear dynamic system:

$$\frac{d}{dt} \mathbf{x}(t) = \mathbf{f}(\mathbf{x}(t)),$$

where  $\mathbf{x}(t) \in \mathbb{R}^D$  denotes the state of system at time  $t$  and function  $\mathbf{f}$  describes the evolution of the system.

- Goal:** Given noisy measurements of the state  $\mathbf{x}(t)$  of the system at several time instances  $t_1, t_2, \dots, t_N$ , our goal is to determine the function  $\mathbf{f}$  consequently discover the underlying dynamical system from data.

- The general form of a linear multistep method with  $M$  steps is

$$\sum_{m=0}^M [\alpha_m \mathbf{x}_{n-m} + \Delta t \beta_m \mathbf{f}(\mathbf{x}_{n-m})] = 0, \quad n = M, \dots, N,$$

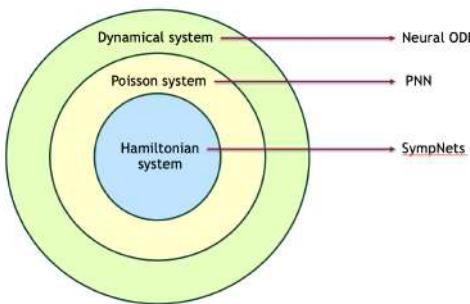
where  $\mathbf{x}_{n-m}$  is state of the system  $\mathbf{x}(t_{n-m})$  at time  $t_{n-m}$ . Different choices for the parameters  $\alpha_m$  and  $\beta_m$  result in specific schemes. For example, trapezoidal rule

$$\mathbf{x}_n = \mathbf{x}_{n-1} + \frac{1}{2} \Delta t (\mathbf{x}_n + \mathbf{f}(\mathbf{x}_{n-1})), \quad n = 1, \dots, N.$$

- The above corresponds to the case where  $M = 1$ ,  $\alpha_0 = -1$ ,  $\alpha_1 = 1$  and  $\beta_0 = \beta_1 = 0.5$ . We put neural netowrk prior on  $\mathbf{f}$ . The parameters of neural network can be learned by optimizing the mean squared loss error function defined as:

$$MSE := \frac{1}{N-M+1} \sum_{n=M}^N \|\mathbf{y}_n\|^2,$$

## Neural ODE, PNN, SympNet



When a model is overfitting, it has large gradient norms compared to the function that we want to reconstruct.

This intuition has been formalized in the work of Oberman et al., 2018

**Idea:** Force the model to have a small Lipschitz constant, while still be close to the observed data (Negrini et al., 2021).

$$\text{Loss}(x, \theta) = \frac{1}{N} \sum_{i=1}^N \|f(x_i) - N(\theta, x_i)\|^2 + \alpha \text{Lip}(N)$$

$$\text{Lip}(N_{int}) = \|\nabla N_{int}\|_{L^\infty(\mathbb{R}^{d+1})}$$

## Hamiltonian Systems: SympNets

### Definition

A matrix  $H \in \mathbb{R}^{2d \times 2d}$  is called symplectic if  $H^T J H = J$        $J := \begin{bmatrix} 0 & I_d \\ -I_d & 0 \end{bmatrix}$

### Definition

A differentiable map  $\phi : U \rightarrow \mathbb{R}^{2d}$  (where  $U \subset \mathbb{R}^{2d}$  is an open set) is called symplectic if the Jacobian matrix  $\frac{\partial \phi}{\partial x}$  is everywhere symplectic, i.e.

$$\left( \frac{\partial \phi}{\partial x} \right)^T J \left( \frac{\partial \phi}{\partial x} \right) = J$$

## Approach 4: Structure-Preserving Neural Networks

The three approaches discussed so far are general purpose models. They are expressive enough to approximate arbitrary dynamical systems, but may not generalize well on specific problems.

### Neural Networks with Prior Physical Knowledge

- Hamiltonian neural networks, symplectic networks

Assume the dynamics are from a Hamiltonian system  $\frac{dx}{dt} = J^{-1} \nabla H^\theta(z)$

- Lagrangian neural network

Assume the dynamics follow the Lagrangian Formalism,

- Poisson neural network

Assume the dynamics are from a Poisson system  $\frac{dx}{dt} = L(z) \nabla H^\theta(z)$

- Graph Neural ODE (Poli et al., 2019)

Assume the RHS of the ODE have certain graph structure (works for molecular dynamics)

- GENERIC formalism informed neural network

Assume the dynamics follow the GENERIC formalism

$$\frac{dx}{dt} = L(z) \nabla E^\theta(z) + M(z) \nabla S^\theta(z)$$

## Summary

- We presented four different approaches to learn dynamical systems from data.
- Approach 1: In Neural ODE the key is to use the adjoint method for backpropagation.
- Approach 2: Use the structure of multi-step methods to learn the coefficients of the RHS from data.
- Approach 3: The recurrent neural network (RNN) family is very effective, especially the Seq2Seq models.
- Approach 4: Neural networks with prior physical knowledge or constraints can be used for Hamiltonian systems, Lagrangian systems, Poisson systems, and general systems by imposing constraints on the structure, e.g. the symplectic structure or the GENERIC formalism of non-equilibrium thermodynamics, which will lead to better generalization.
- Lipschitz regularization of deep neural networks enhances convergence and generalization.

Figure: Neural ODE, PNN and SympNet. Prior knowledge used: SympNet > PNN > Neural ODE; Model expressivity: Neural ODE > PNN > SympNet.

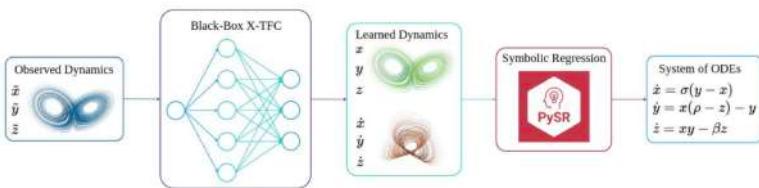
Distillation of NN:

$$\begin{cases} \frac{dx}{dt} = \sigma(y - x) = NN(x, y, t, \theta, \sigma) \\ \frac{dy}{dt} = x(\rho - z) - y = NN(x, y, z, t, \theta, \rho) \\ \frac{dz}{dt} = xy - bz = NN(x, y, z, t, \theta, b) \end{cases}$$

$$\text{given } x(t=0) = x_0, \quad y(t=0) = y_0, \quad z(t=0) = z_0$$

Symbolic Regression (ML) Fit model to predict  $y$  given  $x$ ,  $\mathcal{E}(y, x) = \sum_i^N, y \approx f(x_1, x_2, x_3, \dots)$   
f E Analytic equations

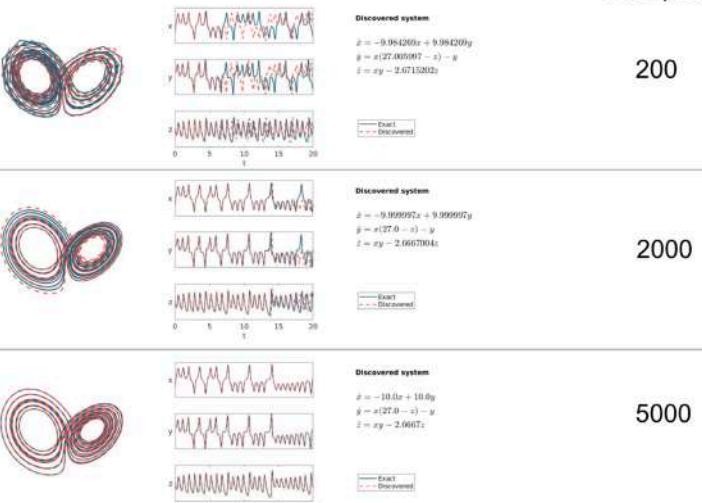
## AI-Lorenz framework



The observed dynamics are used to train the neural network in the black-box X-TFC framework, which learns the dynamics and their rates of change in time. The latter are used as input in the symbolic regression algorithm (PySR), which autonomously distills the explicit mathematical terms that build the system of differential equations.

## Results

### Lorenz system



### Data points

200

2000

5000

## Physics Informed Neural Networks

**Idea:** Encode “physics equations” and “data” into the neural net  $h_\theta(\mathbf{x})$

- **DEF:** Let  $\mathbf{m} = (m_r, m_b)$ . The prototype PINN loss is

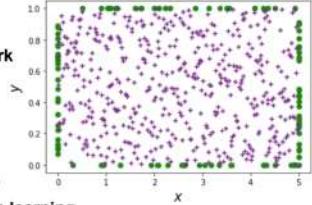
$$\text{Loss}_m^{\text{PINN}}(\theta) = \frac{1}{m_r} \sum_{i=1}^{m_r} \left( \mathcal{L}[h_\theta](\mathbf{x}_r^i) - f(\mathbf{x}_r^i) \right)^2 + \frac{1}{m_b} \sum_{i=1}^{m_b} \left( \mathcal{B}[h_\theta](\mathbf{x}_b^i) - g(\mathbf{x}_b^i) \right)^2$$

- **Method:** minimize  $\text{Loss}_m^{\text{PINN}}(\theta) \implies \theta^*$

$h_m := h_{\theta^*}$  is a physics informed neural network

### Why PINNs?

- Simple and easy to implement it
- Solve PDEs without meshes
- Great potential in solving inverse problems
- Take advantage of rapid development of deep learning



### Approach 3: Dynamic Weights for PINNs

- Loss = Boundary Losses + Residual Losses

$$\mathcal{L} = \lambda_B \mathcal{L}_B + \lambda_f \mathcal{L}_f$$

- Parameters update

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_\theta \mathcal{L}_f - \eta \alpha^{(k+1)} \nabla_\theta \mathcal{L}_B,$$

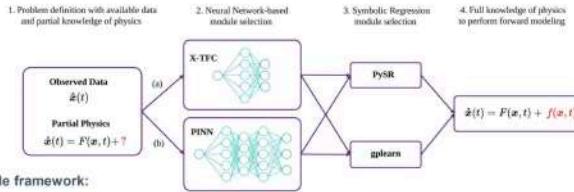
where  $\alpha^{(k+1)} = (1 - \lambda) \alpha^{(k)} + \lambda \hat{\alpha}^{(k+1)}$ , with  $\hat{\alpha}^{(k+1)} = \frac{|\nabla_\theta \mathcal{L}_f|}{|\nabla_\theta \mathcal{L}_B|}$ .

- $|\nabla_\theta \mathcal{L}_B|$  denote the means of  $|\nabla_\theta \mathcal{L}_B|$ .

- $0 \leq \lambda \leq 1$ .

## Gray-Box model Identification AI-Aristotle

Daryakenari, Nazanin Ahmadi, Mario De Florio, Khemraj Shukla, and George Em Karniadakis. \*AI-Aristotle: A Physics-Informed framework for Systems Biology Gray-Box Identification. \* *PLoS computational biology* (2024).

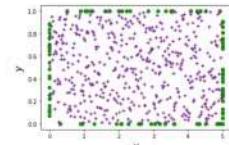
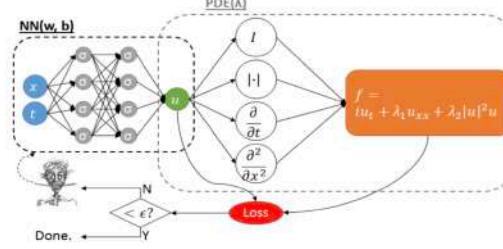


### AI-Aristotle framework:

1. The observed data and the partial knowledge of physics are used to train the selected NN-based module.
2. The selection of the NN-based module needs to be done between (a) X-TFC, recommended for high-resolution data and missing terms discovery, and (b) PINN, recommended for sparse data and parameter estimation. The NN outputs are the time-dependent representations of the missing terms of the dynamical systems, which are fed into the SR algorithm.
3. The selected SR module identifies the mathematical expressions of the missing terms. It is recommended to use both symbolic regressors for cross-validation.
4. The full knowledge of physics is now available, allowing forward modeling performance.

## What is a PINN? Physics-Informed Neural Network

We employ two (or more) NNs that share the same parameters



### Minimize:

$$MSE_u = \frac{1}{N_u} \sum_{i=1}^{N_u} |u(t_u^i, x_u^i) - u^i|^2,$$

$$MSE = MSE_u + MSE_f,$$

$$MSE_f = \frac{1}{N_f} \sum_{i=1}^{N_f} |f(t_f^i, x_f^i)|^2.$$

## Approach 1: PINNs with Soft Constraints

- Loss = Boundary Losses + Residual Losses

$$\mathcal{L} = \lambda_B \mathcal{L}_B + \lambda_f \mathcal{L}_f$$

## Approach 2: Self-Adaptive PINNs

PDE is Defined as

$$\begin{aligned} \mathcal{N}_{\sigma,t}[u(x,t)] &= 0, \quad x \in \Omega, t \in [0, T] \\ u(x,t) &= g(x,t), \quad x \in \partial\Omega, t \in [0, T] \\ u(x,0) &= h(x), \quad x \in \Omega \end{aligned}$$

- Wang S, Yu X, Perdikaris P. PINNs fail to train: A new perspective. *Journal of Physics: Conference Series*. 2022 Jan 15;4

$$\mathcal{L}(\mathbf{w}) = \mathcal{L}_s(\mathbf{w}) + \mathcal{L}_r(\mathbf{w}) + \mathcal{L}_b(\mathbf{w}) + \mathcal{L}_0(\mathbf{w})$$

$$\begin{aligned} \mathcal{L}_s(\mathbf{w}) &= \frac{1}{N_s} \sum_{i=1}^{N_s} \left| u(\mathbf{x}_s^i, t_s^i; \mathbf{w}) - y_s^i \right|^2 \\ \mathcal{L}_r(\mathbf{w}) &= \frac{1}{N_r} \sum_{i=1}^{N_r} r(\mathbf{x}_r^i, t_r^i; \mathbf{w})^2 \\ \mathcal{L}_b(\mathbf{w}) &= \frac{1}{N_b} \sum_{i=1}^{N_b} \left| u(\mathbf{x}_b^i, t_b^i; \mathbf{w}) - g_b^i \right|^2 \\ \mathcal{L}_0(\mathbf{w}) &= \frac{1}{N_0} \sum_{i=1}^{N_0} \left| u(\mathbf{x}_0^i, 0; \mathbf{w}) - h_0^i \right|^2 \end{aligned}$$

## Setup for PINNs

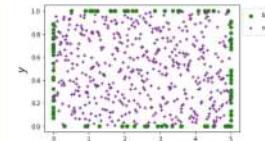
Consider the partial differential equations (PDEs):

$$\begin{aligned} \mathcal{L}[u](\mathbf{x}) &= f(\mathbf{x}) \text{ in } U \subset \mathbb{R}^d \\ \mathcal{B}[u](\mathbf{x}) &= g(\mathbf{x}) \text{ in } \partial U \end{aligned}$$

Goal: Find a neural net  $h_\theta(\mathbf{x})$  that approximates the PDE solution.

What are known to us?

- Pointwise PDE data: i.i.d. samples  $(\mathbf{x}_r^i, f(\mathbf{x}_r^i))$  where  $\mathbf{x}_r^i \in U$
- Residual data:  $(\mathbf{x}_r^i, f(\mathbf{x}_r^i))$  where  $\mathbf{x}_r^i \in \partial U$
- Boundary data:  $(\mathbf{x}_b^i, g(\mathbf{x}_b^i))$  where  $\mathbf{x}_b^i \in \partial U$
- Physics Equations:  $\mathcal{L}[u] - f = 0$  and  $\mathcal{B}[u] - g = 0$



# Physics-Informed Neural Networks

Idea: Encode “physics equations” and “data” into the neural net  $h_\theta(\mathbf{x})$

**DEF:** Let  $\mathbf{m} = (m_r, m_b)$ . The prototype PINN loss is

$$\text{Loss}_m^{\text{PINN}}(\theta) = \frac{1}{m_r} \sum_{i=1}^{m_r} (\mathcal{L}[h_\theta](\mathbf{x}_r^i) - f(\mathbf{x}_r^i))^2 + \frac{1}{m_b} \sum_{i=1}^{m_b} (\mathcal{B}[h_\theta](\mathbf{x}_b^i) - g(\mathbf{x}_b^i))^2$$

**DEF:** The Hölder regularized PINN loss is

$$\text{Loss}_m^{\text{Höld}}(\theta) = \text{Loss}_m^{\text{PINN}}(\theta) + \lambda_r^R \left( [\mathcal{L}[h_\theta]]_{\alpha;U} \right)^2 + \lambda_b^R \left( [\mathcal{B}[h_\theta]]_{\alpha;\partial U} \right)^2$$

Hölder constant  $\rightarrow [u]_{\alpha;U} = \sup_{x \neq y \in U} \frac{\|u(x) - u(y)\|}{\|x - y\|^\alpha}$

Goal: minimize  $\text{Loss}_m^{\text{Höld}}(\theta) \implies \theta^*$

$h_m := h_{\theta_m}$  is a physics informed neural network  $\implies \{h_m\}$

Q: Does  $\{h_m\}$  converge to the PDE solution? In what norm?

## Hölder Regularization: Motivation

Ideally, we wish to minimize the expected PINN loss :

$$\text{Loss}_\infty^{\text{PINN}}(h) = \mathbb{E} [\text{Loss}_m^{\text{PINN}}(h)]$$

**Theorem:** (Shin, Darbon, Karniadakis, 20) Suppose iid samples. Then, with high probability over iid samples,

$$\text{Loss}_\infty^{\text{PINN}}(\theta) \leq C_m \text{Loss}_m^{\text{Höld}}(\theta) + C' \left( m_r^{-\frac{\alpha}{d}} + m_b^{-\frac{\alpha}{d-1}} \right)$$

Expected PINN loss Hölder Regularized empirical loss

- $C_m = \Theta(\max\{\sqrt{m_r}, \sqrt{m_b}\})$
- $C'$  is a universal constant
- $\lambda^R = (\lambda_r^R, \lambda_b^R), \quad \lambda_r^R = \Theta(C_m^{-1} m_r^{-\alpha/d}), \quad \lambda_b^R = \Theta(C_m^{-1} m_b^{-\alpha/(d-1)})$ .

## Summary

- PINN is simple and can be implemented in a few lines of code for any PDE
- Boundary conditions need special treatment and can be implemented via self-adaptive or dynamic weights
- Weighted residual methods can be implemented following a Petrov-Galerkin projection
- Domain decomposition can enhance the accuracy and flexibility of PINNs
- PINN converges for boundary values problems even in the absence of BCs
- Several works have obtained a priori and a posteriori errors for PINNs
- There are three main errors: Approximation, Generalization and Optimization
- PINN is more accurate and faster than the Deep Ritz Method

### Hard Constraints: Dirichlet BC and IC

Dirichlet BC or IC :  $u_i(\mathbf{x}) = g_0(\mathbf{x}), \quad \mathbf{x} \in \Gamma_D$ ,

Trial solution:  $\hat{u}_i(\mathbf{x}; \theta_u) = g(\mathbf{x}) + \ell(\mathbf{x})\mathcal{N}(\mathbf{x}; \theta_u)$ ,

$g(x)$  is a continuous extension of  $g_0(x)$  from  $\Gamma_D$  to  $\Omega$  where  $\mathcal{N}(\mathbf{x}; \theta_u)$  is the network output, and  $\ell$  is a function satisfying the following two conditions:

$$\begin{cases} \ell(\mathbf{x}) = 0, & \mathbf{x} \in \Gamma_D \\ \ell(\mathbf{x}) > 0, & \mathbf{x} \in \Omega - \Gamma_D \end{cases}$$

For example:  $u(0) = 0, u(1) = 1 : g(x) = x, \ell(x) = x(1-x)$

