Neural networks basics

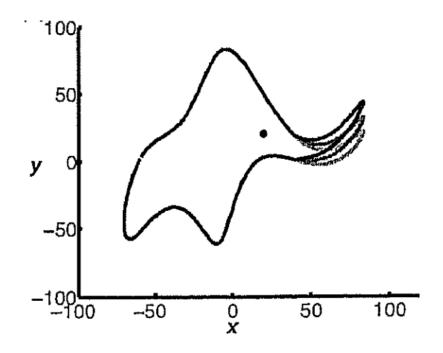
by C. Yao Lai for 2025 Princeton summer school

Outline

- Neural networks in problems governed by PDEs
- NN basics (Lecture 1-3)
 Universal function approximation
 Gradient descent
 Back propagation
- Physics-informed NN (Lecture 4-5)

Von Neumann's elephant

"With four parameters I can fit an elephant, and with five I can make him wiggle his trunk." -Von Neumann



An approximation using five complex parameters was found by Mayer, Khairy, Howard (2010) based on complex Fourier analysis.

(1)

$$x(t) = \sum_{k=0}^{\infty} (A_k^x \cos(kt) + B_k^x \sin(kt)),$$

$$y(t) = \sum_{k=0}^{\infty} (A_k^{\gamma} \cos(kt) + B_k^{\gamma} \sin(kt)),$$

Table I. The five complex parameters p_1, \ldots, p_5 that encode the elephant including its wiggling trunk.

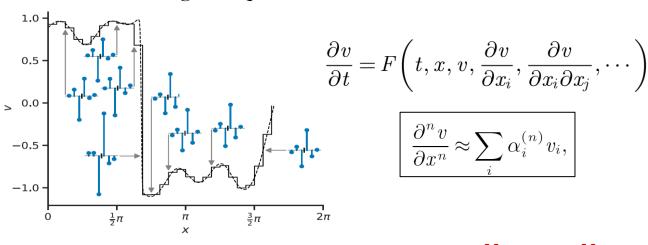
Parameter	Real part	Imaginary par
$p_1 = 50 - 30i$	$B_1^x = 50$	B ^r =−30
$p_2 = 18 + 8i$	$B_{2}^{x}=18$	$B_2^{y}=8$
$p_3 = 12 - 10i$	$A_3^{x} = 12$	$B_3^y = -10$
$p_4 = -14 - 60i$	$A_{5}^{x}=-14$	$A_1^{y} = -60$
$\rho_5 = 40 + 20i$	Wiggle coeff.=40	$x_{\rm eye} = y_{\rm eye} = 20$

Early day ML-PDE

Viscous Burgers eqn

• Leaning how to discretize PDEs (2019)





ML maps discretized solution $(v_1 \dots v_{N_x})$ to the stencil coefficients $(\alpha_1 \dots \alpha_N): \mathbb{R}^{N_x} \to \mathbb{R}^N$

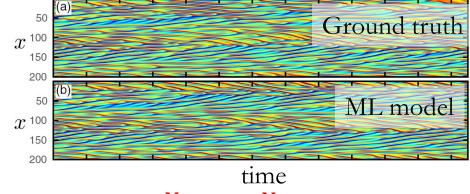
• Learning a ML model to predict simple chaotic system (2018)

KS eqn

PHYSICAL REVIEW LETTERS 120, 024102 (2018)

Model-Free Prediction of Large Spatiotemporally Chaotic Systems from Data: A Reservoir Computing Approach

Jaideep Pathak, 1,2,* Brian Hunt, 3,4 Michelle Girvan, 1,3,2 Zhixin Lu, 1,3 and Edward Ott 1,2,5



ML maps solution from one time step to the subsequent time step: $\mathbb{R}^{N_x} \to \mathbb{R}^{N_x}$



ML weather forecast models

ECMWF is now running a series of data-driven forecasts as part of its experimental suite. These machine-learning based models are very fast, and they produce a 10-day forecast with 6-hourly time steps in approximately one minute. The outputs are available in graphical form.



Explore content Y About the journal Y Publish with us Y

nature > articles > article

Article Open access Published: 04 December 2024

Probabilistic weather forecasting with machine learning

<u>Ilan Price</u> $\stackrel{\ \ \, }{\square}$, <u>Alvaro Sanchez-Gonzalez</u>, <u>Ferran Alet</u>, <u>Tom R. Andersson</u>, <u>Andrew El-Kadi</u>, <u>Dominic Masters</u>, Timo Ewalds, Jacklynn Stott, Shakir Mohamed, Peter Battaglia $\stackrel{\ \ \, }{\square}$, Remi Lam $\stackrel{\ \ \, }{\square}$ & Matthew Willson $\stackrel{\ \ \, }{\square}$

However, these advances have focused primarily on single, deterministic forecasts that fail to represent uncertainty and estimate risk. Overall, MLWP has remained less accurate and reliable than state-of-the-art NWP ensemble forecasts. Here we introduce GenCast, a probabilistic weather model with greater skill and speed than the top operational medium-range weather forecast in the world, ENS, the ensemble forecast of the European Centre for Medium-Range Weather Forecasts⁴. GenCast is an ML weather prediction method, trained on

nature

About the journal ➤ Publish with us ➤

nature > articles > article

Article Open access | Published: 05 July 2023

Accurate medium-range global weather forecasting with 3D neural networks

<u>Kaifeng Bi</u>, <u>Lingxi Xie</u>, <u>Hengheng Zhang</u>, <u>Xin Chen</u>, <u>Xiaotao Gu</u> & <u>Qi Tian</u> □

Nature **619**, 533–538 (2023) | Cite this article

Science

urrent Issue First release papers

pers Archive Ab

ACTICE

RESEARCH ARTICLE | WEATHER FORECASTING

Learning skillful medium-range global weather forecasting



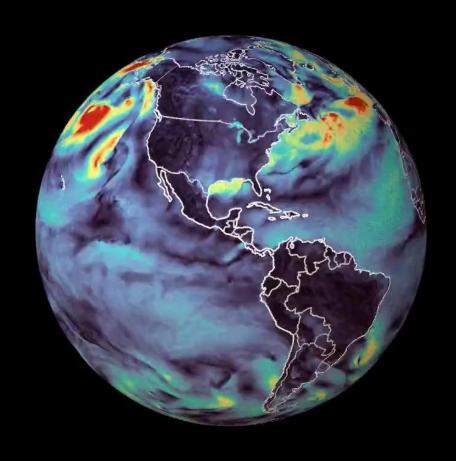
Abstract

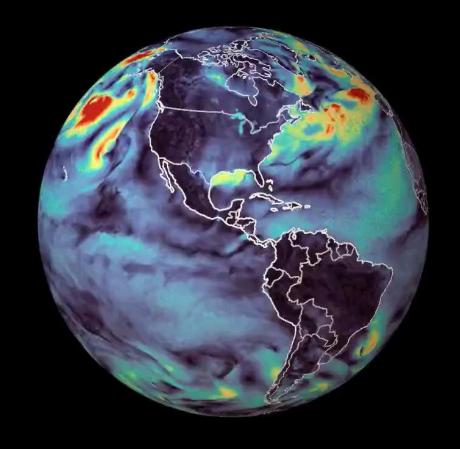
Global medium-range weather forecasting is critical to decision-making across many social and economic domains. Traditional numerical weather prediction uses increased compute resources to improve forecast accuracy but does not directly use historical weather data to improve the underlying model. Here, we introduce GraphCast, a machine learning—based method trained directly from reanalysis data. It predicts hundreds of weather variables for the next 10 days at 0.25° resolution globally in under 1 minute. GraphCast significantly outperforms the most accurate operational deterministic systems on 90% of 1380 verification targets, and its forecasts support better severe event prediction, including tropical cyclone tracking, atmospheric rivers, and extreme temperatures. GraphCast is a key advance in accurate and efficient weather forecasting and helps realize the promise of machine learning for modeling complex dynamical systems.

Pangu-Weather

ML weather forecast models

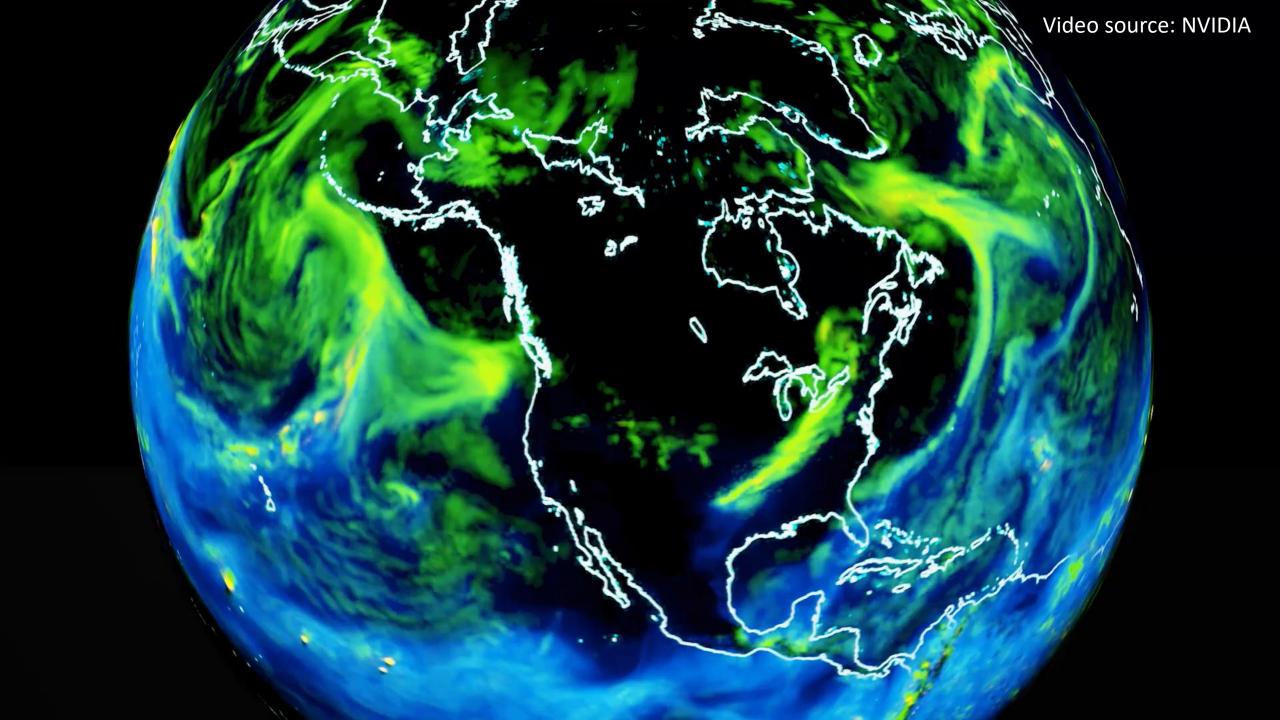
2018-01-01





SFNO

Ground truth



In addition to generating known results faster, can we leverage neural networks to discovery new PDE solutions?

Self-similar blow-up solution for Boussinesq eqn

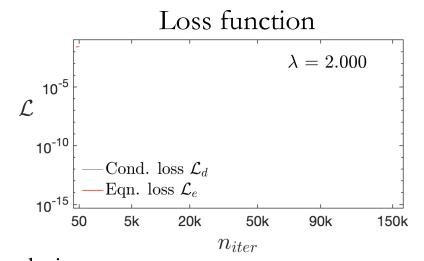
$$f_{1} = \Omega + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla\Omega - \Phi$$

$$f_{2} = (2 + \partial_{y_{1}}U_{1})\Phi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla\Phi + \partial_{y_{1}}U_{2}\Psi$$

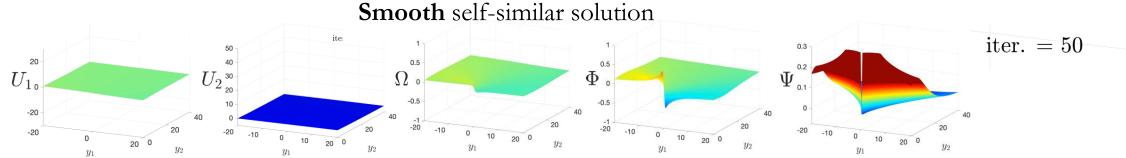
$$f_{3} = (2 + \partial_{y_{2}}U_{2})\Psi + ((1 + \lambda)\mathbf{y} + \mathbf{U}) \cdot \nabla\Psi + \partial_{y_{2}}U_{1}\Phi$$

$$f_{4} = \partial_{y_{1}}U_{1} + \partial_{y_{2}}U_{2}$$

$$f_{5} = \Omega - (\partial_{y_{1}}U_{2} - \partial_{y_{2}}U_{1})$$



Wang-Lai-GómezSerrano-Buckmaster,



Uniform and small equation residues everywhere

Physical Review Letters (2023) f_4 $^{\circ}$ f_3 o f_2 o 10

So...why neural networks?

Very often in math and science we want to know A as a function of B. NN
has the flexibility to fit that function without prior assumptions about the
functional form, particularly useful when the function is high dimensional
(it is a function approximator).

• Leverages GPU (hardware) and open-source software (e.g. tensoflow, Pytorch, JAX). The optimization method is well-developed and user friendly. It can usually be easily adopted for different problems.

We will talk about the challenges after learning about the basics of NN!

Outline

- Neural networks in problems governed by PDEs
- NN basics (Lecture 1-3)
 Universal function approximation
 Gradient descent
 Back propagation
- Physics-informed NN (Lecture 4-5)

What is a neural network?

An analytical model of output y as a function of input x, containing some fitting parameters

1. Linear Regression Model:

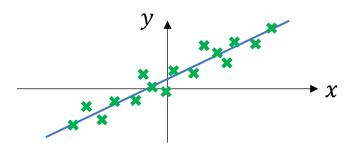
$$y = wx + b$$

$$x \xrightarrow{w} \xrightarrow{b} y$$

Given observations of $\{x_d^i, y_d^i\}_i^n$

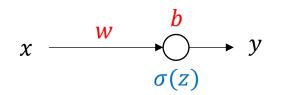
Find the w and b that minimizes

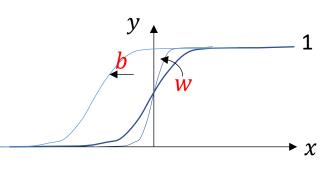
$$J = \sum_{i=1}^{n} (y(x_d^i) - y_d^i)^2$$



2. Logistic Regression Model: make output 0 to 1

$$y = \sigma(wx + b),$$
where $\sigma(z) = \frac{1}{1 + e^{-z}}$ is a sigmoid function





What is a neural network?

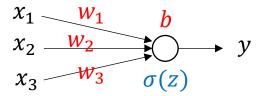
An analytical model of output y as a function of input x, containing some fitting parameters

2. Logistic Regression Model: make output -1 to 1

$$y = \sigma(w_1x_1 + w_2x_2 + w_3x_3 + b),$$

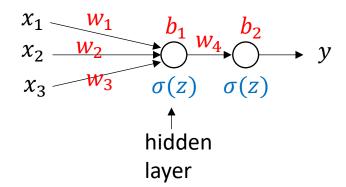
$$where \sigma(z) = \frac{1}{1 + e^{-z}} \text{ is a sigmoid function}$$

$$x_1 \quad w_1 \quad b \quad x_2 \quad w_2 \quad y \quad x_3 \quad \sigma(z)$$



3. Neural network:

$$y = \sigma(w_4\sigma(w_1x_1 + w_2x_2 + w_3x_3 + b_1) + b_2)$$
, where $\sigma(z)$ is a nonlinear activation function



Common choices of $\sigma(z)$

$$\begin{array}{c} sigmoid(z) \\ sin(z) \\ cos(z) \\ tanh(z) \end{array}$$
 Output ranges from -1 to 1

What is a neural network?

An analytical model of output y as a function of input x, containing some fitting parameters

3. Neural network:

More generally...

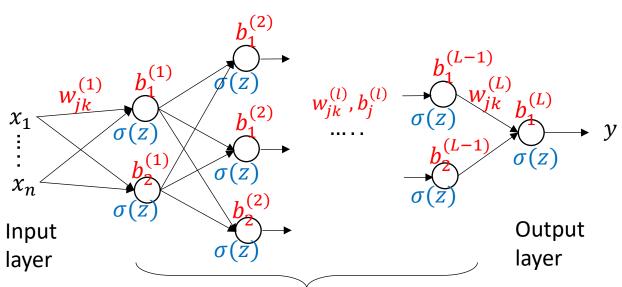
$$y = fn(\mathbf{w}_{jk}^{(l)}, \mathbf{b}_{j}^{(l)}, \mathbf{x}_{i}),$$

) is a nonlinear activation function

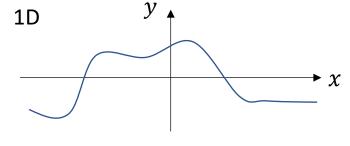
where $\sigma(z)$ is a nonlinear activation function

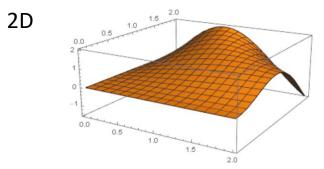
Neural network is a general function approximation

-> Universal function approximation



L-1 hidden layers





n-Dimension surface....

Why are activation function nonlinear?

Common choices of $\sigma(z)$ sigmoid(z) sin(z) cos(z) tanh(z)... $x_1 \qquad w_{jk} \qquad w_{11} \qquad b_{1} \qquad y_{2} \qquad y_{2}$

If activation function is linear,
 NN can only represent a linear function

If activation fn is linear, e.g. $\sigma(z) = z...$

Input x_1, x_2

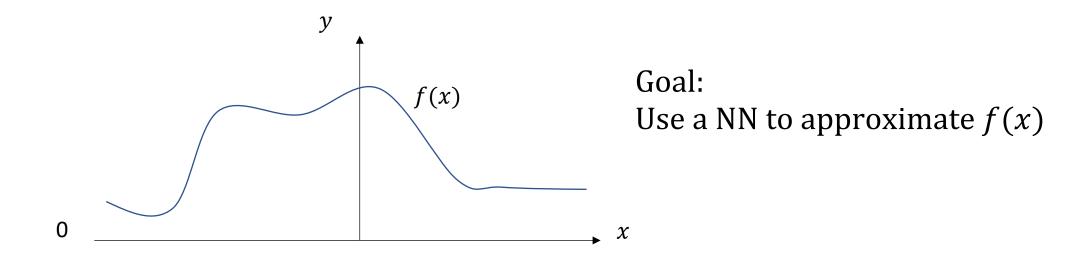
In hidden layer, j=1,2

$$z_j^{(1)} = \sum_{k=1}^{2} w_{jk}^{(1)} x_k + b_j^{(1)}$$
$$a_j^{(1)} = \sigma(z_j^{(1)}) = z_j^{(1)}$$

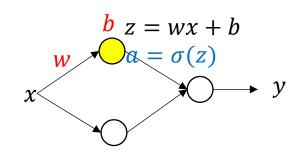
In output layer, j=1

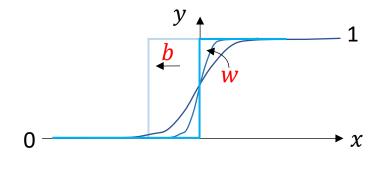
$$z_j^{(2)} = \sum_{k=1}^{2} w_{jk}^{(2)} a_k^{(1)} + b_j^{(2)}$$
$$a_j^{(2)} = \sigma \left(z_j^{(2)} \right) = z_j^{(2)}$$

$$y = a_1^{(2)} = \sum_{k=1}^{2} w_{1k}^{(2)} \left(\sum_{l=1}^{2} w_{kl}^{(1)} x_l + b_k^{(1)} \right) + b_1^{(2)} = Ax_1 + Bx_2 + C$$

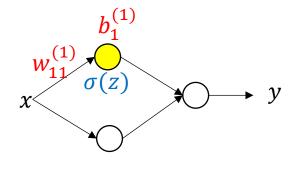


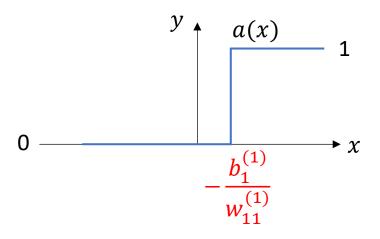
 NN can approximate continuous and smooth functions A visual proof (for sigmoid activation)



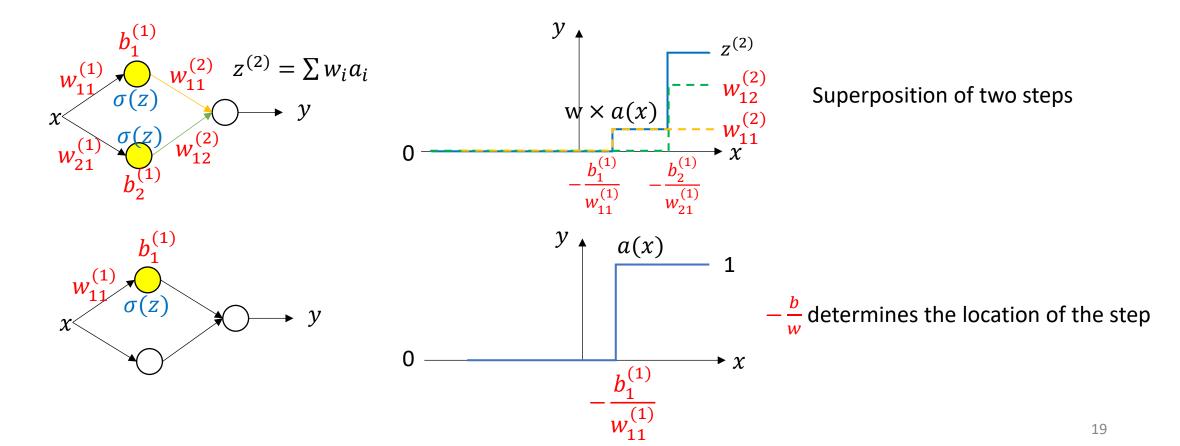


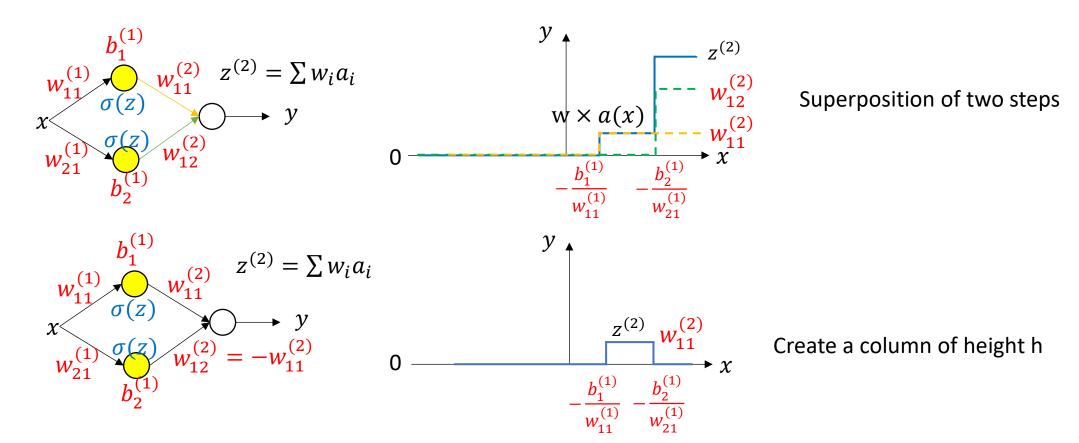
Large w gives a step

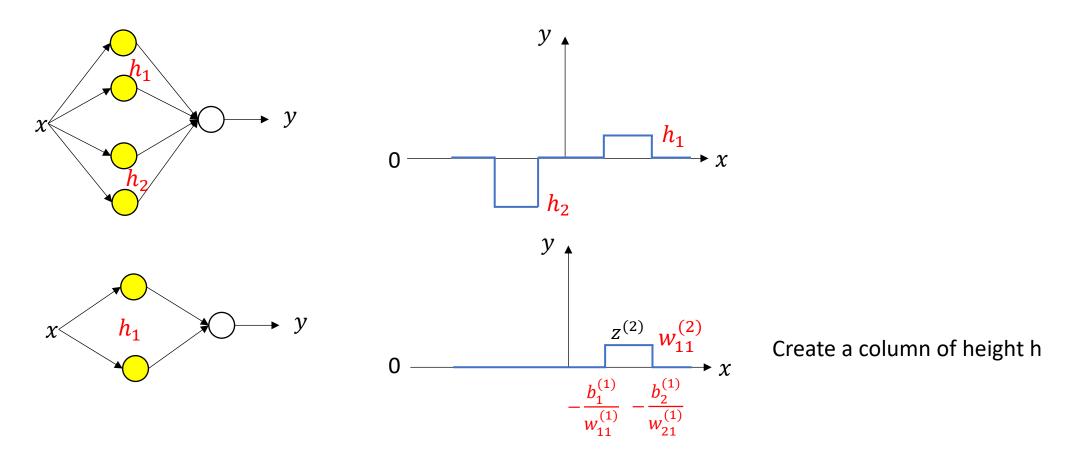


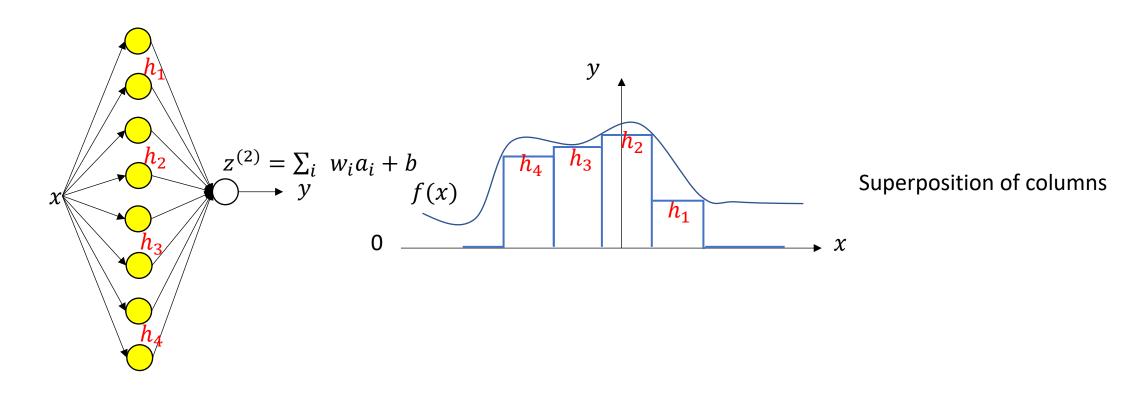


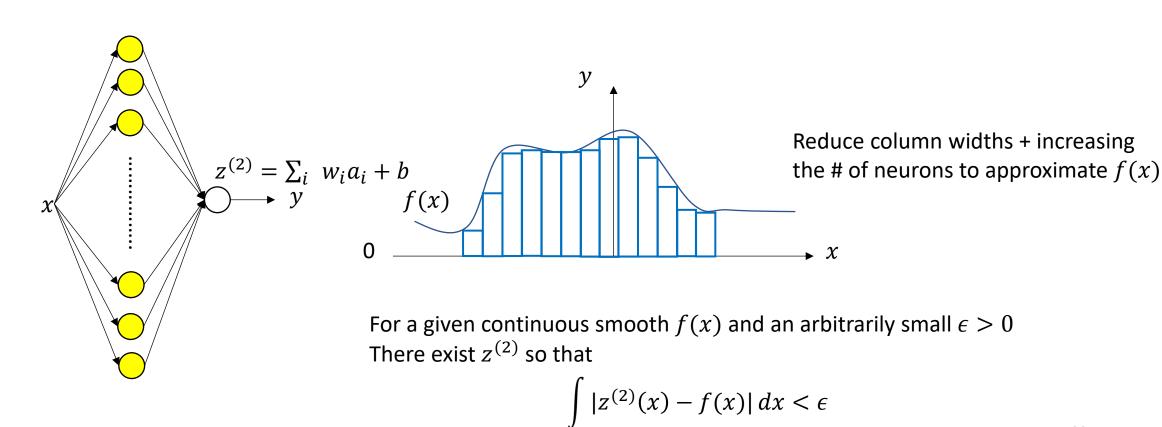
 $-\frac{b}{w}$ determines the location of the step





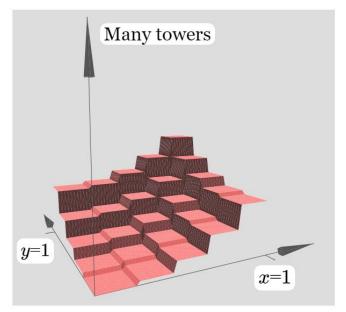






Approximate a 2D surface

What should be the input/output units of NN?





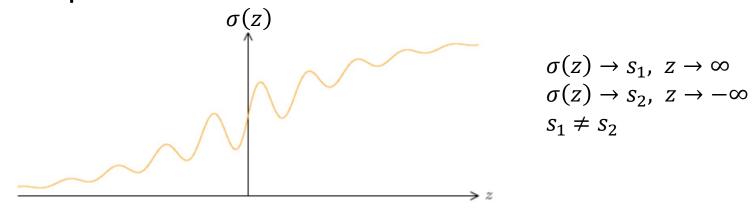
In summary

Dimension of the approximated surface is determined by ...
 Number of input and output units

Column size is determined by ...
 Number of hidden units, values of weights and biases



 Could the following activation function instead of a sigmoid function approximate a step?

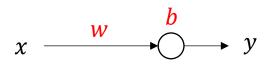


• Could activation function $\sigma(z)=z$ instead of a sigmoid function approximate a step?

Now we know it is possible to tune weights and biases in a NN to approximate functions. We still need an automated method to find the correct weights and biases.

How to find w and b?

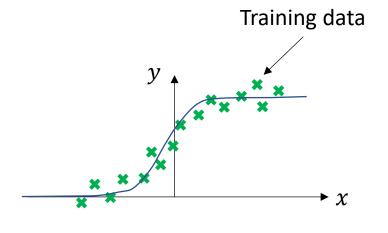
• E.g., Model $y = \sigma(wx + b)$



Given observations of $\{x_d^i, y_d^i\}_i^m$

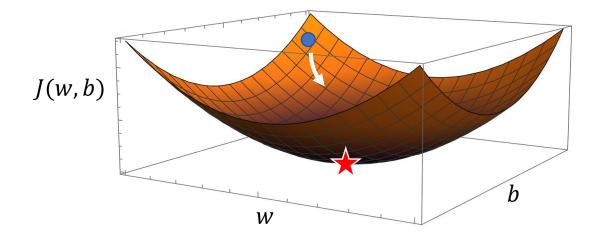
Our goal is to find the model parameters w, b that minimizes the cost function

$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) - y_d^i)^2$$



How to find w and b?

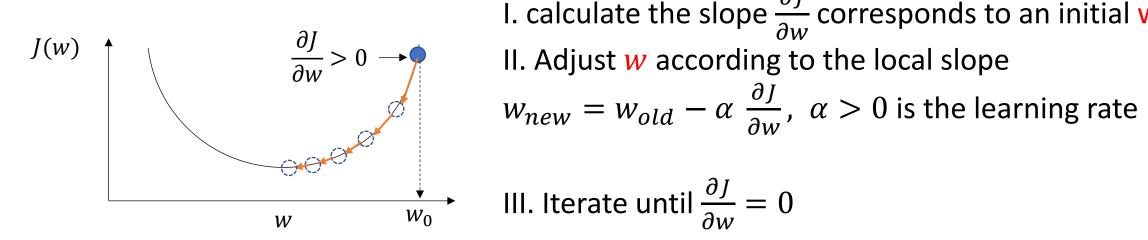
- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$



Find the w, b that minimizes J(w, b)

i.e.
$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial b} = 0$$

- E.g., Model $y = \sigma(\mathbf{w}x + \mathbf{b})$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$

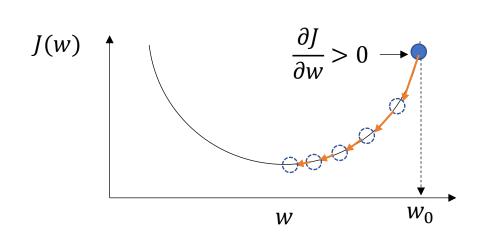


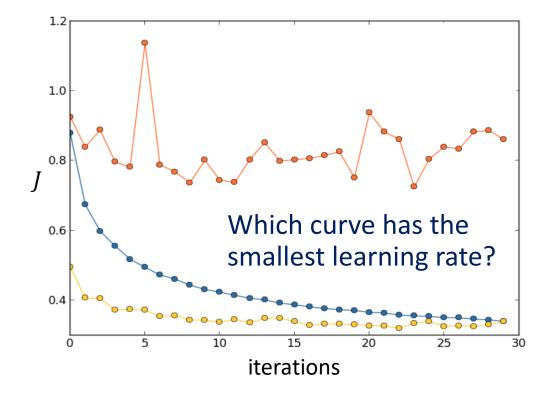
- I. calculate the slope $\frac{\partial J}{\partial w}$ corresponds to an initial w

$$w_{new} = w_{old} - lpha \; rac{\partial J}{\partial w}$$
 , $\; lpha > 0$ is the learning rate

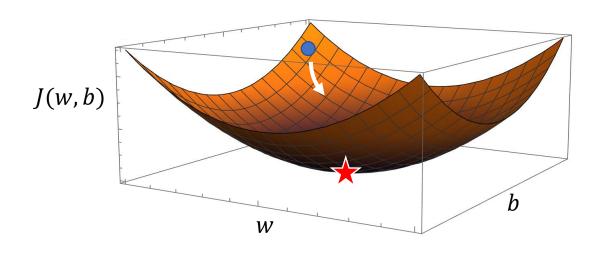
III. Iterate until
$$\frac{\partial J}{\partial w} = 0$$

- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$





- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$



Gradient on a surface

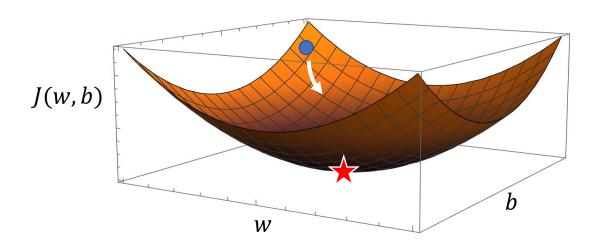
• $-\nabla J$ gives the direction of the steepest decrease of J

$$-\nabla J(w,b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$$

e.g. $-\nabla J(w_0,b_0)=-(10,1)$ What does this mean? Changing w reduces J 10 times faster than changing b

• $-\nabla J$ tells you which weights and biases reduce cost function J the fastest!

- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$



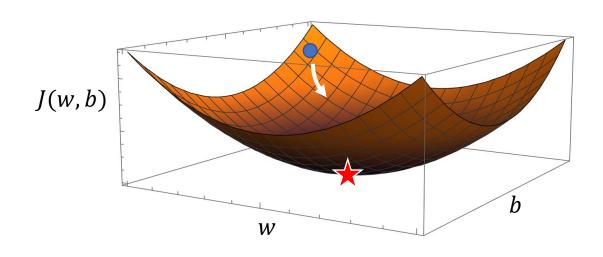
Gradient on a surface

• $-\nabla J$ gives the direction of the steepest decrease of J

$$-\nabla J(w,b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$$

- $(w_{new}, b_{new}) = (w_{old}, b_{old}) \alpha \nabla J(w, b),$ α is the learning rate
- Iterate until $\nabla J = 0$

- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2$



Gradient on a surface

• $-\nabla J$ gives the direction of the steepest decrease of J

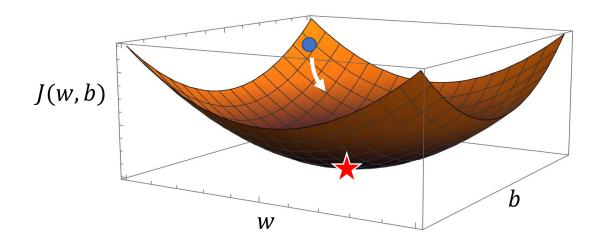
$$-\nabla J(w,b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$$

But how are $\frac{\partial J}{\partial w}$, $\frac{\partial J}{\partial b}$ calculated?

• E.g., Model $y = \sigma(\mathbf{w}x + \mathbf{b})$

Defined for every example i

• Cost function
$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) - y_d^i)^2 = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$$



What would $\frac{dL}{dw}$ be for a linear model $\sigma(z) = z$?

$$L \equiv (y(x) - y_d)^2$$
, $J = \frac{1}{m} \sum_{i=1}^m L \rightarrow \frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial w}$

Back propagation (chain rule)

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial (wx+b)} \frac{\partial (wx+b)}{\partial w} = 2(y-y_d)\sigma'x$$

1) For one example:

$$\frac{\partial L}{\partial w}(w, b, x_d^i, y_d^i) = 2(\sigma(wx_d^i + b) - y_d^i)\sigma' x_d^i$$

2) For the full data set:

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial L}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} 2(\sigma(wx_d^i + b) - y_d^i) \sigma' x_d^i$$

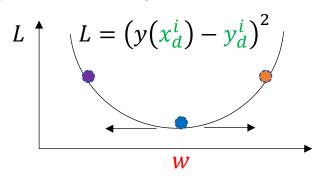
Example:

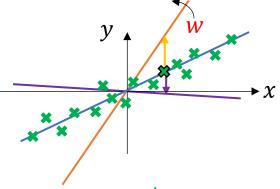
$$x \longrightarrow y$$

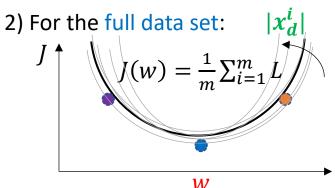
• Linear Regression Model y = wx

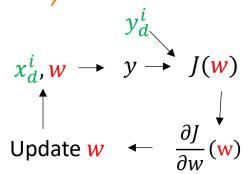
Defined for every example i

- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2 = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$
- 1) For one example:









$$L \equiv (y(x) - y_d)^2$$
, $J = \frac{1}{m} \sum_{i=1}^m L \rightarrow \frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial w}$

Back propagation (chain rule)

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial (wx)} \frac{\partial (wx)}{\partial w} = 2(y - y_d)x$$

1) For one example:

1) For one example:
$$\frac{\partial L}{\partial w}(w, b, x_d^i, y_d^i) = 2(wx_d^i - y_d^i) x_d^i$$

This gradient is defined for every example i, and is a function of w

2) For the full data set:

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial L}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} 2(wx_d^i - y_d^i) x_d^i$$

Example:

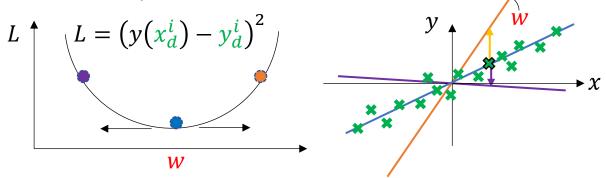


• Linear Regression Model y = wx

Defined for every example i

• Cost function
$$J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) - y_d^i)^2 = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$$

1) For one example:



What would happen if x_d^l are very large/small?

Note that $\frac{\partial J}{\partial w}$ is analytical fn of \mathbf{w}, x_d^i, y_d^i !

$$L \equiv (y(x) - y_d)^2$$
, $J = \frac{1}{m} \sum_{i=1}^m L \rightarrow \frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial w}$

Back propagation (chain rule)

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial (wx)} \frac{\partial (wx)}{\partial w} = 2(y - y_d)x$$

1) For one example:

$$\frac{\partial L}{\partial w}(w, b, x_d^i, y_d^i) = 2(wx_d^i - y_d^i) x_d^i$$

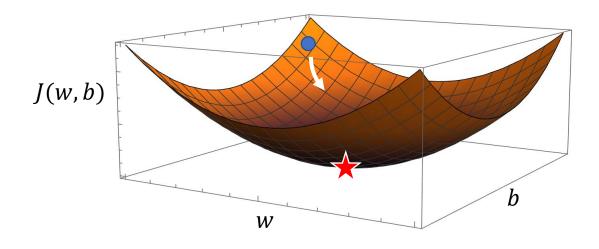
This gradient is defined for every example i, and is a function of w

2) For the full data set:

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} \frac{\partial L}{\partial w} = \frac{1}{m} \sum_{i=1}^{m} 2(wx_d^i - y_d^i) x_d^i$$

How to find w and b? Gradient descent

- E.g., Model $y = \sigma(wx + b)$
- Cost function $J(w,b) = \frac{1}{m} \sum_{i=1}^{m} (y(x_d^i) y_d^i)^2 = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$



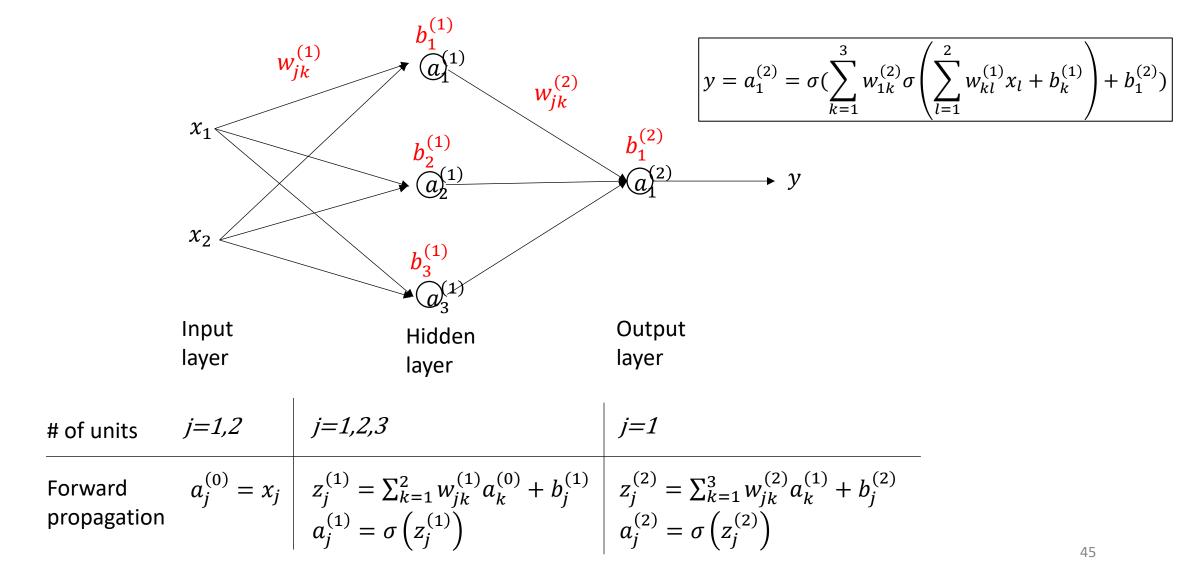
- $-\nabla J(w,b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$ for a given data set at a given w, b is known analytically
- $(w_{new}, b_{new}) = (w_{old}, b_{old}) \alpha \nabla J(w, b),$ α is the learning rate
- Iterate until $\nabla J = 0$

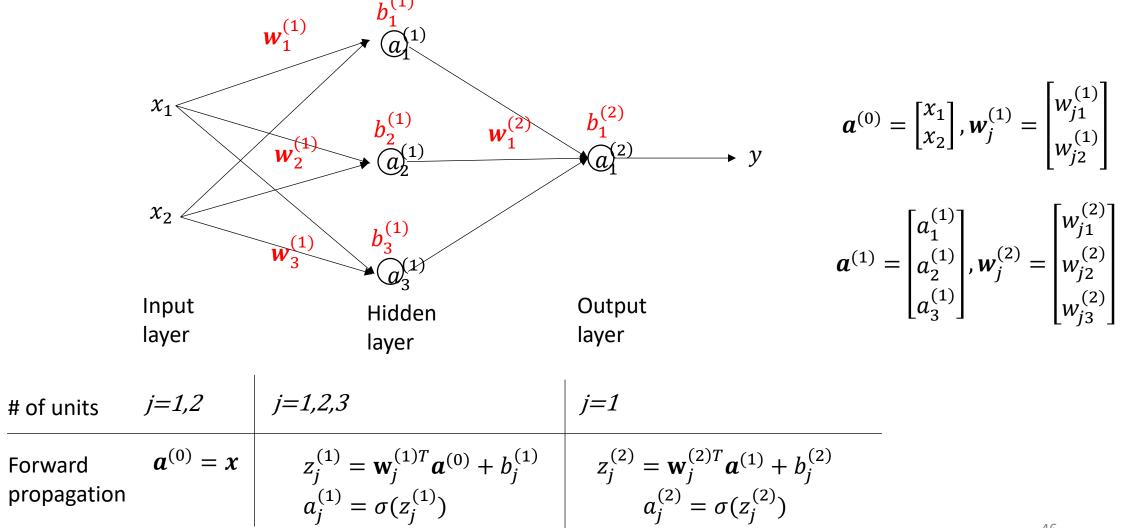
So far we have discussed...

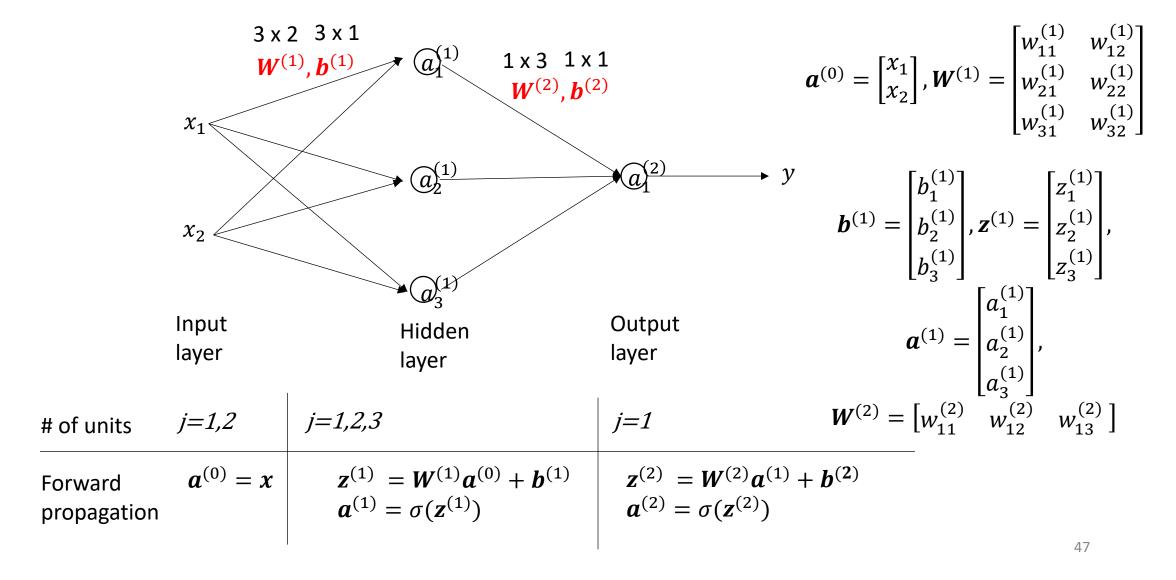
- Universal function approximator
- Gradient descent: a method to find weights and biases that minimize J
- Calculate $-\nabla J(w,b) = -\left(\frac{\partial J}{\partial w}, \frac{\partial J}{\partial b}\right)$ for a simple model $y = \sigma(wx+b)$
 - One example
 - A full dataset

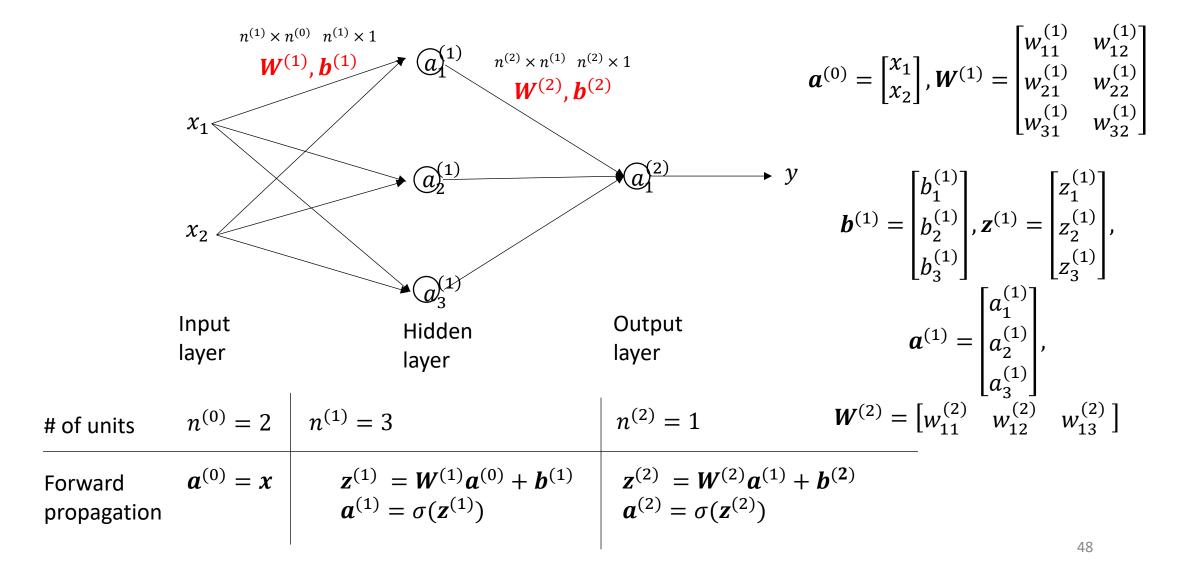
We know how to find w and b for a simple model $y = \sigma(wx + b)$.

What about finding w and b in a neural network?

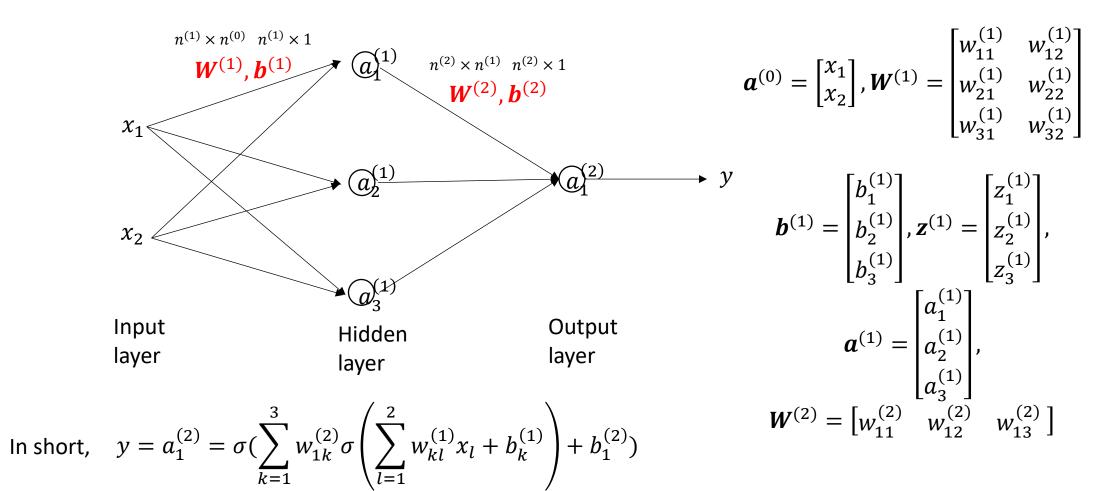








 $\rightarrow y = a^{(2)} = \sigma(W^{(2)}\sigma(W^{(1)}x + b^{(1)}) + b^{(2)})$

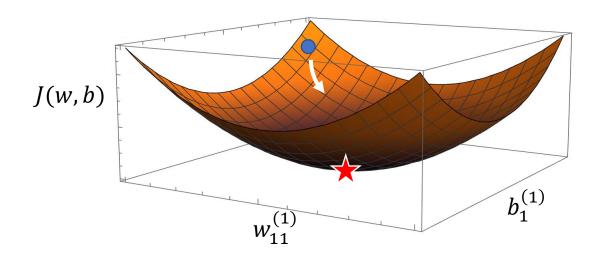


We know how to find w and b for a simple model $y = \sigma(wx + b)$.

What about finding w and b in a neural network? $y = a^{(2)} = \sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$

How to find w and b in a neural network?

- E. g., Model $y = a^{(2)} = \sigma(\mathbf{W}^{(2)}\sigma(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}) + \mathbf{b}^{(2)})$
- Cost function $J(\mathbf{W}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(1)}, \mathbf{b}^{(2)}) = \frac{1}{m} \sum_{i=1}^{m} L(y^i, y_d^i)$

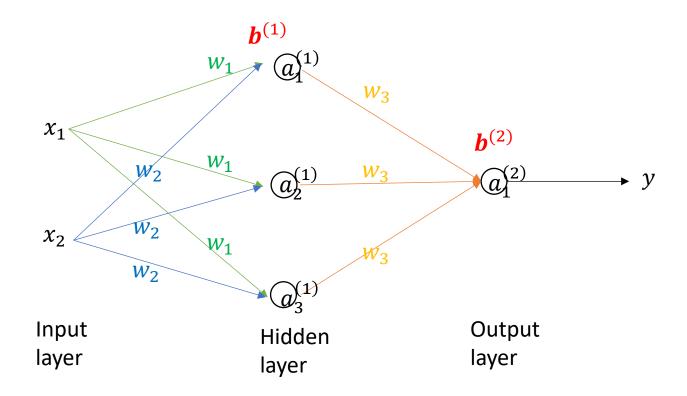


• For the full data set (i=1...m), compute

$$-\nabla J = -\left(\frac{\partial J}{\partial w_{jk}^{(l)}}, \dots, \frac{\partial J}{\partial b_{j}^{(l)}}\right)$$

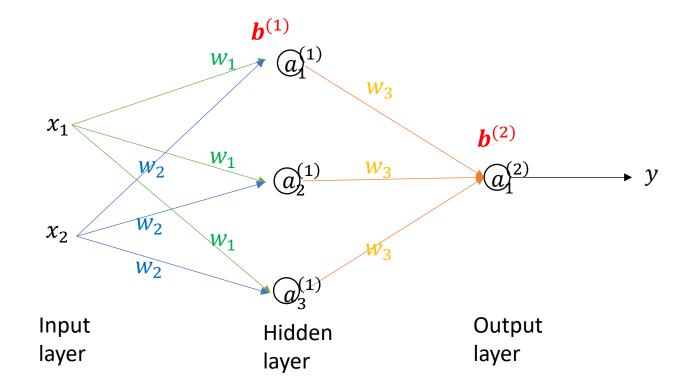
- $w_{jk}{}^{(l)}_{new} = w_{jk}{}^{(l)}_{old} \alpha \frac{\partial J}{\partial w_{jk}^{(l)}}, \ b_{j}{}^{(l)}_{new} = b_{j}{}^{(l)}_{old} \alpha \frac{\partial J}{\partial b_{j}^{(l)}}$ α is the learning rate
- Iterate until $\nabla J = 0$

NN weight Initialization



Q: What would happen if all initial weights are chosen to be the same (e.g., all zeros)?

NN weight Initialization



More specifically, if weights are symmetric:

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_1 & w_2 \\ w_1 & w_2 \\ w_1 & w_2 \end{bmatrix}, \mathbf{W}^{(2)} = \begin{bmatrix} w_3 & w_3 & w_3 \end{bmatrix}$$

(e.g. all zero weights, all constant weights)

Then
$$a_1^{(1)} = a_2^{(1)} = a_3^{(1)}$$

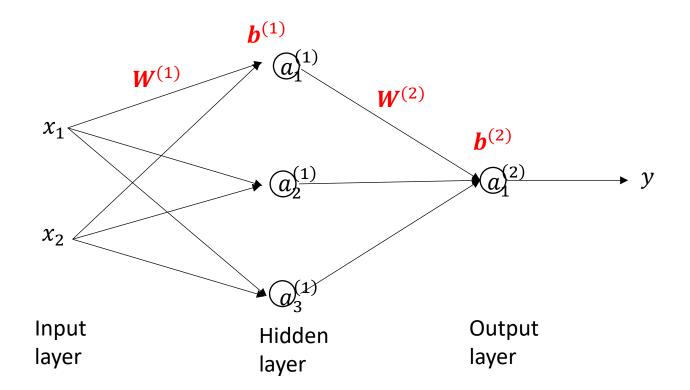
and

$$\begin{bmatrix} w_{1new} & w_{2new} \\ w_{1new} & w_{2new} \\ w_{1new} & w_{2new} \end{bmatrix} = \begin{bmatrix} w_1 & w_2 \\ w_1 & w_2 \\ w_1 & w_2 \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial w_1} & \frac{\partial J}{\partial w_2} \end{bmatrix}$$

The NN never learns to have non-symmetric weights -> $a_1^{(1)}=a_2^{(1)}=a_3^{(1)}$ during training

If the NN weights are **symmetric** \rightarrow Gradient descent will update these hidden units in the same way \rightarrow All hidden units in the same layer will be identical throughout training iterations \rightarrow Equivalent to NN with just 1 hidden unit.

NN weight Initialization



To make the different hidden units approximate different functions

 \rightarrow Initialized weights $w_{ik}^{(l)}$ to random values

What about biases?

→ Biases can just be zeros

$$\mathbf{b}^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \ \mathbf{b}^{(2)} = [0]$$

If the NN weights are **symmetric** \rightarrow Gradient descent will update these hidden units in the same way → All hidden units in the same layer will be identical throughout training iterations \rightarrow Equivalent to NN with just 1 hidden unit.

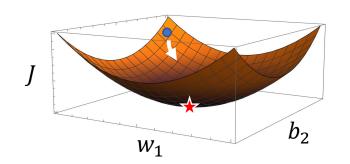
Neural Networks for Pattern Recognition Christopher M. Bishop

Chapter 4 The Multi-layer Perceptron

Forward & Back Propagation: 1. single neuron

Exercise:

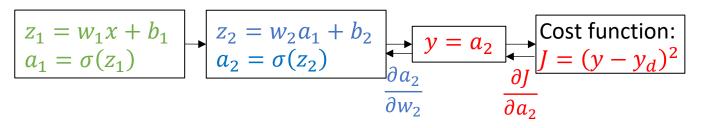
Model: $y = \sigma(w_2\sigma(w_1x + b_1) + b_2)$



 ∇J vary with w_1, w_2, b_1, b_2



Forward propagation →



← Back propagation

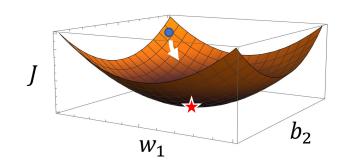
How important is w_2 for changing the cost function J?

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial a_2} \frac{\partial a_2}{\partial w_2} = 2(y - y_d) \sigma'(z_2) a_1$$

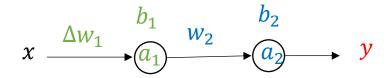
Forward & Back Propagation: 1. single neuron

Exercise:

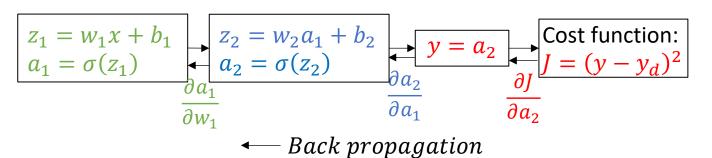
Model: $y = \sigma(w_2\sigma(w_1x + b_1) + b_2)$



 ∇J vary with w_1, w_2, b_1, b_2



Forward propagation →



 ∂J ∂J ∂a_2 ∂a_3

How important is w_1 for

changing the cost function /?

$$\frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial w_1}$$
$$= 2(y - y_d)\sigma'(z_2)w_2\sigma'(z_1)x$$

When training a neural network

Repeat these steps:

- 1. Forward propagate an input
- 2. Compute the cost function
- 3. Compute the gradients of the cost with respect to parameters using backpropagation
- 4. Update each parameter using the gradients, according to the optimization algorithm

Forward & Back Propagation: 1. single neuron

$$x \xrightarrow{w_1} \cdots a \xrightarrow{(L-1)w^{(L)}} a \xrightarrow{(L)} y$$

$$x = a^{(0)} \qquad y = a^{(L)}$$

Forward propagation

Input:
$$a^{(0)} = x$$

...
$$z^{(L-1)} = w^{(L-1)}a^{(L-2)} + b^{(L-1)}$$

$$a^{(L-1)} = \sigma(z^{(L-1)})$$

$$z^{(L)} = w^{(L)}a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$
Output: $y = a^{(L)}$

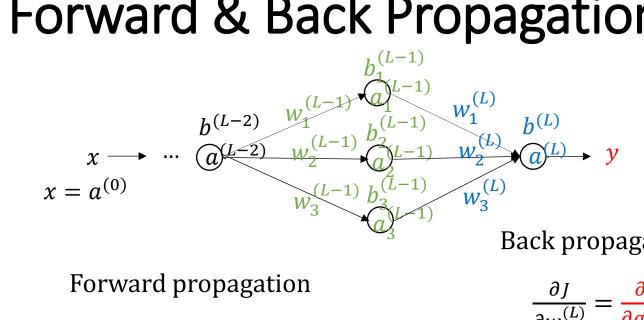
Cost function $J = (y - y_d)^2$

Back propagation (use chain rules)

$$\frac{\partial J}{\partial w^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial b^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial b^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial w^{(L-1)}}$$

$$\frac{\partial J}{\partial b^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a^{(L-1)}} \frac{\partial a^{(L-1)}}{\partial b^{(L-1)}}$$
$$= 2(a^{(L)} - y_d) \sigma'(z^{(L)}) w^{(L)} \sigma'(z^{(L-1)})$$

Forward & Back Propagation: 2. multiple hidden units



Input:
$$a^{(0)} = x$$

$$z_k^{(L-1)} = w_k^{(L-1)} a^{(L-2)} + b_k^{(L-1)}$$

$$a_k^{(L-1)} = \sigma(z_k^{(L-1)})$$

$$z^{(L)} = \sum_{k=1}^3 w_k^{(L)} a_k^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$
Output: $y = a^{(L)}$

Cost function
$$J = (y - y_d)^2$$

Back propagation (use chain rules)

$$\frac{\partial J}{\partial w_k^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial w_k^{(L)}} = 2(a^{(L)} - y_d) \sigma'(\mathbf{z}^{(L)}) a_k^{(L-1)}$$

$$\frac{\partial J}{\partial b^{(L)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial b^{(L)}} = 2(a^{(L)} - y_d) \sigma'(\mathbf{z}^{(L)})$$

$$\frac{\partial J}{\partial w_k^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial w_k^{(L-1)}}$$

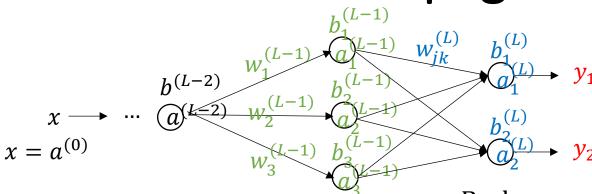
$$= 2(a^{(L)} - y_d) \sigma'(\mathbf{z}^{(L)}) w_k^{(L)} \sigma'(\mathbf{z}_k^{(L-1)}) a^{(L-2)}$$

$$\frac{\partial J}{\partial b_k^{(L-1)}} = \frac{\partial J}{\partial a^{(L)}} \frac{\partial a^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_k^{(L-1)}}{\partial b_k^{(L-1)}}$$

$$= 2(a^{(L)} - y_d) \sigma'(\mathbf{z}^{(L)}) w_k^{(L)} \sigma'(\mathbf{z}_k^{(L-1)})$$

$$= 2(a^{(L)} - y_d) \sigma'(\mathbf{z}^{(L)}) w_k^{(L)} \sigma'(\mathbf{z}_k^{(L-1)})$$

Forward & Back Propagation: 3. multiple outputs



Forward propagation

Input:
$$a^{(0)} = x$$

$$z_{k}^{(L-1)} = w_{k}^{(L-1)} a^{(L-2)} + b_{k}^{(L-1)}$$

$$a_{k}^{(L-1)} = \sigma(z_{k}^{(L-1)})$$

$$z_{j}^{(L)} = \sum_{k=1}^{3} w_{jk}^{(L)} a_{k}^{(L-1)} + b_{j}^{(L)}$$

$$a_{j}^{(L)} = \sigma(z_{j}^{(L)})$$
Output: $y_{j} = a_{j}^{(L)}$

Cost function
$$J = \sum_{j=1}^{2} (y_j - y_{d,j})^2$$

$$\frac{\partial J}{\partial w_{jk}^{(L)}} = \frac{\partial J}{\partial a_{j}^{(L)}} \frac{\partial a_{j}^{(L)}}{\partial w_{jk}^{(L)}} = 2(a_{j}^{(L)} - y_{d,j}) \, \sigma'(z_{j}^{(L)}) a_{k}^{(L-1)}$$

$$\frac{\partial J}{\partial b_{j}^{(L)}} = \frac{\partial J}{\partial a_{j}^{(L)}} \frac{\partial a_{j}^{(L)}}{\partial b_{j}^{(L)}} = 2(a_{j}^{(L)} - y_{d,j}) \, \sigma'(z_{j}^{(L)})$$

$$\frac{\partial J}{\partial w_{k}^{(L-1)}} = \sum_{j=1}^{2} \frac{\partial J}{\partial a_{j}^{(L)}} \frac{\partial a_{j}^{(L)}}{\partial a_{k}^{(L-1)}} \frac{\partial a_{k}^{(L-1)}}{\partial w_{k}^{(L-1)}}$$

$$= \sum_{j=1}^{2} 2(a_{j}^{(L)} - y_{d,j}) \, \sigma'(z_{j}^{(L)}) w_{jk}^{(L)} \, \sigma'(z_{k}^{(L-1)}) a^{(L-2)}$$

$$\frac{\partial J}{\partial b_{k}^{(L-1)}} = \sum_{j=1}^{2} \frac{\partial J}{\partial a_{j}^{(L)}} \frac{\partial a_{j}^{(L)}}{\partial a_{k}^{(L-1)}} \frac{\partial a_{k}^{(L-1)}}{\partial b_{k}^{(L-1)}}$$

$$= \sum_{j=1}^{2} 2(a_{j}^{(L)} - y_{d,j}) \, \sigma'(z_{j}^{(L)}) w_{jk}^{(L)} \, \sigma'(z_{k}^{(L-1)})$$

$$= \sum_{j=1}^{2} 2(a_{j}^{(L)} - y_{d,j}) \, \sigma'(z_{j}^{(L)}) w_{jk}^{(L)} \, \sigma'(z_{k}^{(L-1)})$$

Summary

- Universal function approximator
- Gradient descent: a method to find weights and biases that minimize J
- Calculate \(\nabla J\)
 - One example
 - A full dataset
- Weight initialization
- Forward and backpropagation (a clever way to get gradients of J)

Outline

- Neural networks in problems governed by PDEs
- NN basics (Lecture 1-3)
 Universal function approximation
 Gradient descent
 Back propagation
- Advanced topics: NN's spectral bias
- Physics-informed NN (Lecture 4-5)

Spectral bias: A Fourier analysis of NN

Xu et al. (2022)

For a single-hidden-layer NN with input $x \in \mathbb{R}$ and output $u \in \mathbb{R}$

$$u(x) = \sum_{j=1}^{m} w_j^{(1)} \sigma(w_j^{(0)} x + b_j), \text{ where } \theta_j = \{w_j^{(0)}, w_j^{(1)}, b_j\}$$

The mean square loss measures the distance between the NN output u(x) and ground truth $u_g(x)$

$$J \equiv \int_{-\infty}^{\infty} \frac{1}{2} (u(x) - u_g(x))^2 dx = \int_{-\infty}^{\infty} \frac{1}{2} |\hat{u}(k) - \hat{u}_g(k)|^2 dk$$
, where $\hat{u}(k)$ denotes the Fourier transform of $u(x)$

The loss at frequency k is $J(k) = \frac{1}{2}|\hat{u} - \hat{u}_g|^2$

For tanh activation function, the gradient of loss at J(k) with respect to NN parameters θ_j is

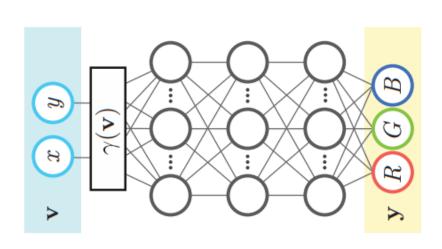
$$\left| \frac{\partial J(k)}{\partial \theta_j} \right| \approx A(k) e^{-\left| \frac{\pi k}{2w_j^{(0)}} \right|} \quad \text{where } A(k) \in [0, +\infty) \text{ is the amplitude of } \hat{u} - \hat{u}_g$$

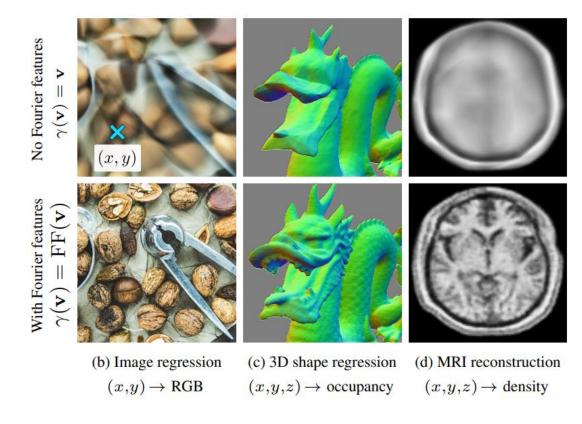
This gradient decay rapidly with frequency k!! This is a simple explanation of the spectral bias.

Spectral bias is bad for multiscale problems

• One solution: Fourier feature $\gamma(\mathbf{v})$ Tancik et al (2020), Wang et al (2021) $\gamma(\mathbf{v})$ maps the input coordinates to a feature space before passing them to NN layers.

$$\gamma(\mathbf{v}) = \left[a_1 \cos(2\pi \mathbf{b}_1^{\mathrm{T}} \mathbf{v}), a_1 \sin(2\pi \mathbf{b}_1^{\mathrm{T}} \mathbf{v}), \dots, a_m \cos(2\pi \mathbf{b}_m^{\mathrm{T}} \mathbf{v}), a_m \sin(2\pi \mathbf{b}_m^{\mathrm{T}} \mathbf{v}) \right]^{\mathrm{T}}$$





• When NN training error is small, the learning becomes really slow. NN learns better when $u \approx 0(1)$. Inspired by the principals of perturbation theory, MSNN includes a superposition of multi-stage NNs, with each stage using a new NN to fit the residue rescaled to O(1) from the previous NN.

$$u_c^{(n)}(x) = u_0(x) + \epsilon_1 u_1(x) + \epsilon_2 u_2(x) + \dots + \epsilon_n u_n(x) = \sum_{j=0}^n \epsilon_j u_j(x)$$

where
$$1 \gg \epsilon_1 \gg \epsilon_2 \gg ... \gg \epsilon_n$$

 $u_0(x)$, $u_1(x)$ and $u_n(x)$ are all of order $\mathcal{O}(1)$ NNs of different stages

Related work with similar flavors: Multi-level neural networks: Aldirany et al. (2023), Multifidelity neural networks: Howard et al. (2023, 2024), Precision Machine Learning: Michaud et al. (2023)...etc 69

Multistage-NN (MSNN) Wang & Lai, J. Comput. Phys. (2024)

