

# **ECE/SIOC 228 Machine Learning for Physical Applications**

## **Lecture 3: Neural Networks and Backpropagation**

**Yuanyuan Shi**

Assistant Professor, ECE

University of California, San Diego

# Logistics & Survey Results

# Homework

UC San Diego

- Assignment 1 released
- Due date: April 21<sup>st</sup> Thursday, 2:00pm (before class)
- Late submission: Late homework within 2 days of the due date is acceptable but the grades of the late homework will be 25% discounted each late day (50% discounted for 2 days). No late homework after 2 days will be accepted without a note from the health center or the student's Resident Dean. In order to ensure fairness, we will be STRICT on the late homework policy.

# Project Website

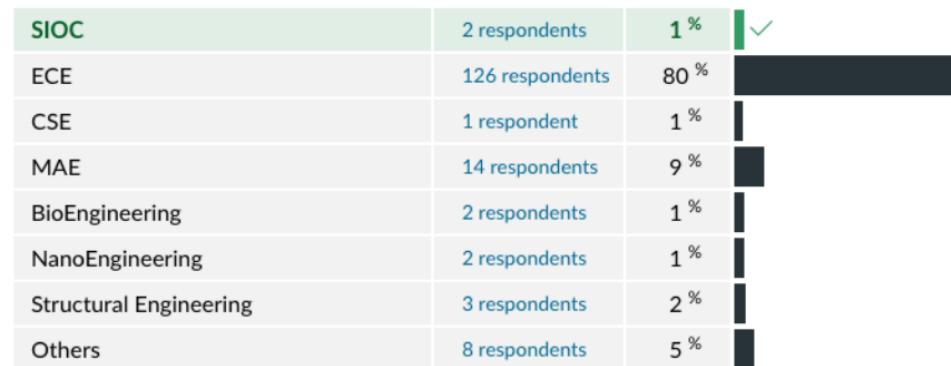
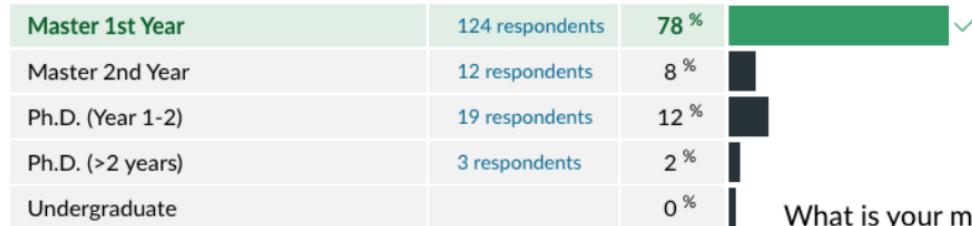
UC San Diego

- Project website: <https://sites.google.com/eng.ucsd.edu/ecesioc228/home>

# Survey Statistics

UC San Diego

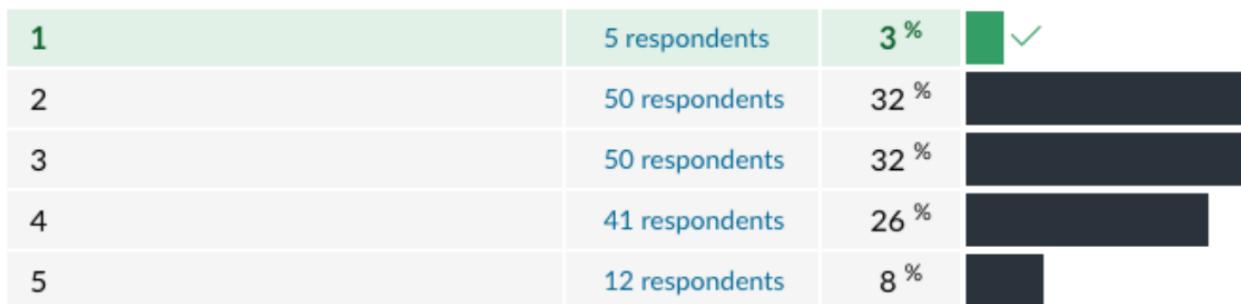
What is degree program are you pursuing and which year are you? (Master, Ph.D., Junior, Senior)



# Topic Familiarity

UC San Diego

1. Deep learning basics: supervised learning, neural networks, gradient descent, linear regression, logistic regression

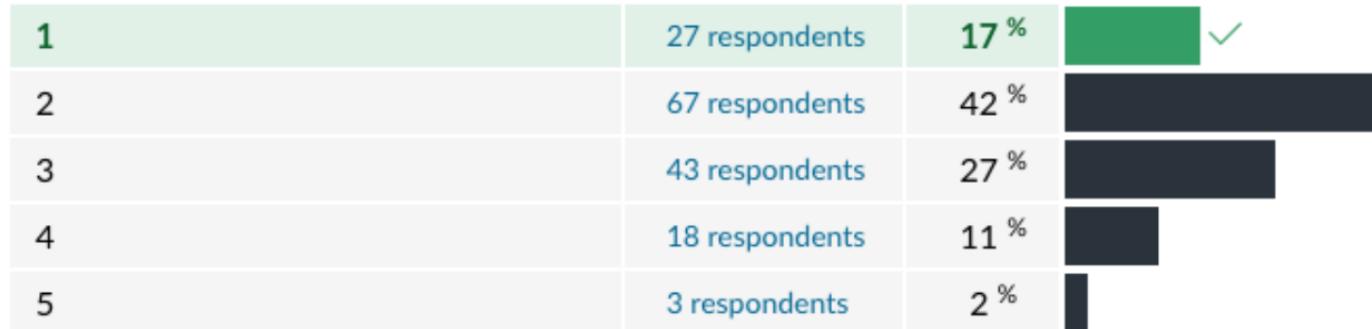


(1) Never heard about it (2) Know a little bit, (3) Manageable, (4) Familiar, (5) Very good at it.

# Topic Familiarity

UC San Diego

2. Advanced deep learning models: convolutional neural networks, recurrent neural networks, graph neural networks

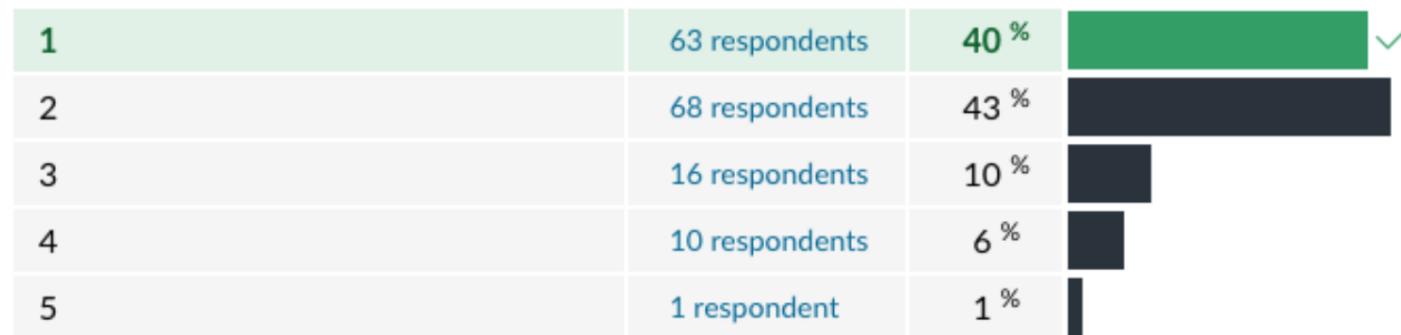


(1) Never heard about it (2) Know a little bit, (3) Manageable, (4) Familiar, (5) Very good at it.

# Topic Familiarity

UC San Diego

3. Reinforcement learning: Markov decision process, policy iteration, value iteration, Q-learning, policy gradient, actor-critic

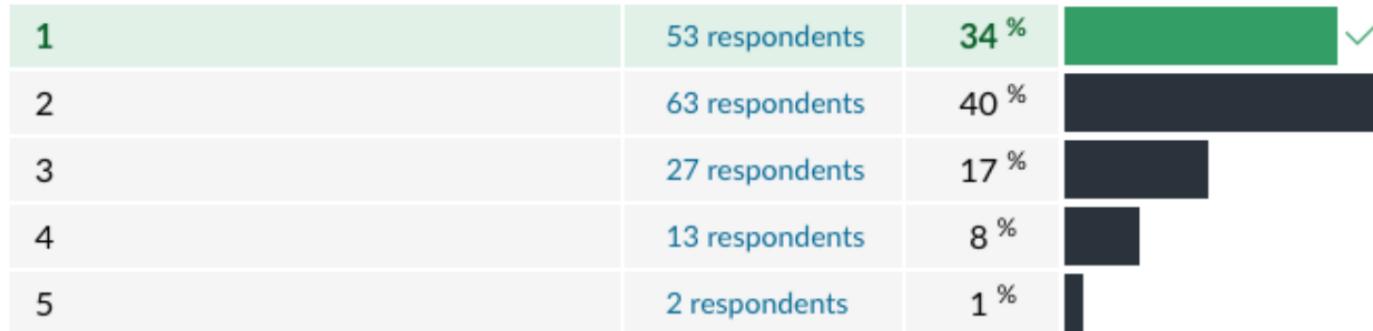


(1) Never heard about it (2) Know a little bit, (3) Manageable, (4) Familiar, (5) Very good at it.

# Topic Familiarity

UC San Diego

## 4. Physics-informed machine learning, partial differential equations

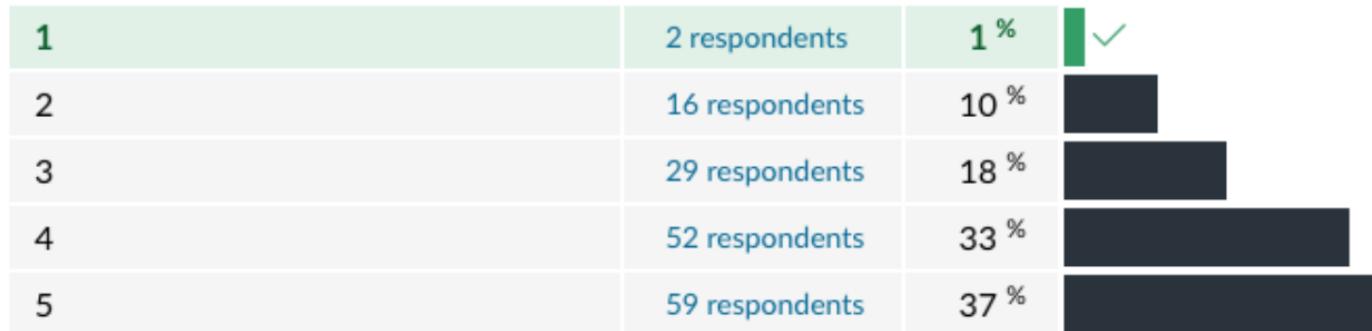


(1) Never heard about it (2) Know a little bit, (3) Manageable, (4) Familiar, (5) Very good at it.

# Topic Familiarity

UC San Diego

## 5. Coding in Python



(1) Never heard about it (2) Know a little bit, (3) Manageable, (4) Familiar, (5) Very good at it.

# Today's Lecture

UC San Diego

- Forward and Backward Computation of Logistic Regression (1-layer Neural Network)
- One Hidden Layer Neural Network *(or called 2-layer NN)*  
*# of hidden layer + output layer*
- Activation Functions
- Forward and Backward Computation of Neural Network with L Hidden Layers
- Universal Approximation Theorem & Why Deep?

# Forward and Backward Computation of Logistic Regression

# Forward Computation of Logistic Regression

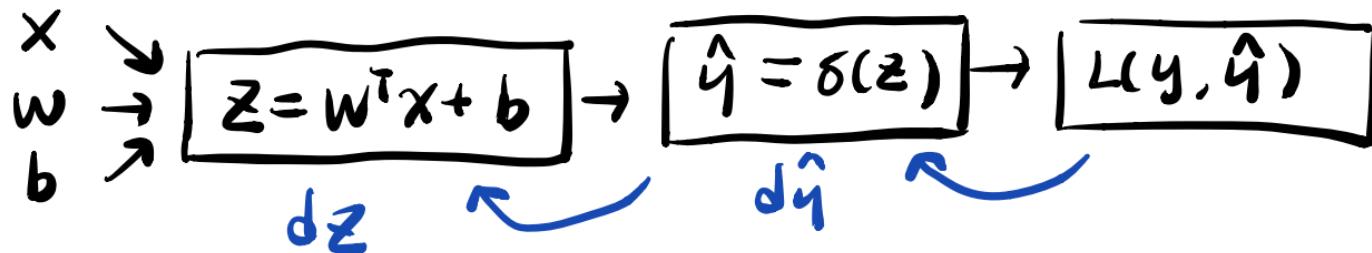
UC San Diego

Logistic regression with 1 sample  $(x, y)$ ,  $x \in \mathbb{R}^n$ ,  $y \in \{0, 1\}$

$$\hat{y} = \sigma(w^T x + b), z = w^T x + b$$

$$L(\hat{y}, y) = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$$

computation graph



# Backward Computation of Logistic Regression

UC San Diego

$$\text{"d}\hat{y} = \frac{\partial L}{\partial \hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

"chain rule"

$$\text{"d}z = \frac{\partial L}{\partial z} = \boxed{\frac{\partial L}{\partial \hat{y}}} \cdot \frac{\partial \hat{y}}{\partial z} = \left[ -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right] \frac{\sigma(z)[1-\sigma(z)]}{\hat{y}(1-\hat{y})}$$

$$= \left[ -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right] \hat{y}(1-\hat{y})$$

$$= -y(1-\hat{y}) + (1-y)\hat{y} = \hat{y} - y$$

# Backward Computation of Logistic Regression

UC San Diego

$$\begin{aligned} "dw" &= \frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w} \\ &= x dz \end{aligned}$$

Gradient update:

$$w \leftarrow w - \eta dw$$

$$b \leftarrow b - \eta db$$

$\eta$ : learning rate

★  $w \in \mathbb{R}^n$  is a vector, so  $\frac{\partial z}{\partial w}$  is

is a vector:  $\frac{\partial z}{\partial w} = \left( \frac{\partial z}{\partial w_1}, \dots, \frac{\partial z}{\partial w_n} \right)$  ↪ also called gradient of  $z$  w.r.t.  $w$

## Clarification on Notations

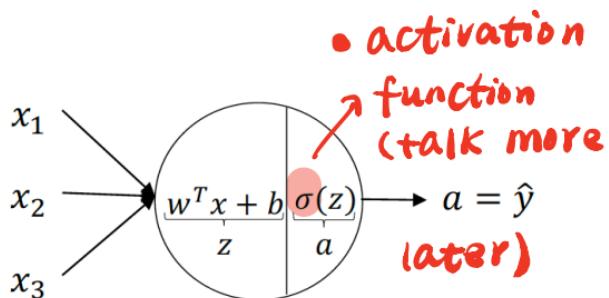
UC San Diego

- \* Both  $\frac{df}{dx}$ ,  $\frac{\partial f}{\partial x}$  are used to denote derivatives.
  - $\frac{df}{dx}$  is more commonly used for function of single variable.
  - If  $f(x_1, x_2, x_3)$  is a multi-variate function,  $\frac{\partial f}{\partial x_1}$  is partial derivative of  $f$  w.r.t.  $x_1$ , given  $\underline{x_2, x_3}$  fixed
  - $\frac{\partial f}{\partial x} = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3} \right)$ ; also called gradient of  $f$
- \* We introduce a short-hand notation in this course
- $\boxed{dx} = \frac{df}{dx}$  (single-variate);  $dx = \underline{\frac{\partial f}{\partial x}} = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3} \right) \leftarrow$  multi-variate

One Hidden Layer Neural  
Network

# One Hidden Layer Neural Network

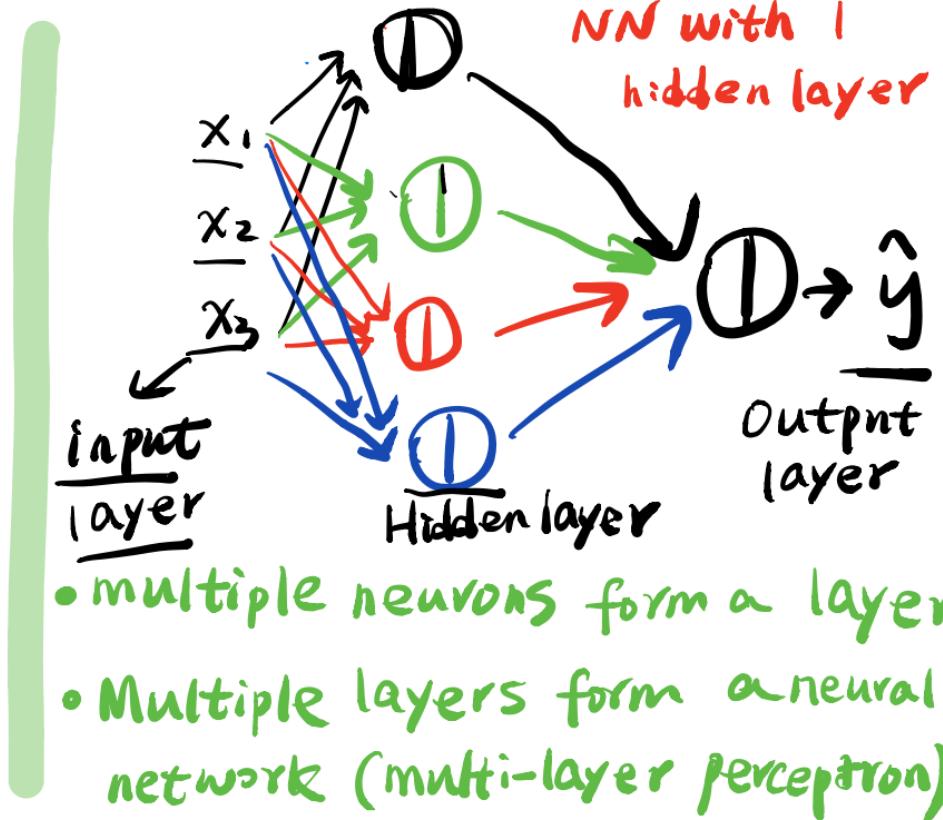
UC San Diego



$$z = w^T x + b$$

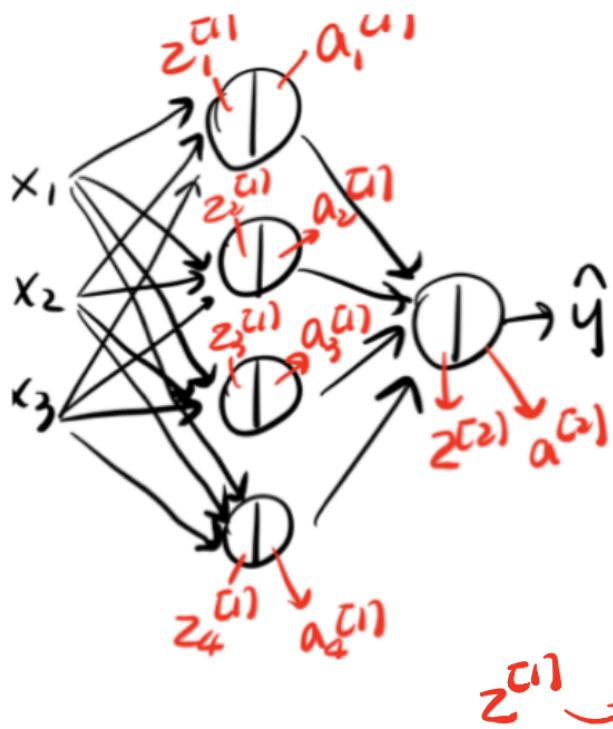
$$a = \sigma(z)$$

⇒ building block for neural network,  
neuron (perceptron)



# Forward Computation

UC San Diego



$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)} \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

$$z_2^{(1)} = w_2^{(1)T} x + b_2^{(1)} \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

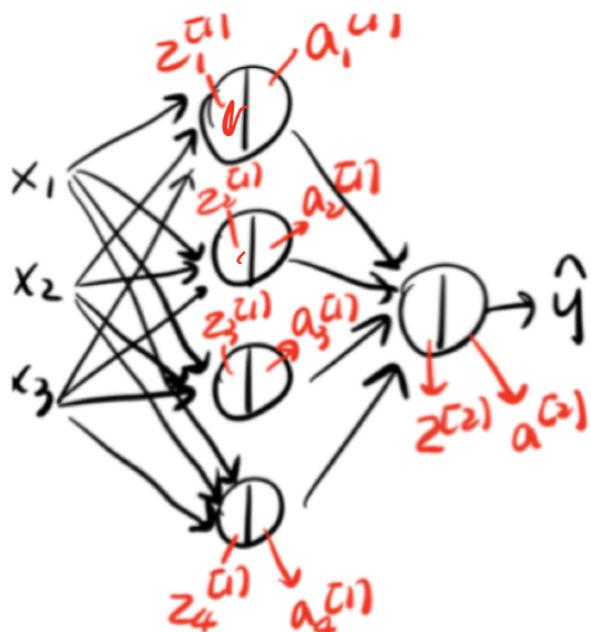
$$z_3^{(1)} = w_3^{(1)T} x + b_3^{(1)} \quad a_3^{(1)} = \sigma(z_3^{(1)})$$

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)} \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

$$\begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix} = \begin{bmatrix} -w_1^{(1)T} & - \\ -w_2^{(1)T} & - \\ -w_3^{(1)T} & - \\ -w_4^{(1)T} & - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} b^{(1)}$$

# Forward Computation

UC San Diego



Given input  $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

$$z^{(1)} = W^{(1)} x + b^{(1)}, a^{(1)} = \sigma(z^{(1)})$$

(4x1)      (4x3) x (3x1)      (4x1)

$$z^{(2)} = W^{(2)} a^{(1)} + b^{(2)}$$

(1x1)      (1x4)      (4x1)      (1x1)

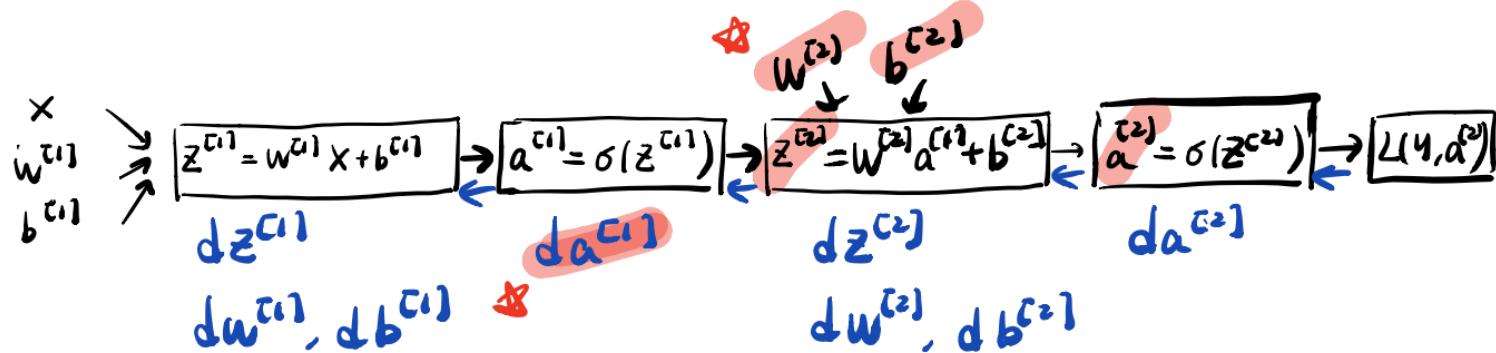
$$\hat{y} = a^{(2)} = \sigma(z^{(2)})$$

# Backward Computation

UC San Diego

- Backpropagation: gradient descent for neural networks

\* next step



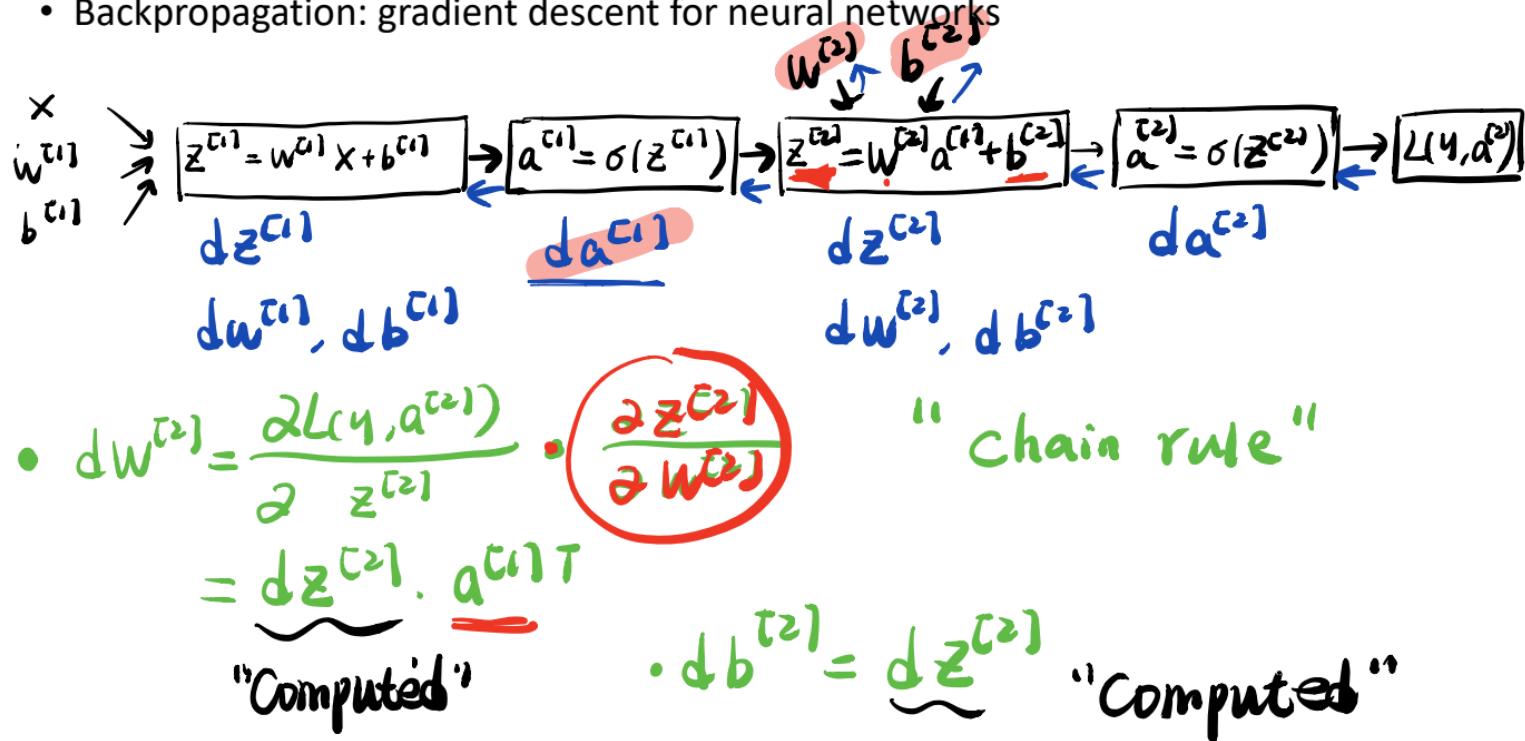
- Suppose:  $L(y, a^{(2)}) = -[y \log a^{(2)} + (1-y) \log(1-a^{(2)})]$

- $da^{(2)} = \frac{\partial L(y, a^{(2)})}{\partial a^{(2)}} = -\frac{y}{a^{(2)}} + \frac{1-y}{1-a^{(2)}}, \quad dz^{(2)} = a^{(2)} - y$

# Backward Computation

UC San Diego

- Backpropagation: gradient descent for neural networks

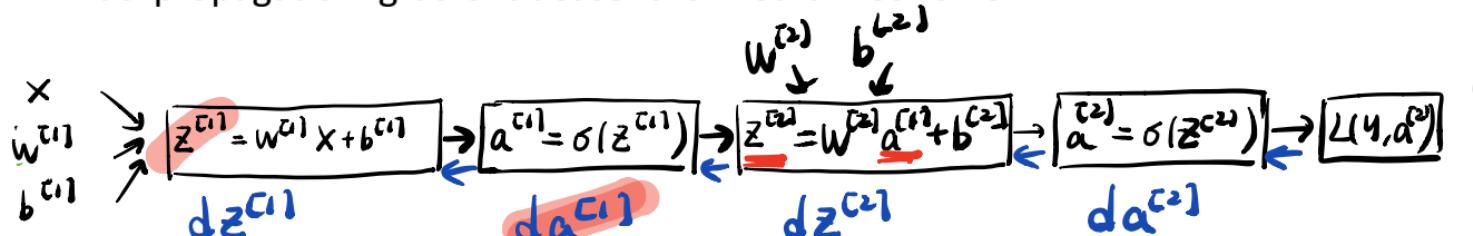


# Backward Computation

UC San Diego

★ next step

- Backpropagation: gradient descent for neural networks



$$\cdot da^{[1]} = W^{[2]T} \cdot dZ^{[2]}$$

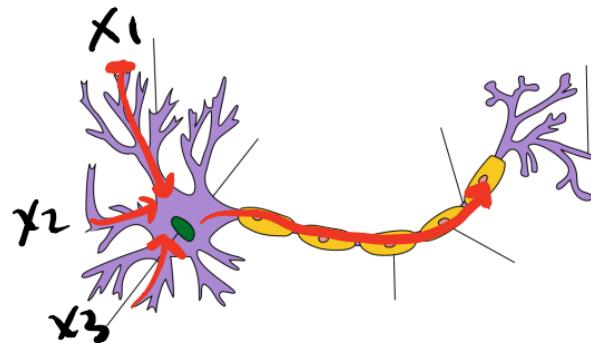
$$\cdot dZ^{[1]} = \underbrace{\frac{da^{[1]}}{(n \times 1)} * \sigma'(z^{[1]})}_{(n \times 1)}, *: \text{element-wise multiplication}$$

$$\cdot dW^{[1]} = dZ^{[1]} \cdot x^T, db^{[1]} = dZ^{[1]}$$

# Activation Functions

# What is an Activation Function?

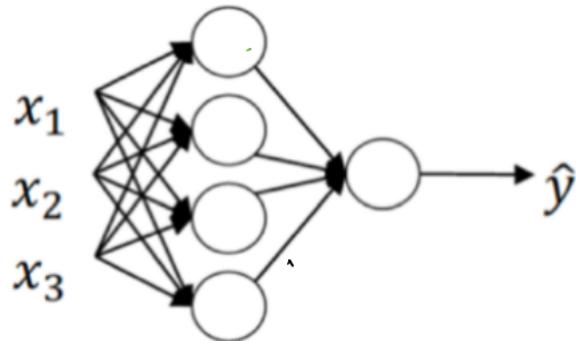
UC San Diego



- \* Inspired by neuron activities in the brain
- \* In "Neural Network", it refers to
  - how weighted sum is transformed to (often nonlinear)  **output**

# Why Do We Need Nonlinear Activation Functions?

UC San Diego



No activation

$$a^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

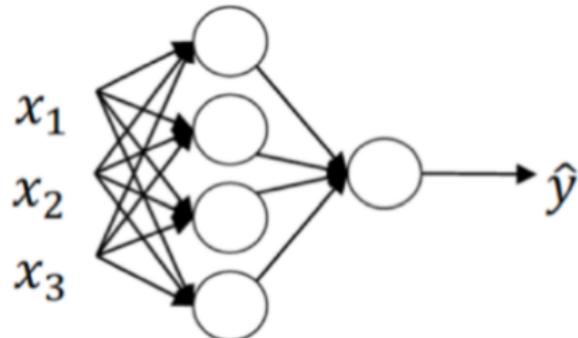
$$= w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]}$$

$$= \underbrace{w^{[2]}w^{[1]}}_{w'} x + \underbrace{w^{[2]}b^{[1]} + b^{[2]}}_{b'}$$

- Without activation func,  
Can only represent  
linear relationship.

# Why Do We Need Nonlinear Activation Functions?

UC San Diego



With activation

$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

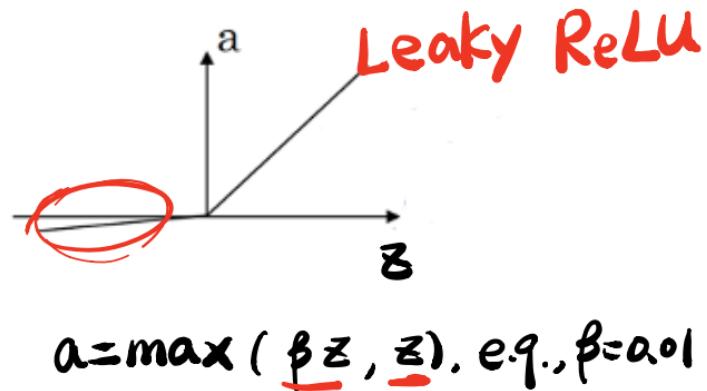
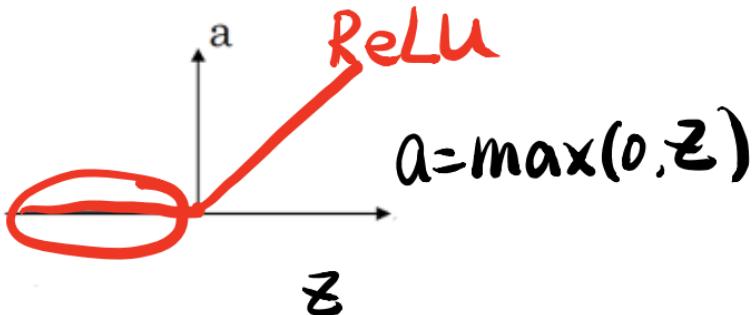
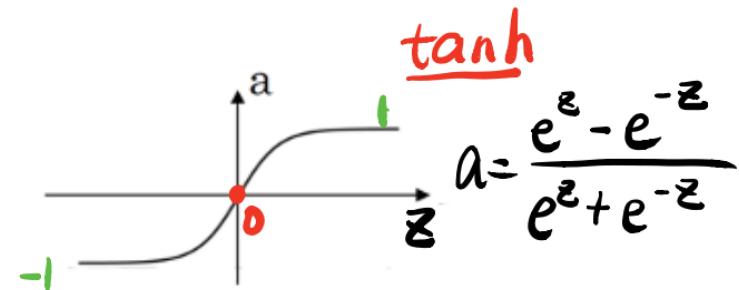
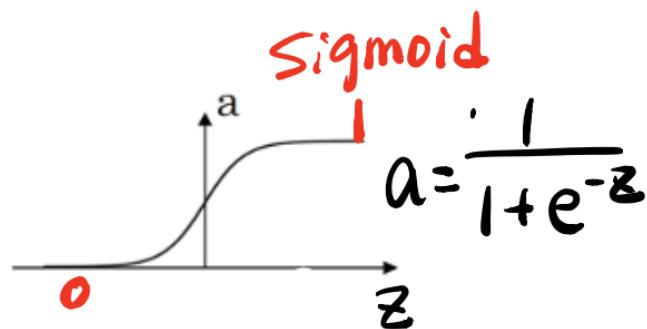
$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$= w^{[2]}\sigma(w^{[1]}x + b^{[1]}) + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

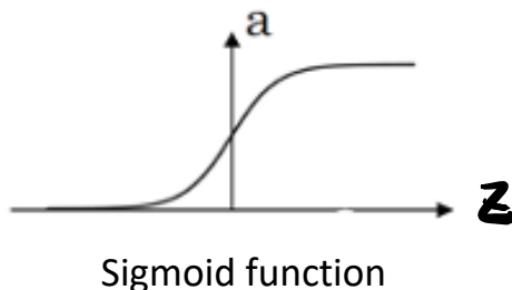
# Different Activation Functions

UC San Diego



# Derivatives of Activation Functions

UC San Diego



$$\Rightarrow \delta(z) = \frac{1}{1 + e^{-z}}$$

$$\Rightarrow \delta'(z) = \delta(z)[1 - \delta(z)]$$

- $z=100, \delta(z) \approx 1$

$$\delta'(z) \approx 1 \cdot (1 - 1) = 0$$

- $z=-100, \delta(z) \approx 0$

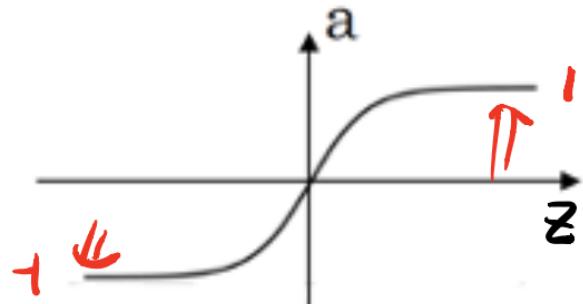
$$\delta'(z) \approx 0[1 - 0] = 0$$

- $z=0, \delta(z)=\frac{1}{2}, \delta'(z)=\frac{1}{4}$

# Derivatives of Activation Functions

UC San Diego

Tanh function



$$\delta(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$\begin{aligned}\delta'(z) &= \frac{d}{dz} \delta(z) \\ &= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\ &= 1 - \tanh^2(z)\end{aligned}$$

- $z=100, \tanh(z) \approx 1, \delta'(z) \approx 0$

Tanh

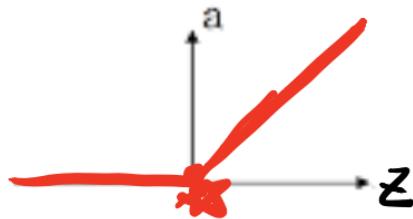
- $z=-100, \tanh(z) \approx -1, \delta'(z) \approx 0$

- $z=0, \tanh(0)=0, \delta'(z)=1$

# Derivatives of Activation Functions

UC San Diego

ReLU

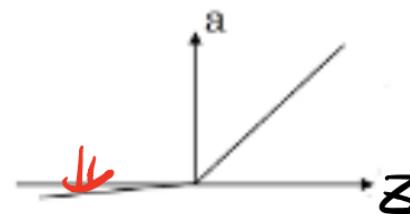


$$\sigma(z) = \max(0, z)$$

$$\sigma'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

$\circled{z=0}$  not differentiable; (empirically), Set it as 0 or 1

Leaky-Relu



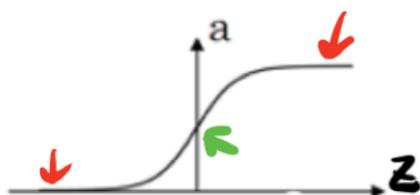
$$\sigma(z) = \max(0.01z, z)$$

$$\sigma'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Set it as 0 or 1

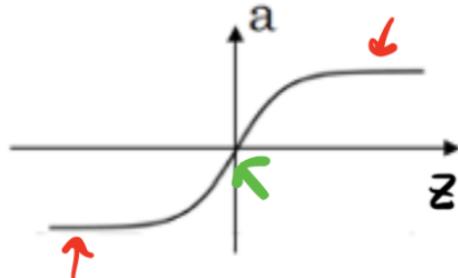
# Pros and Cons of Different Activations

UC San Diego

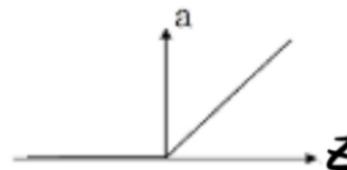


Sigmoid function

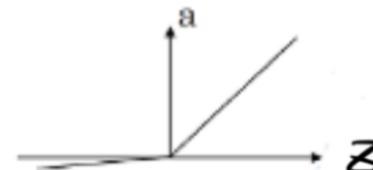
Tanh function



ReLU



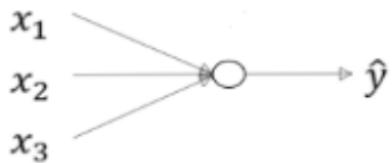
Leaky-Relu



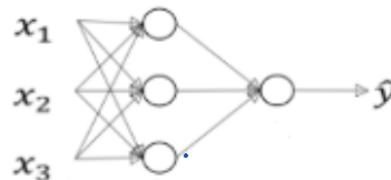
# Forward and Backward Computation of Neural Network

# Neural Networks with Different Number of Layers

UC San Diego



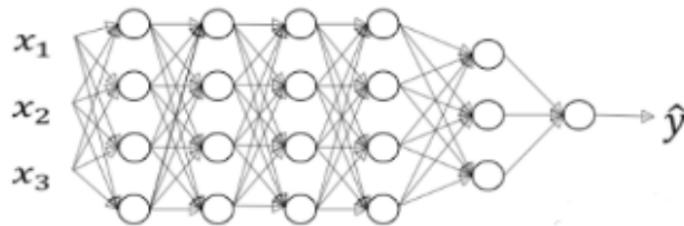
logistic regression



1 hidden layer



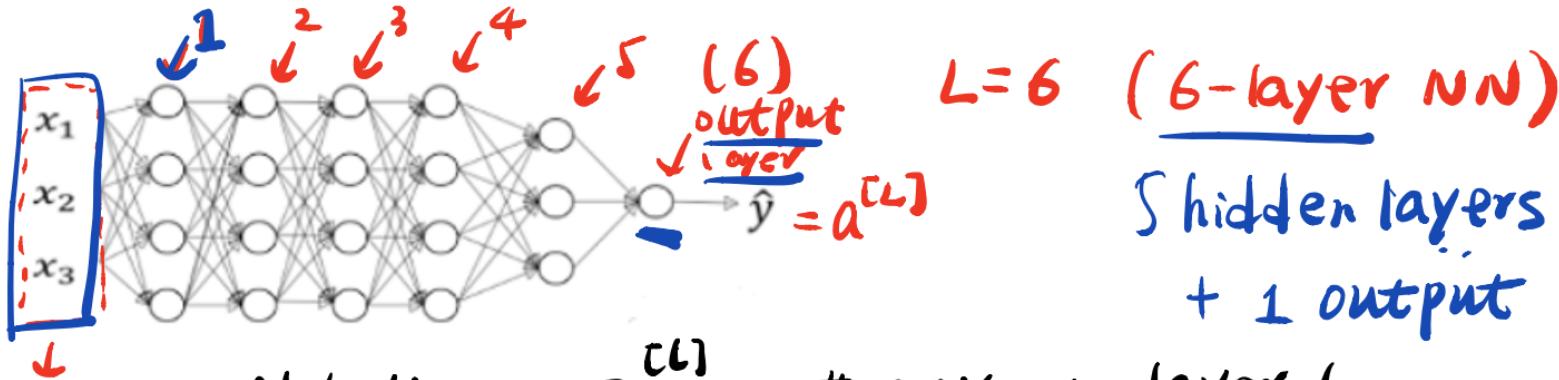
2 hidden layers



5 hidden layers

# Notation

UC San Diego



5 hidden layers  
+ 1 output layer

input layer      Notations:  $n^{[0]} = n_x = 3$

$$n^{[1]} = 4, n^{[2]} = 4$$

$$n^{[3]} = 4, n^{[4]} = 4$$

$$n^{[5]} = 3, n^{[6]} = n_y = 1$$

$n^{[L]}$  = # units in layer L

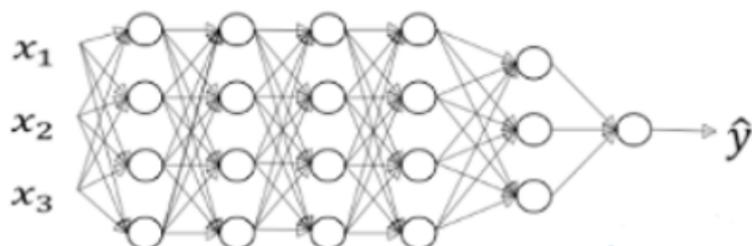
$a^{[L]}$  = activations in layer L

$a^{[L]} = \sigma^{[L]}(z^{[L]})$  weights & bias  
in layer L

$z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]}, a^{[0]} = x$

# Forward Computation

UC San Diego



$$W^{[L]} : (n^{[L]}, n^{[L-1]})$$
$$b^{[L]} : (n^{[L]}, 1)$$

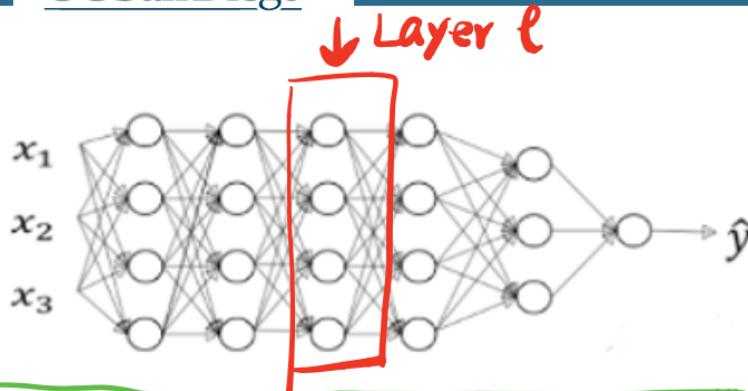
weights & bias  
in layer L

✓  $z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$ ,  $a^{[0]} = x$

✓  $a^{[L]} = \sigma^{[L]}(z^{[L]})$

# Backward Computation

UC San Diego



Layer  $\ell$

forward

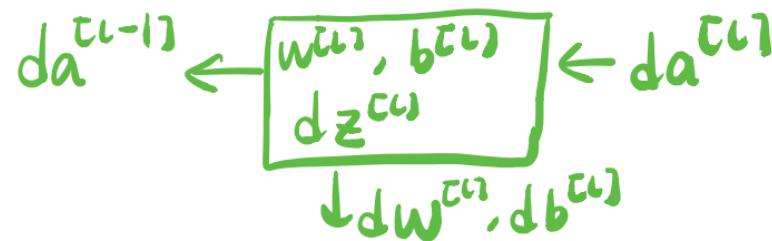
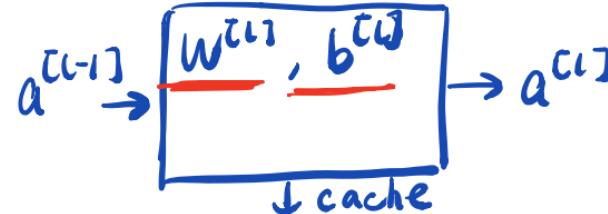
$$\begin{aligned} z^{[l]} &= W^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} &= \sigma(z^{[l]}) \end{aligned}$$

backward elementwise

$$dz^{[l]} = da^{[l]} * \sigma'(z^{[l]}), (z^{[l]})$$

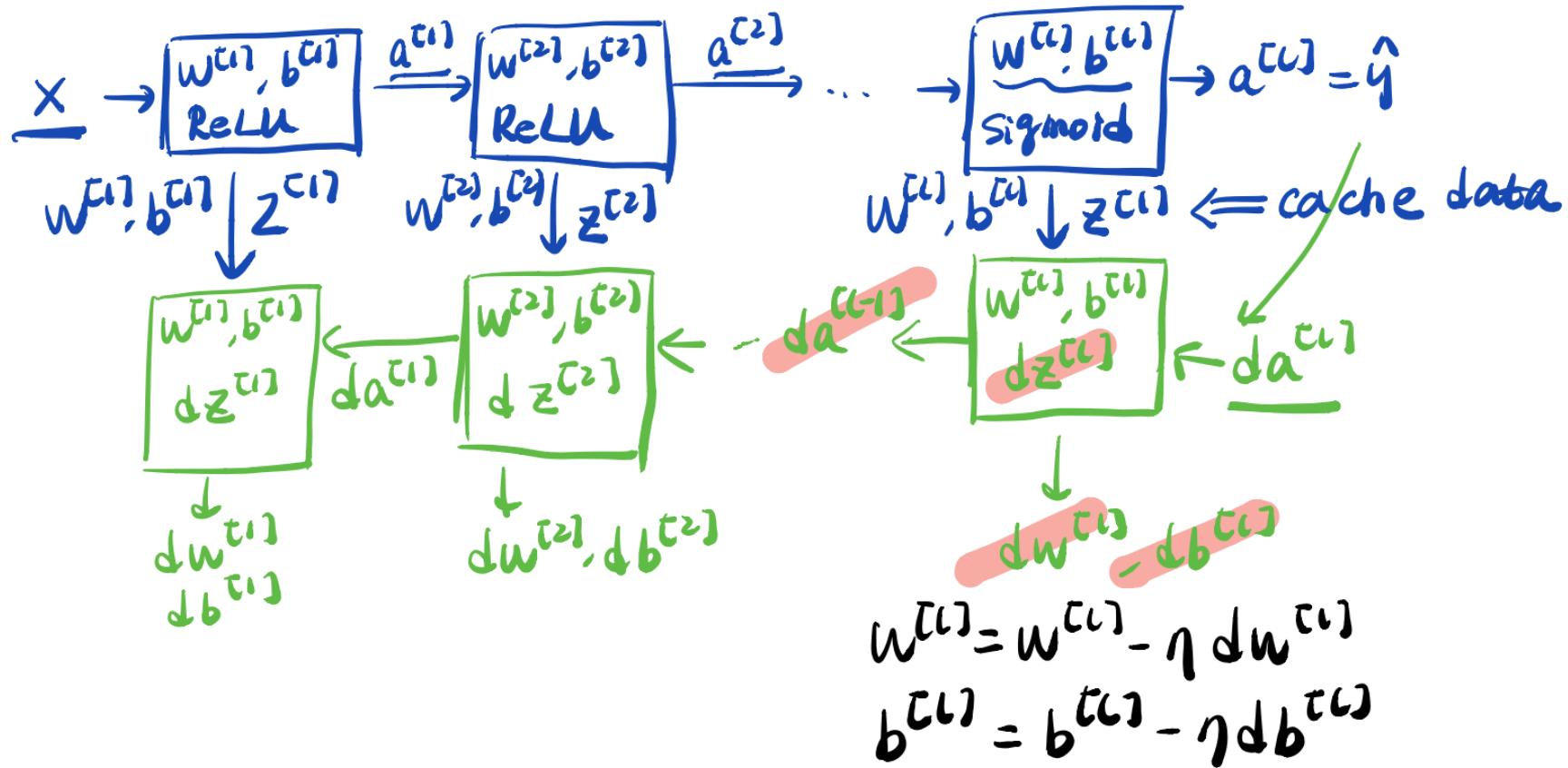
$$dW^{[l]} = dz^{[l]} \cdot a^{[l-1]T}, db^{[l]} = dz^{[l]}$$

$$da^{[l-1]} = W^{[l]T} dz^{[l]}$$



# Forward and Backward Computation

UC San Diego



# Universal Approximation

## Theorem & Why Deep?

# Universal Function Approximation

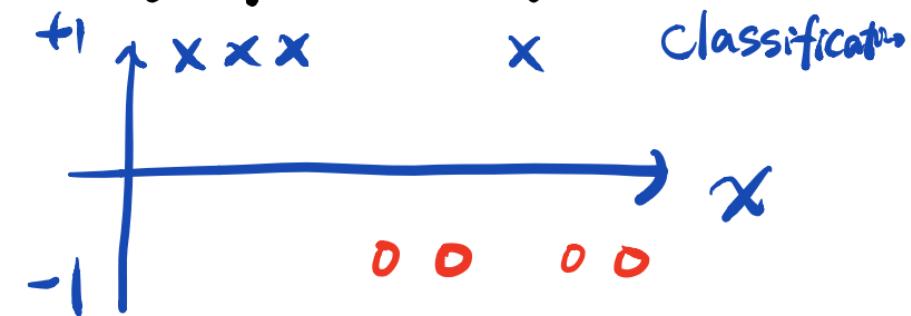
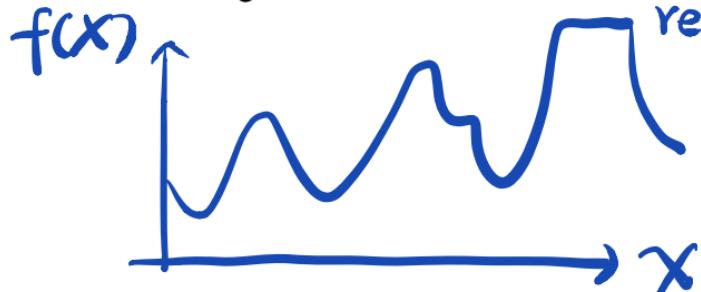
UC San Diego

- Informal theorems:

"a wide enough 1 hidden layer NN can fit any function"

- Empirical findings:

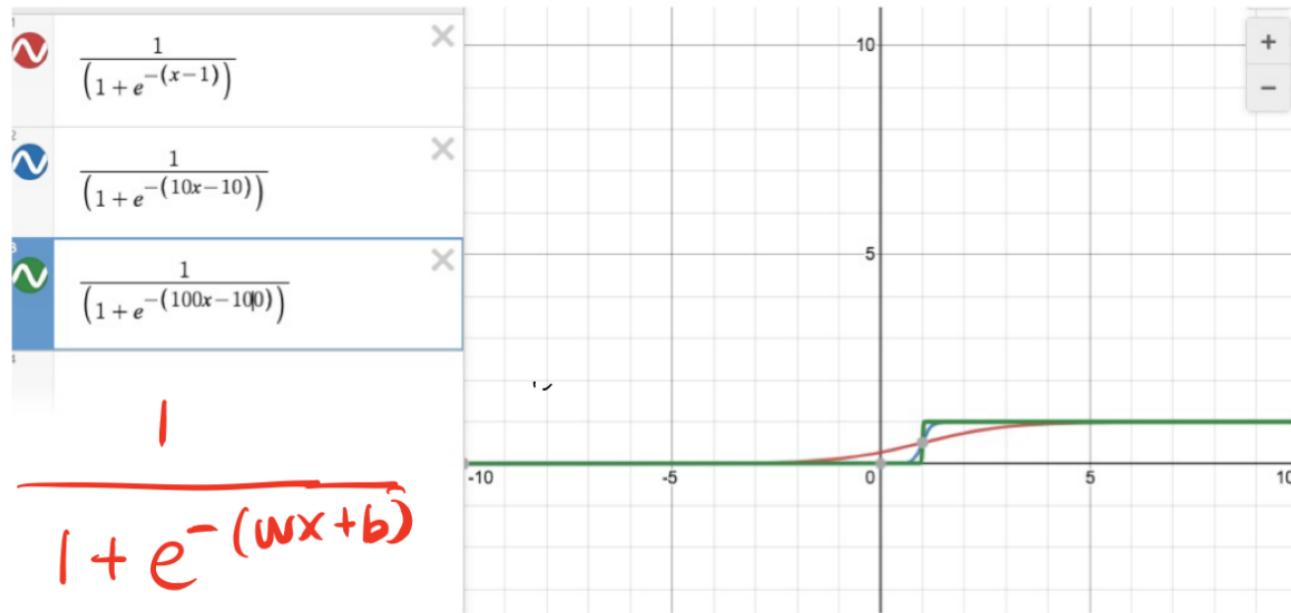
"my NN can fit any labelling of the training data"



# Universal Function Approximation

UC San Diego

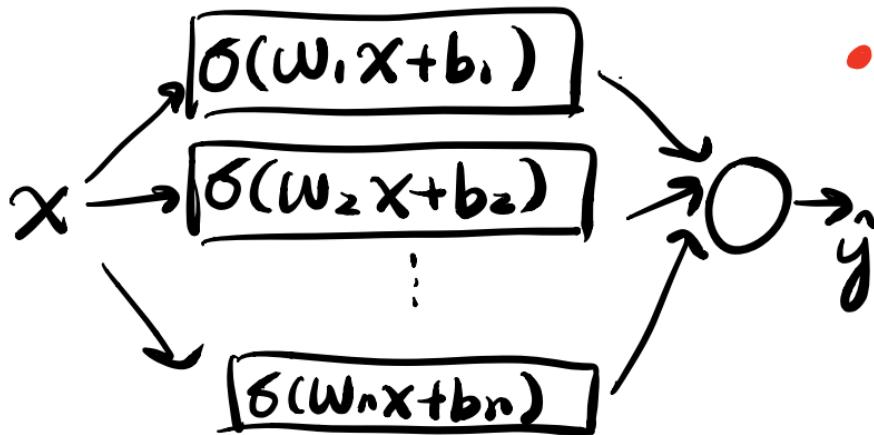
Consider a 1d input  $x \in \mathbb{R}$  & sigmoid activation



# Universal Function Approximation

UCSanDiego

- When  $w$  is very large, Sigmoid functions look like step functions



- Sum of step functions can approximate any function shape:

# Universal Function Approximation

UC San Diego

deep networks are *universal function approximators* (Hornik, 1991)

→ with enough units & layers, can approximate *any* function

## Classic Results

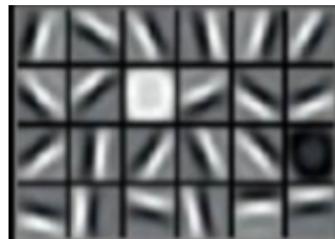
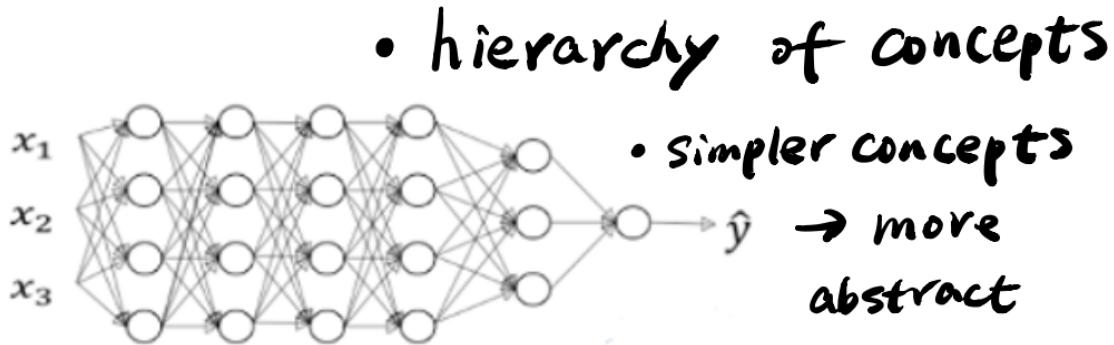
**Theorem:** For one layer neural networks, the error between the estimated network and a target function  $f$  is shown to be bounded by

$$O\left(\frac{C_f^2}{n}\right) + O\left(\frac{nd}{N} \log N\right)$$

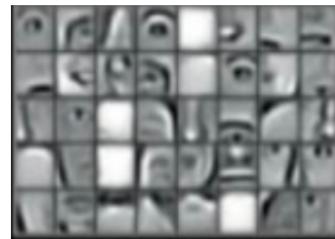
here  $n$  is the number of nodes,  $d$  is the input dimension of the function,  $N$  is the number of training data, and  $C_f$  is the first absolute moment of the Fourier magnitude distribution of  $f$ .

# Why Deep Representation?

UC San Diego



Pixels and  
edges



more complex  
shapes: eyes & mouths



faces