

# Artificial Intelligence and Machine Learning

## Neural Networks

# Lecture Outline

- Logistic Regression Review
- Neural Networks
  - Forward pass
  - Backward pass

# Review: Logistic Regression

Binary

- Sigmoid:

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

- BCE

Multiclass

- Softmax

$$\sigma(z) = \frac{\exp(z)}{\sum_{i=1}^n \exp(z_i)}$$

- Cross ent. loss

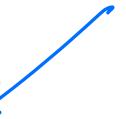


classification

for discrete value

Binary

Multiclass





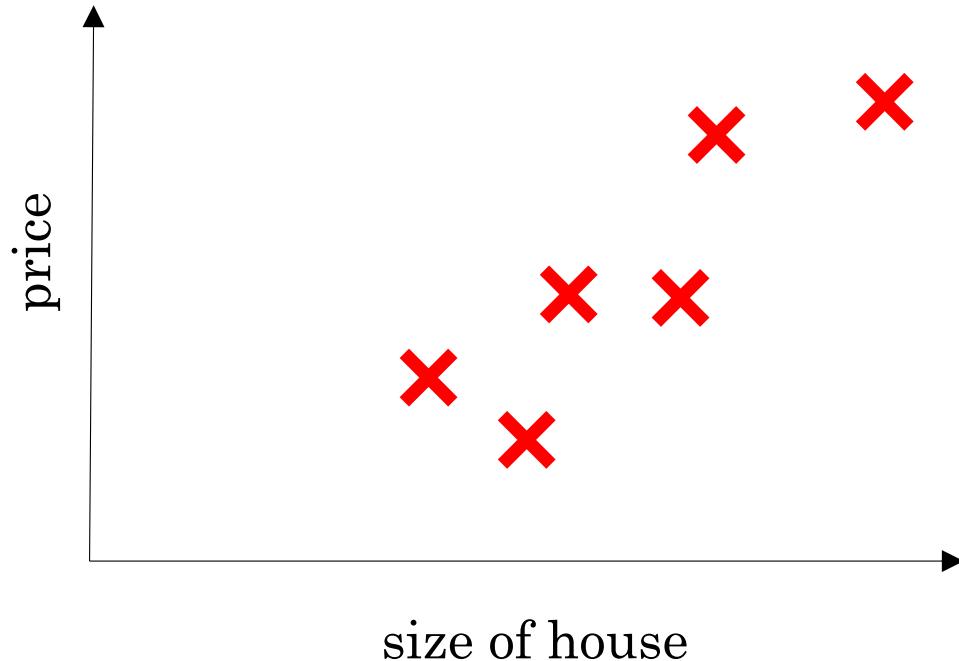
deeplearning.ai

# Introduction to Deep Learning

---

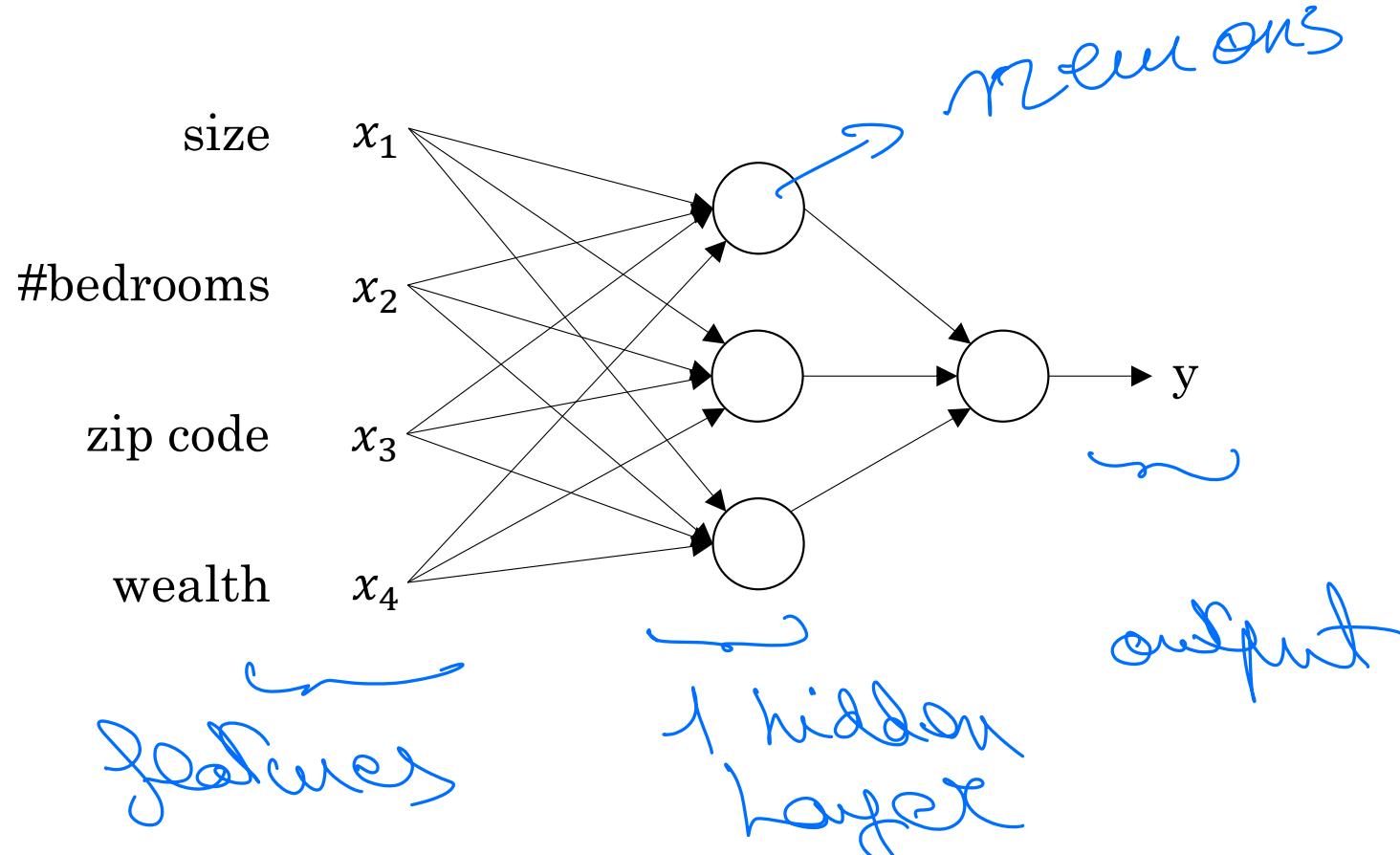
What is a  
**Neural Network?**  
*model inspired by the brain*

# Housing Price Prediction



# Housing Price Prediction

# Housing Price Prediction





deeplearning.ai

# Introduction to Deep Learning

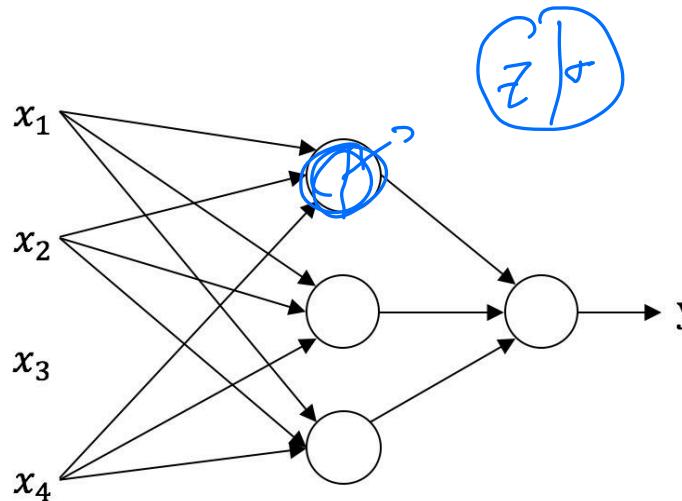
---

## Supervised Learning with Neural Networks

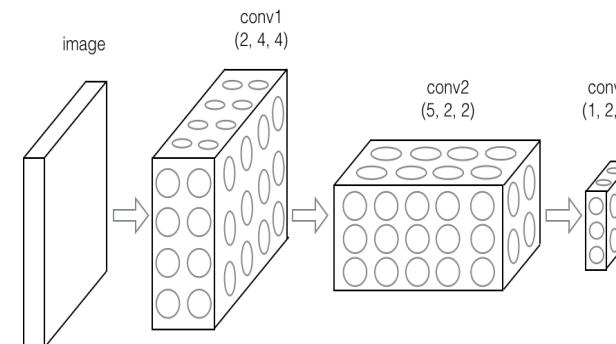
# Supervised Learning

Input(x)	Output (y)	Application
Home features	Price	Real Estate
Ad, user info	Click on ad? (0/1)	Online Advertising
Image	Object (1,...,1000)	Photo tagging
Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Image, Radar info	Position of other cars	Autonomous driving

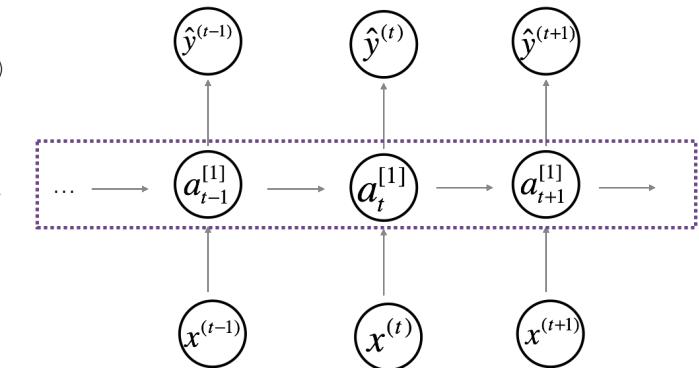
# Neural Network examples



Standard NN



Convolutional NN



Recurrent NN



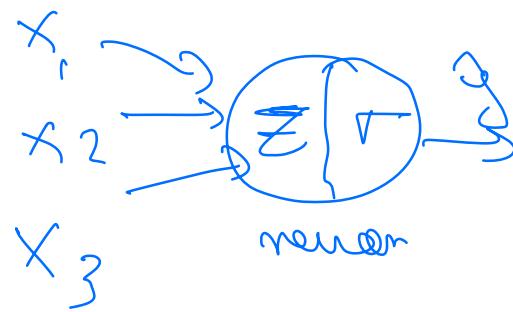
deeplearning.ai

# One hidden layer Neural Network

---

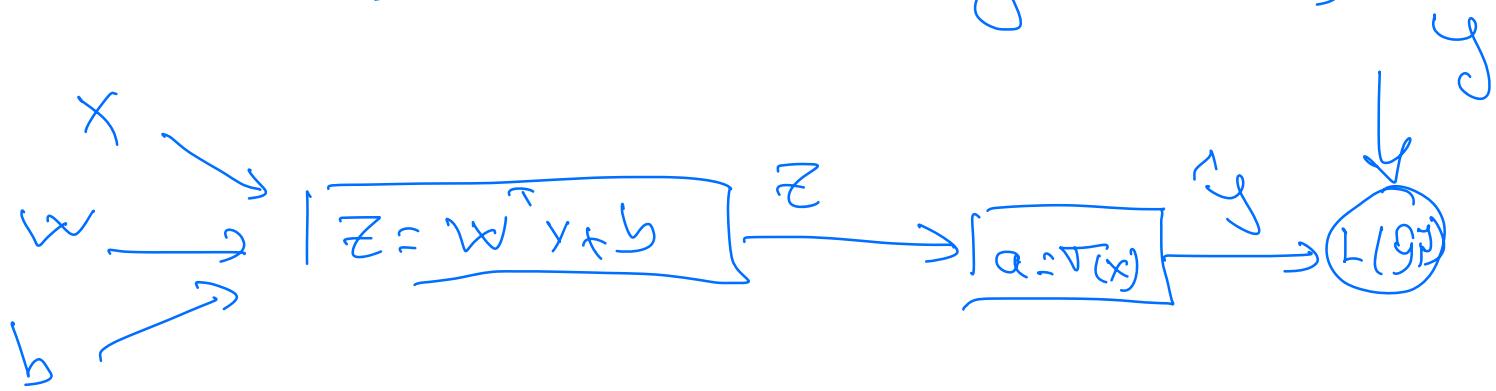
# Neural Networks Overview

\* What inside each neuron?



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

$$\hat{y} = \sigma(Z)$$



\* For each hidden layer or output:

$$Z^{(1)} = w^{(1)T} x + b^{(1)}$$

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

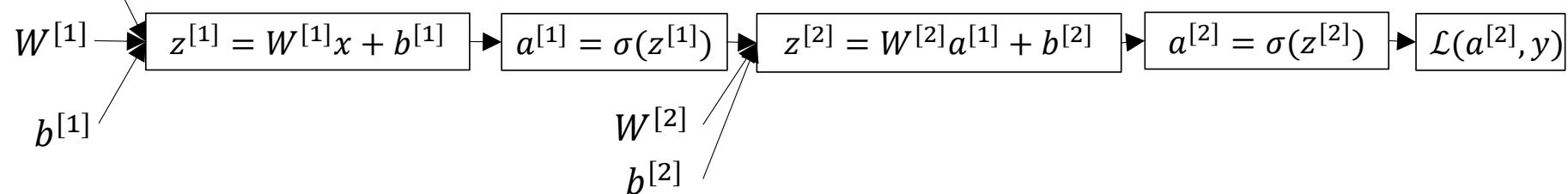
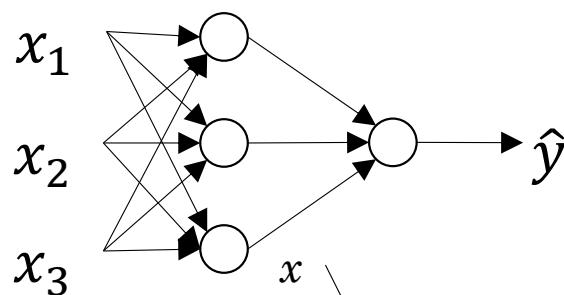
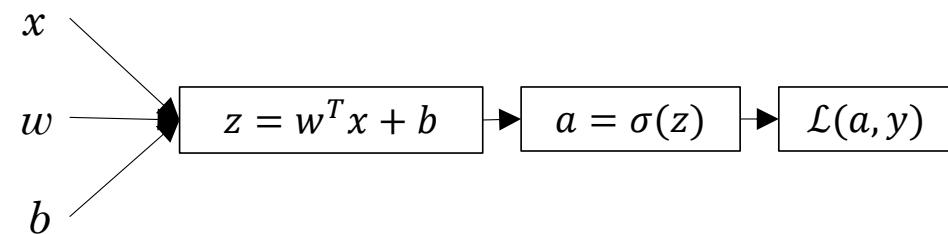
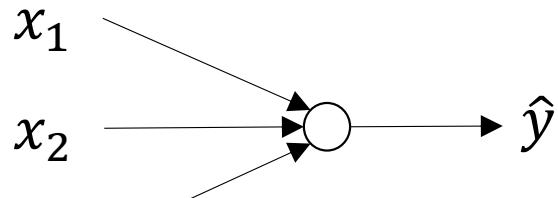
$$Z^{(2)} = w^{(2)T} a^{(1)} + b^{(2)}$$

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

$$Z^{(3)} = w^{(3)T} a^{(2)} + b^{(3)}$$

$$a^{(3)} = \sigma[Z^{(3)}]$$

# What is a Neural Network?





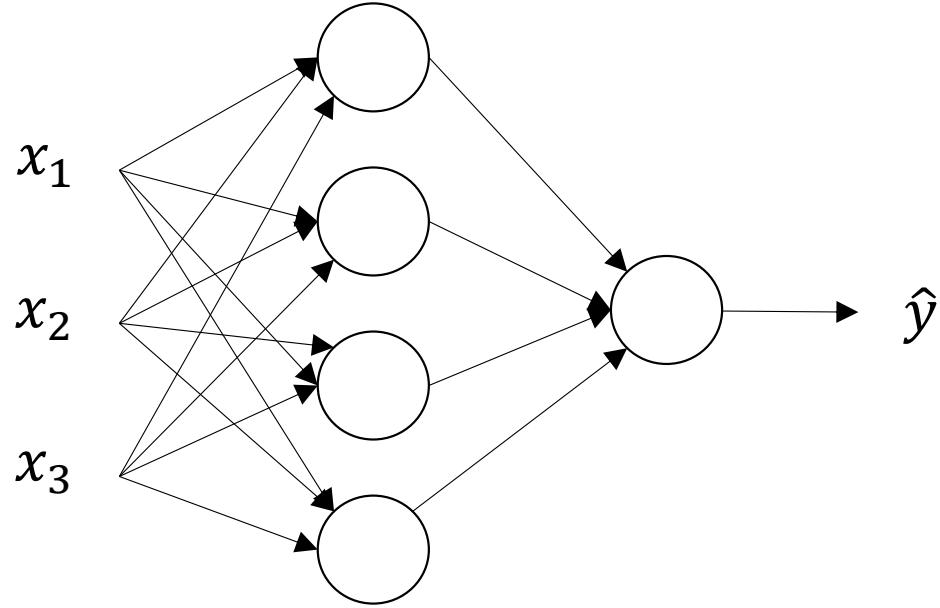
deeplearning.ai

# One hidden layer Neural Network

---

# Neural Network Representation

# Neural Network Representation





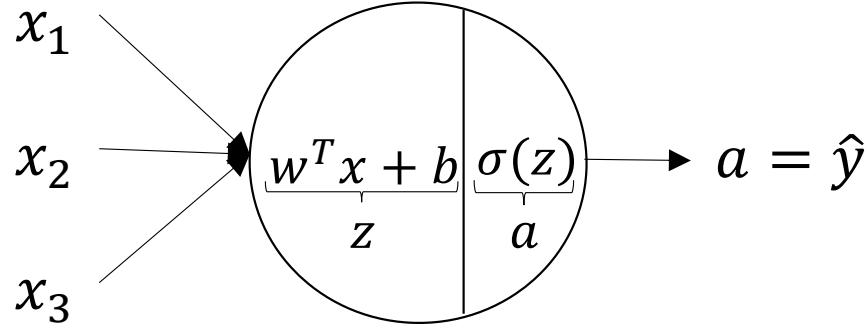
deeplearning.ai

# One hidden layer Neural Network

---

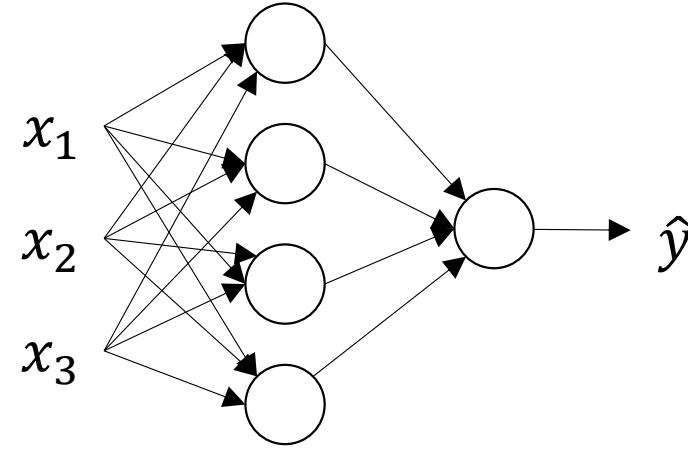
## Computing a Neural Network's Output

# Neural Network Representation

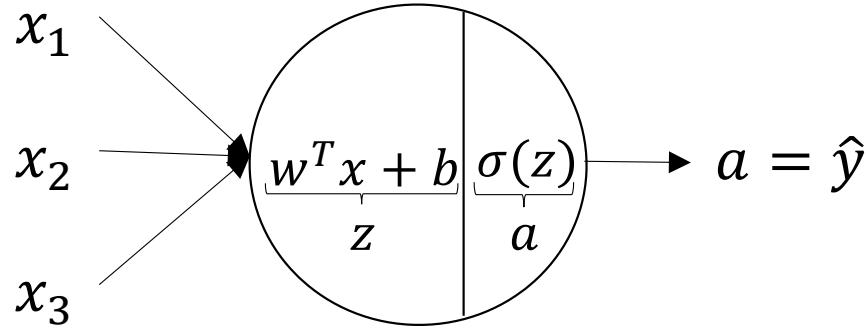


$$z = w^T x + b$$

$$a = \sigma(z)$$

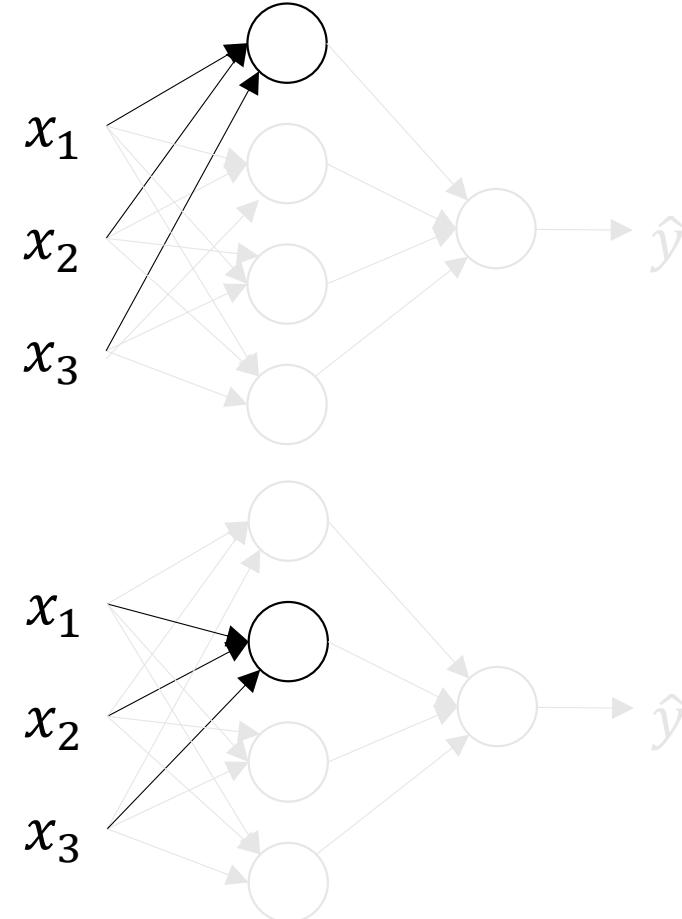


# Neural Network Representation

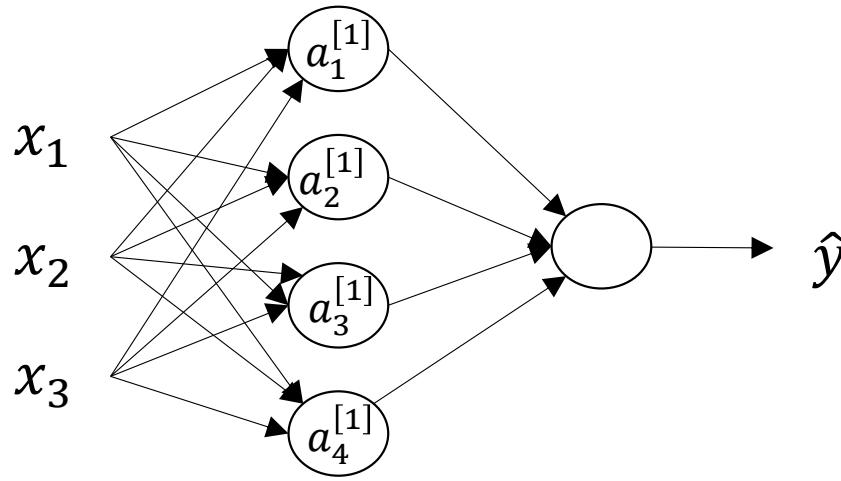


$$z = w^T x + b$$

$$a = \sigma(z)$$



# Neural Network Representation



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

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

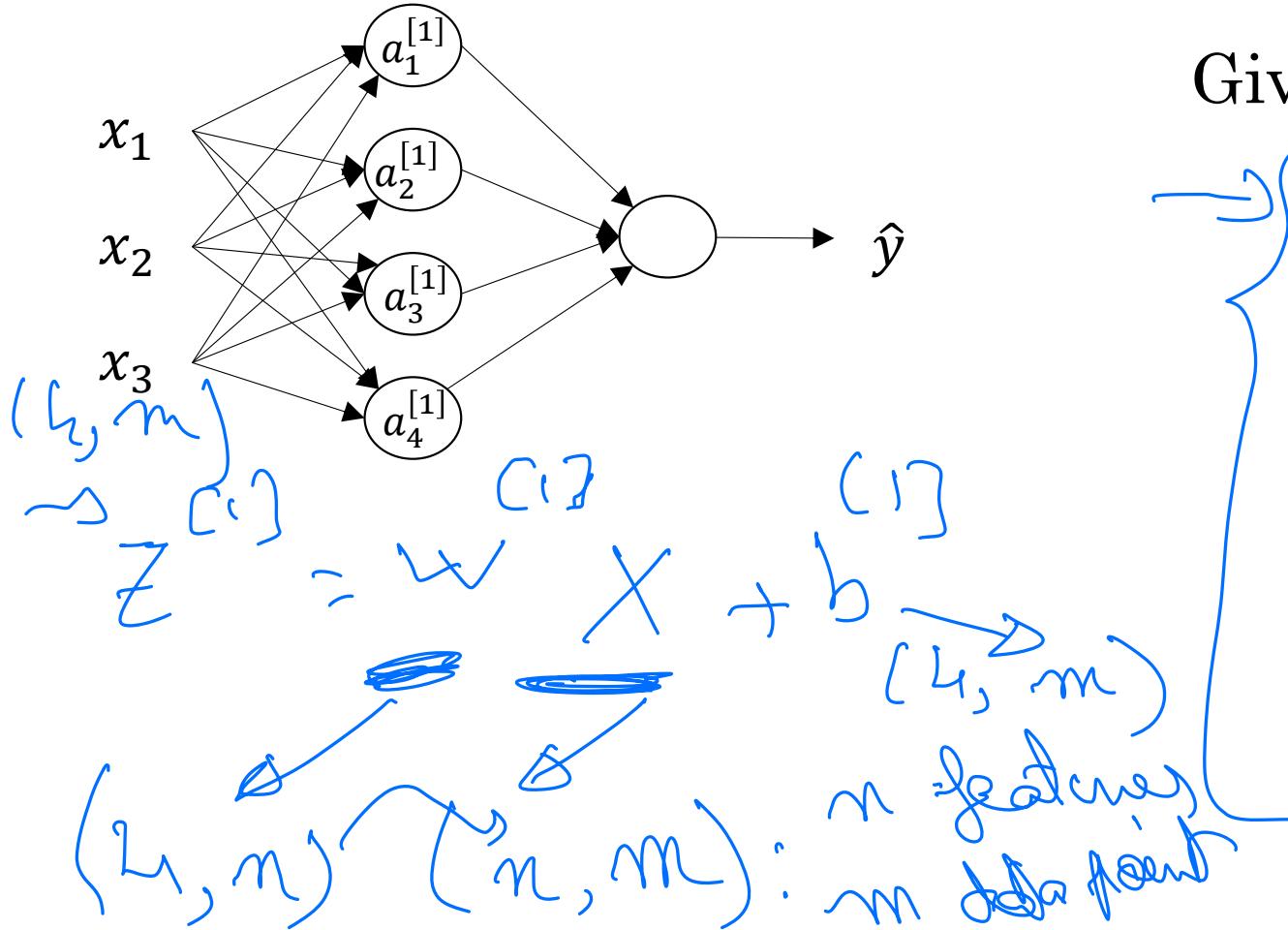
$$z_3^{[1]} = w_3^{[1]T} \underline{x} + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

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

$\Rightarrow$

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

# Neural Network Representation learning

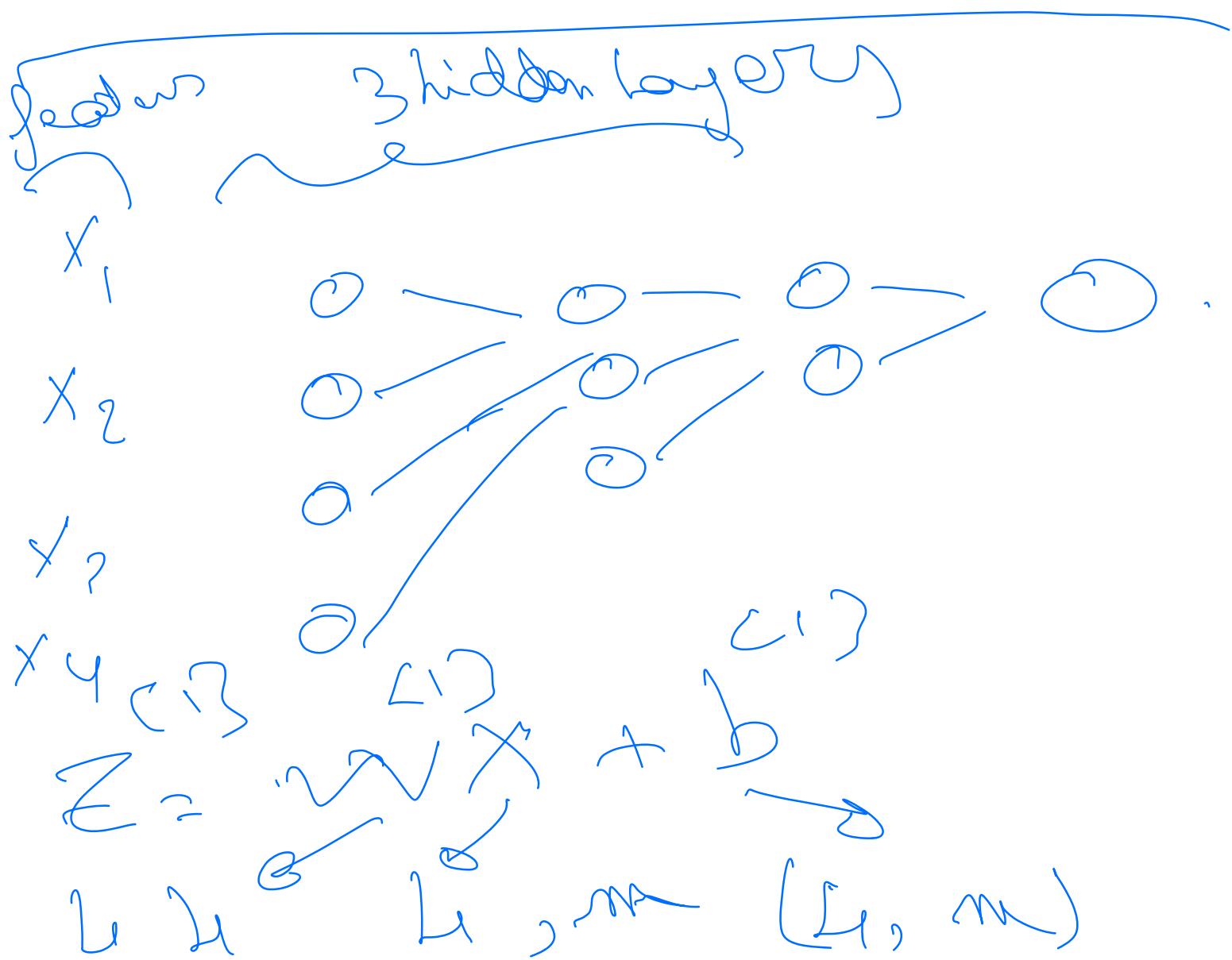


Given input  $x$ :

$a^{(1)} = \tau(z^{(0)}) \rightarrow$  some size  
 $(4, m)$ .

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

$\left\{ \begin{matrix} (4, 4) & (4, m) \\ (4, m) \end{matrix} \right\} \quad (1, m)$



$$Z^{(2)} = w^{(2)(1)} \cancel{a} + b^{(2)}$$

$T(3, 2)$

$$Z = \overset{(1)}{\circlearrowleft} x + b$$

$(3, 1) \rightarrow \begin{matrix} \nearrow \\ \searrow \end{matrix}$

$\begin{matrix} \nearrow \\ \searrow \end{matrix} \rightarrow \begin{matrix} \nearrow \\ \searrow \end{matrix}$

$L_T$



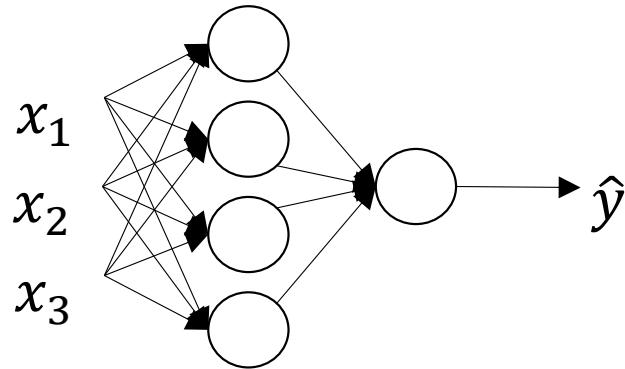
deeplearning.ai

# One hidden layer Neural Network

---

## Vectorizing across multiple examples

# Vectorizing across multiple examples



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

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

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

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

# Vectorizing across multiple examples

for i = 1 to m:

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

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

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

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



deeplearning.ai

# One hidden layer Neural Network

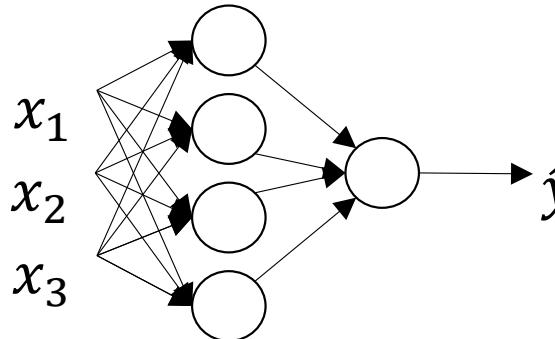
---

## Explanation for vectorized implementation

# Justification for vectorized implementation



# Recap of vectorizing across multiple examples



$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & | \\ a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & & | \end{bmatrix}$$

for i = 1 to m

$$z^{[1](i)} = W^{[1]}x^{(i)} + b^{[1]}$$

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

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

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

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$



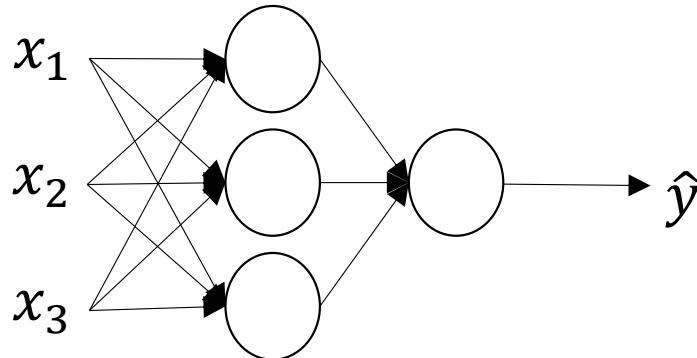
deeplearning.ai

# One hidden layer Neural Network

---

## Activation functions

# Activation functions



Given  $x$ :

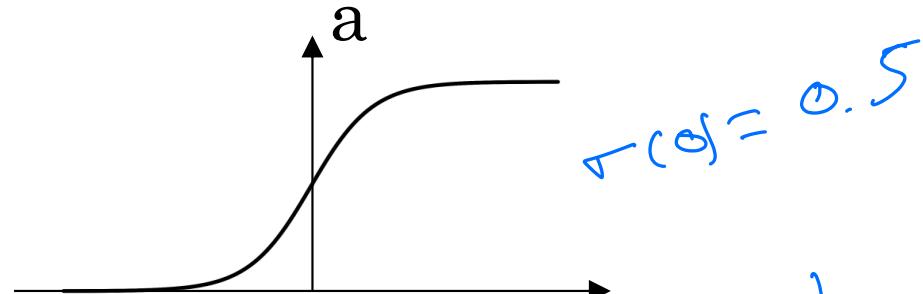
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

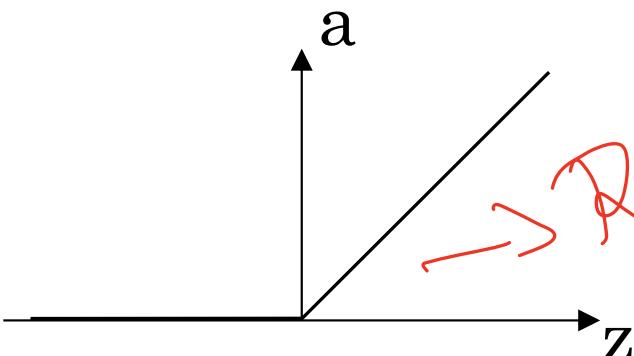
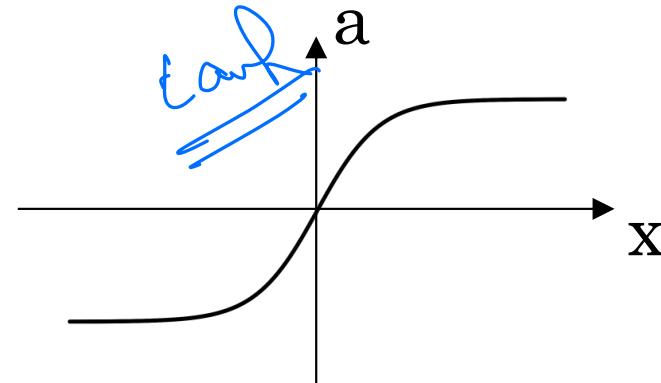
$$a^{[2]} = \sigma(z^{[2]})$$

# Pros and cons of activation functions



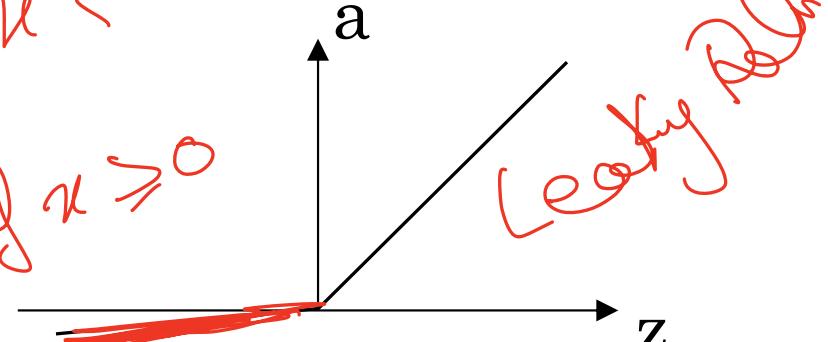
sigmoid: 
$$a = \frac{1}{1 + e^{-z}}$$

$$\frac{1}{1 + e^{-z}} = \frac{1}{2}$$



ReLU :  $\max(0, z)$

$a = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$



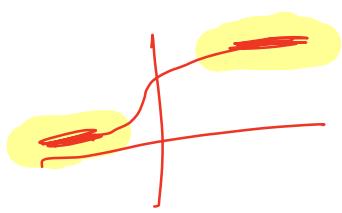
A graph showing the Leaky ReLU function  $a$  versus  $z$ . The function is zero for  $z < 0$  and increases linearly with a small positive slope for  $z \geq 0$ . A red arrow points to the positive part of the graph with the label "Leaky ReLU".

\* Sigmoid:  $\sigma(x) = \frac{1}{1+e^{-x}}$

between (0,1)  $\rightarrow$  like probability

Suitable for binary classification

Derivative:  $\sigma(x) \cdot (1 - \sigma(x))$



gradient become very small for extreme value which slow down the learning process

\* Delta:

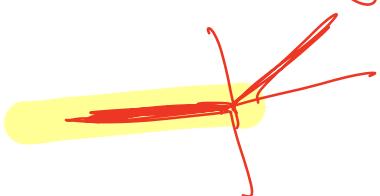
$$\Delta = \begin{cases} 0 & \text{if } a < 0 \\ 1 & \text{if } a \geq 0 \end{cases} \Rightarrow$$

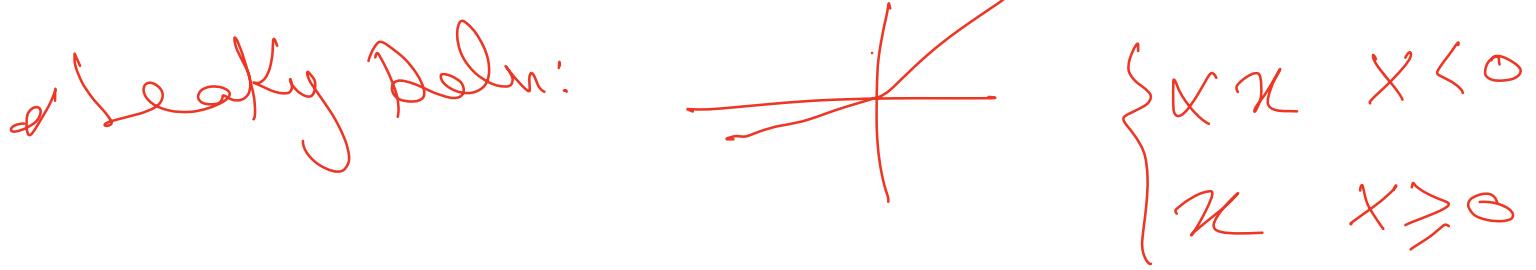
allow fast training  
(fast convergence)

Derivative:

$$\begin{cases} 0 & \text{if } n < 0 \\ 1 & \text{if } n \geq 0 \end{cases}$$

Dying problem: Neurons can become inactive during learning.





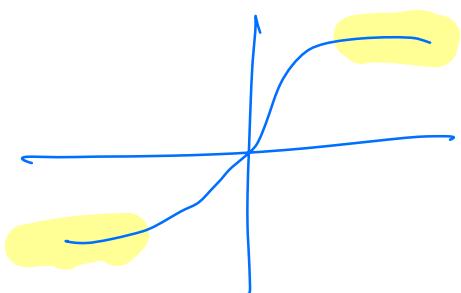
$$\max(0.001x, \kappa)$$

derivative:

$$\begin{cases} x & x < 0 \\ 1 & x \geq 0 \end{cases}$$

→ solve the dying problem.

Tanh:  $\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



$$e^x + e^{-x}$$

range between (-1, 1)

Similar to the sigmoid  
Vanishing gradient Problem at extreme values.



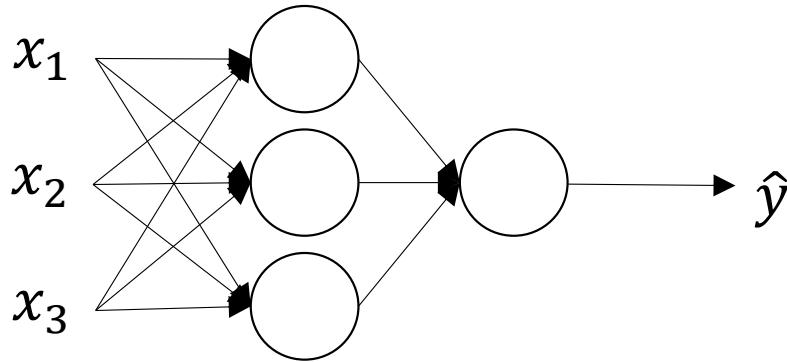
deeplearning.ai

# One hidden layer Neural Network

---

## Why do you need non-linear activation functions?

# Activation function



Given  $x$ :

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

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

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

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



deeplearning.ai

# One hidden layer Neural Network

---

Gradient descent for  
neural networks

# Gradient descent for neural networks

\* Review:

optimization algorithm used to find the weights that minimize the cost (loss) function

$$w_{\text{new}} = w_{\text{old}} - \gamma \cdot \nabla L(\theta)$$

$\xrightarrow{\text{gradient}}$   $\frac{\partial J}{\partial w}$

# Formulas for computing derivatives

$$\rightarrow \frac{\partial L}{\partial w} = ?? \quad \text{or} \quad \frac{\partial L}{\partial b} = ??$$

Forward path:

$$\begin{cases} z' = w^T x + b \\ a' = \sigma(z') \\ z^2 = w^{2T} a' + b \\ a^2 = \sigma(z^2) \end{cases}$$



$L(A^2, y)$  cross entropy:  $L(A^2, y)$

$$\frac{\partial L}{\partial A_2} = \frac{A_2 - y}{A_2(1 - A_2)}$$

$$\boxed{\frac{\partial L}{\partial w_2}, \frac{\partial L}{\partial b_2}, \frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial b_1}}$$

$$\left( \frac{\partial L}{\partial A_2} \right) \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial w_2} = \frac{\partial L}{\partial w_2}$$

$$\left( \frac{\partial L}{\partial A_2} \right) \cdot \frac{\partial A_2}{\partial Z_2} \cdot \frac{\partial Z_2}{\partial b_2} = \frac{\partial L}{\partial b_2}$$

$$\frac{\partial L}{\partial Z_2} = \frac{\partial L}{\partial A_2} \cdot \frac{\partial A_2}{\partial Z_2} = \frac{A_2 - y}{A_2(1 - A_2)} \cdot A_2(1 - A_2)$$

$\cancel{A_2(1 - A_2)}$

$A_2 - y$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2}$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z_2} \cdot A_1^T$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial z_2} \cdot )$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial A_1} \cdot \frac{\partial A_1}{\partial w_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$w_1^T$ ,  $g^{[1]}(z)$

$$\frac{\partial L}{\partial z_1} = w_1^T \frac{\partial L}{\partial z_2} * g^{[1]} \cdot [z]$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_i} \cdot \chi^T$$

$$\frac{\partial L}{\partial b_i} = \frac{\partial L}{\partial z_i} \cdot \frac{\partial z_i}{\partial b_i}$$

$$\frac{\partial L}{\partial b_i} = \frac{\partial L}{\partial z_i} \cdot \dots$$



deeplearning.ai

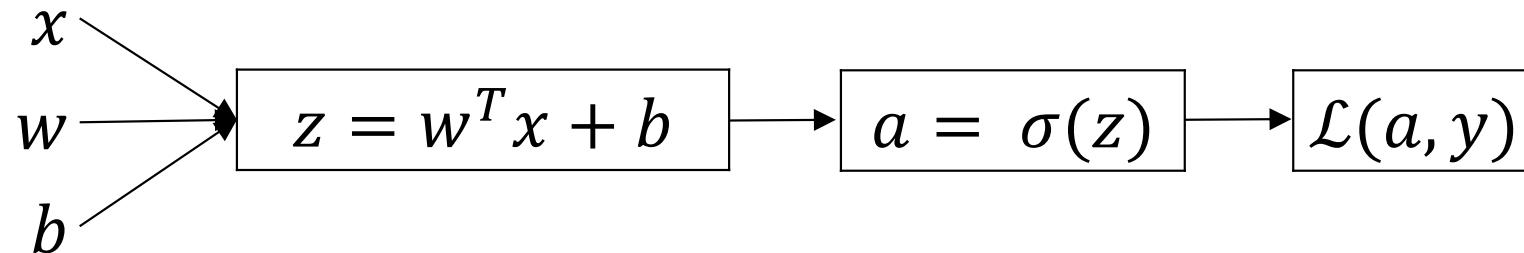
# One hidden layer Neural Network

---

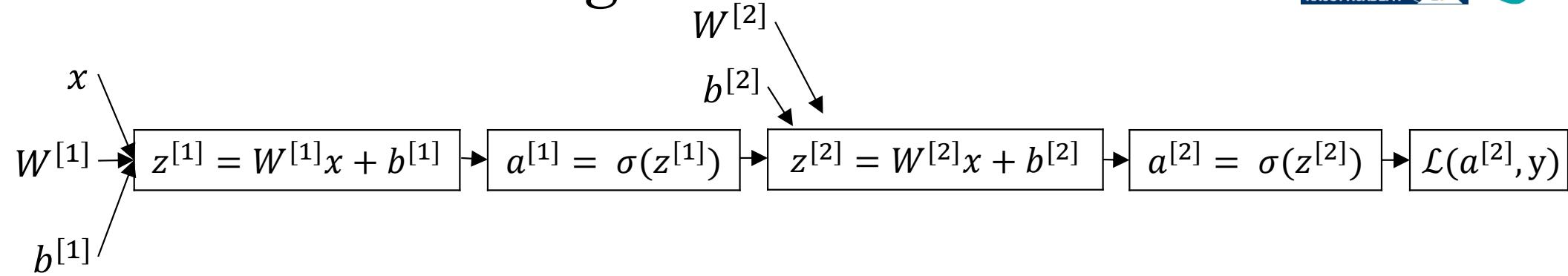
## Backpropagation intuition

# Computing gradients

## Logistic regression



# Neural network gradients



# Summary of gradient descent

$$dz^{[2]} = \cancel{a^{[2]} - y}$$

$$dW^{[2]} = \cancel{dz^{[2]} a^{[1]T}}$$

$$db^{[2]} = \cancel{dz^{[2]}}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = \cancel{dz^{[1]} x^T}$$

$$db^{[1]} = \cancel{dz^{[1]}}$$

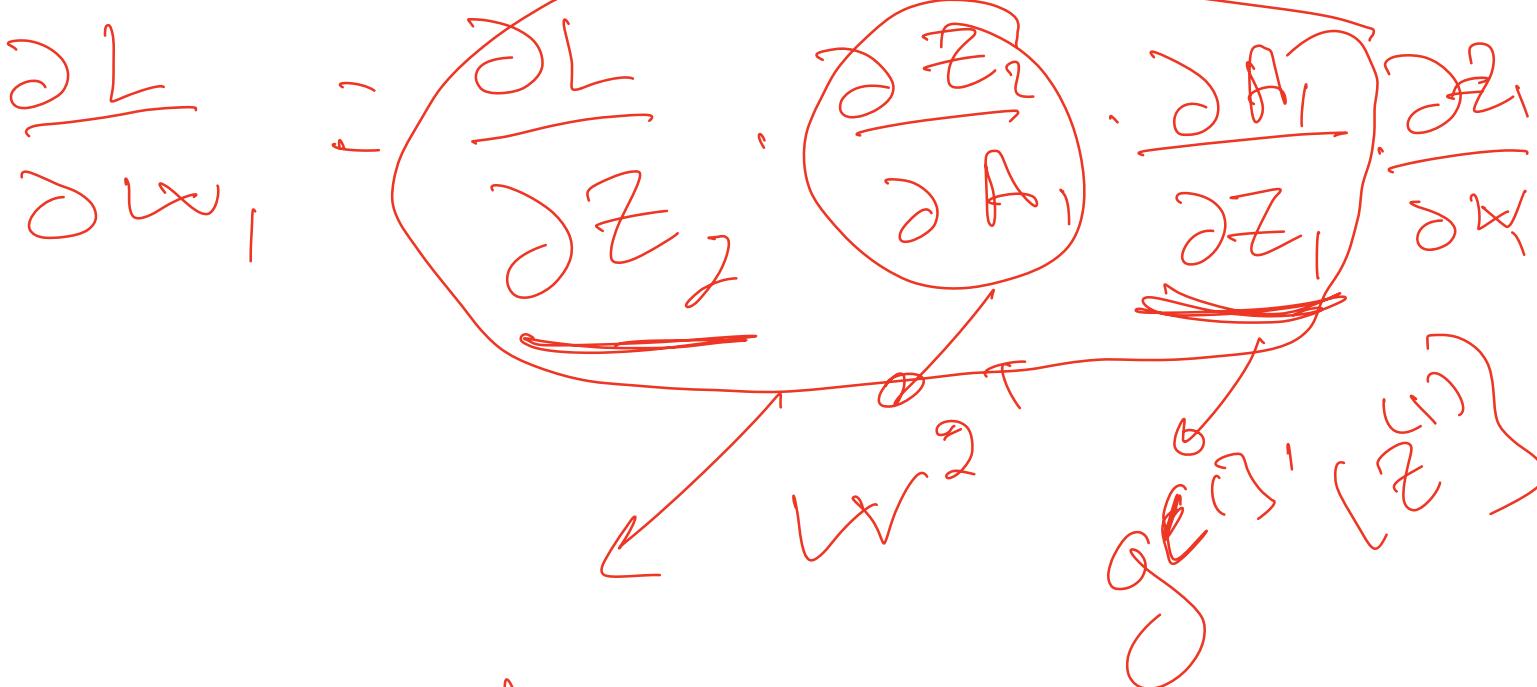
Backpropagation:

$$\frac{\partial L}{\partial A_2} = A^{[2]} - y$$

$$\frac{\partial A_2}{\partial z_2} \rightarrow \frac{\partial L}{\partial z_2}$$

$$\frac{\partial z_2}{\partial A} = A^{[1]T}$$

$$\frac{\partial w_2}{\partial z_2} \rightarrow \frac{\partial z_2}{\partial b_2} = w$$



$$\frac{\partial L}{\partial z^1}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1}$$

$$X \nearrow$$

$$\frac{\partial L}{\partial b_1} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1}$$

# Summary of gradient descent

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$$

$$dZ^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(Z^{[1]})$$

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$$

$$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$$



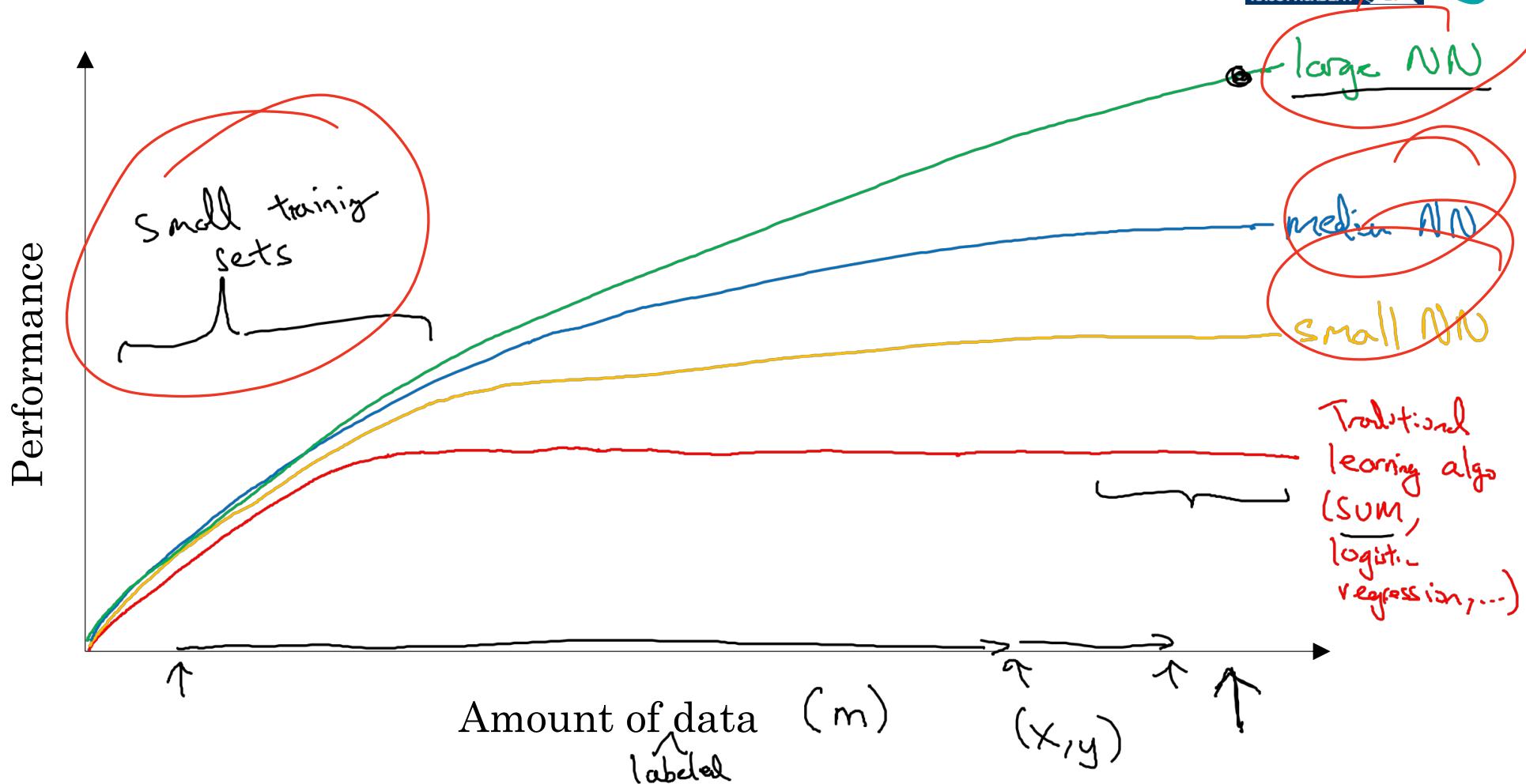
deeplearning.ai

# Introduction to Neural Networks

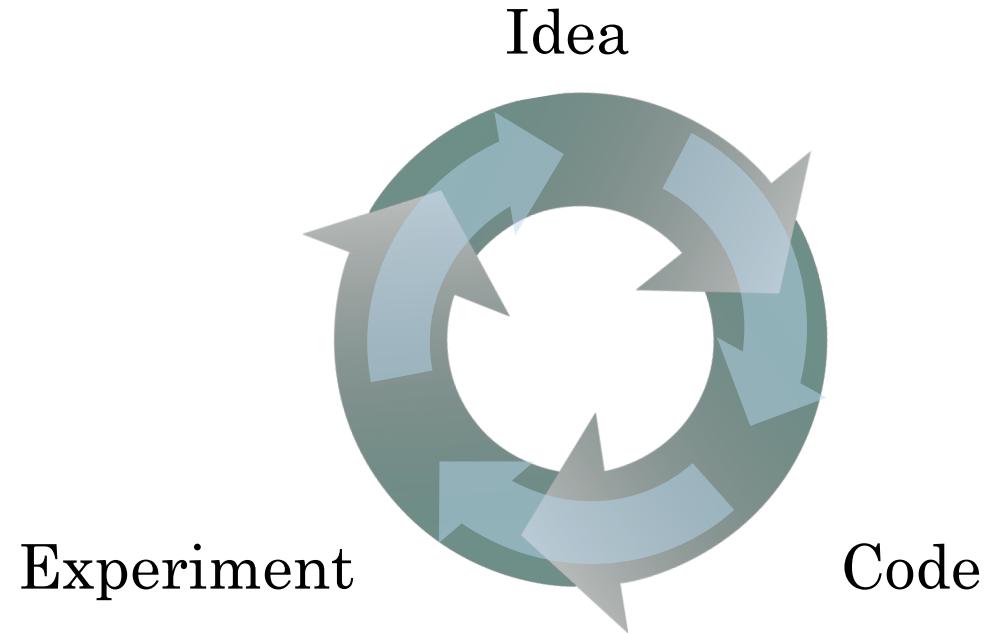
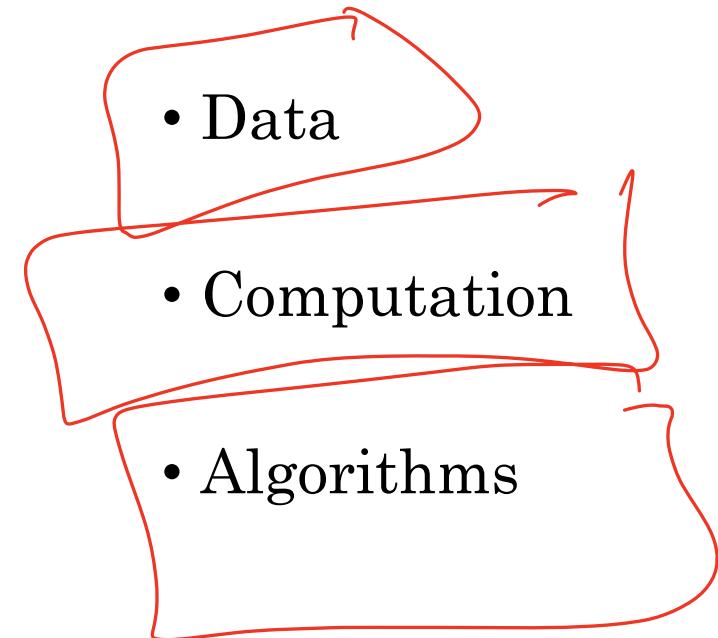
---

## Why is Deep Learning taking off?

# Scale drives deep learning progress



# Scale drives deep learning progress





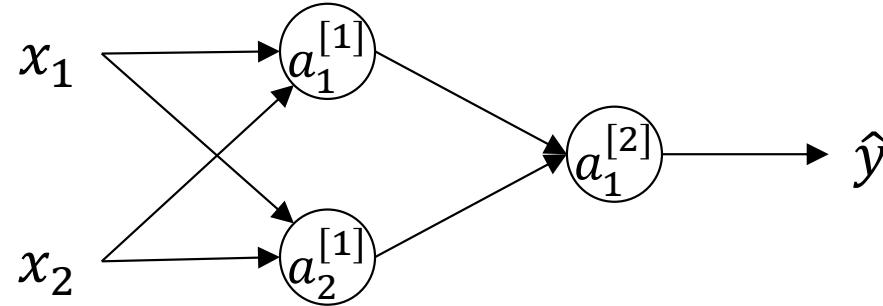
deeplearning.ai

# One hidden layer Neural Network

---

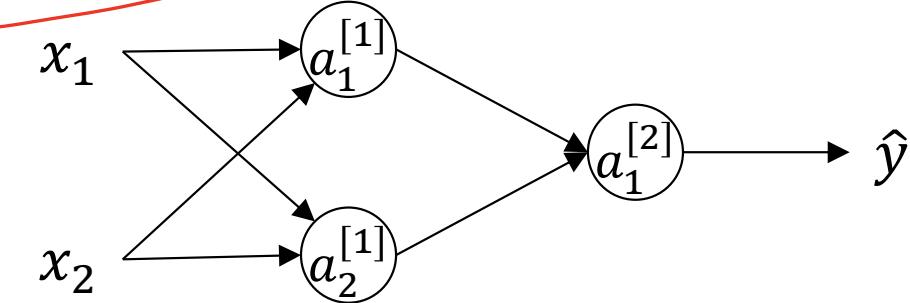
## Random Initialization

# What happens if you initialize weights to zero?



Zero is not generally recommended  
Random initialization is preferred  
to allow the network to learn

# Random initialization



unique feature  
and avoid the  
symmetry problem



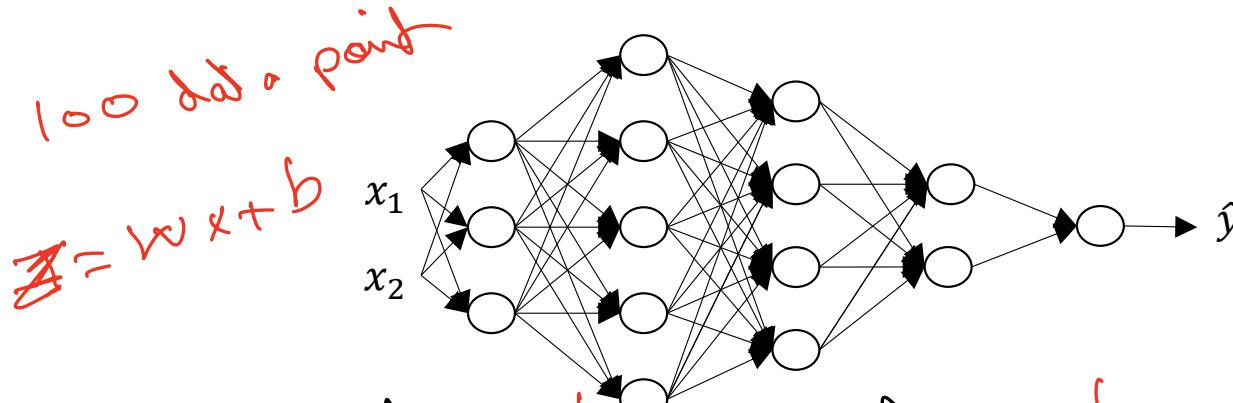
deeplearning.ai

# Deep Neural Networks

---

## Getting your matrix dimensions right

# Parameters $W^{[l]}$ and $b^{[l]}$



$$w^1 (3, 2)$$

$$w^2 (5, 3)$$

$$w^3 (4, 5)$$

$$w^4 (2, 4)$$

$$w^5 (1, 2)$$

$$b^1 (3, 100) \leftarrow x_1$$

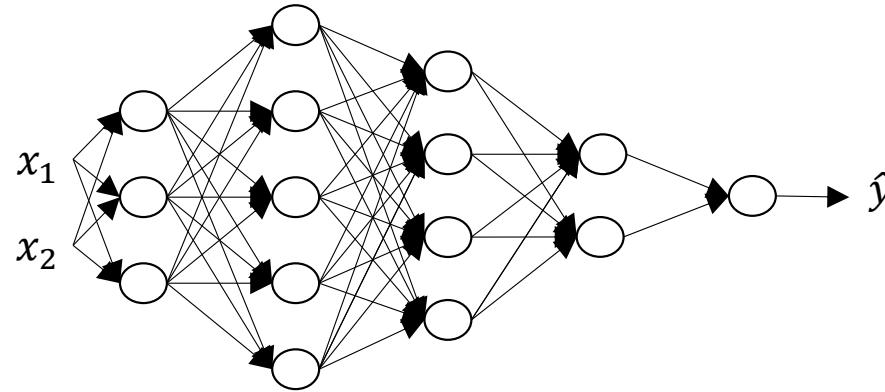
$$b^2 (5, 100) \leftarrow x_2$$

$$b^3 (4, 100) \leftarrow x_3$$

$$b^4 (2, 100) \leftarrow x_4$$

$$b^5 (1, 100) \leftarrow x_5$$

# Vectorized implementation





deeplearning.ai

# Deep Neural Networks

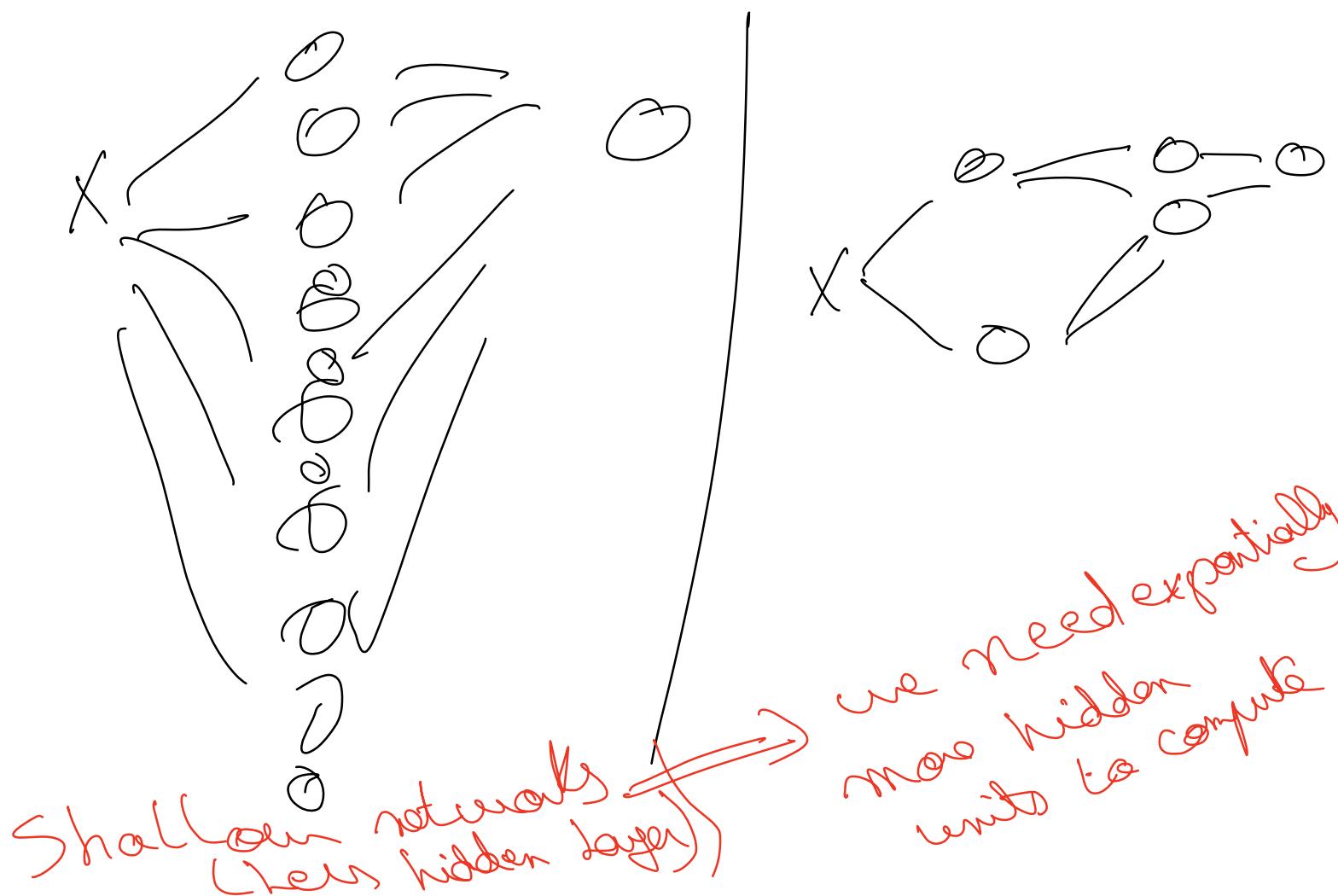
→ increase  
the  
number of  
layers

---

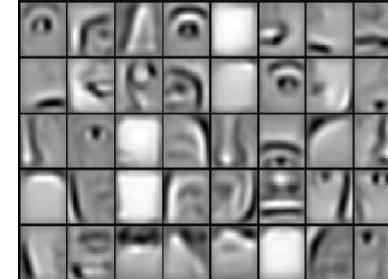
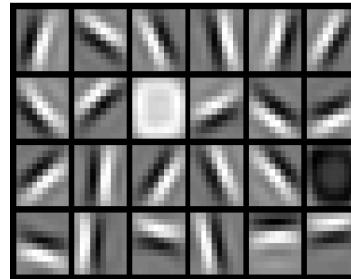
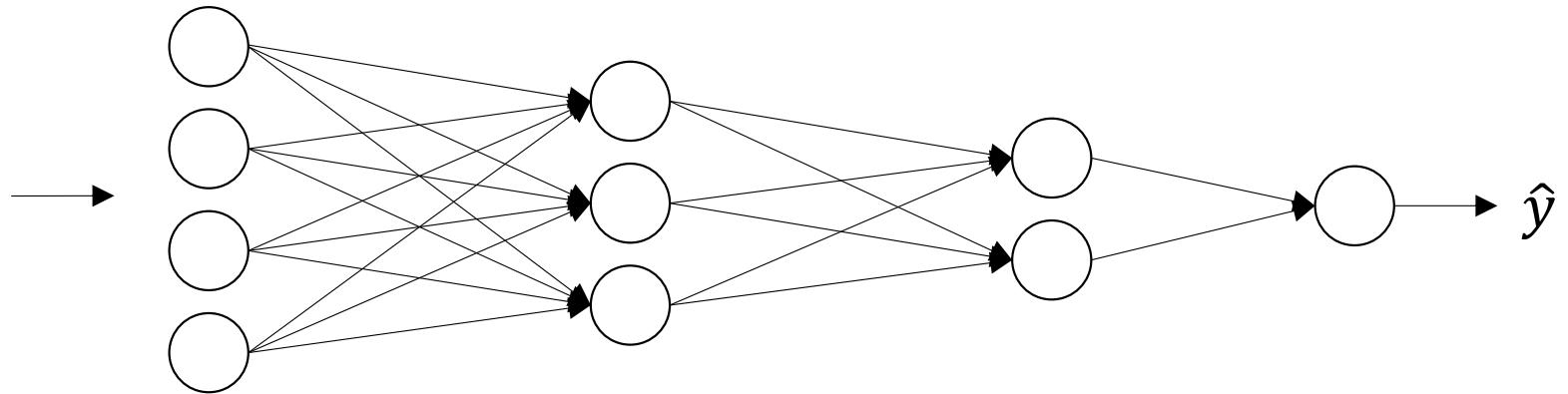
## Why deep representations?

→ increase the # of layers can potentially capture and compute more complex function with a reasonable number of hidden units in each layer.

Increasing the # of layers can potentially capture and compute complex function with a reasonable number of hidden units in each layer



# Intuition about deep representation



# Circuit theory and deep learning

Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.





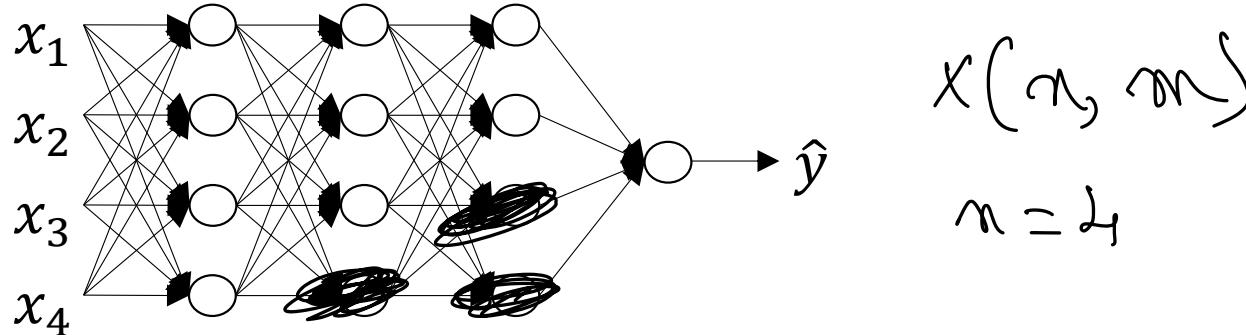
deeplearning.ai

# Deep Neural Networks

---

## Building blocks of deep neural networks

# Forward and backward functions



$w_1(4, h)$

$w_2(3, h)$

$b_1(4, m)$

$b_2(3, m)$

$w_3(2, 3)$

$w_4(1, 2)$

$b_3(2, m)$

$b_4(1, m)$

# Forward and backward functions





deeplearning.ai

# Deep Neural Networks

---

## Forward and backward propagation

# Forward propagation for layer $l$

Input  $a^{[l-1]}$

Output  $a^{[l]}$ , cache ( $z^{[l]}$ )

# Backward propagation for layer $l$

Input  $da^{[l]}$

Output  $da^{[l-1]}, dW^{[l]}, db^{[l]}$

# Summary

Step:

1 - Forward pass

2 - Backward pass

3 - optimizer