

# Artificial Intelligence and Machine Learning

## Neural Networks

# Lecture Outline

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

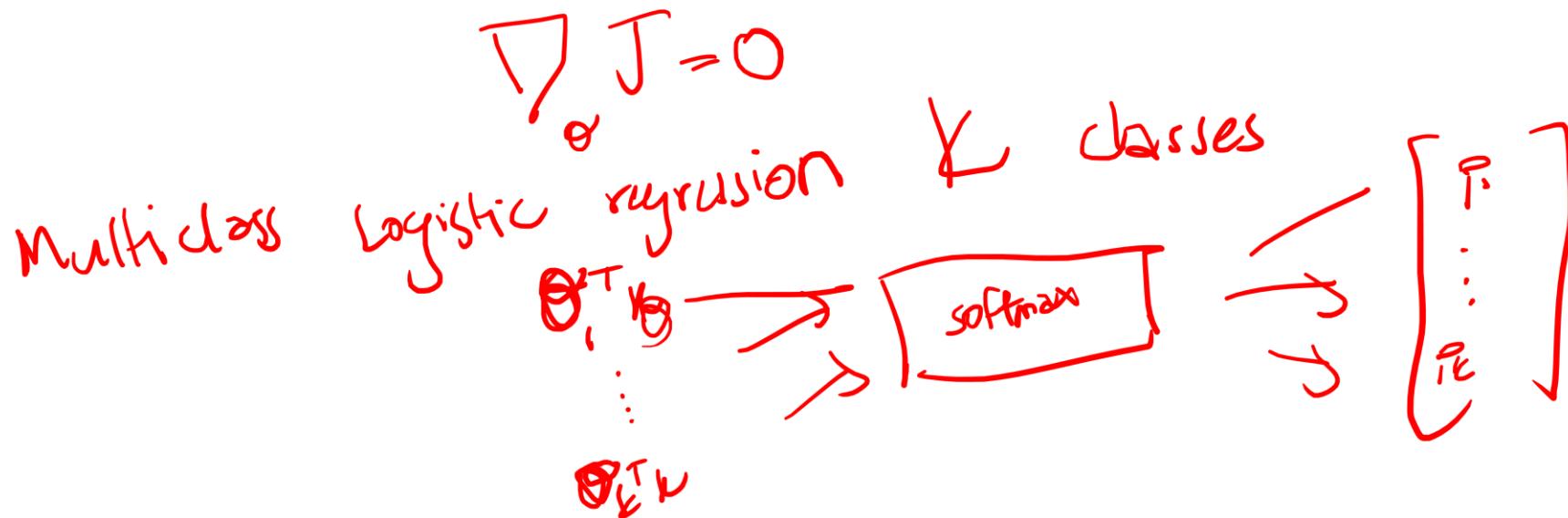
# Review: Logistic Regression

$$p_i = \sigma(\theta^T x)$$

↓  
sigmoid

Loss function: Log Loss

$$\rightarrow J = \frac{1}{N} \sum_{i=1}^N -[y_i \log p_i + (1-y_i) \log(1-p_i)]$$



Gradient descent



learning

$$\theta \leftarrow \theta - \eta \nabla_{\theta} J$$



basic FULL Training



minibatch



Stochastic GD



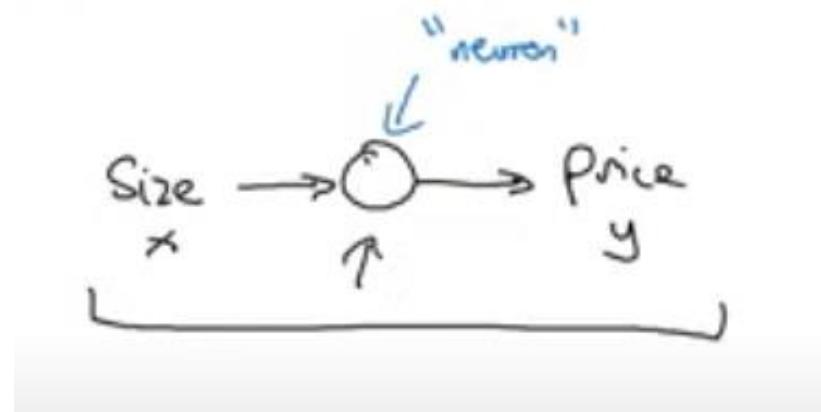
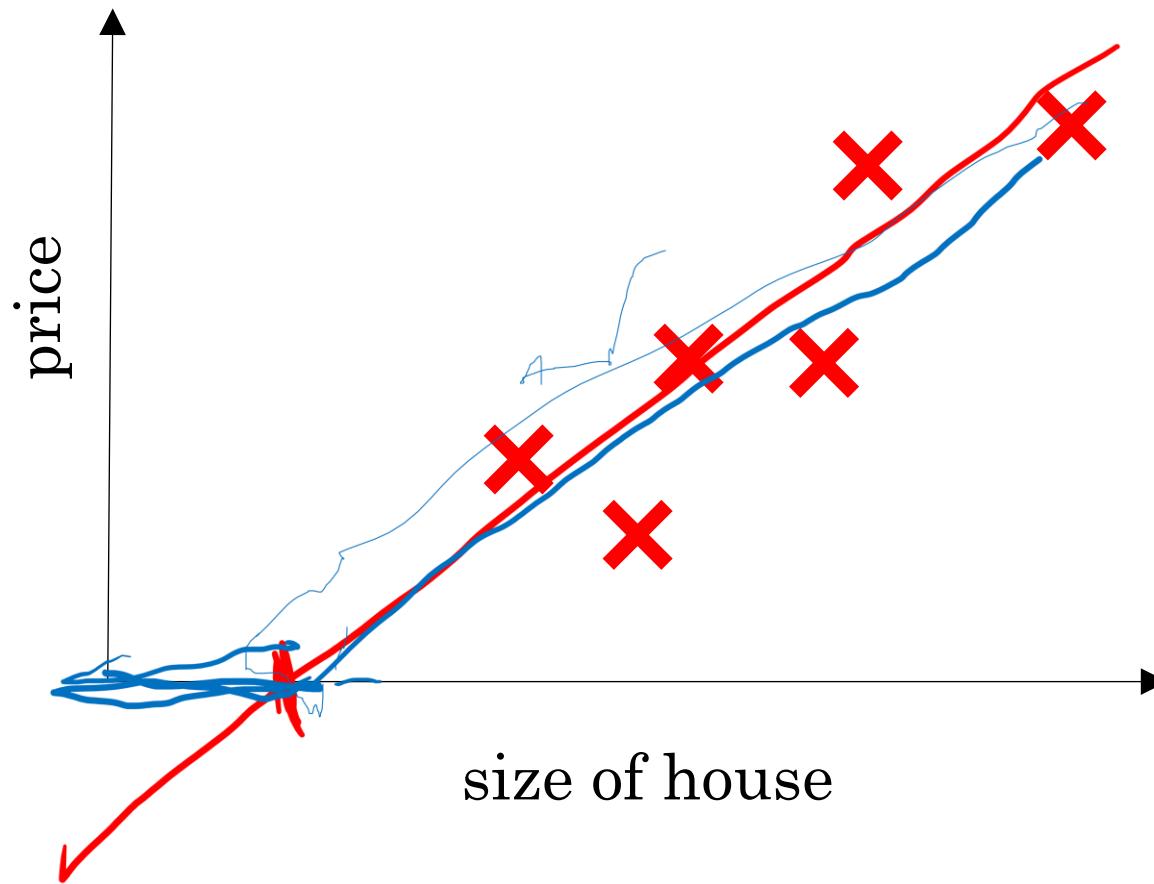
deeplearning.ai

# Introduction to Deep Learning

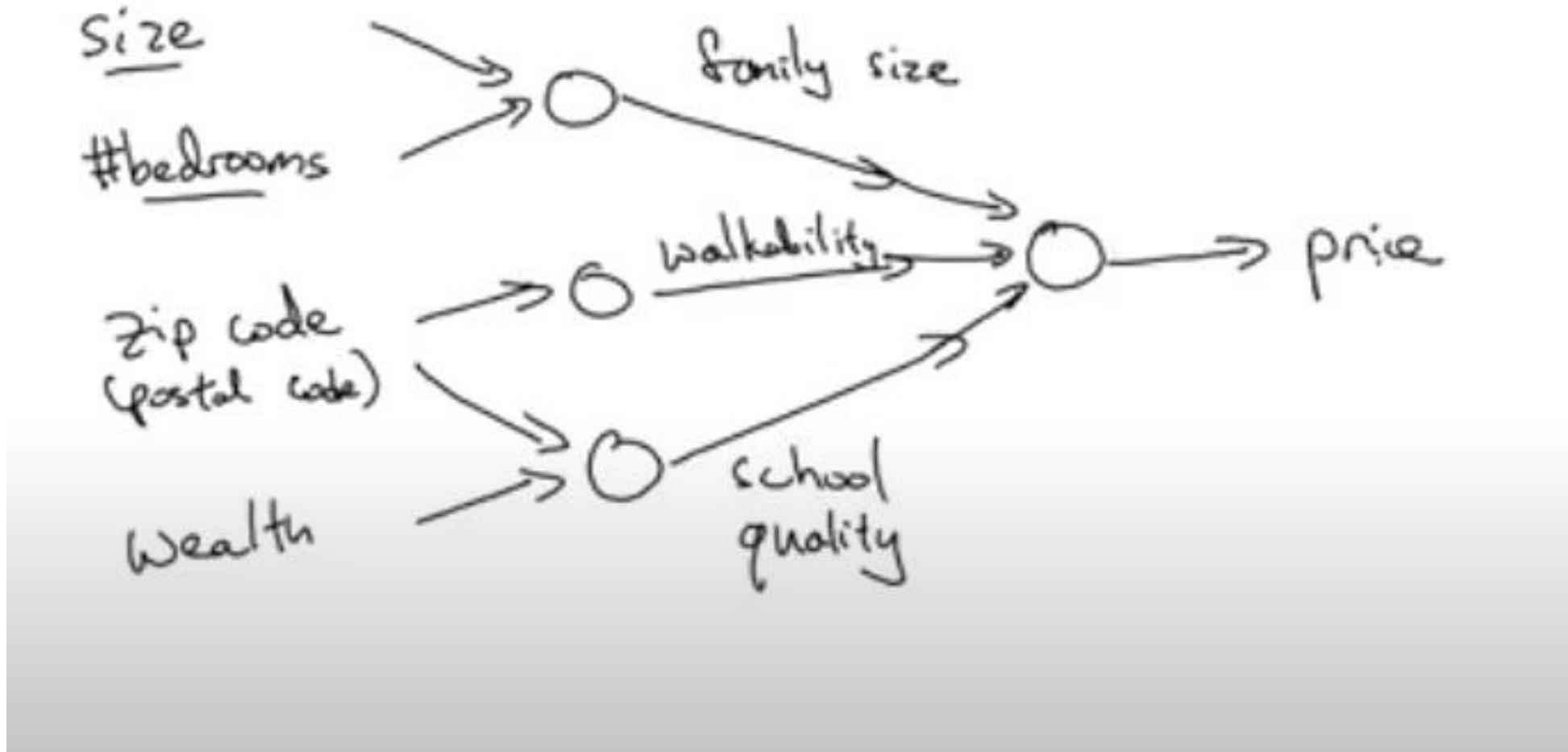
---

## What is a Neural Network?

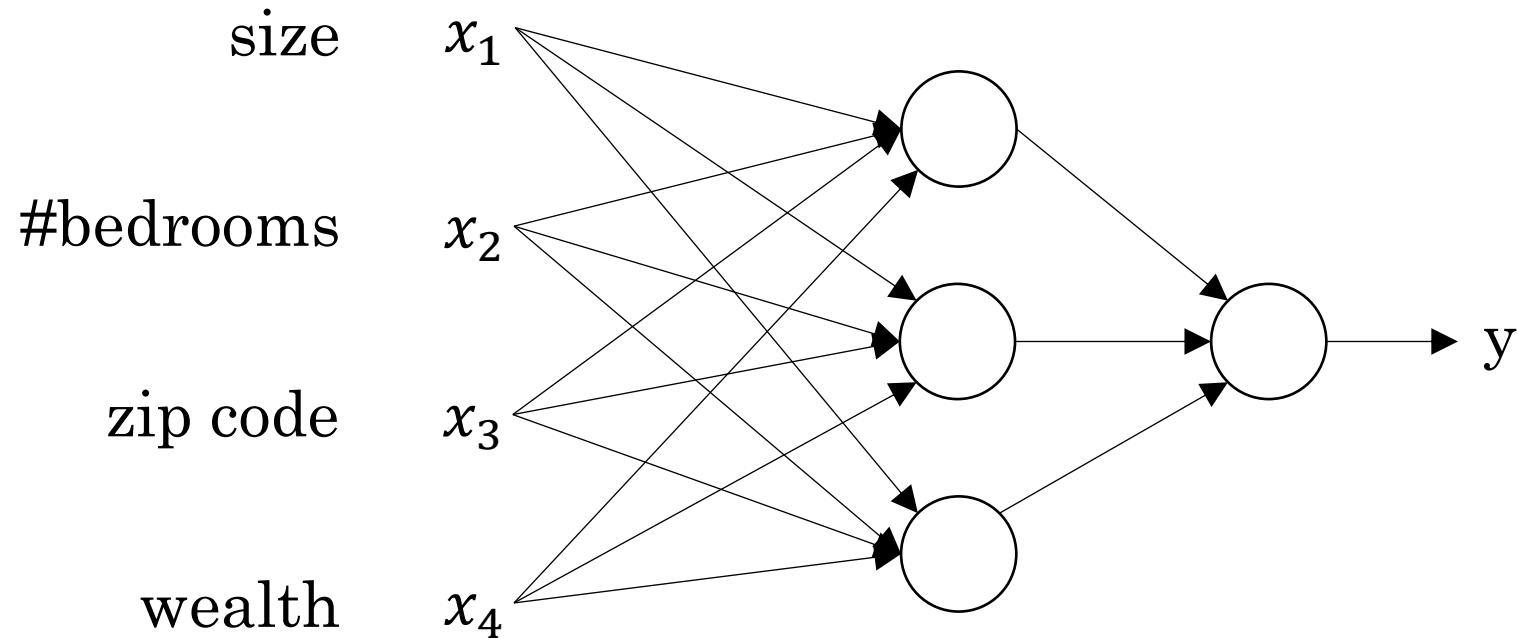
# Housing Price Prediction



# Housing Price Prediction



# Housing Price Prediction





# Introduction to Deep Learning

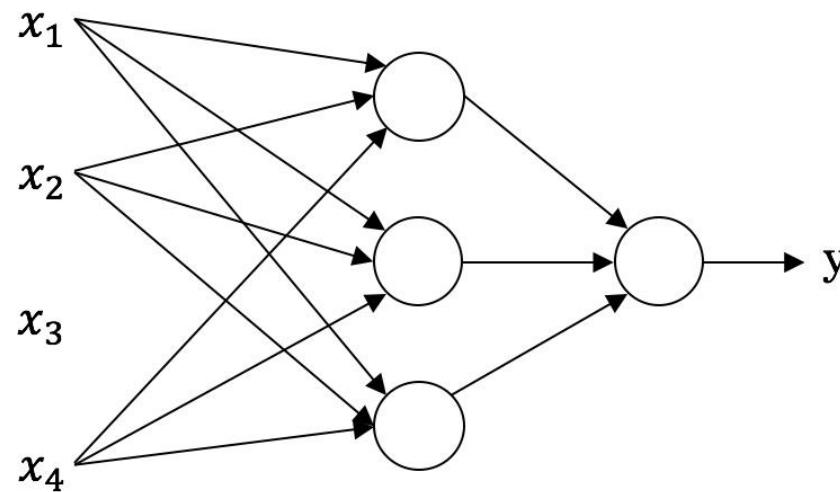
---

## Supervised Learning with Neural Networks

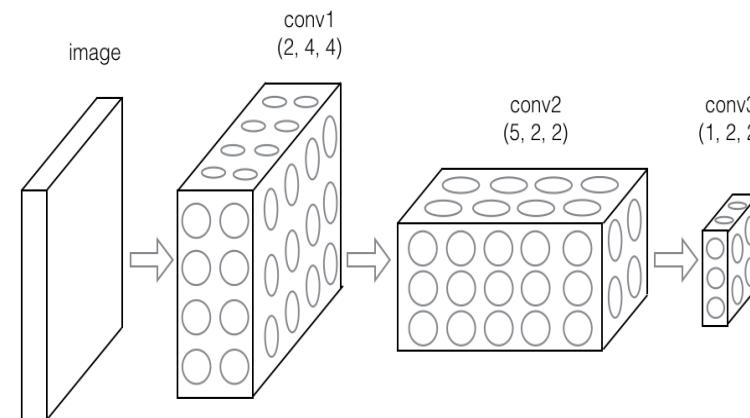
# Supervised Learning

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

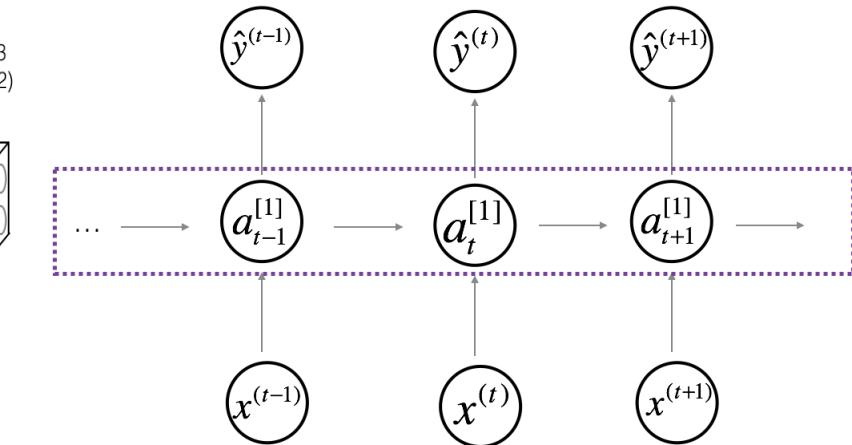
# Neural Network examples



Standard NN



Convolutional NN



Recurrent NN



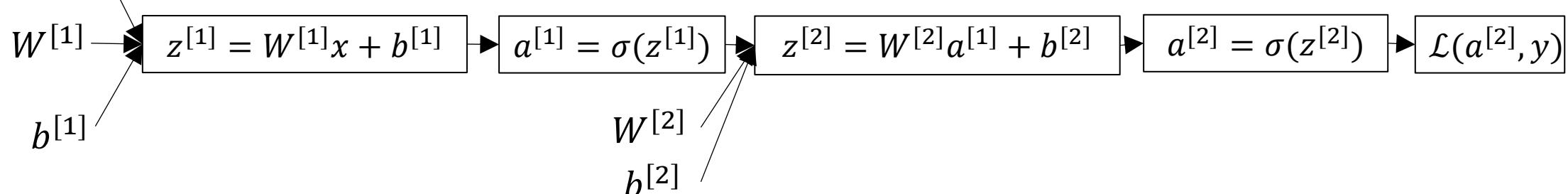
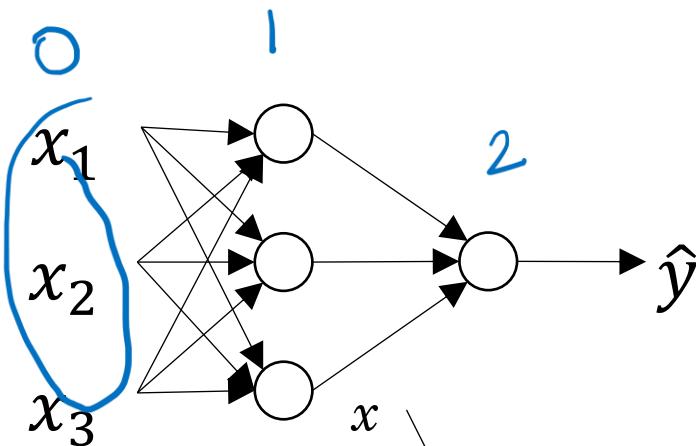
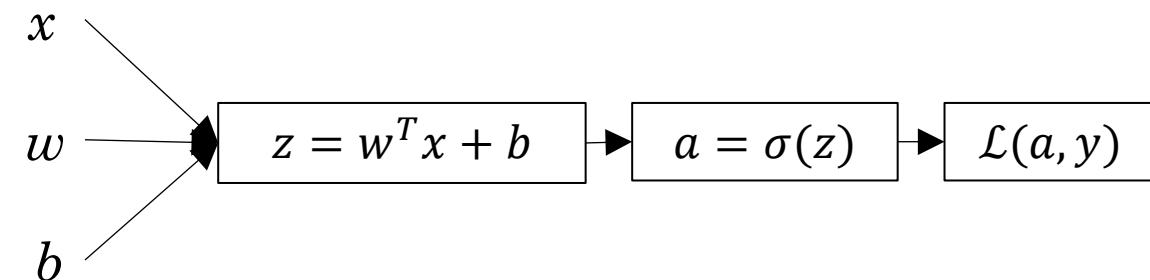
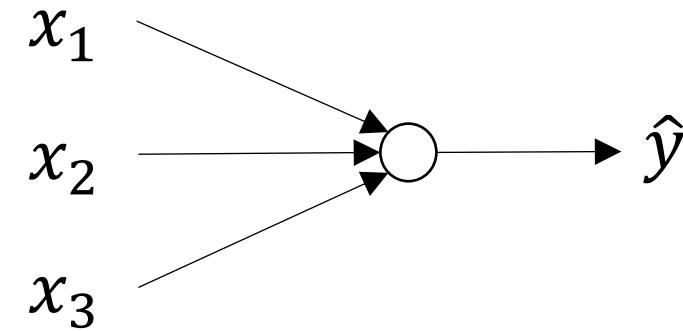
deeplearning.ai

# One hidden layer Neural Network

---

# Neural Networks Overview

# What is a Neural Network?





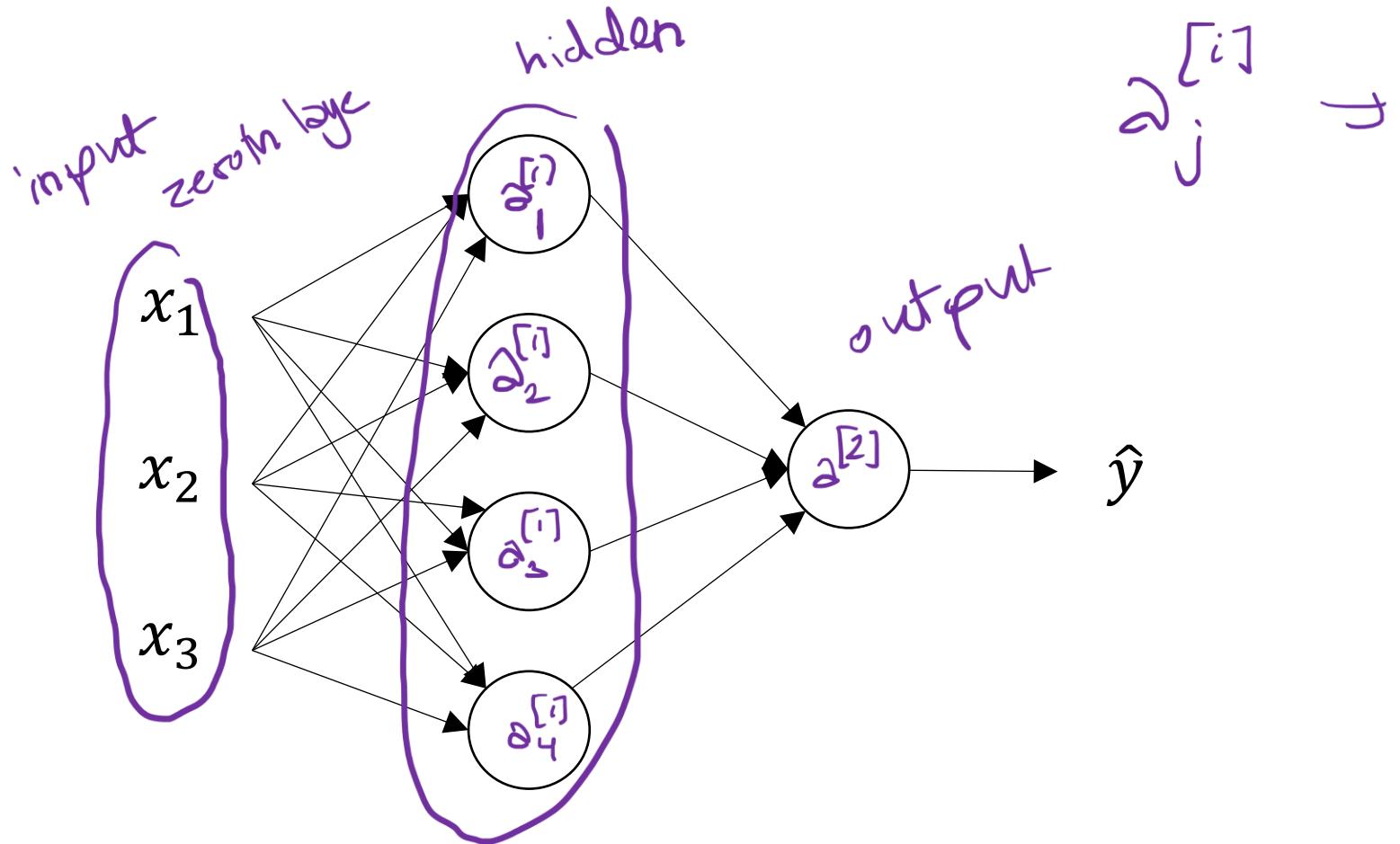
deeplearning.ai

# One hidden layer Neural Network

---

## Neural Network Representation

# Neural Network Representation





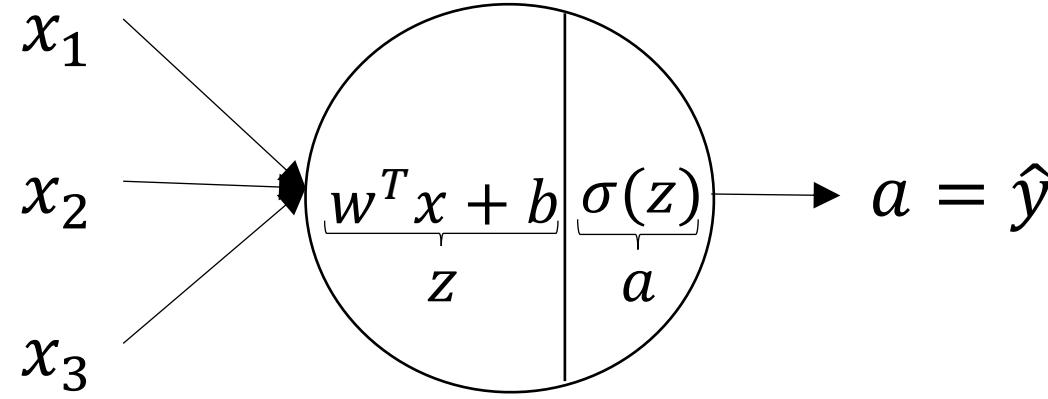
deeplearning.ai

# One hidden layer Neural Network

---

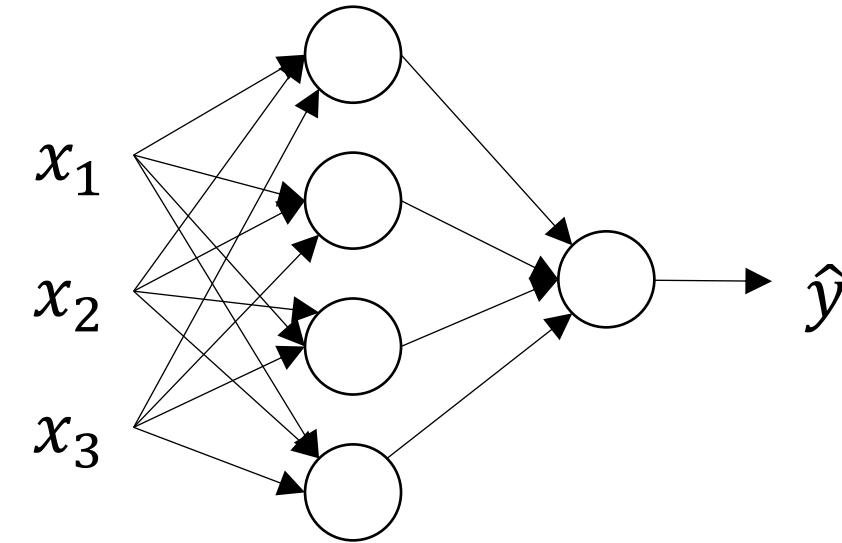
## Computing a Neural Network's Output

# Neural Network Representation

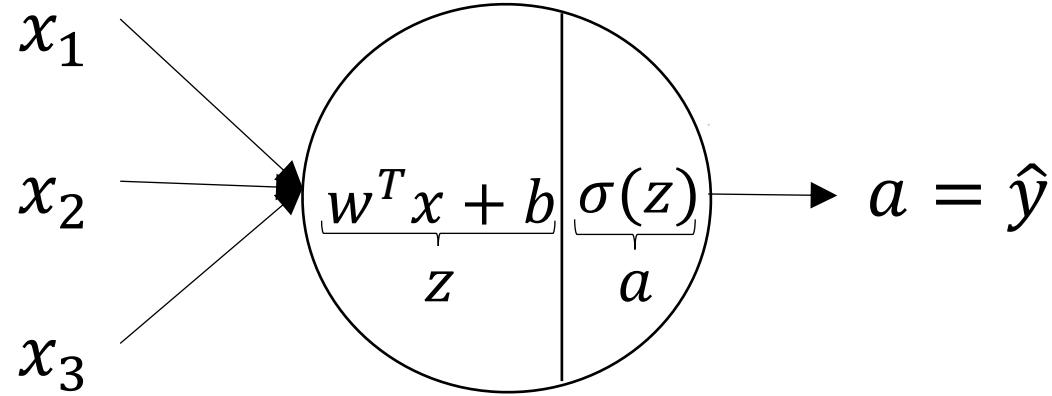


$$z = w^T x + b$$

$$a = \sigma(z)$$

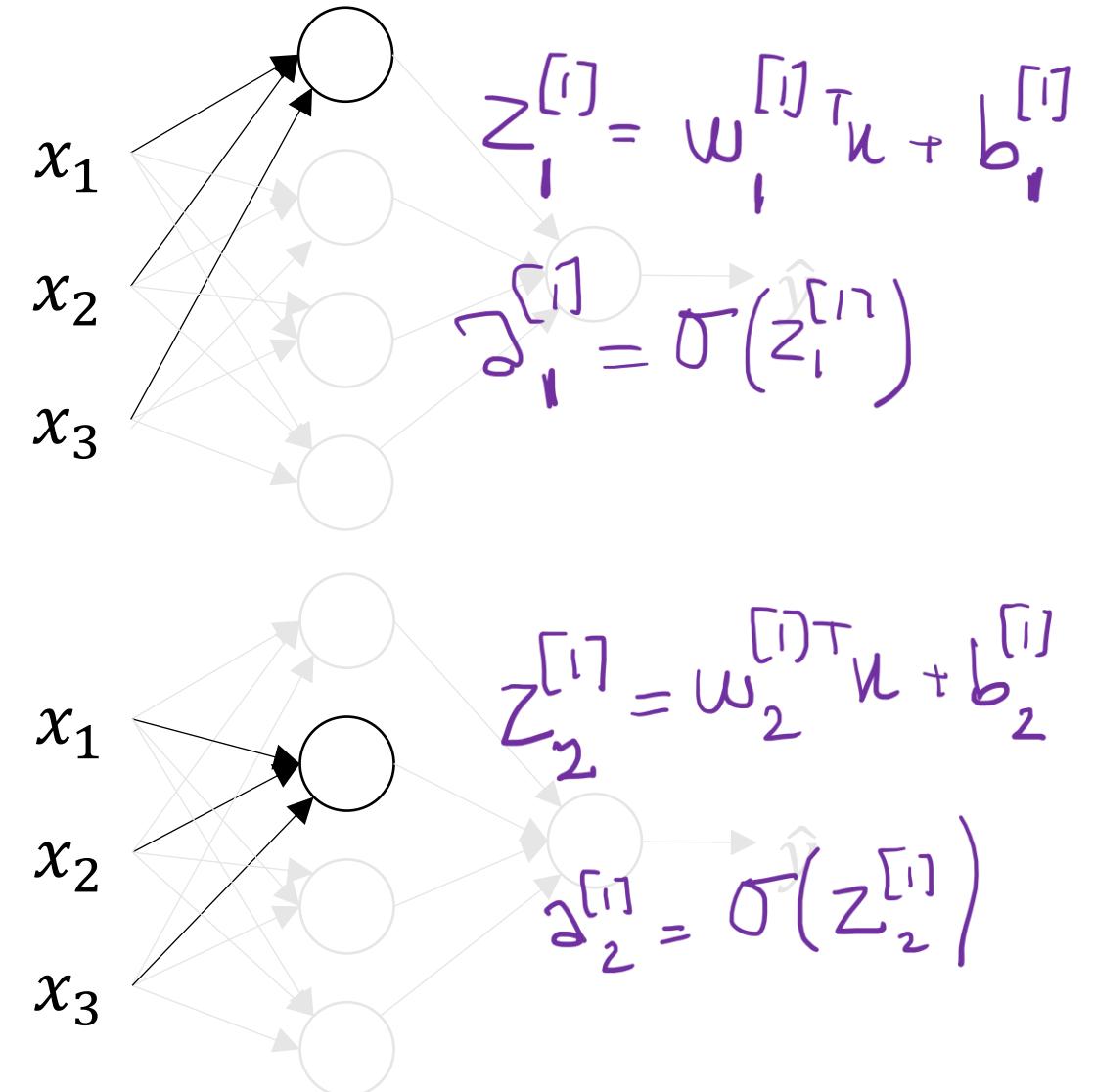


# Neural Network Representation

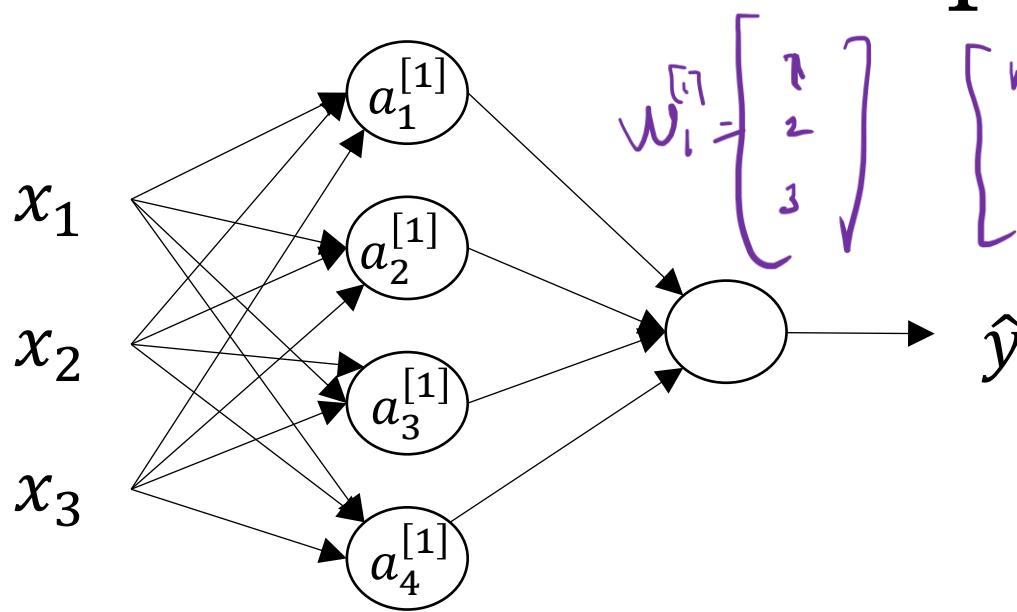


$$z = w^T x + b$$

$$a = \sigma(z)$$



# Neural Network Representation



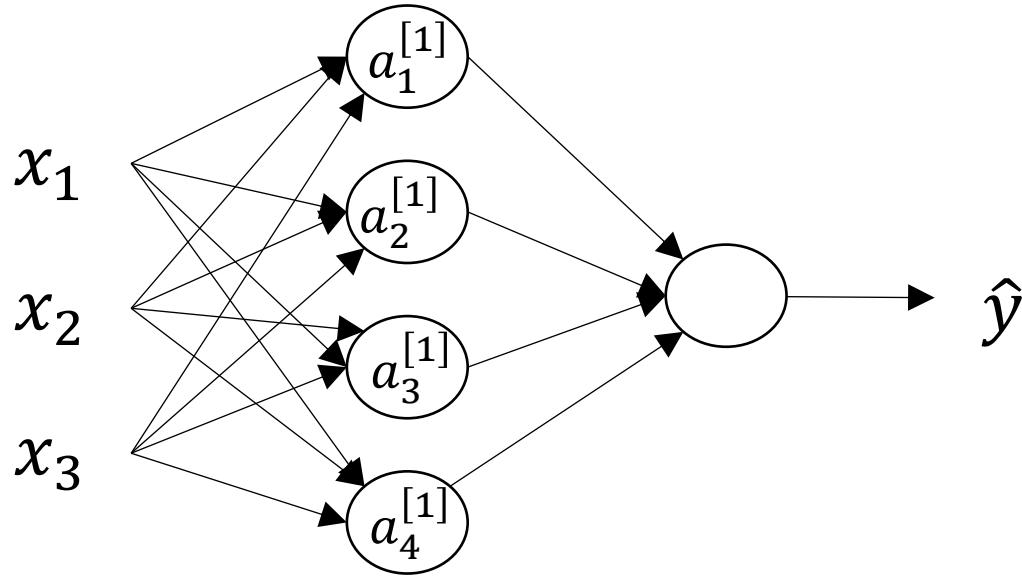
$$\left\{ \begin{array}{l} z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \\ z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \\ z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \\ z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \end{array} \right. \quad \left\{ \begin{array}{l} a_1^{[1]} = \sigma(z_1^{[1]}), \\ a_2^{[1]} = \sigma(z_2^{[1]}), \\ a_3^{[1]} = \sigma(z_3^{[1]}), \\ a_4^{[1]} = \sigma(z_4^{[1]}) \end{array} \right.$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}$$

$$W^{[1]} = \begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \\ w_4^{[1]T} \end{bmatrix}$$

# Neural Network Representation learning



Given input  $x$ :

$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \end{array} \right.$$

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

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



deeplearning.ai

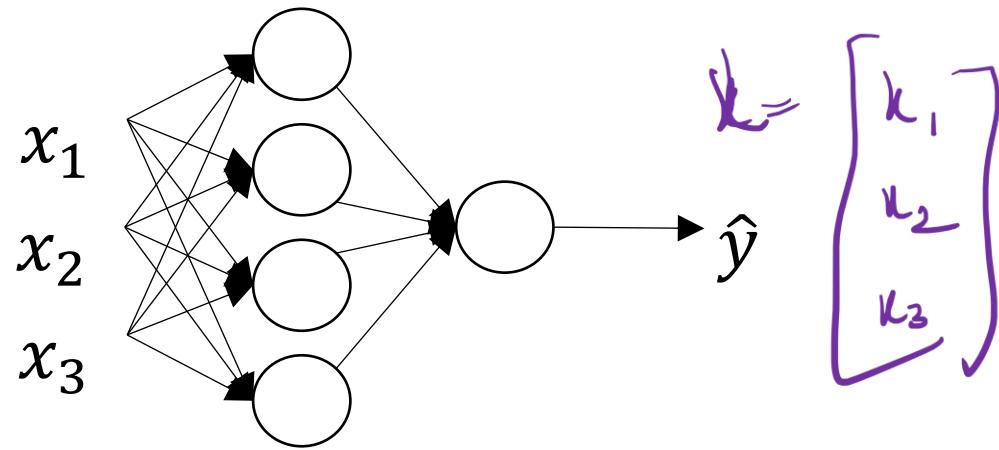
$$\{k_i, y_i\}_{i=1}^m$$

# One hidden layer Neural Network

---

## Vectorizing across multiple examples

# Vectorizing across multiple examples



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

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

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

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

$$X = [x^{(1)} \quad k^{(2)} \quad \dots \quad k^{(m)}]$$

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

$$Z^{[1]} = [z^{(1)} \quad z^{(2)} \quad \dots \quad z^{(m)}]$$

$$A^{[1]} = [a^{[1](1)} \quad a^{[1](2)} \quad \dots \quad a^{[1](m)}]$$

# Vectorizing across multiple examples



for i = 1 to m:

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

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

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

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

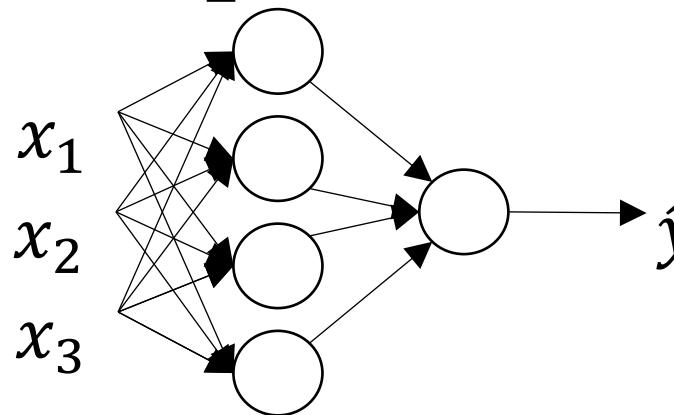
$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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

# 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

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

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

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

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

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

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

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

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



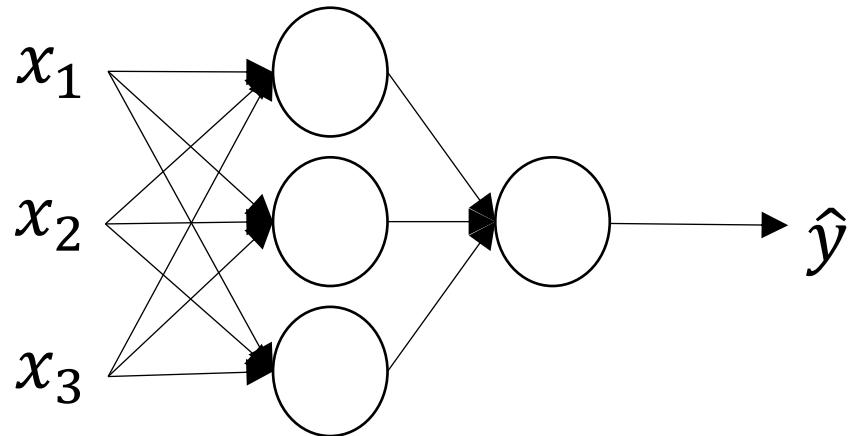
deeplearning.ai

# One hidden layer Neural Network

---

## Activation functions

# Activation functions



Given  $x$ :

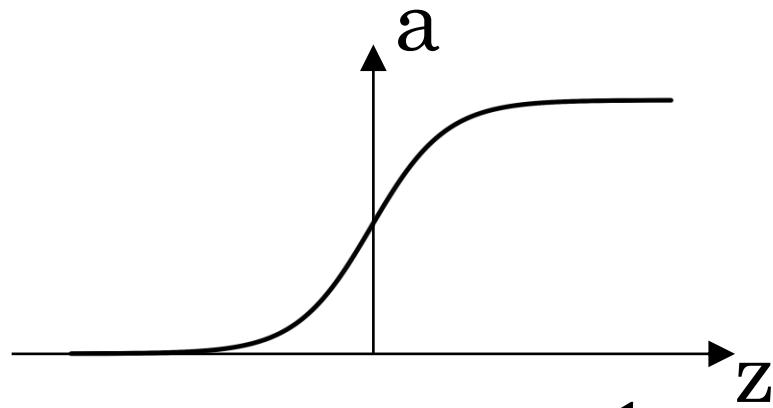
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

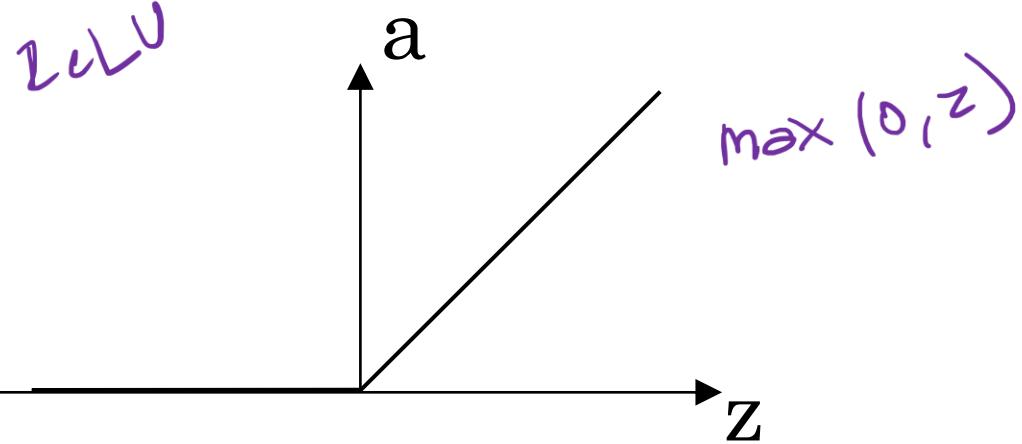
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

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

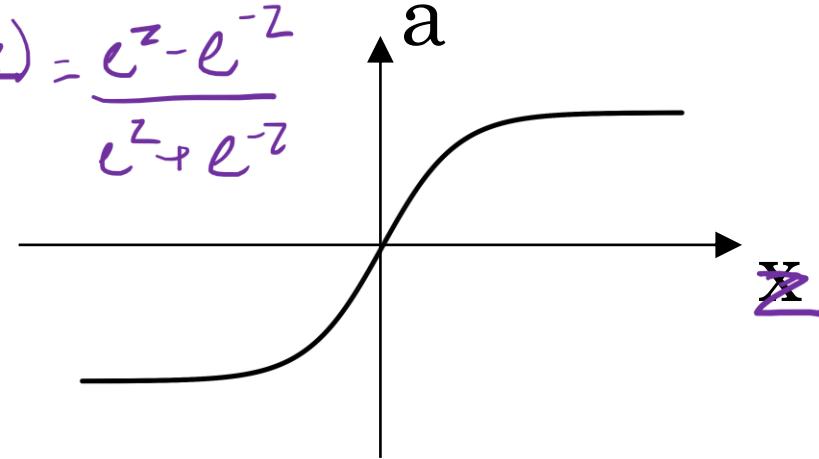
# Pros and cons of activation functions



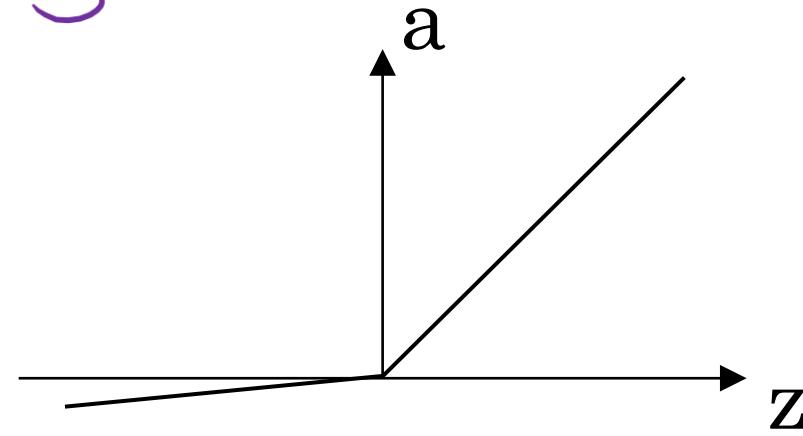
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



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



Leaky ReLU





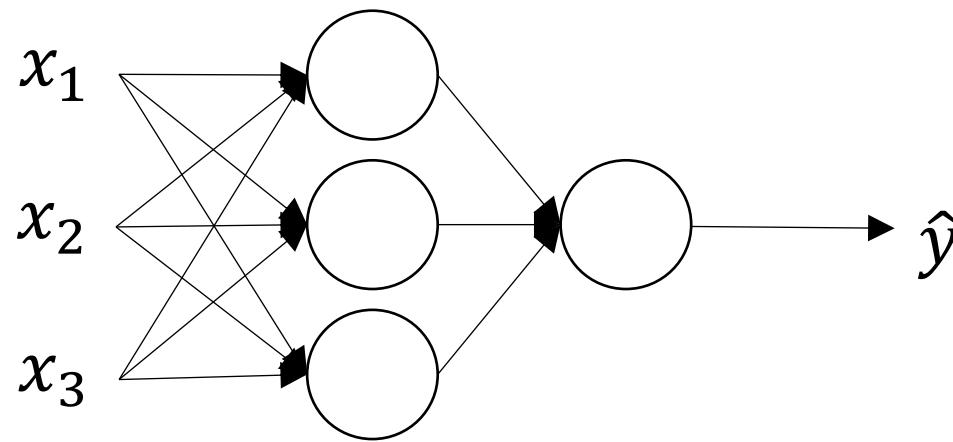
deeplearning.ai

# One hidden layer Neural Network

---

Why do you  
need non-linear  
activation functions?

# Activation function



Given  $x$ :

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

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

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

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

$$z^{[2]} = Z^{[2]}$$

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

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

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

$$g(z) = z$$

linear activation

identity act.

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

$$z^{[2]} = W'x + b'$$

$$z^{[2]} = \underline{W'x + b'}$$



deeplearning.ai

# One hidden layer Neural Network

---

# Gradient descent for neural networks

# Gradient descent for neural networks

Parameters :  $w^{[1]}, w^{[2]}, b^{[1]}, b^{[2]}$

$$J = \text{LogLoss} = \frac{1}{m} \sum_{i=1}^m -(y_i \log p_i + (1-y_i) \log (1-p_i))$$

$$w^{[1]} \leftarrow w^{[1]} - \alpha \frac{\partial J}{\partial w^{[1]}}$$

$$b^{[1]} \leftarrow b^{[1]} - \alpha \frac{\partial J}{\partial b^{[1]}}$$

$$w^{[2]} \leftarrow w^{[2]} - \alpha \frac{\partial J}{\partial w^{[2]}}$$

$$b^{[2]} \leftarrow b^{[2]} - \alpha \frac{\partial J}{\partial b^{[2]}}$$

Repeat

forward propagation

Compute predictions

$(\hat{y}^{(i)})_{i=1, \dots, m}$

Compute  $dW^{[1]}, db^{[1]}, dW^{[2]}, db^{[2]}$

backward propagation

$$w^{[1]} \leftarrow w^{[1]} - \alpha dw^{[1]}, b^{[1]} \leftarrow b^{[1]} - \alpha db^{[1]}$$

$$w^{[2]} \leftarrow w^{[2]} - \alpha dw^{[2]}, b^{[2]} \leftarrow b^{[2]} - \alpha db^{[2]}$$

# Forward propagation

$g$

$$Z^{[1]} = W^{[1]} \times \neg b^{[1]}$$

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

$$Z^{[2]} = W^{[2]} A^{[1]} \times \neg b^{[2]}$$

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

# Formulas for computing derivatives

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

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

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

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

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

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



deeplearning.ai

# One hidden layer Neural Network

---

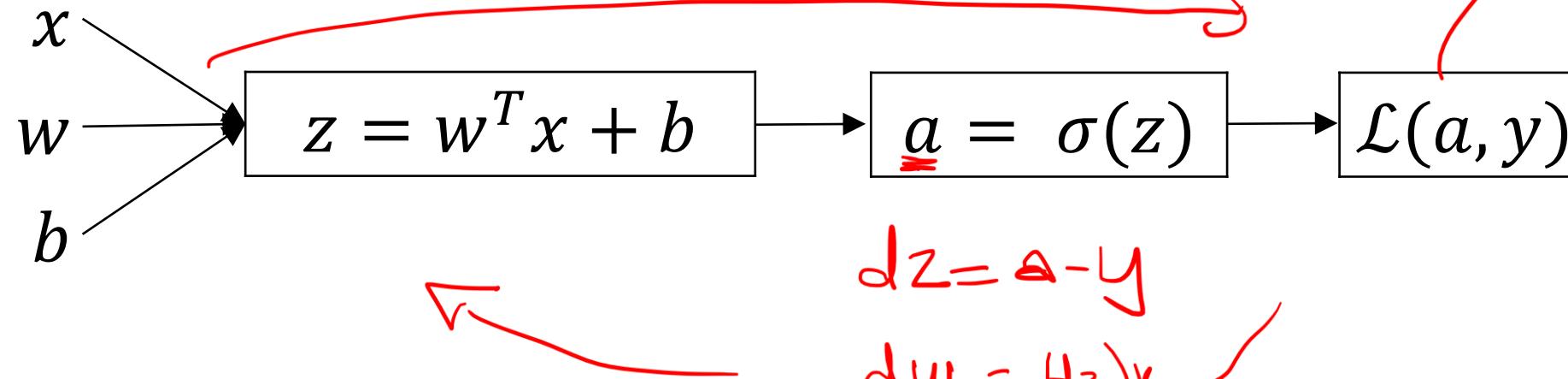
## Backpropagation intuition

# Computing gradients

Logistic regression

$$\frac{dw}{J} = \frac{\partial J}{\partial w}$$

$$\frac{db}{J} = \frac{\partial J}{\partial b}$$



$$w \leftarrow w - \alpha \frac{\partial J}{\partial w}$$

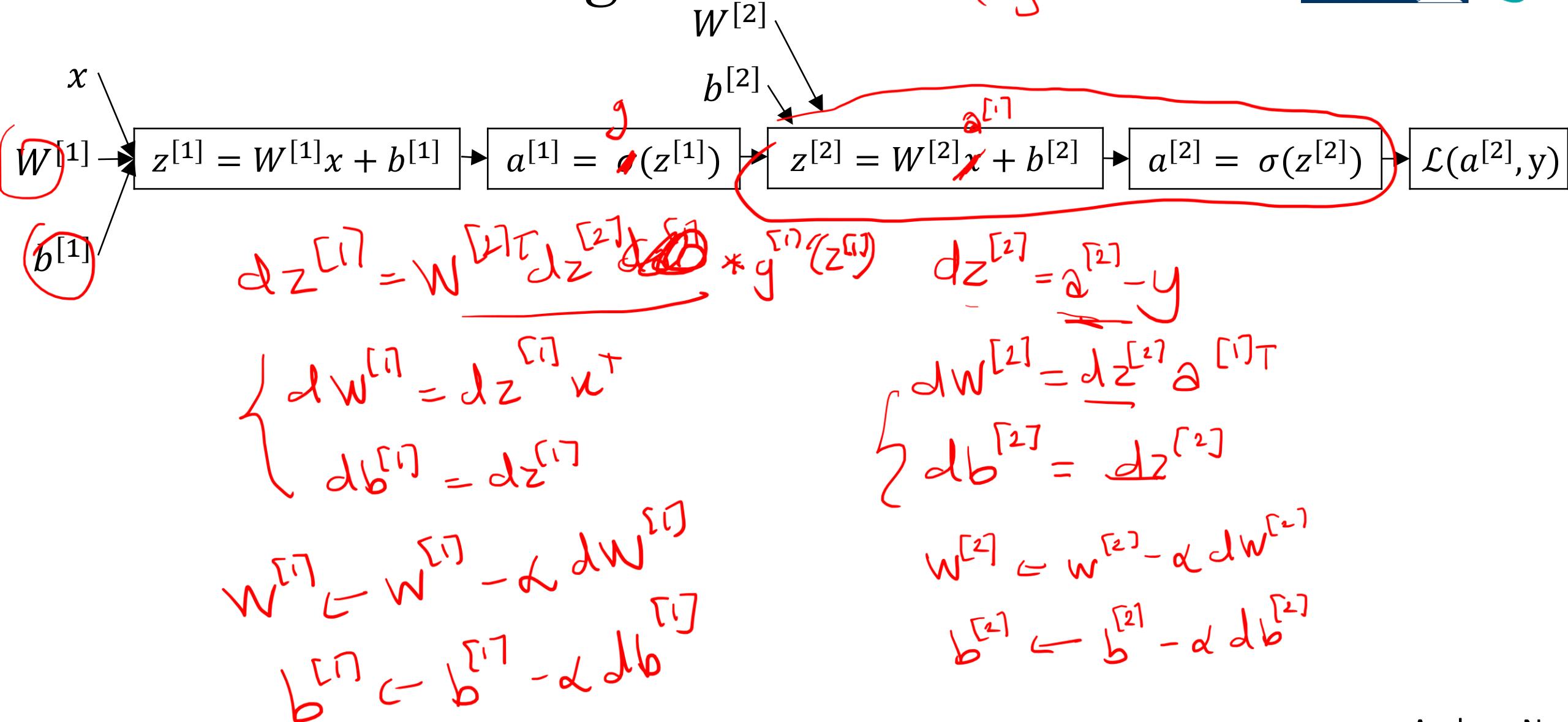
$$b \leftarrow b - \alpha \frac{\partial J}{\partial b}$$

$$\frac{\partial z}{\partial w} = \sigma(z)(1 - \sigma(z))x$$

$$\frac{\partial J}{\partial w} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m (\sigma(z) - y_i)x_i$$

$$\frac{\partial z}{\partial b} = \sigma(z)(1 - \sigma(z))$$

# Neural network gradients



# Summary of gradient descent

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

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

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

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]\prime}(z^{[1]})$$

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

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

# Summary of gradient descent

one training example  $\underline{k}$

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

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

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

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

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

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

vectorized implementation of one training example

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

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

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

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

$$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T \quad \frac{1}{m} \begin{bmatrix} dZ^{[1](1)} & \dots & dZ^{[1](m)} \end{bmatrix} \begin{bmatrix} u^{(1)} \\ \vdots \\ u^{(m)} \end{bmatrix}$$

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

$$= \frac{1}{m} \sum_{i=1}^m dZ^{[1](i)} u^{(i)}$$



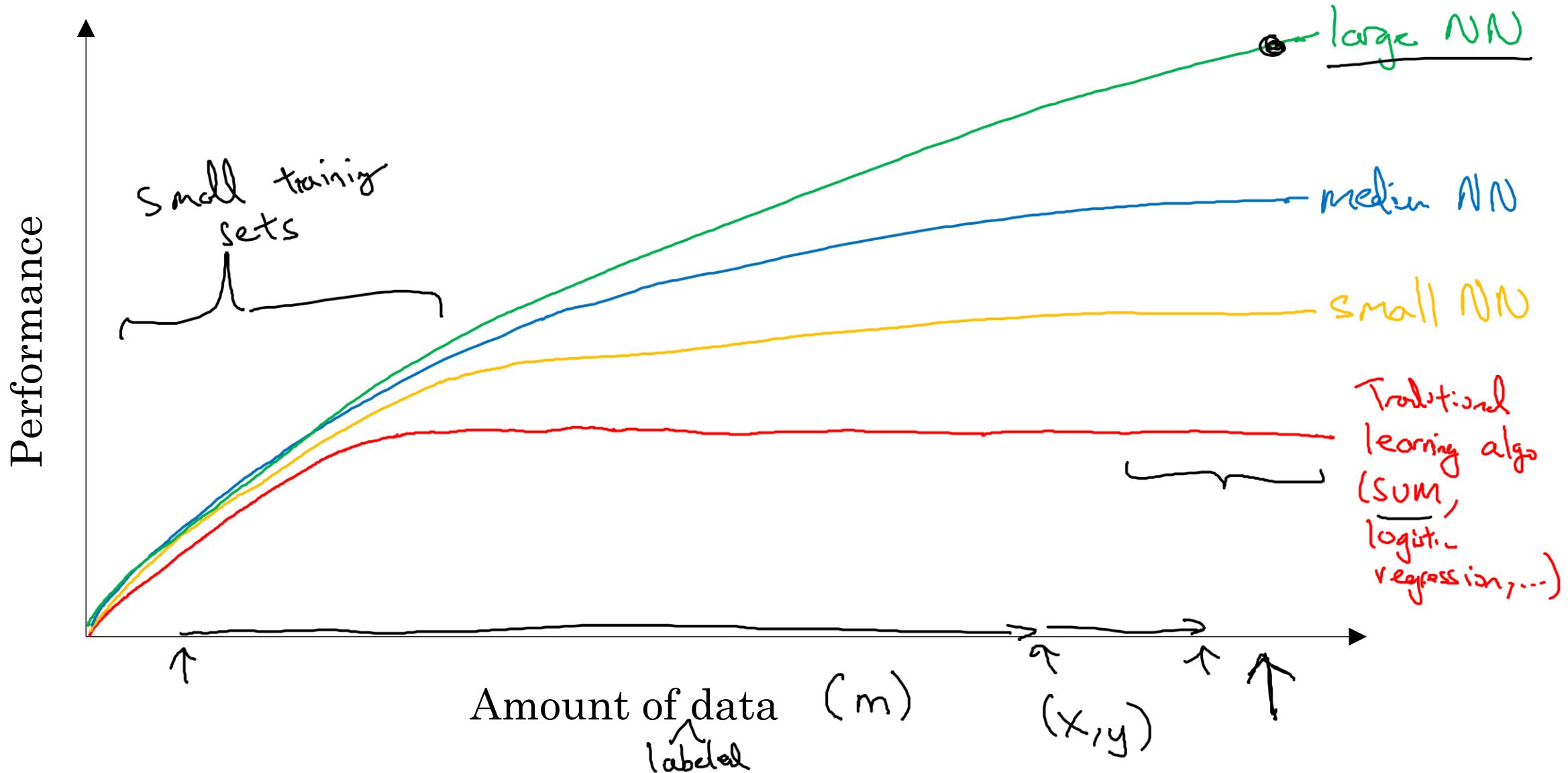
deeplearning.ai

# Introduction to Neural Networks

---

## Why is Deep Learning taking off?

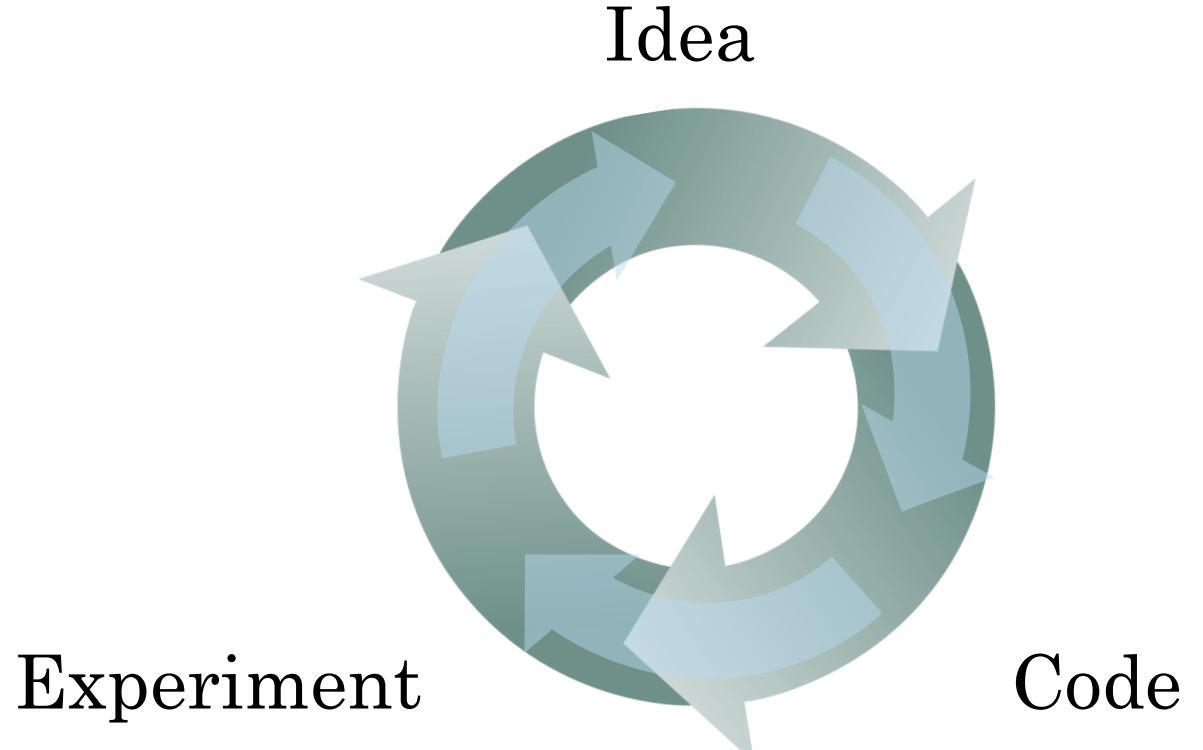
# Scale drives deep learning progress



# Scale drives deep learning progress



- Data
- Computation
- Algorithms





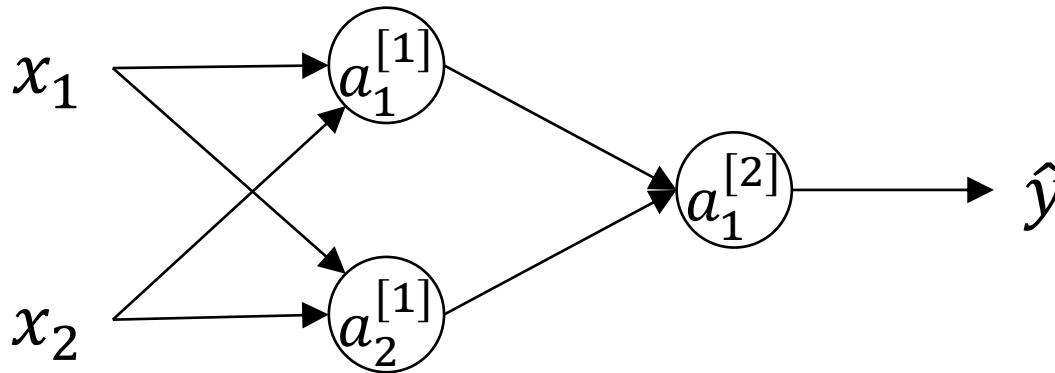
deeplearning.ai

# One hidden layer Neural Network

---

## Random Initialization

# What happens if you initialize weights to zero?



Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$

$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

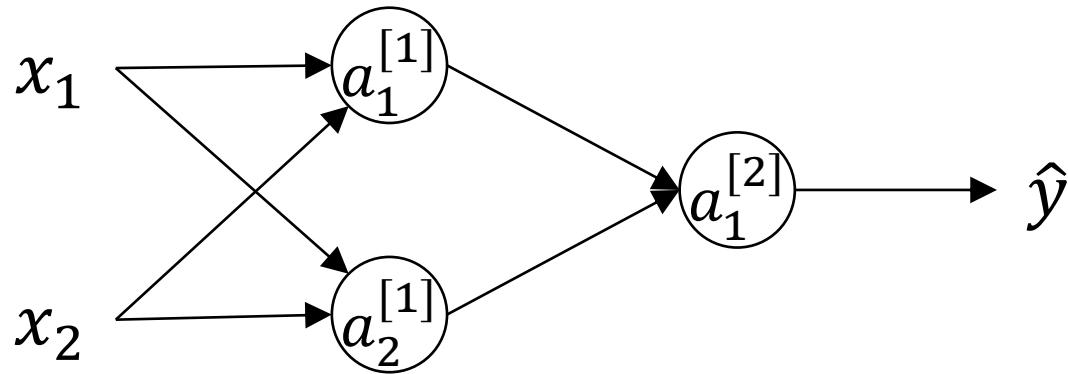
$$W^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix} \quad b^{[2]} = \begin{bmatrix} 0 \end{bmatrix}$$

$$\underbrace{a_1^{[1]} = a_2^{[1]} = g^{[1]}(0)}$$

$$\underbrace{dz_1^{[1]} = dz_2^{[1]} = 0}$$

$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix} \quad W^{[1]} \leftarrow W^{[1]} - dW$$

# Random initialization



$$W^{[1]} = \text{np.random.randn}((2, 2) * 0.01)$$

$$b^{[1]} = \text{np.zeros}((2, 1))$$

$$W^{[2]} = \text{np.random.randn}((1, 2) * 0.01)$$

$$b^{[2]} = \text{np.zeros}((1, 1))$$



deeplearning.ai

# Deep Neural Networks

---

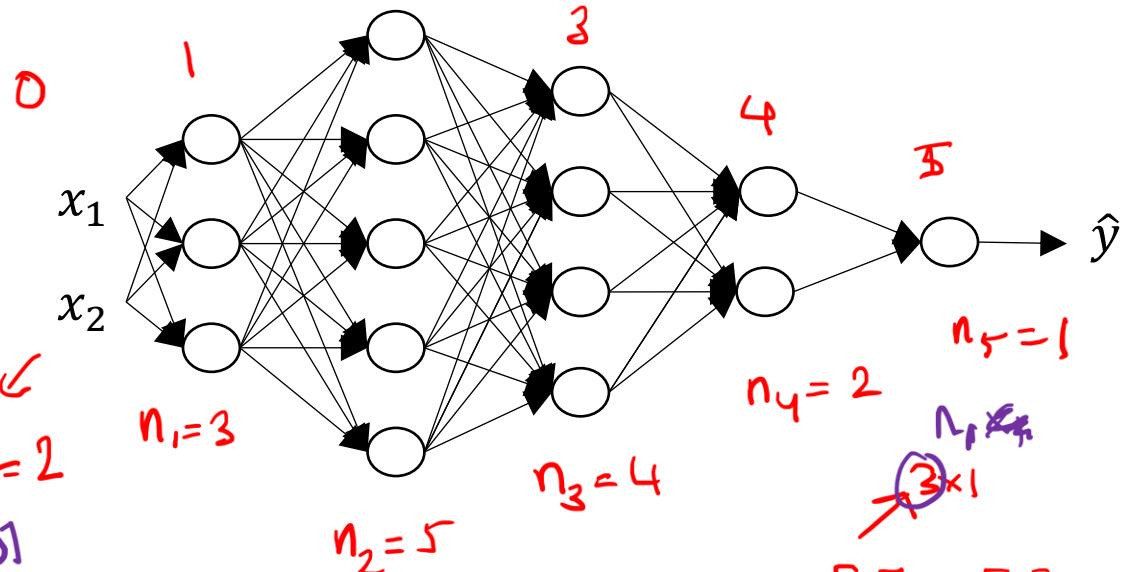
## Getting your matrix dimensions right

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

$$a^{[0]} = \nu$$

$$\nu = \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$$

$$n_0 = n_L = 2$$



$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

Dimensions:  $a^{[0]}: 3 \times 1$ ,  $W^{[1]}: 3 \times 2$ ,  $b^{[1]}: 2 \times 1$ ,  $z^{[1]}: 3 \times 1$

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

Dimensions:  $g^{[1]}: 3 \times 1$ ,  $a^{[1]}: 3 \times 1$

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

Dimensions:  $a^{[1]}: 3 \times 1$ ,  $W^{[2]}: 5 \times 3$ ,  $b^{[2]}: 5 \times 1$ ,  $z^{[2]}: 5 \times 1$

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

Dimensions:  $g^{[2]}: 5 \times 1$ ,  $a^{[2]}: 5 \times 1$

$i^{\text{th}} \text{ layer}$

$z^{[i]} \rightarrow n_i \times 1$

$a^{[i]} \rightarrow n_i \times 1$

$b^{[i]} \rightarrow n_i \times 1$

$W^{[i]} \rightarrow n_i \times n_{i-1}$

# Vectorized implementation

$$\{(x_i, y_i)\}_{i=1}^m$$

$n_i \times 1$

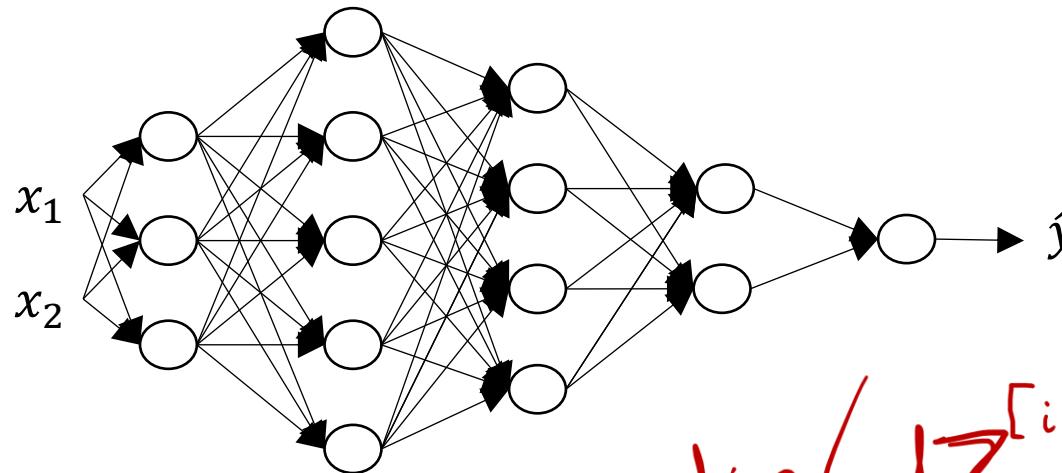
i<sup>th</sup> layer

$$Z^{[i]}$$

$n_i \times m$

$$A^{[i]} = [z^{[i](1)} \dots z^{[i](m)}]$$

$n_i \times m$



$$z^{[i]} = [z^{[i](1)} \dots z^{[i](m)}]$$

$$\dim(\underline{dZ}^{[i]}) = \dim(Z^{[i]}) \\ = n_i \times m$$

$$\dim(\underline{db}^{[i]}) = \dim(b^{[i]}) \\ = n_i \times 1$$



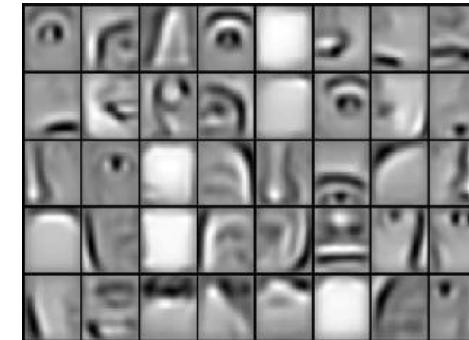
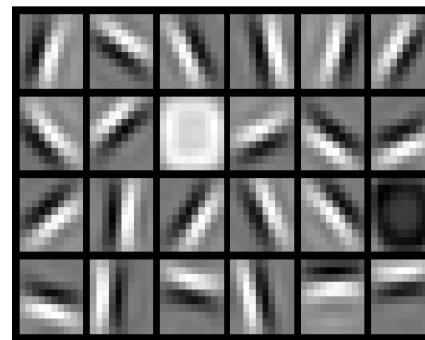
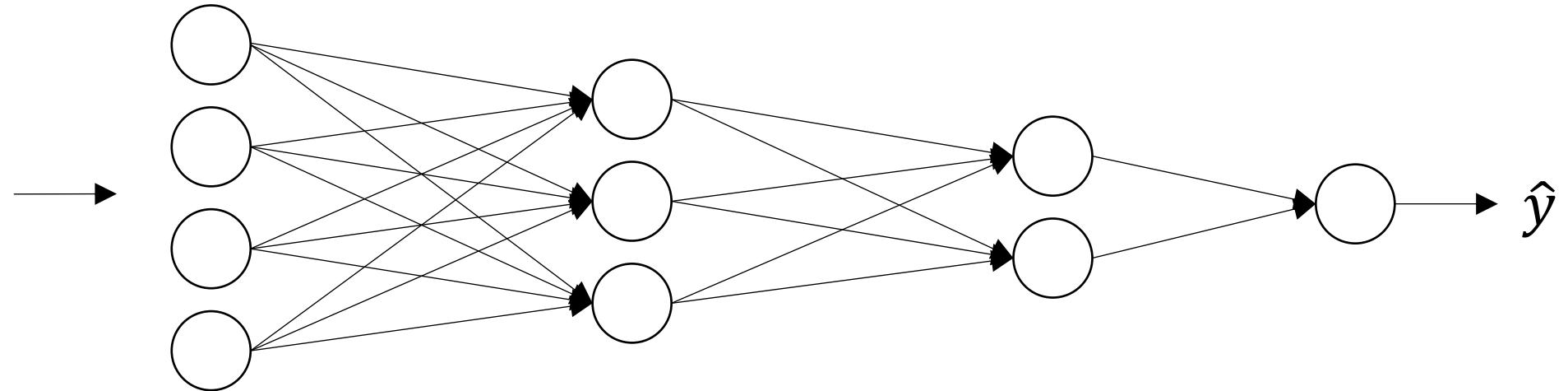
deeplearning.ai

# Deep Neural Networks

---

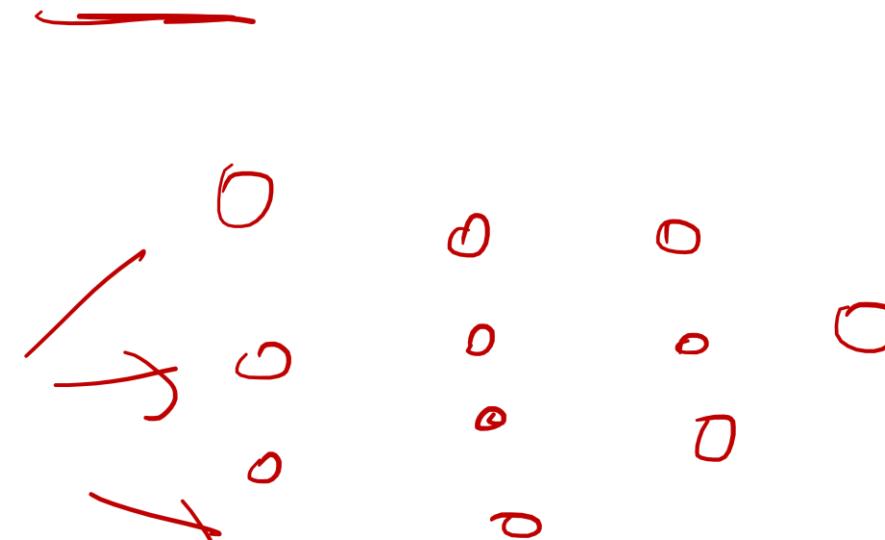
## Why deep representations?

# 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.





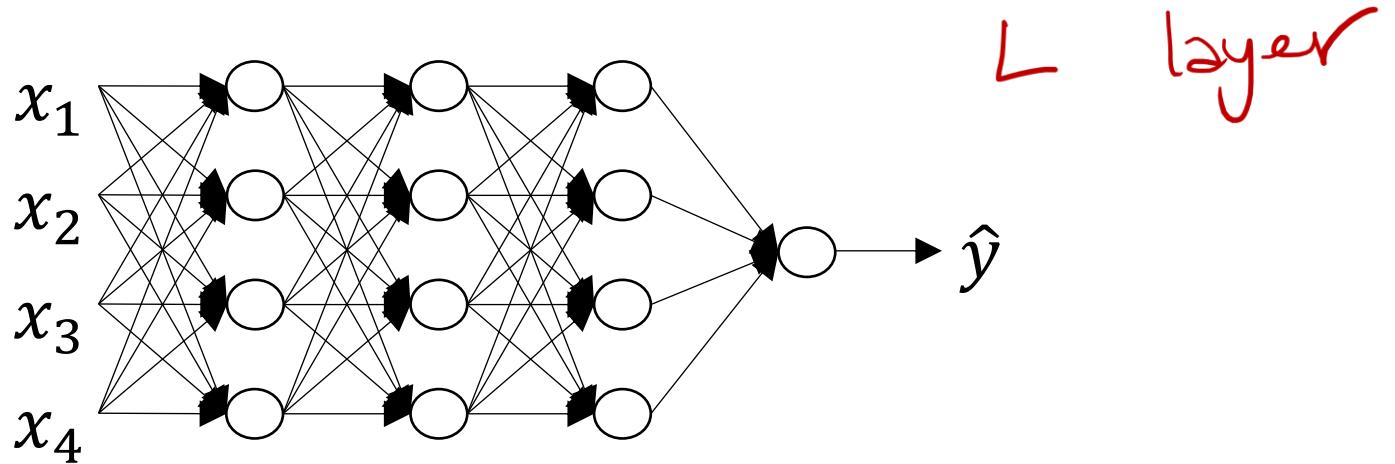
deeplearning.ai

# Deep Neural Networks

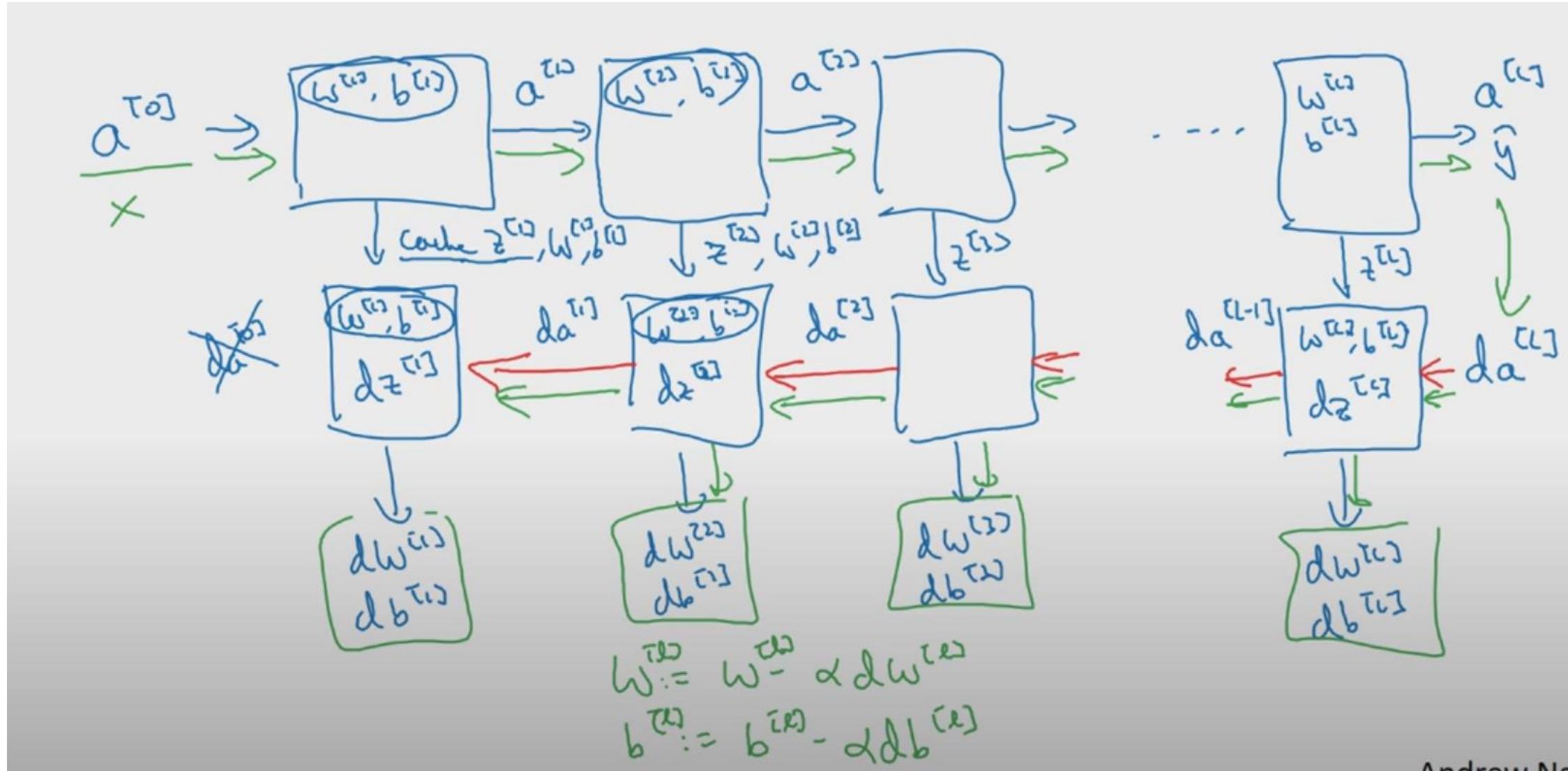
---

Building blocks of  
deep neural networks

# Forward and backward functions



# Forward and backward functions





deeplearning.ai

# Deep Neural Networks

---

## Forward and backward propagation

# Forward propagation for layer $l$

$L$

Input  $a^{[l-1]}$        $l^{\text{th}}$  layer

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

one sample

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

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

vectorized

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

$$A^{[l]} = g^{[l]}(Z^{[l]})$$

# Backward propagation for layer $l$

Input  $da^{[l]}$

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

*single sample*

$$dz^{[l]} = da^{[l]} * g^{[l]'}(z^{[l]})$$

$$dW^{[l]} = dz^{[l]} \cdot \underline{a^{[l-1]}}$$

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

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

*Vectorized*

$$dz^{[l]} = \underline{dA^{[l]}} * g^{[l]'}(z^{[l]})$$

$$dW^{[l]} = \left(\frac{1}{m}\right) dz^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} np \cdot \text{sum}(dz^{[l]}, \text{axis}=1, \text{keepdims=True})$$

$$dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

# Summary

