

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



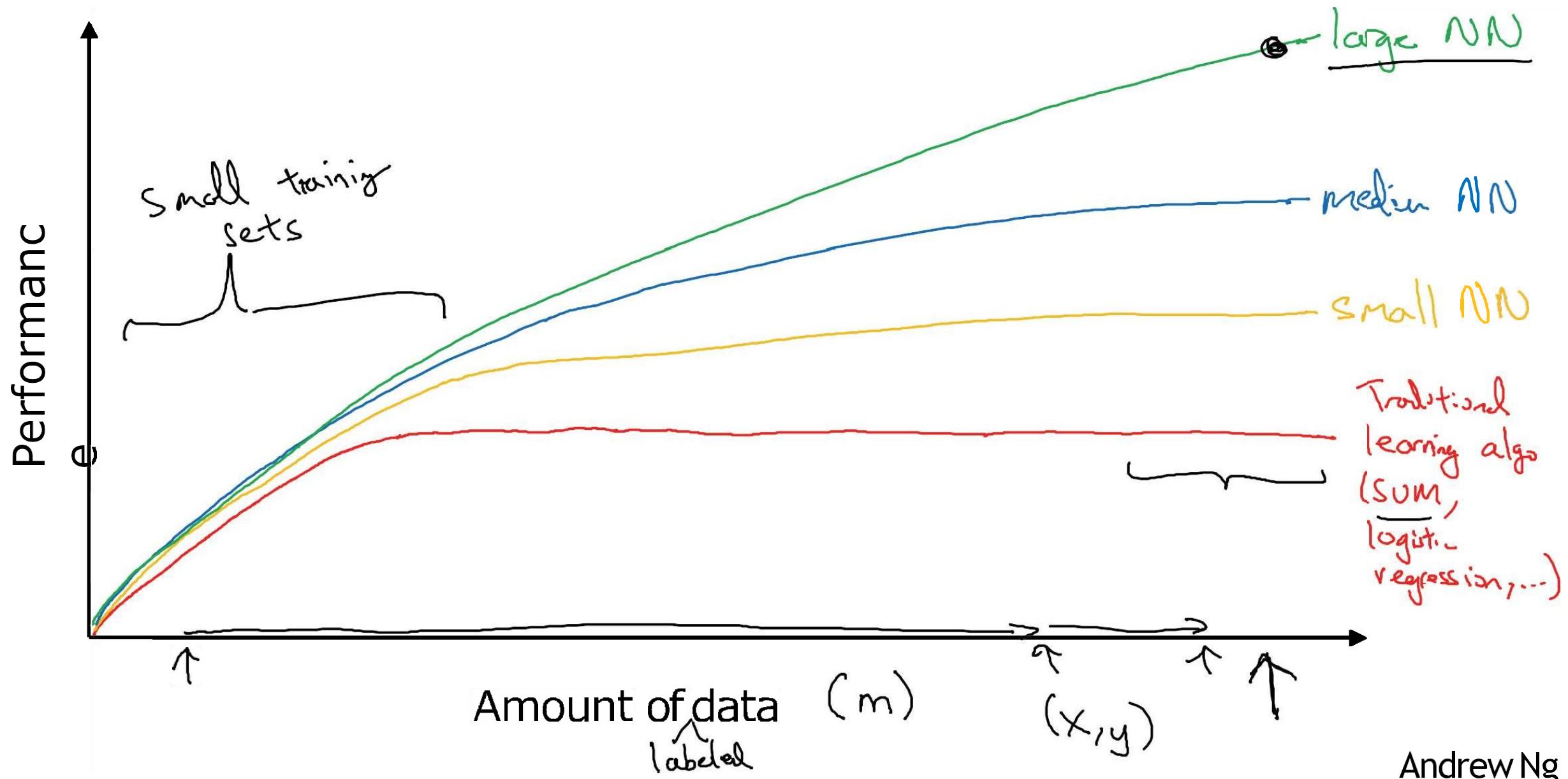
deeplearning.ai

# Introduction to Neural Networks

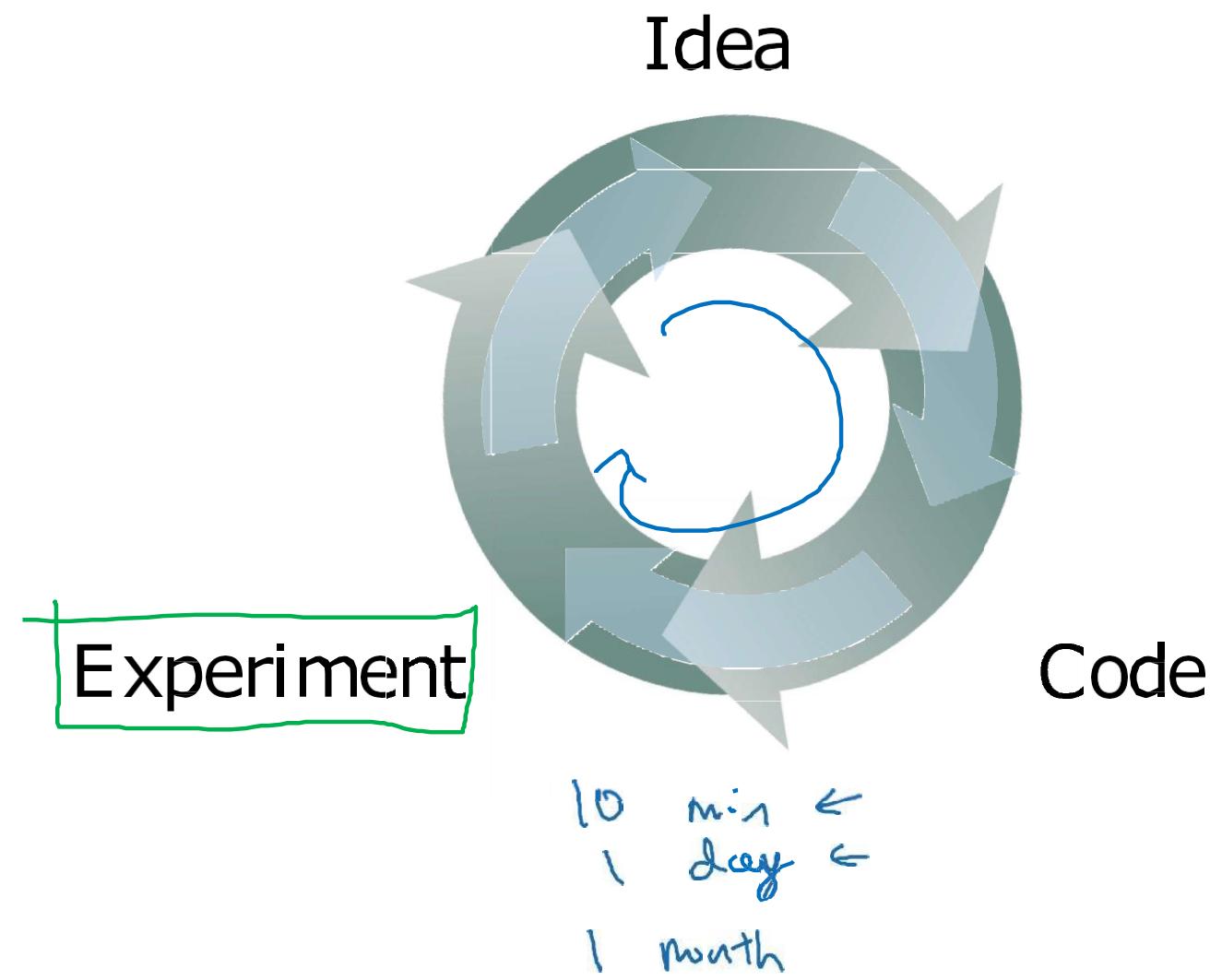
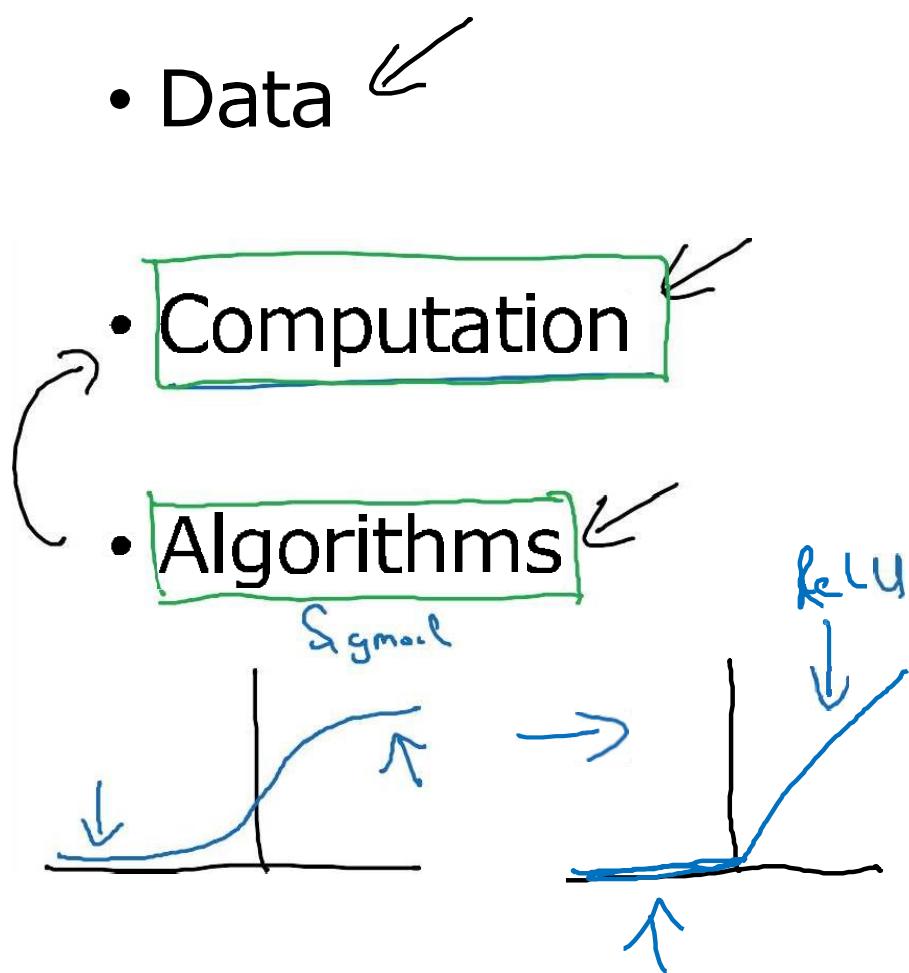
---

## Why is Deep Learning taking off?

# Scale drives deep learning progress



# Scale drives deep learning progress



Andrew Ng



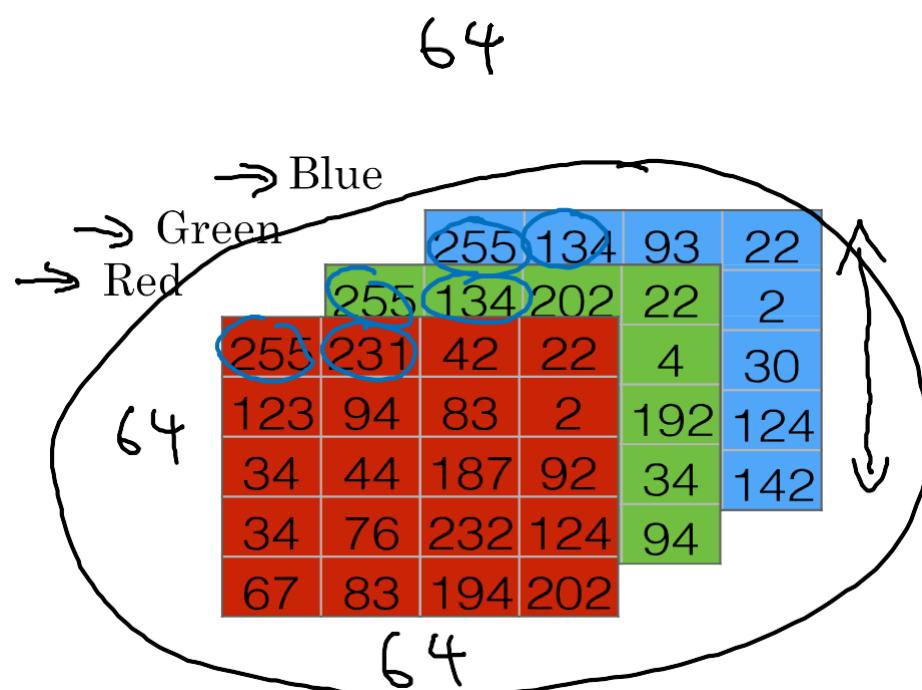
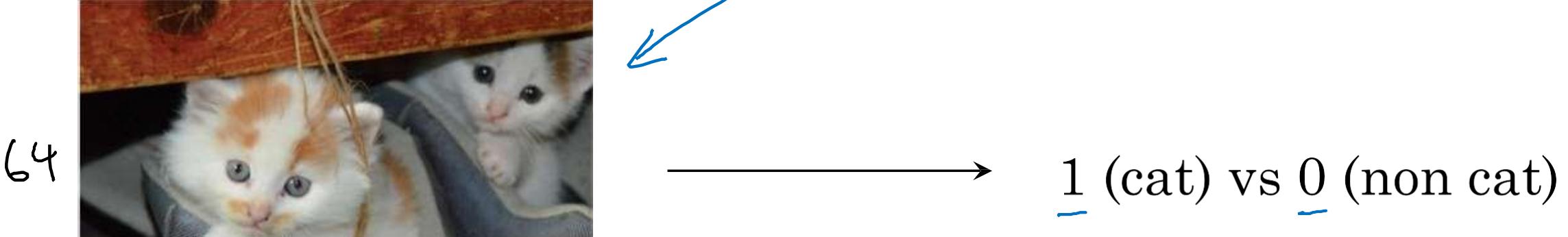
deeplearning.ai

# Basics of Neural Network Programming

---

## Binary Classification

# Binary Classification



$$X = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix}$$

$$64 \times 64 \times 3 = 12288$$

$$n = n_x = 12288$$

$$X \longrightarrow y$$

# Notation

$$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$$

$m$  training examples :  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$$M = M_{\text{train}}$$

$$M_{\text{test}} = \# \text{test examples.}$$

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

$X \in \mathbb{R}^{n_x \times m}$        $X.\text{shape} = (n_x, m)$

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$Y \in \mathbb{R}^{1 \times m}$$

$$Y.\text{shape} = (1, m)$$



deeplearning.ai

# Basics of Neural Network Programming

---

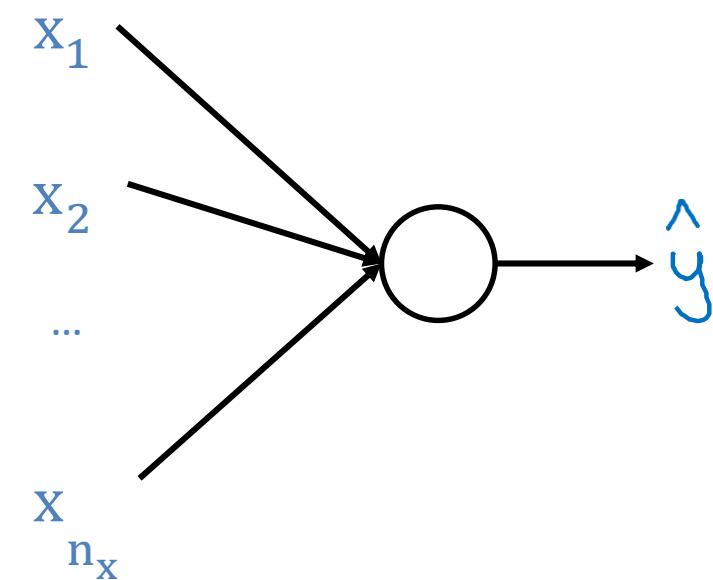
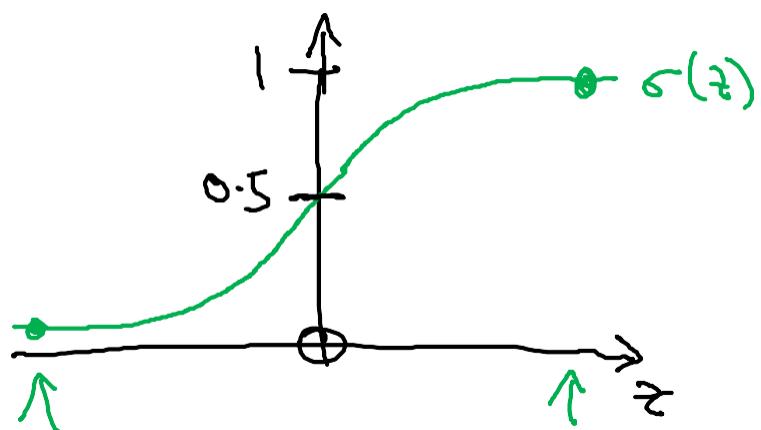
## Logistic Regression

# Logistic Regression

Given  $x$ , want  $\hat{y} = P(y=1|x)$   
 $x \in \mathbb{R}^{n_x}$

Parameters:  $\underline{\omega} \in \mathbb{R}^{n_x}$ ,  $b \in \mathbb{R}$ .

Output  $\hat{y} = \sigma(\underbrace{\omega^T x + b}_z)$



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

If  $z$  large  $\sigma(z) \approx \frac{1}{1+0} = 1$

If  $z$  large negative number

$$\sigma(z) = \frac{1}{1+e^{-z}} \approx \frac{1}{1+\text{Bignum}} \approx 0$$

Andrew Ng



deeplearning.ai

# Basics of Neural Network Programming

---

## Logistic Regression cost function

# Logistic Regression cost function

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}} \quad z^{(i)} = w^T x^{(i)} + b$$

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , want  $\hat{y}^{(i)} \approx y^{(i)}$ .

$x^{(i)}$   
 $y^{(i)}$   
 $z^{(i)}$

i-th  
example.

**Loss** (error) function:

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

If  $y=1$ :  $L(\hat{y}, y) = - \log \hat{y} \leftarrow$  want  $\log \hat{y}$  large, want  $\hat{y}$  large

If  $y=0$ :  $L(\hat{y}, y) = - \log(1-\hat{y}) \leftarrow$  want  $\log 1-\hat{y}$  large ... want  $\hat{y}$  small

**Cost** function:  $J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$



deeplearning.ai

# Basics of Neural Network Programming

---

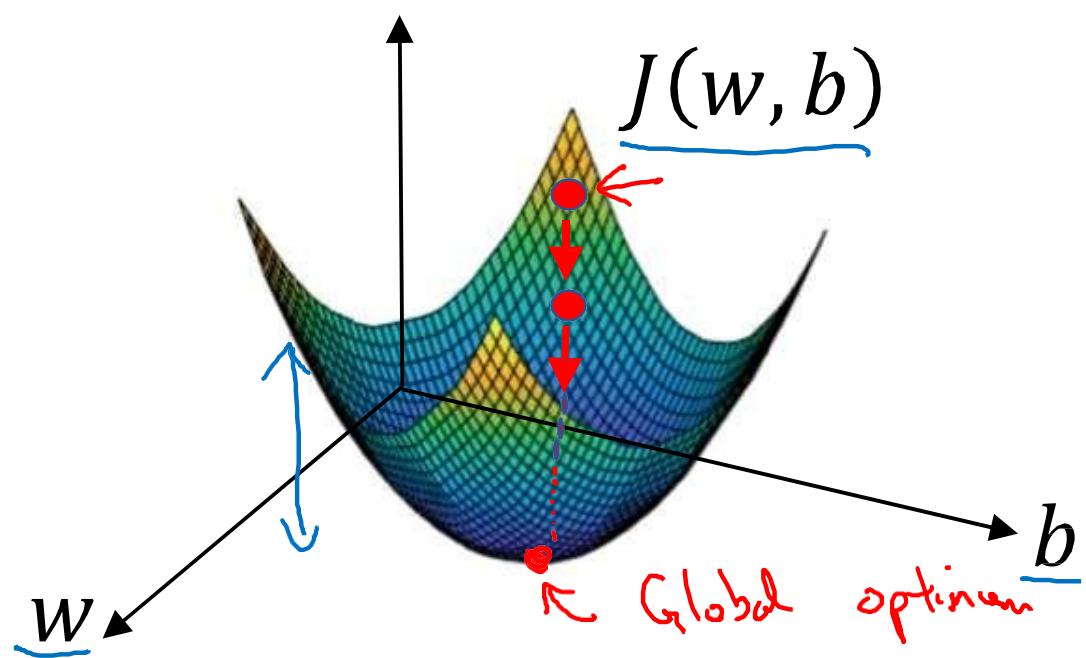
## Gradient Descent

# Gradient Descent

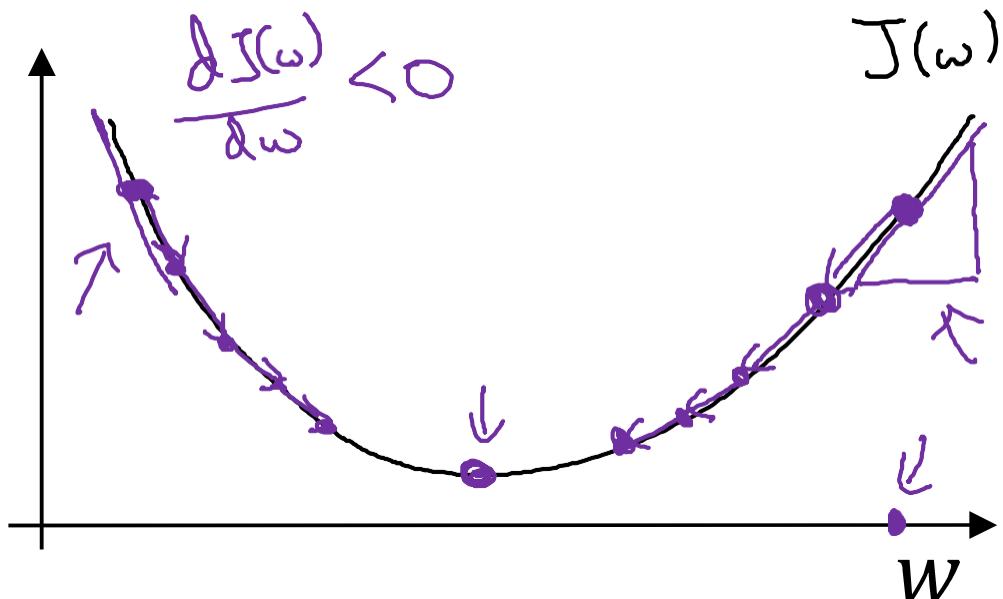
Recap:  $\hat{y} = \sigma(w^T x + b)$ ,  $\sigma(z) = \frac{1}{1+e^{-z}}$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

Want to find  $w, b$  that minimize  $J(w, b)$



# Gradient Descent



Repeat {  
 $w := w - \alpha$   
 $w' = w - \underline{\alpha \frac{dJ(w)}{dw}}$   
 $\frac{dJ(w)}{dw} = ?$

↑ ↑ learning rate "dw"

---

$J(w, b)$

$w := w - \alpha \frac{dJ(w, b)}{dw}$

$b := b - \alpha \frac{dJ(w, b)}{db}$

$\frac{dJ(w, b)}{dw}$        $\frac{dJ(w, b)}{db}$

"partial derivative"  $J$

$dJ(w, b)$        $dJ(w, b)$

$dw$        $db$



deeplearning.ai

# Basics of Neural Network Programming

---

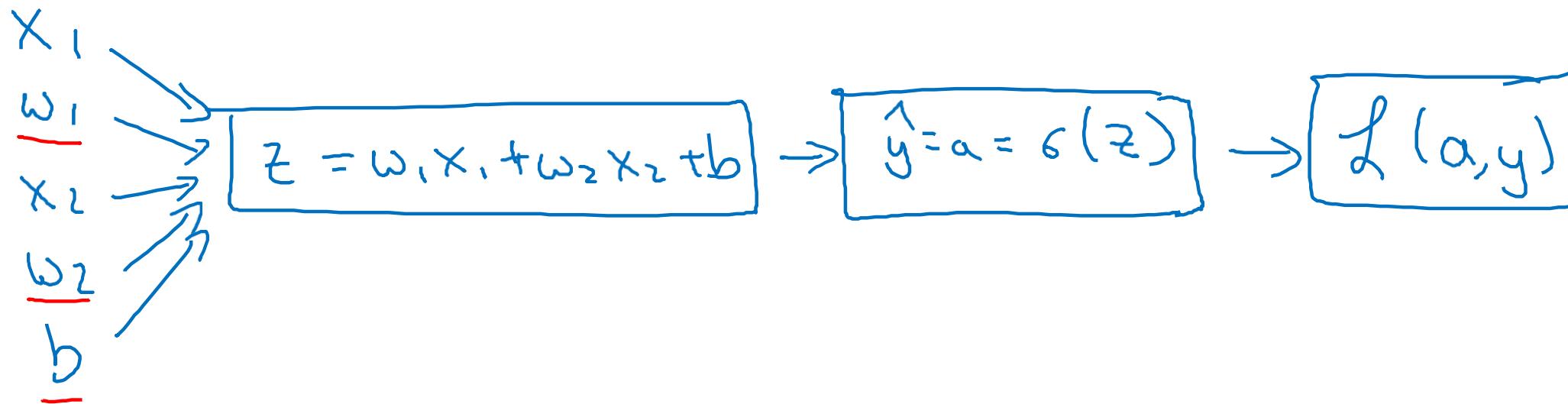
## Logistic Regression Gradient descent

# Logistic regression recap

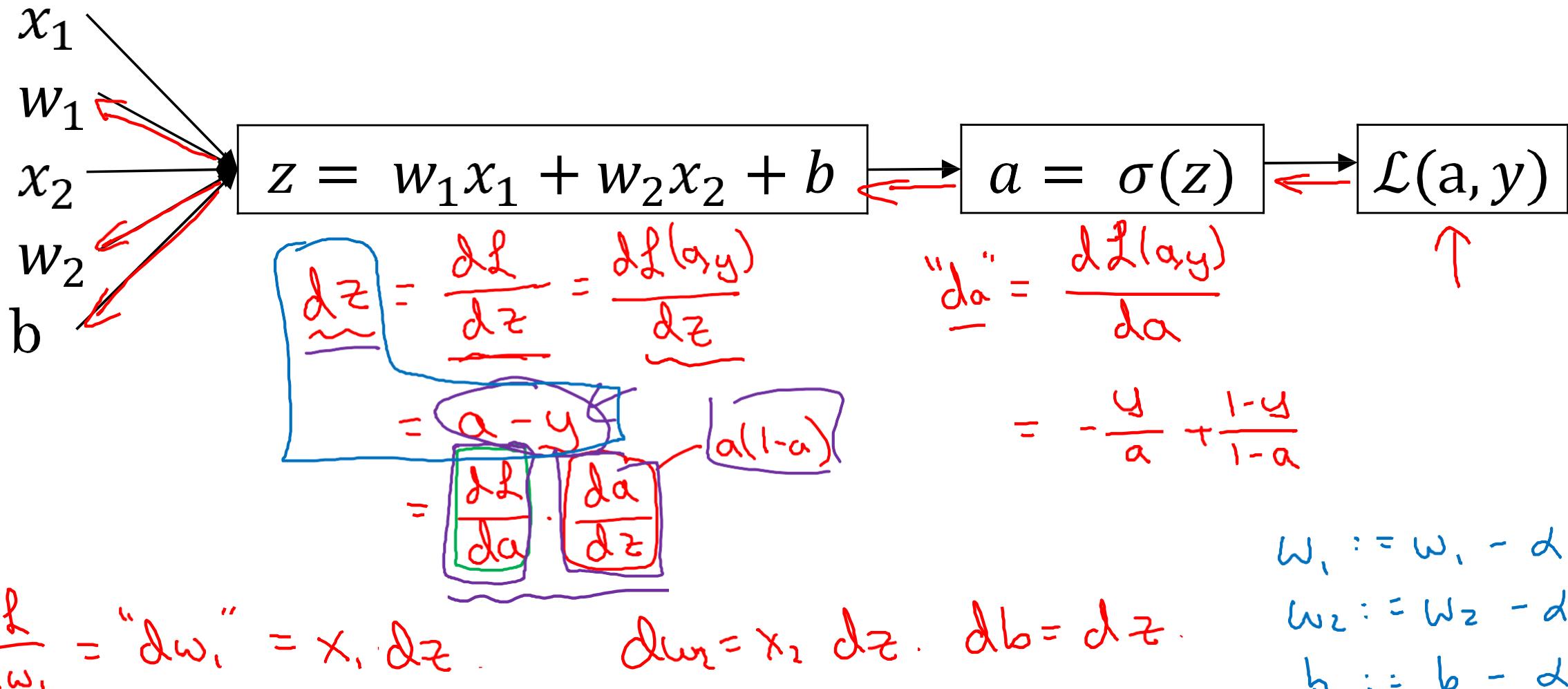
$$z = w^T x + b$$

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

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$



# Logistic regression derivatives





deeplearning.ai

# Basics of Neural Network Programming

---

Gradient descent  
on  $m$  examples

# Logistic regression on $m$ examples

$$\underline{J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m l(a^{(i)}, y^{(i)})$$
$$a^{(i)} = \hat{y}^{(i)} = g(z^{(i)}) = g(\omega^\top x^{(i)} + b)$$
$$(x^{(i)}, y^{(i)})$$
$$\underline{dw_1^{(i)}}, \underline{dw_2^{(i)}}, \underline{db^{(i)}}$$

$$\underline{\frac{\partial}{\partial \omega_j} J(\omega, b)} = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\partial}{\partial \omega_j} l(a^{(i)}, y^{(i)})}_{dw_j^{(i)}} - (x^{(i)}, y^{(i)})$$

# Logistic regression on $m$ examples

$$J = 0; \underline{\partial w_1} = 0; \underline{\partial w_2} = 0; \underline{\partial b} = 0$$

$\rightarrow$  For  $i = 1$  to  $m$

$$z^{(i)} = \omega^\top x^{(i)} + b$$

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

$$J_t = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)})]$$

$$\underline{\partial z^{(i)}} = a^{(i)} - y^{(i)}$$

$$\begin{aligned} \uparrow & \left[ \begin{array}{l} \partial w_1 += x_1^{(i)} \partial z^{(i)} \\ \partial w_2 += x_2^{(i)} \partial z^{(i)} \\ \partial b += \partial z^{(i)} \end{array} \right] \quad \downarrow n=2 \\ \partial w_3 \\ \vdots \\ \partial w_n \end{aligned}$$

$$J / m \leftarrow$$

$$\partial w_1 / m; \quad \partial w_2 / m; \quad \partial b / m. \quad \leftarrow$$

$$\partial w_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \underline{\partial w_1}$$

$$w_2 := w_2 - \alpha \underline{\partial w_2}$$

$$b := b - \alpha \underline{\partial b}$$

Vectorization



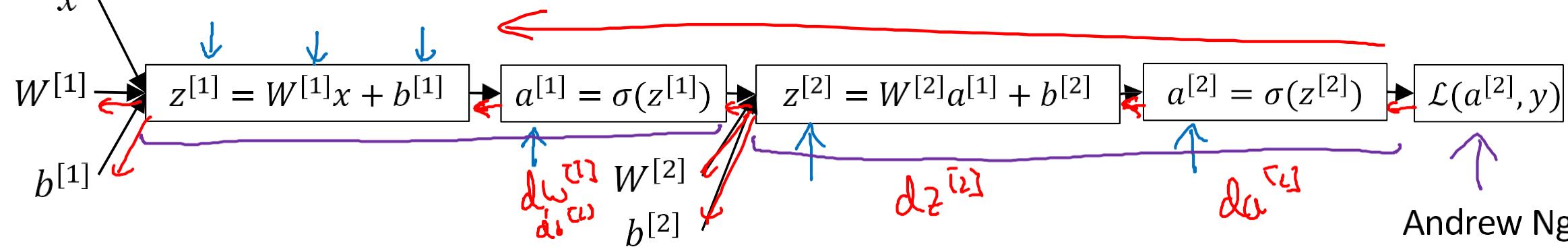
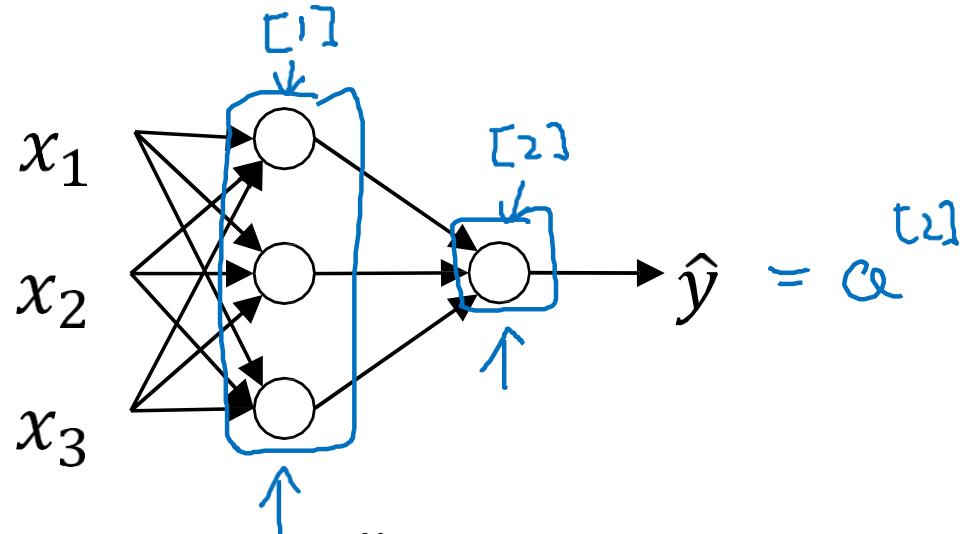
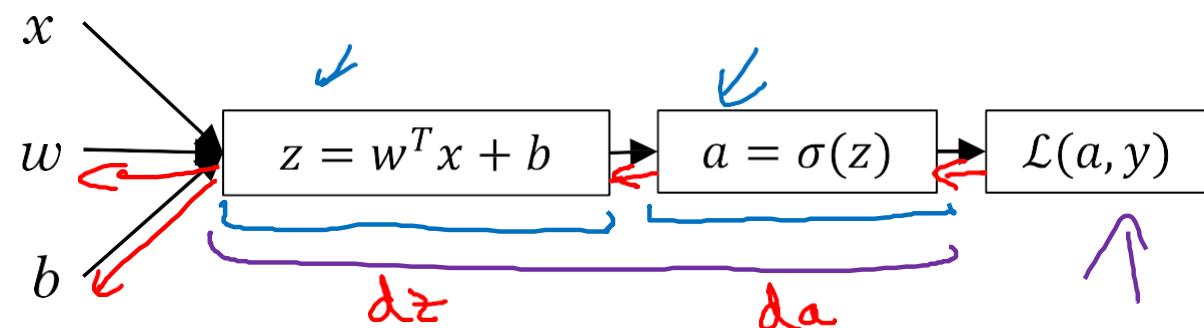
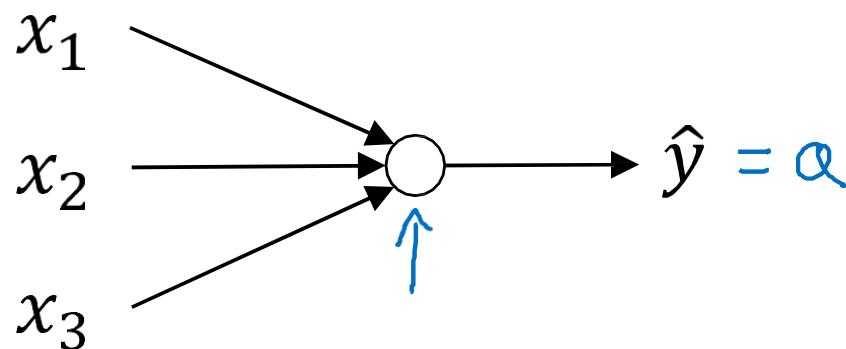
deeplearning.ai

One hidden layer  
Neural Network

---

Neural Networks  
Overview

# What is a Neural Network?



Andrew Ng



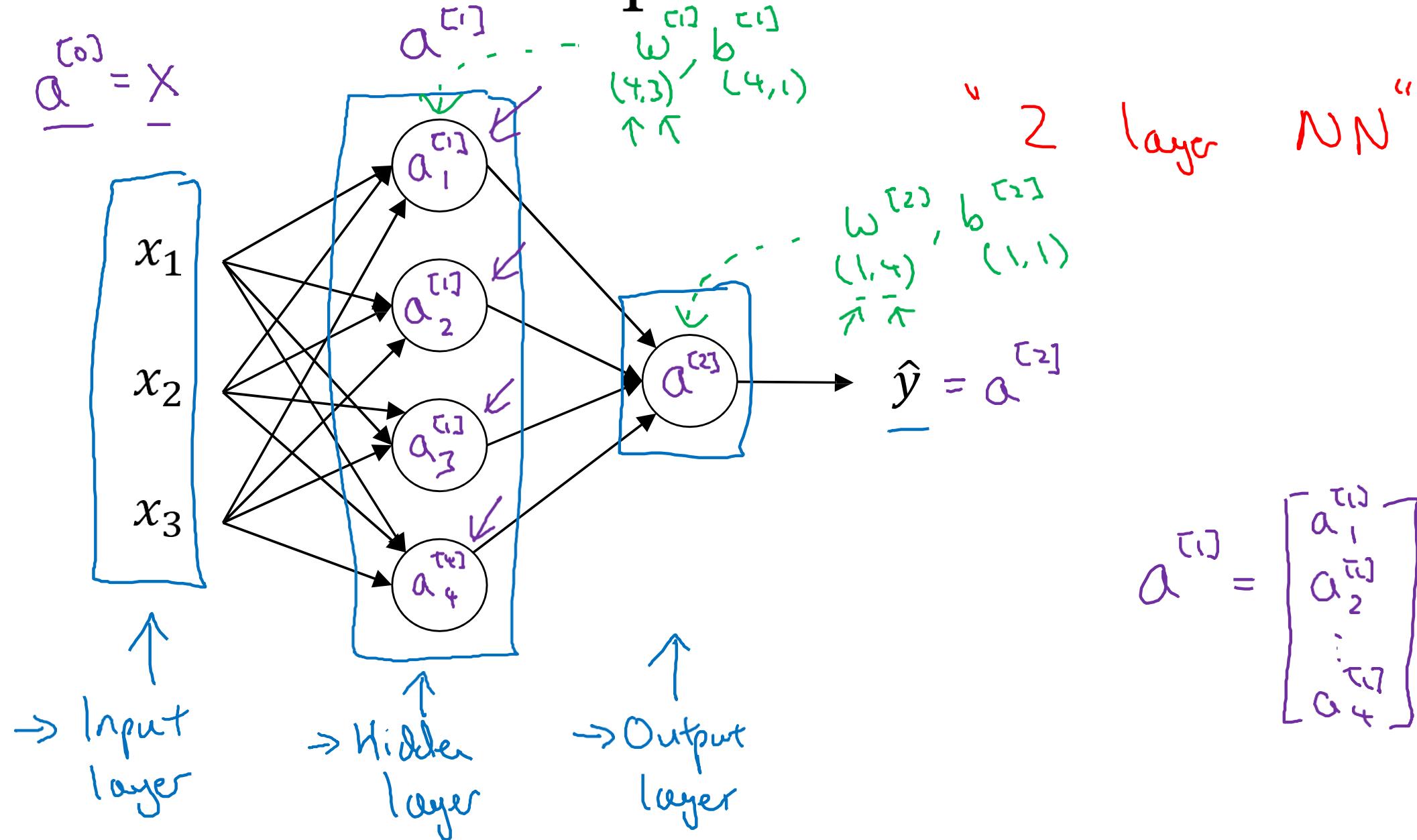
deeplearning.ai

One hidden layer  
Neural Network

---

Neural Network  
Representation

# Neural Network Representation





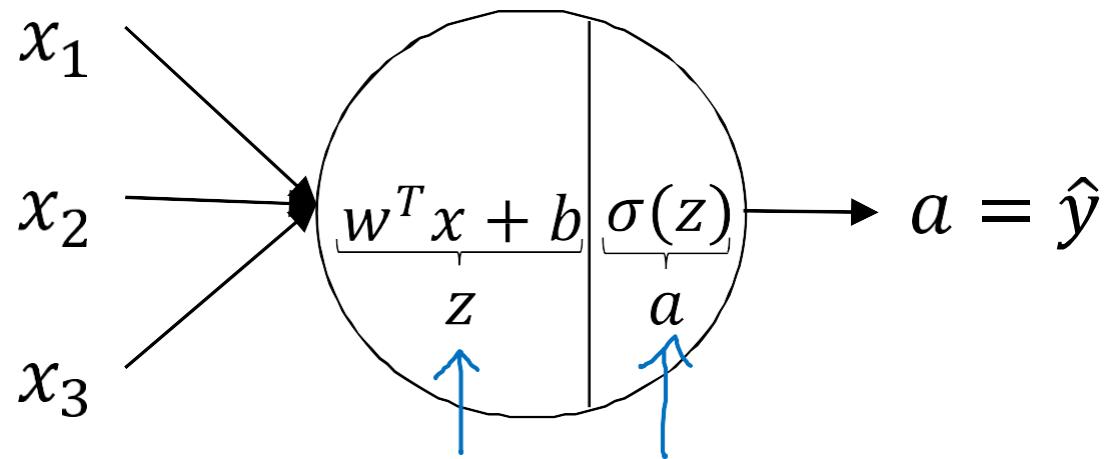
deeplearning.ai

One hidden layer  
Neural Network

---

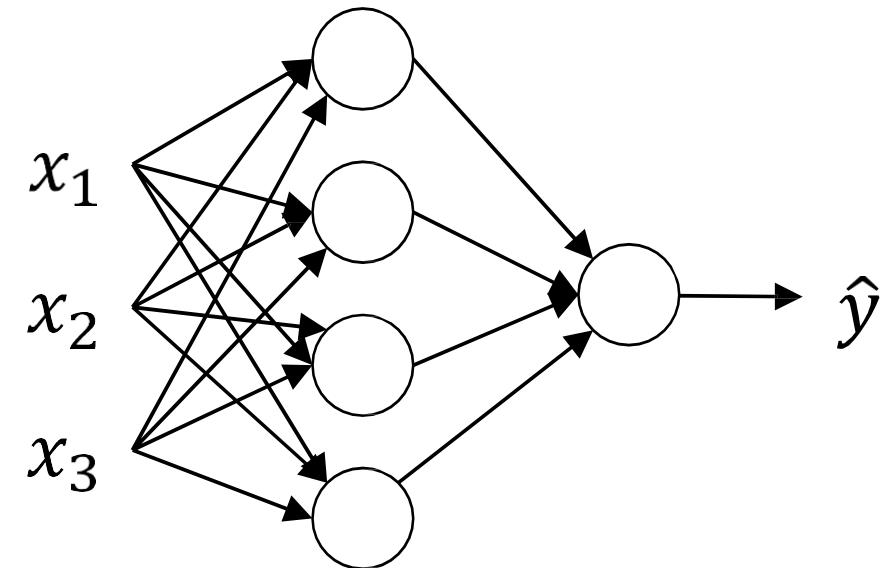
Computing a  
Neural Network's  
Output

# Neural Network Representation

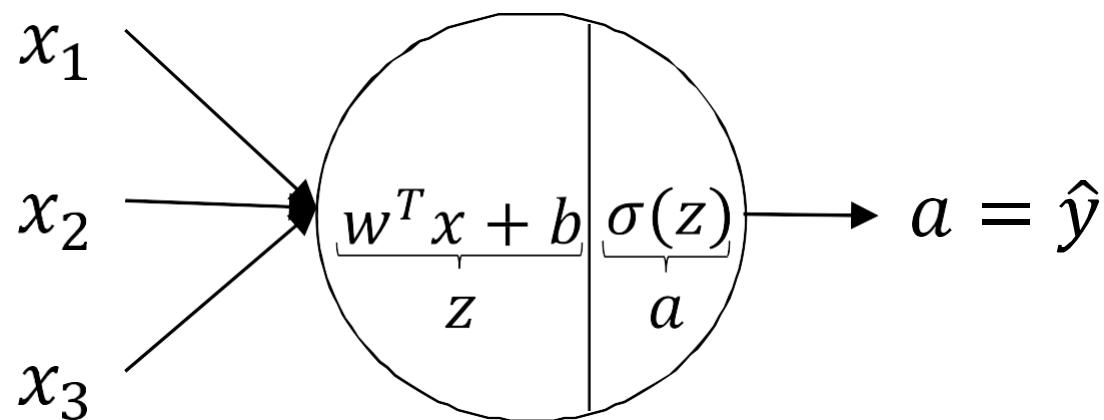


$$z = w^T x + b$$

$$a = \sigma(z)$$

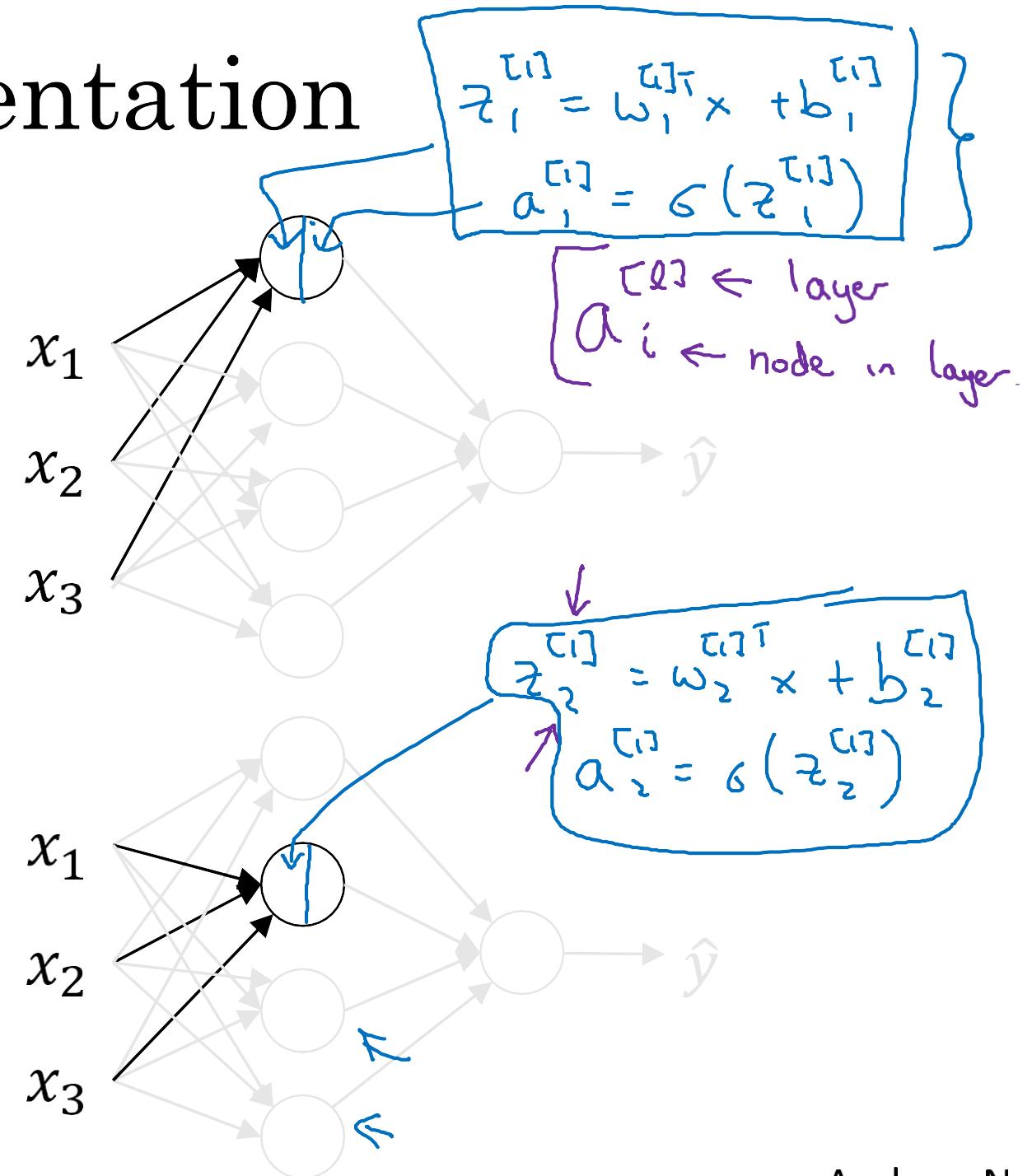


# Neural Network Representation



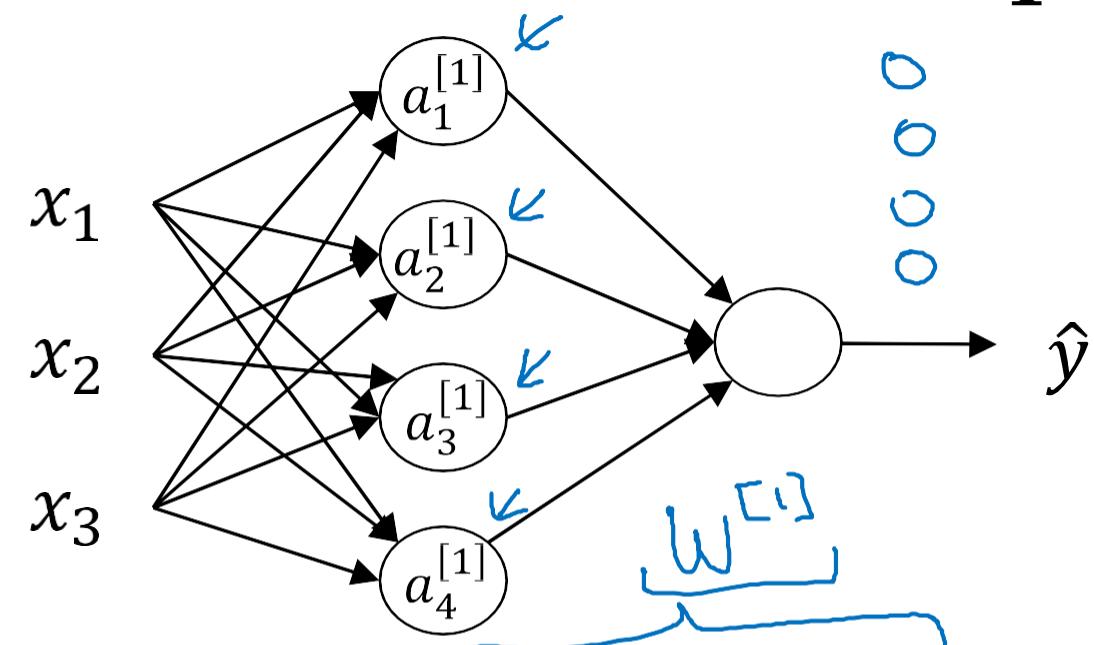
$$z = w^T x + b$$

$$a = \sigma(z)$$



Andrew Ng

# Neural Network Representation



$$\rightarrow z^{[1]} = \begin{bmatrix} -w_1^{T,1} \\ -w_2^{T,1} \\ -w_3^{T,1} \\ -w_4^{T,1} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = g(z^{[1]})$$

$$\rightarrow a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ \vdots \\ a_4^{[1]} \end{bmatrix} = g(z^{[1]})$$

Handwritten notes for the forward pass calculation:

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

$$a_3^{[1]} = \sigma(z_3^{[1]})$$

$$a_4^{[1]} = \sigma(z_4^{[1]})$$

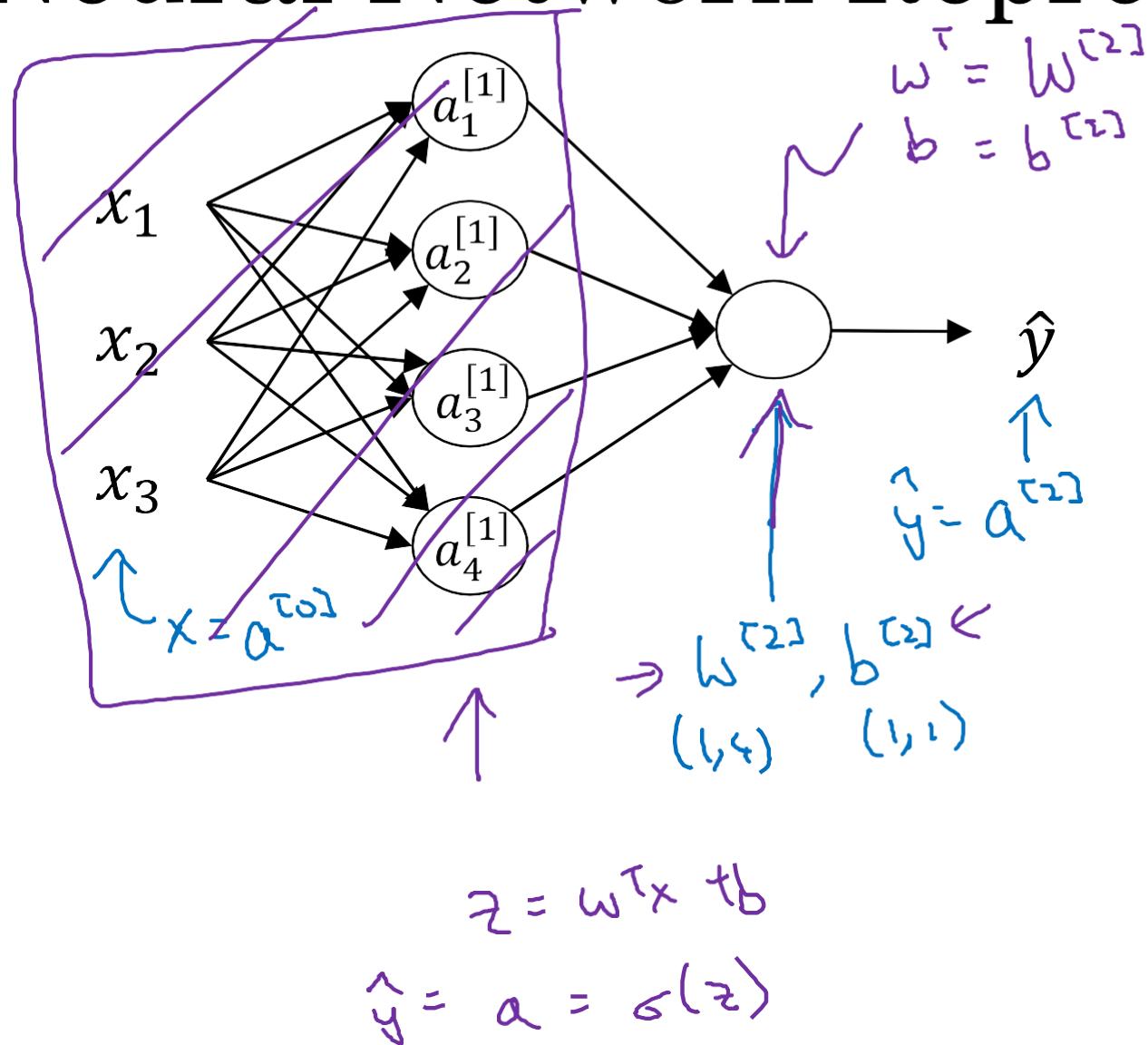
$$(w_i^{[1]T} x + b_i^{[1]}) \quad \sigma(\cdot)$$

$$z_1^{[1]} \quad z_2^{[1]} \quad z_3^{[1]} \quad z_4^{[1]}$$

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

Andrew Ng

# Neural Network Representation learning



Given input  $x$ :

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

$(4, 1) \quad (4, 3) \quad (3, 1) \quad (4, 1)$

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

$(4, 1) \quad (4, 1)$

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

$(1, 1) \quad (1, 4) \quad (4, 1) \quad (1, 1)$

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

$(1, 1) \quad (1, 1)$



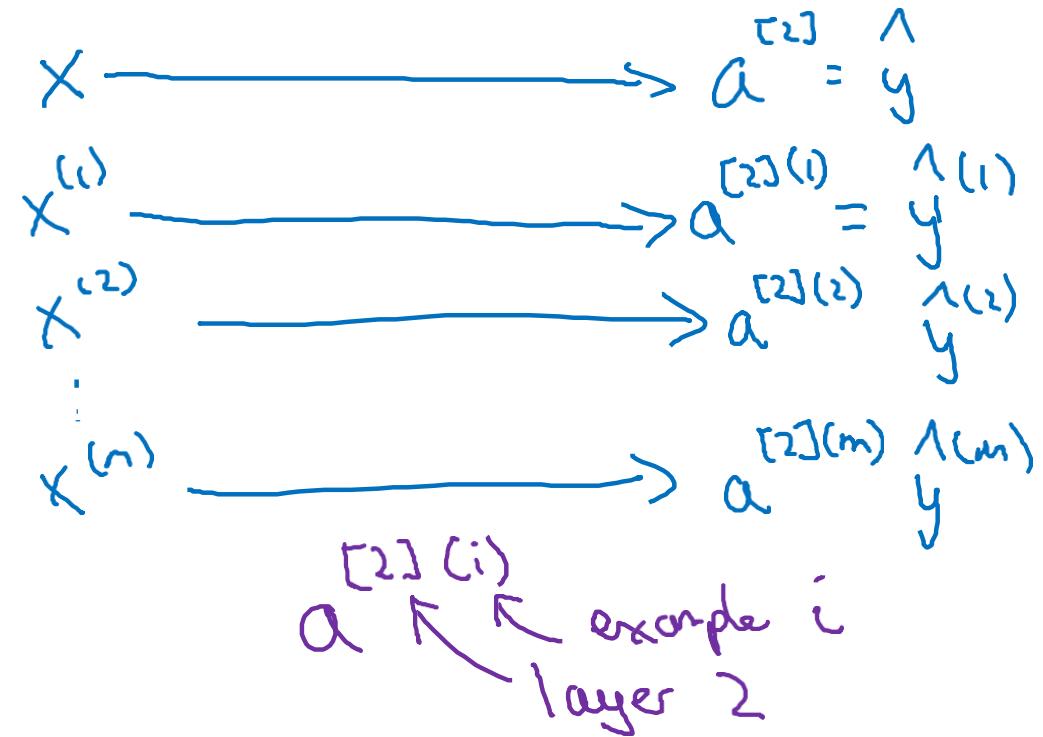
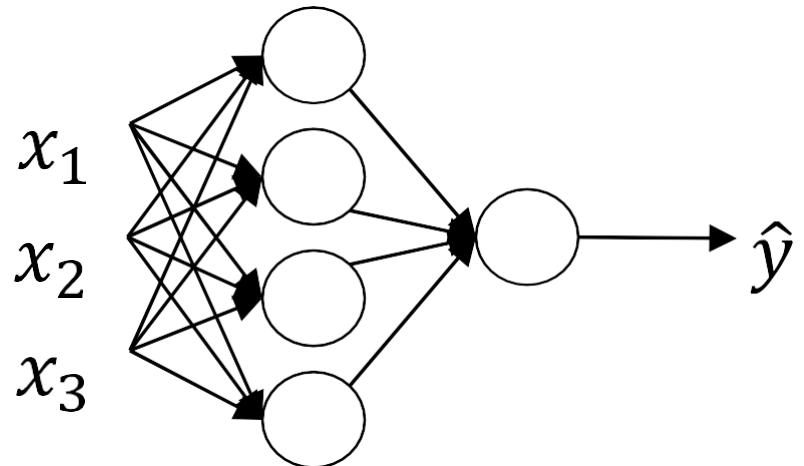
deeplearning.ai

One hidden layer  
Neural Network

---

Vectorizing across  
multiple examples

# Vectorizing across multiple examples



$$\left\{ \begin{array}{l} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = \sigma(z^{[1]}) \\ z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = \sigma(z^{[2]}) \end{array} \right.$$

for  $i = 1$  to  $n$ ,

$$\begin{aligned} z^{[1](i)} &= w^{[1]}x^{(i)} + b^{[1]} \\ a^{[1](i)} &= \sigma(z^{[1](i)}) \\ z^{[2](i)} &= w^{[2]}a^{[1](i)} + b^{[2]} \\ a^{[2](i)} &= \sigma(z^{[2](i)}) \end{aligned}$$

Andrew Ng

# Vectorizing across multiple examples

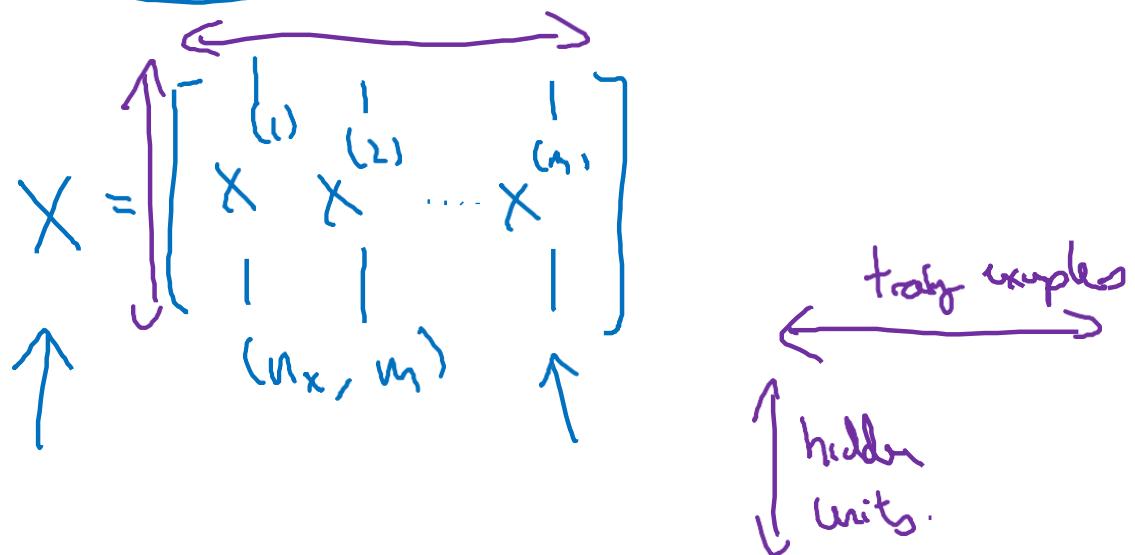
for i = 1 to m:

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

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

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

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

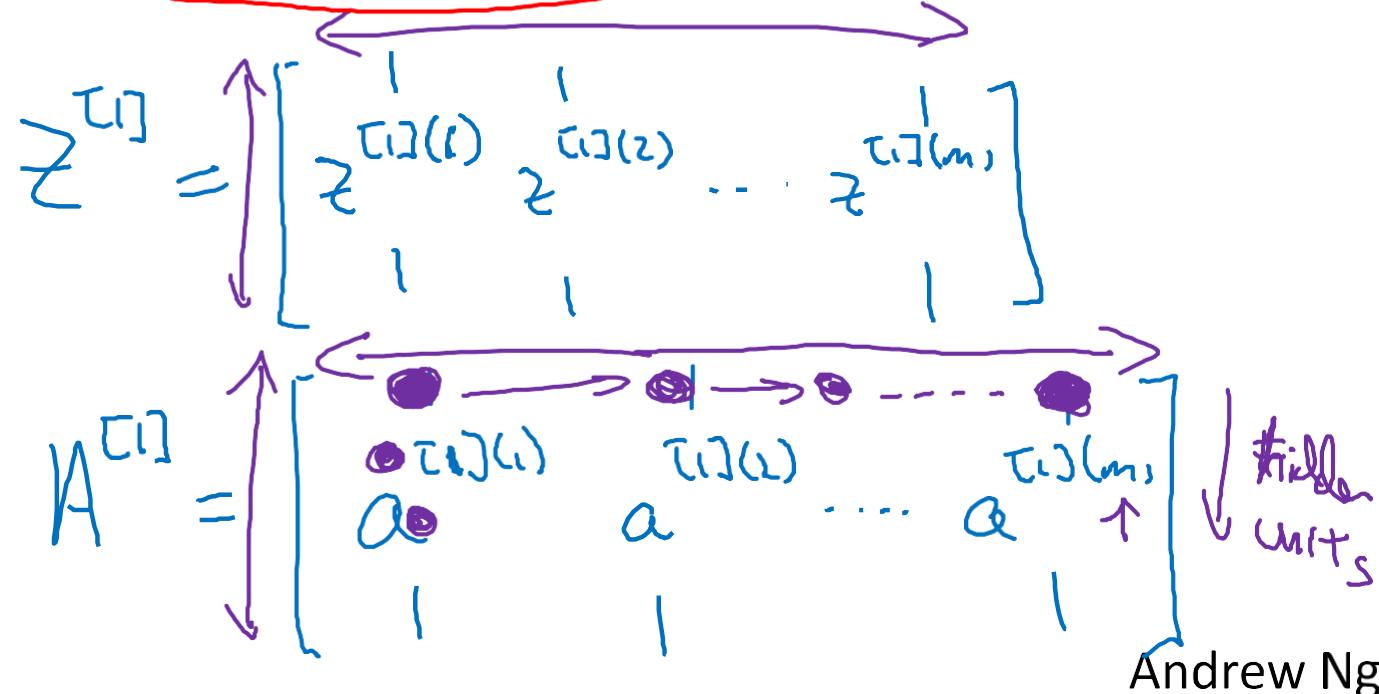


$$z^{[1]} = W^{[1]}X + b^{[1]}$$

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

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

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





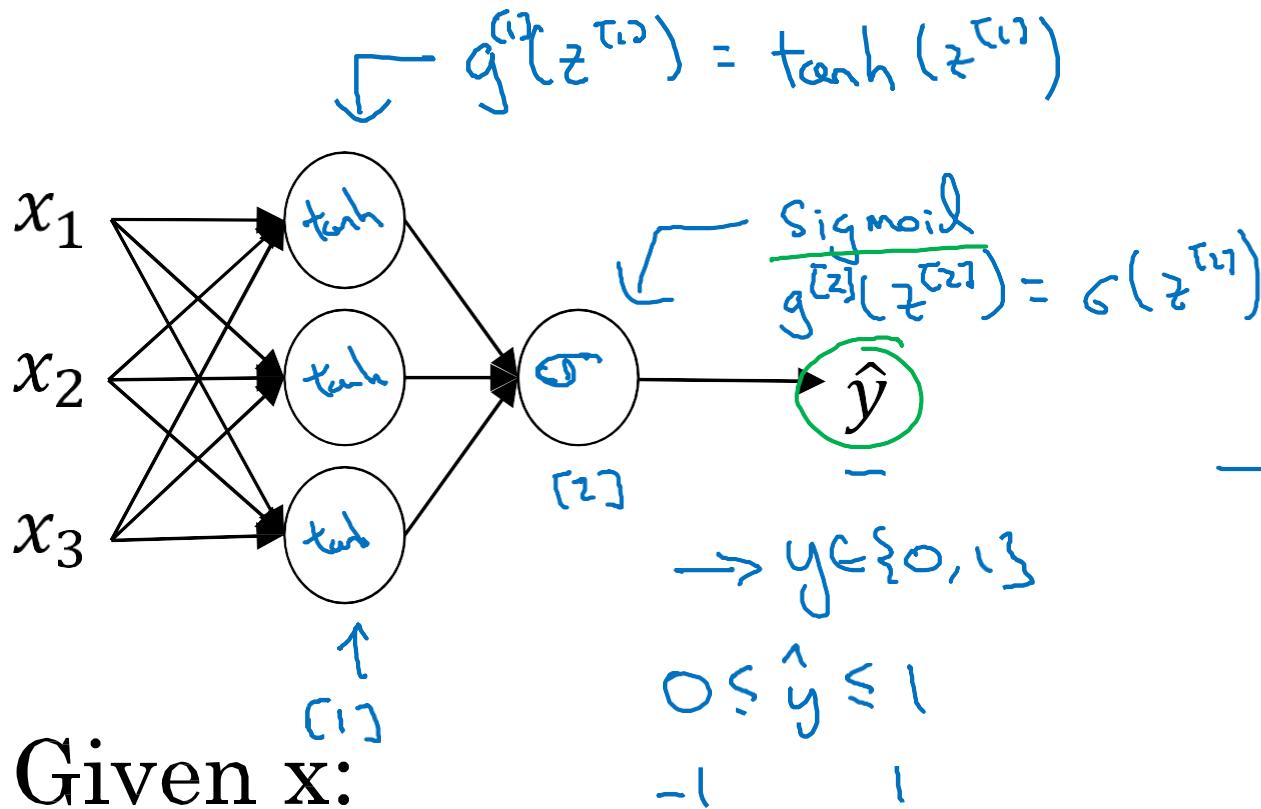
deeplearning.ai

One hidden layer  
Neural Network

---

Activation functions

# Activation functions



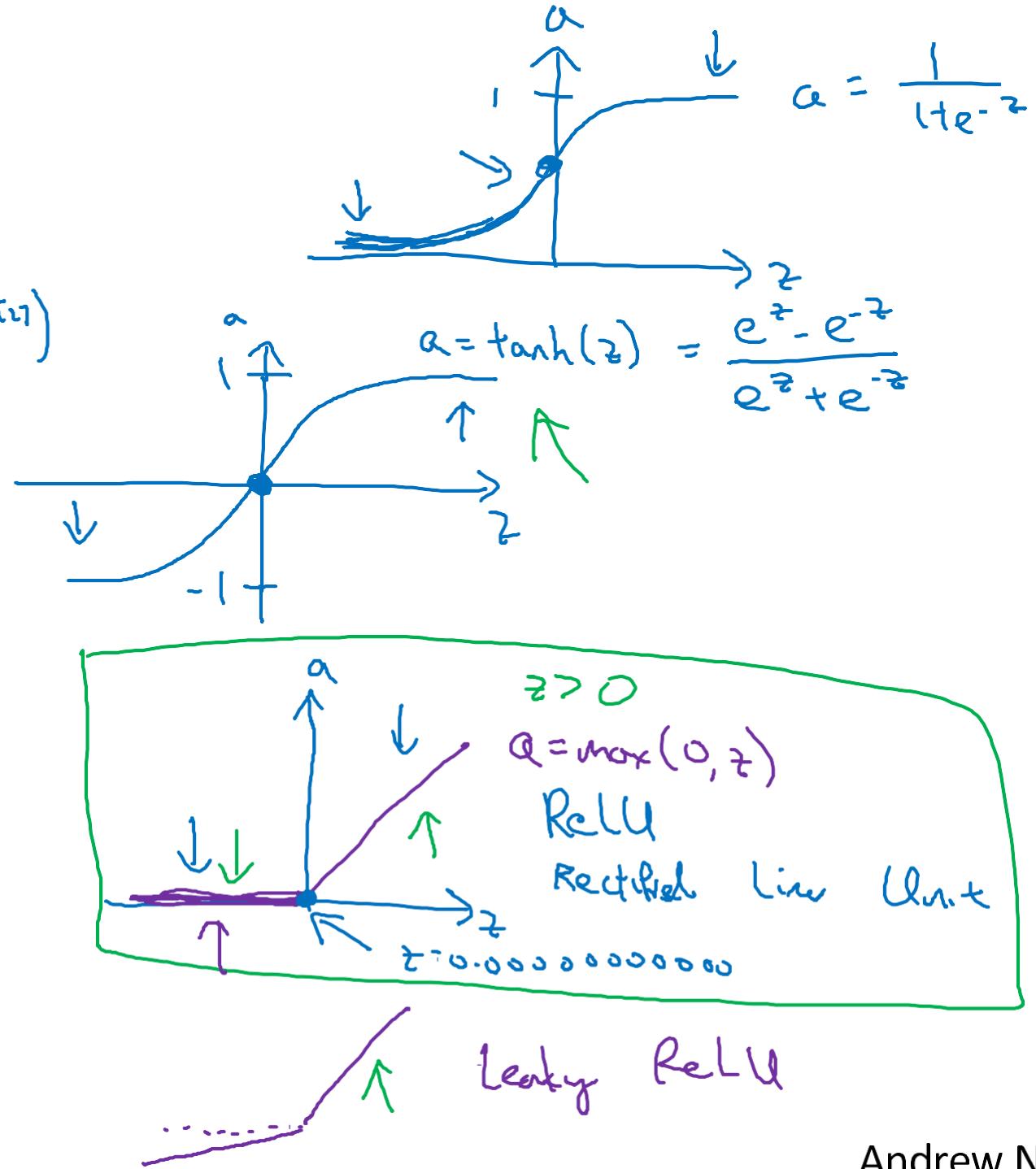
Given x:

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

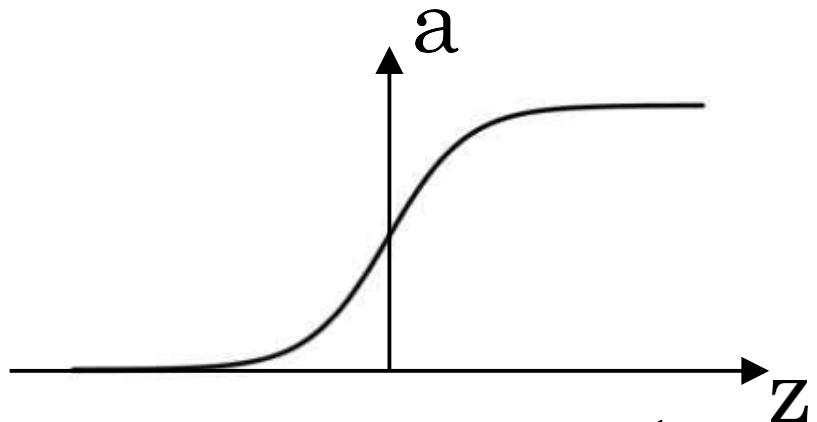
$$\rightarrow a^{[1]} = \sigma(z^{[1]}) \quad g^{[1]}(z^{[1]})$$

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

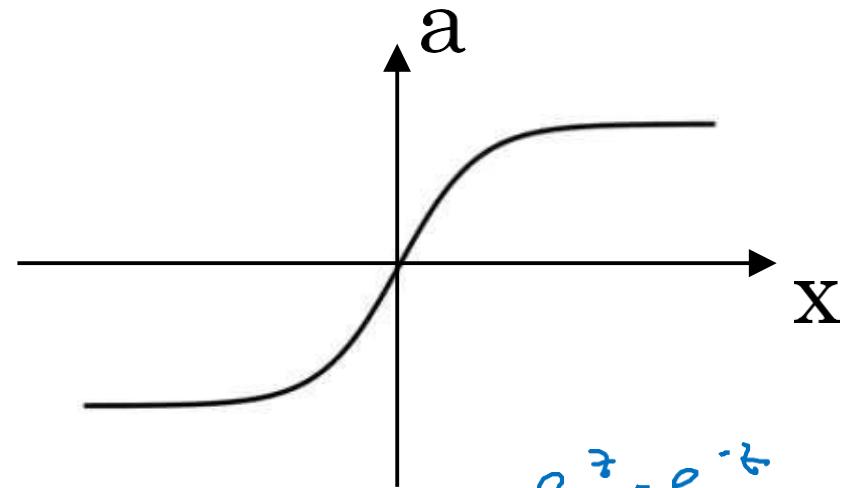
$$\rightarrow a^{[2]} = \sigma(\cancel{z^{[2]}}) \quad \cancel{g^{[2]}(z^{[2]})}$$



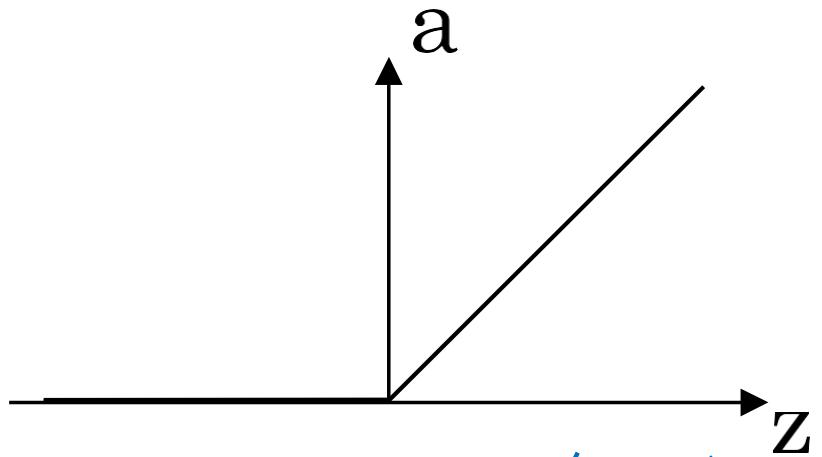
# Pros and cons of activation functions



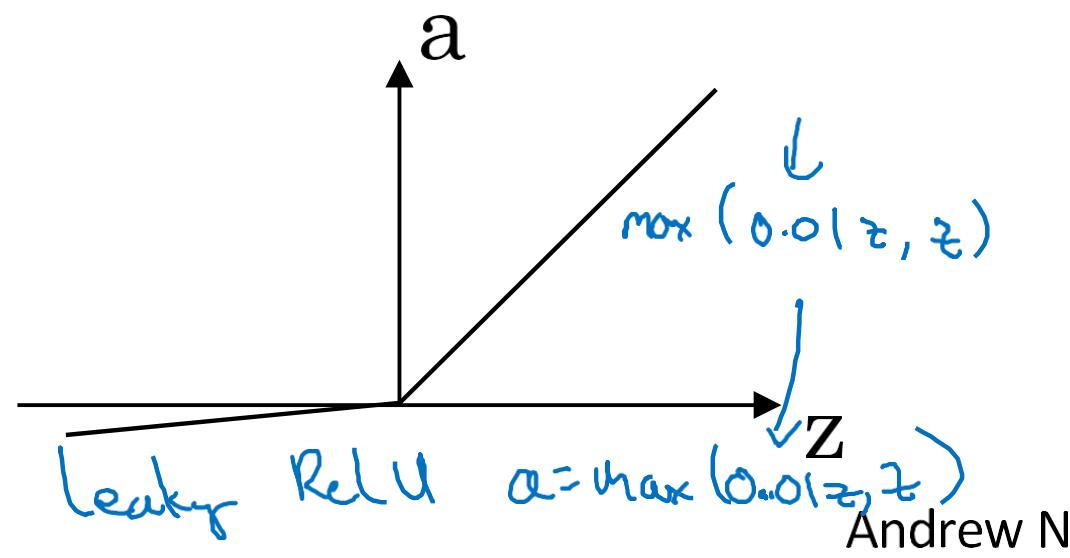
$$\text{sigmoid: } a = \frac{1}{1 + e^{-z}}$$



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



$$\text{ReLU} \quad a = \max(0, z)$$



$$\text{Leaky ReLU} \quad a = \max(0.01z, z)$$

Andrew Ng



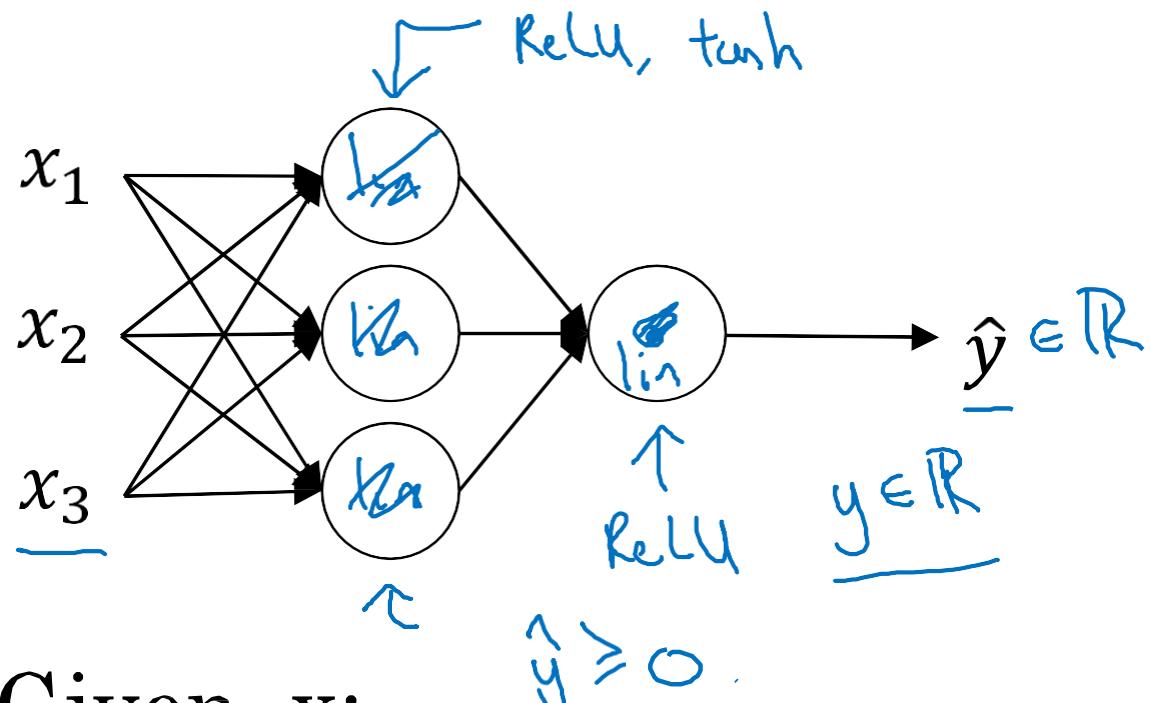
deeplearning.ai

# One hidden layer Neural Network

---

Why do you  
need non-linear  
activation functions?

# Activation function



Given  $x$ :

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

$$\rightarrow a^{[1]} = g^{[1]}(z^{[1]}) \geq^{C_1}$$

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

$$\rightarrow a^{[2]} = g^{[2]}(z^{[2]}) \geq^{C_2}$$

$g(z) = z$   
"linear activation  
function"

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

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

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

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

$$= \underline{W'x + b'}$$

$$g(x) = z$$



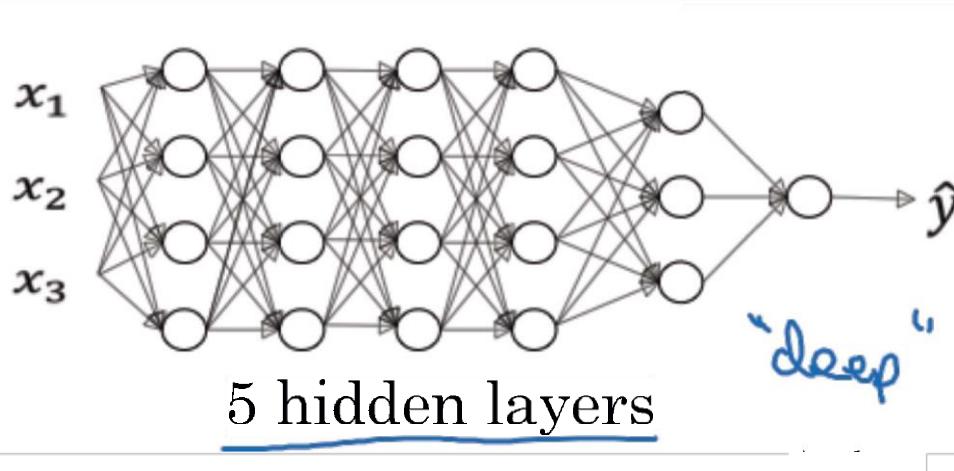
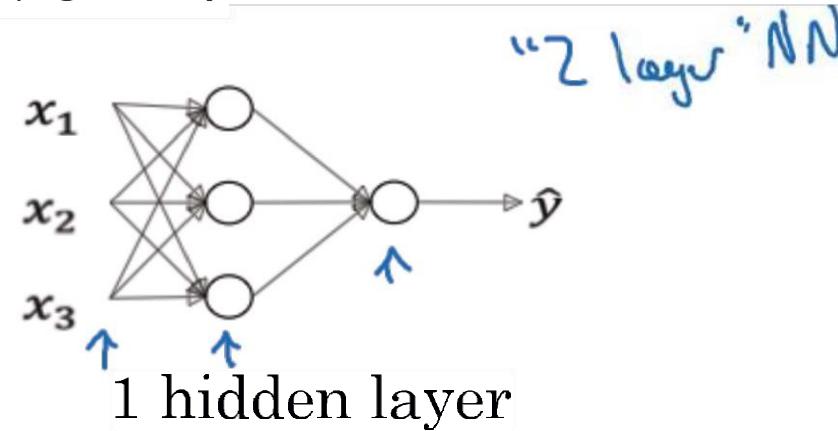
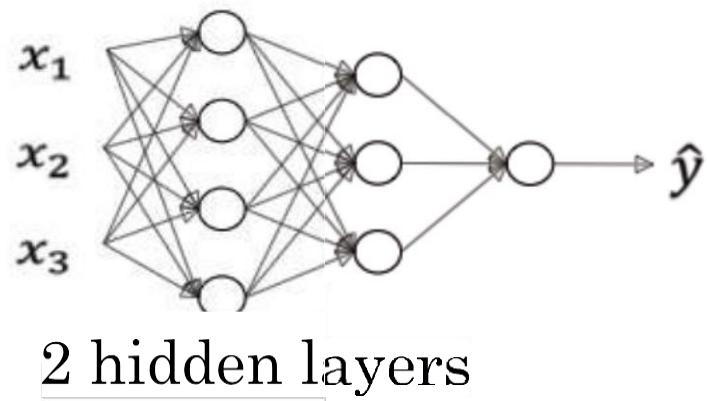
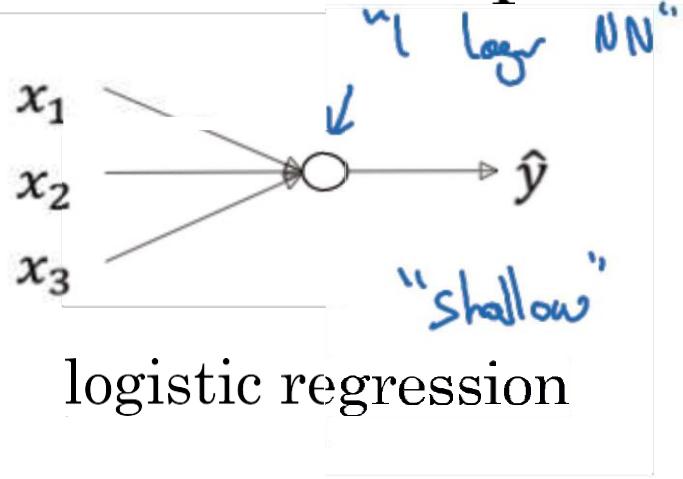
deeplearning.ai

# Deep Neural Networks

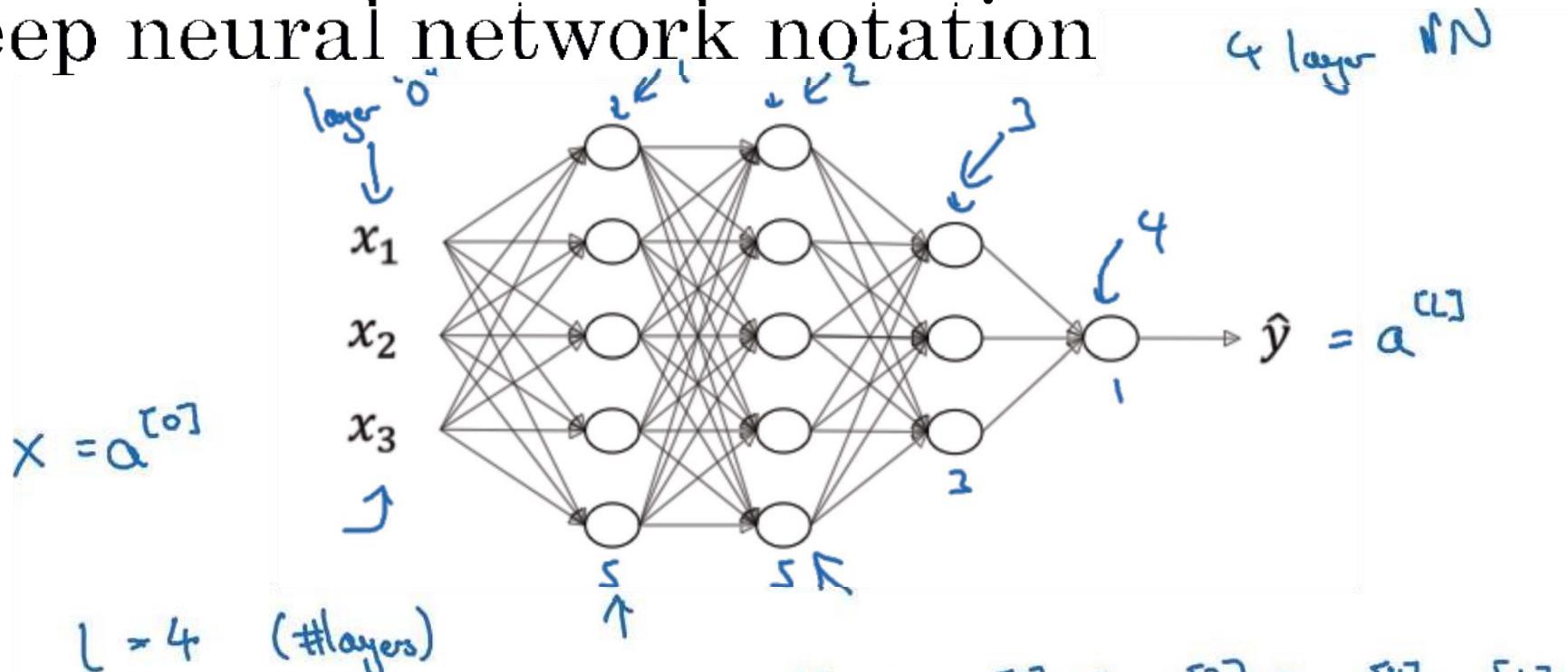
---

Deep L-layer  
Neural network

# What is a deep neural network?



# Deep neural network notation



$$l = 4 \quad (\# \text{layers})$$

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

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

$a^{[l]} = g^{[l]}(z^{[l]})$ ,  $w^{[l]}$  = weights for  $z^{[l]}$

$$n^{[1]} = 5, n^{[2]} = 5, n^{[3]} = 3, n^{[4]} = n^{[L]} = 1$$

$$n^{[0]} = n_x = 3$$



deeplearning.ai

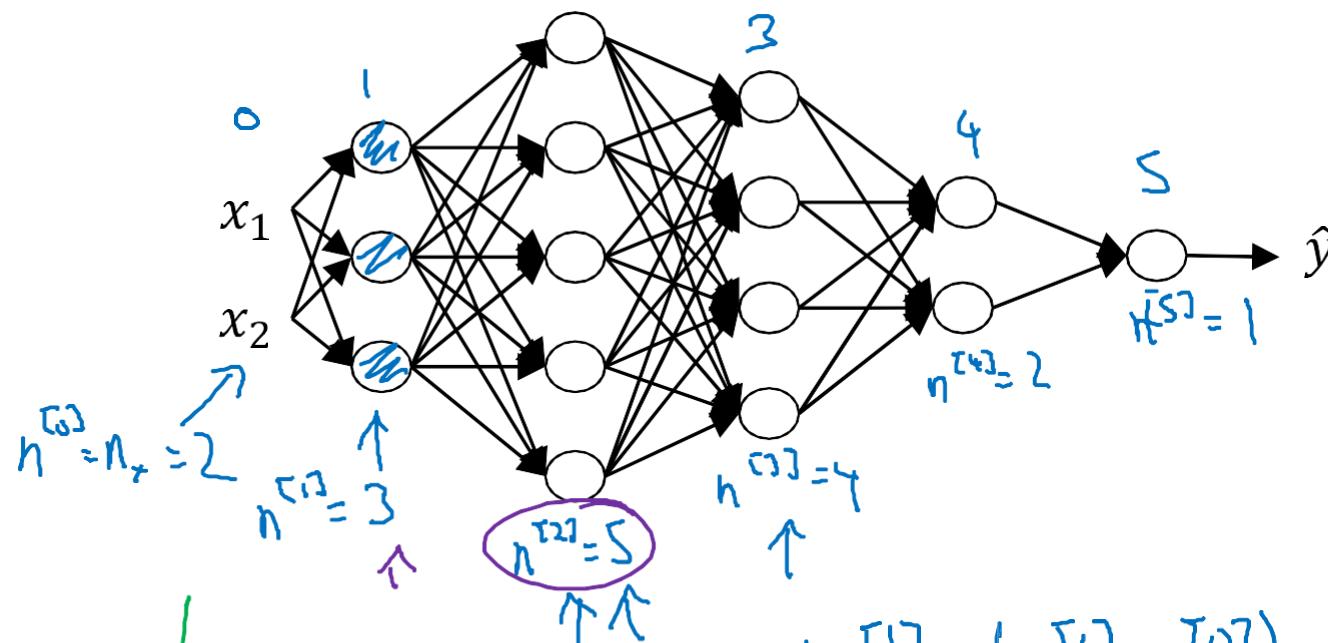
# Deep Neural Networks

---

## Getting your matrix dimensions right

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

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



$$z^{[1]} = \underbrace{w^{[1]} \cdot x}_{(3,1) \leftarrow (3,2) \begin{pmatrix} 2 \\ 1 \end{pmatrix}} + \underbrace{b^{[1]}_{(3,1)}}_{(n^{[1]}, 1)} = \underbrace{(n^{[1]}, 1)}_{\text{---}}$$

$$\begin{bmatrix} \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$$w^{[1]}: (n^{[1]}, n^{[0]})$$

$$w^{[2]}: (5, 3) \quad (n^{[2]}, n^{[1]})$$

$$z^{[2]} = \underbrace{w^{[2]} \cdot a^{[1]}}_{\rightarrow (5,1)} + \underbrace{b^{[2]}_{(5,1)}}_{(n^{[2]}, 1)}$$

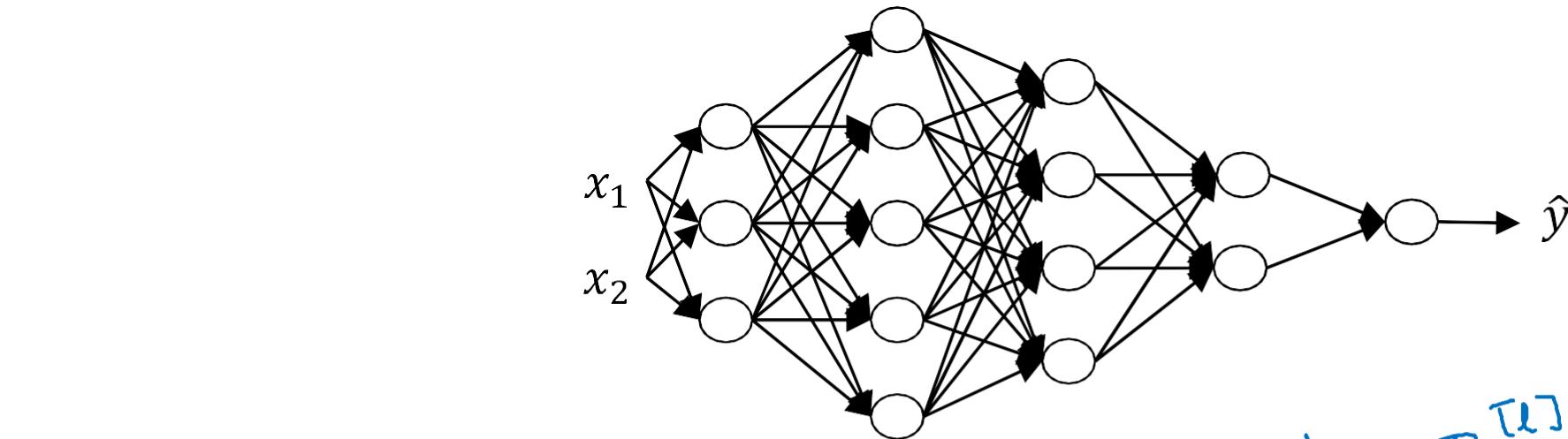
$$w^{[3]}: (4, 5)$$

$$w^{[4]}: (2, 4), \quad w^{[5]}: (1, 2)$$

$L=5$

$$\begin{cases} w^{[l]}: (n^{[l]}, n^{[l-1]}) \\ b^{[l]}: (n^{[l]}, 1) \\ dw^{[l]}: (n^{[l]}, n^{[l-1]}) \\ db^{[l]}: (n^{[l]}, 1) \end{cases}$$

# Vectorized implementation



$$z^{[l]} = w^{[l]} \cdot x + b^{[l]}$$

$(n^{[l]}, 1)$      $(n^{[l]}, n^{[l+1]})$      $\underbrace{(n^{[l]}, n^{[l]})}_{(n^{[l]}, 1)}$

$$\begin{bmatrix} z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \end{bmatrix}$$

$$\overbrace{\begin{bmatrix} z^{[l]} \\ \vdots \\ z^{[l]} \end{bmatrix}}^{\uparrow} = \underbrace{w^{[l]} \cdot X + b^{[l]}}_{\uparrow} \quad \underbrace{(n^{[l]}, m)}_{\uparrow} \quad \underbrace{(n^{[l]}, n^{[l+1]})}_{\uparrow} \quad \underbrace{(n^{[l]}, m)}_{\uparrow} \quad \underbrace{(n^{[l]}, 1)}_{\uparrow}$$

$$z^{[l]}, a^{[l]} : (n^{[l]}, 1)$$

$$z^{[l]}, A^{[l]} : (n^{[l]}, m)$$

$$l=0 \quad A^{[0]} = X = (n^{[0]}, m)$$

$$dz^{[l]}, dA^{[l]} : (n^{[l]}, m)$$



deeplearning.ai

# Deep Neural Networks

---

## Parameters vs Hyperparameters

# What are hyperparameters?

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \dots$

Hyperparameters:

learning rate  $\frac{\alpha}{n}$   
#iterations

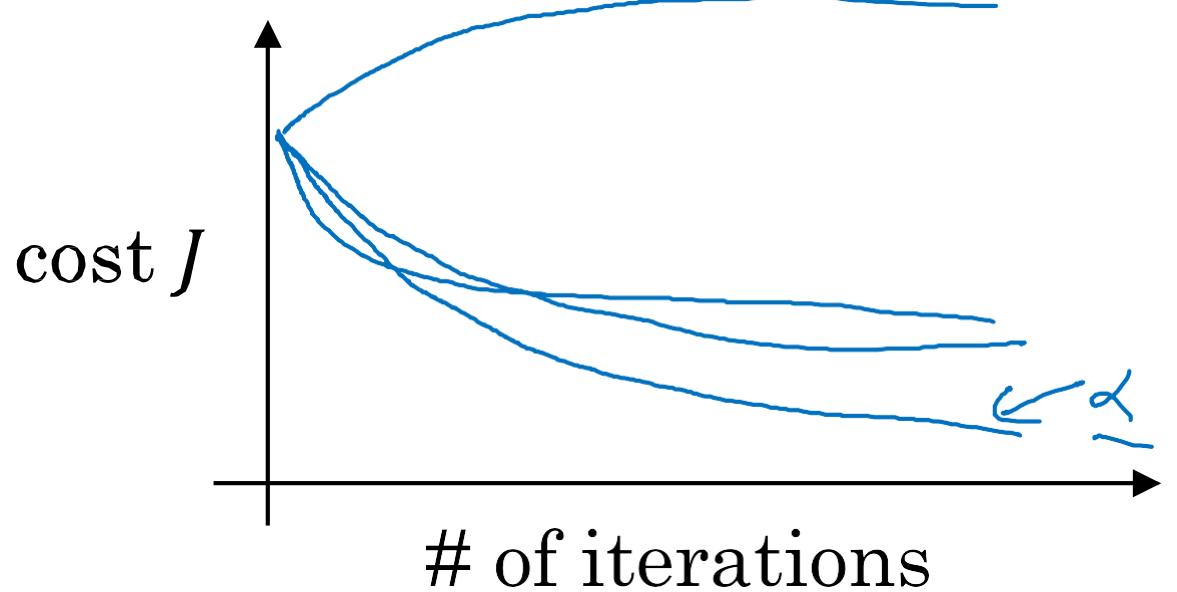
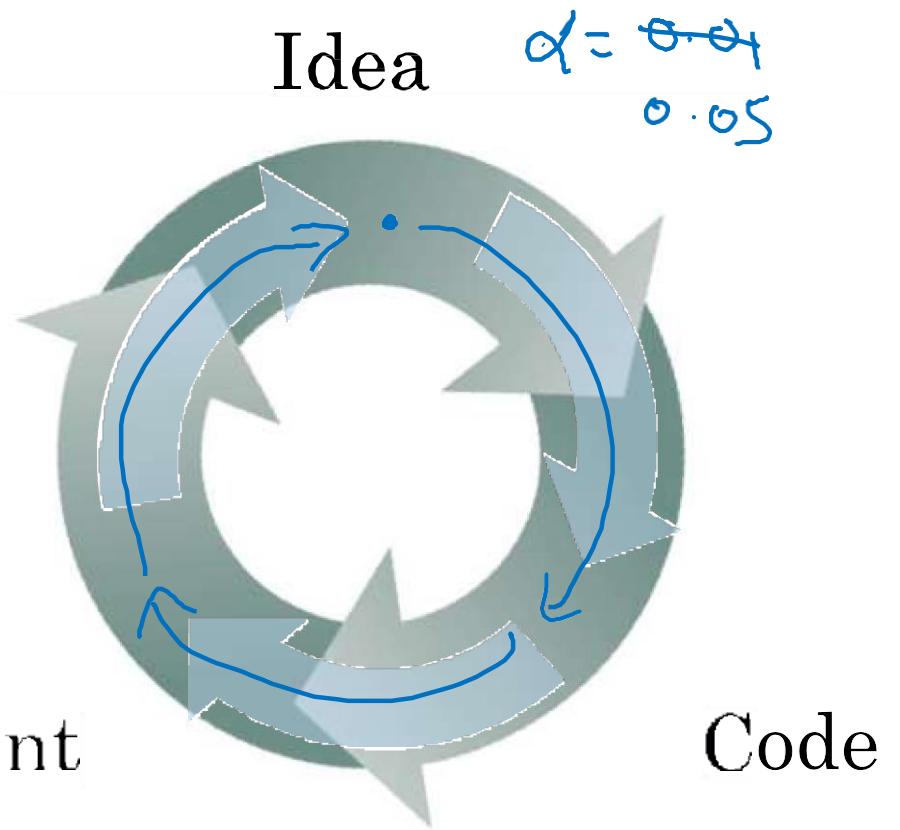
#hidden layers  $L$

#hidden units  $n^{[1]}, n^{[2]}, \dots$

choice of activation function

Later: Momentum, minibatch size, regularizations, ...

# Applied deep learning is a very empirical process





deeplearning.ai

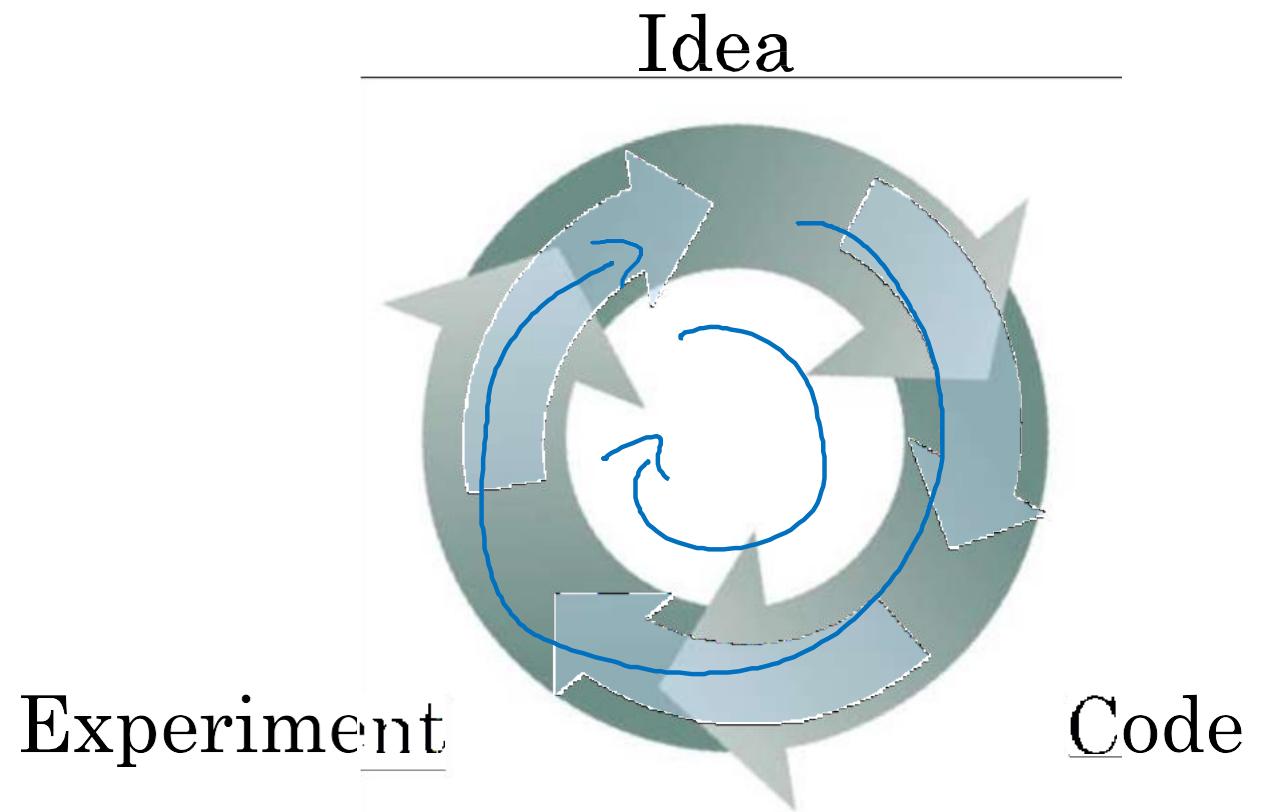
Setting up your  
ML application

---

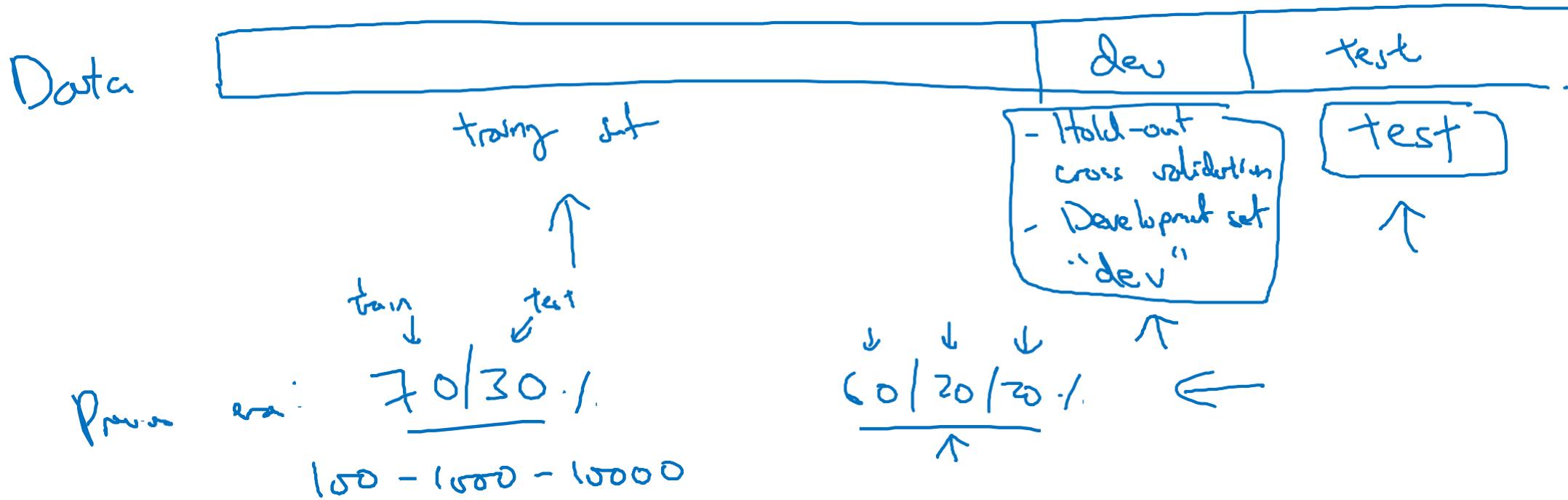
Train/dev/test  
sets

# Applied ML is a highly iterative process

{ # layers  
# hidden units  
learning rates  
activation functions  
...



# Train/dev/test sets



Big data:

1,000,000

10,000

10,000

98 / 1 / 1 . . .

99.5 { 25 { 25  
    .4   .1 . . .

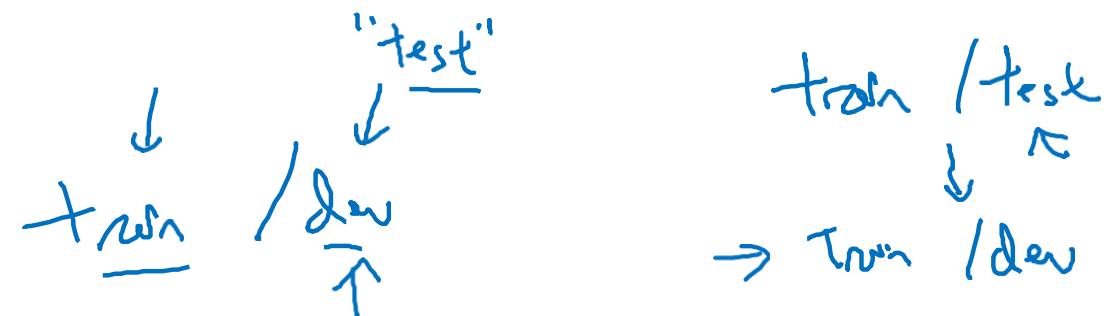
# Mismatched train/test distribution

Conts

Training set:  
Cat pictures from }  
webpages

Dev/test sets:  
Cat pictures from }  
users using your app }

→ Make sure dev and test come from same distribution.



Not having a test set might be okay. (Only dev set.)



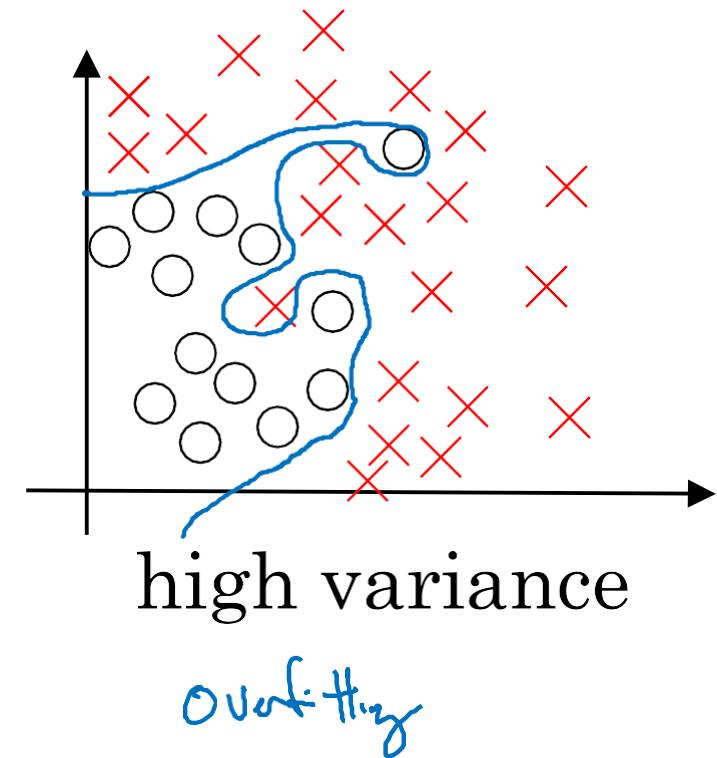
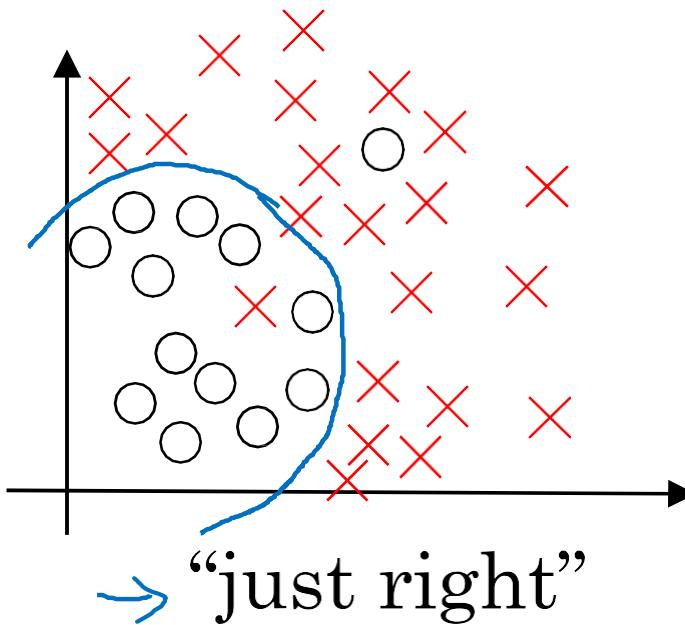
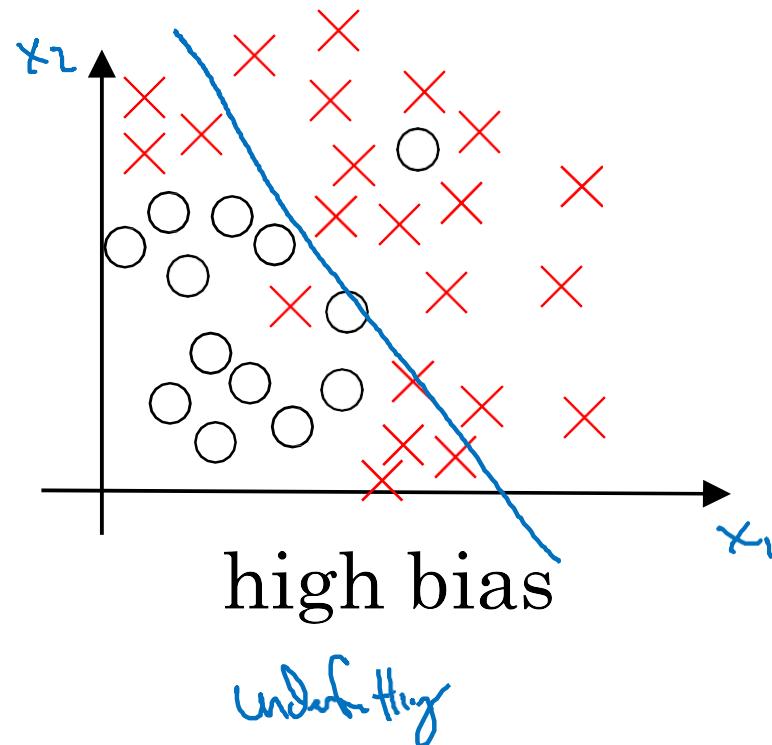
deeplearning.ai

Setting up your  
ML application

---

Bias/Variance

# Bias and Variance



# Bias and Variance

Cat classification



Train set error:

1%

Dev set error:

11%

Human  
 $\approx 0\%$

high variance



15% ↙

16% ↙

high bias



15%  
30%  
high bias  
& high varan

0.5%

1%

low bias  
low variance  
↑

Optimal (Bayes) error:  ~~$\approx 0\%$~~  15% Blury images



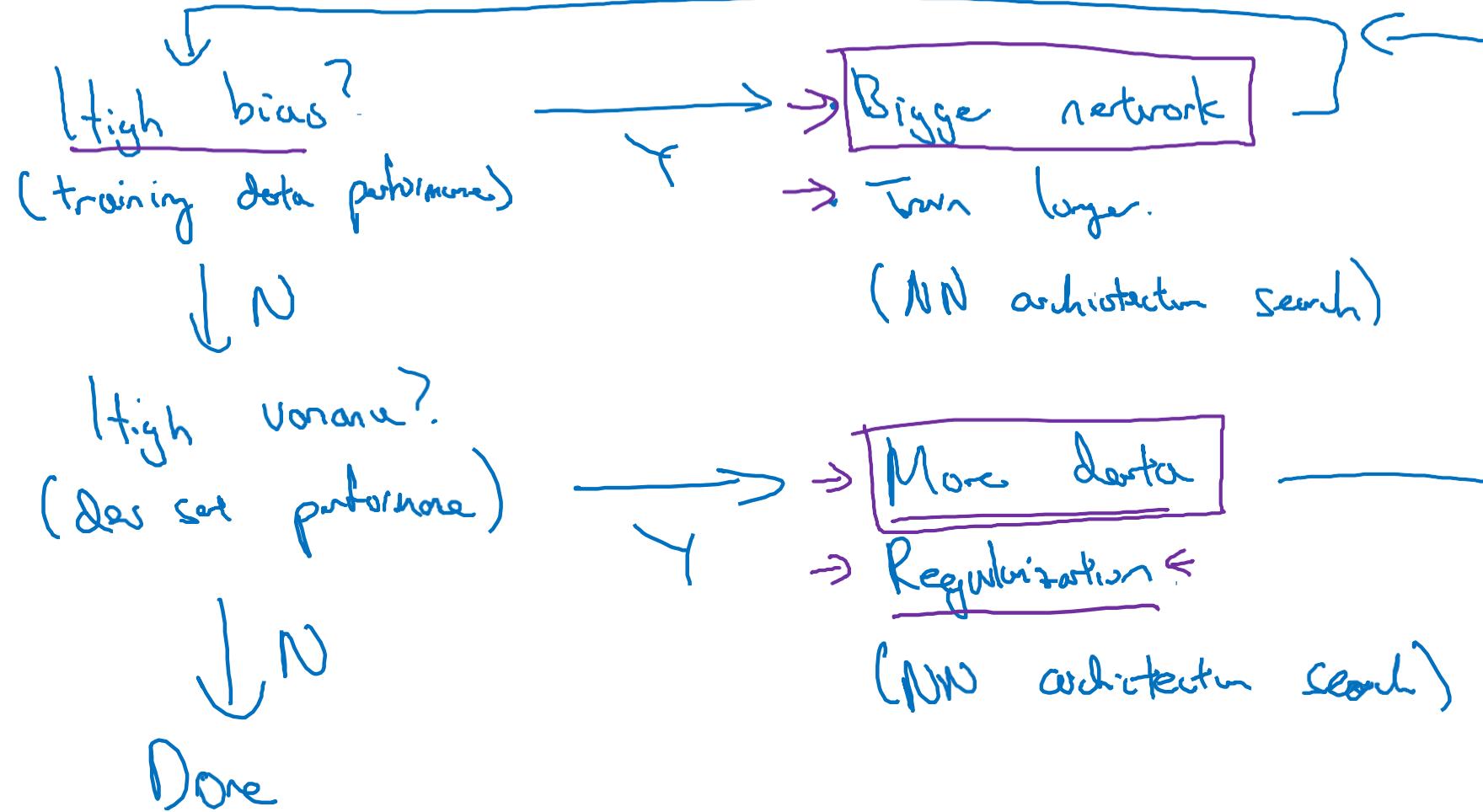
deeplearning.ai

Setting up your  
ML application

---

Basic “recipe”  
for machine learning

# Basic recipe for machine learning





deeplearning.ai

# Regularizing your neural network

---

## Regularization

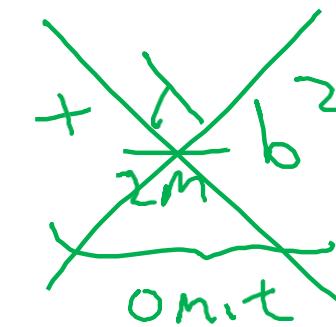
# Logistic regression

$$\min_{w,b} J(w,b)$$

$$w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

$\lambda$  = regularization parameter  
lambda

$$J(w,b) = \underbrace{\frac{1}{m} \sum_{i=1}^m l(y^{(i)}, y_{(i)})}_{\text{L}_2 \text{ regularization}} + \frac{\lambda}{2m} \|w\|_2^2$$



$$\|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$$

# Neural network

$$\rightarrow J(\omega^{(0)}, b^{(0)}, \dots, \omega^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(y^{(i)}, \hat{y}^{(i)})}_{\text{loss function}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L \|\omega^{(l)}\|_F^2}_{\text{regularization term}}$$

$$\|\omega^{(l)}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{(l-1)}} (\omega_{ij}^{(l)})^2$$

$\omega^{(l)} : (n^{(l)}, n^{(l-1)})$

"Frobenius norm"

$$\|\cdot\|_2^2$$

$$\|\cdot\|_F^2$$

$$d\omega^{(l)} = \boxed{(\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)}}$$

$$\rightarrow \omega^{(l)} := \omega^{(l)} - \alpha d\omega^{(l)}$$

$$\frac{\partial J}{\partial \omega^{(l)}} = d\omega^{(l)}$$

"Weight decay"

$$\omega^{(l)} := \omega^{(l)} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} \omega^{(l)} \right]$$

$$= \omega^{(l)} - \frac{\alpha \lambda}{m} \omega^{(l)} - \alpha (\text{from backprop})$$

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right) \omega^{(l)}}_{\leq 1} - \alpha \underbrace{(\text{from backprop})}$$



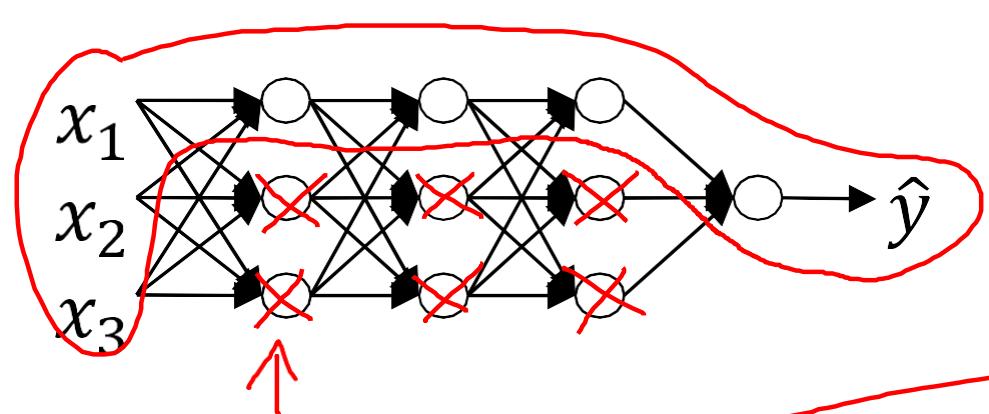
deeplearning.ai

# Regularizing your neural network

---

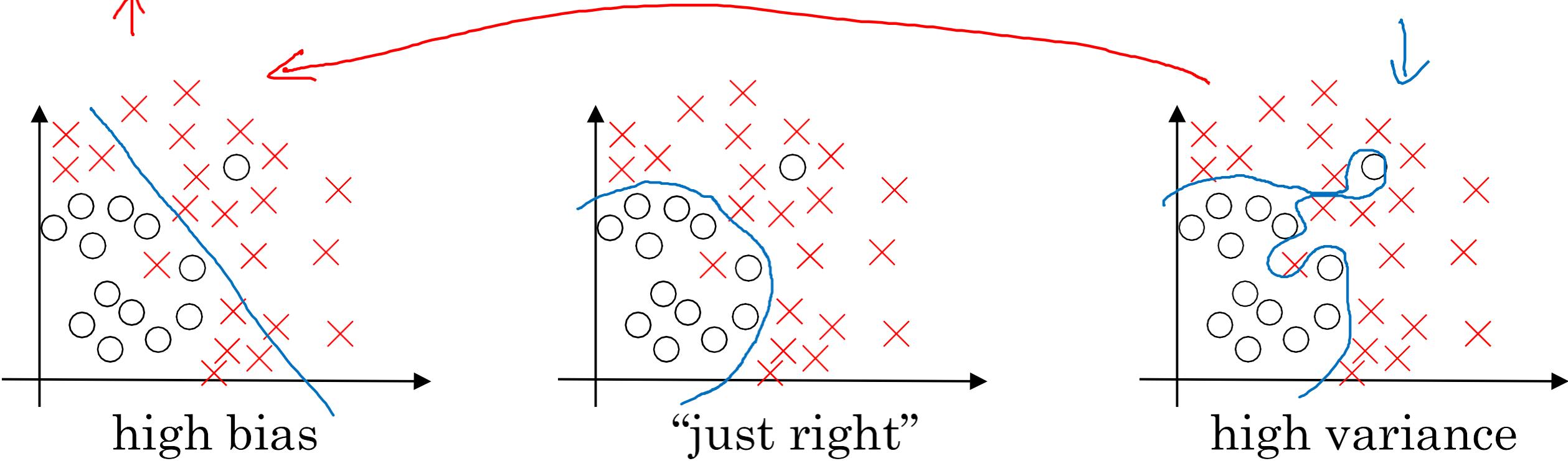
## Why regularization reduces overfitting

# How does regularization prevent overfitting?

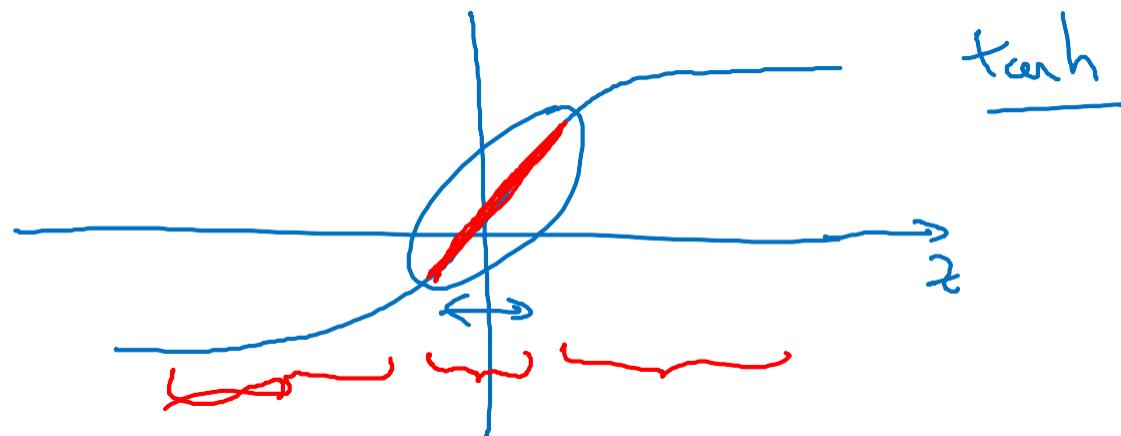


$$J(\boldsymbol{w}^{(m)}, \boldsymbol{b}^{(m)}) = \frac{1}{m} \sum_{i=1}^m L(y_i^{(i)}, \hat{y}_i^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

$\underline{\boldsymbol{w}^{(l)} \approx 0}$



# How does regularization prevent overfitting?

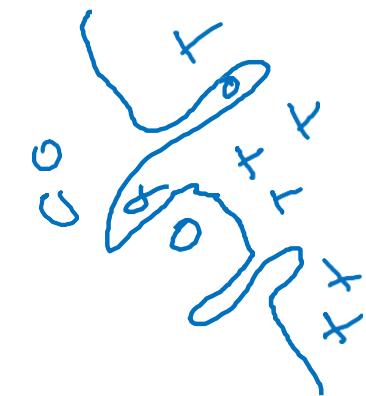


$$g(z) := \tanh(z)$$

$$\lambda \uparrow$$

$$w^{[l]} \downarrow$$

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



Every layer  $\approx$  linear

$$J(\dots) = \left[ \sum_i L(\hat{y}^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2m} \sum_l \|w^{[l]}\|_F^2$$





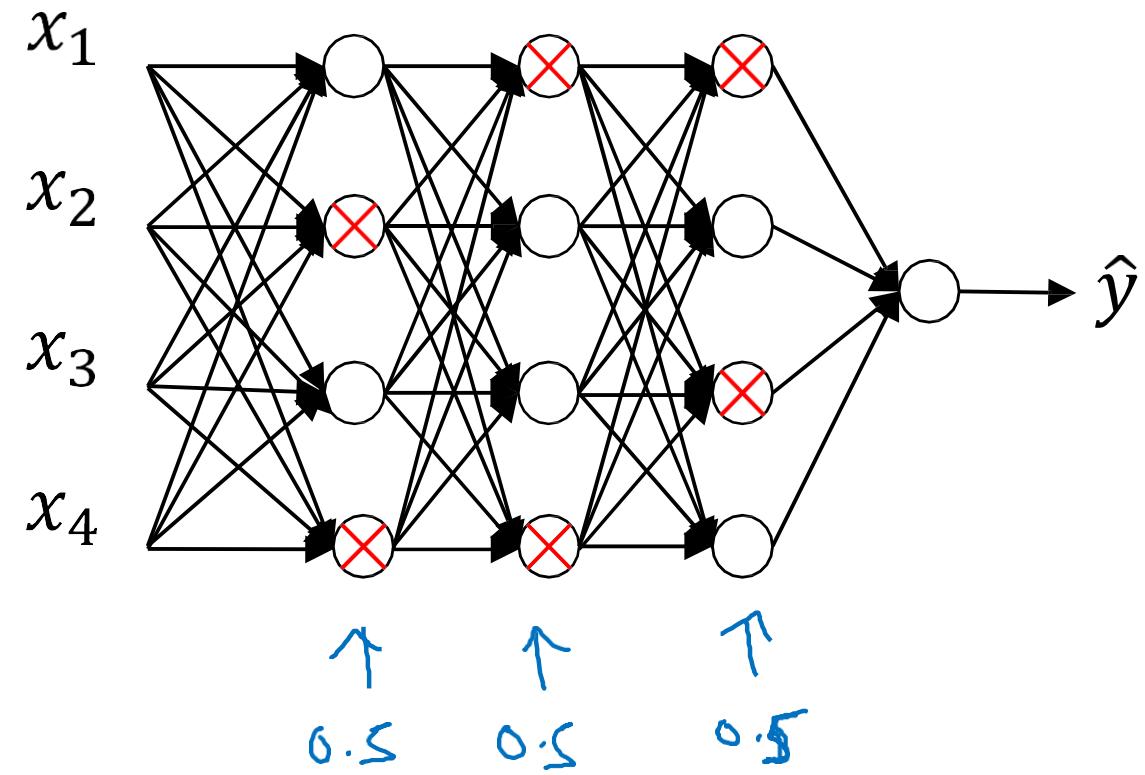
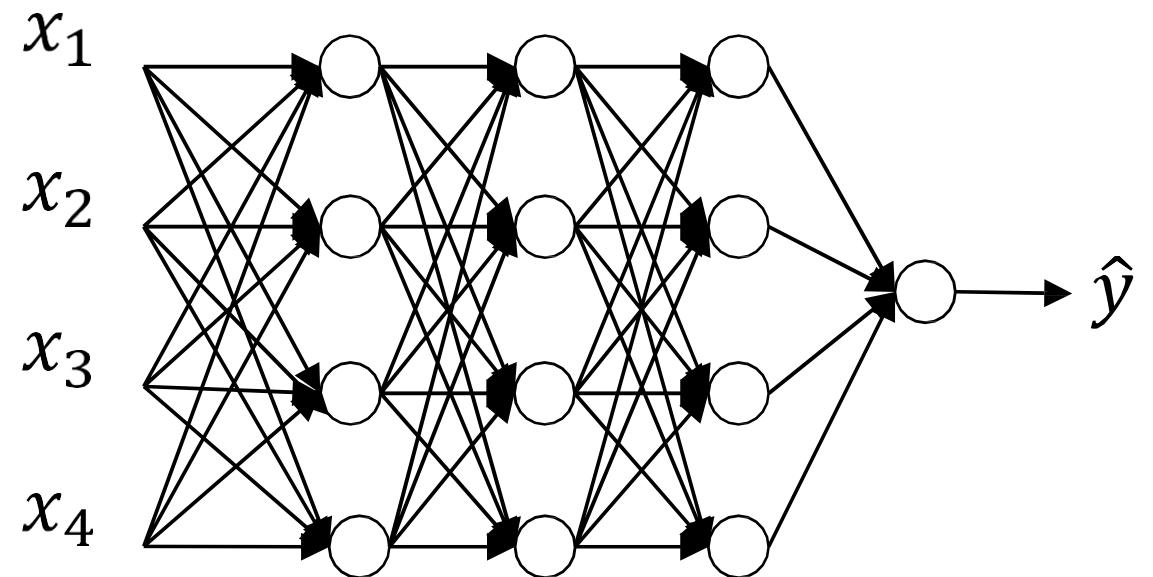
deeplearning.ai

# Regularizing your neural network

---

## Dropout regularization

# Dropout regularization



# Implementing dropout (“Inverted dropout”)

Illustrate with layer  $\ell=3$ .  $\text{keep-prob} = \frac{0.8}{5}$  0.2

$\rightarrow d_3 = \underbrace{\text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep-prob}}$

$\underbrace{a_3}_{\text{a}_3} = \text{np.multiply}(a_3, d_3)$  #  $a_3 * d_3$ .

$\rightarrow a_3 /= \cancel{0.8} \text{ keep-prob}$   $\leftarrow$

50 units  $\rightsquigarrow$  10 units shut off

$$z^{[4]} = w^{[4]} \cdot \underbrace{a^{[3]}}_{\text{50 units}} + b^{[4]}$$

$\mathcal{J}$  reduced by 20%.

Test

$$1 = \underline{0.8}$$

# Making predictions at test time

$$a^{(0)} = X$$

No drop out

$$z^{(1)} = w^{(1)} \underline{a^{(0)}} + b^{(1)}$$

$$a^{(1)} = g^{(1)}(\underline{z^{(1)}})$$

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

$$a^{(2)} = \dots$$

$$\downarrow \hat{y}$$

$\lambda$  = keep-prob



deeplearning.ai

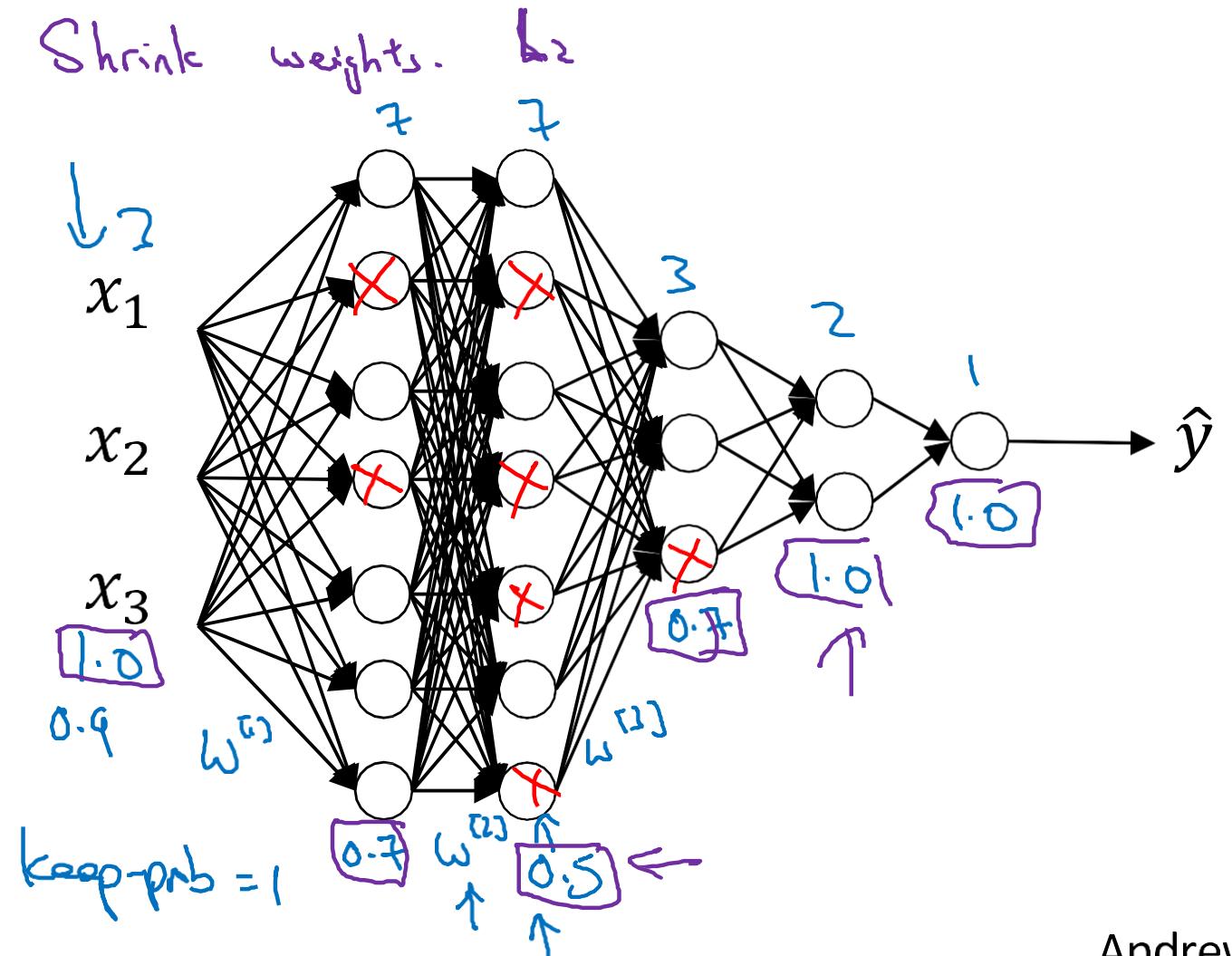
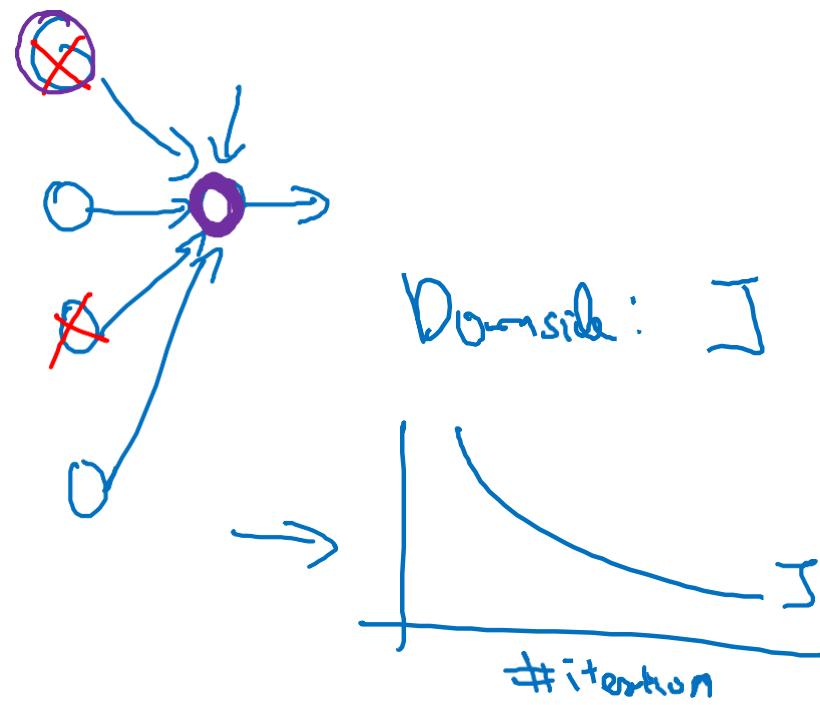
# Regularizing your neural network

---

## Understanding dropout

# Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights.  $\rightsquigarrow$  Shrink weights.





deeplearning.ai

# Regularizing your neural network

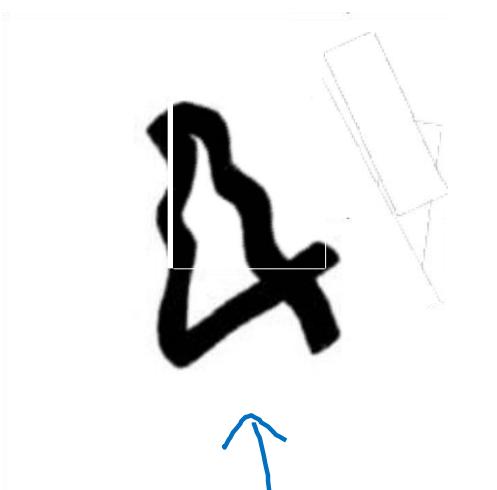
---

## Other regularization methods

# Data augmentation



4

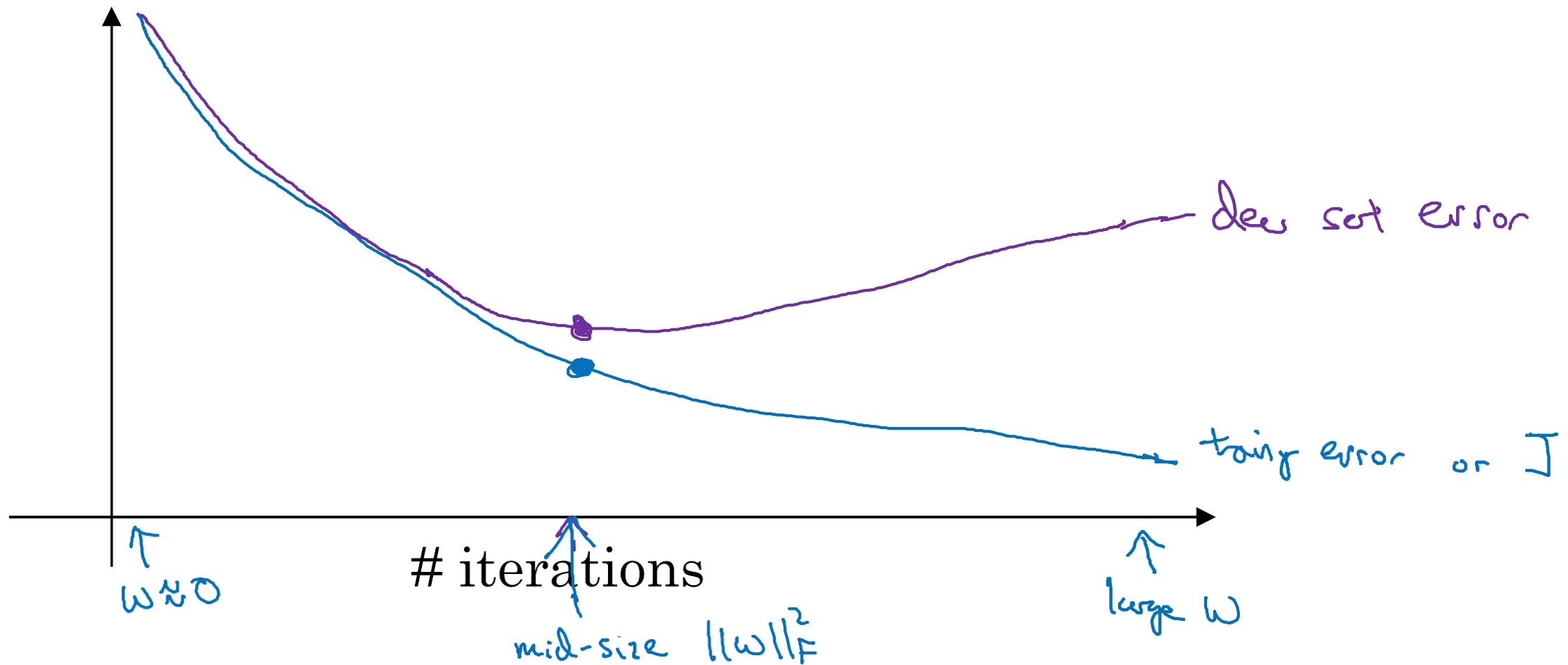


4



Andrew Ng

# Early stopping





deeplearning.ai

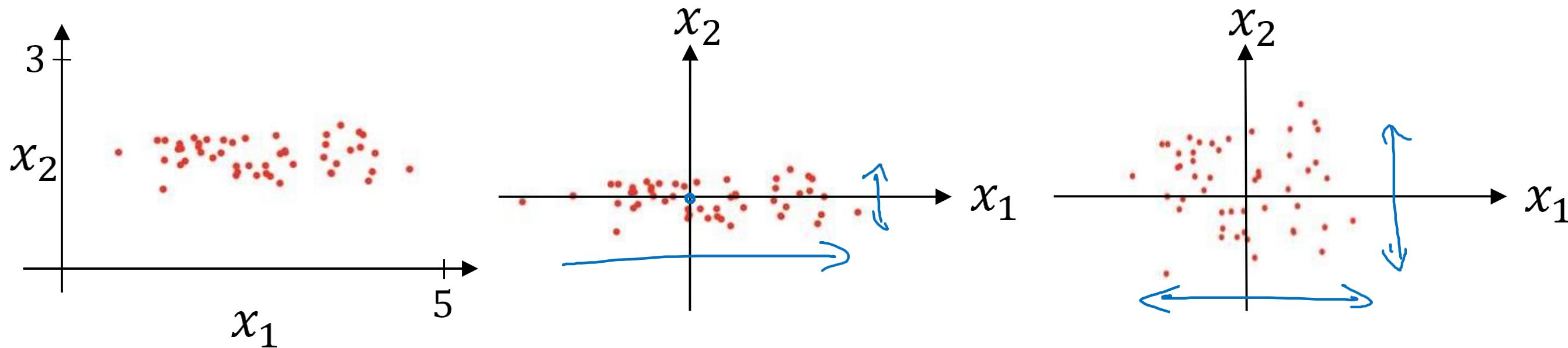
# Setting up your optimization problem

---

## Normalizing inputs

# Normalizing training sets

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$



Subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$x := x - \mu$$

Normalize variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m x^{(i)} * x^{(i)}$$

~ element-wise

$$x / \sigma$$

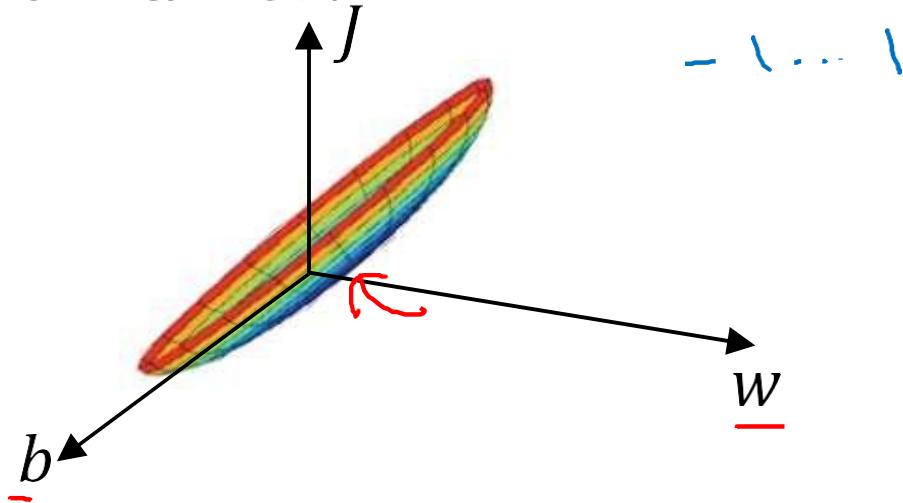
Use same  $\mu, \sigma$  to normalize test set.

# Why normalize inputs?

$w_1$   $x_1$ : 1...1000  $\leftarrow$

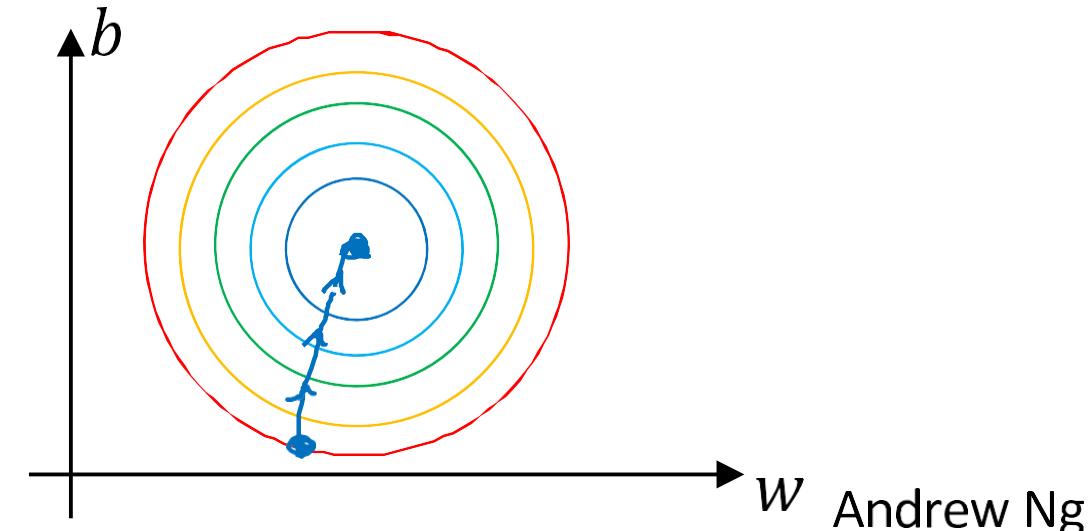
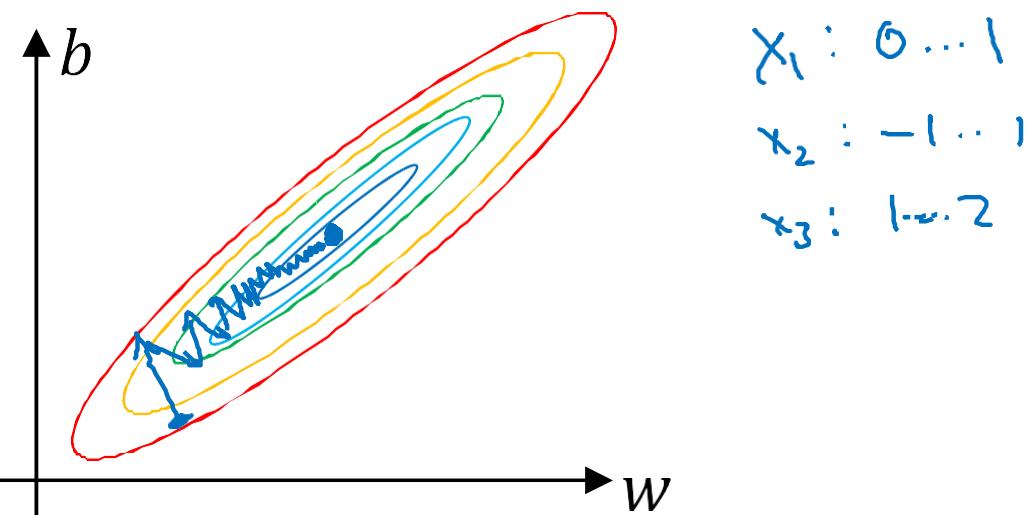
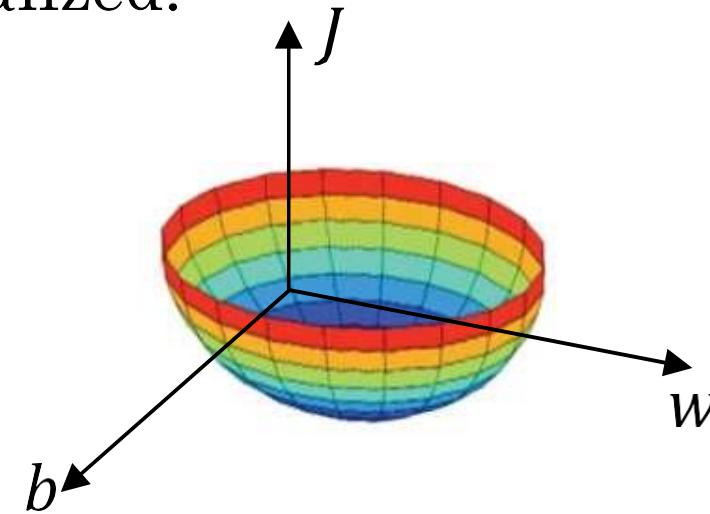
$w_2$   $x_2$ : 0...1  $\leftarrow$

Unnormalized:



$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Normalized:





deeplearning.ai

# Optimization Algorithms

---

## Mini-batch gradient descent

# Batch vs. mini-batch gradient descent

$X, Y$

$X^{\{t\}}, Y^{\{t\}}$

Vectorization allows you to efficiently compute on  $m$  examples.

$$X = [x^{(1)} \ x^{(2)} \ x^{(3)} \ \dots \ x^{(1000)} \ | \ x^{(1001)} \ \dots \ x^{(2000)} \ | \ \dots \ | \ \dots \ x^{(m)}]$$

$\underbrace{x^{(1)}}_{(n_x, m)} \quad (n_x, 1000) \quad \underbrace{x^{(2)}}_{(n_x, 1000)} \quad \dots \quad \underbrace{x^{(5,000)}}_{(n_x, 1000)}$

$$Y = [y^{(1)} \ y^{(2)} \ y^{(3)} \ \dots \ y^{(1000)} \ | \ y^{(1001)} \ \dots \ y^{(2000)} \ | \ \dots \ | \ \dots \ y^{(m)}]$$

$\underbrace{y^{(1)}}_{(1, m)} \quad (1, 1000) \quad \underbrace{y^{(2)}}_{(1, 1000)} \quad \dots \quad \underbrace{y^{(5,000)}}_{(1, 1000)}$

What if  $m = 5,000,000$ ?

5,000 mini-batches of 1,000 each

Mini-batch  $t$ :  $X^{\{t\}}, Y^{\{t\}}$

$x^{(i)}$   
 $z^{[l]}$   
 $X^{\{t\}}, Y^{\{t\}}$

# Mini-batch gradient descent

repeat {  
for  $t = 1, \dots, 5000$  {

Forward prop on  $X^{\{t\}}$

$$z^{(l)} = w^{(l)} X^{\{t\}} + b^{(l)}$$

$$A^{(l)} = g^{(l)}(z^{(l)})$$

:

$$A^{(L)} = g^{(L)}(z^{(L)})$$

Compute cost  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|w^{(l)}\|_F^2$ .

Backprop to compute gradients wrt  $J^{\{t\}}$  (using  $(X^{\{t\}}, Y^{\{t\}})$ )

$$w^{(l)} := w^{(l)} - \alpha \nabla J^{\{t\}} \quad b^{(l)} := b^{(l)} - \alpha \nabla b^{(l)}$$

} }

"1 epoch"  
pass through training set.

1 step of gradient descent  
using  $\frac{X^{\{t\}}, Y^{\{t\}}}{1000}$   
(as if 1000)

$X, Y$



deeplearning.ai

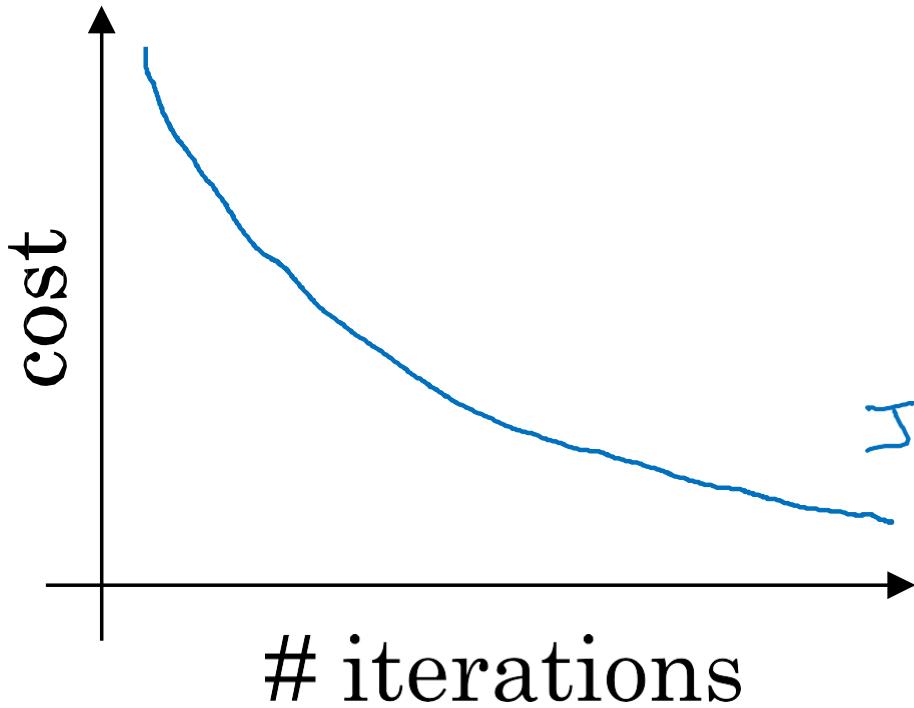
# Optimization Algorithms

---

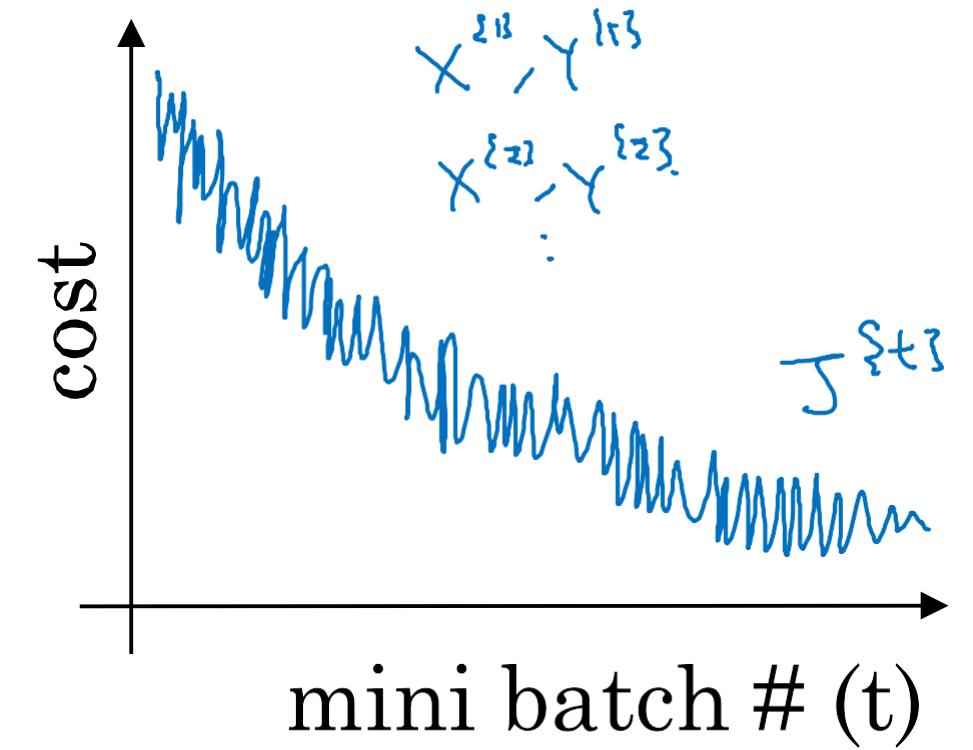
## Understanding mini-batch gradient descent

# Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent

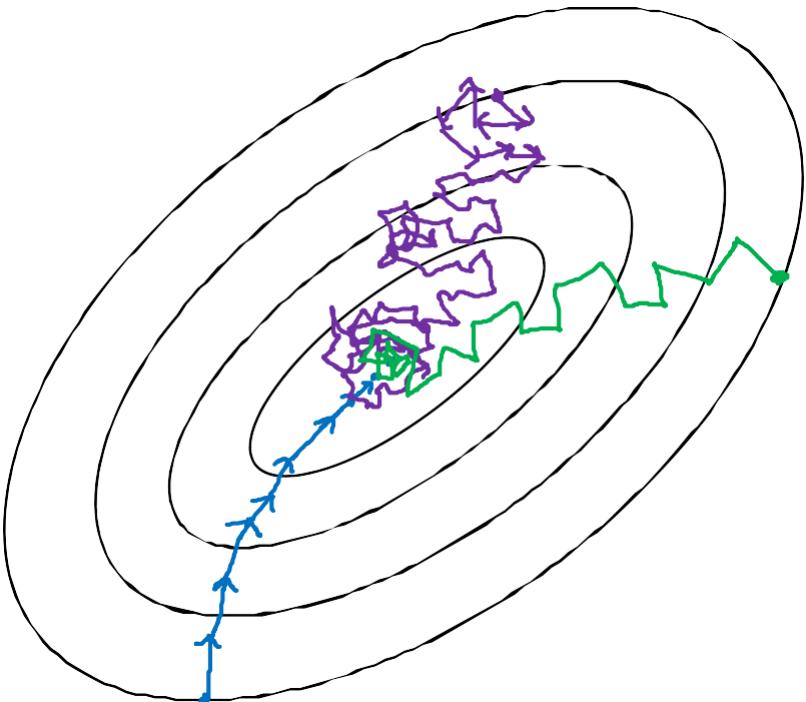


Plot  $J^{st}$  w.r.t.  $x^{st}, y^{st}$

# Choosing your mini-batch size

- If mini-batch size =  $m$  : Batch gradient descent.  $(X^{\{1\}}, Y^{\{1\}}) = (X, Y)$ .
- If mini-batch size = 1 : Stochastic gradient descent. Every example is it own  
 $(X^{S13}, Y^{S13}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.

In practice: Somewhere in-between 1 and  $m$



Stochastic  
gradient  
descent

{  
Use speedups  
from vectorization

In-between  
(mini-batch size  
not too big/small)

}  
Fastest learning

- Vectorization.  
(~1000)
- Mini batches without  
processing entire train set.

Batch  
gradient descent  
(mini-batch size =  $m$ )

}  
Two long  
per iteration

# Choosing your mini-batch size

If small toy set : Use batch gradient descent.  
 $(m \leq 2000)$

Typical mini-batch sizes:

$$\rightarrow 64, 128, 256, 512 \quad \frac{1024}{2^{10}}$$

$2^6 \quad 2^7 \quad 2^8 \quad 2^9$



Make sure mini-batch fits in CPU/GPU memory.

$$X^{\{t\}}, Y^{\{t\}}$$

here



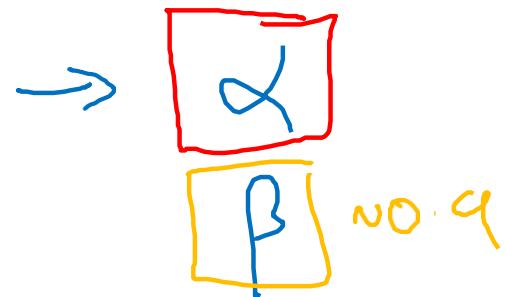
deeplearning.ai

# Hyperparameter tuning

---

## Tuning process

# Hyperparameters



$\beta_1, \beta_2, \epsilon$   
0.9 0.999  $10^{-8}$

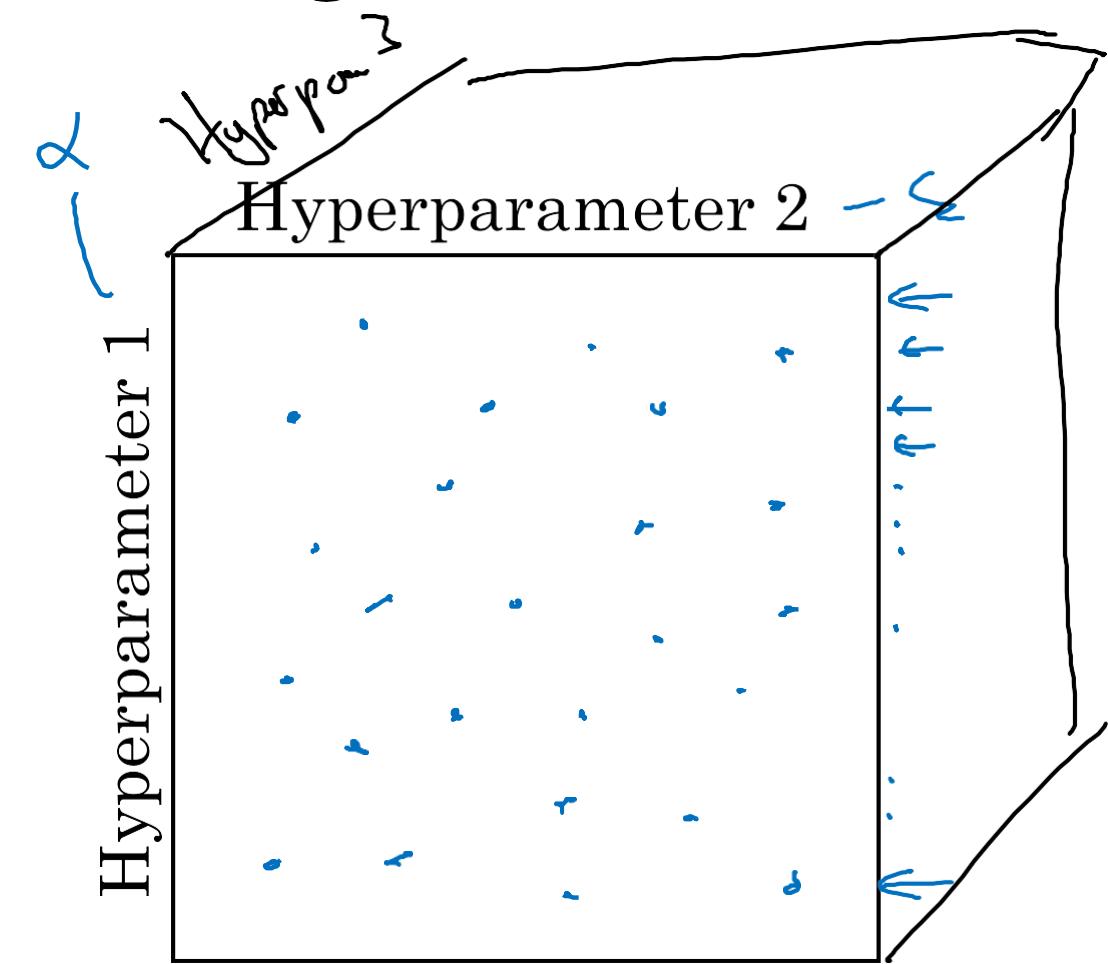
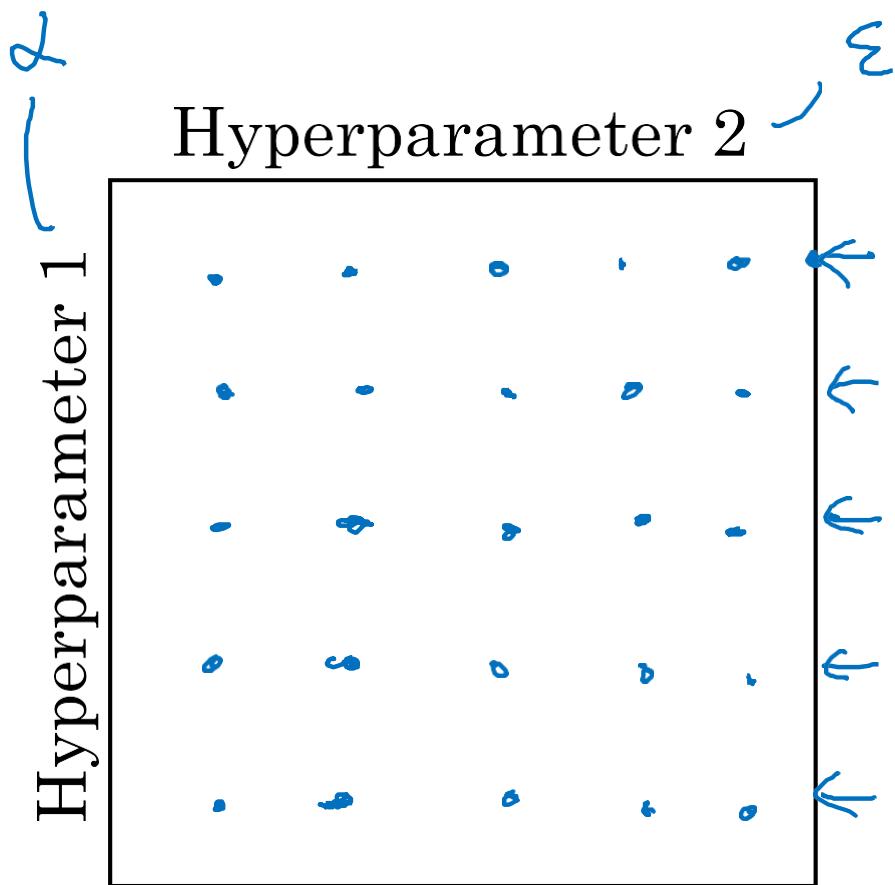
# layers

# hidden units

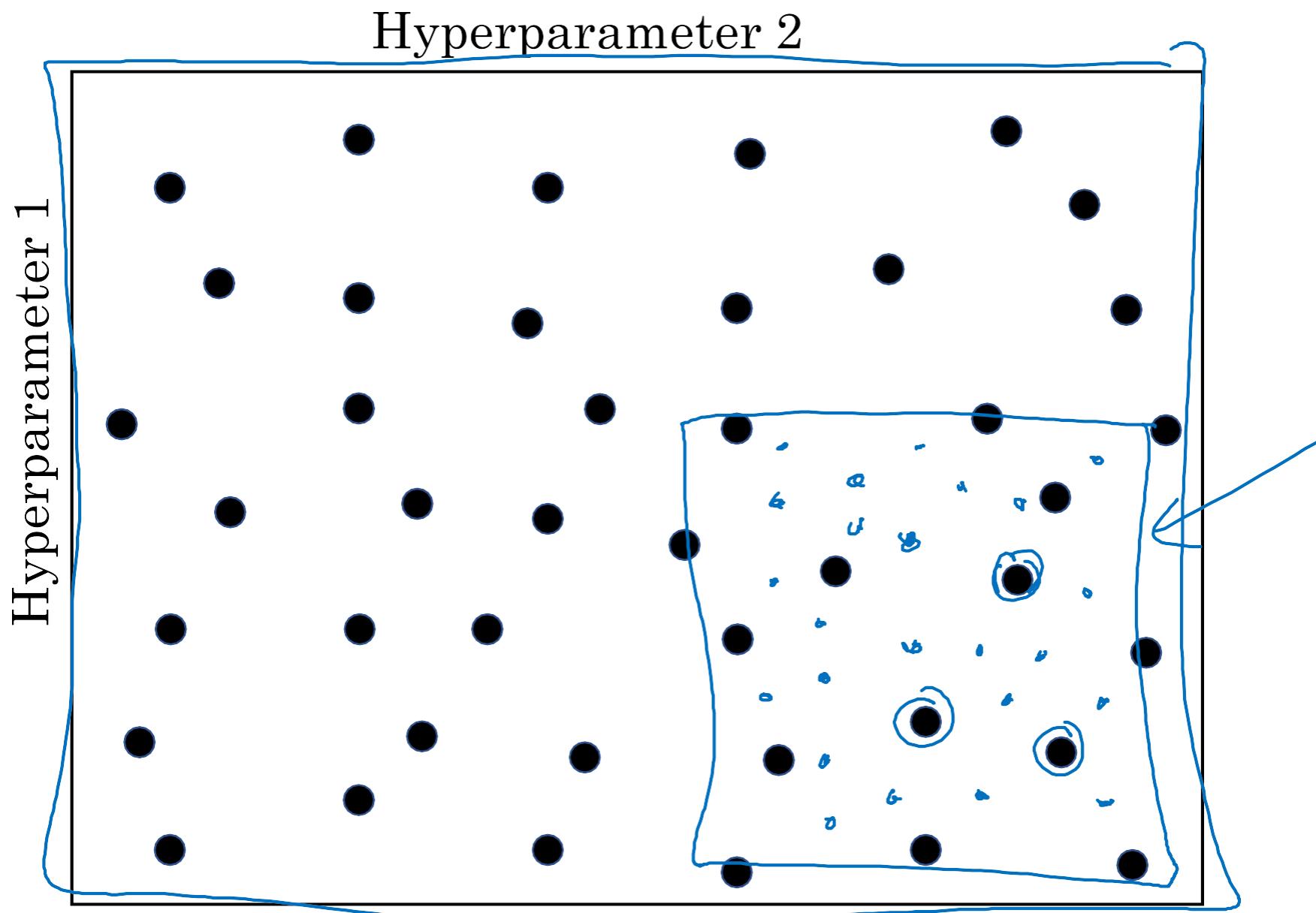
learning rate decay

mini-batch size

# Try random values: Don't use a grid



# Coarse to fine





deeplearning.ai

Multi-class  
classification

---

Softmax regression

# Recognizing cats, dogs, and baby chicks , other



3

1

2

0

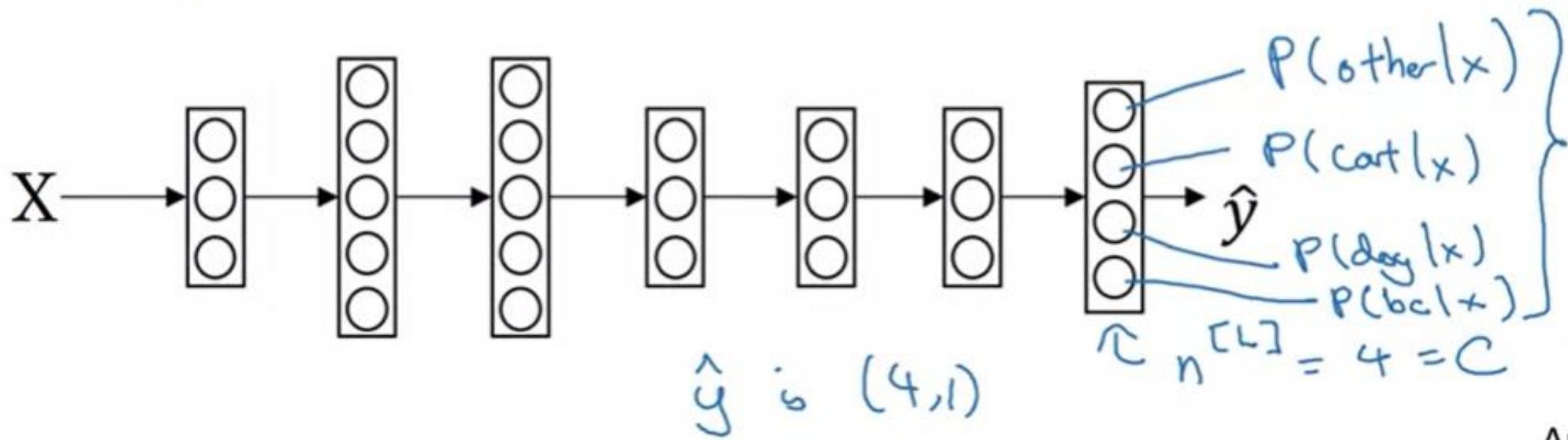
3

2

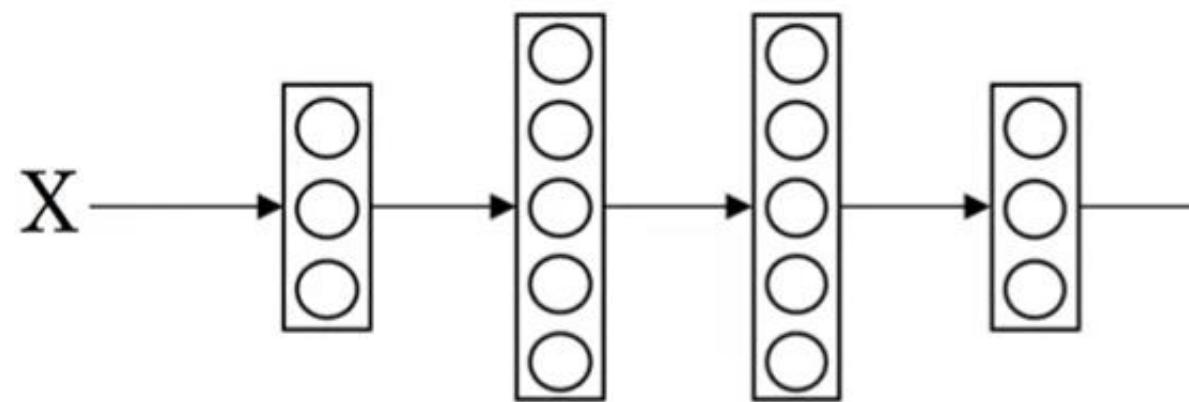
0

1

$$C = \# \text{classes} = 4 \quad (0, \dots, 3)$$



# Softmax layer



$$\rightarrow z^{[L]} = \omega^{[L]} a^{[L-1]} + b^{[L]}$$

$(4,1)$

Activation function:

$$\rightarrow t = e^{(z^{[L]})}$$

$(4,1)$

$$\rightarrow a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}, \quad a_i^{[L]} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

$$a^{[L]} = g^{[L]}(z^{[L]}) \quad (4,1)$$

layer L

$$\begin{aligned} & \frac{e^5}{176.3} = 0.842 \\ & \frac{e^2}{176.3} = 0.042 \\ & \frac{e^{-1}}{176.3} = 0.002 \\ & \frac{e^3}{176.3} = 0.114 \end{aligned}$$

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \sum_{j=1}^4 t_j = 176.3$$

$$a^{[L]} = \frac{t}{176.3}$$



deeplearning.ai

## Multi-class classification

---

# Training a softmax classifier



0:00 / 10:07



# Loss function

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \text{cat} \quad y_2 = 1$$

$y_1 = y_3 = y_4 = 0$

$$\underline{\ell(\hat{y}, y)} = \frac{-\sum_{j=1}^4 y_j \log \hat{y}_j}{\text{small}}$$

$$a^{(i)} \Rightarrow \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$$J(w^{(i)}, b^{(i)}, \dots) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$$

$$-\underline{y_2 \log \hat{y}_2} = -\underline{\log \hat{y}_2}$$

Make  $\hat{y}_2$  big.

2	3	7	4	6	2	9
6	6	9	8	7	4	3
3	4	8	3	8	9	7
7	8	3	6	6	3	4
4	2	1	8	3	4	6
3	2	4	1	9	8	3
0	1	3	9	2	1	4

This is the input matrix in 17 th slide of  
CMPE362\_W13\_CNN-part1.pdf

1	3	2	1	3
2	9	1	1	5
1	3	2	3	2
8	3	5	1	0
5	6	1	2	9

This is the input matrix in 33 rd slide of  
CMPE362\_W13\_CNN-part1.pdf