

# 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

-  $\frac{1}{n} \sum_{i=1}^n (h_\theta(x^{(i)}) - y^{(i)})^2$  -  
Setting up your  
ML application

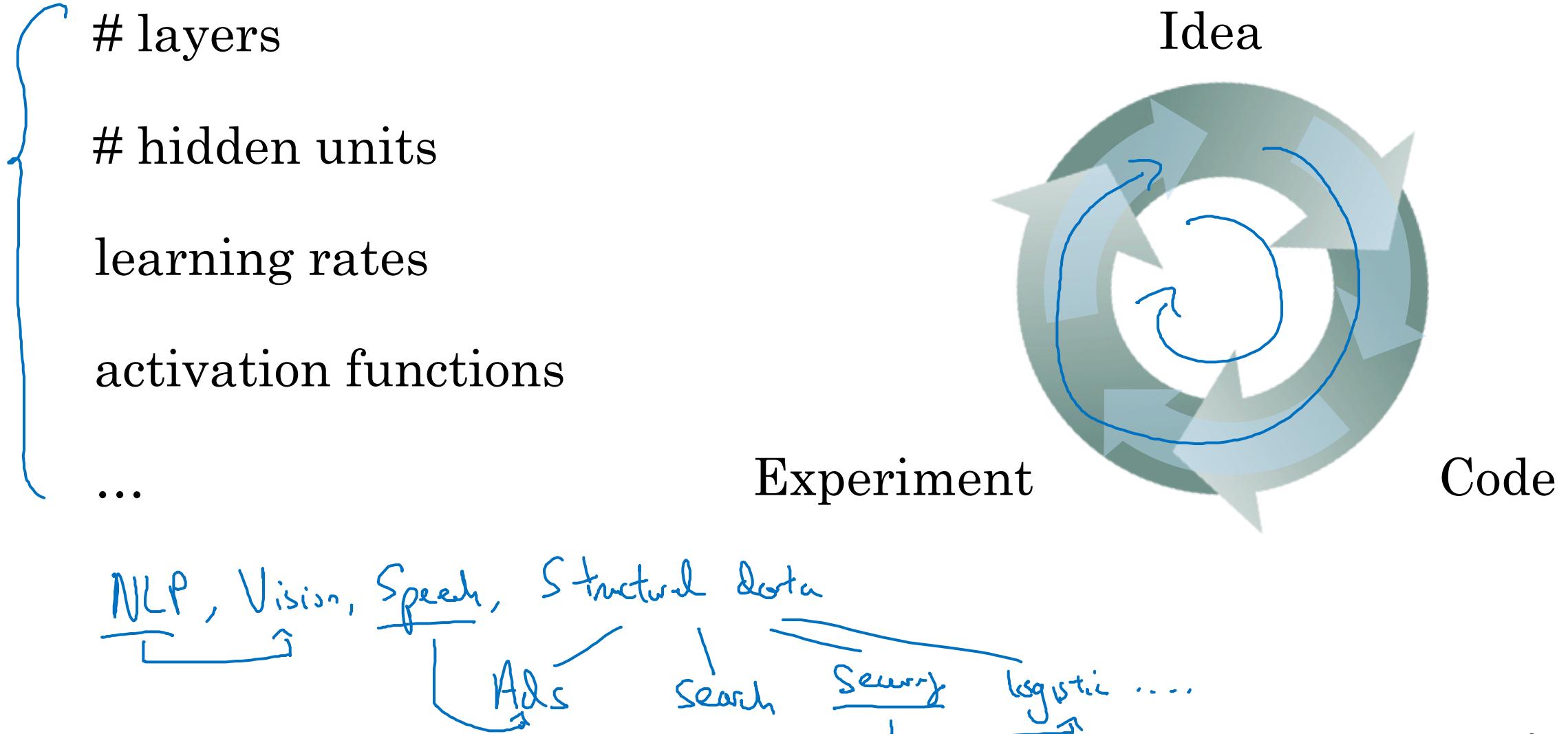
---

Train/dev/test  
sets

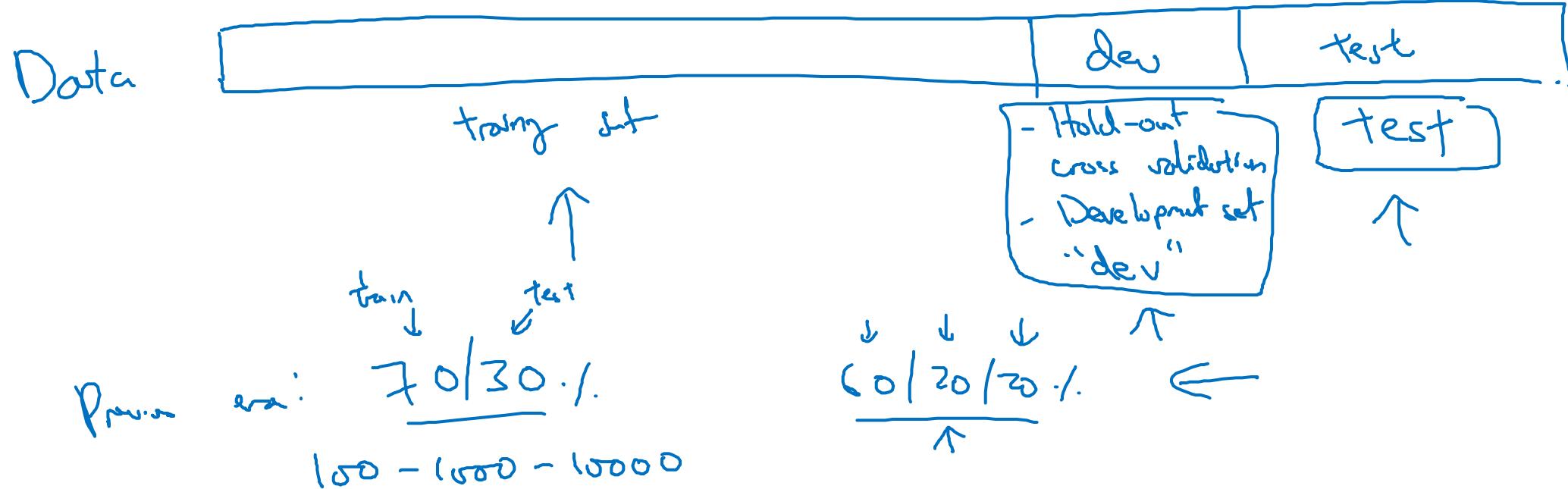
validation



# Applied ML is a highly iterative process



# Train/dev/test sets



Big data! 1,000,000

10,000 10,000

98 / 1 / 1 %.

99.5 { 25 / 25 %.

· 4 { - 1 · 1.

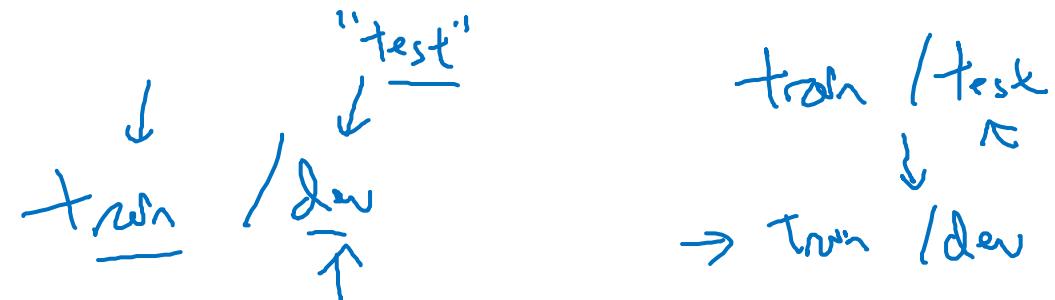
# Mismatched train/test distribution

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.)



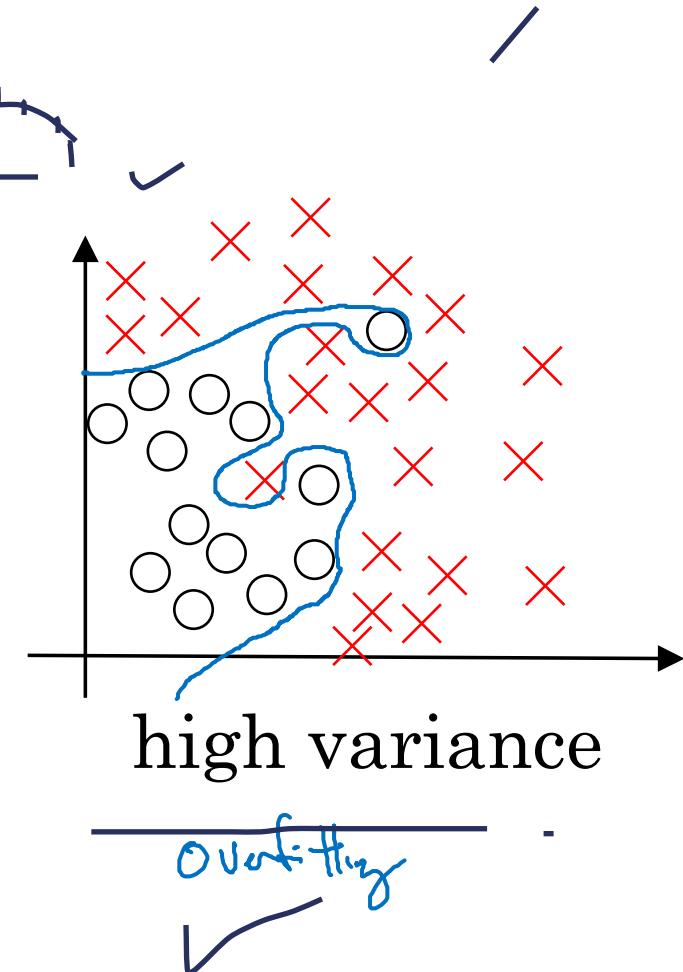
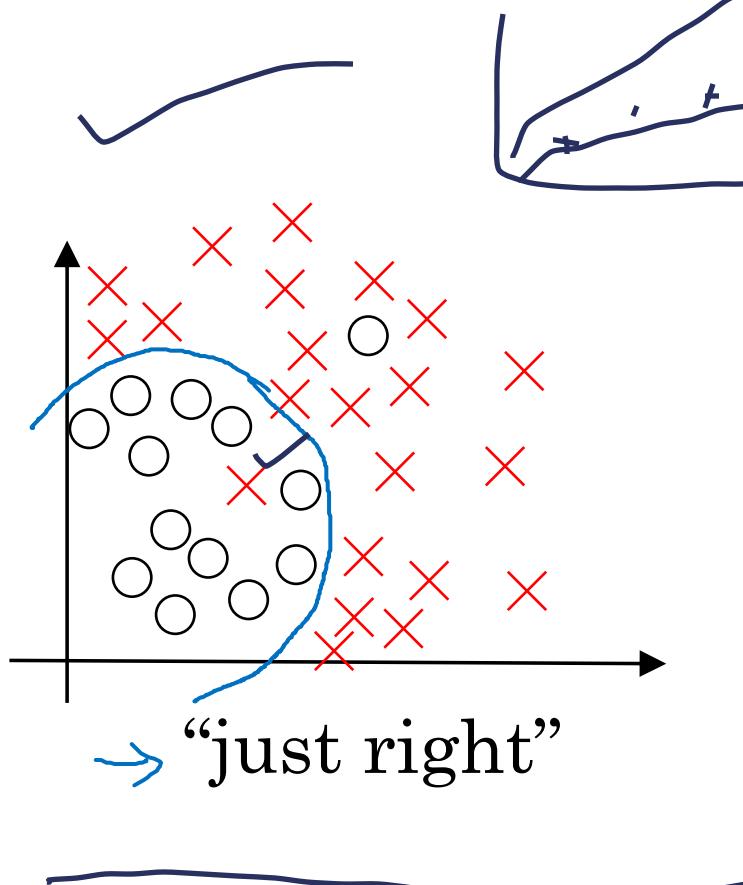
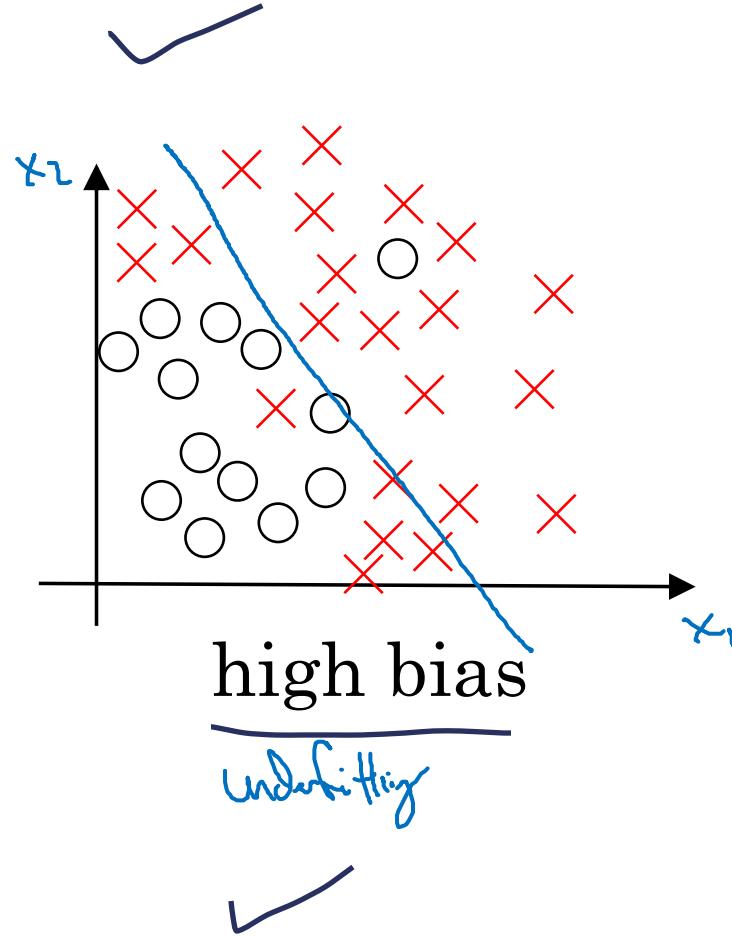
deeplearning.ai

Setting up your  
ML application

---

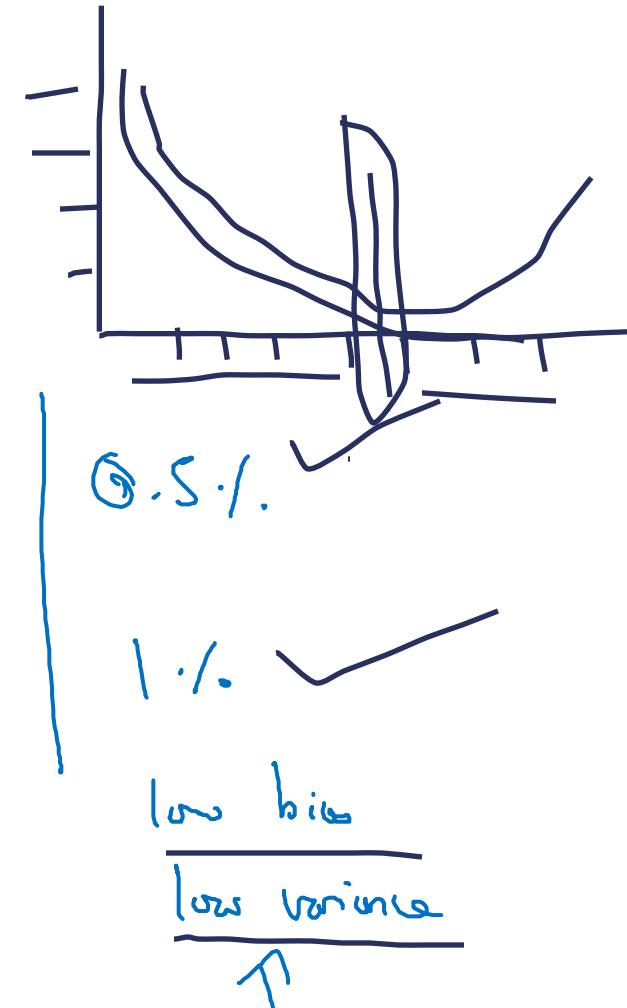
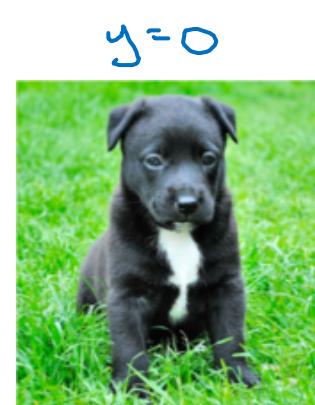
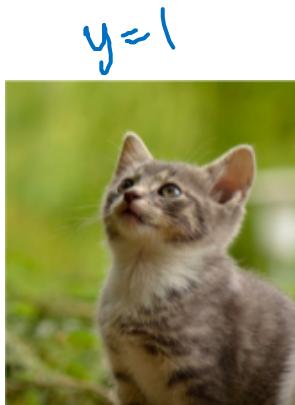
Bias/Variance

# Bias and Variance



# Bias and Variance

## Cat classification



Train set error:

✓  
1%

Dev set error:

11%

high variance  
↑

15% ↙

16% ↘

high bias  
↑  
↑

15%

30%

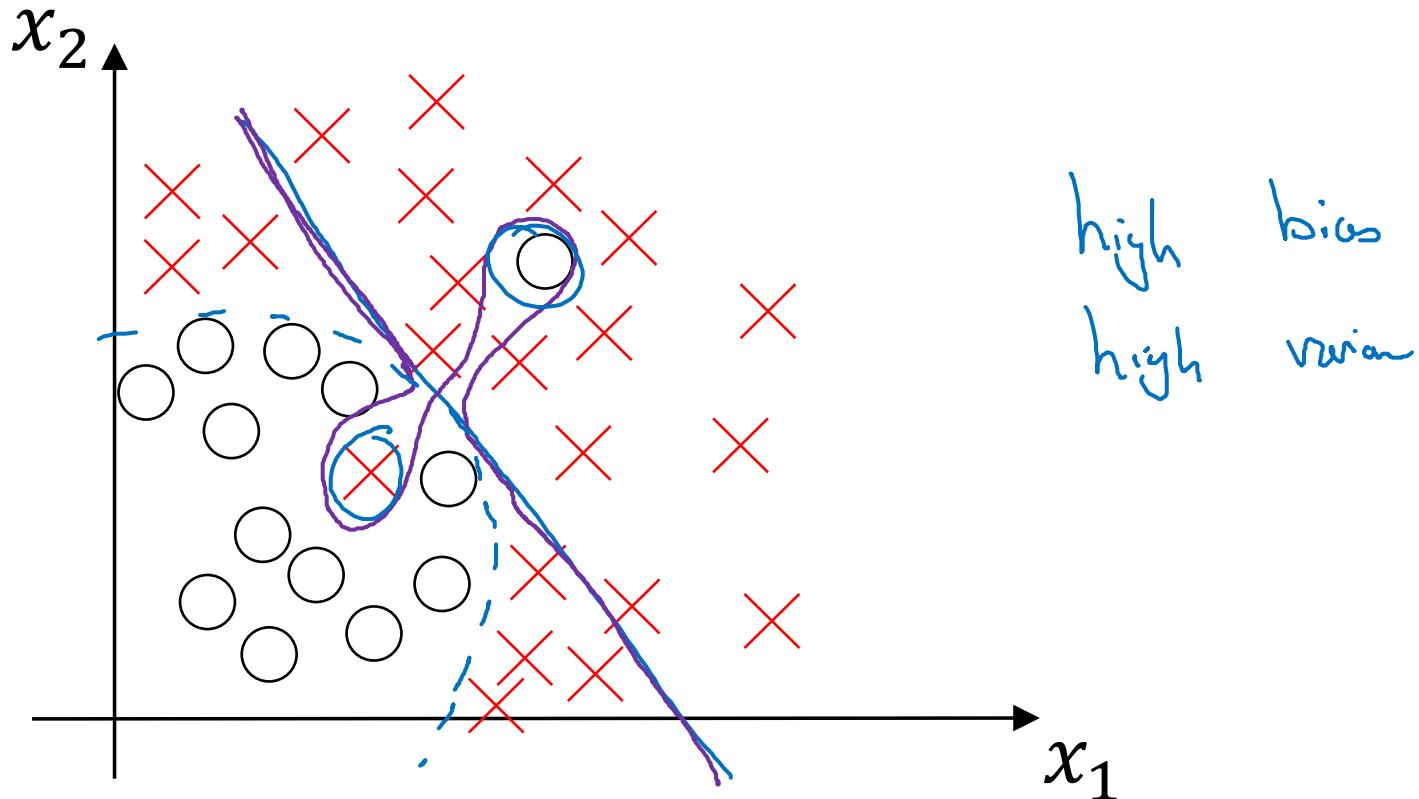
high bias  
& high varan

Human: ~20%

Optimal (Bayes) error: ~20% to 15%

Blurry images

# High bias and high variance





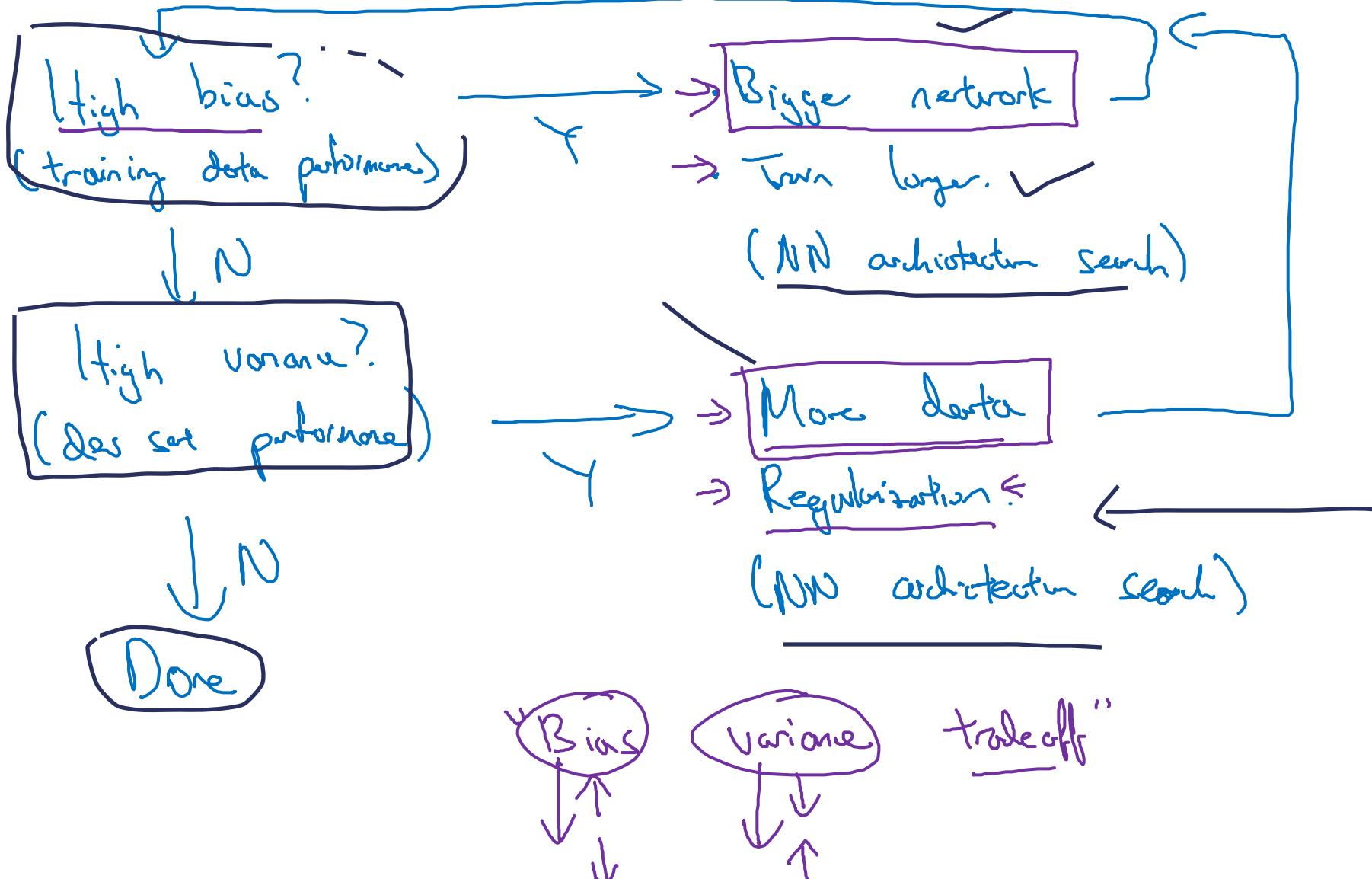
deeplearning.ai

# 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  
lambd

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

$$+ \frac{\lambda}{2m} \|w\|_2^2$$

$$+ \cancel{\frac{\lambda}{2m} b^2}$$

omit

L<sub>2</sub> regularization

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

L<sub>1</sub> regularization

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$$

w will be sparse

# Neural network

$$\rightarrow J(w^{(1)}, b^{(1)}, \dots, w^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(\hat{y}^{(i)}, y^{(i)})}_{\text{"Frobenius norm"}}$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l)}} \sum_{j=1}^{n^{(l+1)}} (w_{ij}^{(l)})^2$$

$$w^{(l)} : (n^{(l)}, n^{(l+1)})$$

$$\frac{\lambda}{2m} \left[ \sum_{l=1}^L \|w^{(l)}\|_F^2 \right] + \sum_{l=1}^L \sum_{i=1}^{n^{(l)}} |w_{i1}^{(l)}|^2$$

Weight decay

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

$$\rightarrow \underline{w^{(l)}} := w^{(l)} - \alpha \Delta w^{(l)}$$

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

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

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



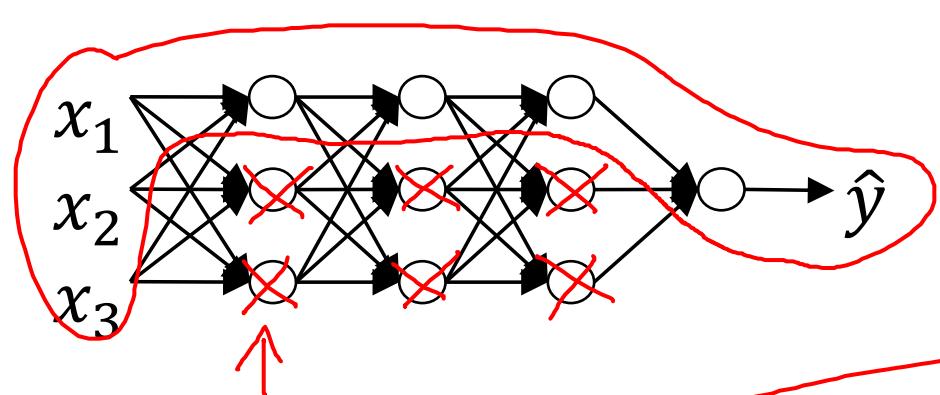
deeplearning.ai

# Regularizing your neural network

---

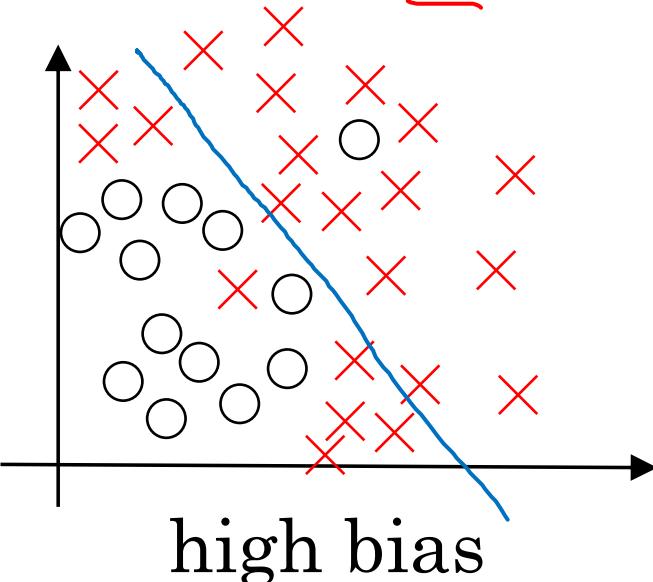
## Why regularization reduces overfitting

# How does regularization prevent overfitting?

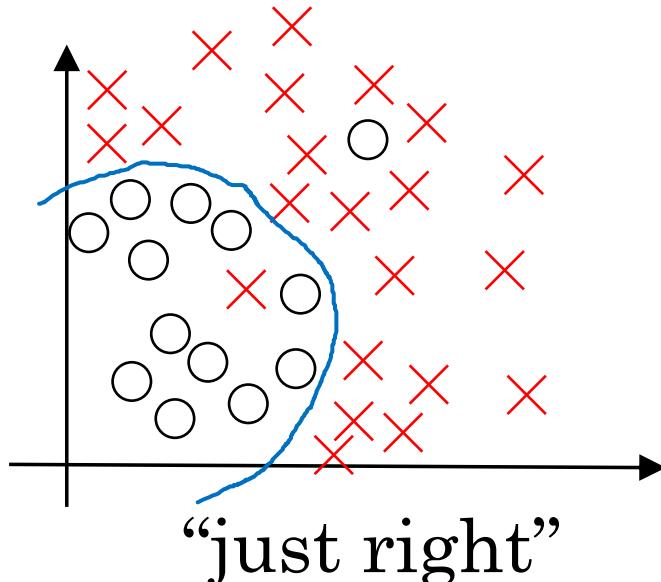


$$J(\boldsymbol{\theta}^{(m)}, \boldsymbol{b}^{(m)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

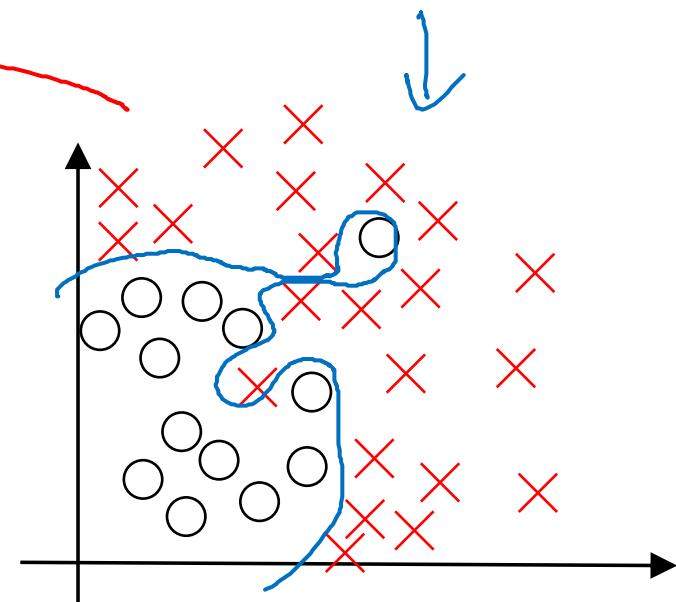
$\boldsymbol{w}^{(l)} \approx 0$



high bias

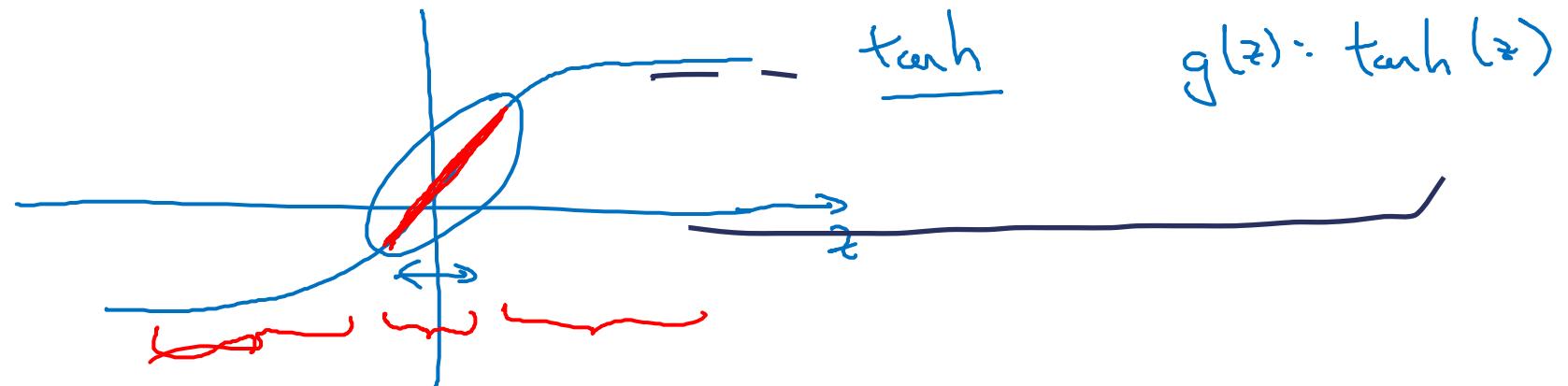


"just right"



high variance

# How does regularization prevent overfitting?

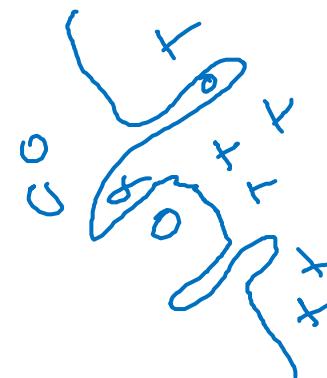


$$\lambda \uparrow$$

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

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

Every layer  $\approx$  linear.



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





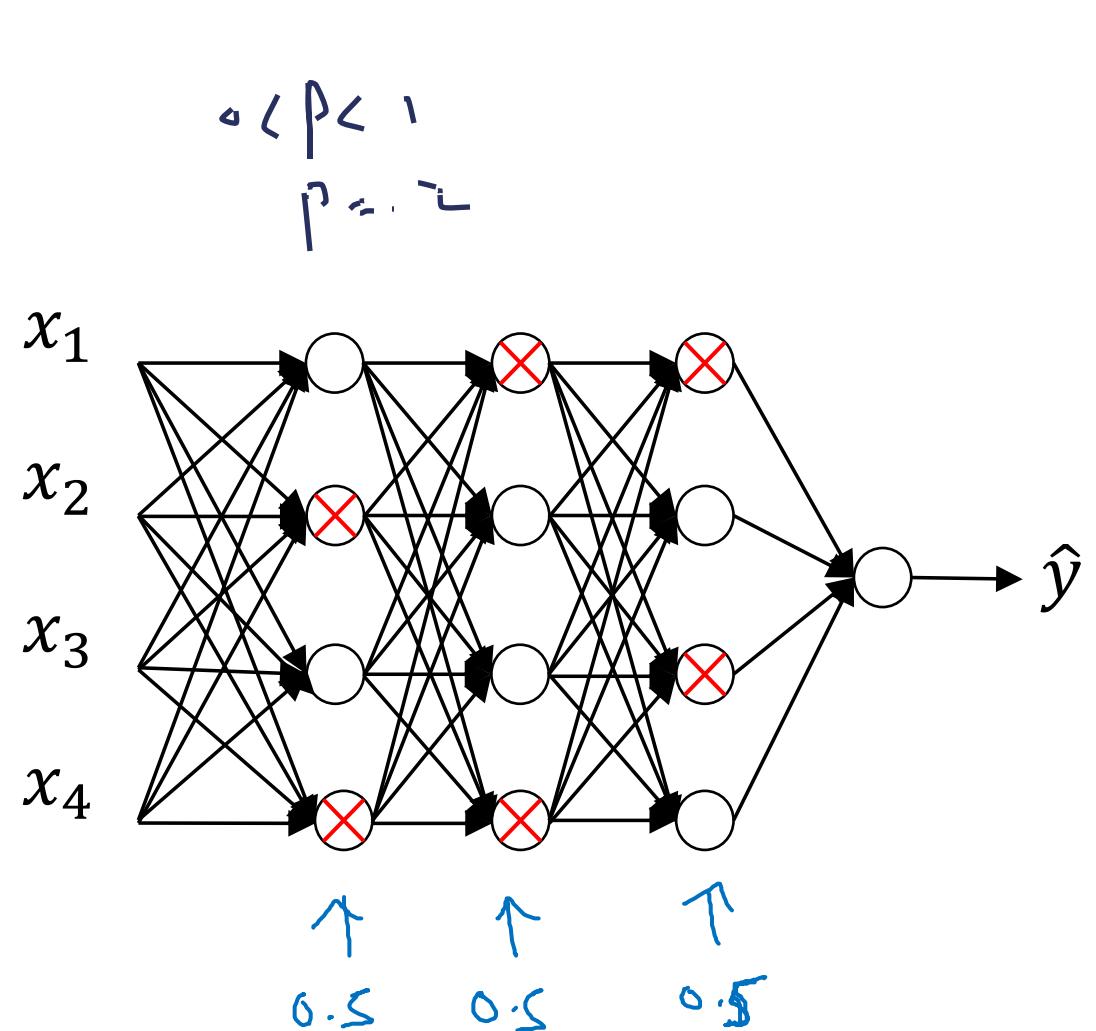
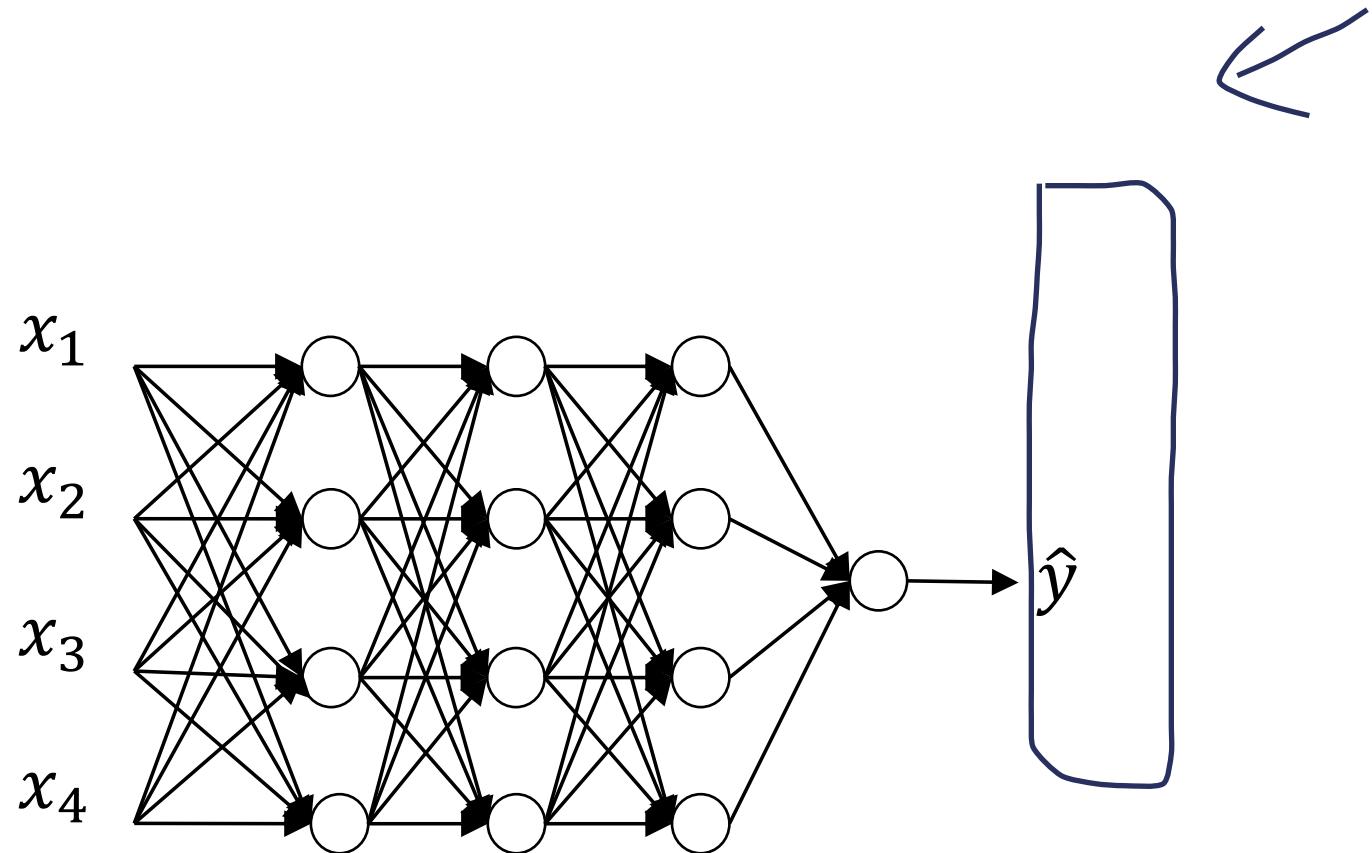
deeplearning.ai

Regularizing your  
neural network

---

Dropout  
regularization

# Dropout regularization



# Implementing dropout (“Inverted dropout”)

Illustrate with layer l=3. keep-prob = 0.8 0.2

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

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

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

50 units.  $\rightsquigarrow$  10 units shut off

$$\overline{z^{(4)}} = w^{(4)} \cdot \frac{a^{(3)}}{\cancel{C}} + b^{(4)}$$

$\cancel{C}$  reduced by 20%.

Test

$$1 = \underline{0.8}$$

# Making predictions at test time

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

No drop out.

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

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

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

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

$$\downarrow \hat{y}$$

$\lambda$  = keep-prob



deeplearning.ai

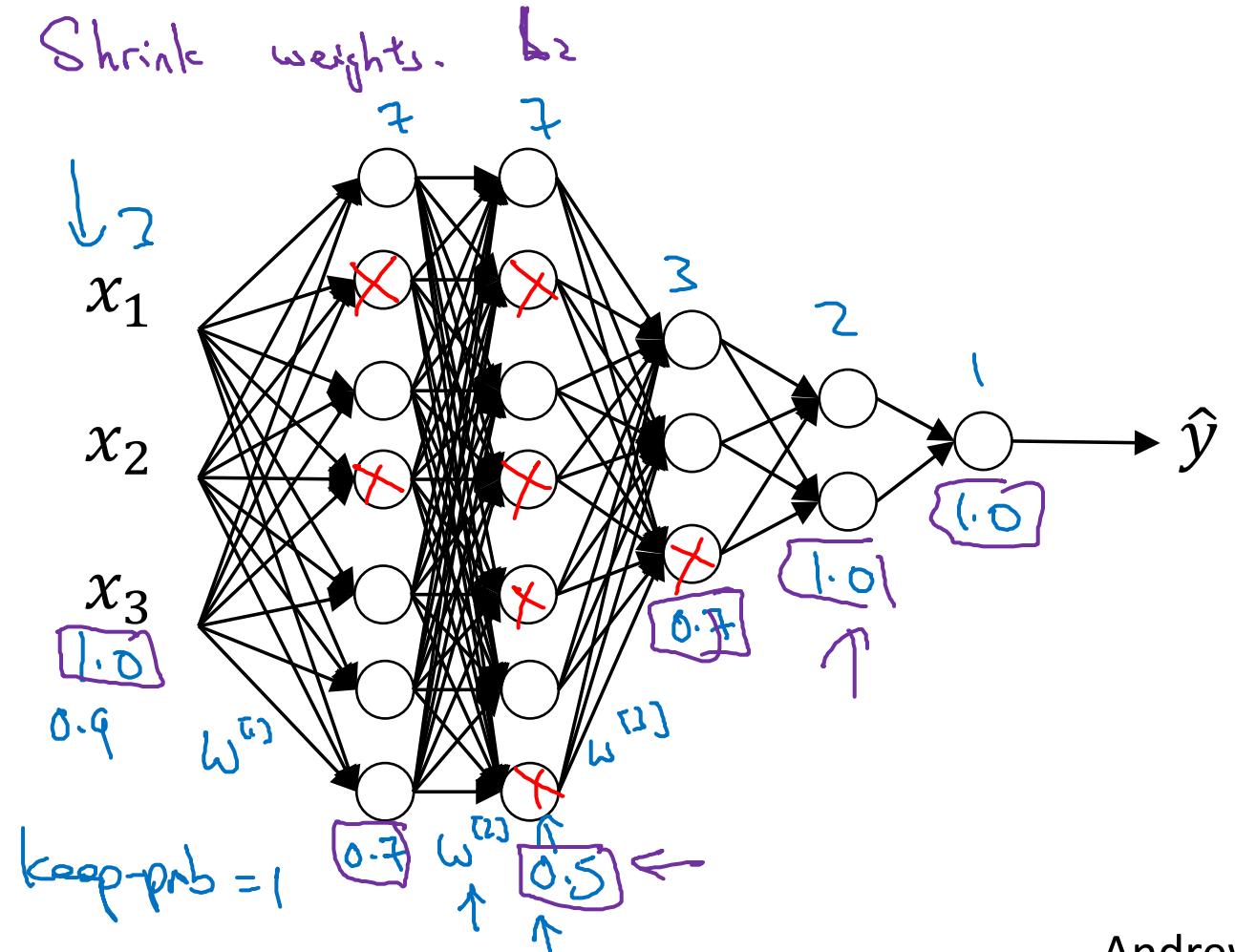
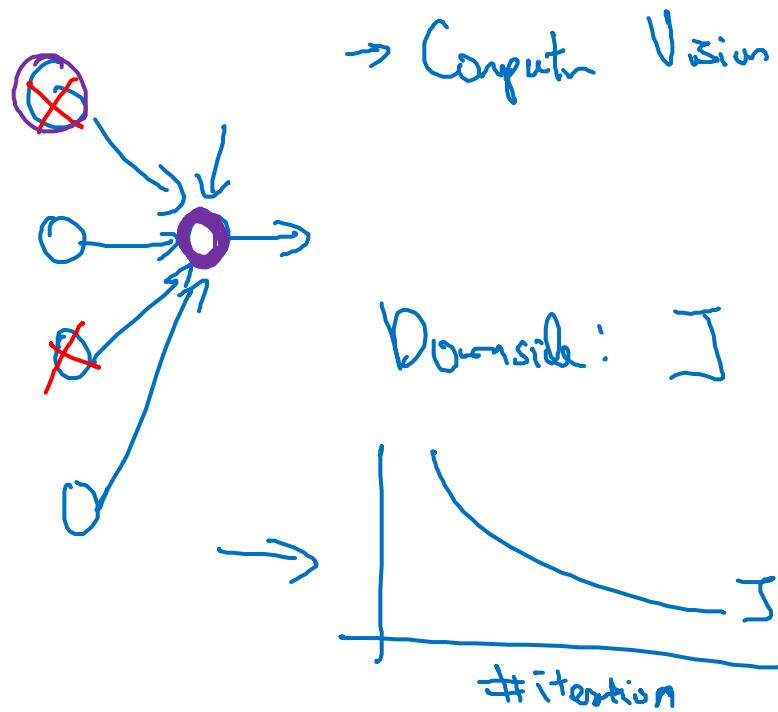
Regularizing your  
neural network

---

Understanding  
dropout

# Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights.  $\rightsquigarrow$  Shrink weights.  $b_2$





deeplearning.ai

# Regularizing your neural network

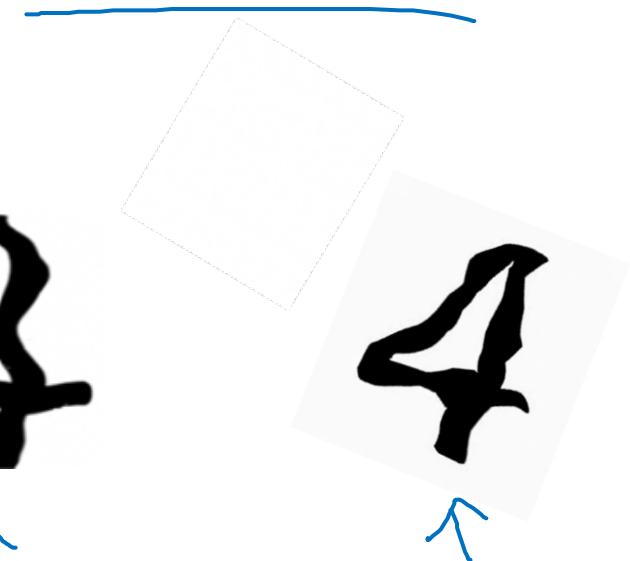
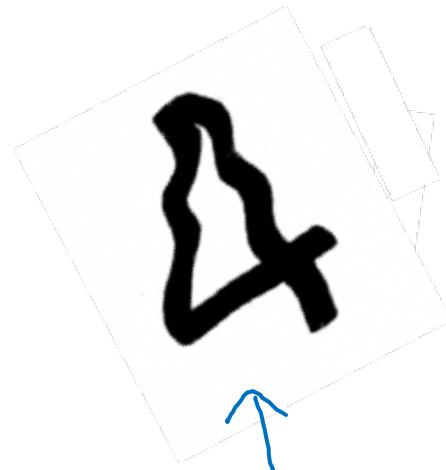
---

## Other regularization methods

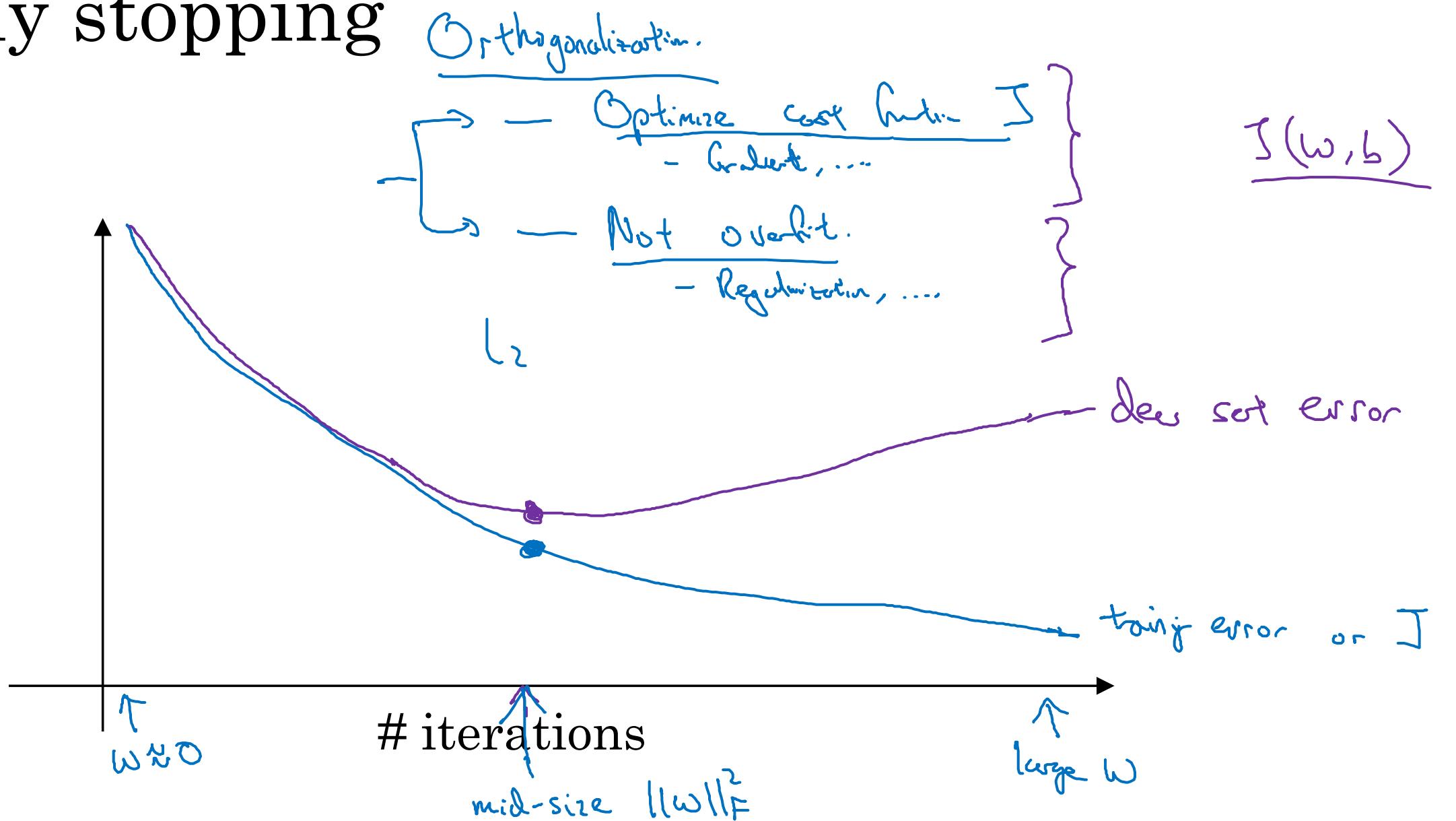
# Data augmentation



4



# Early stopping





deeplearning.ai

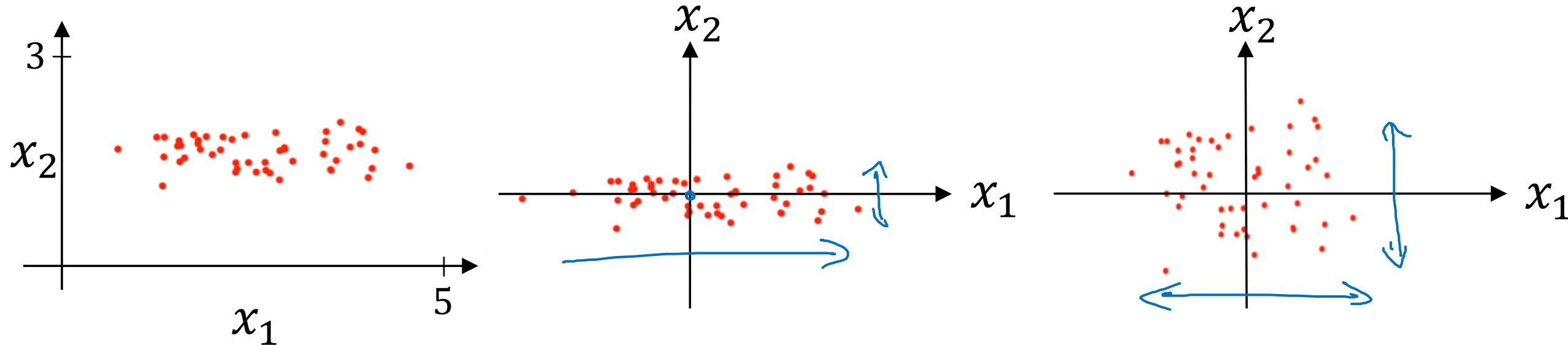
Setting up your  
optimization problem

---

Normalizing inputs

# Normalizing training sets

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



Subtract mean:

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

$$\underline{x := x - \mu}$$

Normalize variance

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

~ element-wise

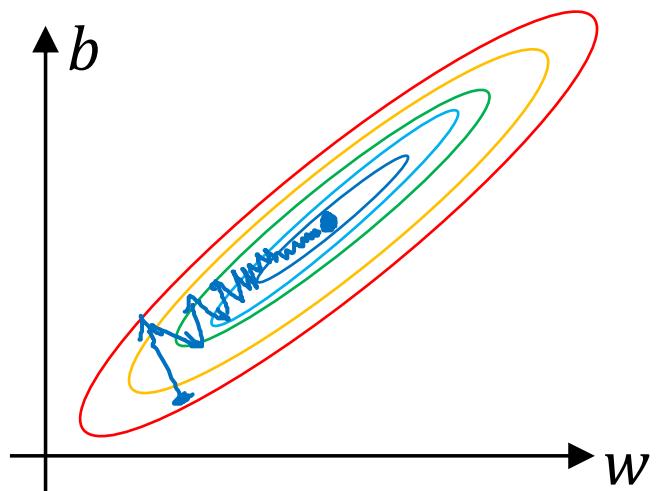
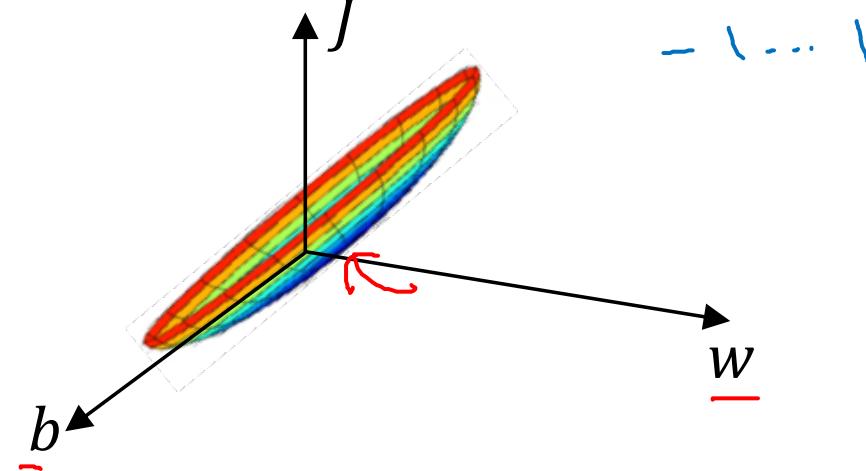
$$\underline{x / \sigma^2}$$

Use same  $\mu, \sigma^2$  to normalize test set.

# Why normalize inputs?

$\omega_1 \quad x_1: \frac{1 \dots 1000}{0 \dots 1} \leftarrow$   
 $\omega_2 \quad x_2: \frac{0 \dots 1}{-1 \dots 1} \leftarrow$

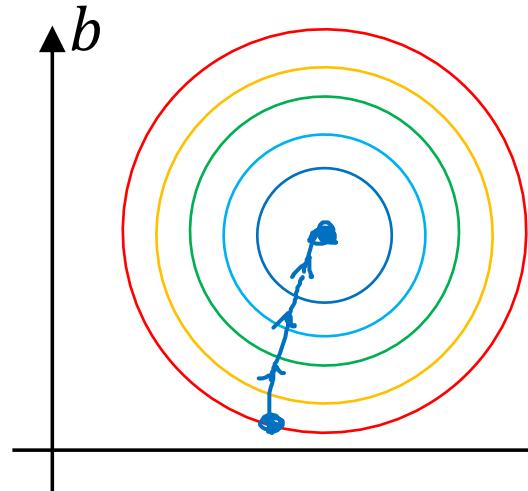
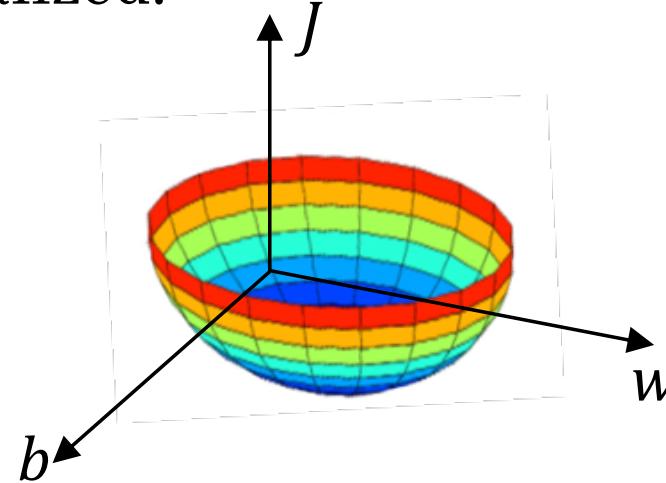
Unnormalized:



$x_1: 0 \dots 1$   
 $x_2: -1 \dots 1$   
 $x_3: 1 \dots 2$

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

Normalized:



$w$  Andrew Ng



deeplearning.ai

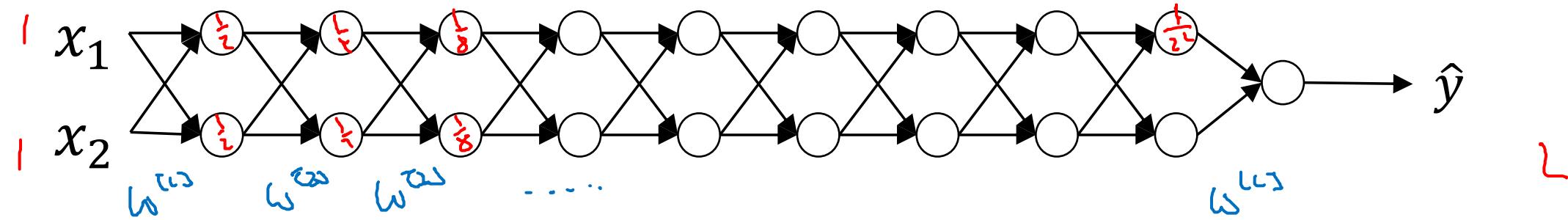
Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

# Vanishing/exploding gradients

$L=150$



$$\underline{g(z) = z} . \quad b^{[L]} = 0 .$$



$$w^{[1]} > I$$

$$w^{[2]} < I \quad \begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$$

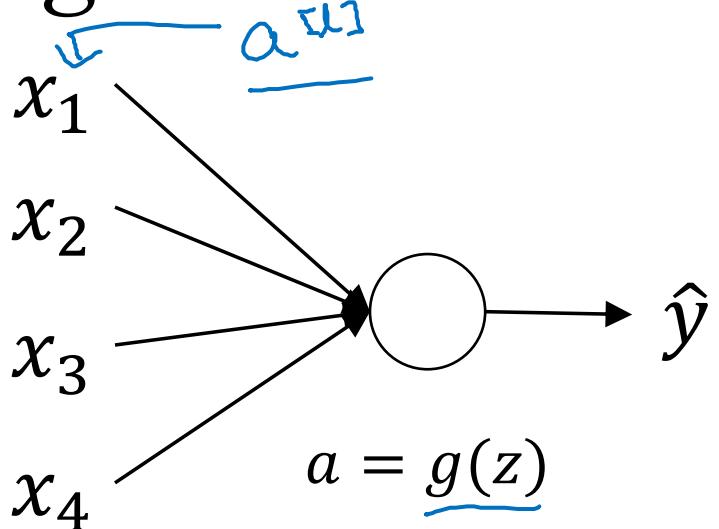
$$w^{[L]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 6.5 \end{bmatrix}$$

$$\hat{y} = w^{[L]} \begin{bmatrix} 0.5 & \\ & 1.5 & 0 \\ 0 & 0.5 & 6.5 \end{bmatrix}^{L-1} x$$

$$z^{[L-1]} = w^{[L]} x$$

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

# Single neuron example



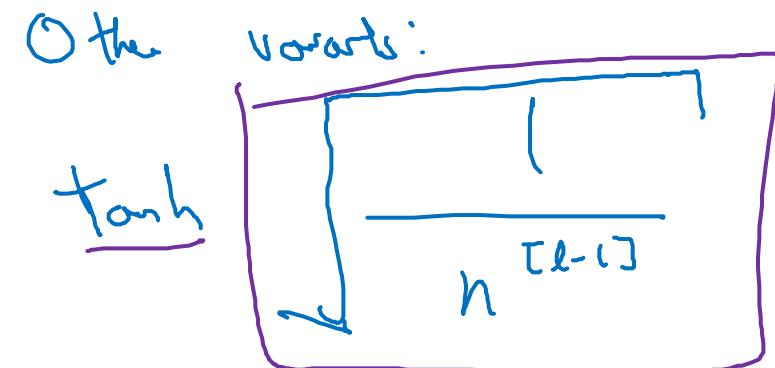
$$z = \underline{w_1}x_1 + \underline{w_2}x_2 + \cdots + \underline{w_n}x_n \quad \cancel{\text{if } n \text{ is large}}$$

$\downarrow$   
Large  $n \rightarrow$  Smaller  $w_i$

$$\text{Var}(w_i) = \frac{1}{n} \frac{2}{n}$$

$$\underline{w^{[l]}} = \text{np.random.randn}(\text{shape}) * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

$\downarrow$   
 $\text{ReLU}$        $\underline{g^{[l]}}(z) = \text{ReLU}(z)$



Xavier initialization  $\uparrow$

$$\frac{2}{n^{[l-1]} + n^{[l]}}$$

$\uparrow$



deeplearning.ai

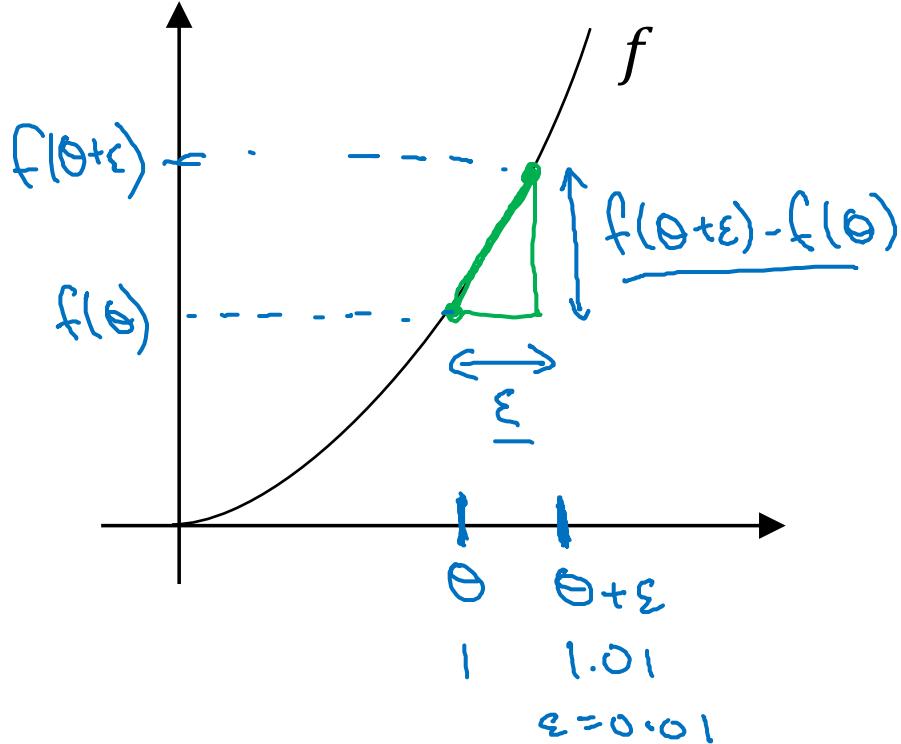
# Setting up your optimization problem

---

## Numerical approximation of gradients

# Checking your derivative computation

$$\begin{aligned} f(\theta) &= \underline{\theta^3} \\ \theta &\in \mathbb{R}. \\ \text{I} \end{aligned}$$



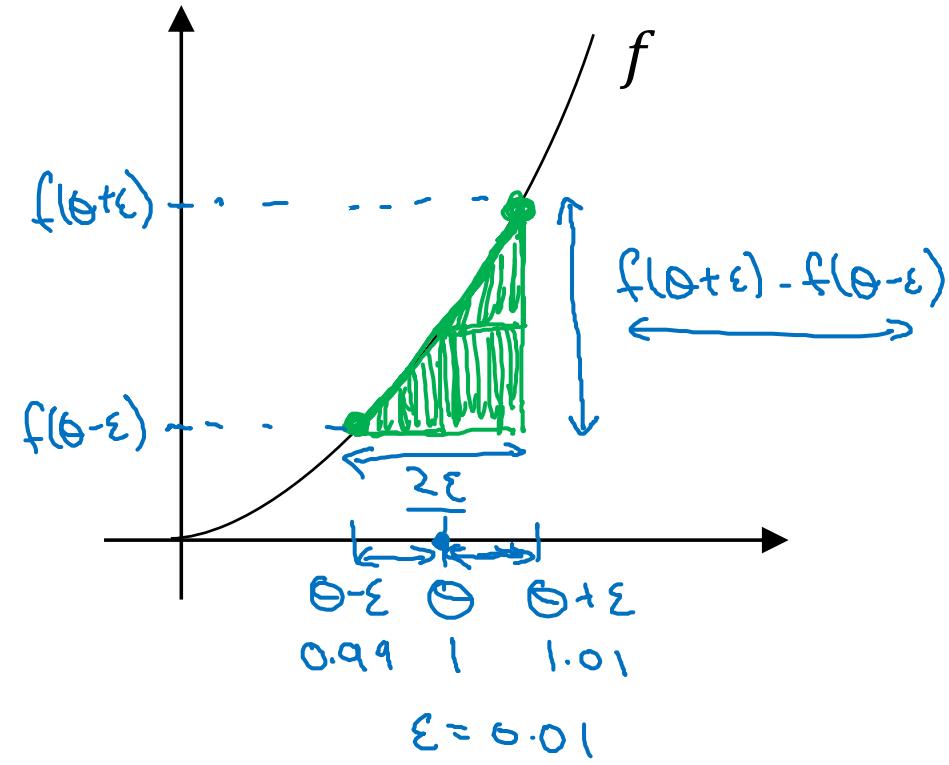
$$\begin{aligned} g(\theta) &= \frac{d}{d\theta} f(\theta) = f'(\theta) \\ g(\theta) &= 3\theta^2. \\ g(1) &= 3 \cdot (1)^2 = 3 \\ \text{when } \theta &= 1 \\ \frac{dw}{db} \end{aligned}$$

$$\begin{aligned} \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} &\approx g(\theta) \\ \frac{(1.01)^3 - 1^3}{0.01} &= 3.0301 \\ \frac{3.0301}{0.0301} &\approx 3 \end{aligned}$$

$$\begin{aligned} \theta &= 1 \\ \theta + \epsilon &= 1.01 \end{aligned}$$

# Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



$$\left[ \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \right] \approx g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

$\left\{ f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon}$	$\frac{\mathcal{O}(\epsilon^2)}{0.01} = \underline{0.0001}$	$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}$ $\uparrow \qquad \uparrow$	$\text{error: } \mathcal{O}(\frac{\epsilon}{0.01}) = 0.01$
--	---	---	--



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking

# Gradient check for a neural network

Take  $\underline{W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}}$  and reshape into a big vector  $\underline{\theta}$ .

$$J(\underline{w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}}) = J(\underline{\theta})$$

Take  $\underline{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}}$  and reshape into a big vector  $\underline{d\theta}$ .

Is  $d\theta$  the gradient of  $J(\theta)$ ?

# Gradient checking (Grad check)

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$$

for each  $i$ :

$$\rightarrow \underline{d\theta_{\text{approx}}[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \varepsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \varepsilon, \dots)}{2\varepsilon}$$

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i}$$

$$d\theta_{\text{approx}} \stackrel{?}{\approx} d\theta$$

Check

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$

$$\varepsilon = 10^{-7}$$

$$\approx \boxed{10^{-7} - \text{great!}} \leftarrow 10^{-5}$$

$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$



deeplearning.ai

# Setting up your optimization problem

---

## Gradient Checking implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{\partial \theta_{\text{approx}}^{[i]}}{\uparrow} \longleftrightarrow \frac{\partial \theta^{[i]}}{\uparrow}$$

- If algorithm fails grad check, look at components to try to identify bug.

$$\frac{\partial b^{[l]}}{\uparrow} \quad \frac{\partial w^{[l]}}{\uparrow}$$

- Remember regularization.

$$J(\theta) = \frac{1}{m} \sum_i f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_l \|w^{(l)}\|_F^2$$

$\frac{\partial \theta}{\partial \theta} = \text{gradient of } J \text{ wrt. } \theta$

- Doesn't work with dropout.

$\exists$

keep-prob = 1.0

- Run at random initialization; perhaps again after some training.

$w, b \text{ NO}$