

# 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>



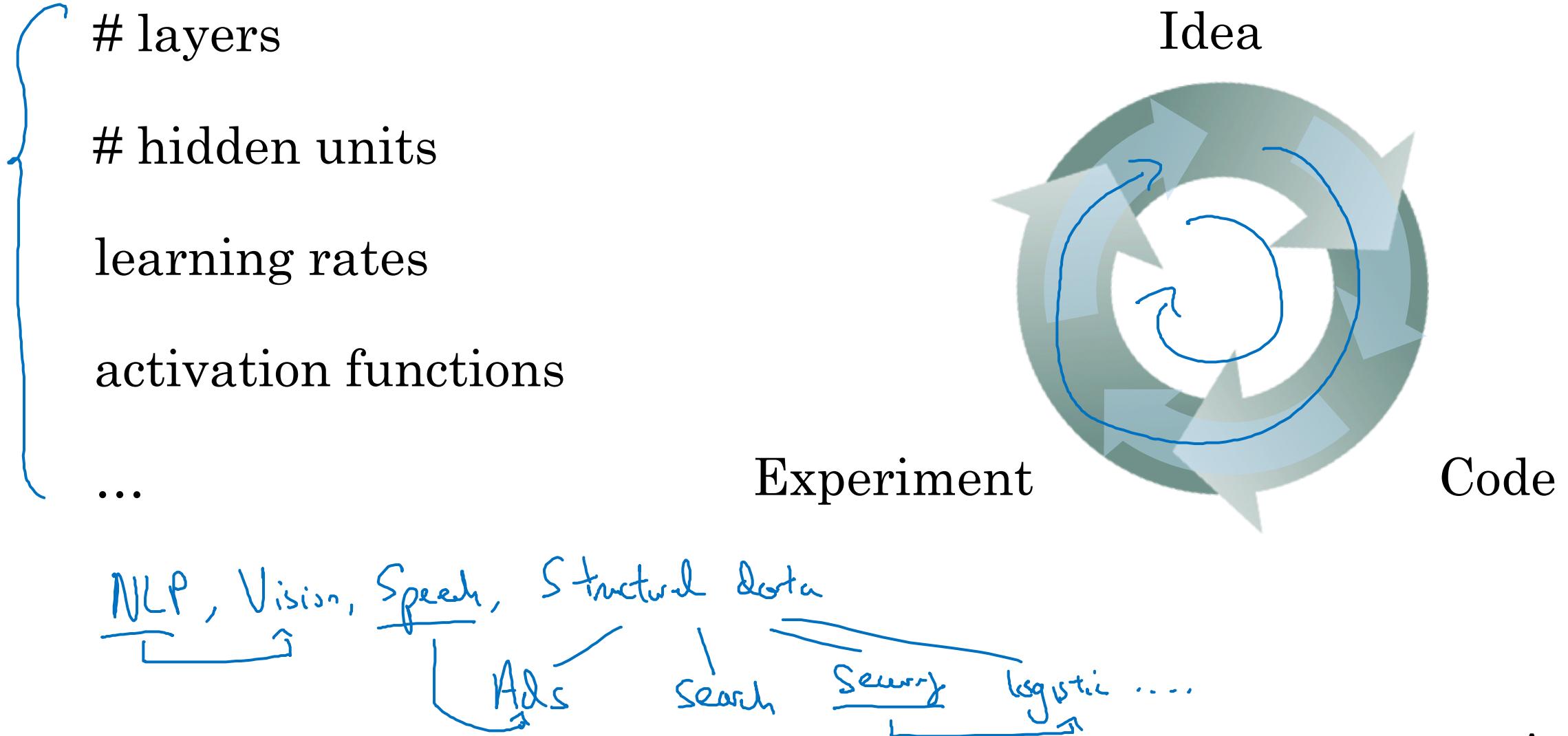
deeplearning.ai

Setting up your  
ML application

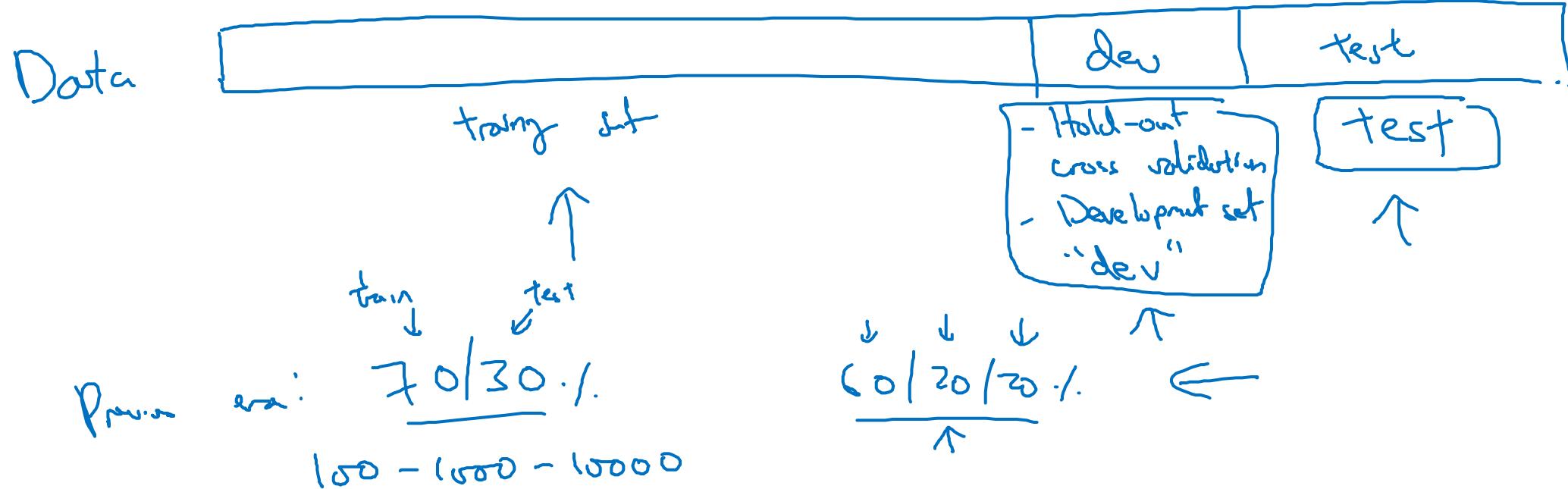
---

Train/dev/test  
sets

# 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 . . .

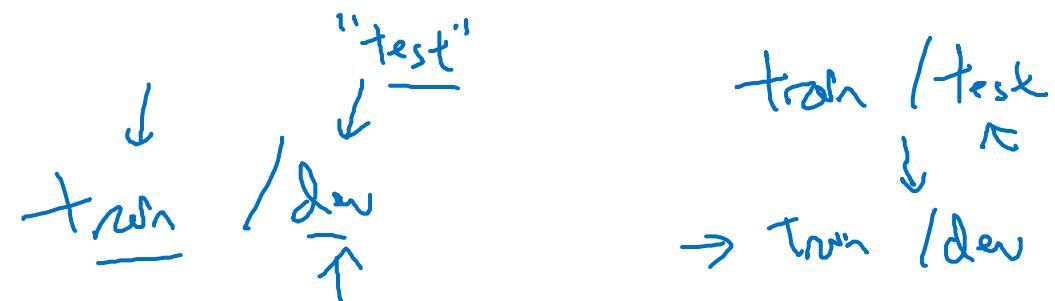
# Mismatched train/test distribution

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.)



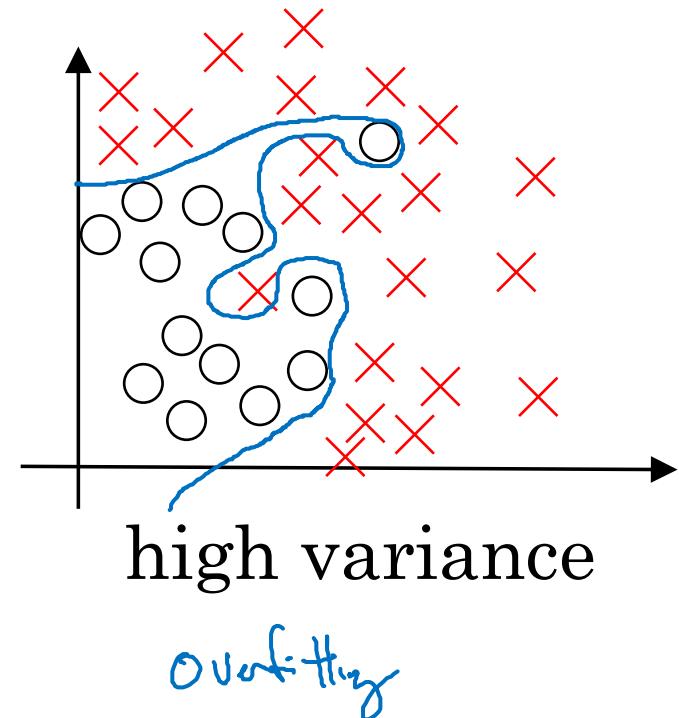
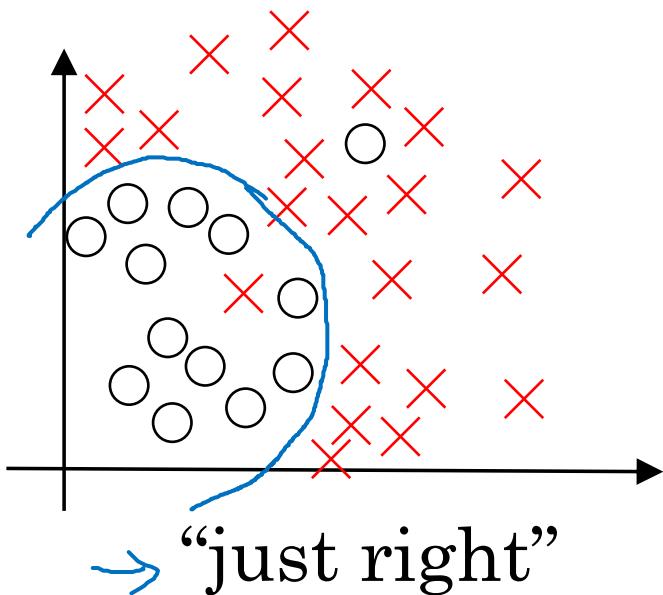
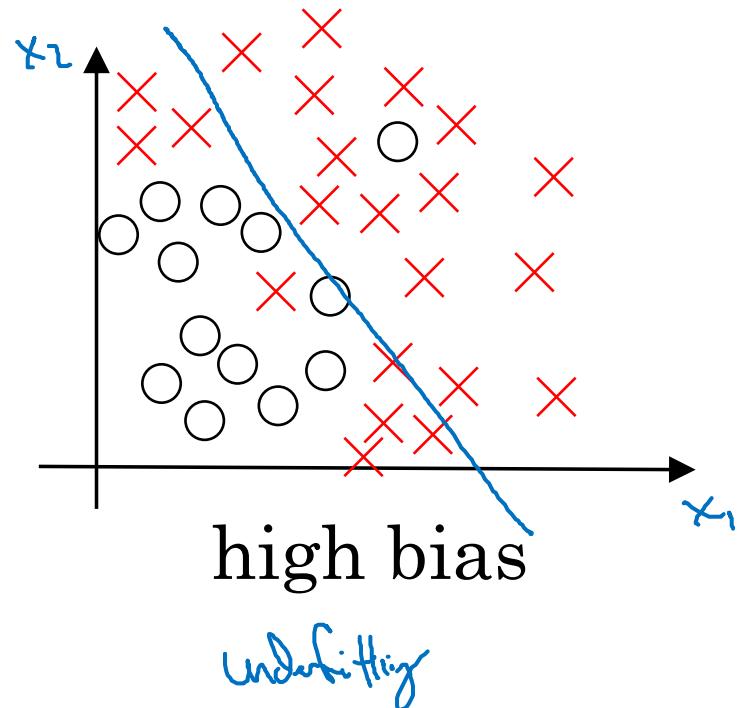
deeplearning.ai

Setting up your  
ML application

---

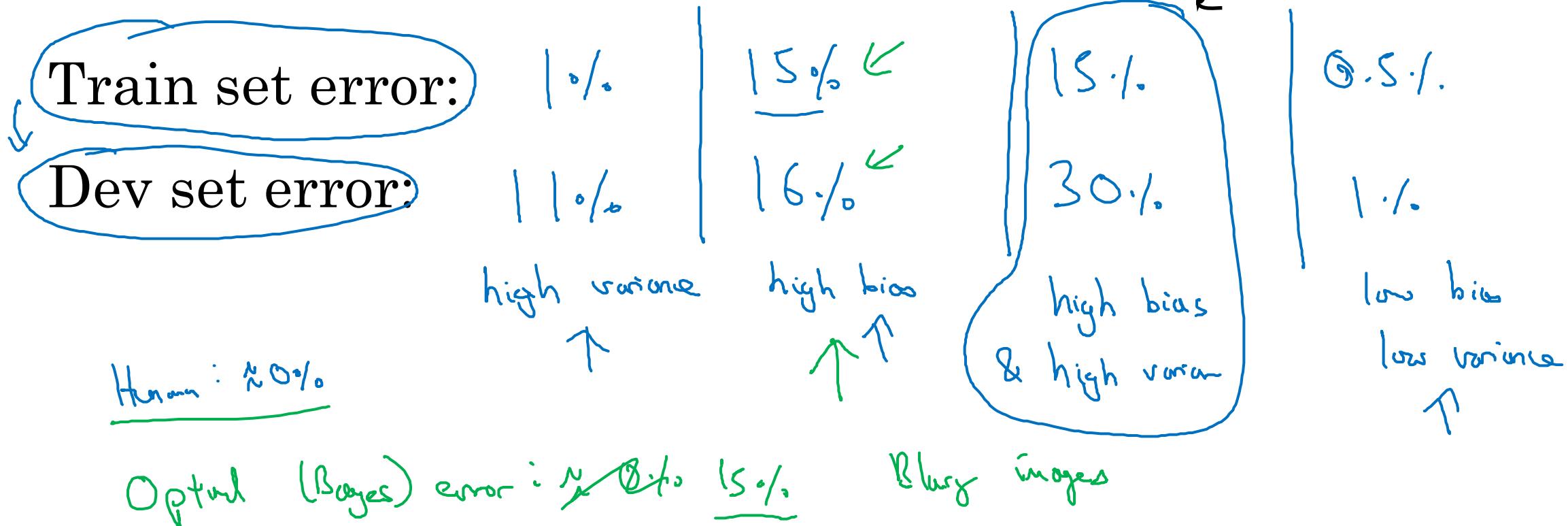
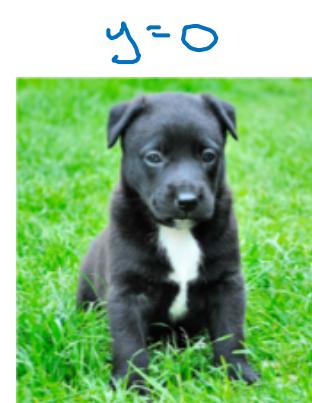
Bias/Variance

# Bias and Variance

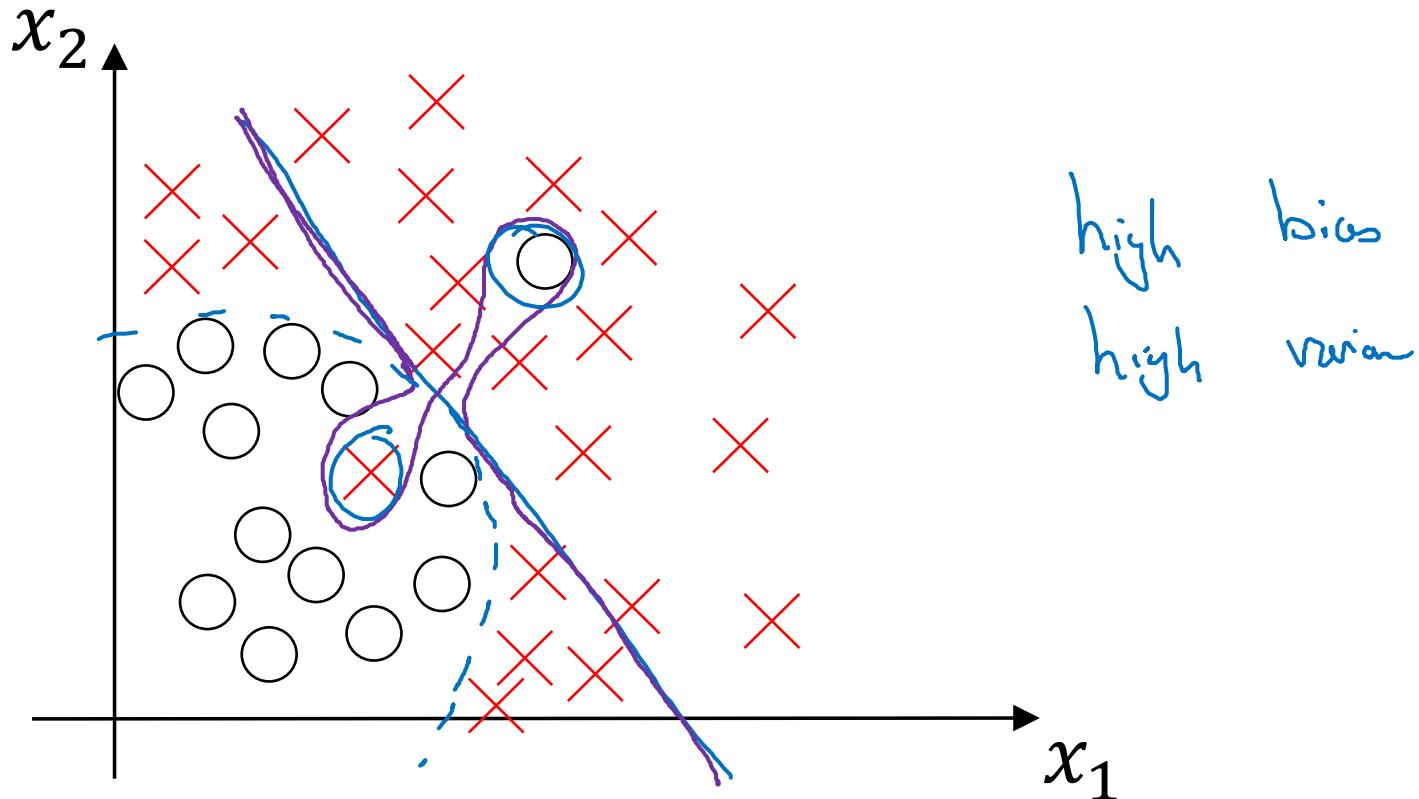


# Bias and Variance

## Cat classification



# High bias and high variance





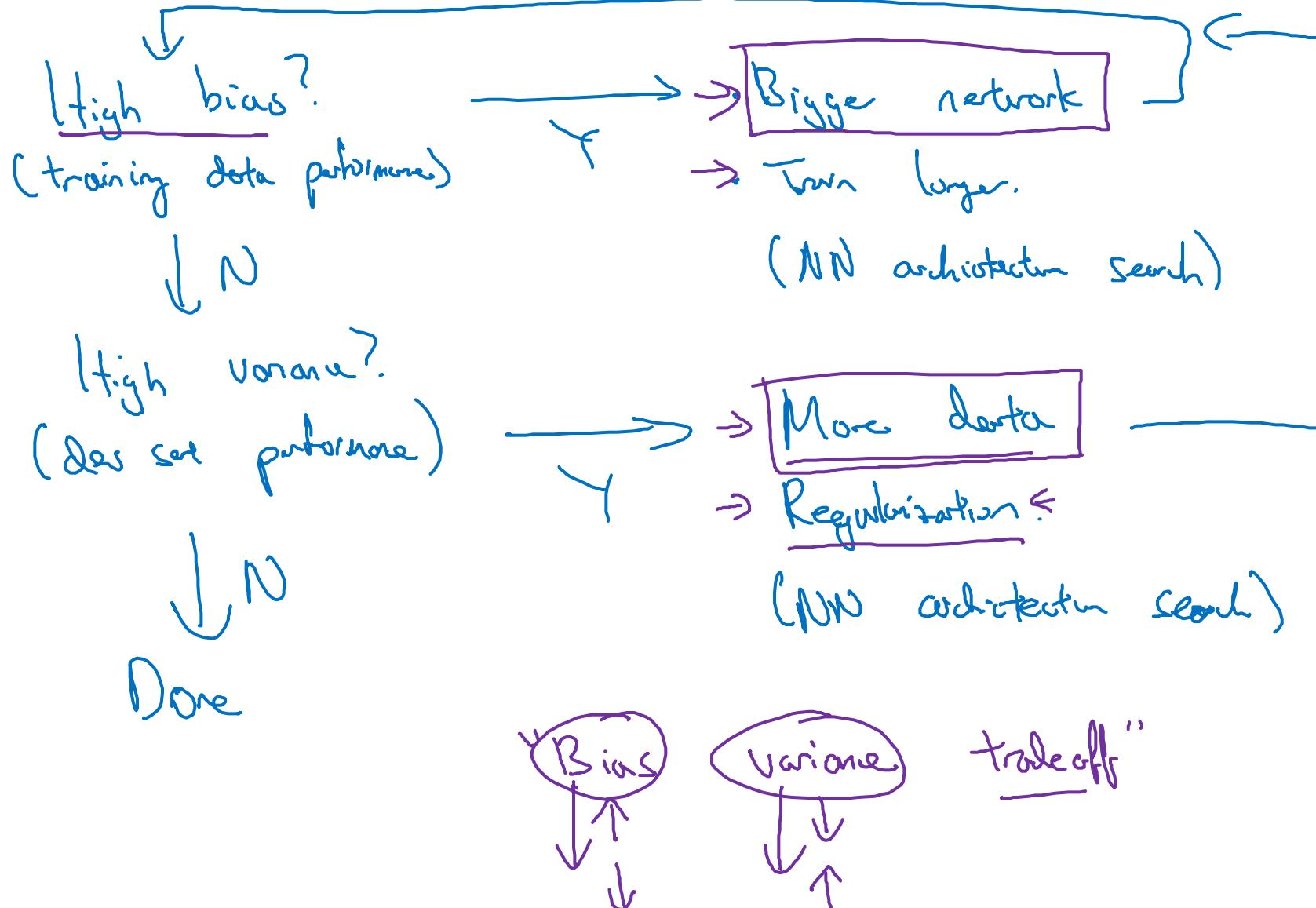
deeplearning.ai

# 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  
lambda lambd

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

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

omit

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

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

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

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

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

"Frobenius norm"

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

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

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

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

$$\frac{\partial J}{\partial w^{(l)}} = dw^{(l)}$$

"Weight decay"

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

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

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{< 1} \underbrace{w^{(l)}}_{> 0} - \alpha (\text{from backprop})$$



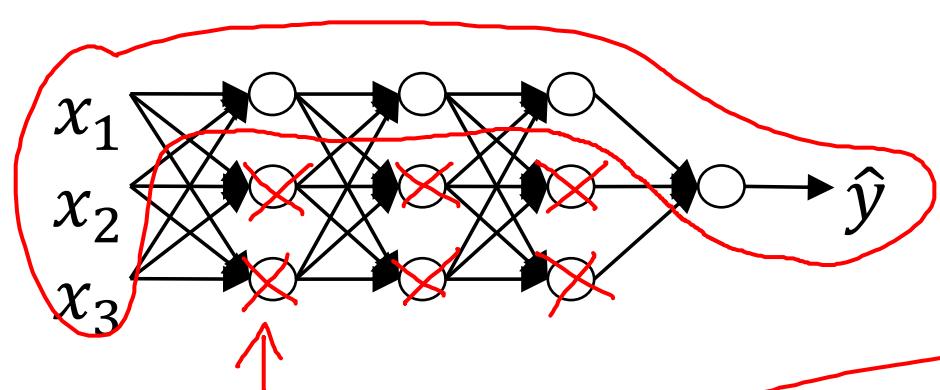
deeplearning.ai

# Regularizing your neural network

---

## Why regularization reduces overfitting

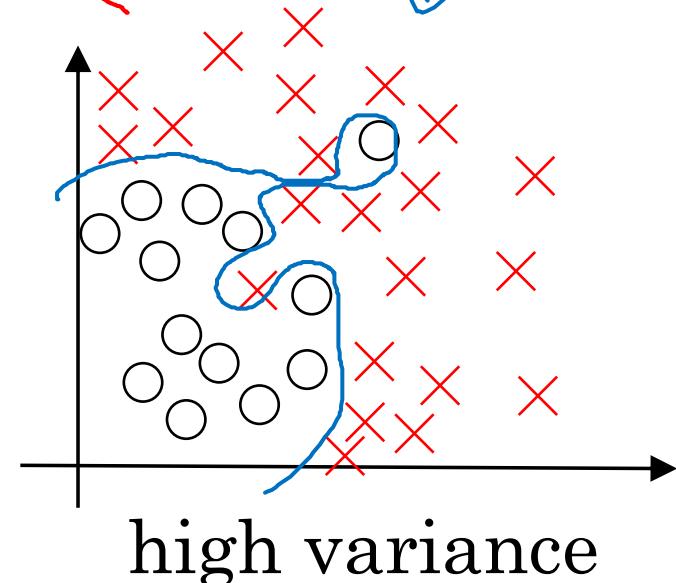
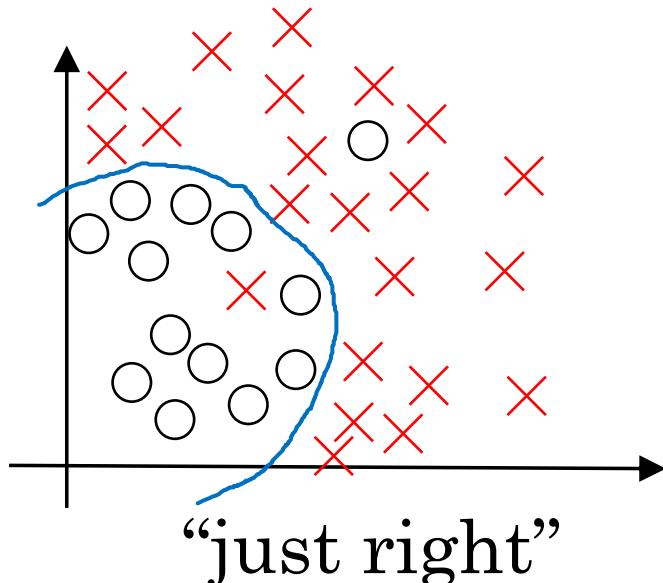
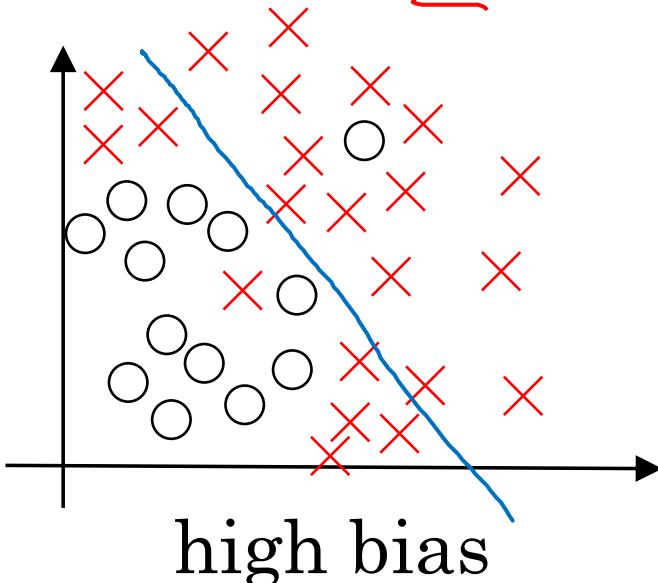
# How does regularization prevent overfitting?



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

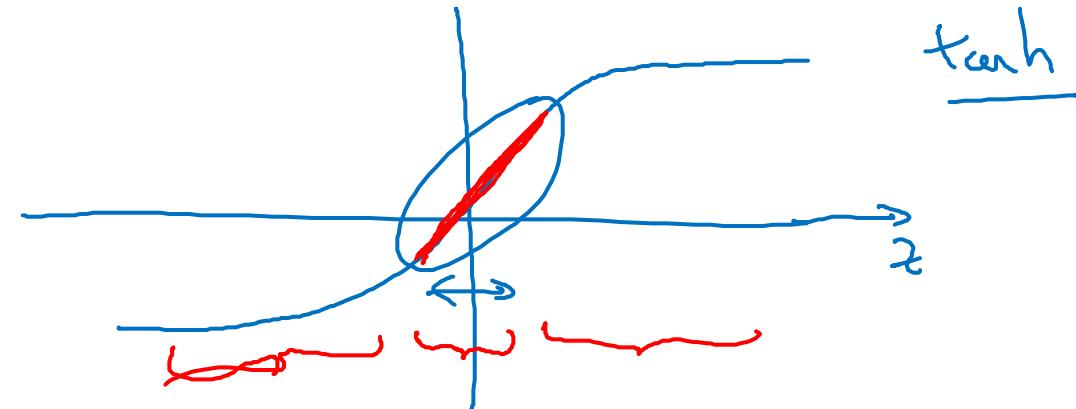
$\lambda \approx 0$

lambda 1 e  
yaklaşınca w 0  
a yaklaşır



# How does regularization prevent overfitting?

w nun küçük değerleri alması linear bir problem çözümümüzü sağlar



tanh

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

$$\lambda \uparrow$$

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

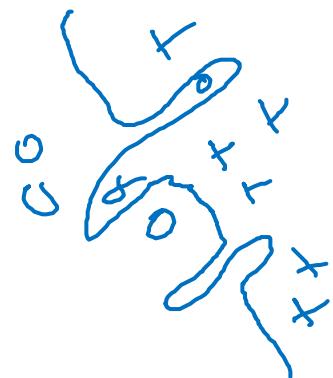
$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

Every layer  $\approx$  linear.

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



weight decay olduğu için cost sürekli azalıyor olmasa bir yerden sonra düz devam ederdi





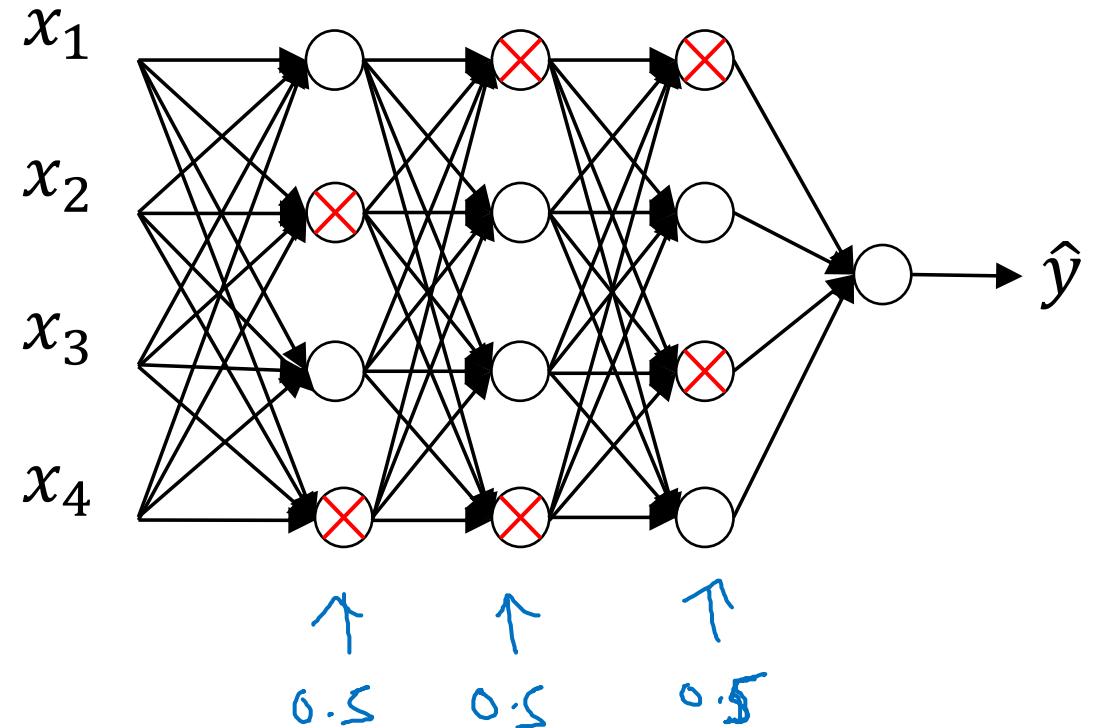
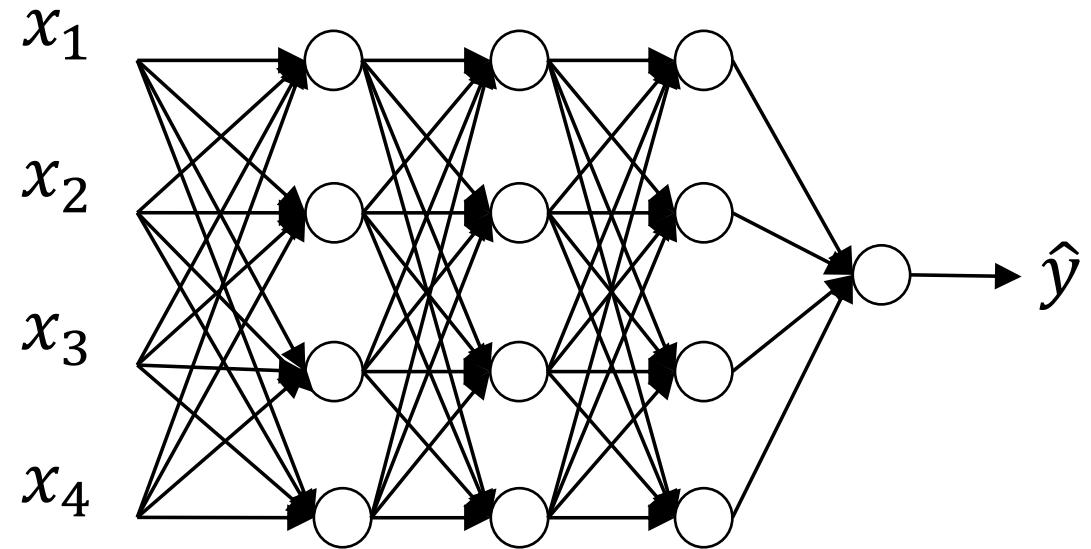
deeplearning.ai

Regularizing your  
neural network

---

Dropout  
regularization

# Dropout regularization



# Implementing dropout (“Inverted dropout”)

Illustrate with layer  $l=3$ .  $\text{keep-prob} = \frac{0.8}{x}$   $\underline{\underline{0.2}}$

$$\rightarrow d_3 = \underbrace{\text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keep-prob}}_{\text{d3}}$$

$$\underbrace{a_3}_{\text{a3}} = \text{np.multiply}(a_3, d_3) \quad \# a_3 * d_3.$$

$$\rightarrow \underbrace{a_3 /=\cancel{0.8} \text{ keep-prob}}_{\text{a3}} \leftarrow$$

50 units.  $\rightsquigarrow$  10 units shut off

$$z^{(4)} = w^{(4)} \cdot \underbrace{\frac{a^{(3)}}{x}}_{\text{reduced by } 20\%} + b^{(4)}$$

Test

$$x = \underline{\underline{0.8}}$$

# Making predictions at test time

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

No drop out.

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

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

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

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

$$\downarrow \hat{y}$$

we dont want randomness in test time that would just add noise to our output

$\lambda$  = keep-prob



deeplearning.ai

Regularizing your  
neural network

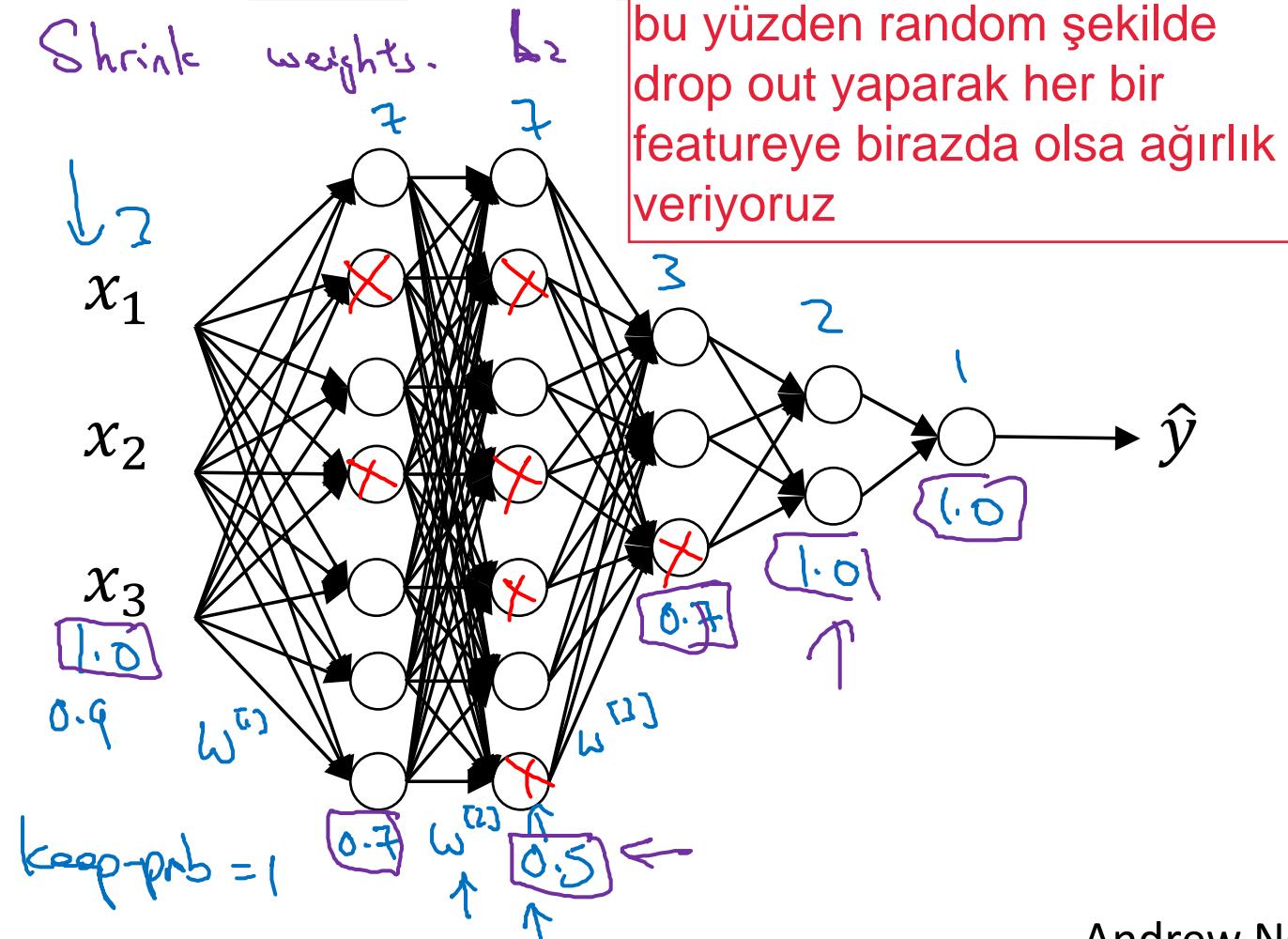
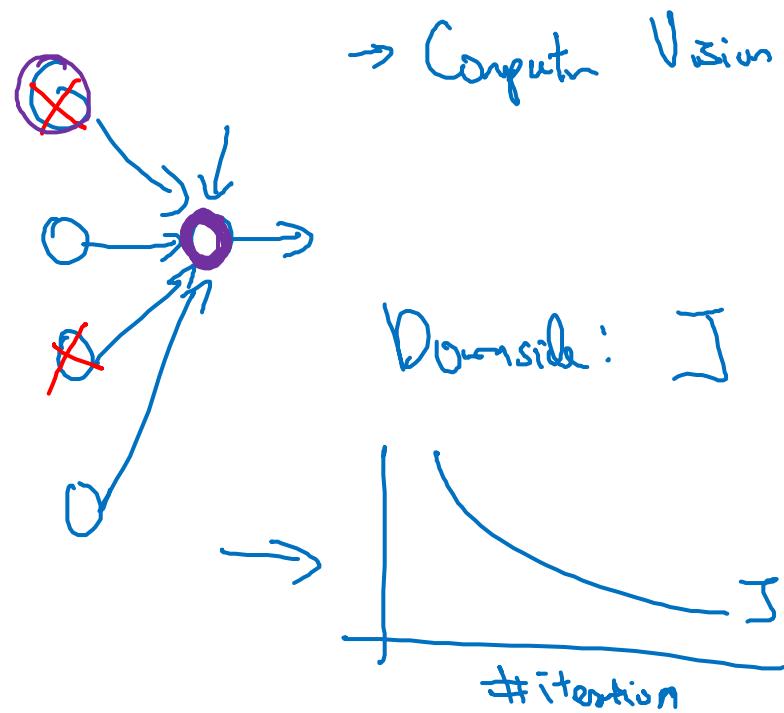
---

Understanding  
dropout

# Why does drop-out work?

mantık şöyle: modelimizdeki weighlerimizin tek bir featureye bağlı olmasını istemiyoruz dağınık olmaları daha mantıklı

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





deeplearning.ai

# Regularizing your neural network

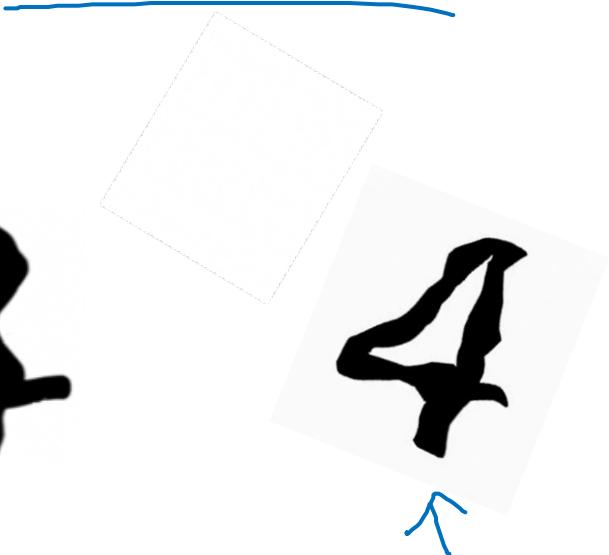
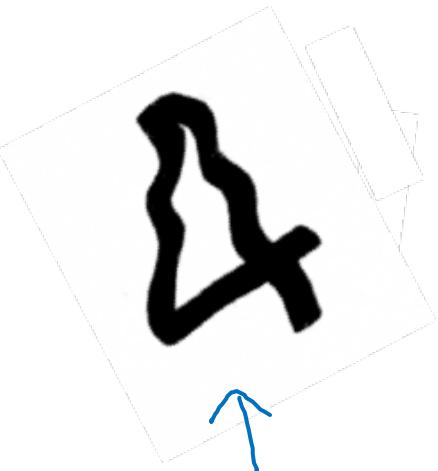
---

## Other regularization methods

# Data augmentation



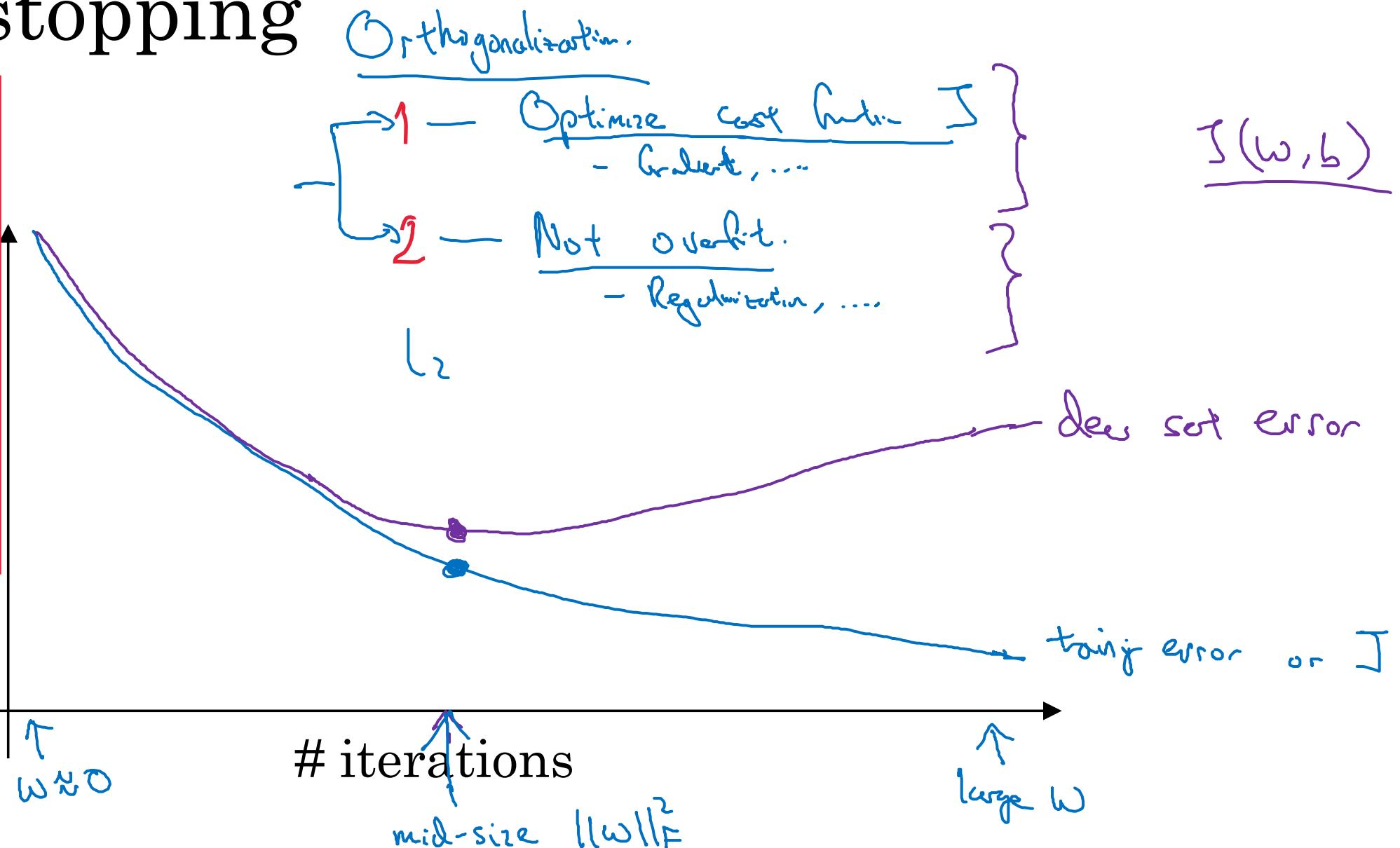
4



# Early stopping

1 ve 2 işlerini aynı anda yapamayız bu sebeple en optimal yerde gradient descend durdukmak early stopping oluyor ama bu cost functionumu daha fazla optimize etmemizi engelliyor.

L2 regularizationunda bu dezavantaj yok ama onunda lambda değerini bulmak maliyetli





deeplearning.ai

Setting up your  
optimization problem

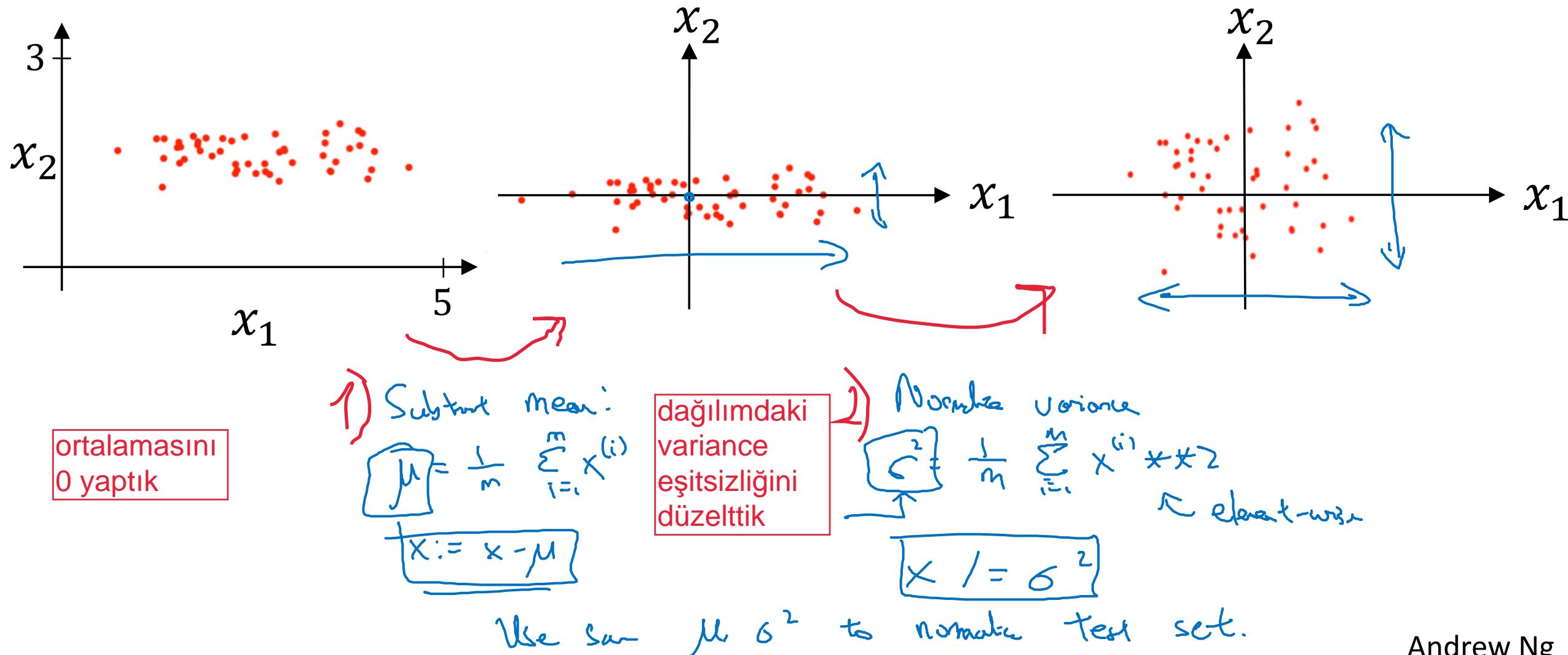
---

Normalizing inputs

# Normalizing training sets

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

normalize derken aynı m ve var değerlerini kullan train ve test için ve bunları train datası üzerinden hesapla

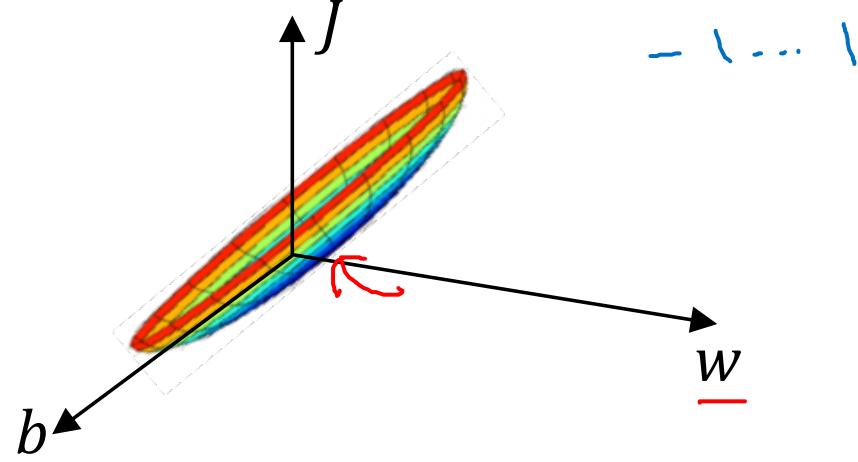


# Why normalize inputs?

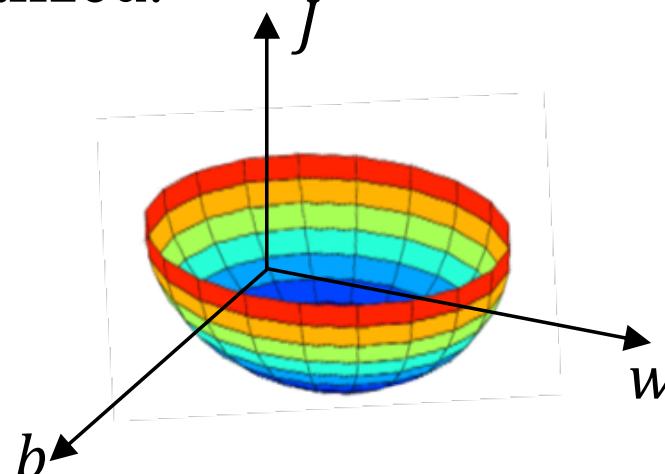
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$w_1$   $x_1: \underline{1 \dots 1000} \leftarrow$   
 $w_2$   $x_2: \underline{0 \dots 1} \leftarrow$

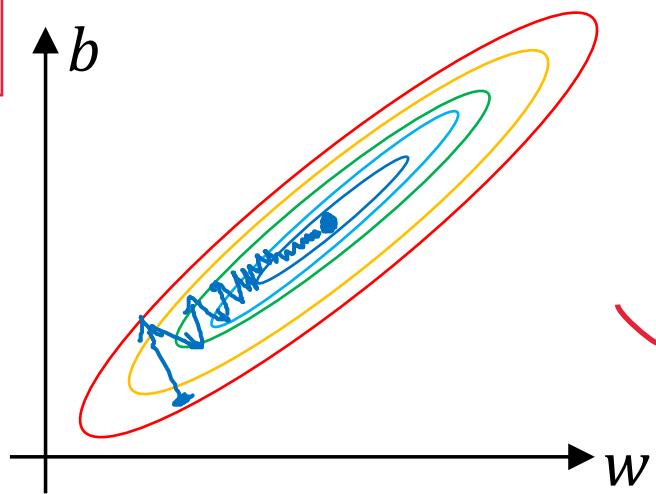
Unnormalized:



Normalized:



weightlerimi  
zin aralıkları  
birbirinden  
çok farklı  
olması  
optimization  
algoritmaları  
için büyük  
bir sorun

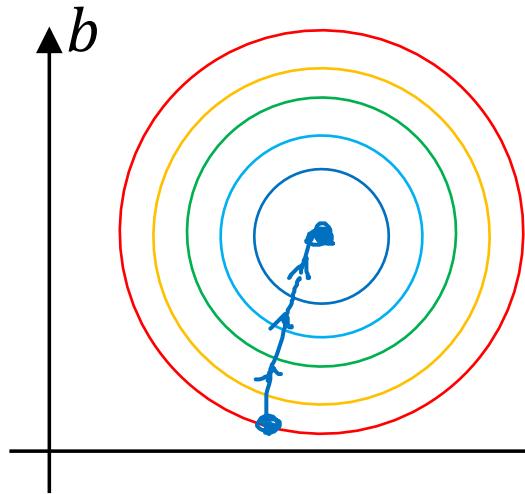


$x_1: 0 \dots 1$

$x_2: -1 \dots 1$

$x_3: 1 \dots 2$

benzer aralıklara  
normalize ettiğimizde  
GD daha hızlı şekilde  
converge olucaktır



$w$  Andrew Ng



deeplearning.ai

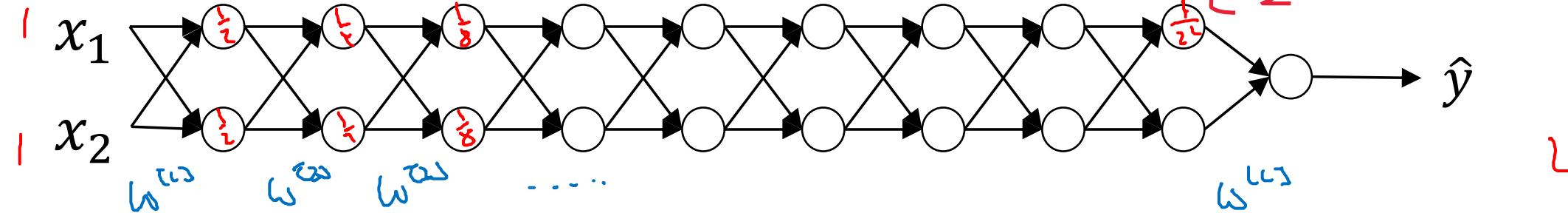
Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

# Vanishing/exploding gradients

$L = 150$



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

$$\hat{y} = w^{[L]} \underbrace{\left( w^{[L-1]} \underbrace{\left( w^{[L-2]} \dots \right)}_{\text{Identity matrix}} \right)}_{\text{Identity matrix}} \times a^{[L-1]}$$

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

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

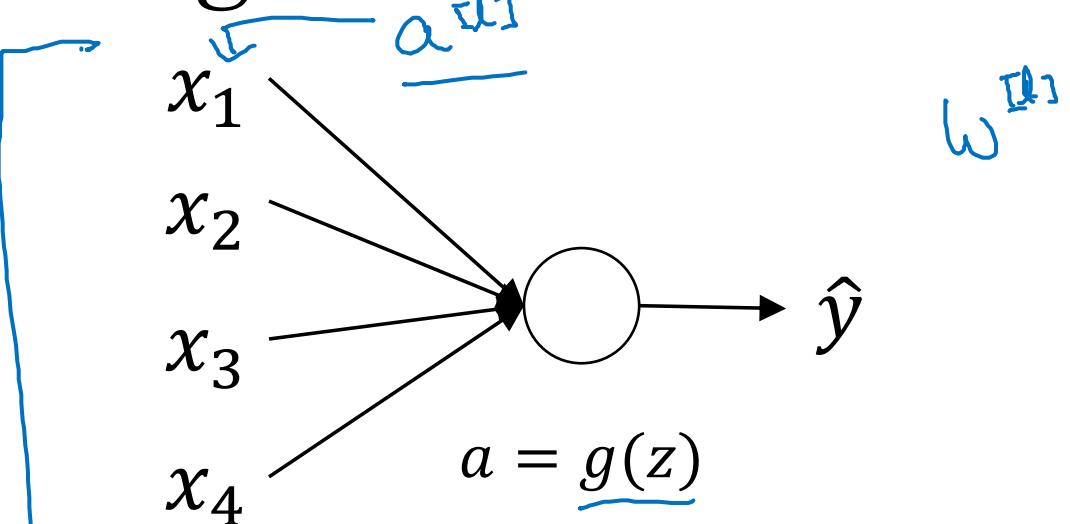
$$w^{[L]} = \begin{bmatrix} 1.5 & 0 & 0 \\ 0 & 0 & 1.5 \\ 0 & 1.5 & 6.5 \end{bmatrix}$$

$$\begin{aligned} z^{[L]} &= w^{[L]} x \\ a^{[L]} &= g(z^{[L]}) = z^{[L]} \\ \hat{y} &= w^{[L]} \underbrace{\begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}}_{I-1} \times \underbrace{\begin{bmatrix} 1.5 & & & \\ & 1.5 & & \\ & & 1.5 & \\ & & & 6.5 \end{bmatrix}}_{I-1} x \end{aligned}$$

eğer  $w$  muz identity matrixten birazcık büyük veya küçük olursa layerlarımızın çok fazla olduğu durumlarda bu katlanarak arttığı için gradient explode veya gradient vanishing durumlarıyla karşılaşırız.

w lerimiz her layerde aynı

# Single neuron example



x lerimizin sayısı  
artıkça w larımızın  
değerlerinin  
azalmasını istiyoruz  
ki exploding/  
vanishing problemi  
ortadan tam  
olmasada kalksın

$$z = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$$

Large  $n \rightarrow$  Smaller  $w_i$

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

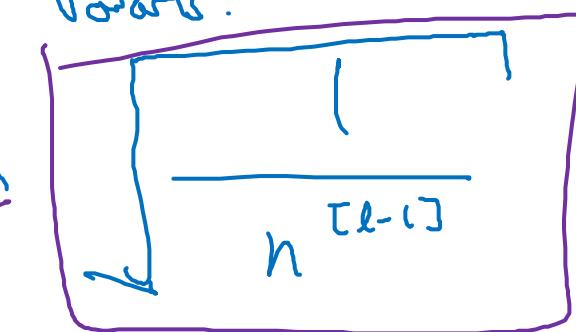
$$w^{[l]} = \text{np.random.rand}(\text{shape}) * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

ReLU

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

Other variants:

Tanh



Xavier initialization

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

- So lets say when we initialize w's like this (better to use with tanh activation):

```
np.random.rand(shape) * np.sqrt(1/n[1-1])
```

or variation of this (Bengio et al.):

```
np.random.rand(shape) * np.sqrt(2/(n[1-1] + n[1]))
```

- Setting initialization part inside sqrt to  $2/n[1-1]$  for ReLU is better:

```
np.random.rand(shape) * np.sqrt(2/n[1-1])
```



deeplearning.ai

# Setting up your optimization problem

---

## Numerical approximation of gradients

# Checking your derivative computation

$$f(\theta) = \theta^3$$

$\theta \in \mathbb{R}$

I

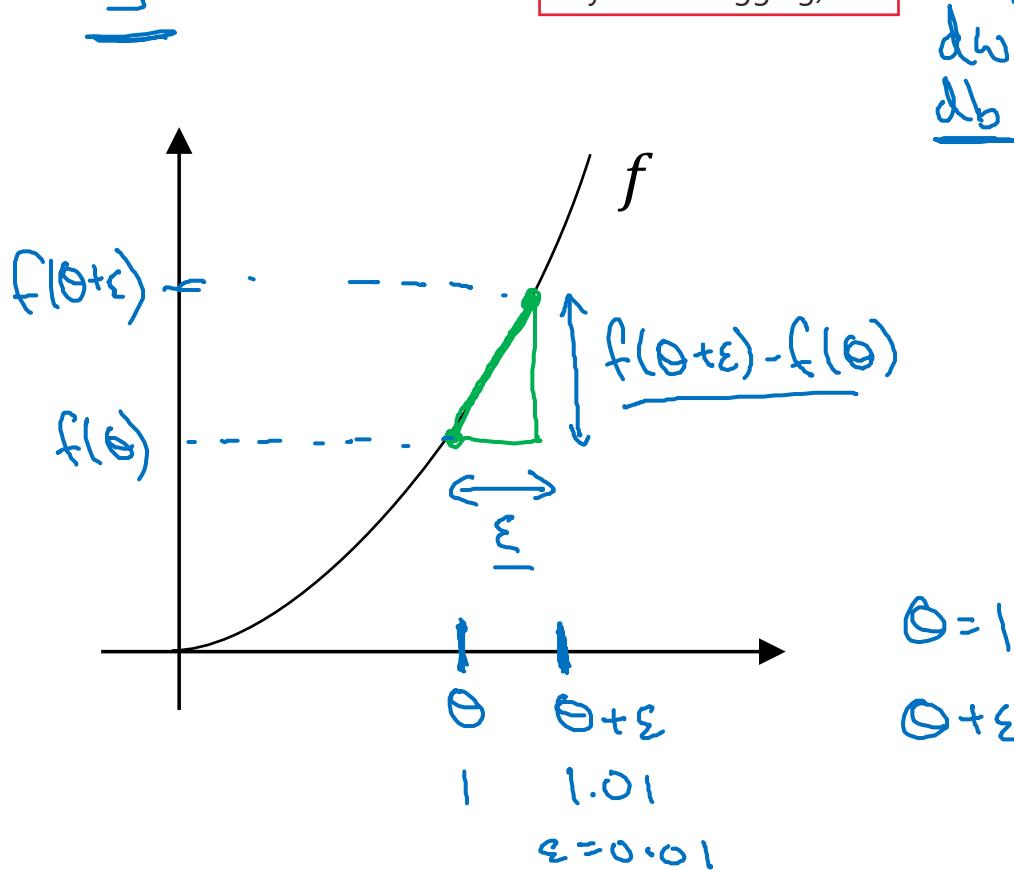
gradient checking which tells you if your implementation of backpropagation is correct.

but it's slower than gradient descent (so use only for debugging).

$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

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

$$g(\theta) = 3 - (1)^2 = 3 \text{ when } \theta = 1$$



$$\frac{(1.01)^3 - 1^3}{0.01} = 3.0301$$

$\approx g(\theta)$

$\approx 3$

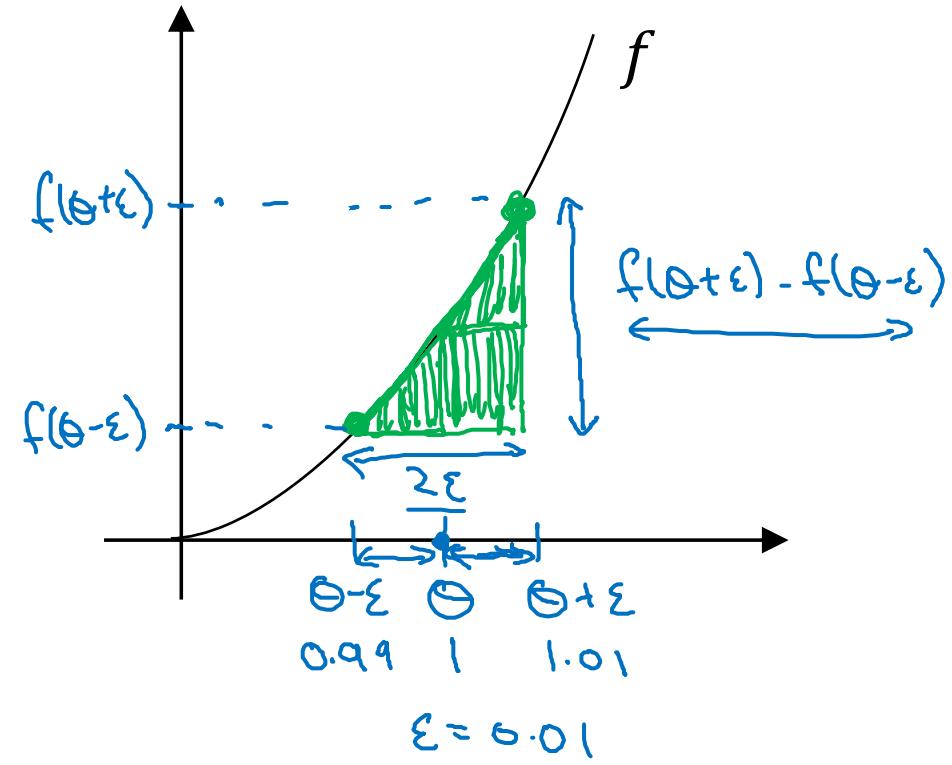
$0.0301$

$3.1$

$3.2$

# 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  $i$ :

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

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

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

Check

öklid

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

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

$\approx$

$10^{-7}$  - great!

$10^{-5}$

$\rightarrow 10^{-3}$  - worry. debug



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking  
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

debug için kullan değerler birbirine çok yakınsa kapatıp devam et değilse hata ara

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

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

hatayı birbirinden çok farklı olan değerlerde aramamız lazım bunlarla Q yi veren W ve B lere bakarak yapacağız

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

- 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.

$$J \quad \underline{\text{keep\_prob} = 1.0}$$

drop out için kullanıcksak  $\text{keep\_prob} = 1.0$  olarak kontrol et sonra drop out uygula

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

$$\underline{w, b \text{ NO}}$$