



deeplearning.ai

Setting up your  
ML application

---

Train/dev/test  
sets

structured  
data

advertising  
search  
security  
logistics

# Applied ML is a highly iterative process

# layers

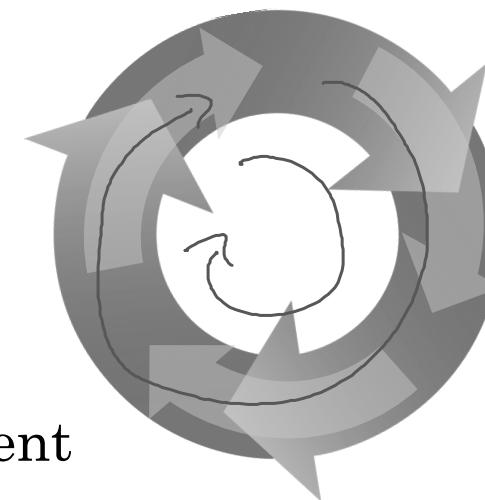
# hidden units

learning rates

activation functions

...

Idea



Experiment

Code

NLP, Vision, Speech, Structural data

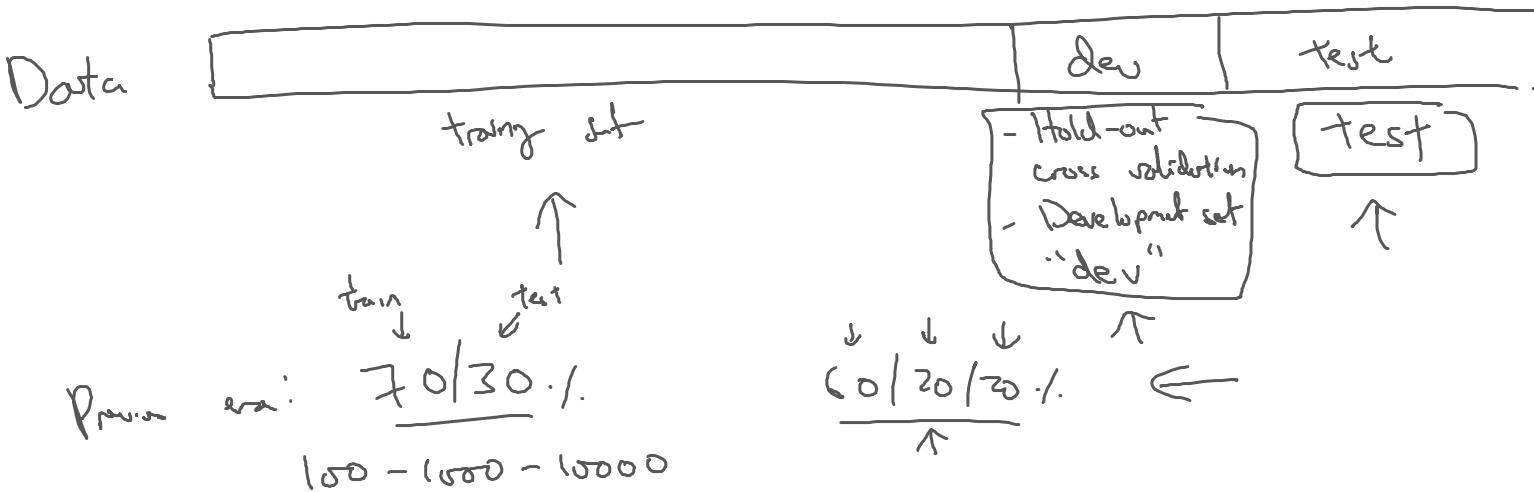
Ads      Search      Security      logistic ...

Data : ① Training ② Development ③ Testing

Dev set = hold-out = cross validation = development  $\rightsquigarrow$

بيان مفهوم عملية جندين للعريقة جامع است.

## Train/dev/test sets



Big data: 1000,000      10,000      10,000

98 / 1 / 1 .%

99.5 { .25 / .25 { - 1 .%

جائز: مربع صالح dev/test دلالة solution  $\Leftarrow$  ملائمة dev/Test والـ Train ليسا Mismatched Distributions

# Mismatched train/test distribution

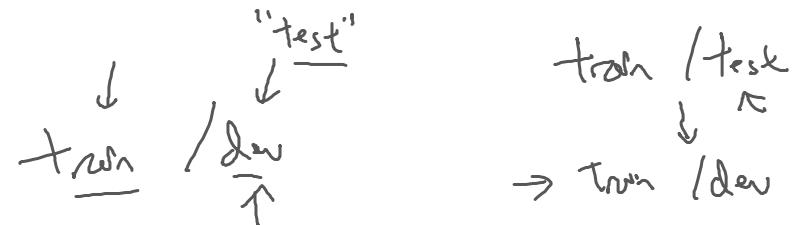
Certs

Training set:  
Cat pictures from  
webpages



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

→ Make sure derived test come from some distributions.



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

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

(میتوان) dev کردن Test ، ابتدا dev set میباشد (چون ما شنیده بیلی go) unbiased estimate logar

← وقتی سوچو گوئی Train/Test می‌شود، Train/Dev منظور



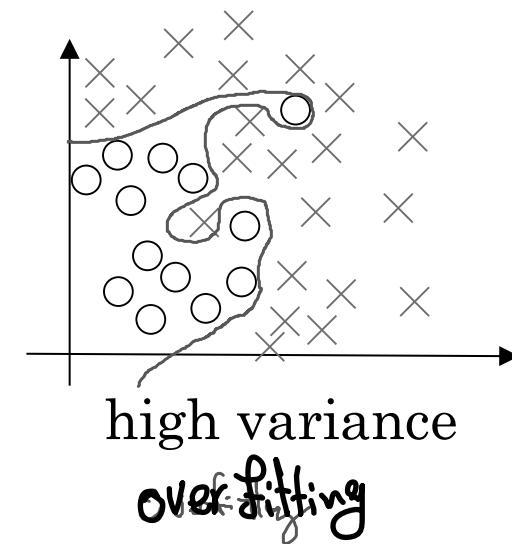
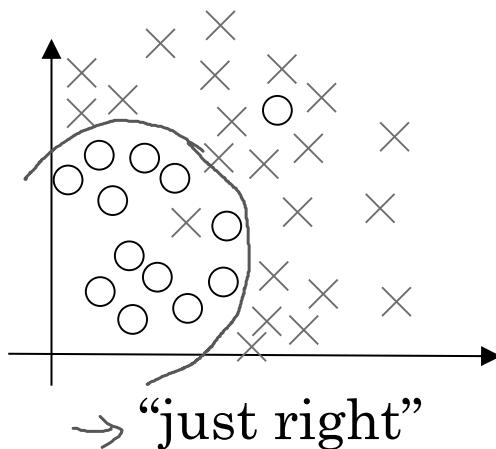
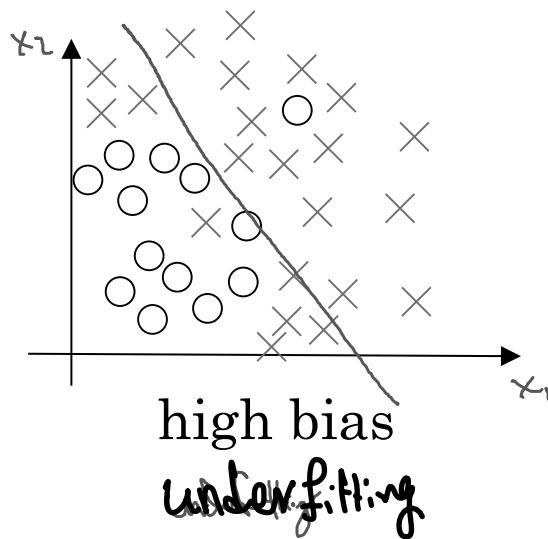
deeplearning.ai

Setting up your  
ML application

---

Bias/Variance

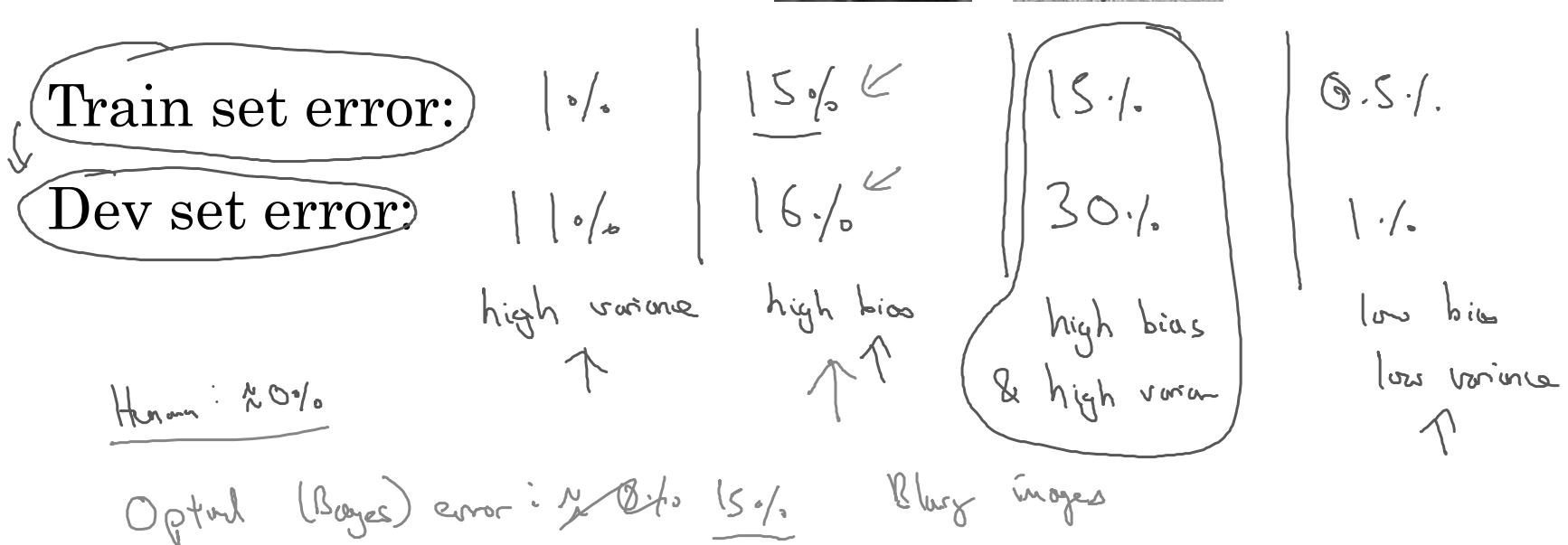
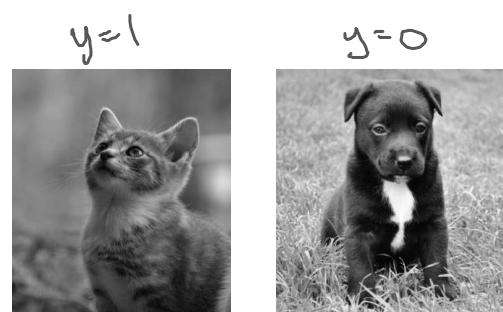
# Bias and Variance



آریک انسان بخواهد بین سه/کدام classification کند سهاید جا خطا ۱۵٪ است.  $\Rightarrow$  این عدد معیاری برای قضاوت بالایا پاکیزه optimal (Bayes) Error خامیت خطا Train, Dev, Train خطا داشت.

# Bias and Variance

## Cat classification

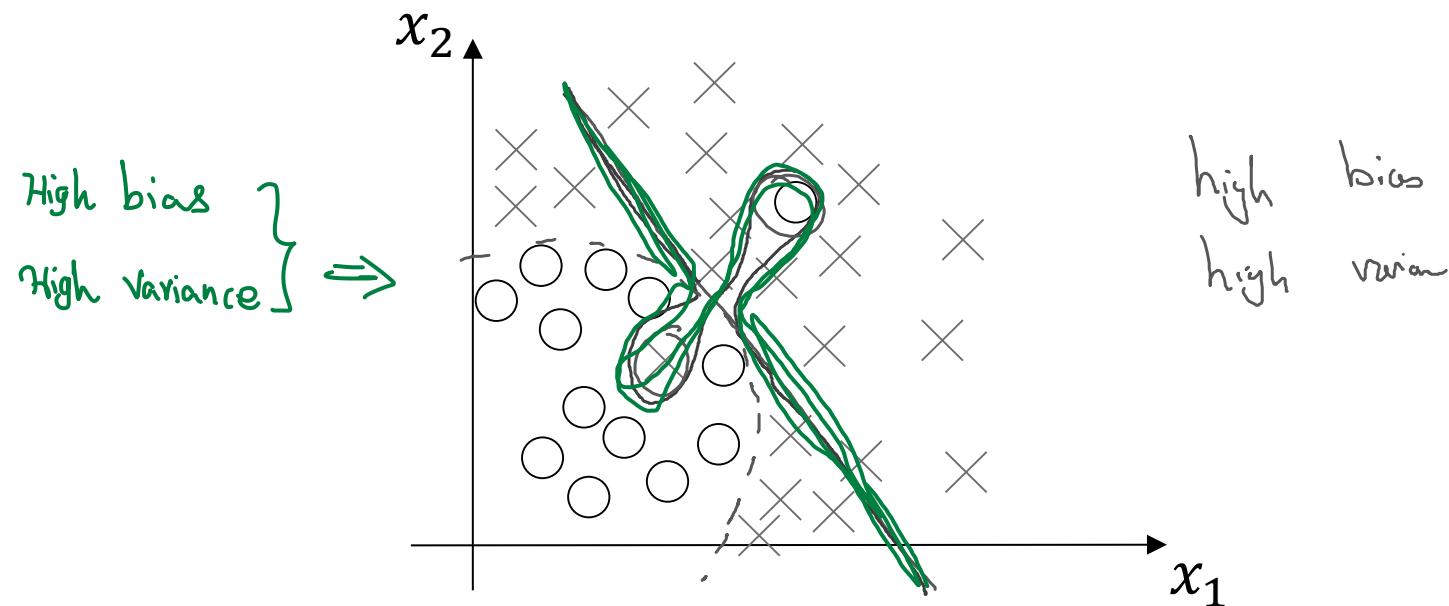


بهترین حالت زمانی است که هر دو minimal variance, Bias و maximum bias باشند.

High Bias: Training error  $\gg$  Bayes error  $\Rightarrow$  underfitting  
 High Variance: Dev error  $\gg$  Training error  $\Rightarrow$  overfitting

Can we have under/over fitting simultaneously?

# High bias and high variance





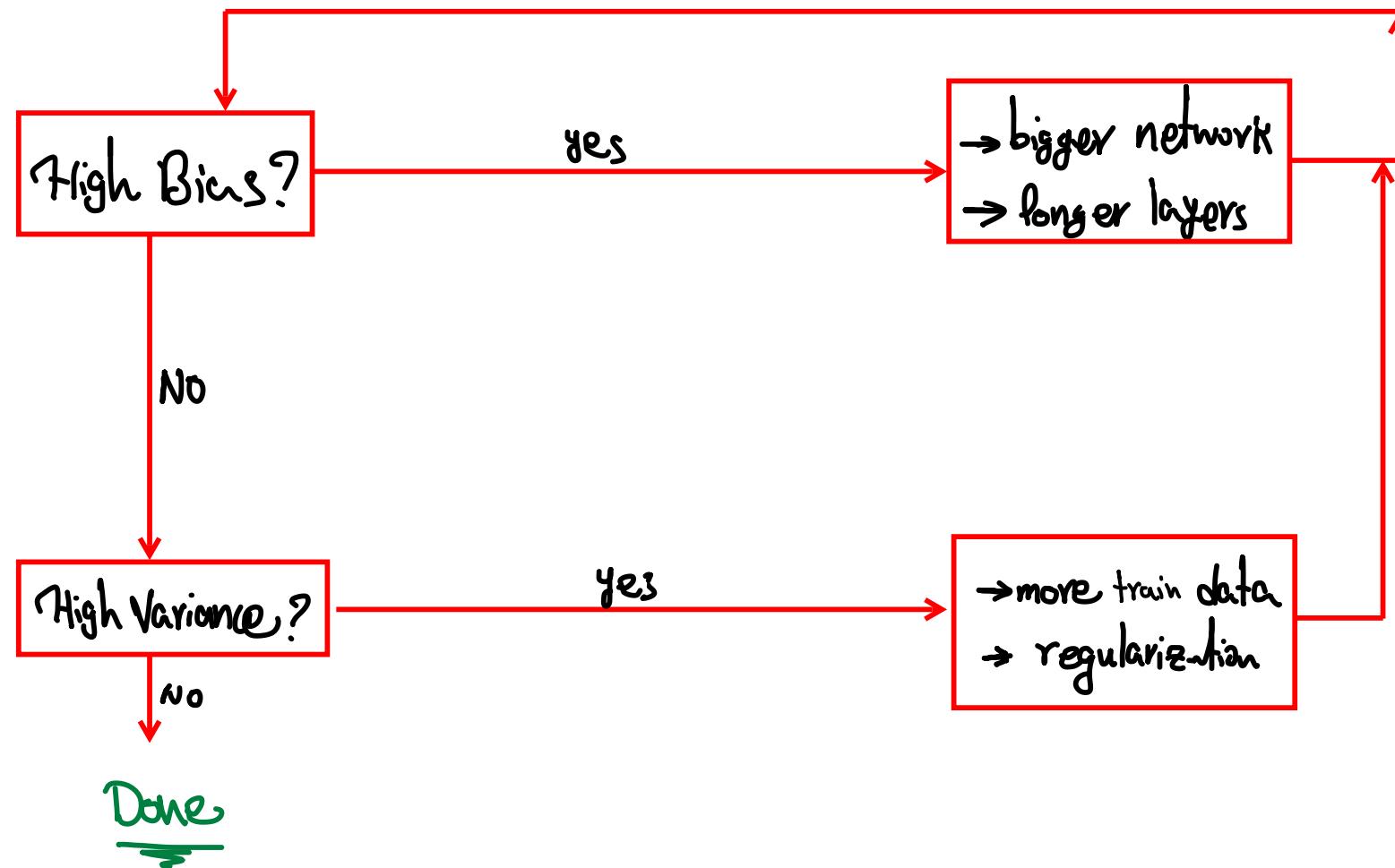
deeplearning.ai

Setting up your  
ML application

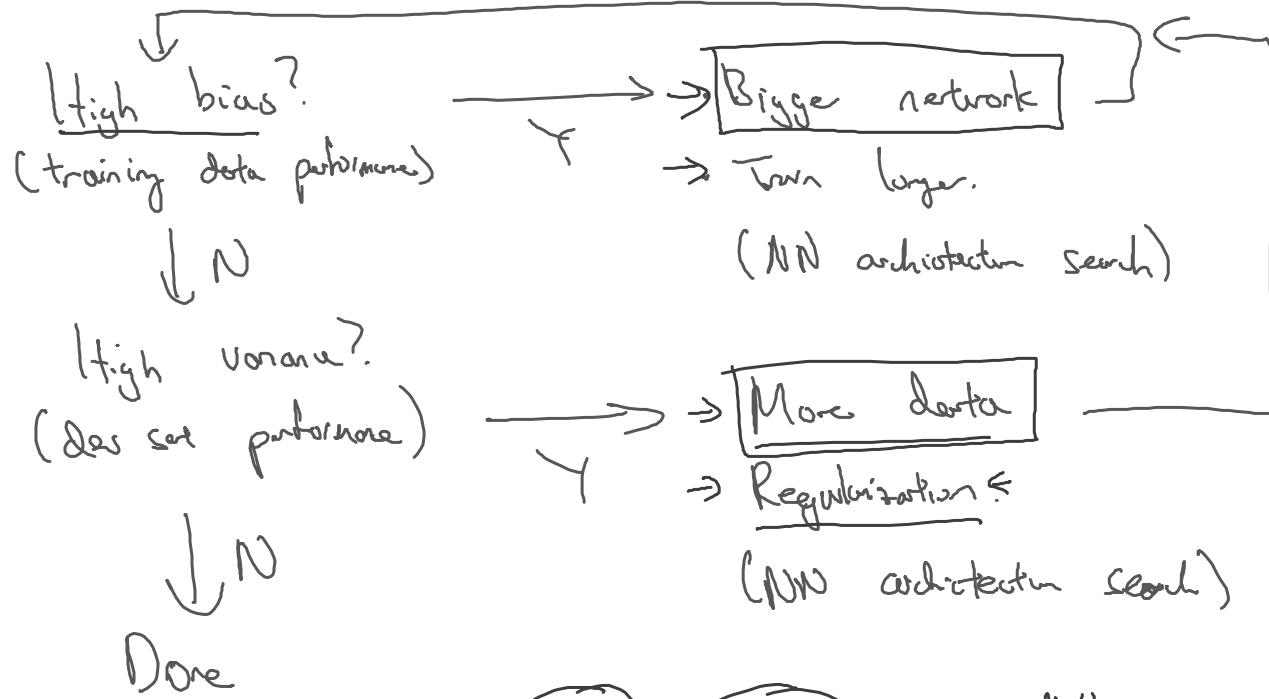
---

Basic “recipe”  
for machine learning

# Basic “recipe” for machine learning



# Basic recipe for machine learning





deeplearning.ai

Regularizing your  
neural network

---

Regularization

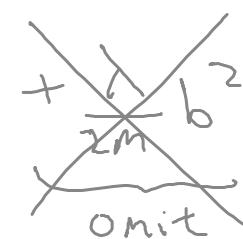
# Logistic regression

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

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

$\lambda$  = regularization parameter  
 lambda  
lambd

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



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

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

w will be sparse

L<sub>2</sub> Regularization  $\frac{\lambda}{2m} \|w\|_2^2$  (weight decay)

For vectors

L<sub>1</sub> Regularization  $\frac{\lambda}{2m} \|w\|_1 \Rightarrow$  sparse w (more zero components)

$$J = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_{l=1}^L \|W^{[l]}\|_F^2$$

Frobenius norm (matrix)  $\|W\|_F^{[l]} = \sqrt{\sum_{i=1}^n \sum_{j=1}^{n^{[l]}} (W_{i,j}^{[l]})^2}$  where  $W^{[l]}: (n^{[l]}, n^{[l-1]})$

## Neural network

$$\rightarrow J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(y^{(i)}, \hat{y}^{(i)})}_{\text{"Frobenius norm"}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2}_{w^{[l]}: (n^{[l]}, n^{[l-1]})}$$

$$\|w^{[l]}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^{n^{[l]}} (w_{i,j}^{[l]})^2$$

"Frobenius norm"  $\| \cdot \|_2^2$   $\| \cdot \|_F^2$

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

$$\frac{\partial J}{\partial w^{[l]}} = dW^{[l]}$$

$$\rightarrow w^{[l]} := w^{[l]} - \alpha \frac{\partial J}{\partial w^{[l]}}$$

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

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

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{< 1} \underbrace{w^{[l]}}_{\text{without regularization}} - \alpha (\text{from backprop})$$

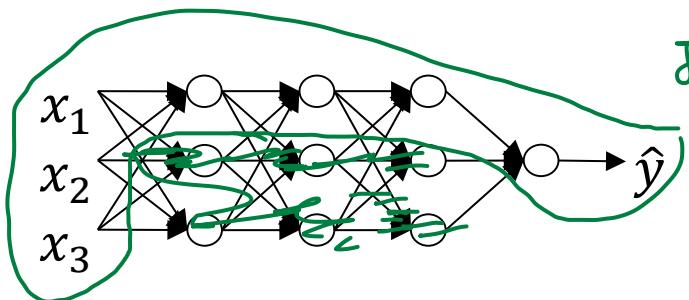
$L_2$  Regularization = "Weight decay"

$$w^{[l]} = \left(1 - \frac{\alpha \lambda}{m}\right) w^{[l]} - \alpha (\text{from backprop})$$

with

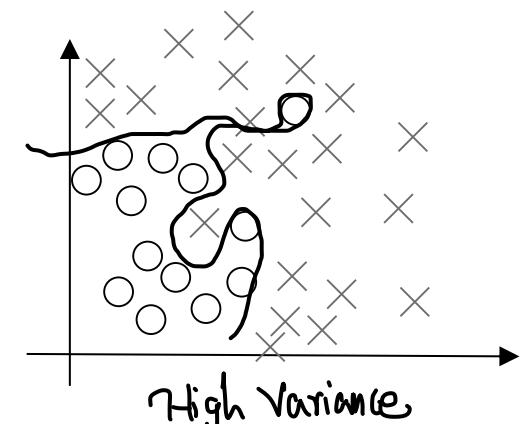
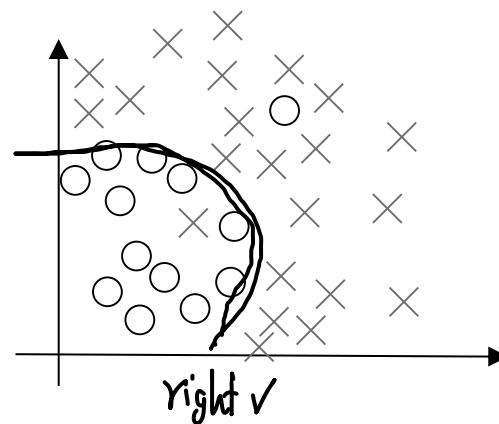
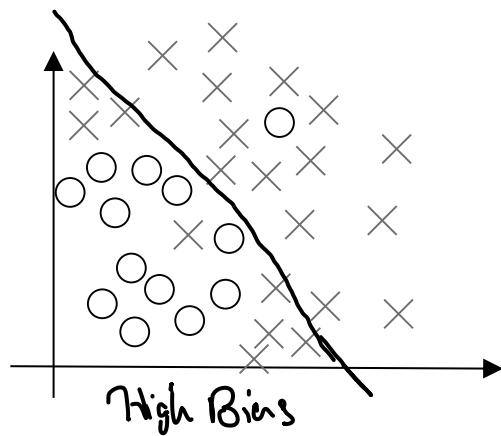
$$w^{[l]} = w^{[l]} - \alpha (\text{from backprop}) \rightsquigarrow \text{without regularization}$$

# How does regularization prevent overfitting?

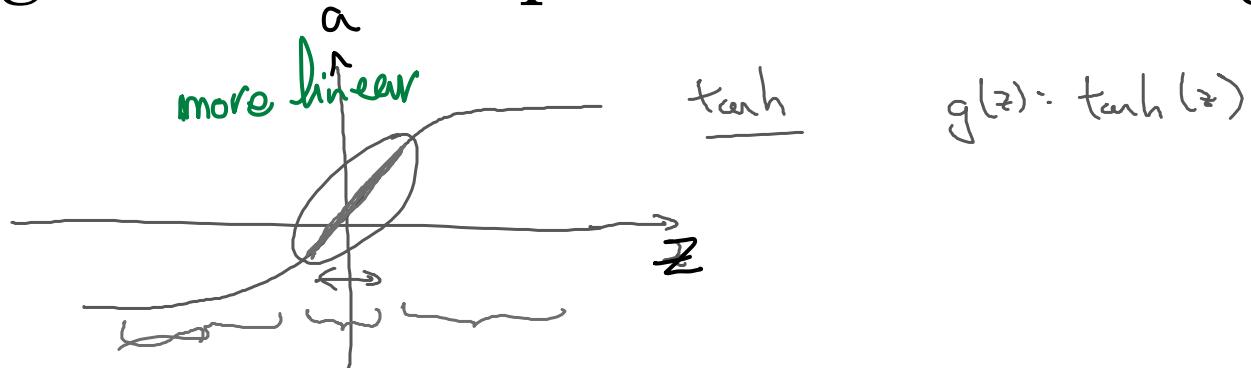


$$\mathcal{J}(\mathbf{w}^{[l]}, \mathbf{b}^{[l]}) = \frac{1}{m} \sum_i \ell(\hat{y}_i, y_i) + \frac{\lambda}{2m} \sum_{j=0}^L \|\mathbf{w}^{[l]}\|_F^2$$

If  $\lambda \uparrow \Rightarrow \mathcal{J}_{\min}$  leads to  $\mathbf{w}^{[l]} \approx 0 \Rightarrow$  simpler model  $\Rightarrow$  overfitting



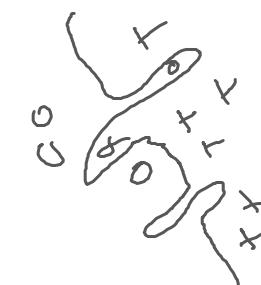
# How does regularization prevent overfitting?



$$\lambda \uparrow$$

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

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



Every layer  $\approx$  linear.

---

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





deeplearning.ai

Regularizing your  
neural network

---

Dropout  
regularization

# Dropout regularization

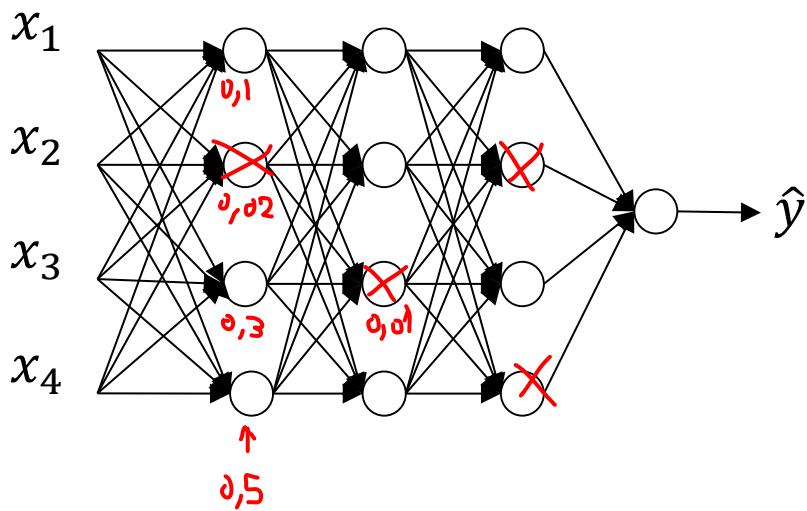
صلهای جریه Performance  $\Leftarrow$  گنجینه حذف می‌کنیم.

keep\_prob

اول طبقه جانشی overfitting

بصورت تصادفی هر hidden node  $\uparrow$  احتمال صحتی، کسریها را حذف می‌کنیم.

لوب: بازیگردانی iteration را انجام می‌دهیم



$\uparrow$   
0.5  
 $\uparrow$   
0.5  
 $\uparrow$   
0.5

Andrew Ng

معروف ترین نوع از regularization technique: Invert dropout

## Implementing dropout ("Inverted dropout")

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

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

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

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

50 units.  $\rightsquigarrow$  10 units shut off

$$z^{[4]} = w^{[4]} \cdot \underline{\frac{a^{[3]}}{5}} + b^{[4]}$$

$\cancel{5}$  reduced by 20%.

Test

$$1 = \underline{0.8}$$

1.  $d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{Keep_Prob} \Rightarrow$  جزوی  $0, 2, 0, 8$  کهای  $F, T$  باشد  $d_3$
2.  $a_3 = \text{np.multiply}(a_3, d_3)$
3.  $a_3 = a_3 / \text{keep-prob}$

• مهم انجام dropout  $\Rightarrow$  Test set  $\rightarrow$

## Making predictions at test time

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

No drop out.

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

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

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

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



$\lambda$  = keep-prob

Andrew Ng



deeplearning.ai

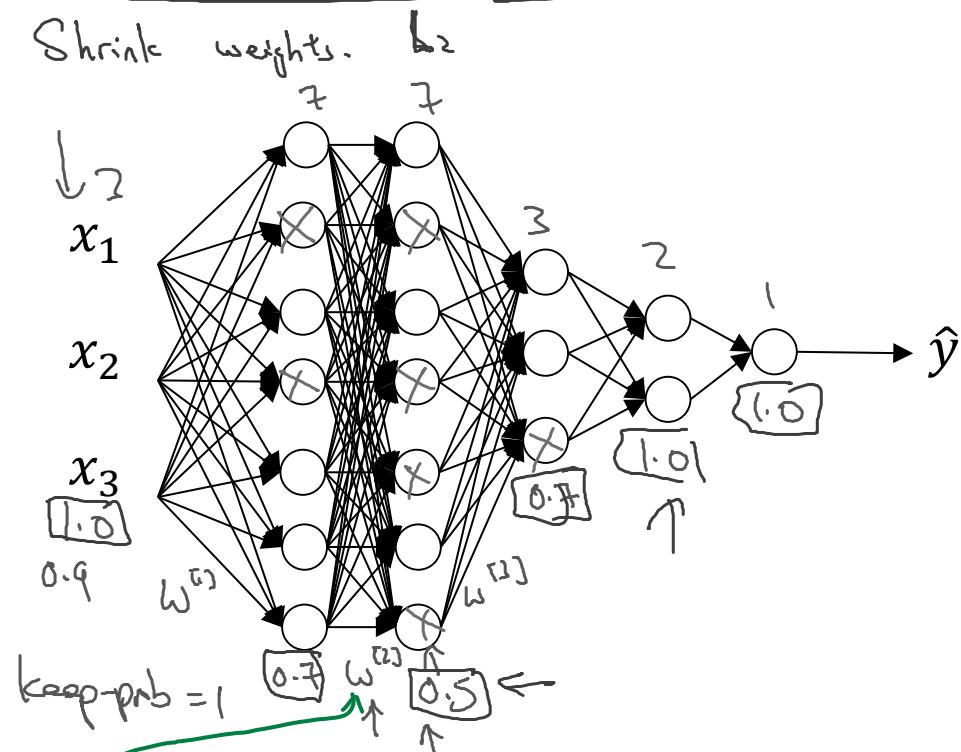
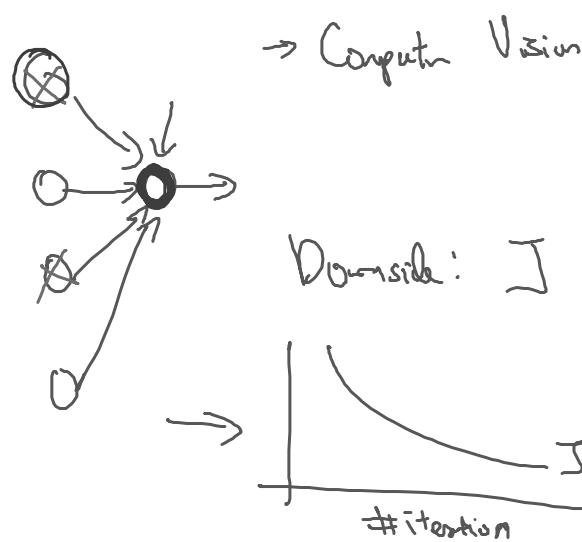
Regularizing your  
neural network

---

Understanding  
dropout

Why does drop-out work?  $\Rightarrow$  Similar to L<sub>2</sub> & shrinking the weights!

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



برای هر لایه صورت  $Keep\text{-}prob$  خاص خود را تعریف کرد  $\Rightarrow$  مخصوصاً برای لایه‌هایی که بزرگ‌شوند حتی صورت بخوبی انجام نمایند.

الن من در dropout این است که cost function  $\rightarrow$  ادکننده‌اند و این مترادف باشد.

Regularization: ①  $L_2$  ② dropout ③ data augmentation ④ early stopping



deeplearning.ai

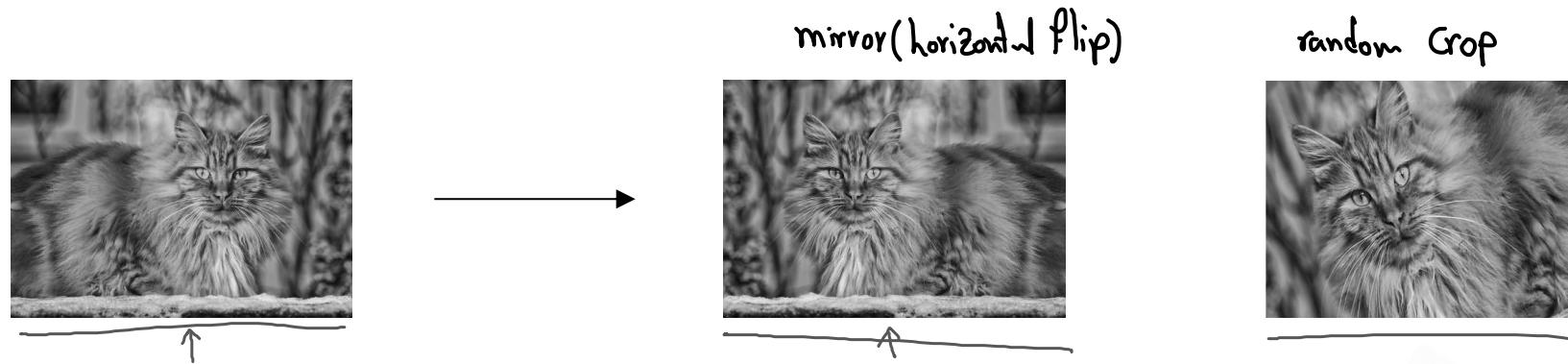
# Regularizing your neural network

---

## Other regularization methods

بهای جمع کنیم  $\leftarrow$  augmentation more training data  $\rightarrow$  regularization

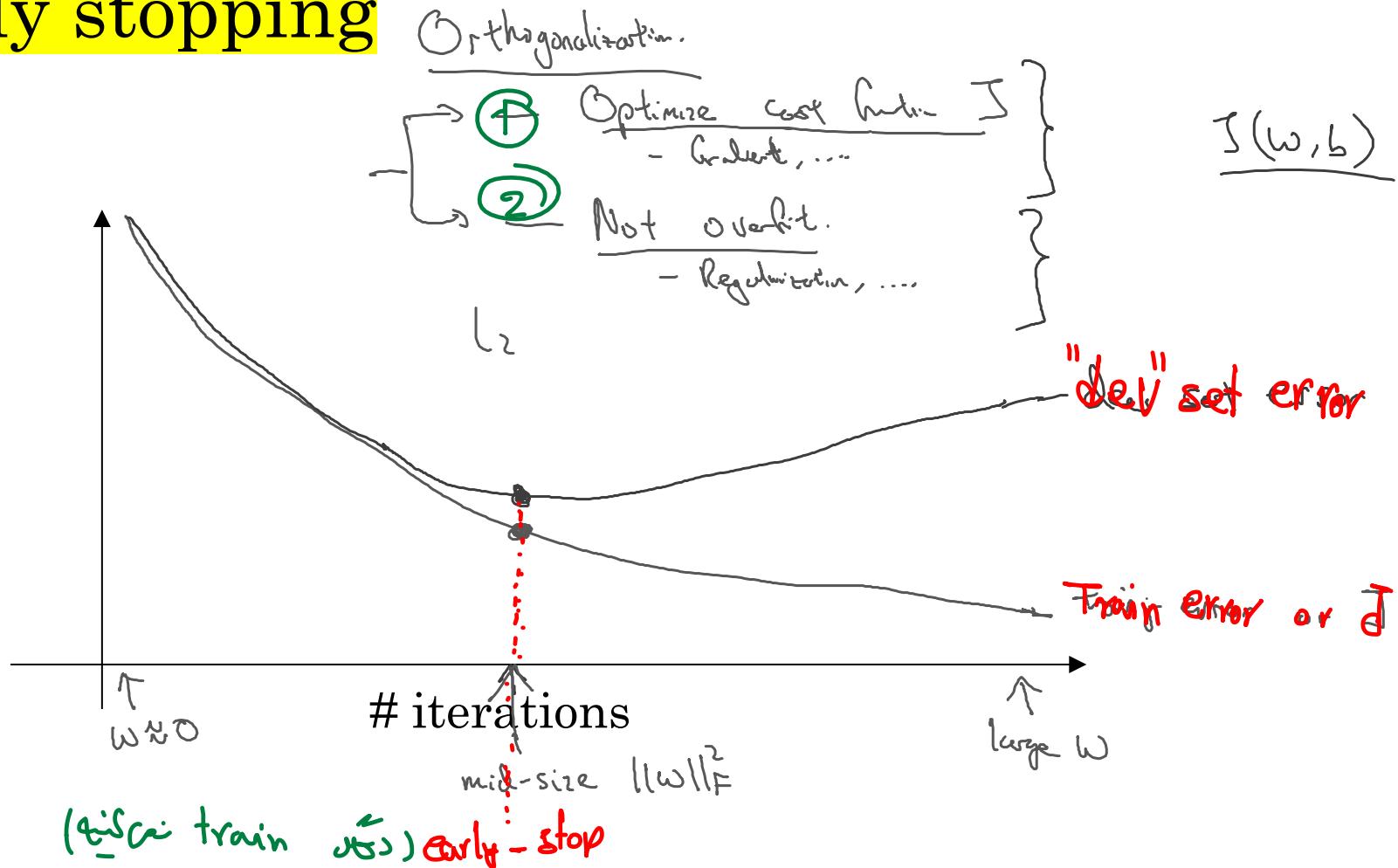
## Data augmentation



orthogonalization & to independently do the two tasks.

• if one task, 1, orthogonalization: early stopping  $\leftarrow$

## Early stopping





deeplearning.ai

Setting up your  
optimization problem

---

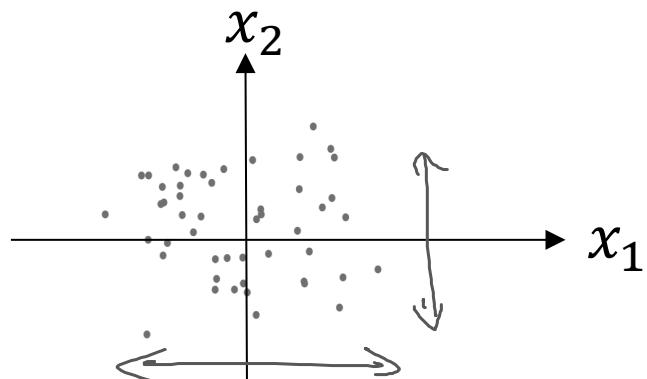
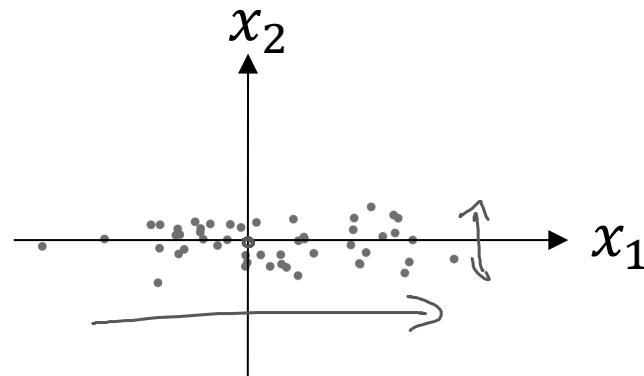
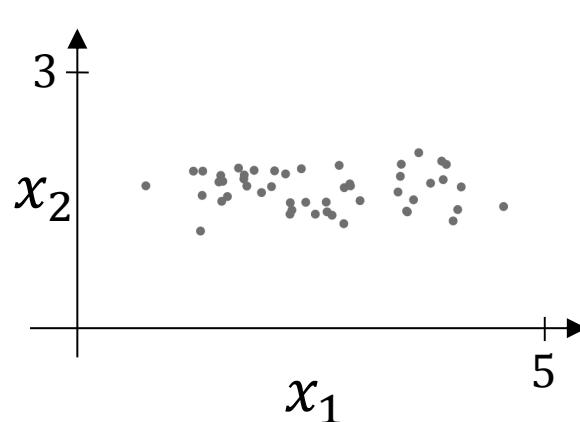
Normalizing inputs

همیانگی روابط رانی داده ها صفر دارند

دل جاید مدل هم باشد و جداگانه بناشد  
برای رادکهای normalization

## Normalizing training sets $x = \frac{x - \mu}{\sigma}$

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



Subtract mean:

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

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

Normalize variance

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

$\sigma$  element-wise

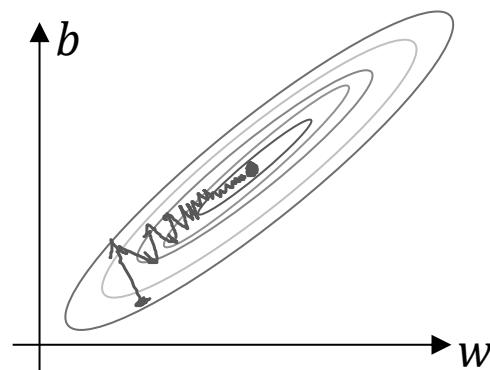
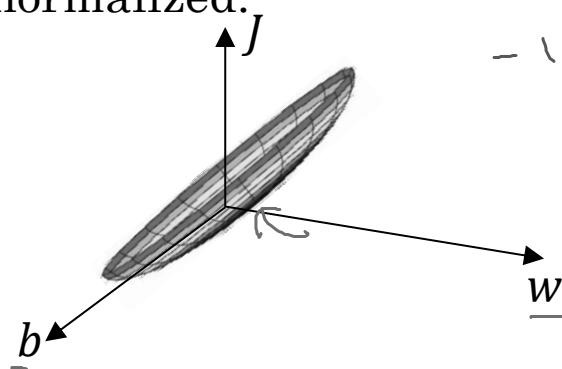
$$\underline{x / \sigma}$$

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

Andrew Ng

# Why normalize inputs?

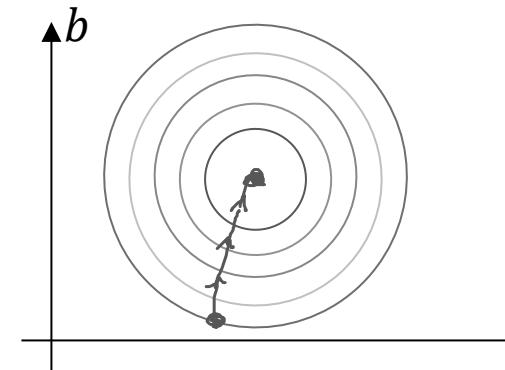
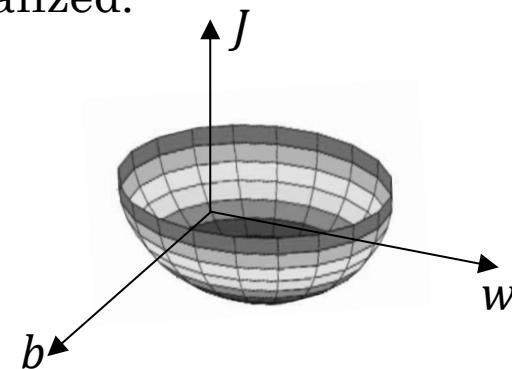
Unnormalized:  
 $w_1$   $x_1$ : 1...1000 ←  
 $w_2$   $x_2$ : 0...1 ←  
- 1...1



$x_1$ : 0...1  
 $x_2$ : -1...1  
 $x_3$ : 1...2

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

Normalized:



Andrew Ng



# Setting up your optimization problem

## Vanishing/exploding gradients

deeplearning.ai

$\frac{\partial J}{\partial w} \rightarrow 0$  or  $\rightarrow \infty$  in deep NNs  $\Leftrightarrow (w^{[l]} < 1 \text{ or } > 1) \text{ or } w^{[l]} \text{ کمتر یا بیشتر از یک جاست}$ . آن مقدار اولیه  $w^{[l]}$  هوتمند انتخاب کرد.  $\leftarrow$   $\text{initial values}$   $\leftarrow$   $\text{جایزه}$ .

$$\text{Relu: } \sqrt{\frac{2}{n^{[l-1]}}}$$

$$\text{Tanh: } \frac{1}{\sqrt{n^{[l-1]}}}$$

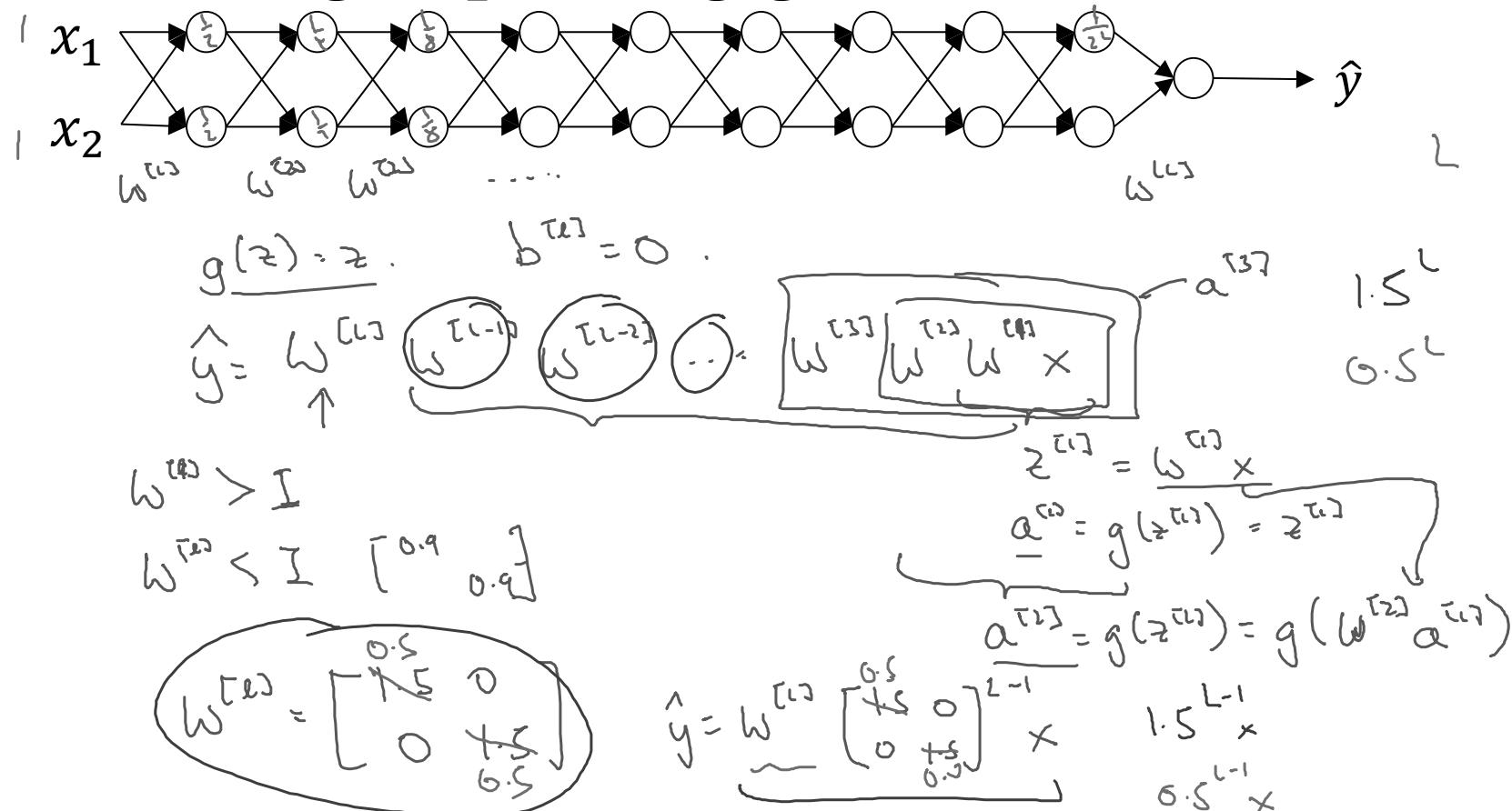
$$\text{Or: } \sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

آن جانشی سود کر  $\leftarrow \text{Var}(w^{[l]})$  بعورت  $\frac{2}{n}$  سود  $\leftarrow$   $w^{[l]}$  خیلی کمتر یا بیشتر از یک نیست.

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

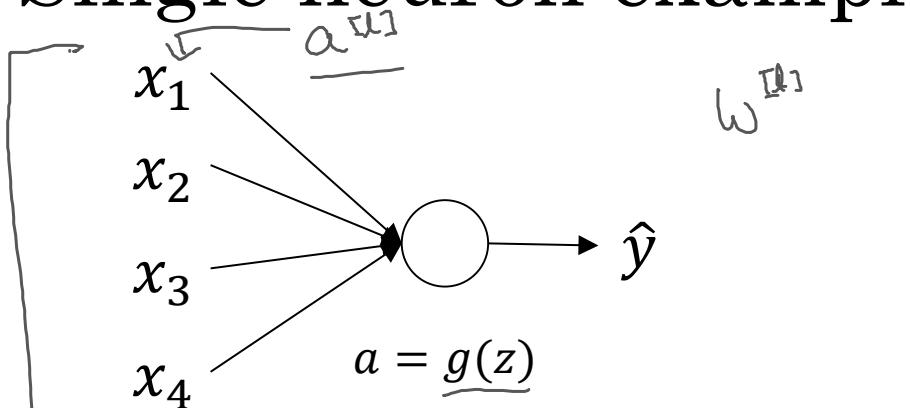
# Vanishing/exploding gradients

$L = 150$



Andrew Ng

# Single neuron example

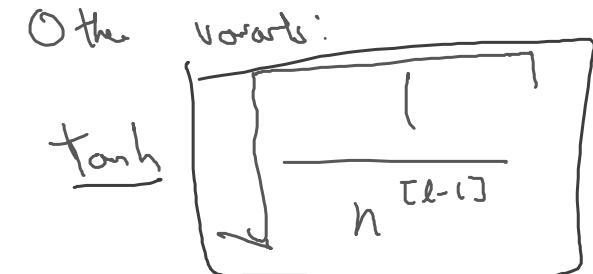


$$z = \underline{w_1}x_1 + \underline{w_2}x_2 + \cdots + \underline{w_n}x_n \quad \cancel{\text{for}}$$

Large  $n \rightarrow$  Smaller  $w_i$

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

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



Xavier initialization ↑

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

↑

Andrew Ng



deeplearning.ai

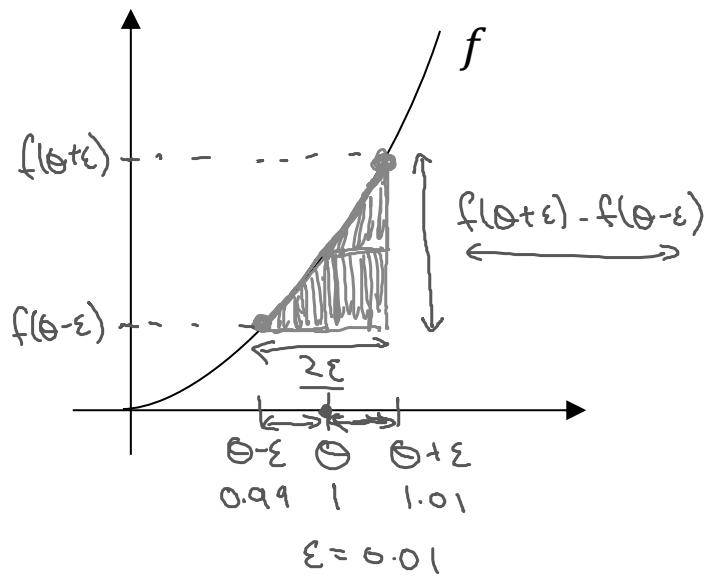
Setting up your  
optimization problem

---

Numerical approximation  
of gradients

# Checking your derivative computation

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



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

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

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

approx error:  $0.0001$

(prev slide:  $3.0301$ , error:  $0.03$ )

$$\left\{ f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \right. \quad \begin{array}{l} O(\epsilon^2) \\ 0.01 \\ \underline{0.0001} \end{array} \quad \left| \quad \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} \quad \text{error: } O(\epsilon) \right. \quad \begin{array}{l} 0.01 \end{array}$$

Andrew Ng



deeplearning.ai

Setting up your  
optimization problem

---

**Gradient Checking**

## Gradcheck

# Gradient check for a neural network

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

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

- ② Take  $\underline{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}}$  and reshape into a big vector  $\underline{d\theta}$ .

concatenate

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

# Gradient checking (Grad check)

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

③ for each  $i$ :

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

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{\text{approx}} \stackrel{?}{\approx} d\theta$$

④ Check

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

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

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

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



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.

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

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

- Doesn't work with dropout.

$$J \quad \underline{\text{keep-prob} = 1.0}$$

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

$$\underline{w, b \approx 0}$$

Andrew Ng