



deeplearning.ai

Setting up your  
ML application

---

Train/dev/test  
sets

Applied ML is a highly iterative process

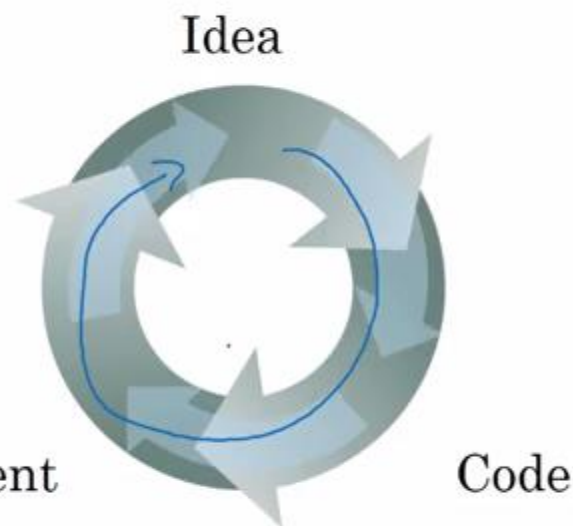
# layers

# hidden units

learning rates

activation functions

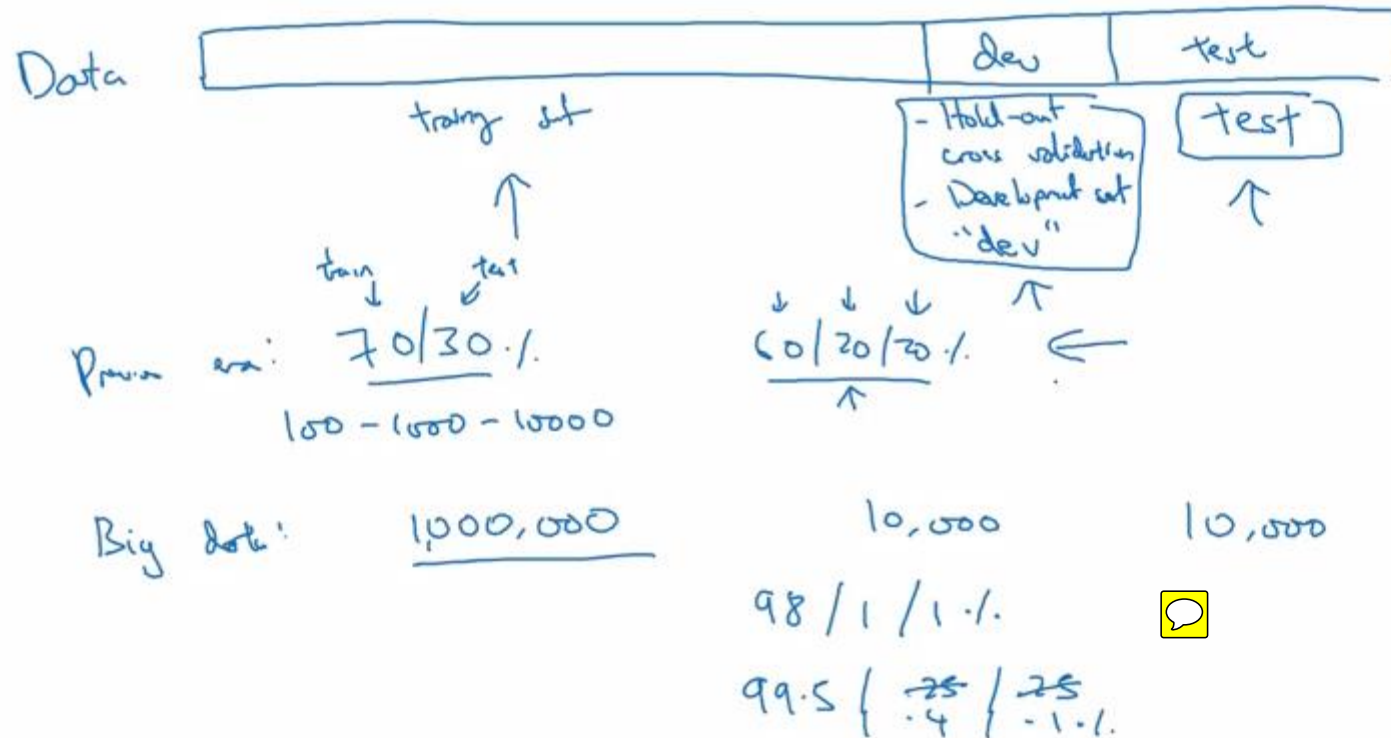
...



NLP, Vision, Speech, Structured Data

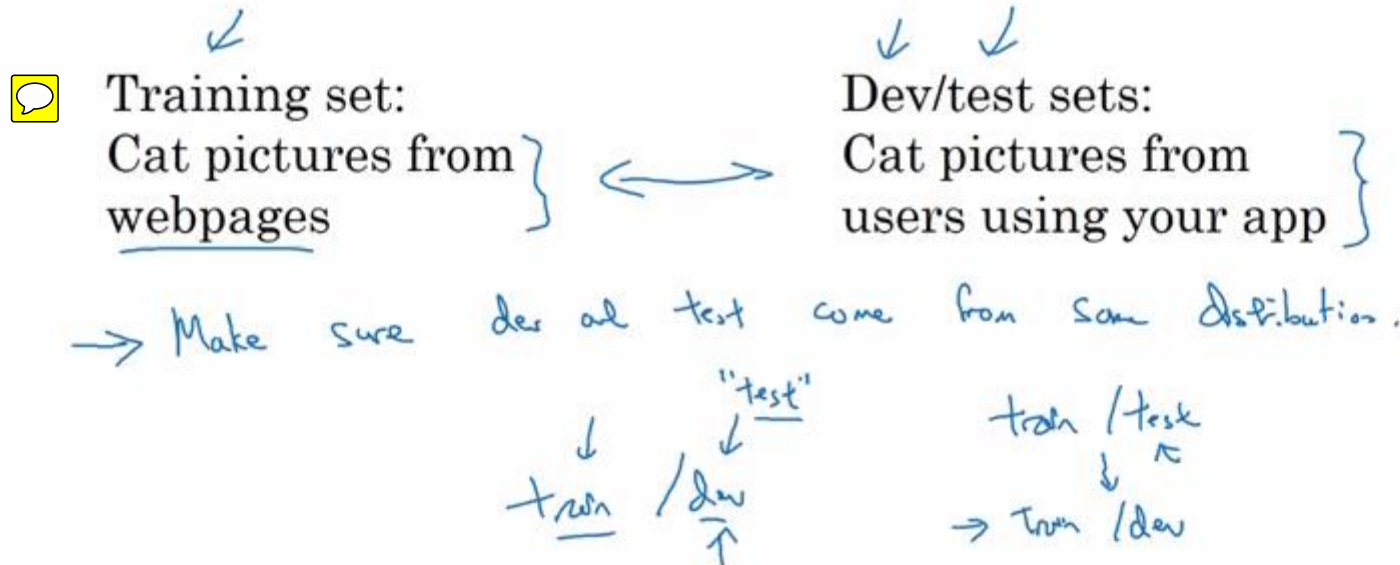
└─┬─┐  
Ads Search Security Logistic ...

# Train/dev/test sets



# Mismatched train/test distribution

Certs



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



deeplearning.ai

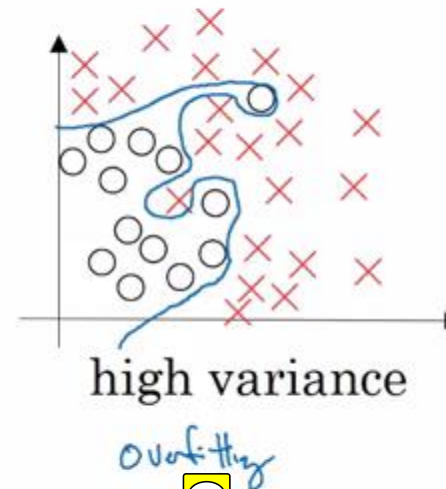
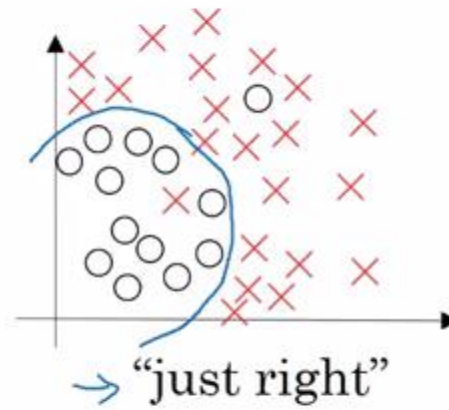
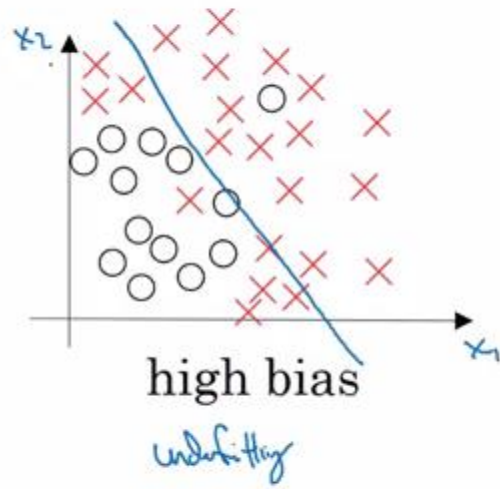
Setting up your  
ML application

---

**Bias/Variance**

---

## Bias and Variance



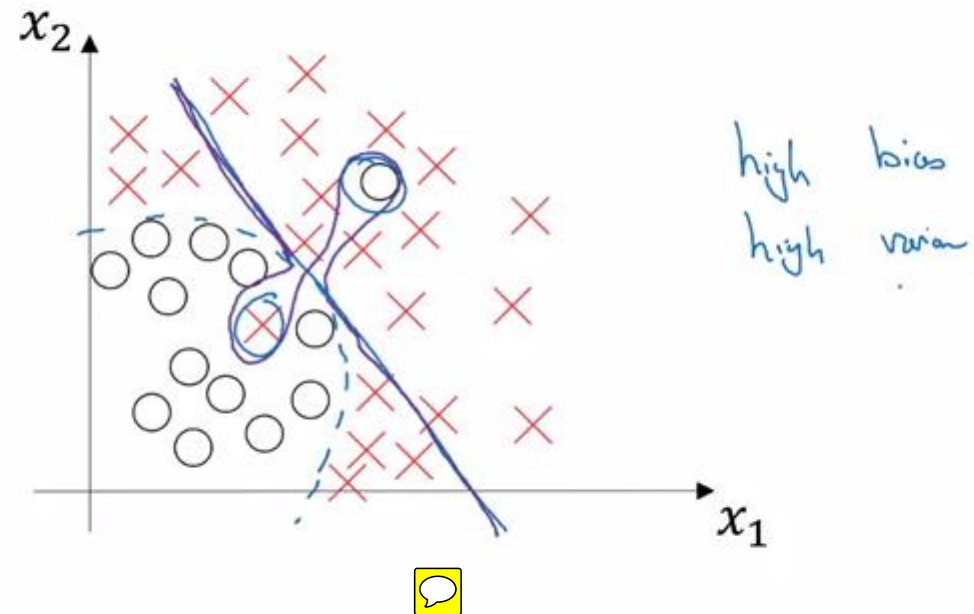
# Bias and Variance

Cat classification



Train set error:	1%	15% ↙	15%	0.5%
Dev set error:	11%	16% ↙	30%	1%
	high variance ↑	high bias ↑↑	high bias & high variance	low bias low variance ↑
Human: ~0%				
Optimal (Bayes) error: ~0%	15%	Blurry images		

# High bias and high variance







deeplearning.ai

# Setting up your ML application

---

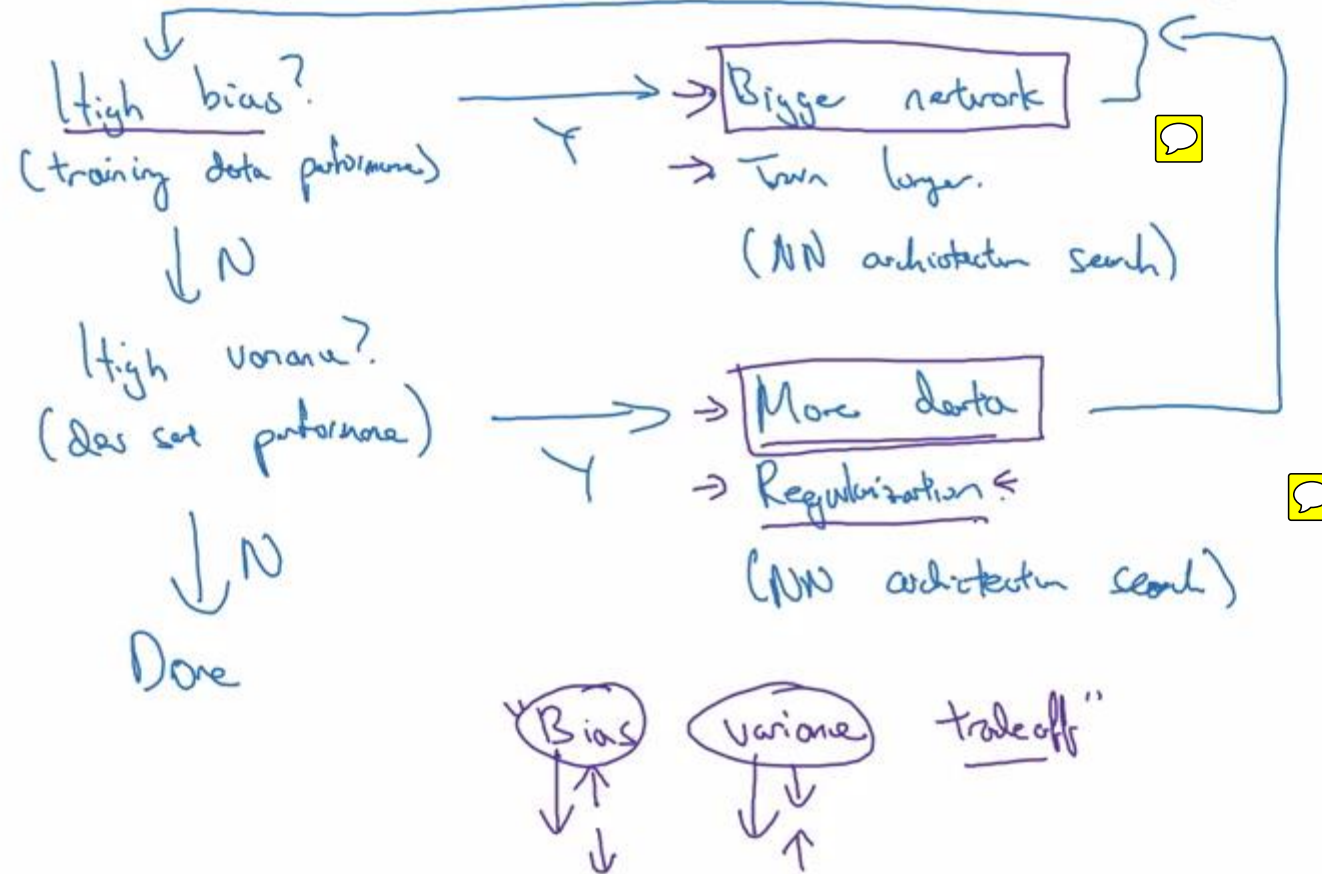
## Basic “recipe” for machine learning



0:21 / 6:21



# Basic recipe for machine learning





deeplearning.ai

Regularizing your  
neural network

---

**Regularization**

# Logistic regression

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

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

$\lambda$  = regularization parameter  
lambda lambda

$$\mathcal{J}(w,b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})}_{\text{cost function}} + \frac{\lambda}{2m} \underbrace{\|w\|_2^2}_{\text{L2 regularization}}$$

~~$$+ \frac{\lambda}{2m} b^2$$~~  
 omit

$L_2$  regularization  $\underline{\|w\|_2^2} = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$

$L_1$  regularization  $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$

$w$  will be sparse

# Neural network

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

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

$$w: \begin{pmatrix} n^{[l]} & n^{[l-1]} \\ \uparrow & \uparrow \end{pmatrix}$$

"Frobenius norm"

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

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

$$dw^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

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

$$\rightarrow w^{[l]} := w^{[l]} - \alpha dw^{[l]}$$

"Weight decay"

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

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



deeplearning.ai

Regularizing your  
neural network

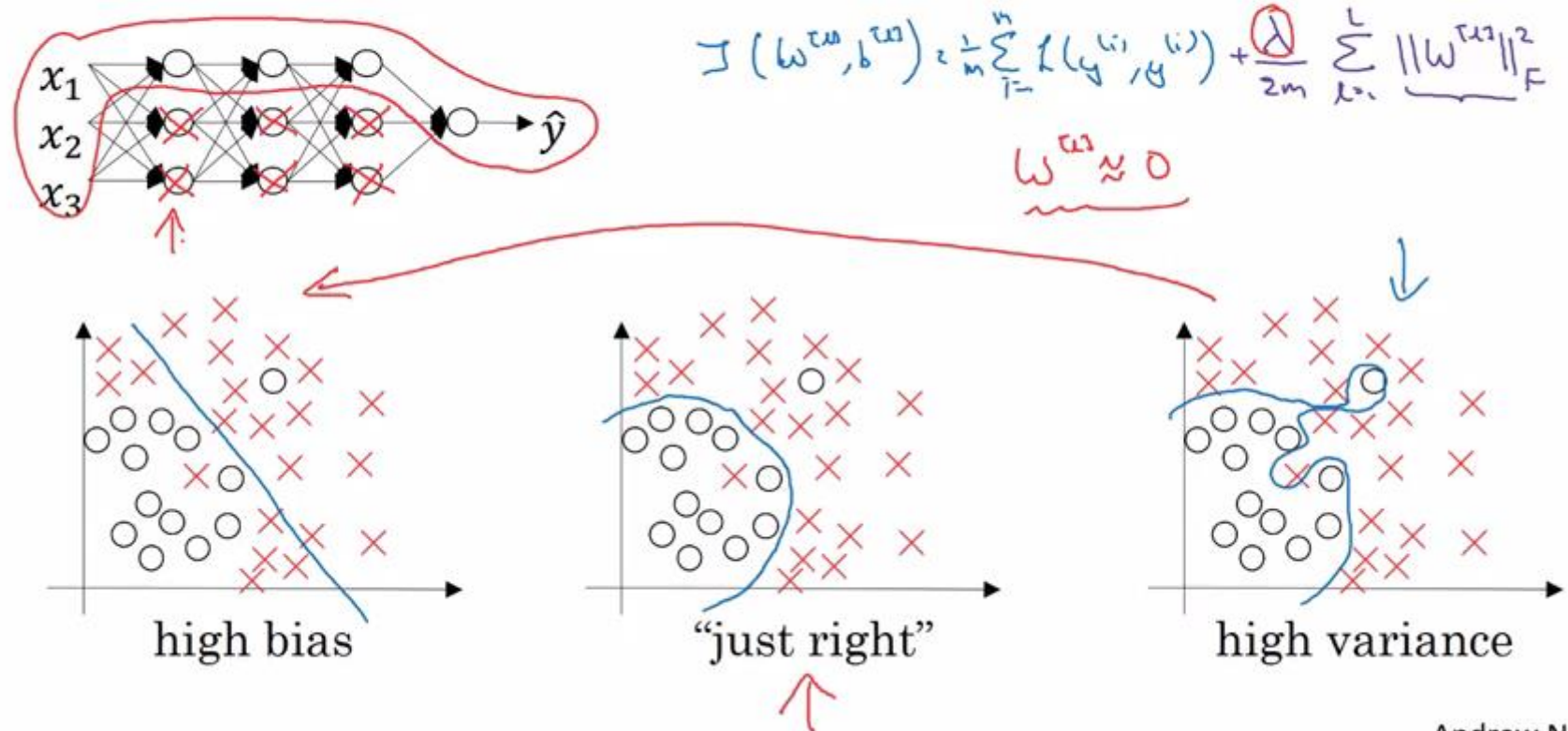
---

Why regularization  
reduces overfitting

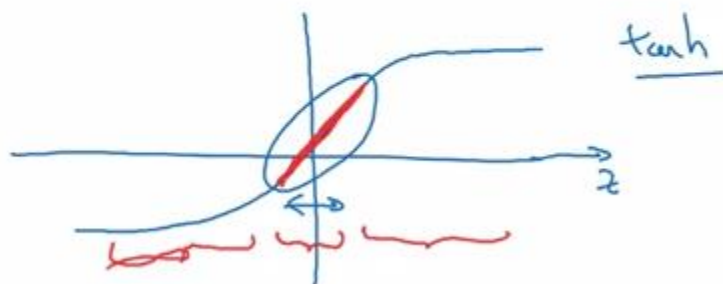




# How does regularization prevent overfitting?



# How does regularization prevent overfitting?



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



$\lambda \uparrow$

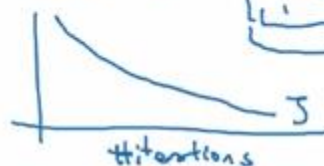
$W^{[L]} \downarrow$

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



Every layer  $\approx$  linear.

$$J(\dots) = \underbrace{\sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)})}_{\text{Loss}} + \underbrace{\frac{\lambda}{2} \sum_l \|W^{[l]}\|_F^2}_{\text{Regularization}}$$







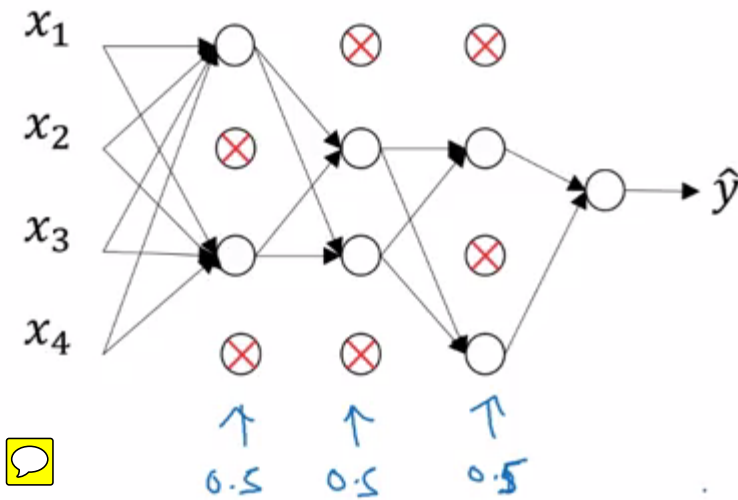
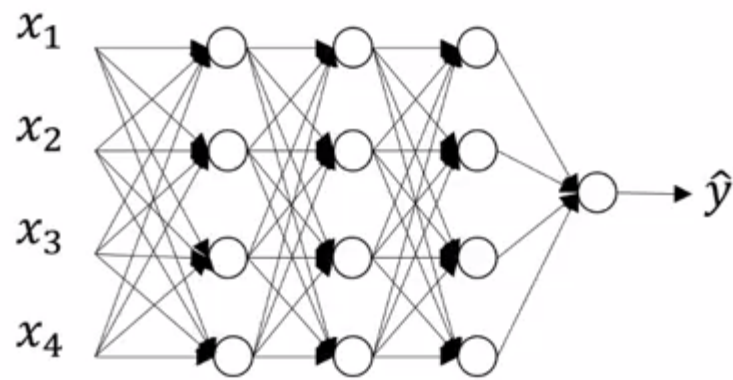
deeplearning.ai

# Regularizing your neural network

---

## Dropout regularization

# Dropout regularization



## Implementing dropout (“Inverted dropout”)

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

→  $d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep-prob}$

$a3 = \text{np.multiply}(a3, d3)$  #  $a3 \neq d3$ .

→  $a3 /= \text{keep-prob}$  ←

50 units.  $\leadsto$  10 units shut off

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

↑

↑ reduced by 20%.

$$/= 0.8$$

Test

## Making predictions at test time

$$a^{[0]} = X$$

No drop out.

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

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

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

$$a^{[2]} = \dots$$

↓  
 $\hat{y}$

$\neq$  keep-prob



deeplearning.ai

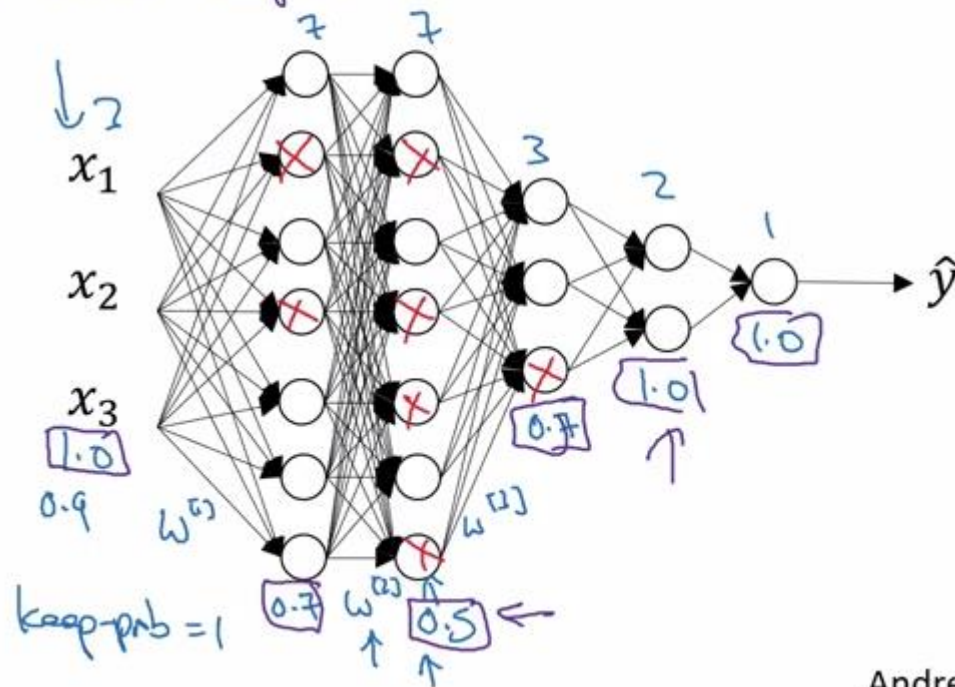
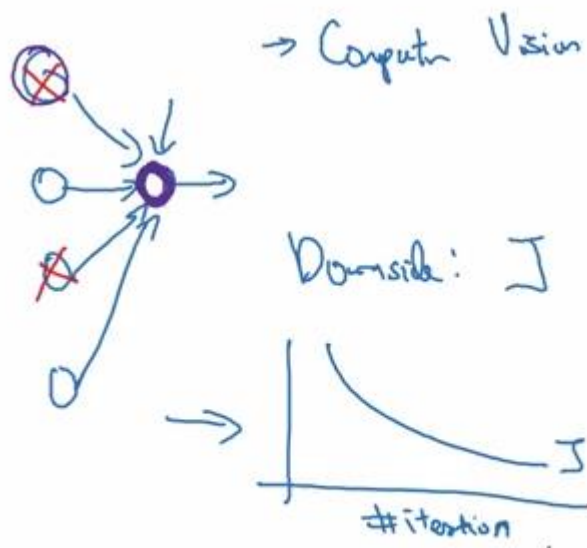
Regularizing your  
neural network

---

Understanding  
dropout

# Why does drop-out work?

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





deeplearning.ai

Regularizing your  
neural network

---

Other regularization  
methods

# Data augmentation



4



4

4

4



2:40 / 8:23



Andrew Ng

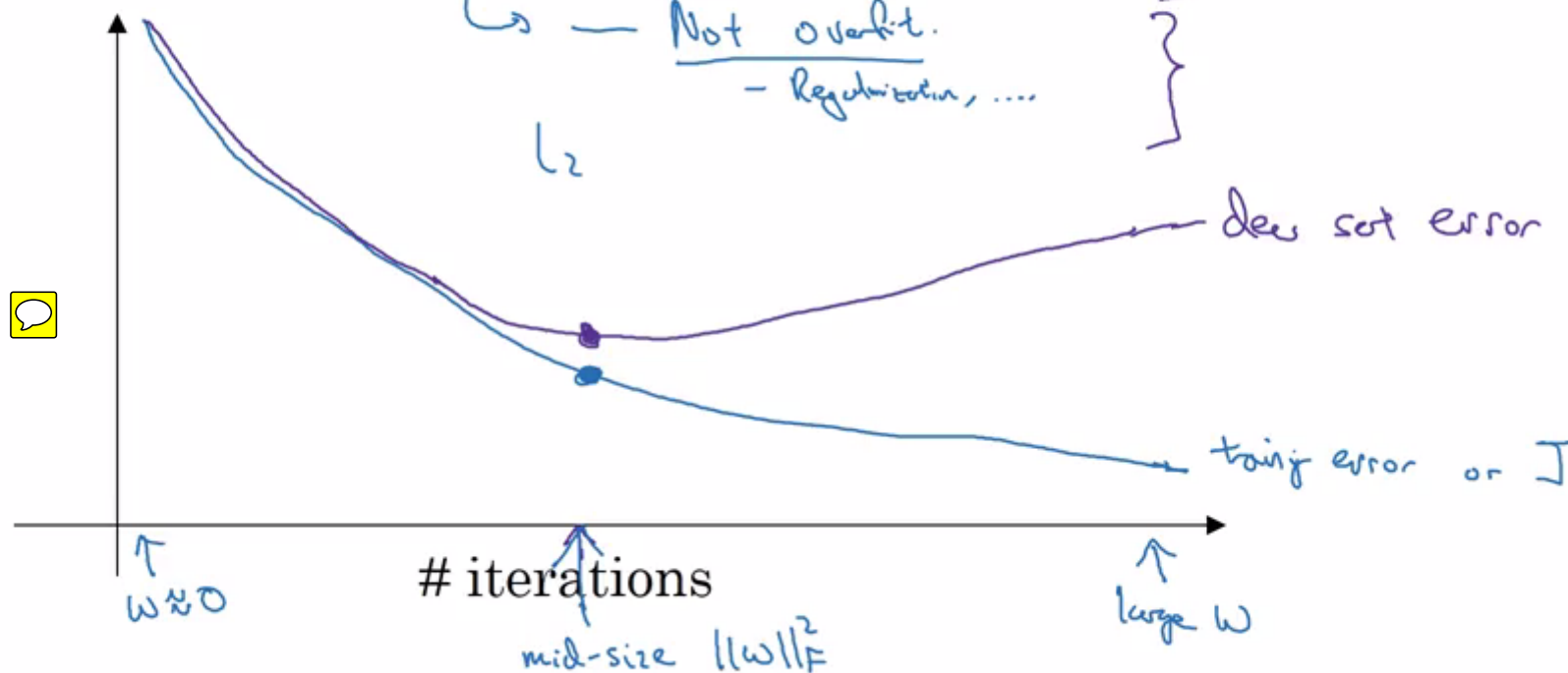


# Early stopping

Orthogonalization.

Optimize cost function  $J$   
- Gradient, ...  
Not overfit.  
- Regularization, ...

$J(w, b)$





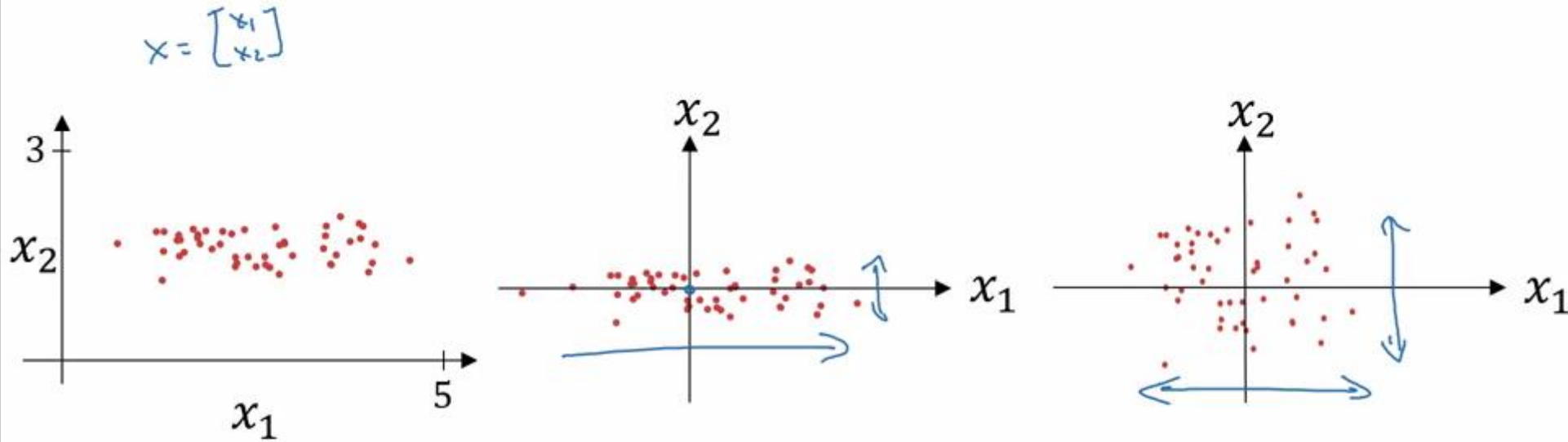
deeplearning.ai

Setting up your  
optimization problem

---

**Normalizing inputs**

# Normalizing training sets



Subtract mean:

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



$$x := x - \mu$$

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

Normalize variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^m x^{(i)} * x^{(i)} \quad \leftarrow \text{element-wise}$$

$$x /= \sigma$$



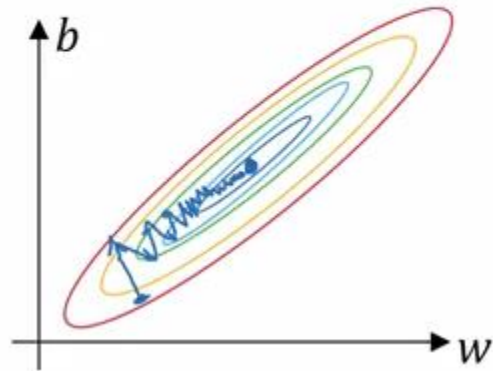
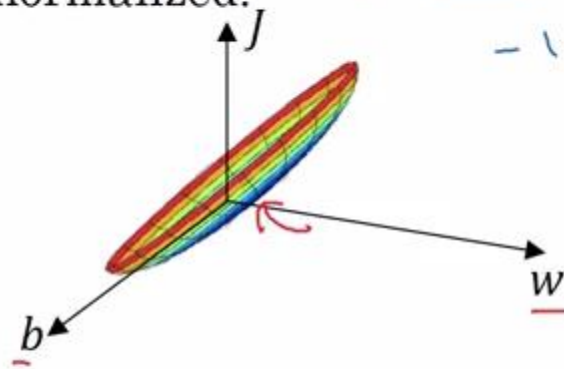
# Why normalize inputs?

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



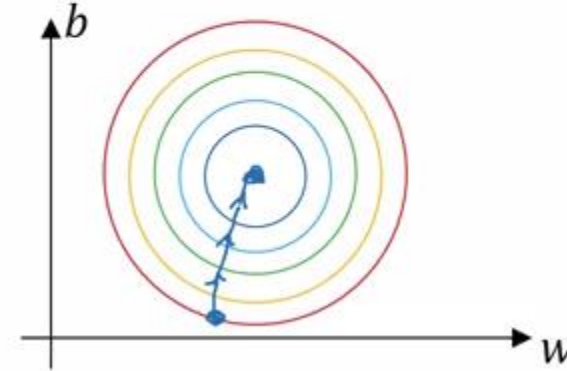
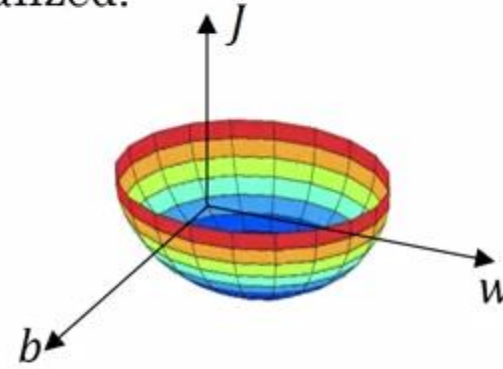
Unnormalized:

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



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

Normalized:





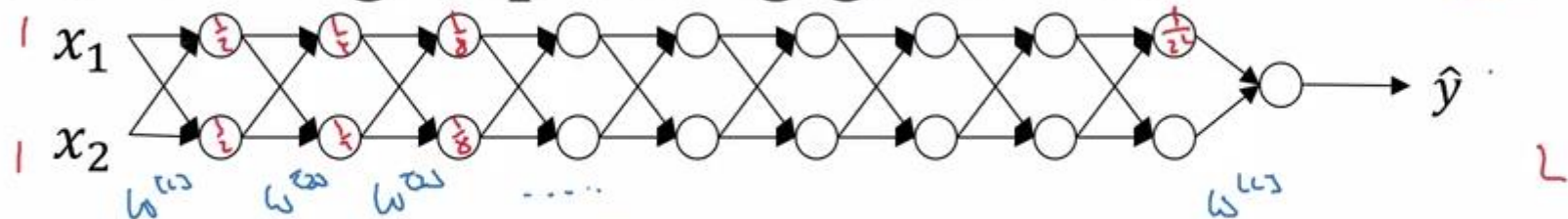
deeplearning.ai

Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

# Vanishing/exploding gradients



$g(z) = z$   
 $\hat{y} = w^{(L)} \left( w^{(L-1)} \left( w^{(L-2)} \dots \left( w^{(1)} x \right) \right) \right)$   
 $b^{(L)} = 0$   
 $1.5^L$   
 $0.5^L$



$w^{(1)} > I$



$w^{(2)} < I \begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$



$w^{(2)} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$   
 $0.5$   
 $0.5$

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

$z^{(1)} = w^{(1)} x$   
 $a^{(1)} = g(z^{(1)}) = z^{(1)}$   
 $a^{(2)} = g(z^{(2)}) = g(w^{(2)} a^{(1)})$





deeplearning.ai

Setting up your  
optimization problem

---

Weight initialization  
for deep networks

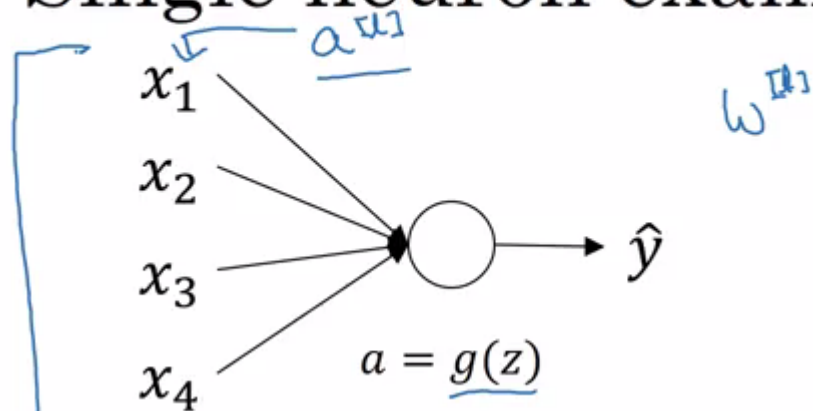


0:03 / 6:11





# Single neuron example



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

☐ large  $n \rightarrow$  Smaller  $w_i$

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

$\underline{w}^{[1]} = n.p. \cdot \text{random} \cdot \text{random}(\text{shape}) * n.p. \cdot \text{sqr}t\left(\frac{2}{n^{[1-1]}}\right)$   
 $\text{ReLU}$   $g^{[1]}(z) = \text{ReLU}(z)$

Other variants:

tanh

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

Xavier initialization ↑

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

↑





deeplearning.ai

Setting up your  
optimization problem

---

Numerical approximation  
of gradients

# Checking your derivative computation 云课堂

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



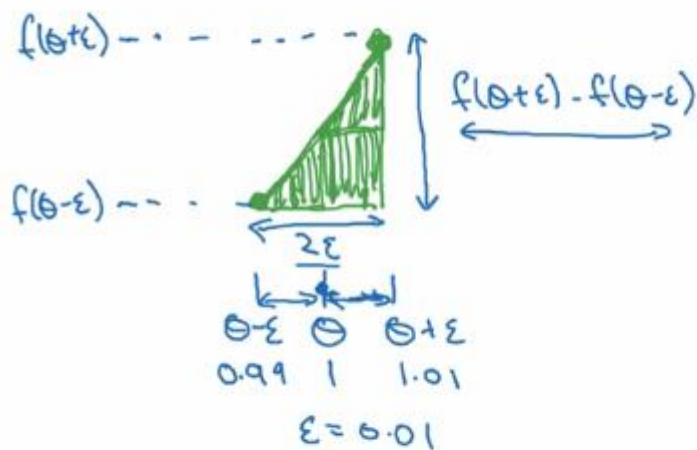
$$\frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \approx \underline{g(\theta)}$$

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

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

approx error: 0.0001

(prev slide: 3.0301, error: 0.03)



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



deeplearning.ai


---

Setting up your  
optimization problem

---


## Gradient Checking

# Gradient check for a neural network

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

*concentrate*

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

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

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

# Gradient checking (Grad check)

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each  $i$ :

☐  $\rightarrow \underline{d\theta_{\text{approx}}[i]} = \frac{J(\theta_1, \theta_2, \dots, \overset{\downarrow}{\theta_i + \epsilon}, \dots) - J(\theta_1, \theta_2, \dots, \overset{\downarrow}{\theta_i - \epsilon}, \dots)}{2\epsilon}$

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



Checks

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

$\epsilon = 10^{-7}$

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

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



6:18 / 6:34



Andrew Ng



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking  
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

$$\underline{db^{[L]}} \quad \underline{dw^{[L]}}$$

- Remember regularization.

$$\underline{J(\theta)} = \frac{1}{n} \sum_i \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2n} \sum_l \|w^{[l]}\|_F^2$$

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

- Doesn't work with dropout.

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

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

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