

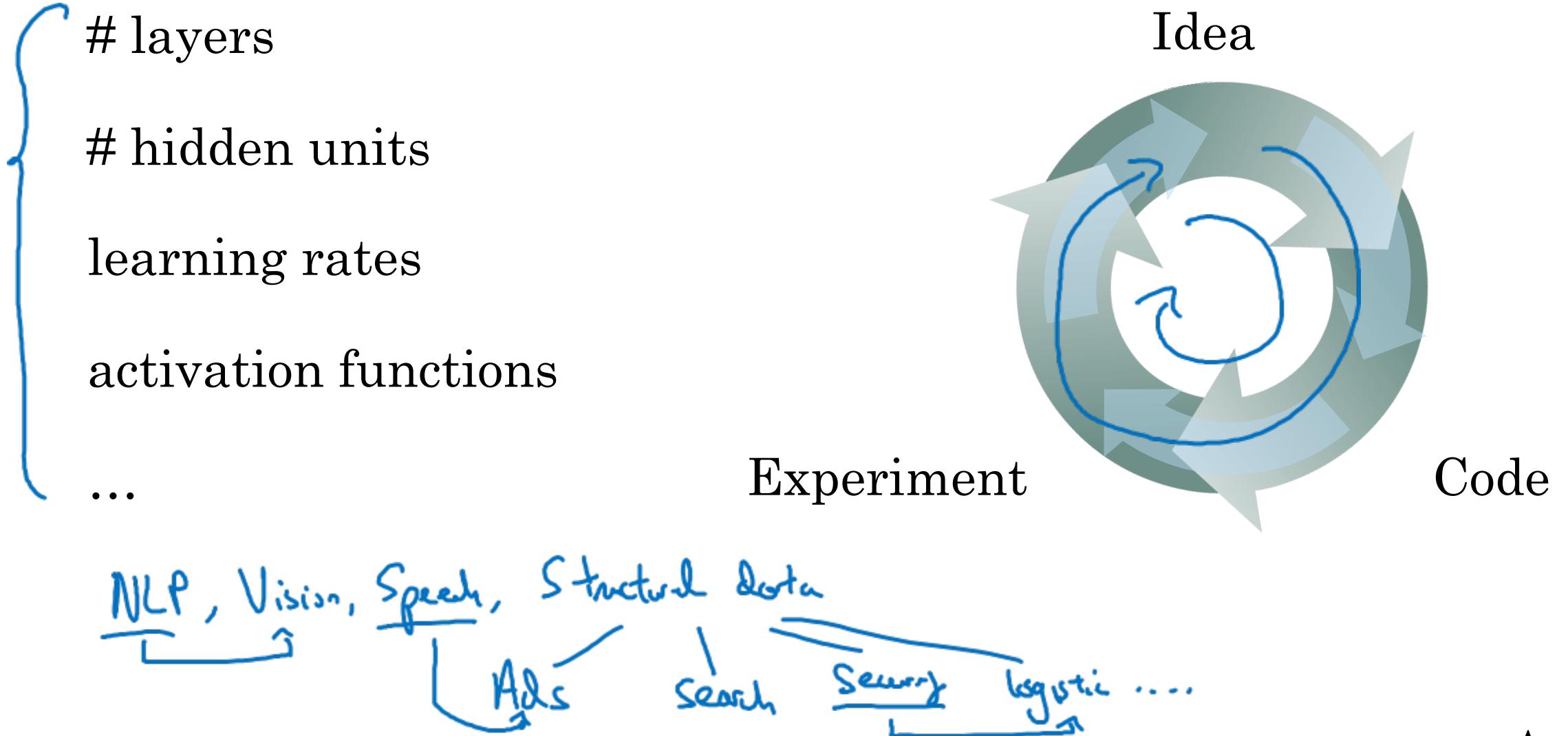


deeplearning.ai

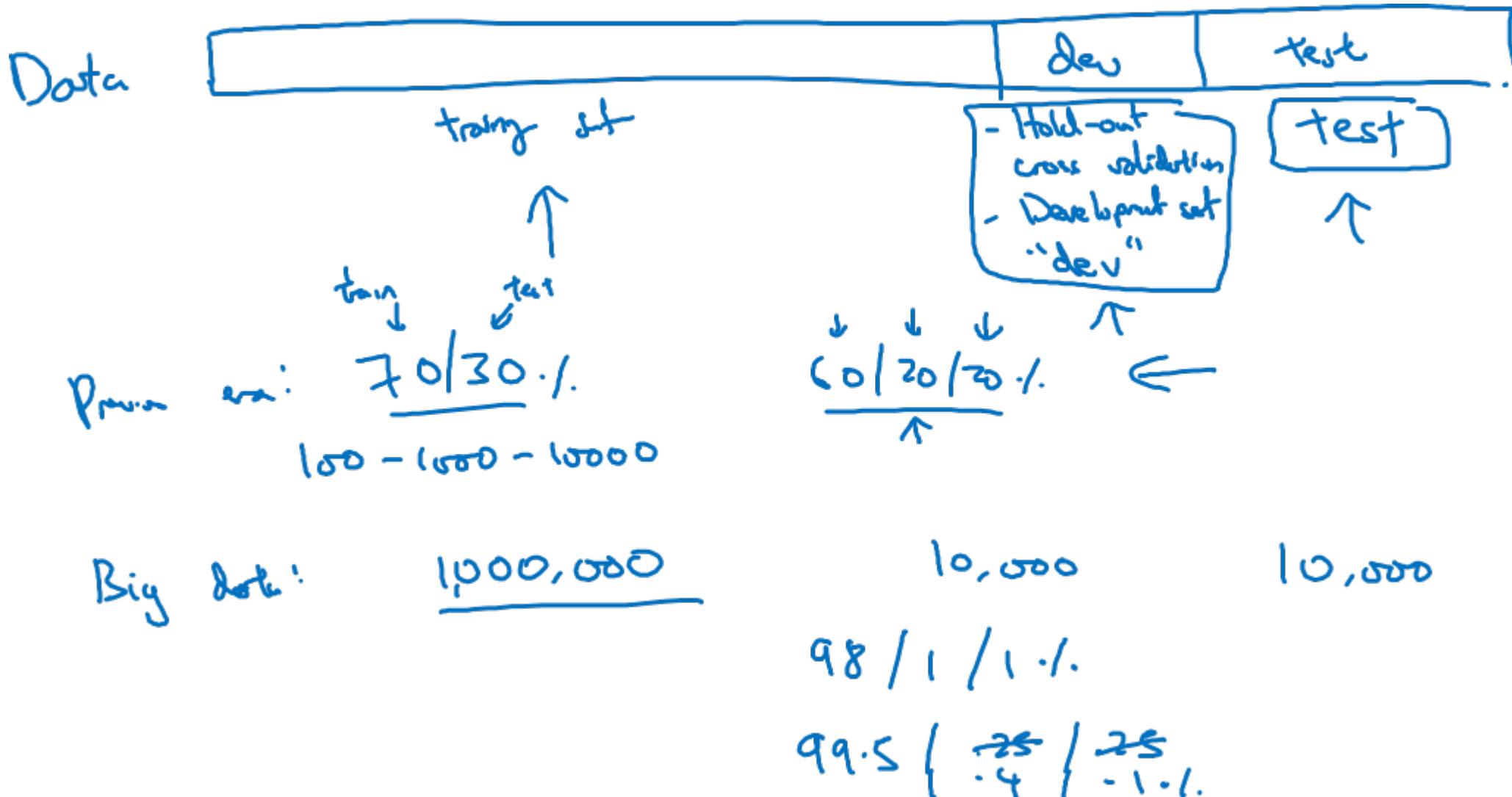
Setting up your ML application

Train/dev/test sets

Applied ML is a highly iterative process



Train/dev/test sets



Mismatched train/test distribution

Conts

Training set:

Cat pictures from }
webpages

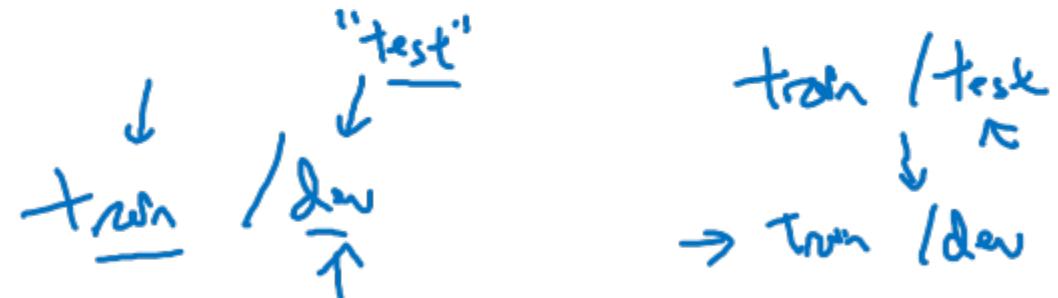


Dev/test sets:

Cat pictures from }
users using your app



→ Make sure dev and test come from same distribution.



train / test
↓
→ train / dev

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

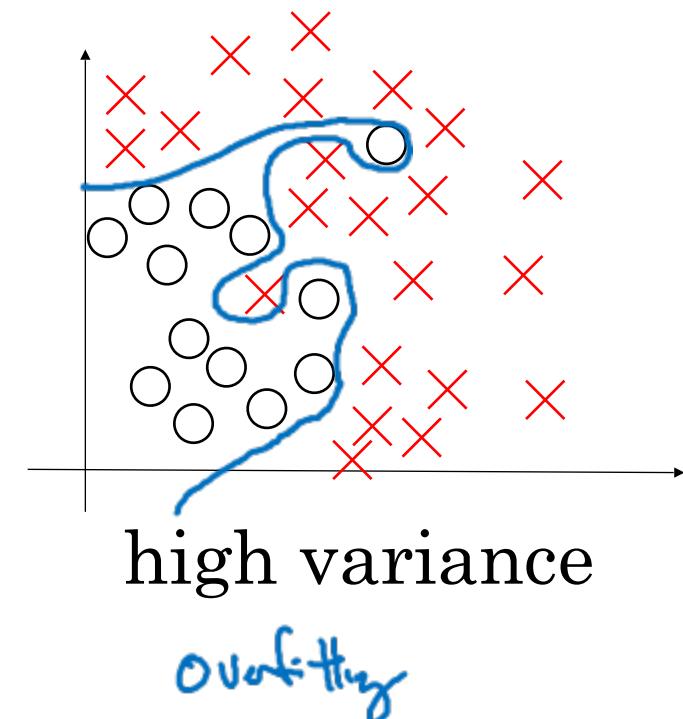
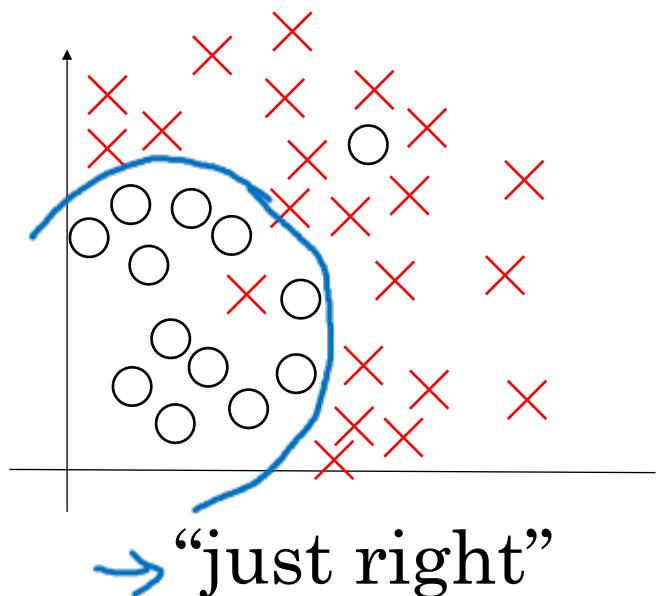
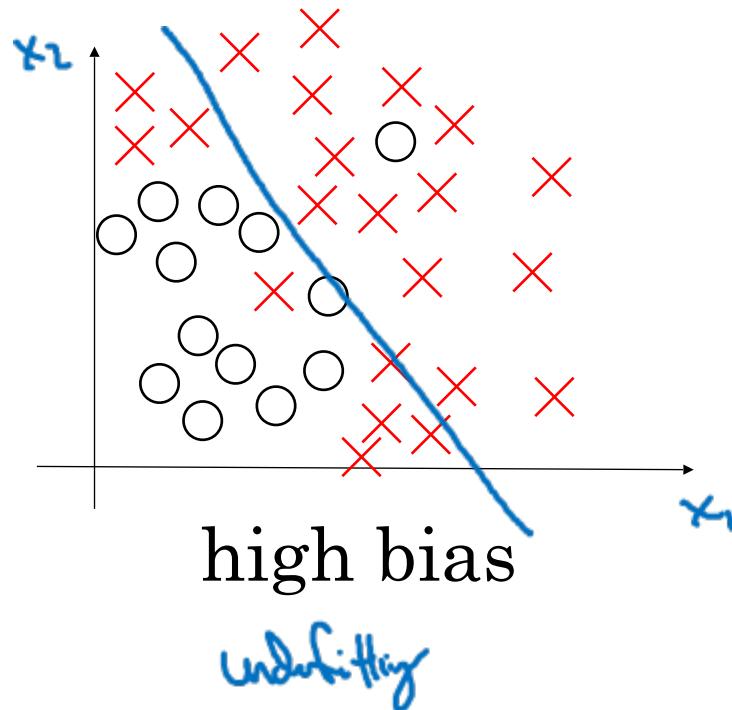


deeplearning.ai

Setting up your ML application

Bias/Variance

Bias and Variance



Bias and Variance

Cat classification

Train set error:

1%

$y=1$



Dev set error:

11%

15% ↗

$y=0$



Human: 10%

high variance



high bias



0.5%

1%

15%

30%

high bias
& high varan

low bias

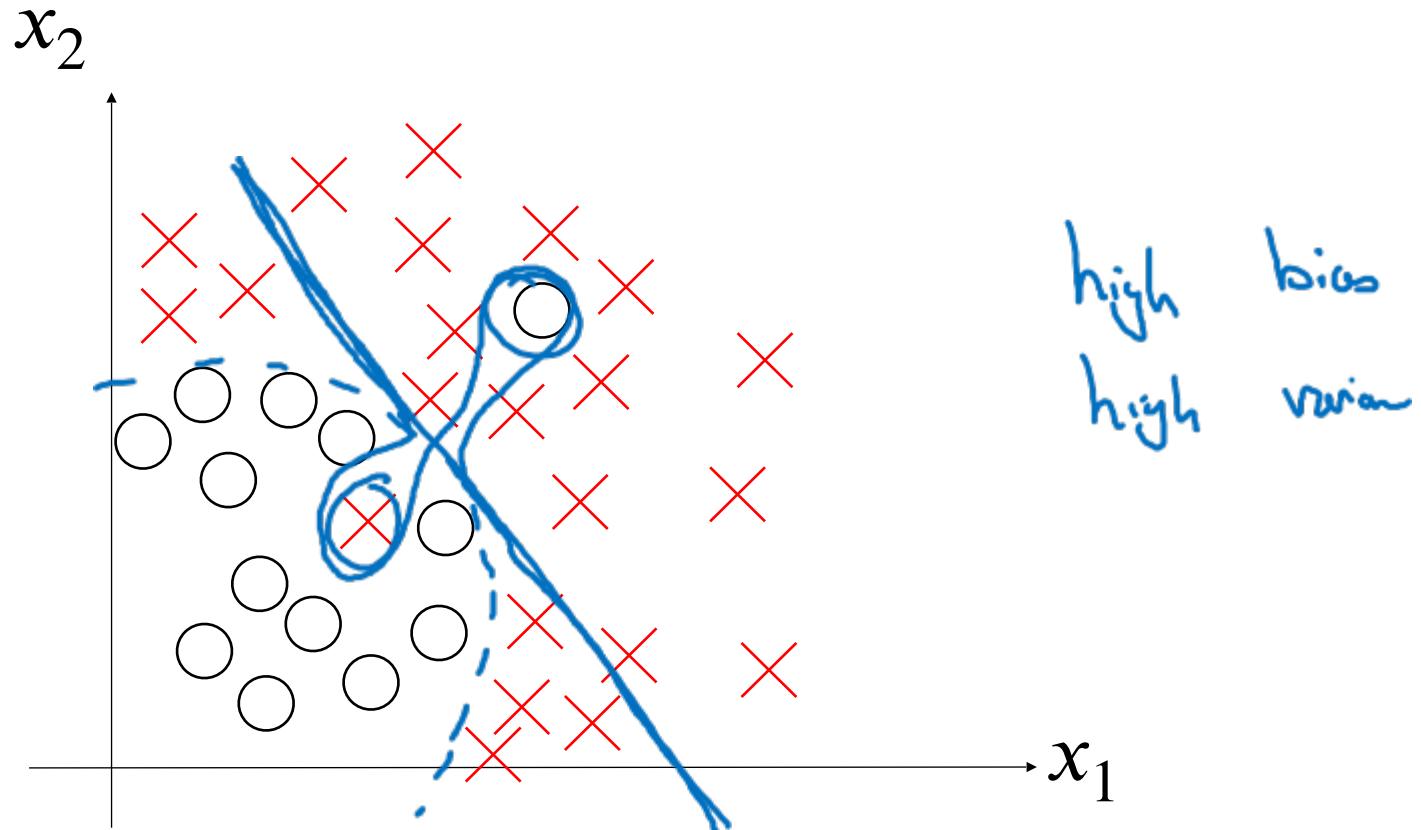
low variance



Optimal (Bayes) error: ~~10%~~ to 15%

Blurry images

High bias and high variance





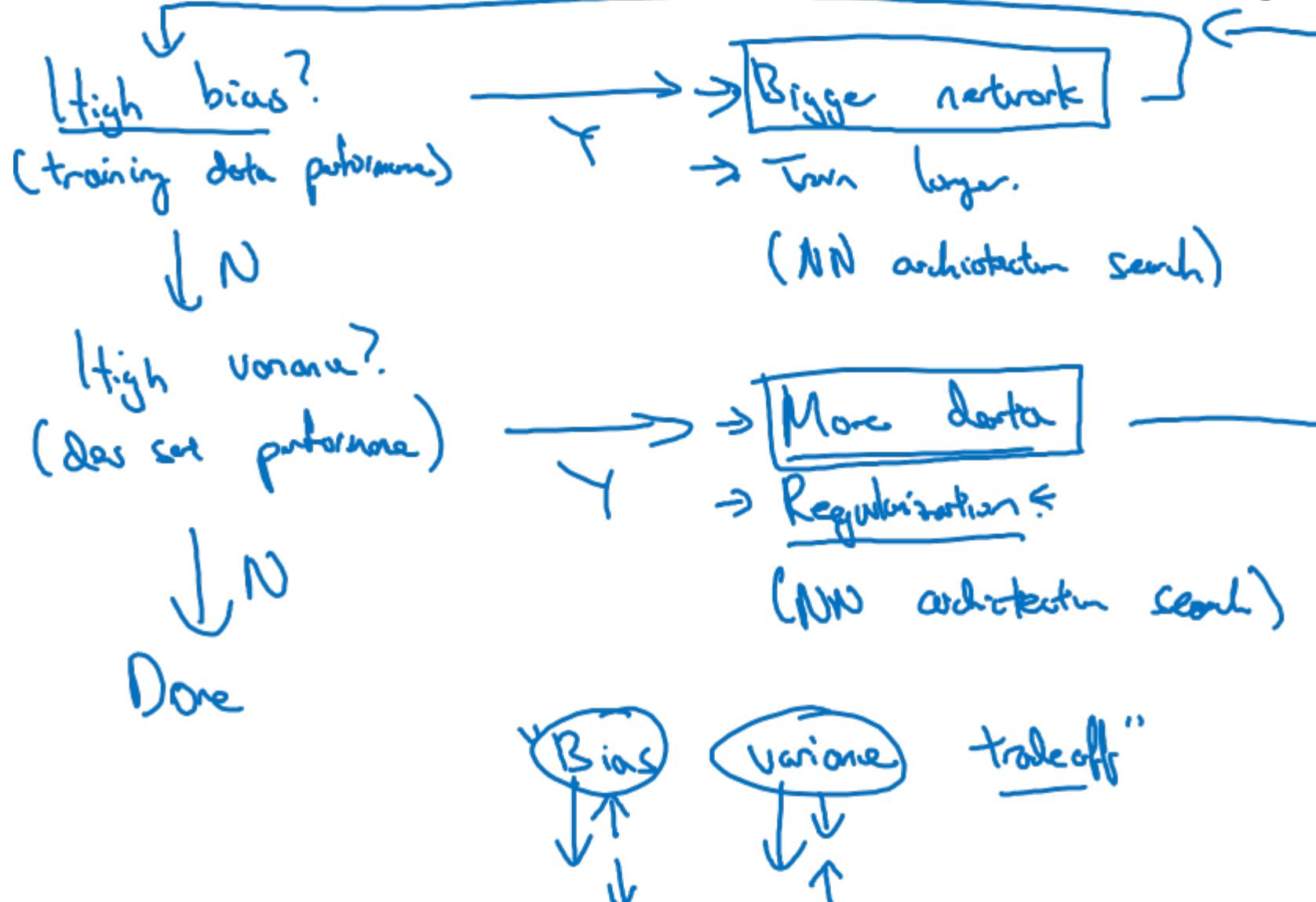
deeplearning.ai

Setting up your ML application

Basic “recipe” for machine learning

Basic “recipe” for machine learning

Basic recipe for machine learning





deeplearning.ai

Regularizing your
neural network

Regularization

Logistic regression

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

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

λ = regularization parameter
lambda lambd

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

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

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

L₁ regularization

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

w will be sparse

Neural network

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

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

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

"Frobenius norm"

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

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

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

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

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

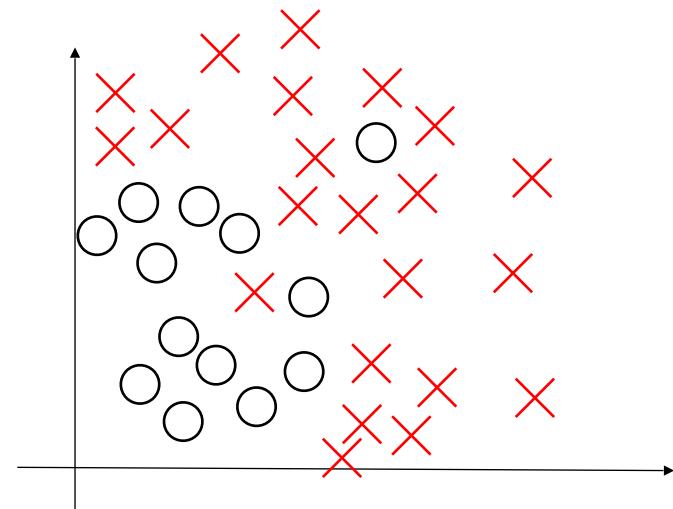
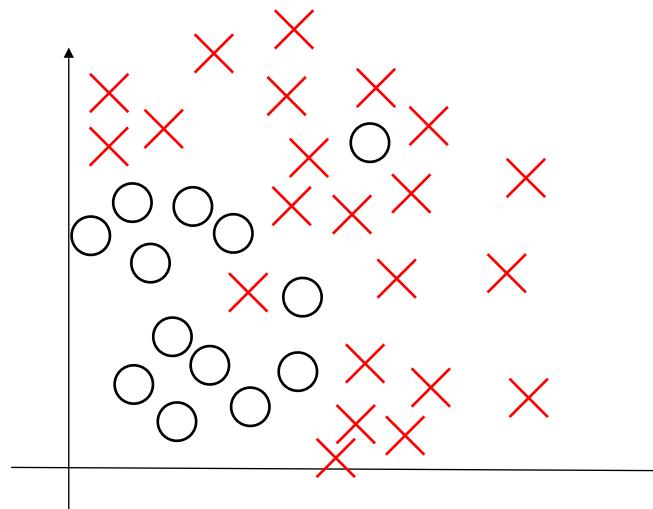
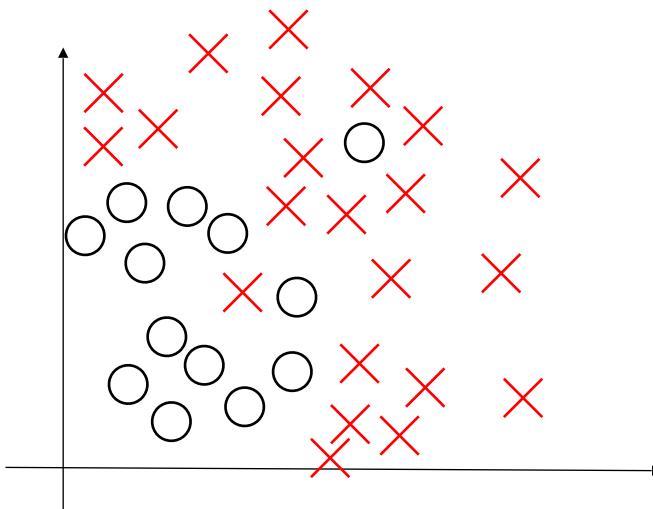
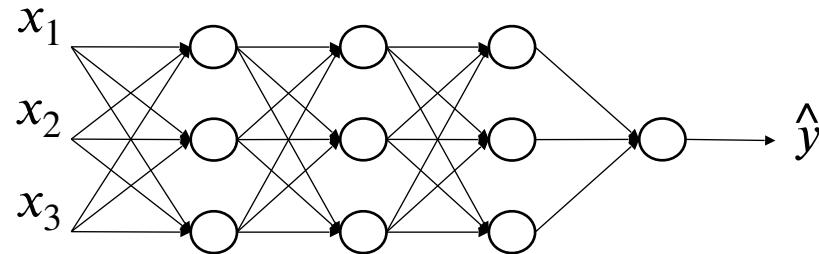
"Weight decay"

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

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

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

How does regularization prevent overfitting?



How does regularization prevent overfitting?

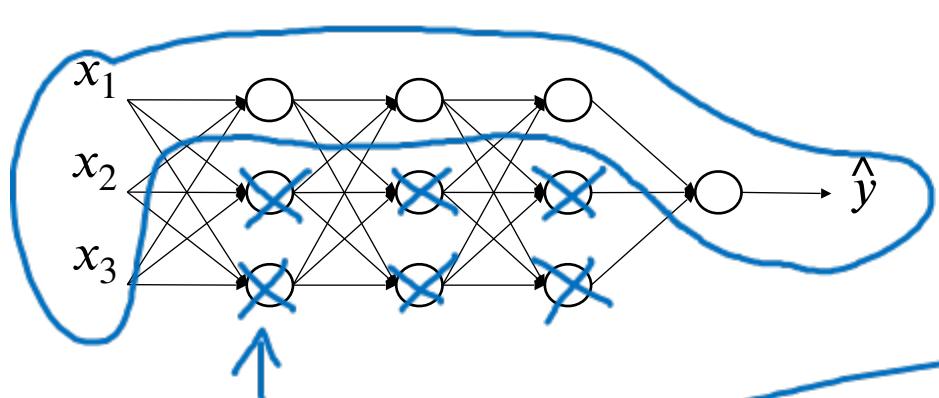


deeplearning.ai

Regularizing your neural network

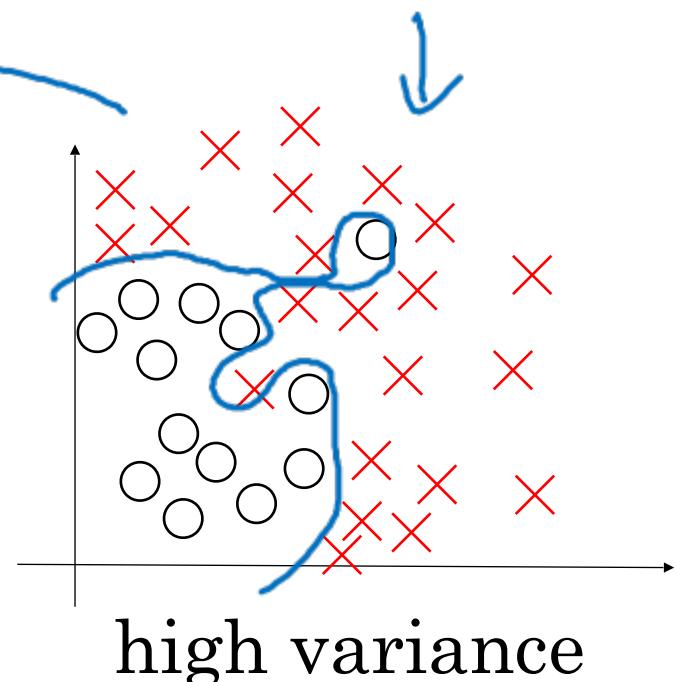
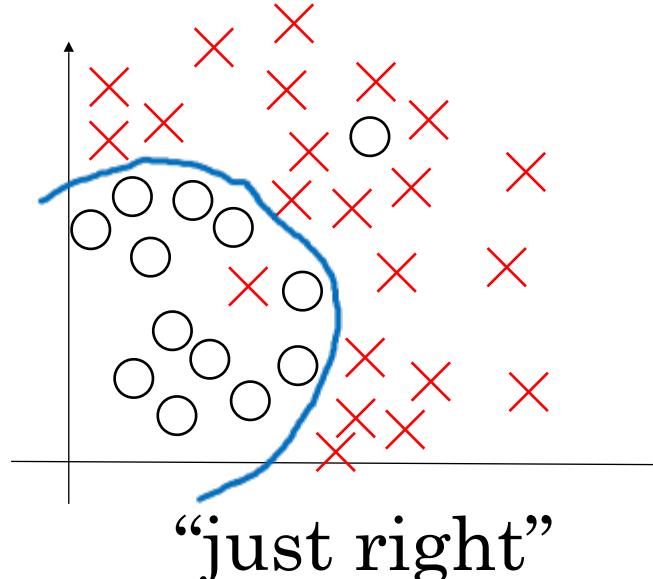
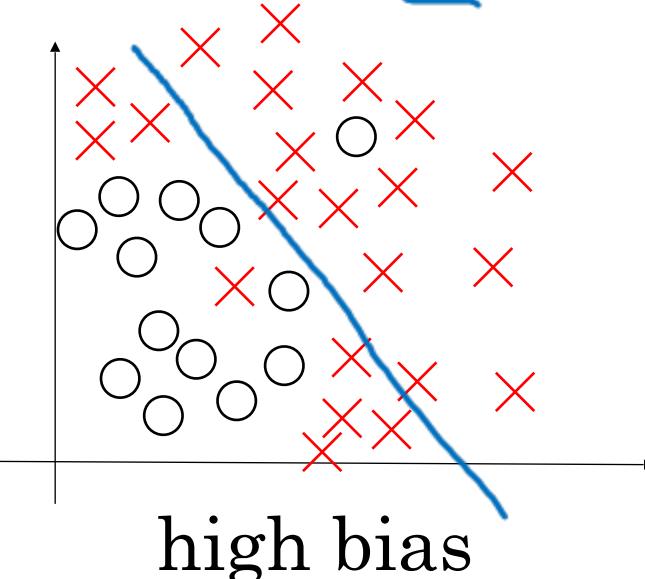
Why
regularization
reduces

How does regularization prevent overfitting?

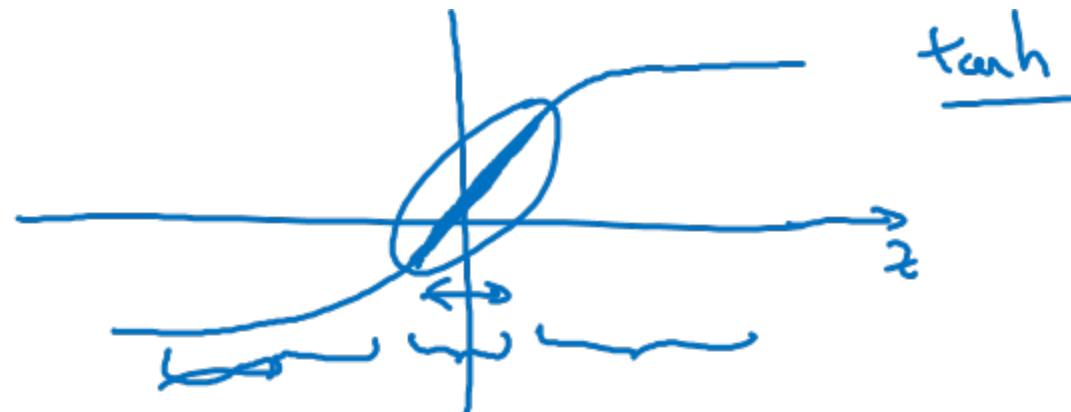


$$J(\omega^w, b^w) = \frac{1}{m} \sum_{i=1}^m (l(y^{(i)}, \hat{y}^{(i)})) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

$\omega^w \approx 0$



How does regularization prevent overfitting?



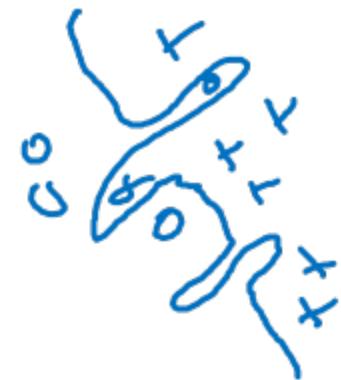
$$g(z) = \tanh(z)$$

$$\lambda \uparrow$$

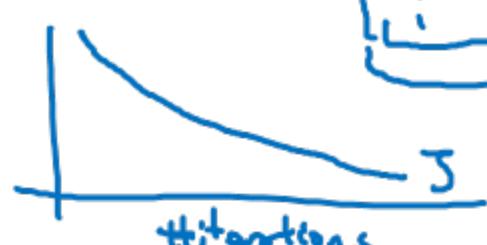
$$w^{[k]} \downarrow$$

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

Every layer \approx linear.



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



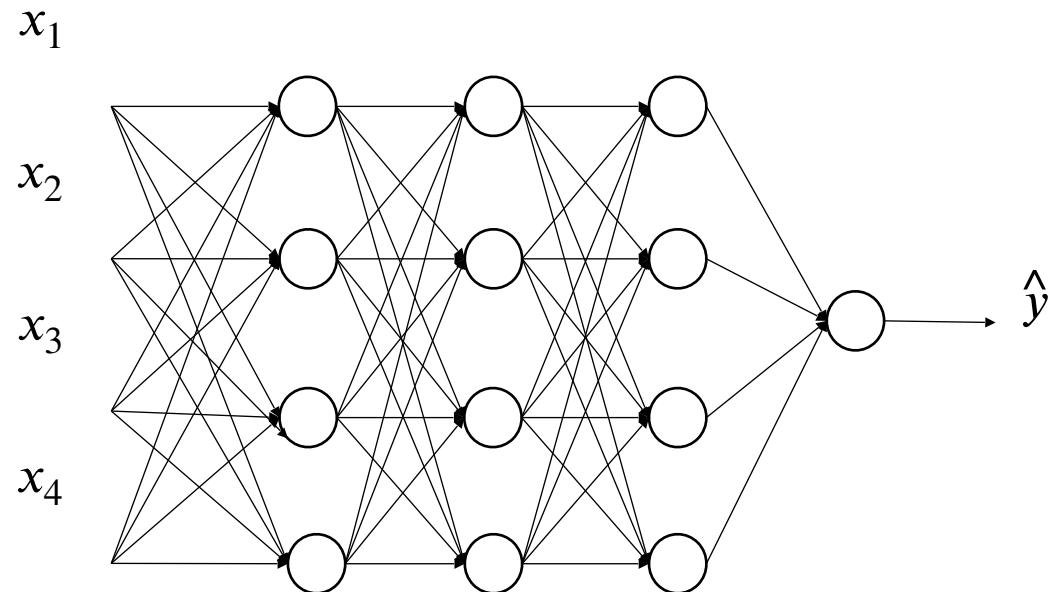


deeplearning.ai

Regularizing your neural network

Dropout regularization

Dropout regularization



\uparrow
0.5 \uparrow
0.5 \uparrow
0.5

Implementing dropout (“Inverted dropout”)

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

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

$$\overbrace{a_3}^{\text{a3}} = \text{np.multiply}(a_3, d_3) \quad \# a_3 * d3.$$

$$\rightarrow \overbrace{a_3 /=\cancel{\text{keep-prob}}}^{\text{a3}} \leftarrow$$

50 units. \rightsquigarrow 10 units shut off

$$z^{(4)} = w^{(4)} \cdot \overbrace{a^{(3)}}^T + b^{(4)}$$

T reduced by 20%.

Test

$$1 = \underline{\underline{0.8}}$$

Making predictions at test time

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

No drop out.

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

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

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

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

$$\downarrow \hat{y}$$

λ = keep-prob



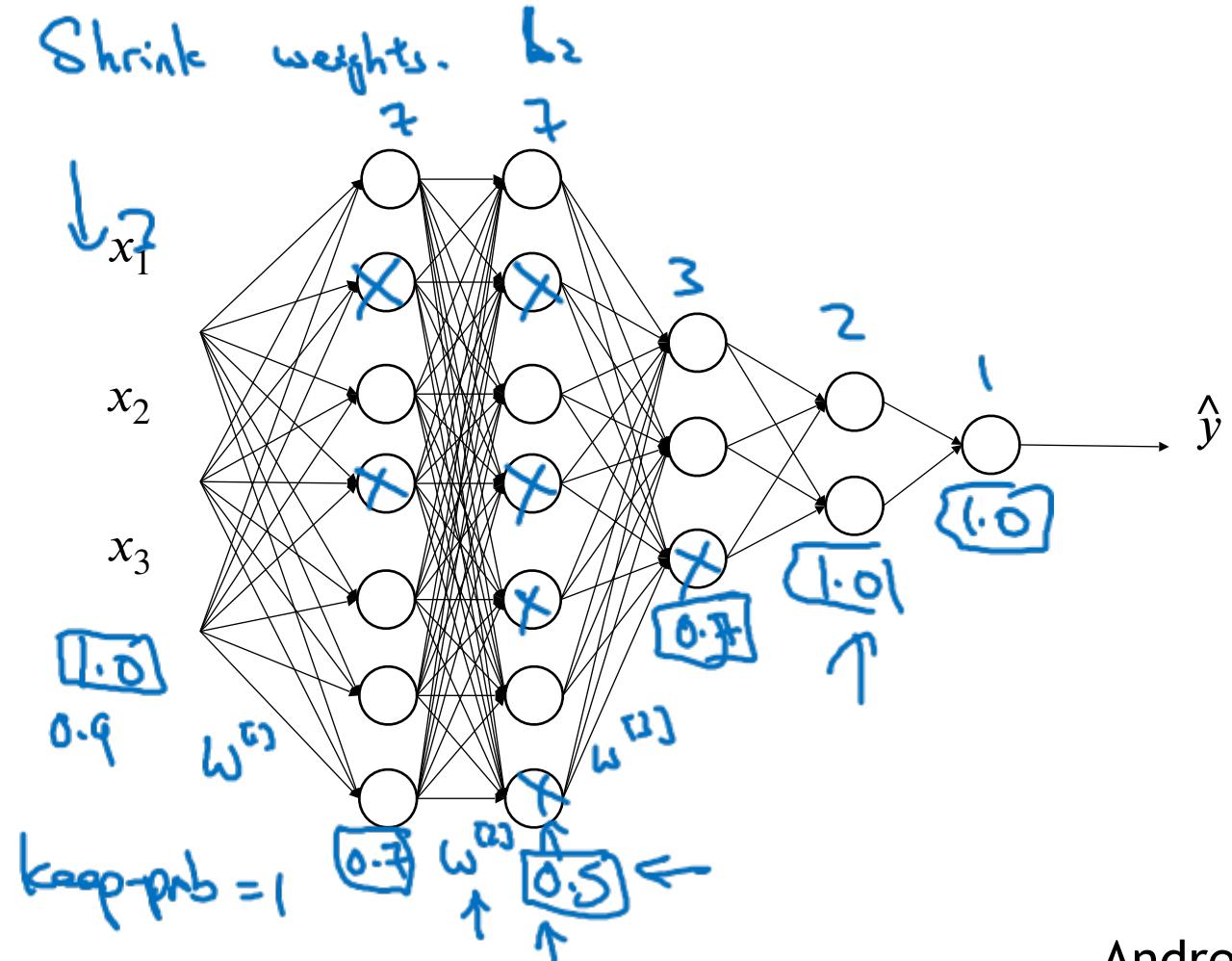
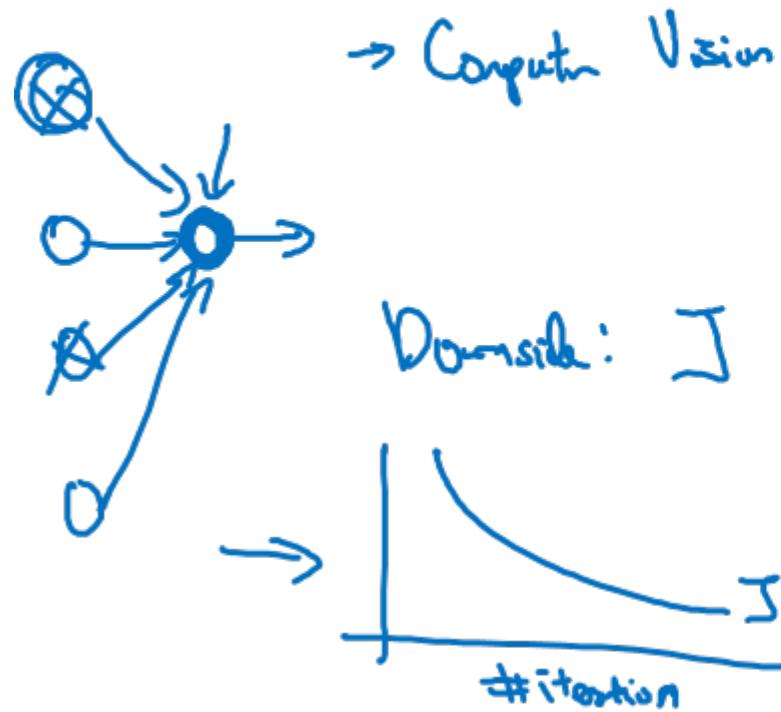
deeplearning.ai

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

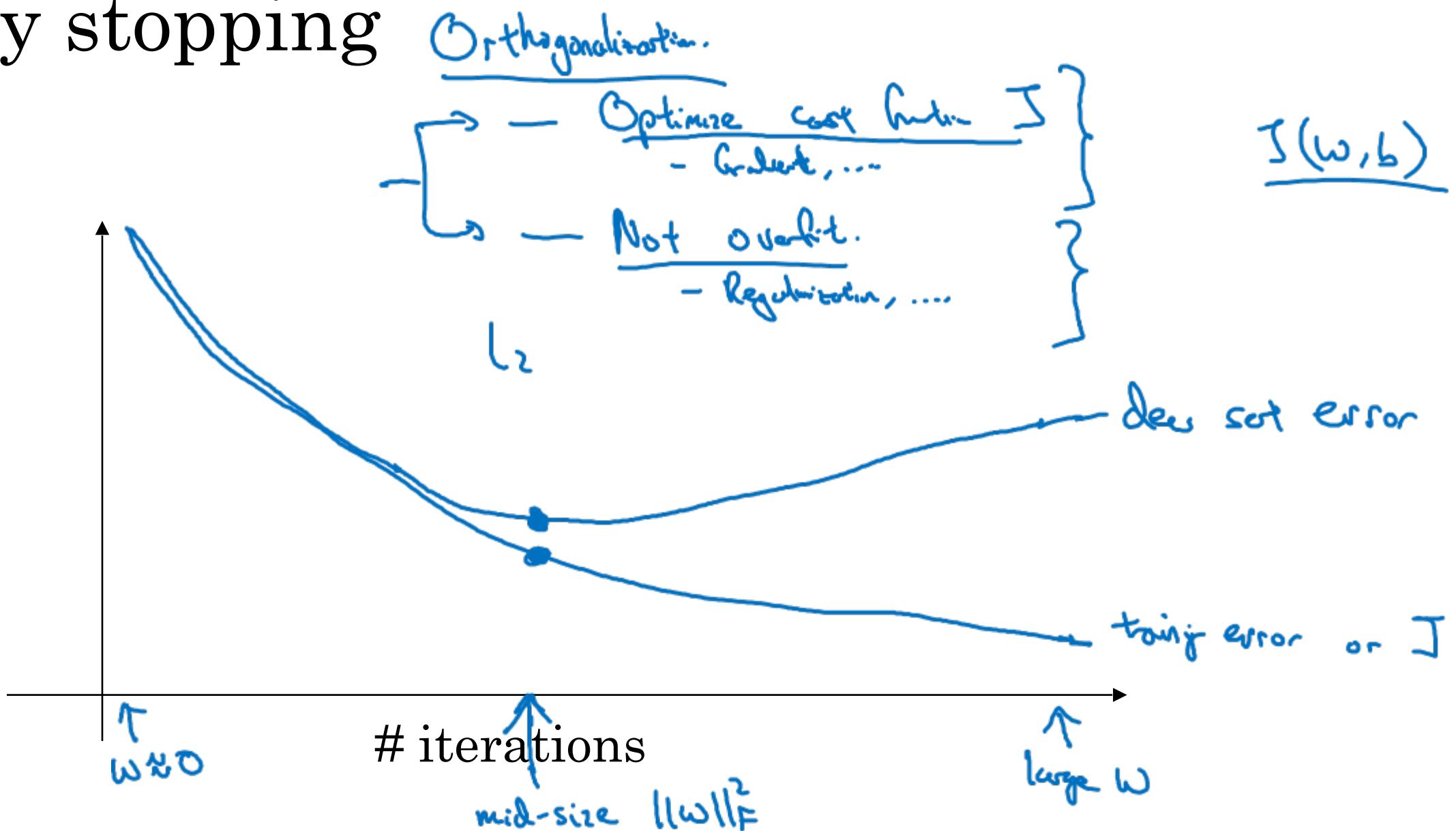
Data augmentation



4

A large, bold black digit '4' centered on the page.

Early stopping





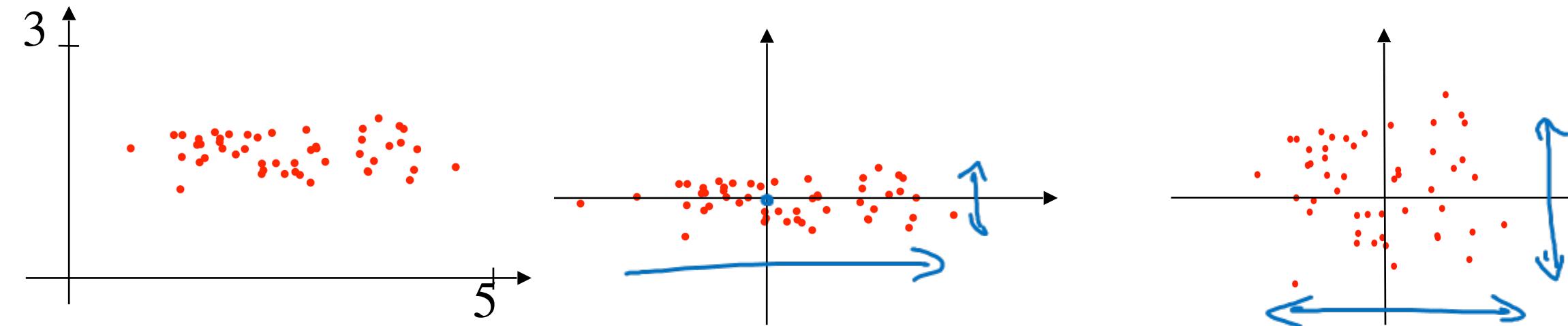
deeplearning.ai

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

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

\sim element-wise

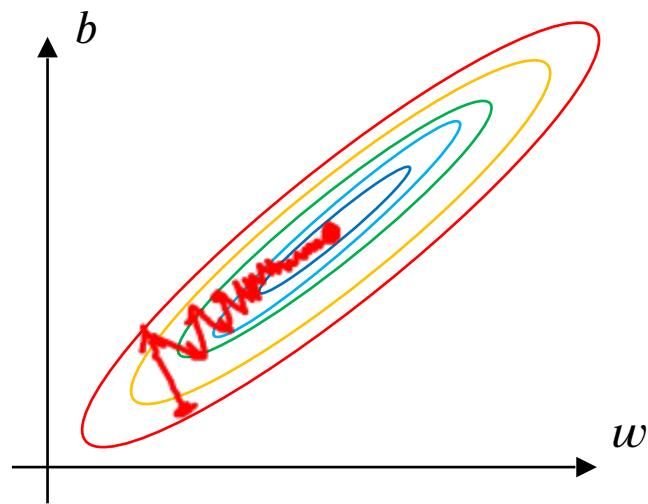
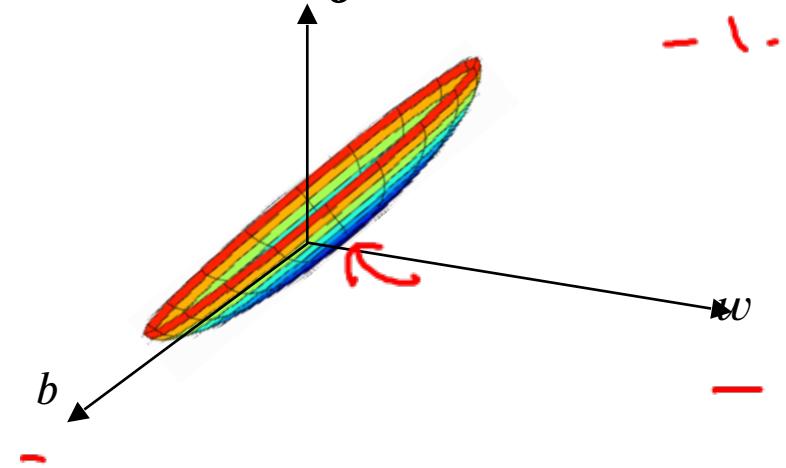
$$\underline{x / \sigma^2}$$

Use same μ, σ^2 to normalize test set.

Why normalize inputs?

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

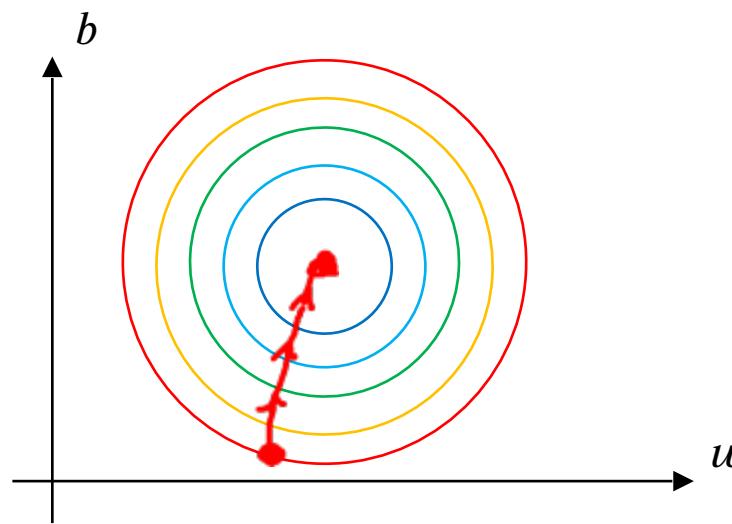
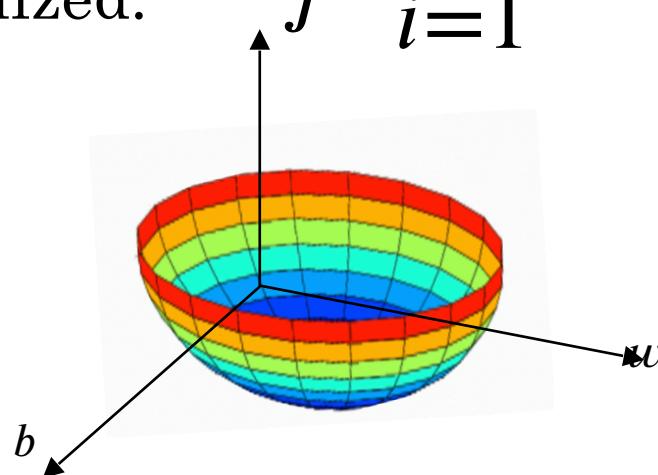
Unnormalized:



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

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

Normalized:





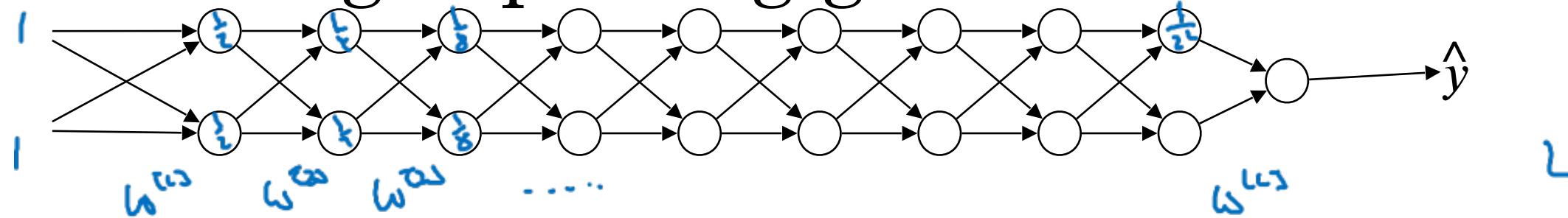
deeplearning.ai

Setting up your
optimization
problem

Vanishing/
exploding
gradients

Vanishing/exploding gradients

$L=150$



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



$$w^{[0]} > I$$

$$w^{[2]} < I \quad [0.9 \quad 0.9]$$

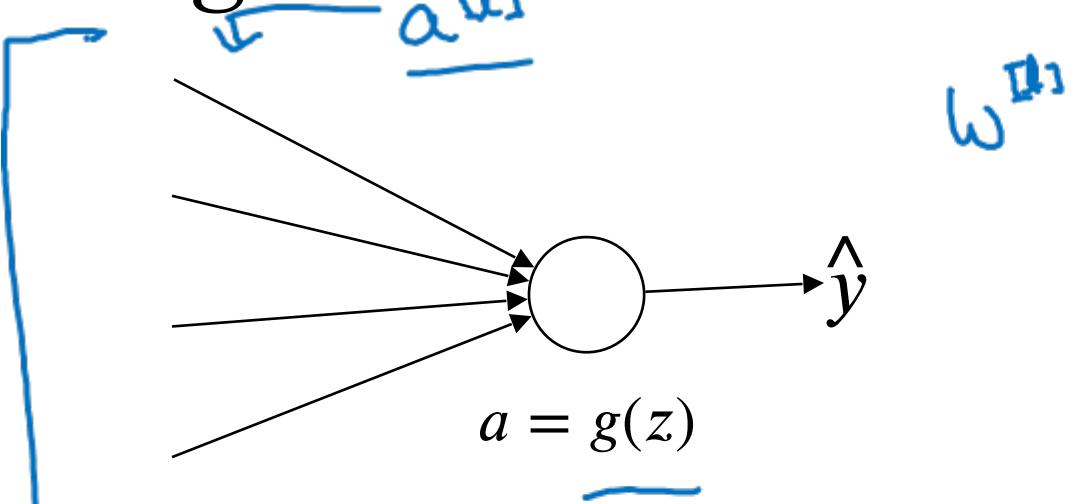
$$w^{[0]} = \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}$$

$$\begin{aligned} z^{[1]} &= \underline{w^{[0]} x} \\ a^{[0]} &= g(z^{[1]}) = z^{[1]} \end{aligned}$$

$$\begin{aligned} a^{[2]} &= g(z^{[1]}) = g(w^{[2]} a^{[1]}) \\ \hat{y} &= \underline{w^{[L-1]} \begin{bmatrix} 0.5 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0.5 \end{bmatrix}^{L-1} x} \end{aligned}$$

$$\begin{aligned} 1.5^{L-1} x \\ 0.5^{L-1} x \end{aligned}$$

Single neuron example



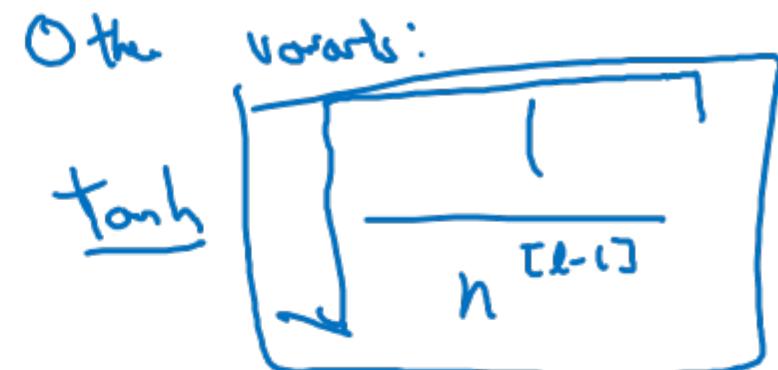
$$z = \underline{w_1 x_1 + w_2 x_2 + \dots + w_n x_n} \quad \times$$

Large $n \rightarrow$ Smaller w_i

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

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

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



Xavier initialization ↑

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

↑



deeplearning.ai

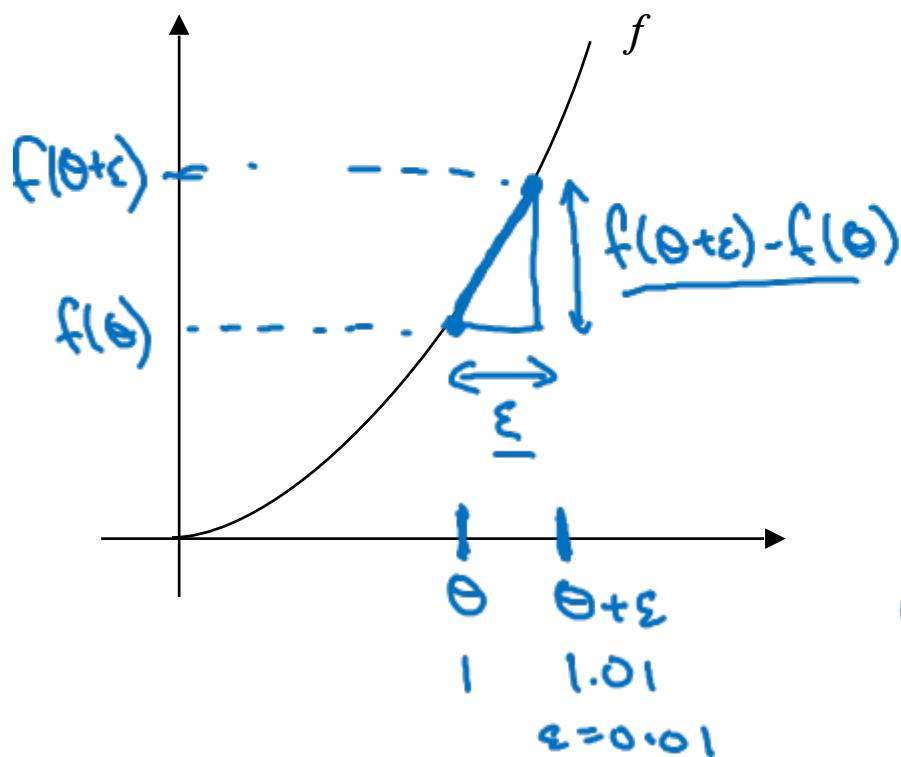
Setting up your optimization problem

Numerical approximation of gradients

Checking your derivative computation

$$\frac{f(\theta)}{\theta \in \mathbb{R.}}$$

I



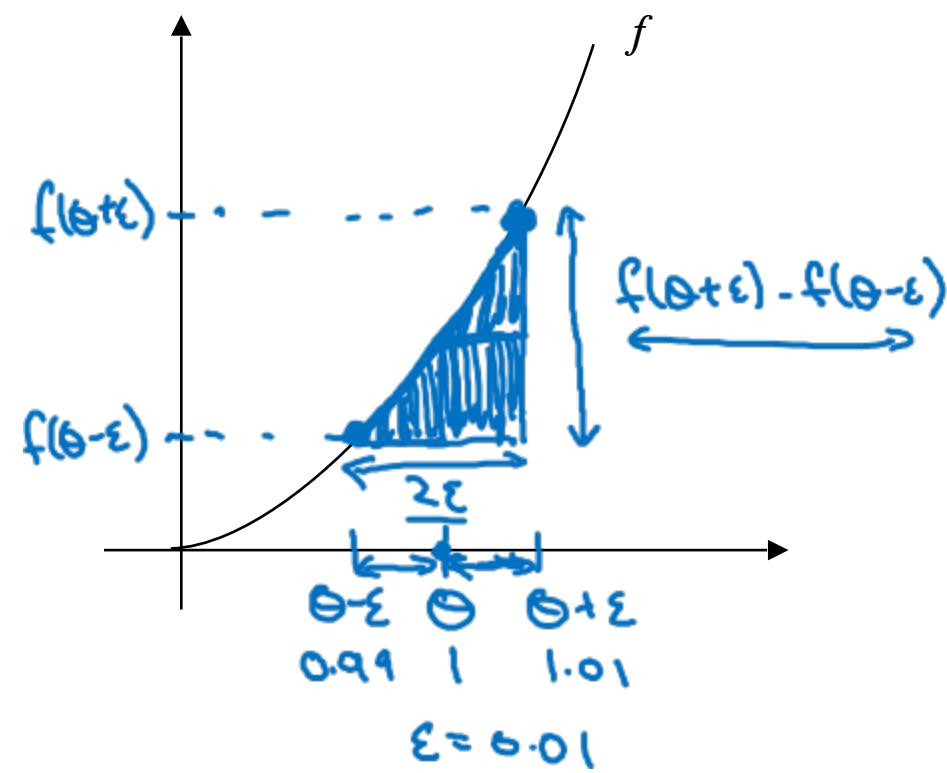
$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$
$$\frac{dw}{db}$$
$$\frac{g(\theta)}{\theta = 1}$$
$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

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

$\theta = 1$
 $\theta + \epsilon = 1.01$
 $\epsilon = 0.01$

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)

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



deeplearning.ai

Setting up your
optimization problem

Gradient Checking

Gradient check for a neural network

Take $\boxed{\text{and}}$ and $\boxed{\text{reshape}}$ into a big vector

concatenate

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

Take $\boxed{\text{and}}$ and $\boxed{\text{reshape}}$ into a big vector

concatenate

J

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 :

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

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

Checks

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

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

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

$$10^{-5}$$

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



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.

$$\frac{\partial b}{\uparrow} \quad \frac{\partial w}{\uparrow}$$

- Remember regularization.

$$J(\theta) = \frac{1}{m} \sum_i \ell(y_i, \hat{y}_i) + \frac{\lambda}{2m} \sum_j \|w^{(j)}\|_F^2$$

$\Delta\theta = \text{gradt of } J \text{ wrt. } \theta$

- Doesn't work with dropout.

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

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

$w, b \text{ no}$