

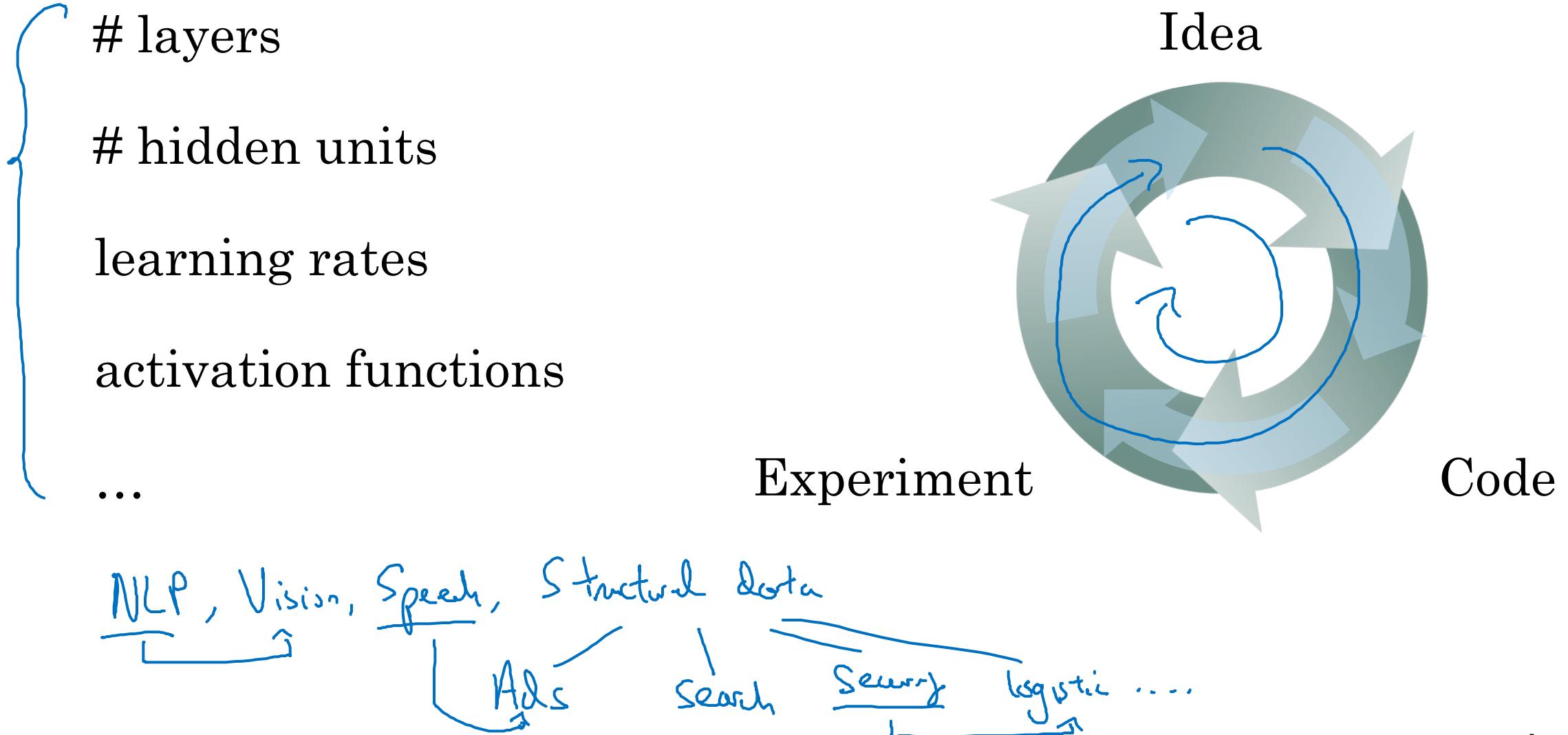


deeplearning.ai

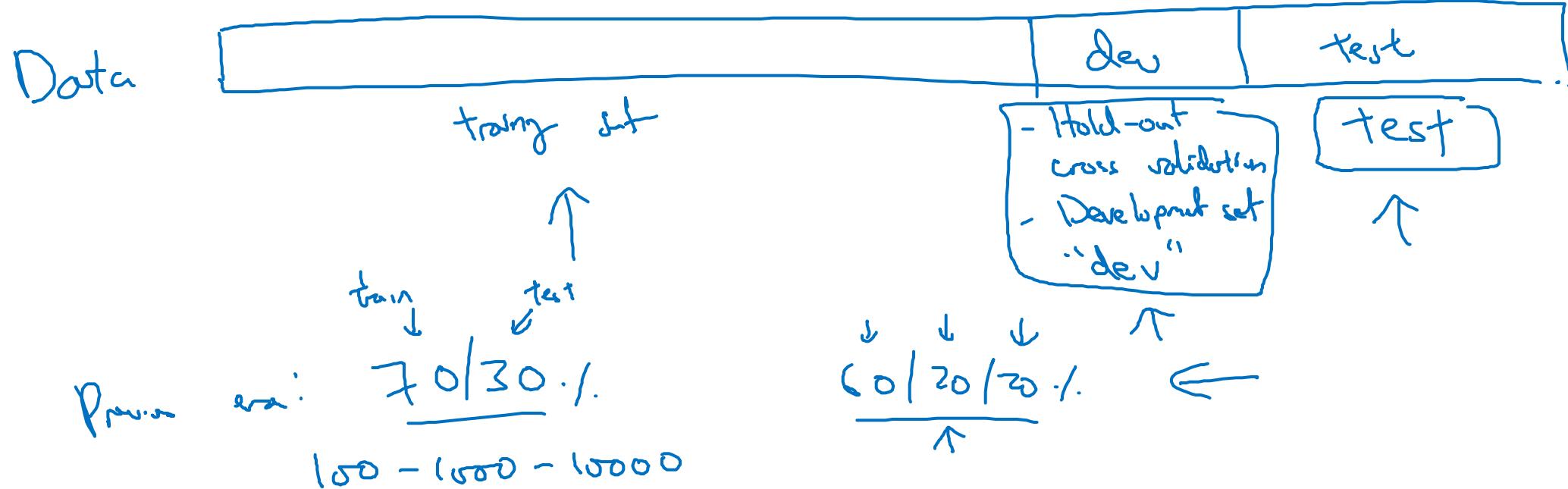
Setting up your
ML application

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 · 1

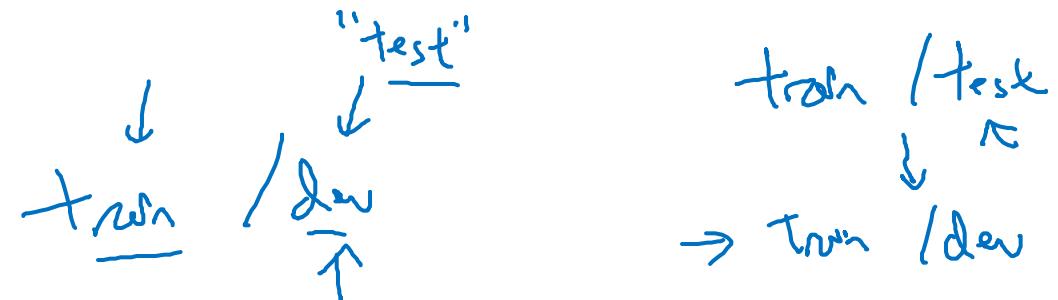
Mismatched train/test distribution

Conts

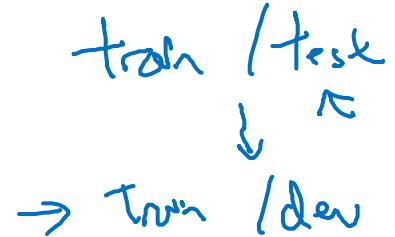
Training set:
Cat pictures from }
webpages

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.)



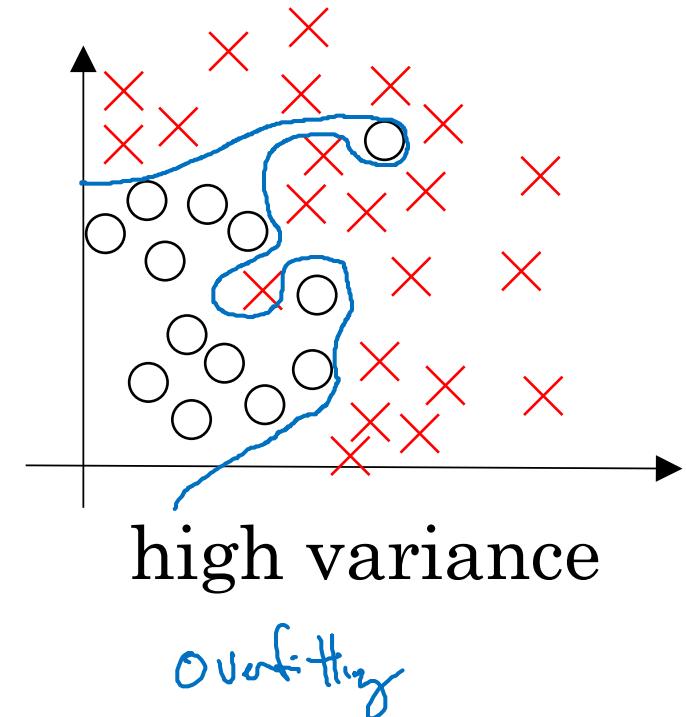
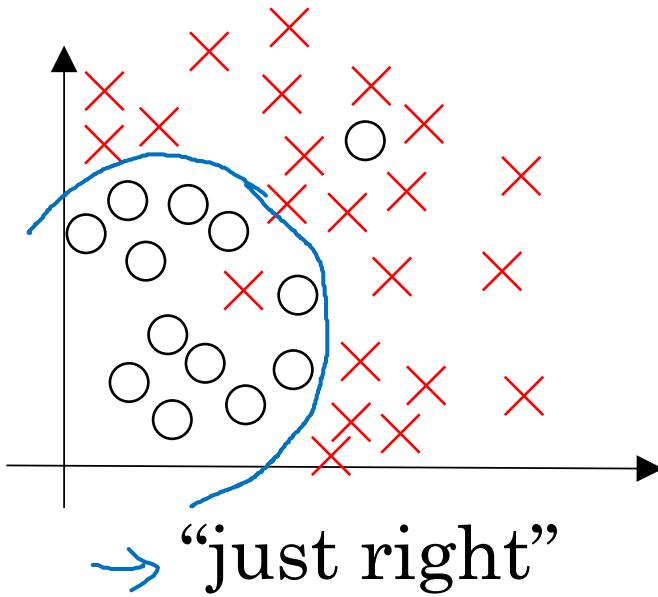
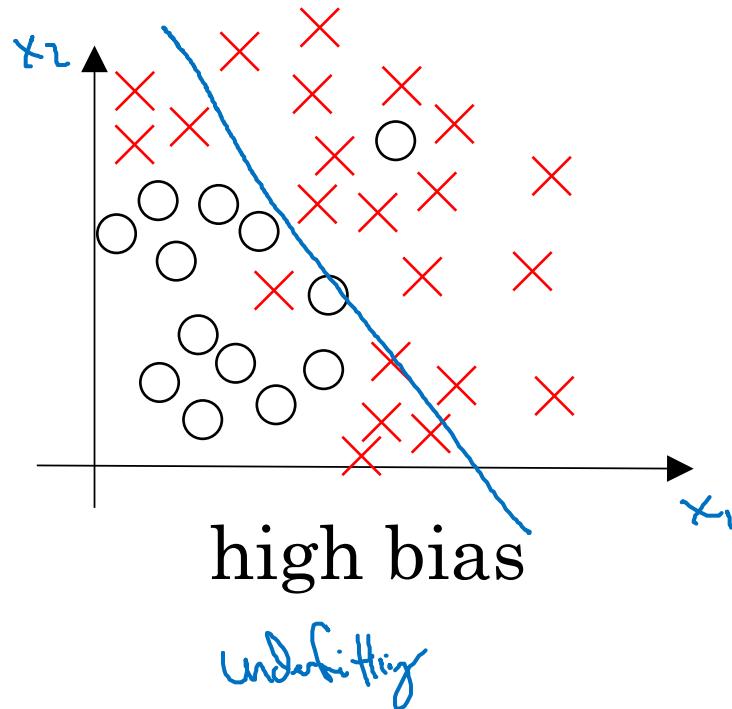


deeplearning.ai

Setting up your
ML application

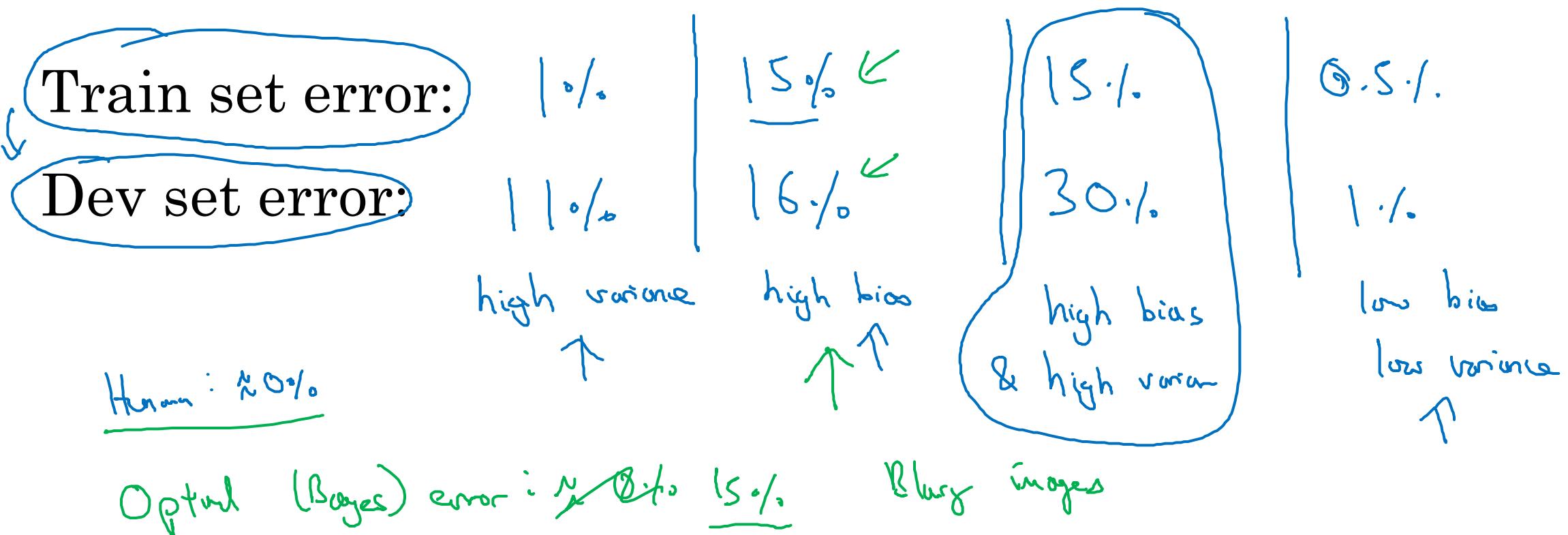
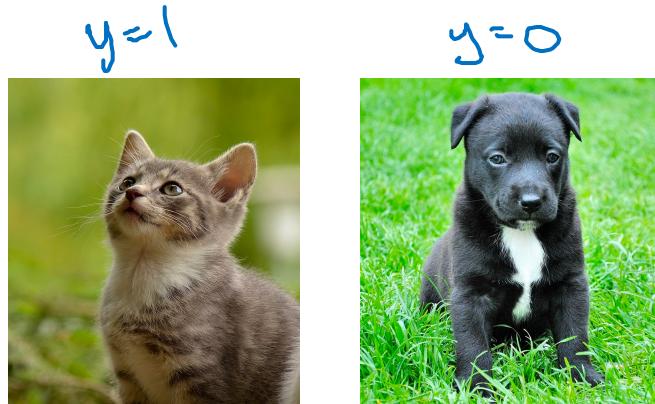
Bias/Variance

Bias and Variance

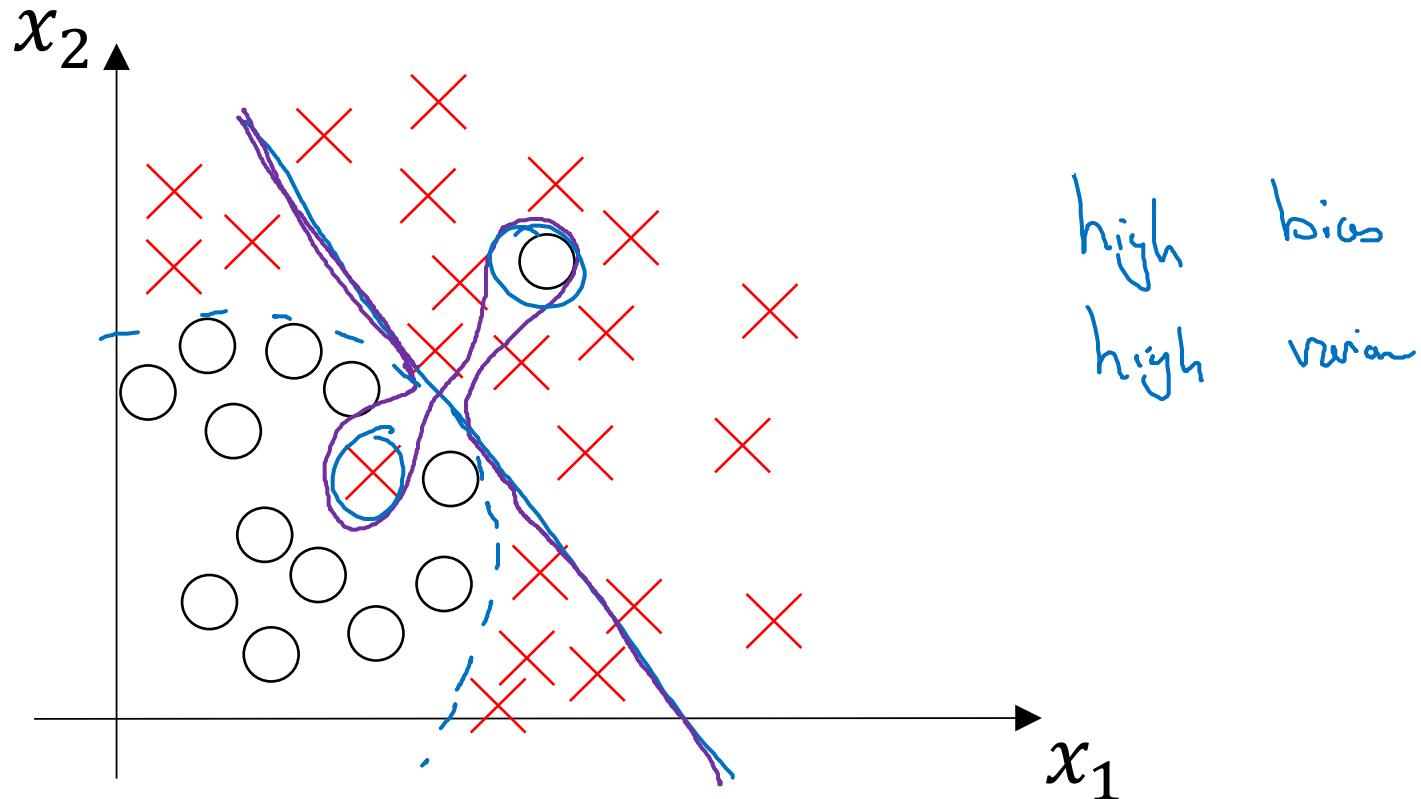


Bias and Variance

Cat classification



High bias and high variance





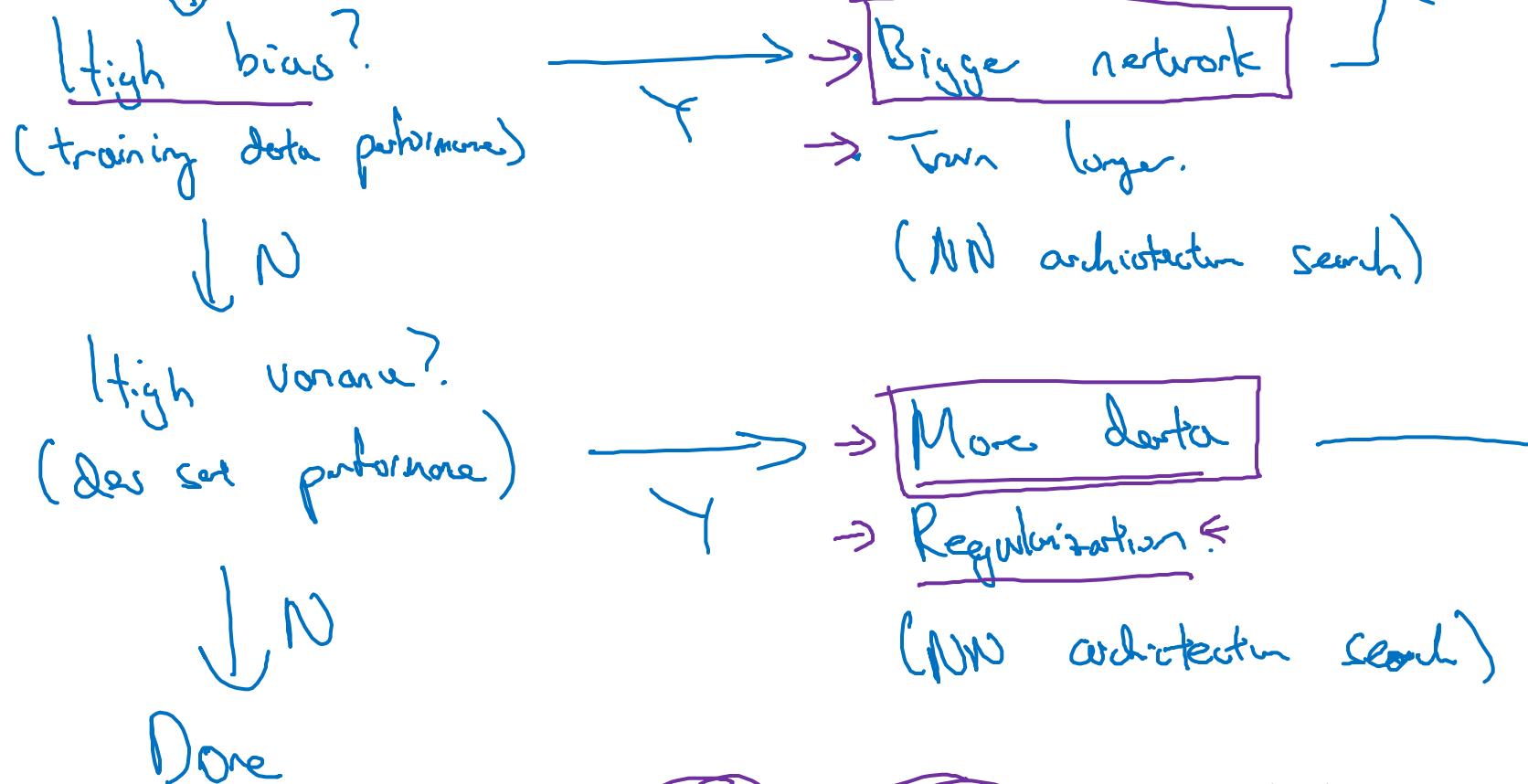
deeplearning.ai

Setting up your
ML application

Basic “recipe”
for machine learning

Basic “recipe” for machine learning

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

$$J(w, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})}_{\text{L1 regularization}} + \frac{\lambda}{2m} \|w\|_2^2$$

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

omit

L₂ regularization

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

L₁ regularization

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

w will be sparse

Neural network

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

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

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

"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"

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

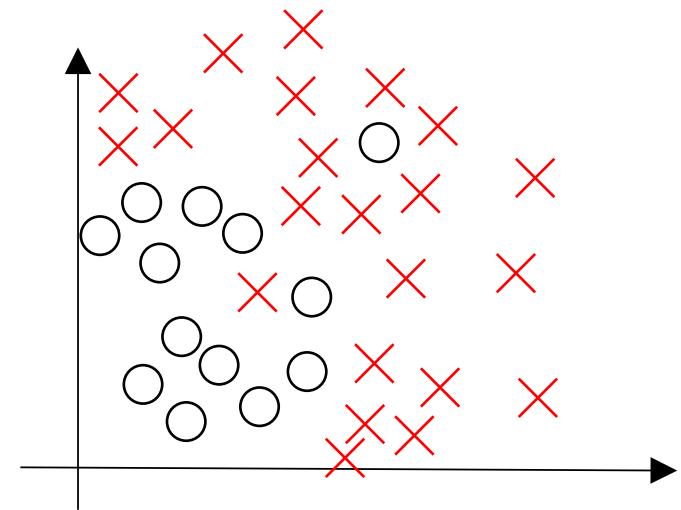
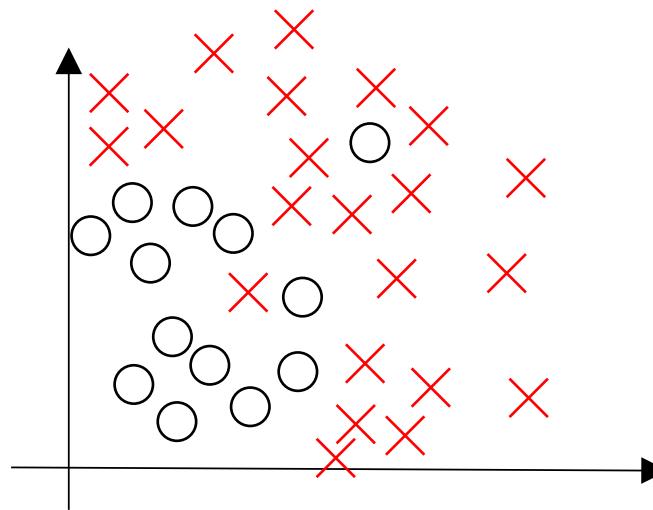
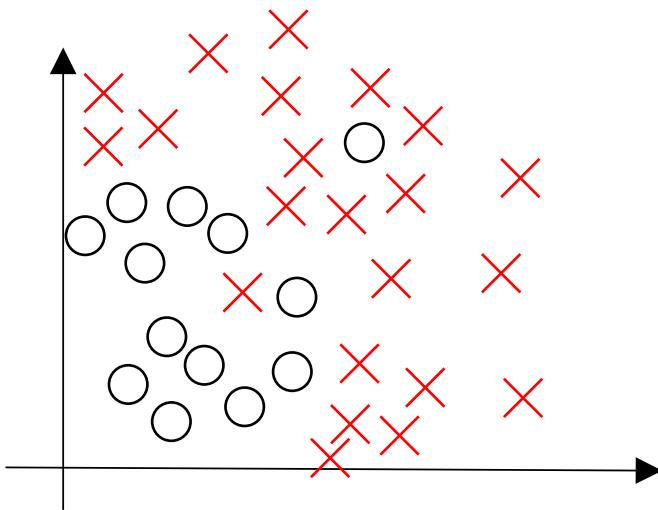
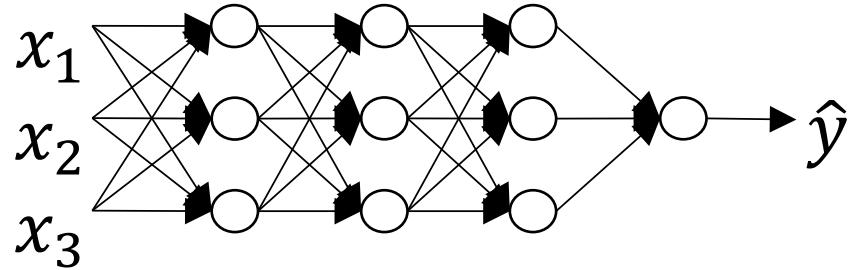
$$= w^{(l)} - \frac{\alpha \lambda}{m} w^{(l)} - \alpha (\text{from backprop})$$

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

Neural network

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

How does regularization prevent overfitting?



How does regularization prevent overfitting?

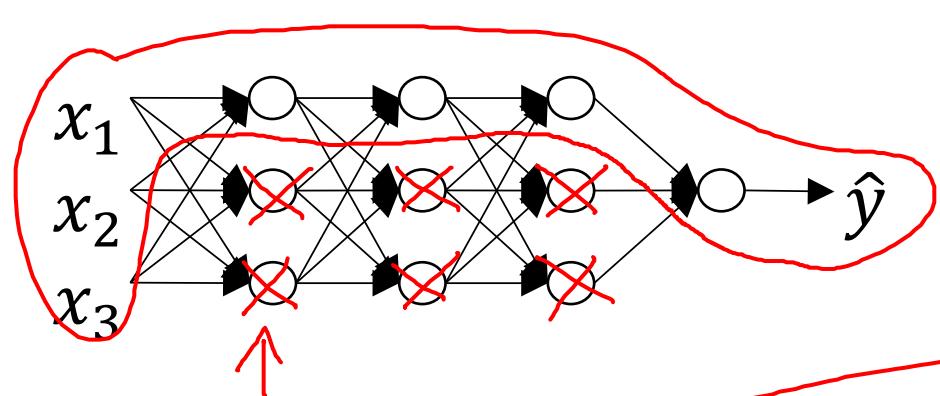


deeplearning.ai

Regularizing your neural network

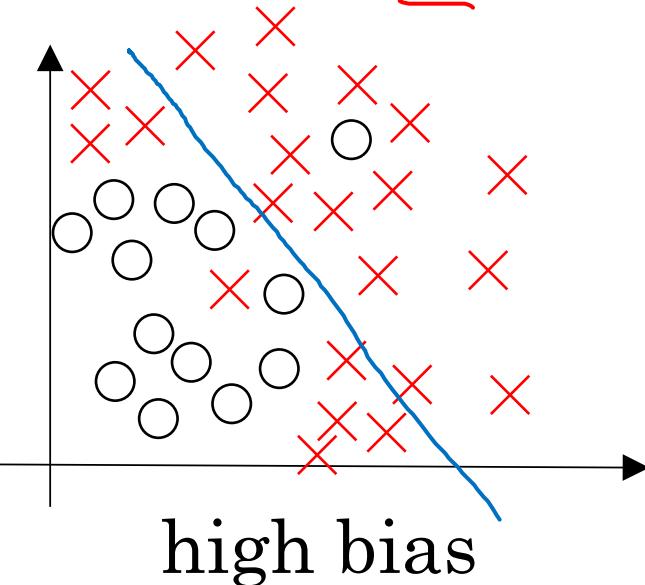
Why regularization reduces overfitting

How does regularization prevent overfitting?

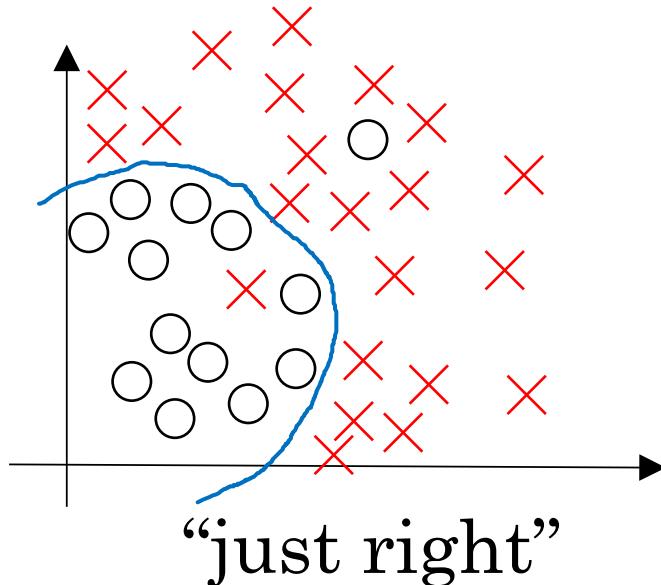


$$J(\boldsymbol{w}^{(1)}, \boldsymbol{b}^{(1)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

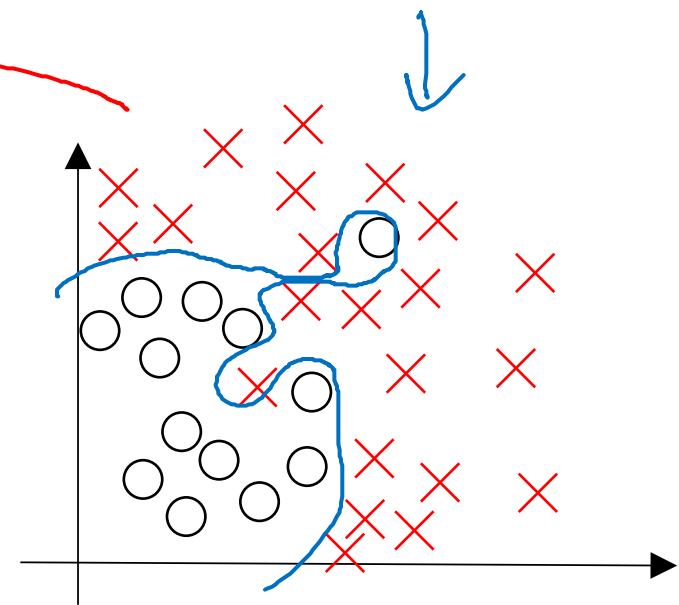
$\boldsymbol{w}^{(1)} \approx 0$



high bias

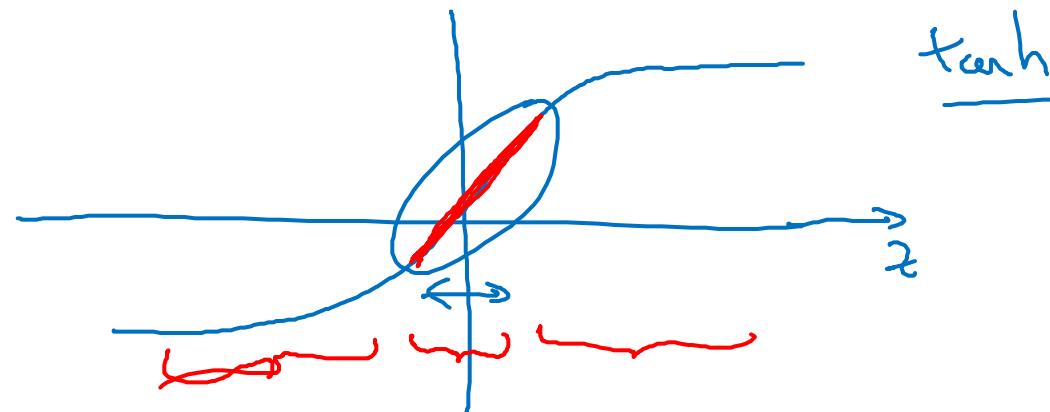


"just right"



high variance

How does regularization prevent overfitting?



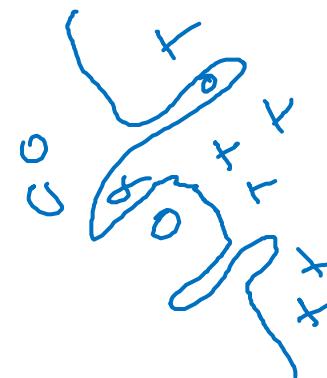
$$\lambda \uparrow$$

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

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

Every layer \approx linear.

$$J(\dots) = \boxed{\sum_i L(\hat{y}^{(i)}, y^{(i)})} + \lambda \sum_{l=2}^L \|\underline{w^{[l]}}\|_F^2$$



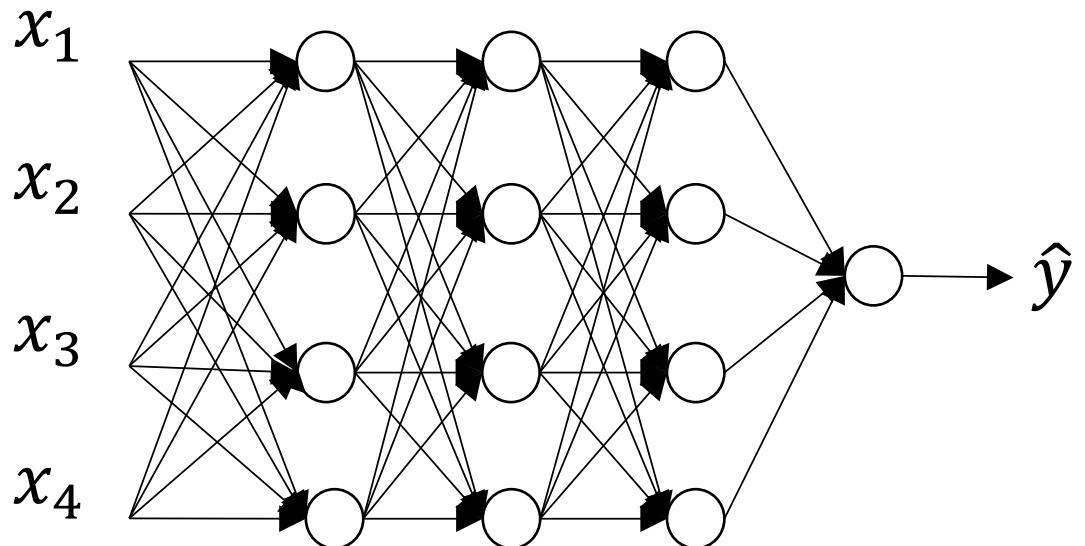


deeplearning.ai

Regularizing your
neural network

Dropout
regularization

Dropout regularization



\uparrow
0.5 \uparrow
0.5 \uparrow
0.5

Implementing dropout (“Inverted dropout”)

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

$\rightarrow d^3 = \underbrace{\text{np.random.rand}(a^3.\text{shape}[0], a^3.\text{shape}[1]) < \text{keep-prob}}$

a^3 = $\text{np.multiply}(a^3, d^3)$ $\# a^3 * d^3$.

$\rightarrow a^3 /= \cancel{\text{keep-prob}}$ \leftarrow

50 units. \rightsquigarrow 10 units shut off

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

x reduced by 20%.

Test

$$1 = \underline{0.8}$$

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



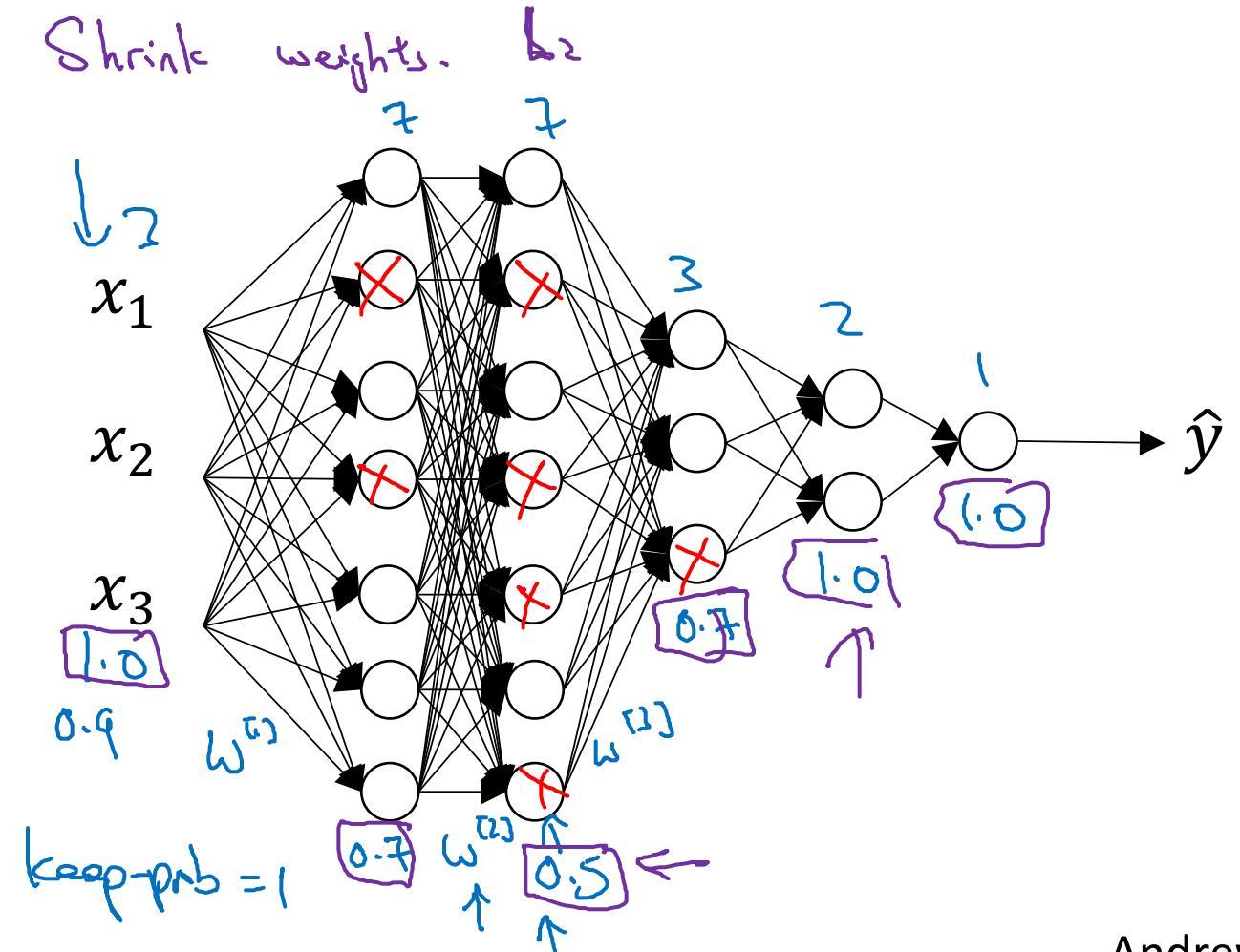
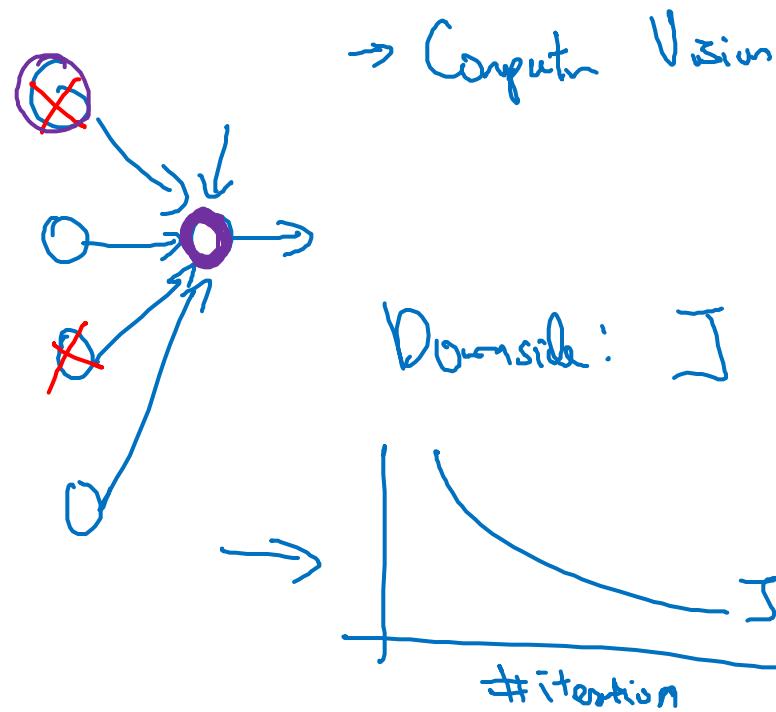
deeplearning.ai

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. b_2





deeplearning.ai

Regularizing your neural network

Other regularization methods

Data augmentation



4

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

4

A black silhouette of the digit '4' with a wavy, hand-drawn style.

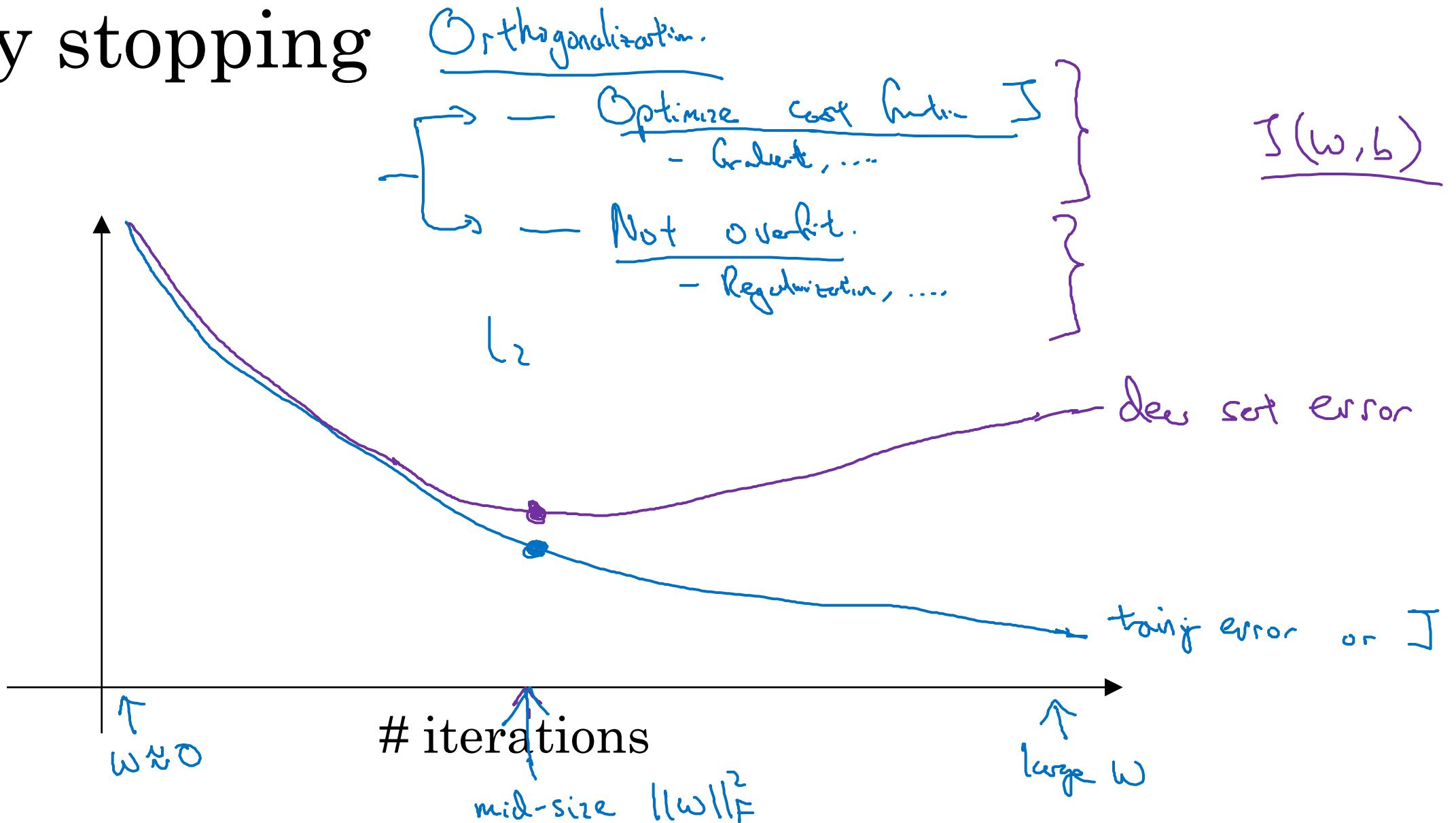
4

A black silhouette of the digit '4' with a more standard, clean font style.

4

A black silhouette of the digit '4' with a wavy, hand-drawn style, enclosed within a white rectangular box.

Early stopping





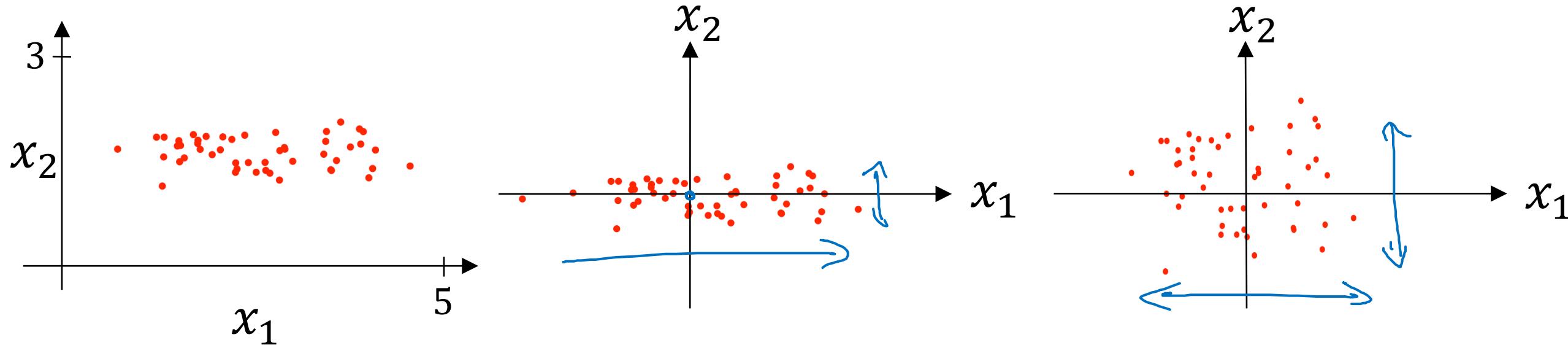
deeplearning.ai

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

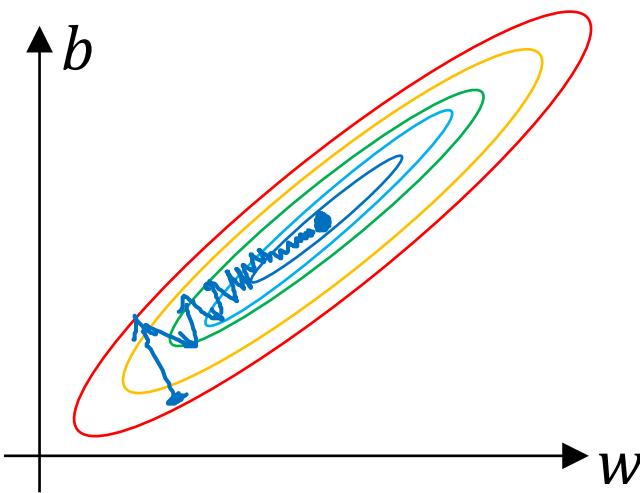
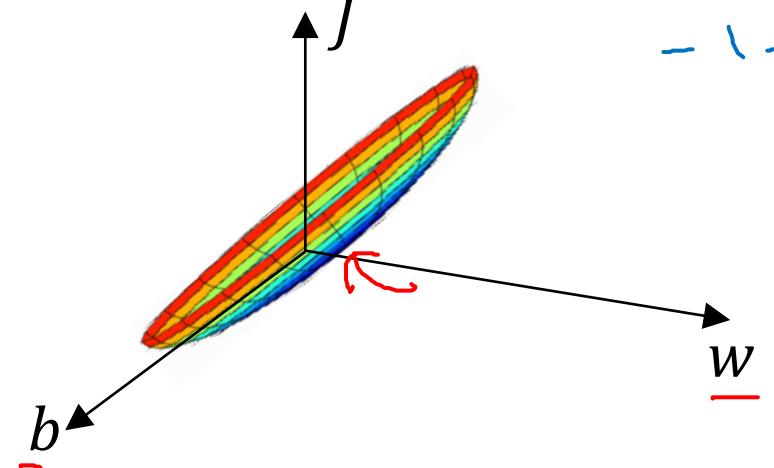
$$\underline{x / \sigma^2}$$

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

Why normalize inputs?

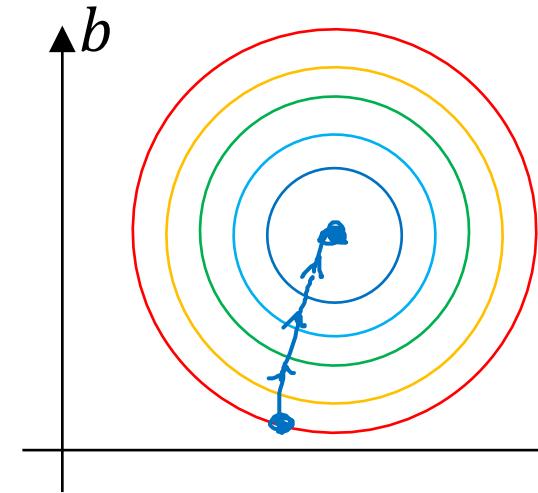
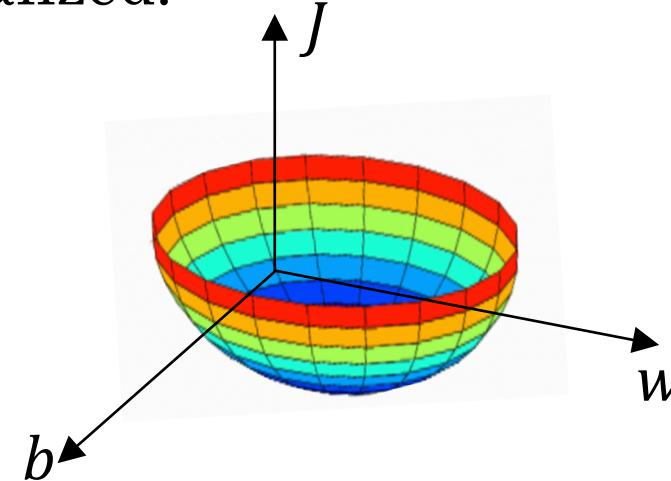
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:
 ω_1 $x_1: \underline{1...100}$ ←
 ω_2 $x_2: \underline{0...1}$ ←
- ... -



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

Normalized:





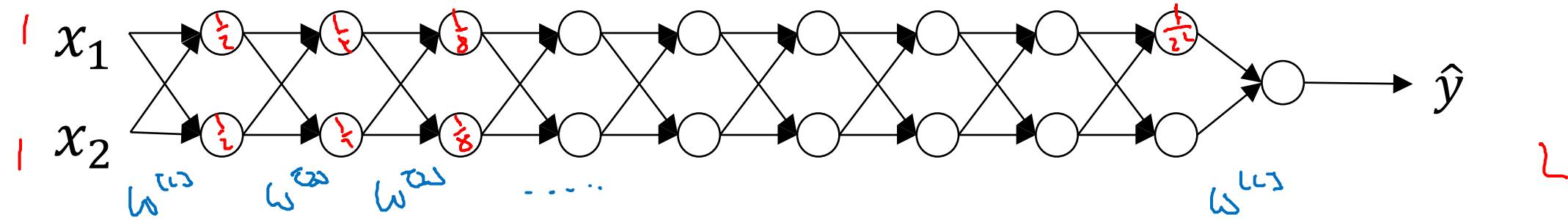
deeplearning.ai

Setting up your
optimization problem

Vanishing/exploding
gradients

Vanishing/exploding gradients

$L=150$



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



$$w^{(1)} > I$$

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

$$w^{(2)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 6.5 \end{bmatrix}$$

$$z^{(1)} = \underline{w^{(1)} x}$$

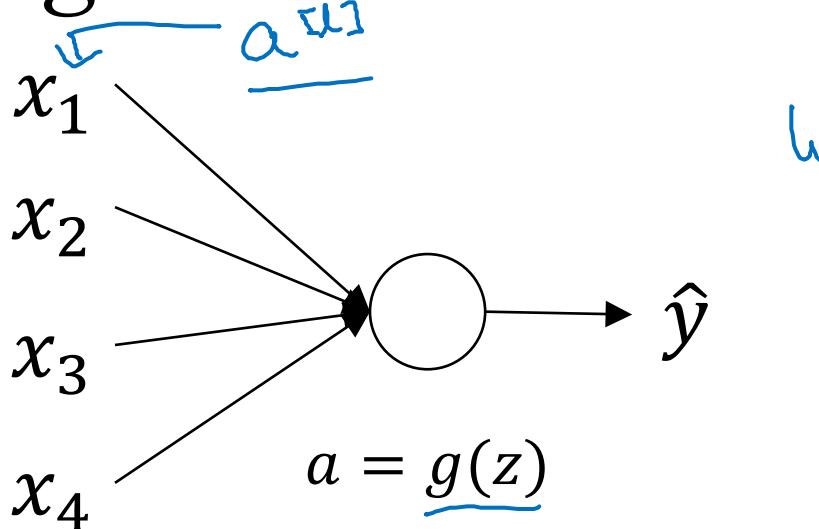
$$a^{(1)} = g(z^{(1)}) = z^{(1)}$$

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

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

$$1.5^{L-1} x \\ 6.5^{L-1} x$$

Single neuron example



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

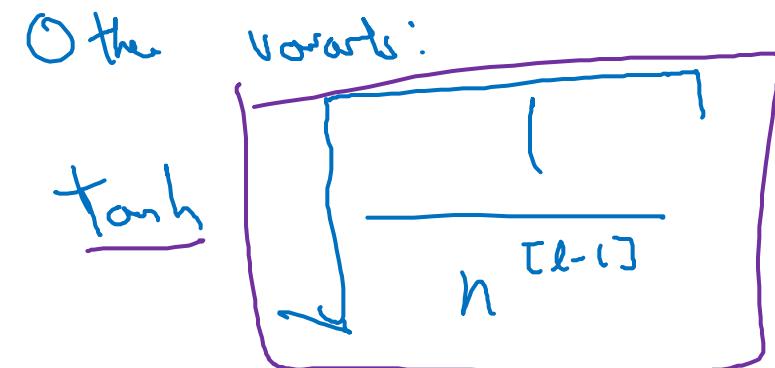
Large $n \rightarrow$ Smaller w_i

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

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

ReLU

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



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

↑



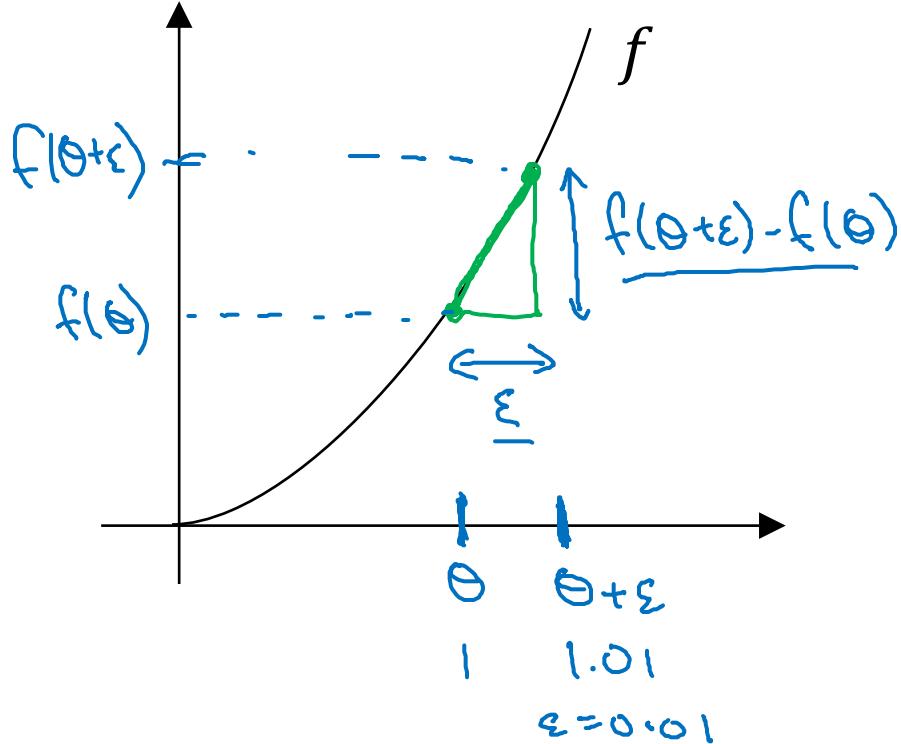
deeplearning.ai

Setting up your optimization problem

Numerical approximation of gradients

Checking your derivative computation

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



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

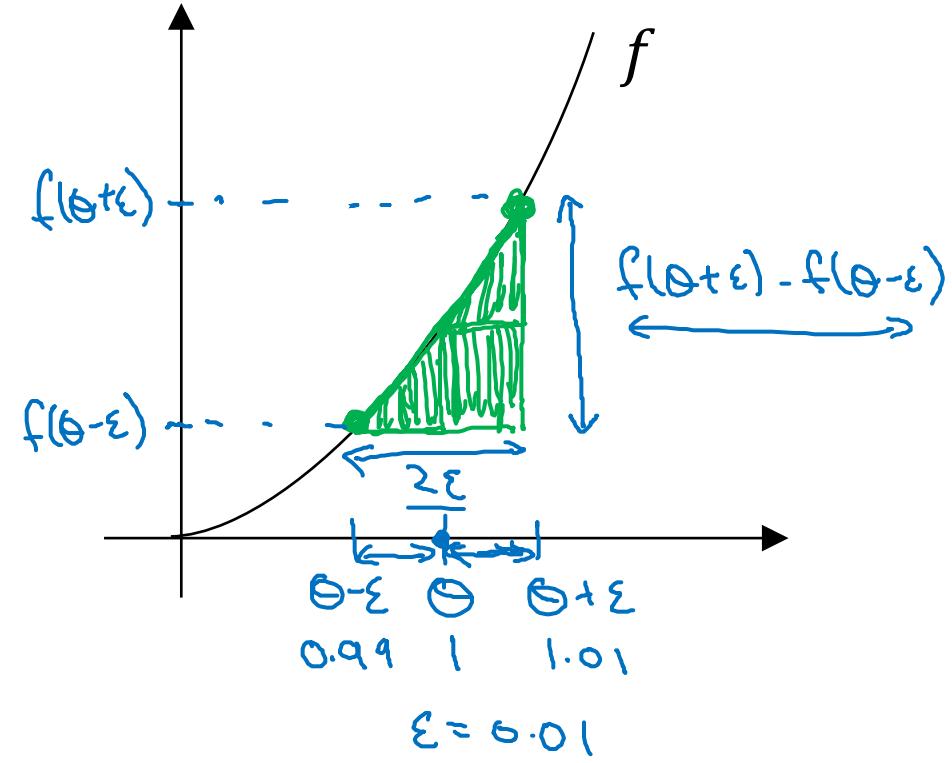
$$\begin{aligned} \frac{f(\theta+\epsilon) - f(\theta)}{\epsilon} &\approx g(\theta) \\ \frac{(1.01)^3 - 1^3}{0.01} &= \frac{3.0301}{0.01} \\ &\approx 3 \end{aligned}$$

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

$$\begin{aligned} 3.0301 \\ 3.1 \\ 3.2 \end{aligned}$$

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)

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



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 $\underline{\theta}$.

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

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 $\mathcal{J}(\theta)$?

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

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

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

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

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



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.

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

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

- Doesn't work with dropout.

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

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

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