



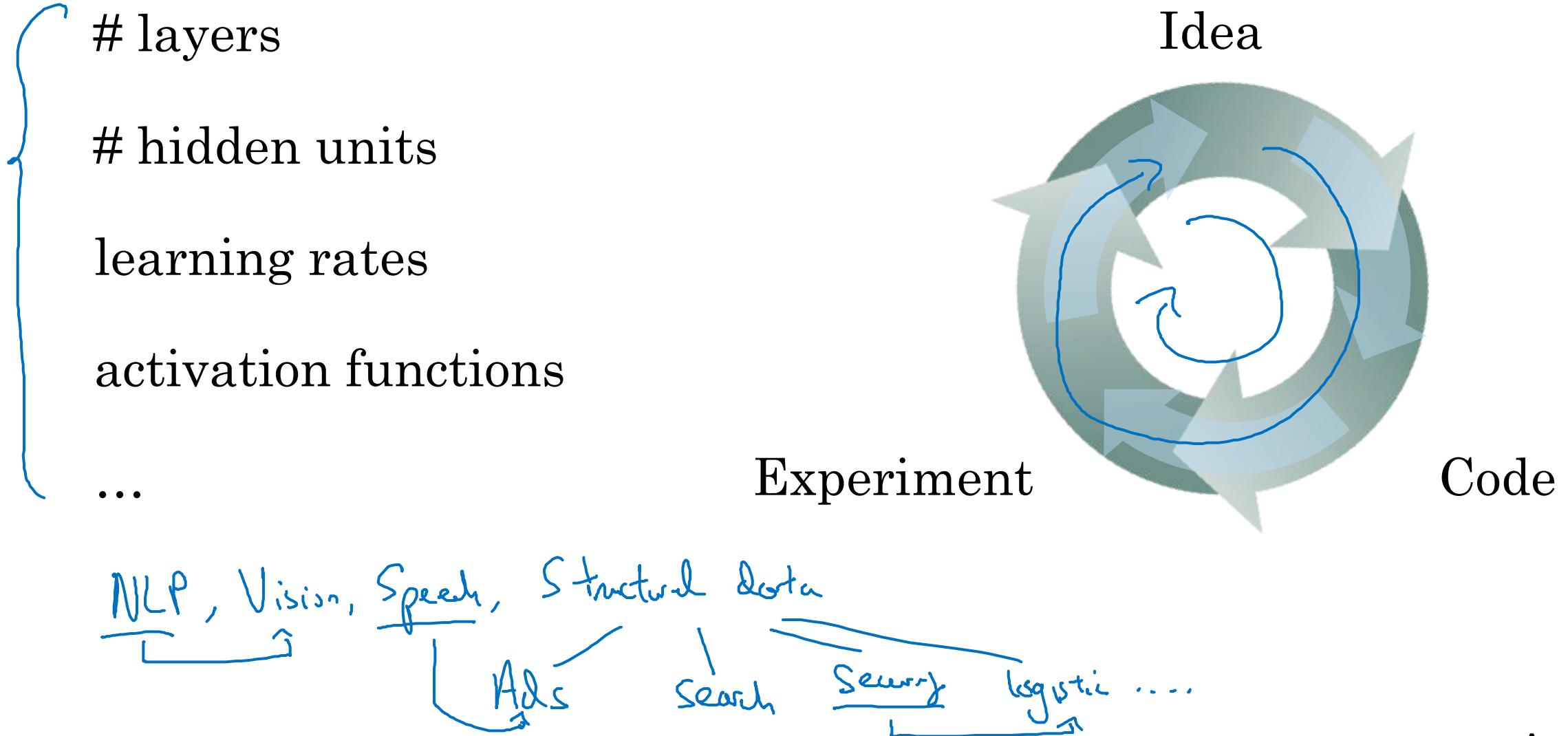
deeplearning.ai

Setting up your  
ML application

---

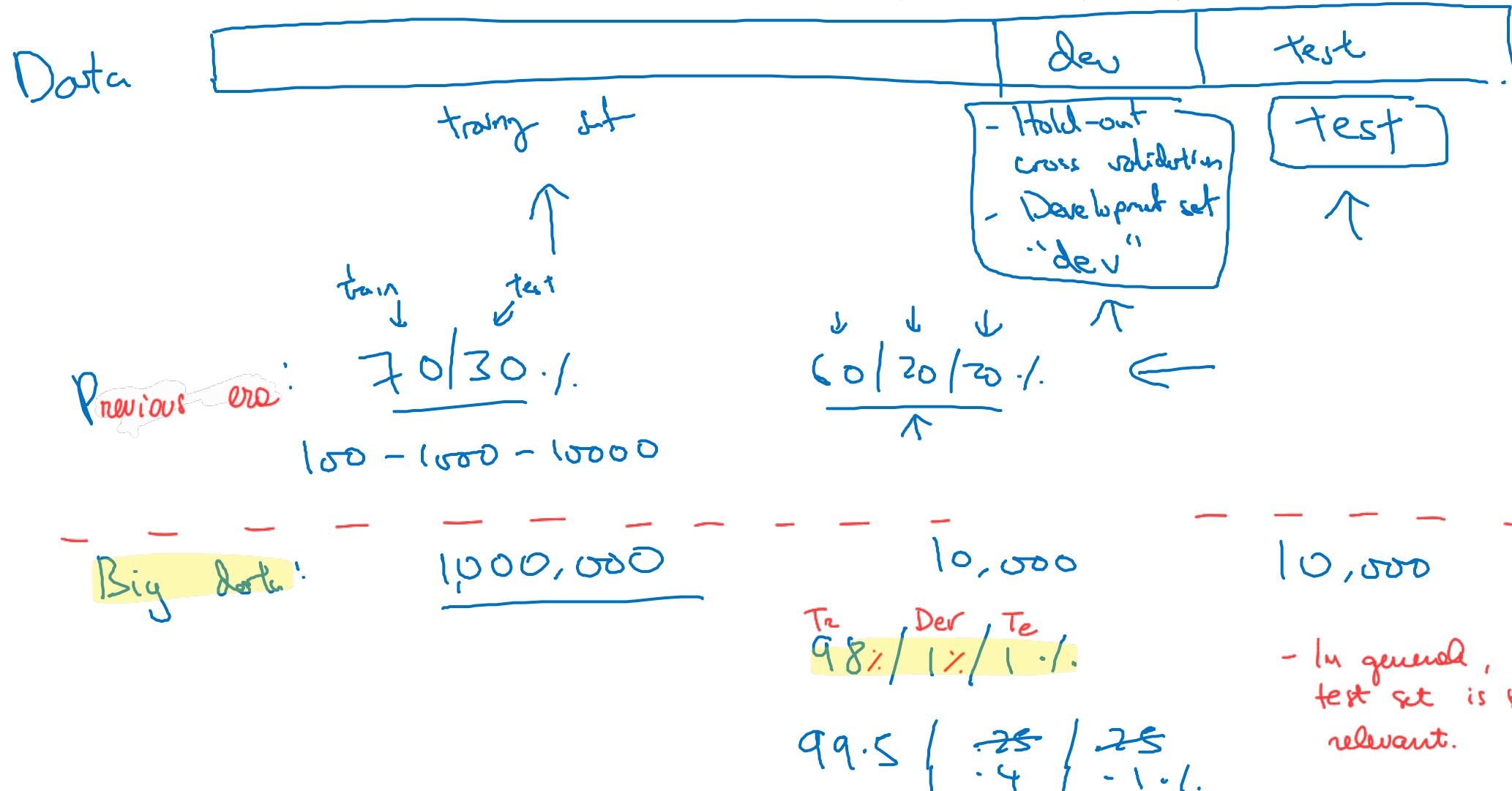
Train/dev/test  
sets

# Applied ML is a highly iterative process



# Train/dev/test sets

- In a traditional "small data" setting a typical train/test split is 70/30%.
- In a big data setting it makes more sense to have the test & development (= validation) sets much smaller  $\leq 1\%$  each.



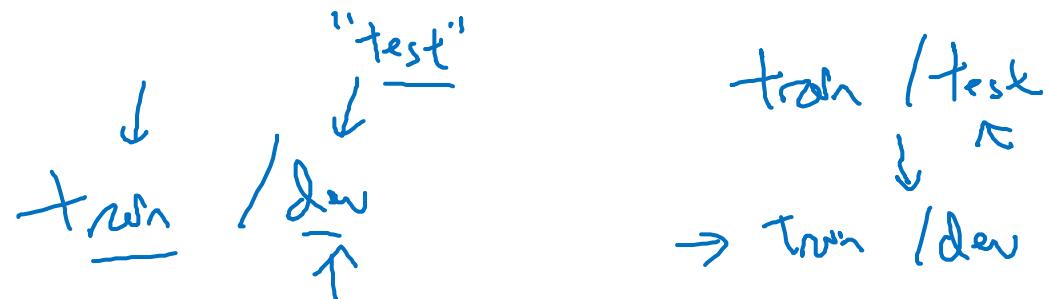
- In general, in big data the test set is the least relevant.

# Mismatched train/test distribution

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.)

IF ONE DOESN'T NEED AN UNBIASED ESTIMATE OF THE MODEL!!

↳ Sounds like BS to me ...

- It is common to have non i.d. TRAIN & TEST sets. This is fine as long as the validation set is i.d. to the test set (because the network is optimized, using the dev set, to work on the target distribution).

- WEIRD. Some groups not always use a separate validation (dev) and test set. They only have train & dev (called "hold out dev set" or, increasingly, test set) and to train & test on those, although this is improper because they are overfitting on the test set.



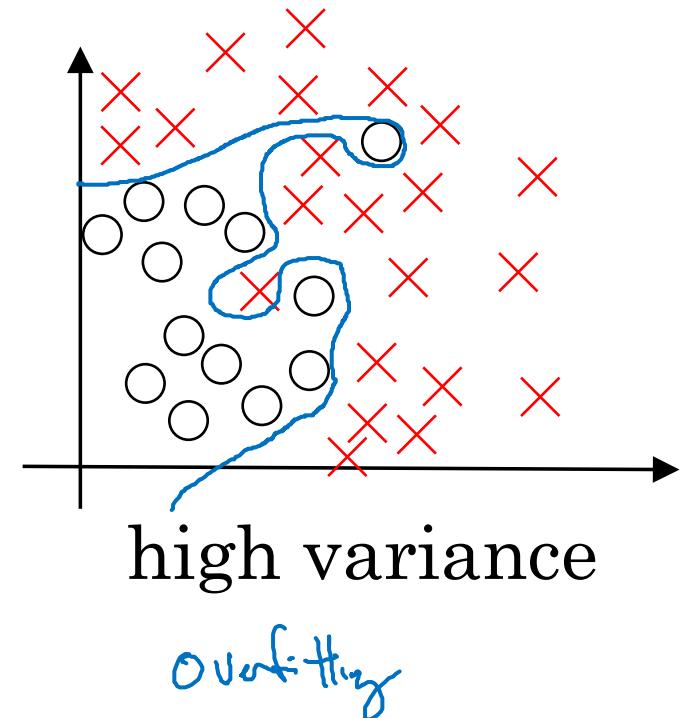
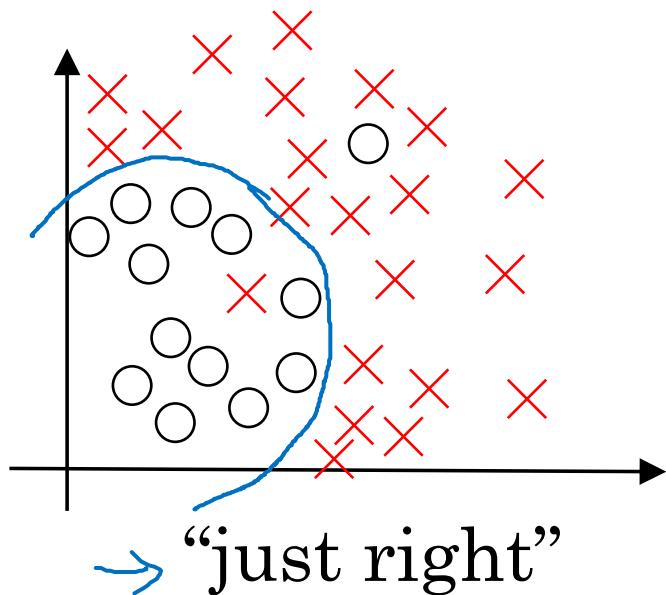
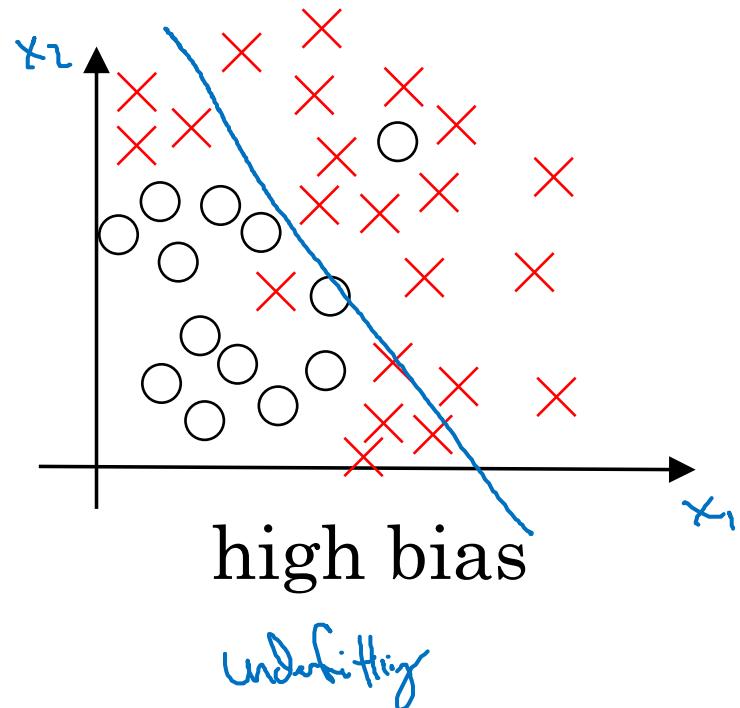
deeplearning.ai

Setting up your  
ML application

---

Bias/Variance

# Bias and Variance



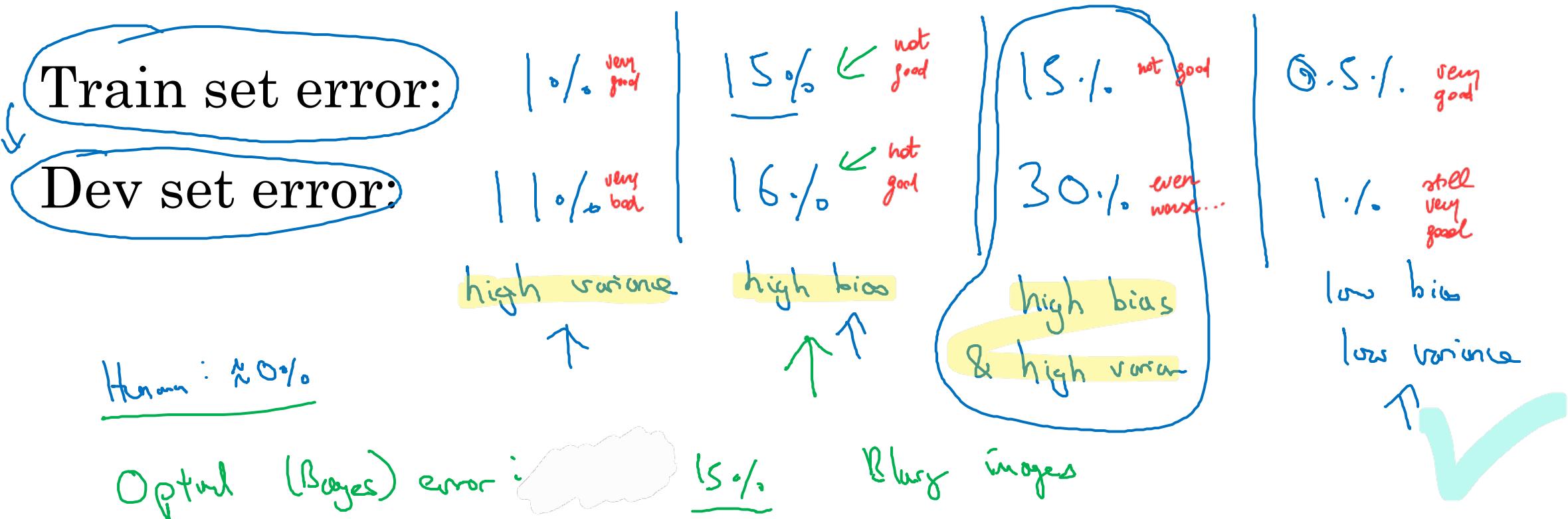
# Bias and Variance

## Cat classification

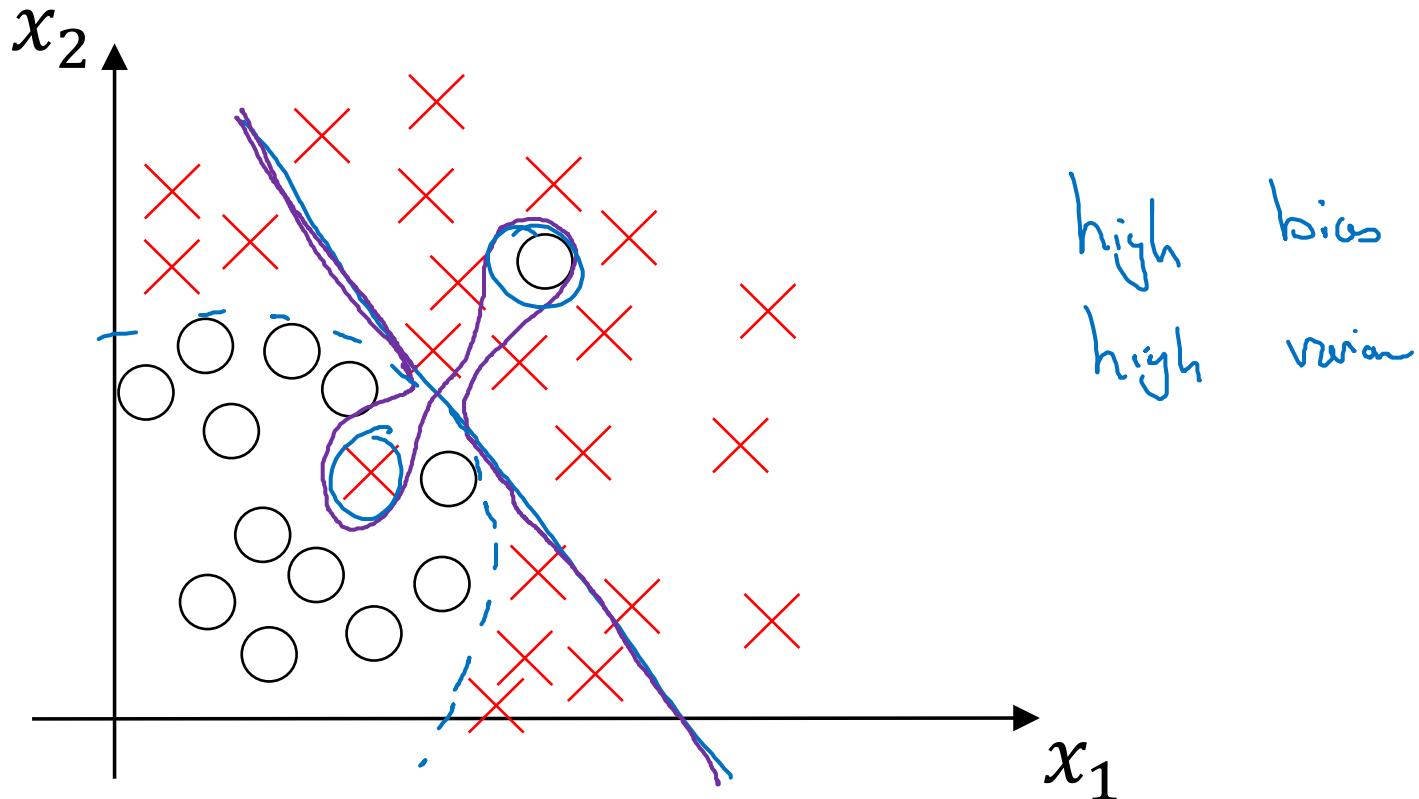


→ Metrics to evaluate **BIAIS & VARIANCE** are the TRAIN & DEV SET ERRORS.

→ We want both low train error and dev set error.



# High bias and high variance





deeplearning.ai

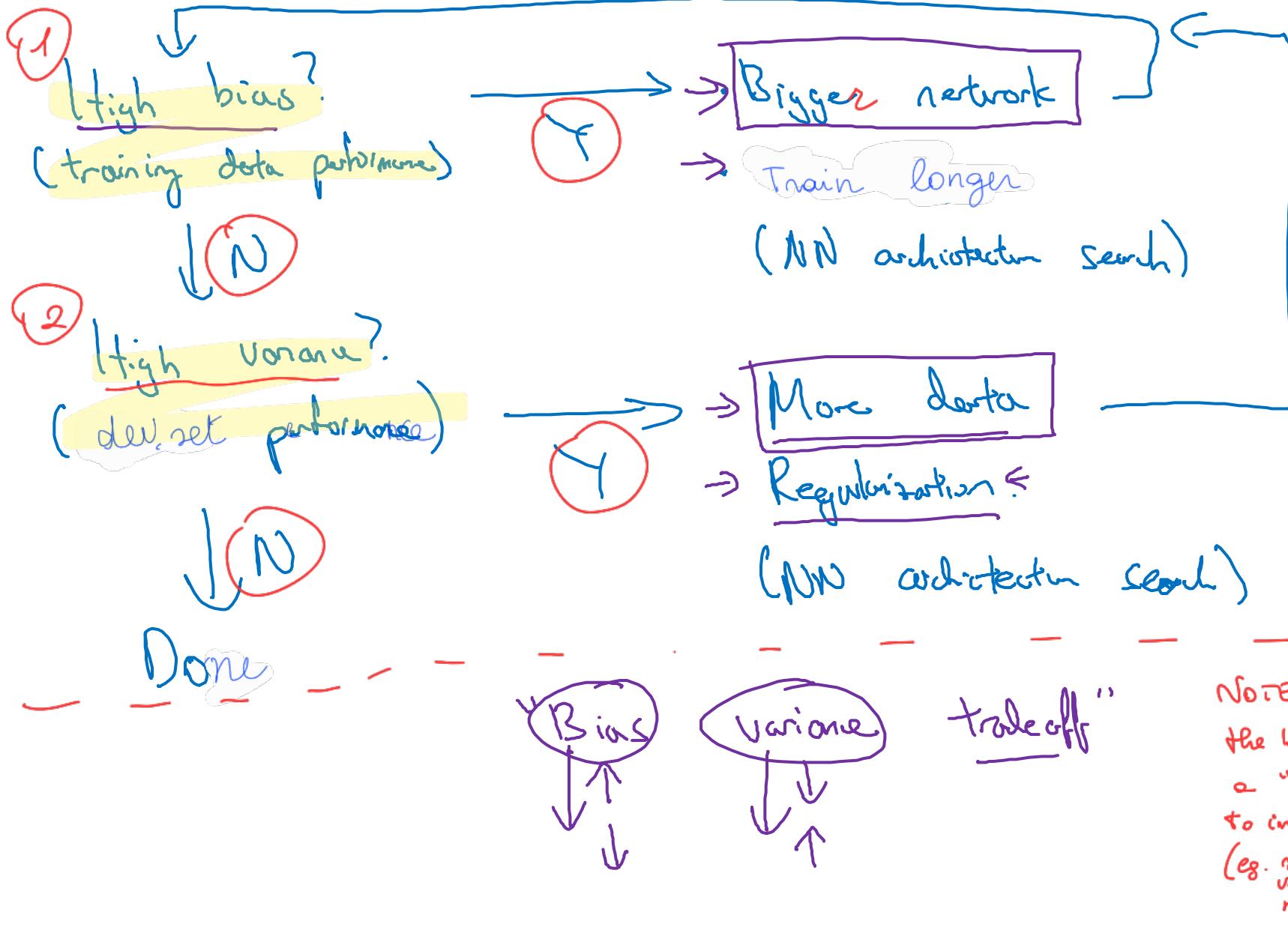
Setting up your  
ML application

---

Basic “recipe”  
for machine learning

# Basic recipe for machine learning

How to reduce both  
BIAS AND VARIANCE



NOTE: for deep learning applications,  
the bias/variance tradeoff is less of  
a "tradeoff" because there are options  
to improve (reduce) one without hurting the other.  
(eg. more depth  $\rightarrow$  reduces bias but does not impact  
variance that much.  
more data  $\rightarrow$  reduces variance without impacting  
bias (not much))

# Andrew Ng



deeplearning.ai

# Regularizing your neural network

---

## Regularization

Is another way to reduce variance  
(other than increasing the data)

example: regularization in logistic regression

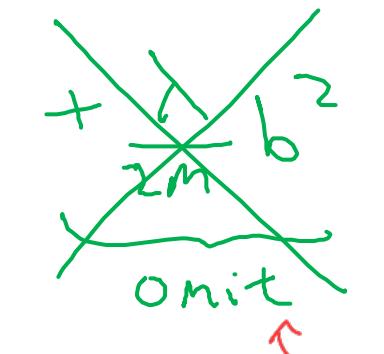
# Logistic regression

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

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

$\lambda$  = regularization parameter  
lambd

$$J(w, b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})}_{\text{L1 regularization}} + \underbrace{\frac{\lambda}{2m} \|w\|_2^2}_{\text{L2 regularization}}$$



often the bias is left unregularized because it is only a scalar, while  $w$  can be high dimensional and contribute more to the variance.

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

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

$w$  will be sparse

# Neural network

- If we apply the Frob. regularization to the weights (new term in the cost function) the gradients of the cost function are : - SAME  $\frac{\partial J}{\partial w^{[l]}}$ ,  $\frac{\partial J}{\partial w^{[l]}} = \frac{\partial J}{\partial w^{[l]}} + \frac{\lambda}{m} w^{[l]}$ .

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

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

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

$$\textcircled{1} \quad \text{"Frobenius norm"} \quad \| \cdot \|_2^2 \quad \| \cdot \|_F^2$$

$$\textcircled{2} \quad \text{GRADIENT DESCENT with regularized weights}$$

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

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

$$\textcircled{3} \quad \text{a.k.a. "Weight decay"} \quad \rightarrow w^{[l]} := w^{[l]} - \alpha \frac{\partial w^{[l]}}{\partial w^{[l]}}$$

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

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

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

intuition for naming  
 "weight decay": at each update the weights are updated by a factor  $< 1$ .



deeplearning.ai

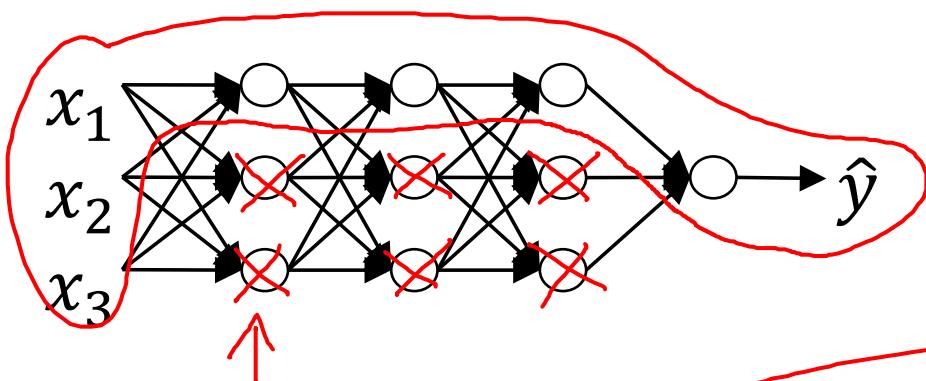
# Regularizing your neural network

---

## Why regularization reduces overfitting

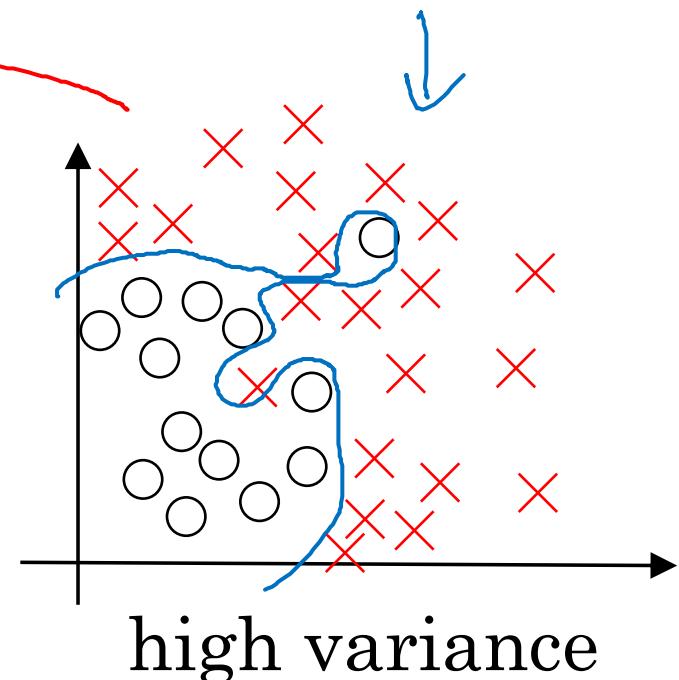
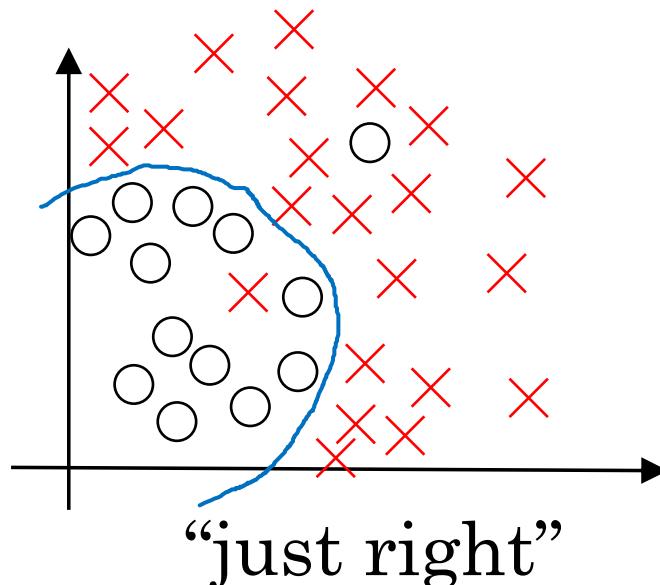
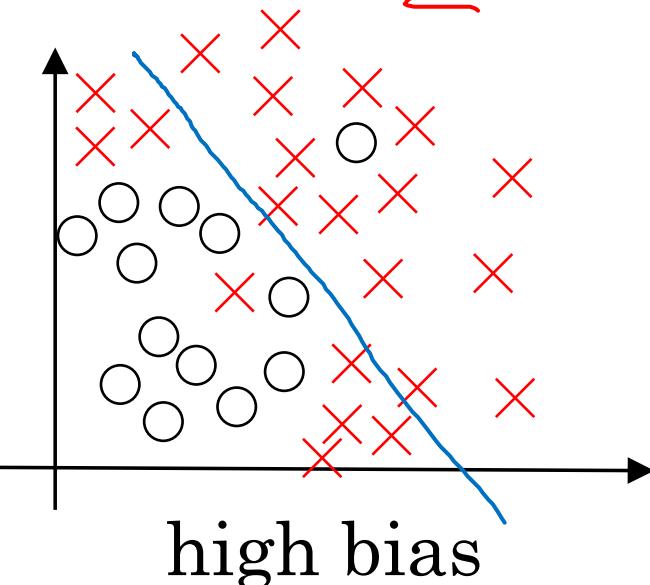
# How does regularization prevent overfitting?

- INTUITION 1: regularization reduces nonlinearity by dampening or killing some neurons.



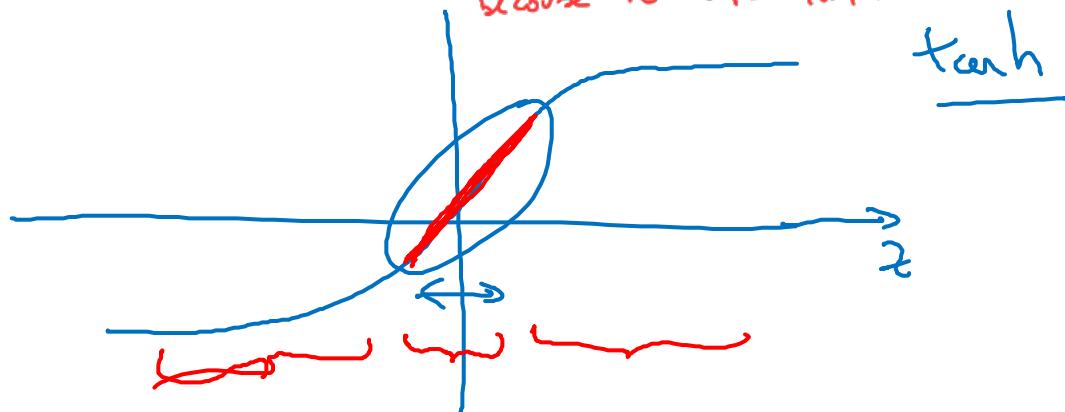
$$J(\boldsymbol{w}^{(1)}, \boldsymbol{b}^{(1)}) = \frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|\boldsymbol{w}^{(l)}\|_F^2$$

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



# How does regularization prevent overfitting?

INTUITION 2: regularization reduces nonlinearity in the case of sigmoidal activations because it keeps  $\|w\|$  low  $\Rightarrow |z|$  low  $\Rightarrow$  sigmoid activations remain in the "linear regime".

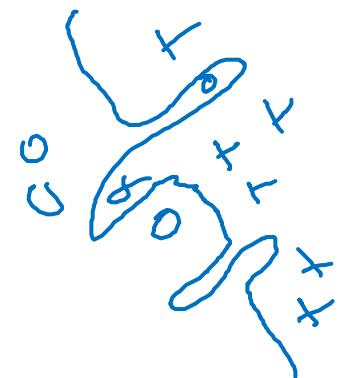


$$\lambda \uparrow$$

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

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

$\rightarrow$  Every layer  $\approx$  linear



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





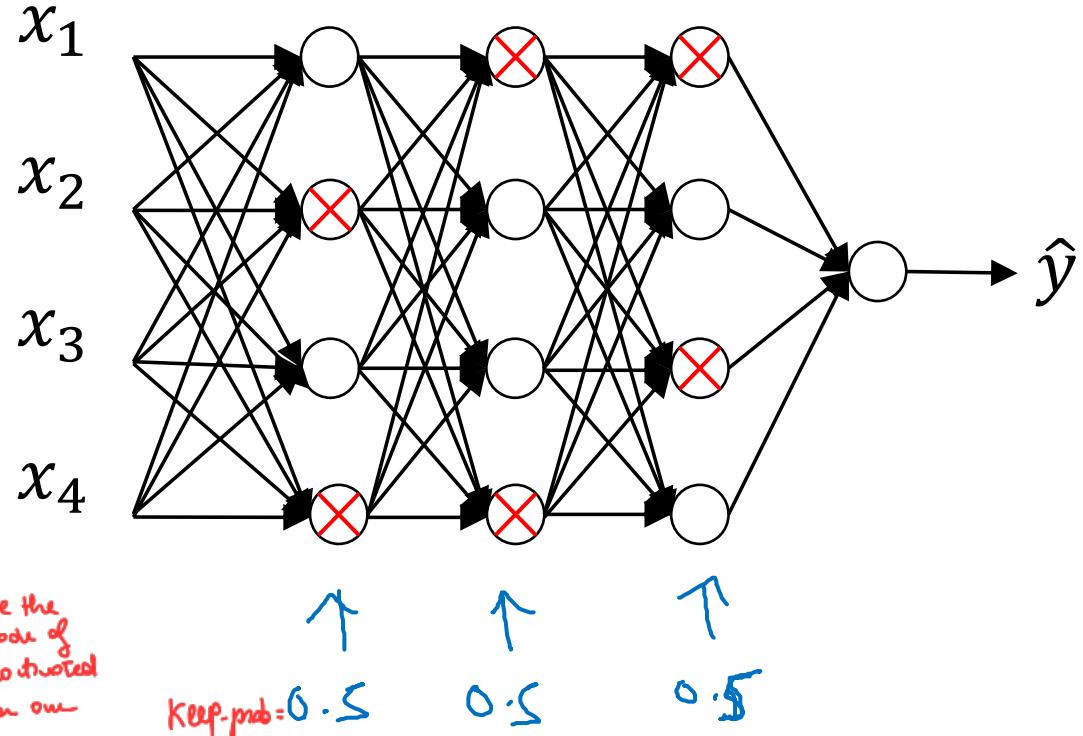
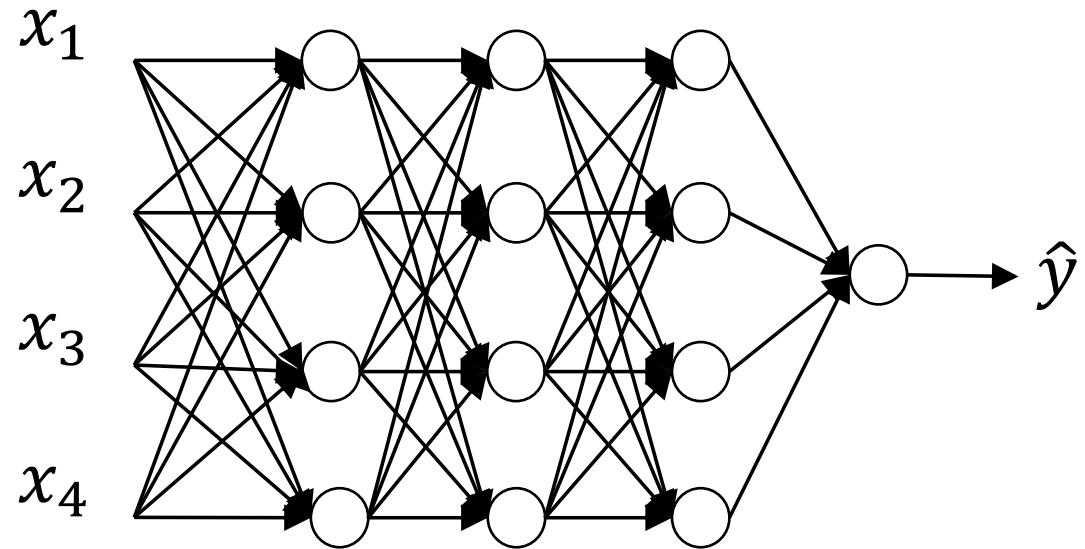
deeplearning.ai

Regularizing your  
neural network

---

Dropout  
regularization

# Dropout regularization



- For each layer, define the probability that any node of the layer could be deactivated at training time (for our minibatch).

- At training time for each sample/minibatch randomly select nodes to deactivate, conduct model update, reactivate the nodes and continue with the next minibatch.

# Implementing dropout (“Inverted dropout”)

HENCE USED IN PRACTICE THAN STANDARD DROPOUT

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

0.2

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

$a_3 = \text{np.multiply}(a_3, d_3)$

#  $a_3 * d_3$ .

→  $a_3 /= \text{keep-prob}$

CHARACTERISTIC OF THE INVERTED DROPOUT: during training, the activations of a layer are scaled by the inverse of the keep-prob  $\approx$  to keep the expected value (average magnitude) of the layer's activations  $\approx$  unchanged.

50 units.  $\rightsquigarrow$  10 units shut off

$$z^{[4]} = w^{[4]} \cdot \frac{a^{[3]}}{\text{keep-prob}} + b^{[4]}$$

$\mathcal{J}$  reduced by  $20\%$ .

$$\text{keep-prob} = \underline{0.8}$$

Test

This technique reduces scaling problems at test time. (If we didn't scale at training time, we should add an extra scaling at test time).

# Making predictions at test time

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

→ Dropout is NOT used at test time.

No drop out.

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

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

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

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

$$\downarrow \hat{y}$$

~~$\lambda$  = keep-prob~~  
no scaling is needed  
at test time if  
we used inverted  
dropout.



deeplearning.ai

Regularizing your  
neural network

---

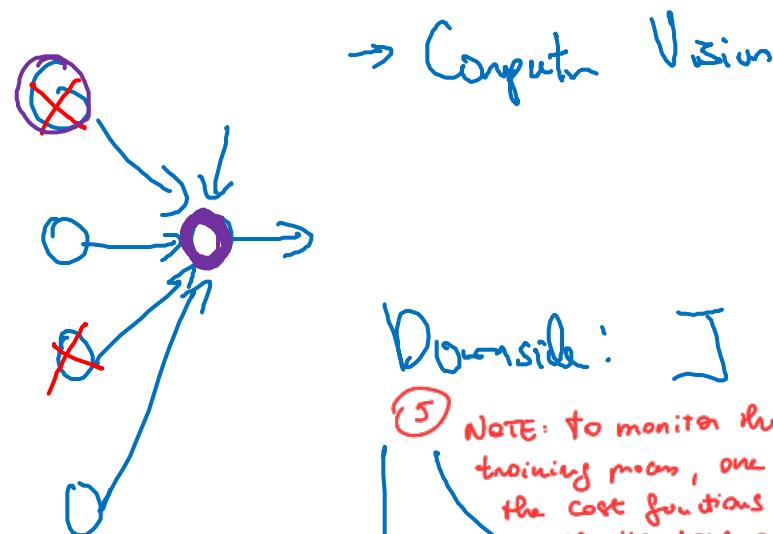
Understanding  
dropout

# Why does drop-out work?

① → "DROPOUT SPREADS OUT THE WEIGHTS"

=> intuitively, by spreading the weights it shrinks the weights (it prevents that any weight becomes too big)  
=> it has a regularizing effect

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

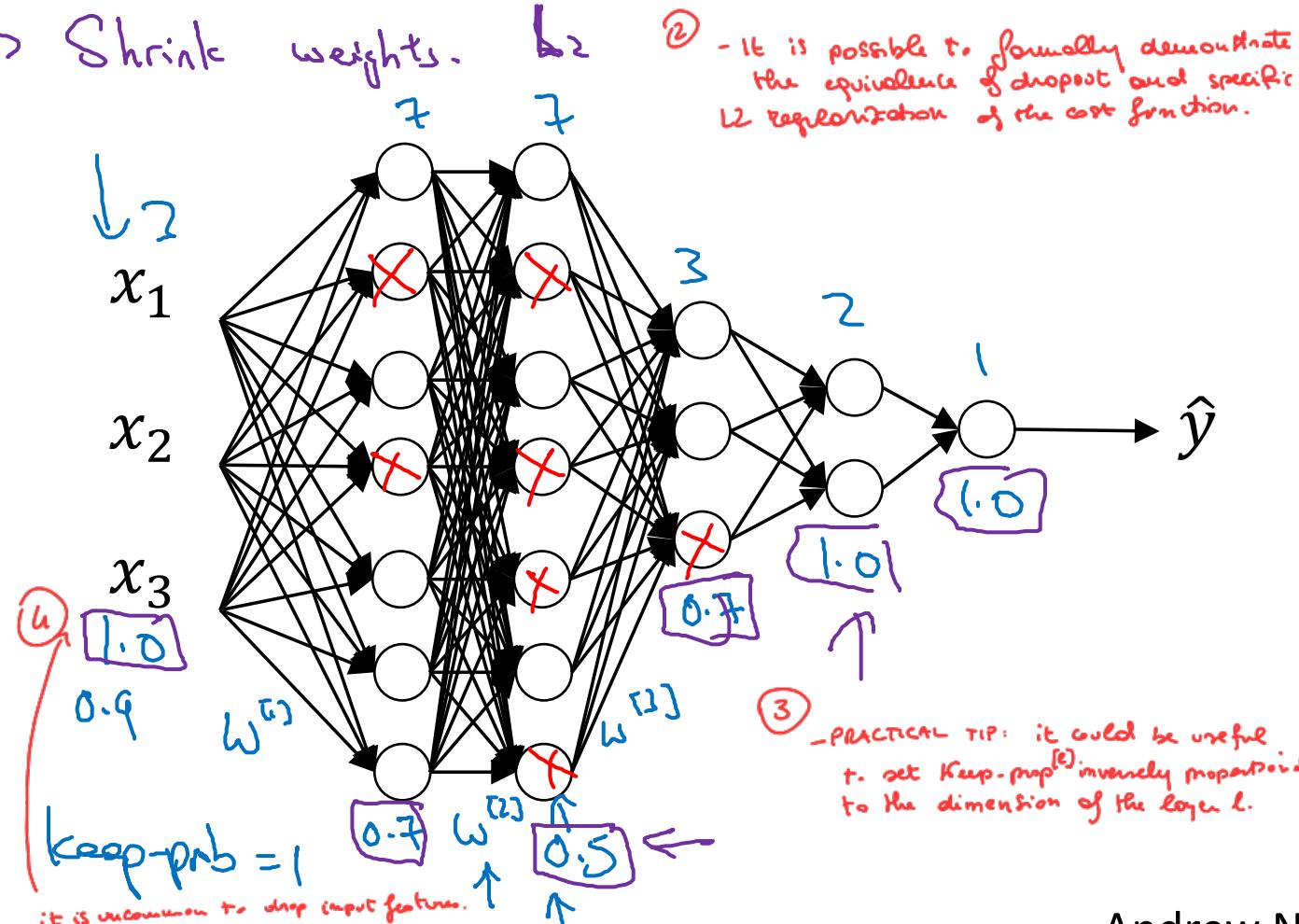


Downside: ↴

⑤ NOTE: to monitor the training process, one plots the cost function over the iterations and checks that it is monotonically decreasing.

#iteration However, monotonicity may not be observed if dropout is used. The solution is to periodically deactivate dropout, compute/plot the loss, then reactivate it and continue the training.

→ Shrink weights.





deeplearning.ai

# Regularizing your neural network

---

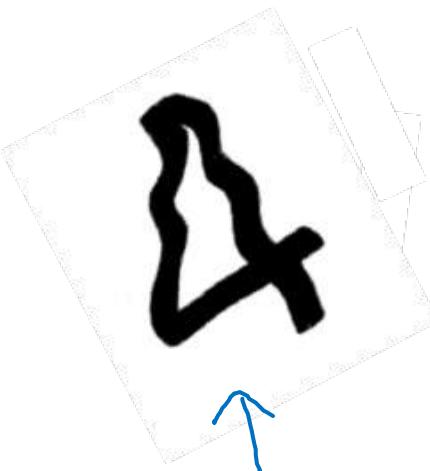
## Other regularization methods

# Data augmentation

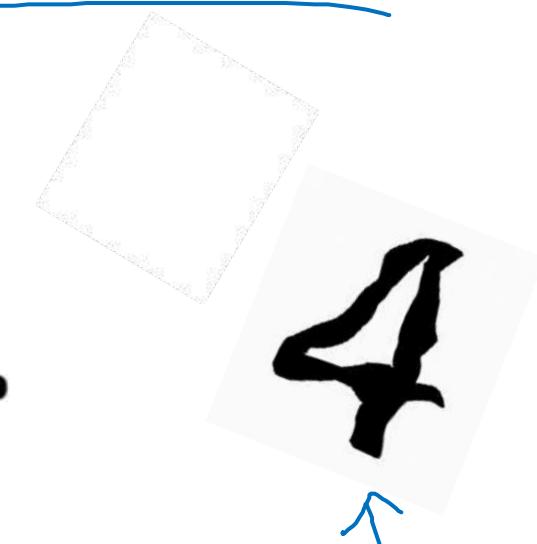
- e.g. for images: rotation, scaling, flipping, zooming



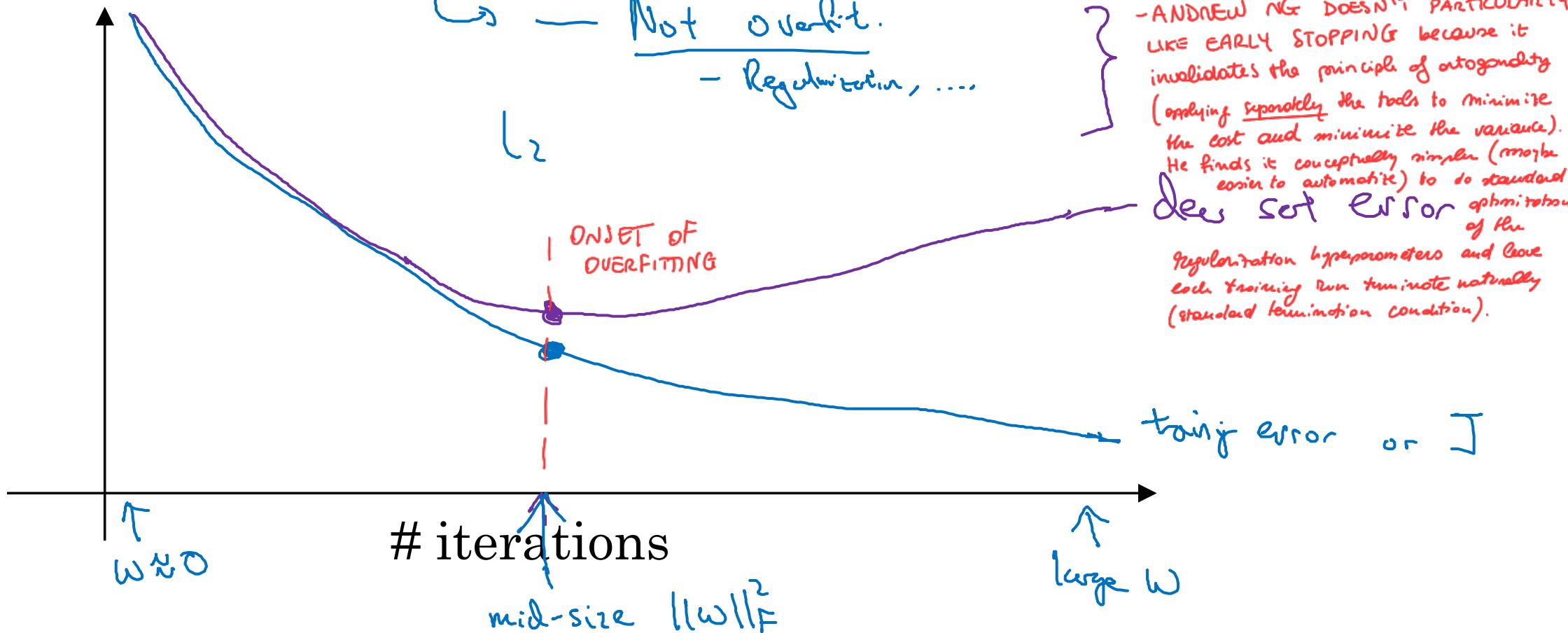
4



4



# Early stopping





deeplearning.ai

Setting up your  
optimization problem

---

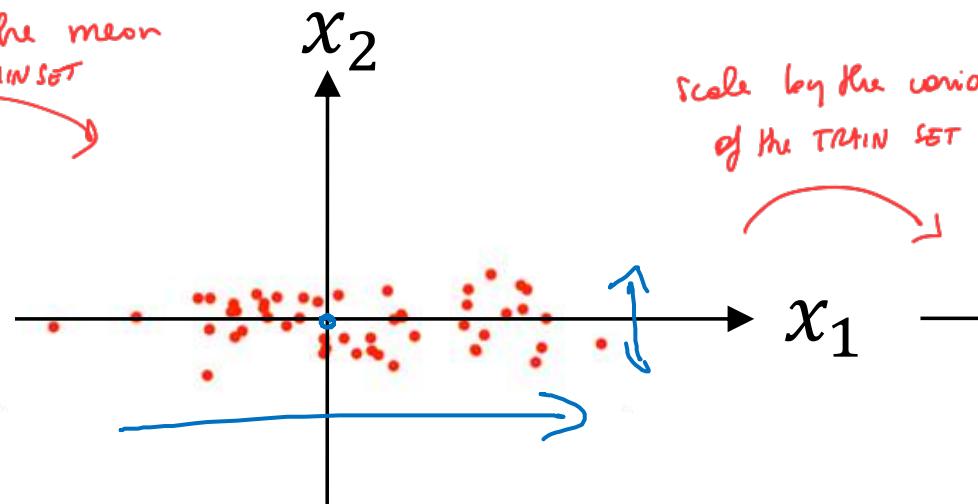
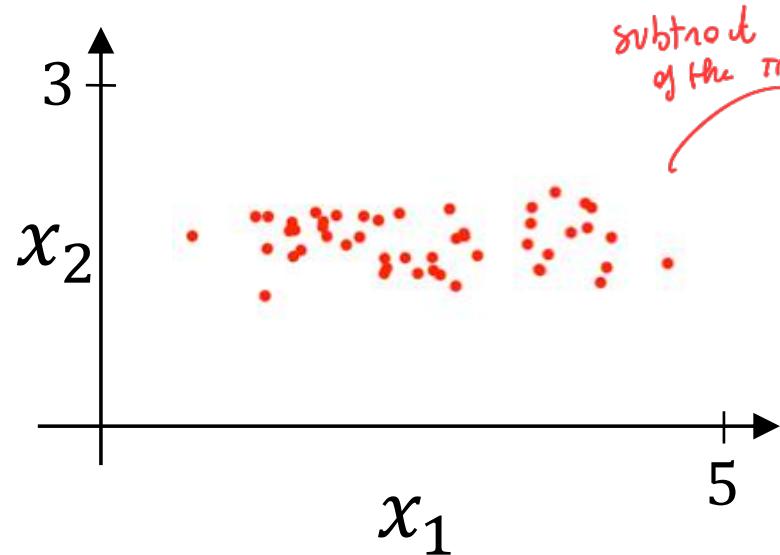
Normalizing inputs

# Normalizing training sets

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

STANDARDIZING

subtract the mean  
of the TRAIN set



- Technique to speed up training
- Scale the training set so it has  $\mu = 0$ ,  $\sigma^2 = 1$ .
- obtained by computing  $\mu_{\text{TRAIN}}$ ,  $\sigma_{\text{TRAIN}}$  and scaling BOTH training AND test sets by then ( $\frac{x_i - \mu_{\text{train}}}{\sigma_{\text{train}}}$ ).
- NOTE! This operation removes the magnitude information from the data and it can be used ONLY when the magnitude of the variables is not an interesting information!



Subtract mean:

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

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

Normalize variance

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

element-wise

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

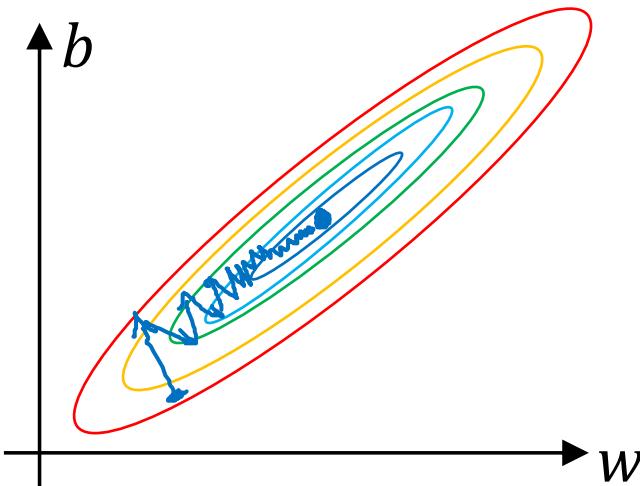
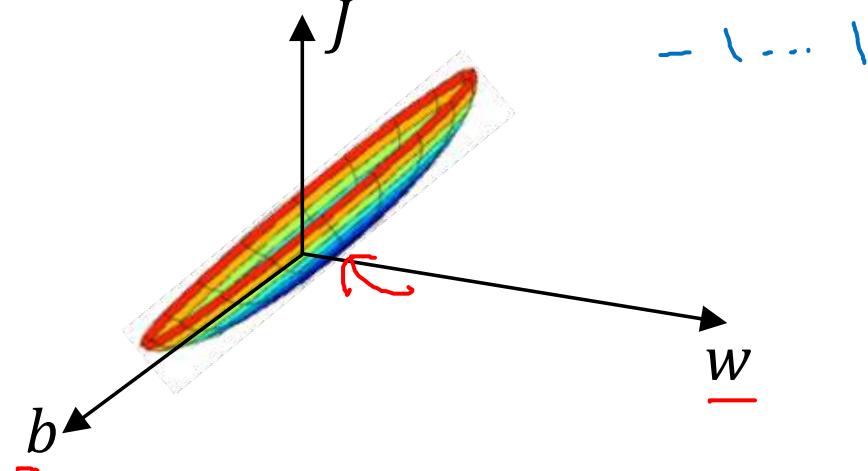
Use same  $\mu, \sigma^2$  to normalize test set.

# Why normalize inputs?

$w_1 \quad x_1: \frac{1 \dots 1000}{\text{...}} \leftarrow$

$w_2 \quad x_2: \frac{0 \dots 1}{\text{...}} \leftarrow$

Unnormalized:



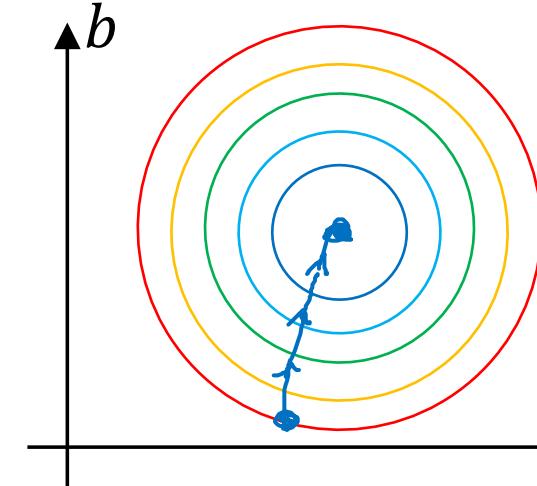
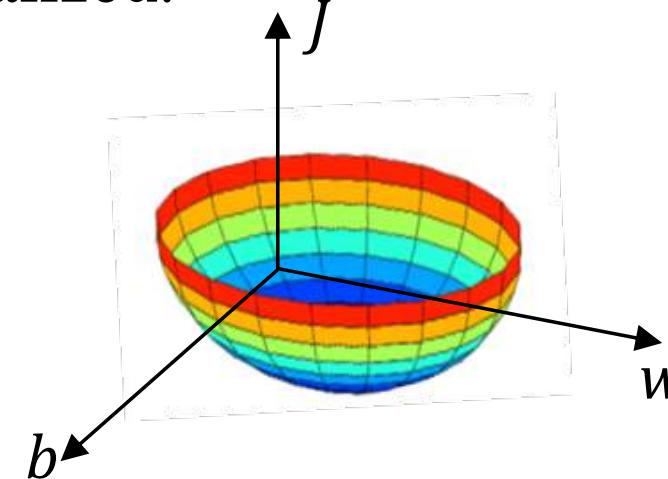
$x_1: 0 \dots 1$

$x_2: -1 \dots 1$

$x_3: 1 \dots 2$

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

Normalized:



*Ideas: the cost function is easier to optimize with GD if all the features have similar ranges.*



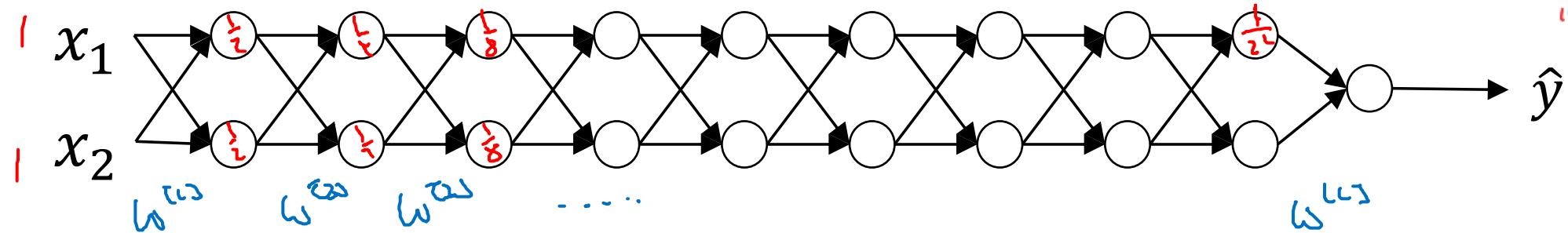
deeplearning.ai

Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

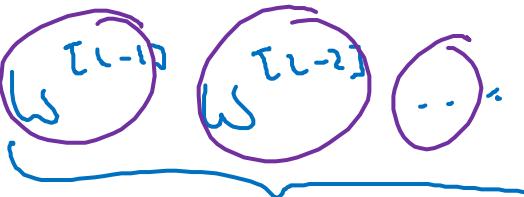
# Vanishing/exploding gradients



assume linear activation and  $b=0$

then the output

$$\hat{y} = \omega^{[L]}$$



$$\omega^{[l]} > 1$$

$$\omega^{[l]} < 1 \quad [0.9 \quad 0.9]$$

$$\omega^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \\ 0 & 0.5 \end{bmatrix}$$

$$\hat{y} = \omega^{[L]} \begin{bmatrix} 0.5 \\ 1.5 \\ 0 \\ 0.5 \end{bmatrix} x$$

$$1.5^{L-1} x \\ 0.5^{L-1} x$$

$$z^{[l]} = \omega^{[l]} x$$

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

$$a^{[l]} = g(z^{[l]}) = g(\omega^{[l]} a^{[l-1]})$$

$$1.5 \\ 0.5^L$$

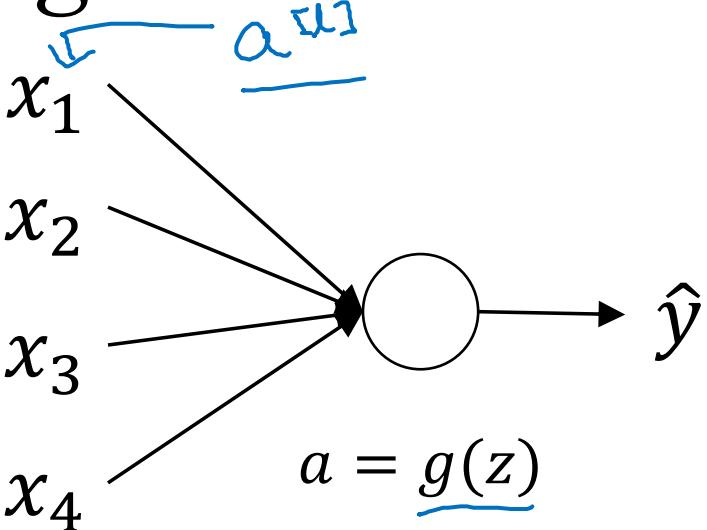
$L=150$

INTUITION: for very deep networks the activations in later layers (or similarly the gradients in earlier layers) can increase/decrease exponentially if the weights are  $>1$ / $<1$ .

- This is because the activation of one layer is, in practice, multiplied by the activation of the next one.
- This can be prevented with an intelligent initialization of the weights.

# Single neuron example

How to initialize the weights



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

Large  $n \rightarrow$  Smaller  $w_i$

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

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

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

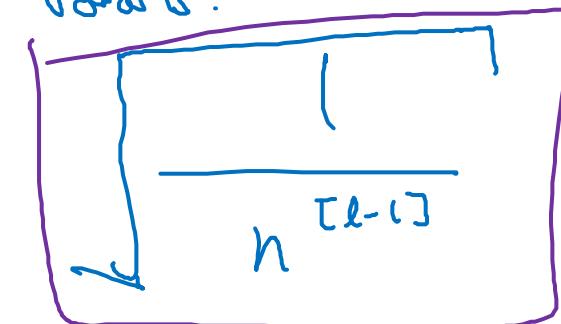
$w^{[l]}$  "HE INITIALIZATION"

- commonly used for neurons with ReLU activations.
- assuming that consecutive layers are densely connected (each neuron on layer  $l$  has  $l-1$  inputs), it makes sense (to limit the output of that neuron) to set each arc's weight to  $\frac{1}{n-\text{inputs}} = \frac{1}{n^{[l-1]}}$ . In practice one initializes the weights of neurons on layer  $l$  to

$$w^l \sim N\left(0, \frac{2}{n^{[l-1]}}\right)$$

Other variants:

Tanh



Xavier initialization

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

↑

- In practice, this is not one of the most influential hyperparameters.



deeplearning.ai

# Setting up your optimization problem

---

## Numerical approximation of gradients

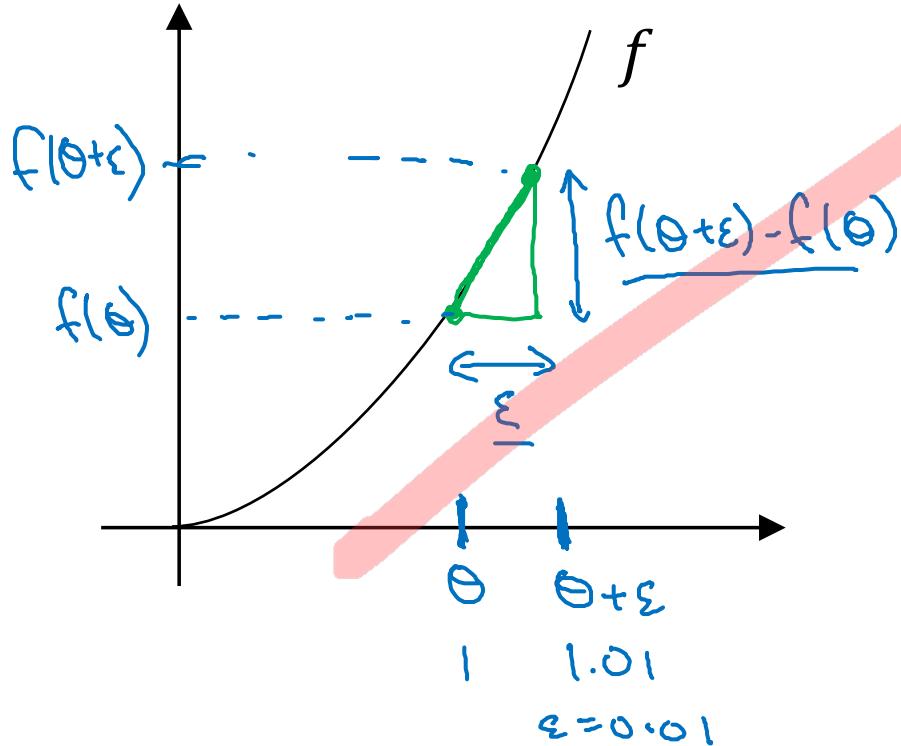
One can approximate a gradient numerically as

$$f'(θ) = \frac{f(θ+ε) - f(θ-ε)}{2ε}$$

This approximation has an error of order  $O(ε^2)$ , and is more precise than the standard "one-sided" derivative, which has an error with  $O(ε)$ .

# Checking your derivative computation

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



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

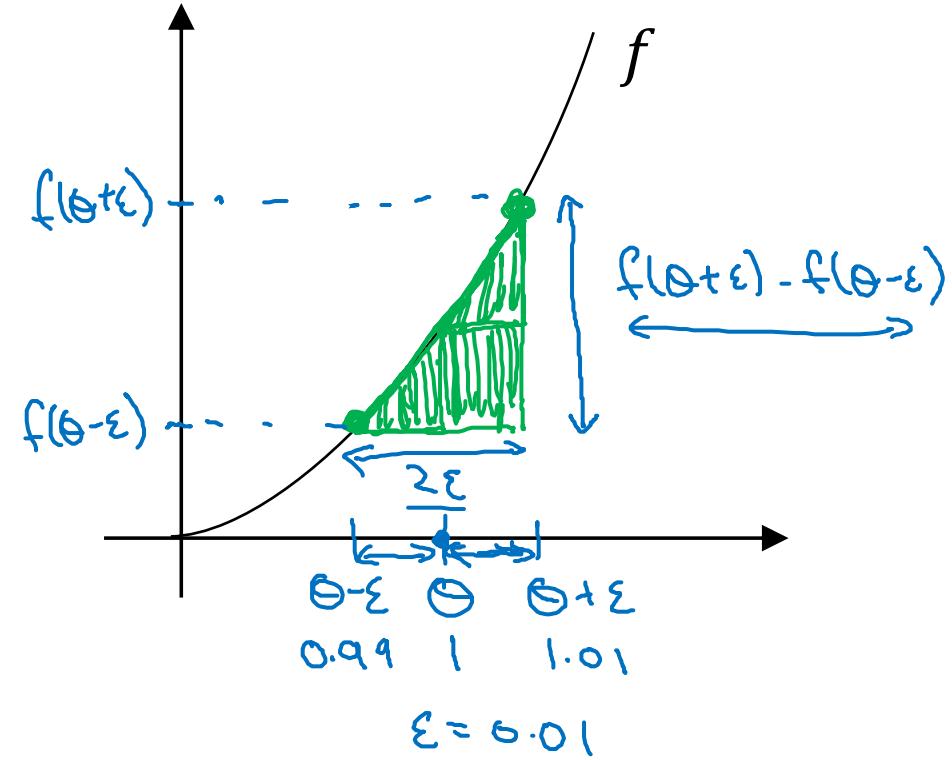
$$\begin{aligned} g(\theta) &= \frac{d}{d\theta} f(\theta) = f'(\theta) \\ g(\theta) &= 3\theta^2 \\ g(1) &= 3 \cdot (1)^2 = 3 \\ \text{when } \theta &= 1 \end{aligned}$$

$$\frac{dw}{db}$$

$$\begin{aligned} &\boxed{\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon}} \approx g(\theta) \\ &\frac{(1.01)^3 - 1^3}{0.01} = \frac{3.0301}{0.01} \approx 3 \\ &\frac{3.0301}{3.1} \approx 3 \end{aligned}$$

# Checking your derivative computation

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



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

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

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

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

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

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

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)

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

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

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

- compute  $d\theta_{\text{approx}}[i]$  independently  
for each  $\theta[i]$ .

- verify that  $d\theta_{\text{approx}}$  and  $d\theta$  (analytical)  
are not too dissimilar.

- In practice:  
- use  $\varepsilon = 10^{-7}$   
- if error  $\sim 10^{-7}$  ok  
 $\sim 10^{-5}$  maybe ok  
 $\sim 10^{-3}$  problem

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

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

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



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking  
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{d\theta_{approx}[i]}{d\theta[i]} \leftarrow \frac{d\theta[i]}{\Delta}$$

$$+ \frac{d\theta_{approx}[i] - d\theta[i]}{\Delta}$$

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

$$\frac{db^{[l]}}{\Delta} \quad \frac{dw^{[l]}}{\Delta}$$

- Remember regularization.

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

$d\theta$  = gradient of  $J$  wrt.  $\theta$

- Doesn't work with dropout.

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

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

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

- it's possible that the backprop implementation works well for some values of the parameters and badly for others  
↳ e.g.  $w, b \sim 0$  ↳ e.g.  $w, b \gg 0$ .

- Therefore you could run the gradcheck at the beginning of the training (params are in a certain range), then update for a bit, and then repeat the gradient check (params in different range).