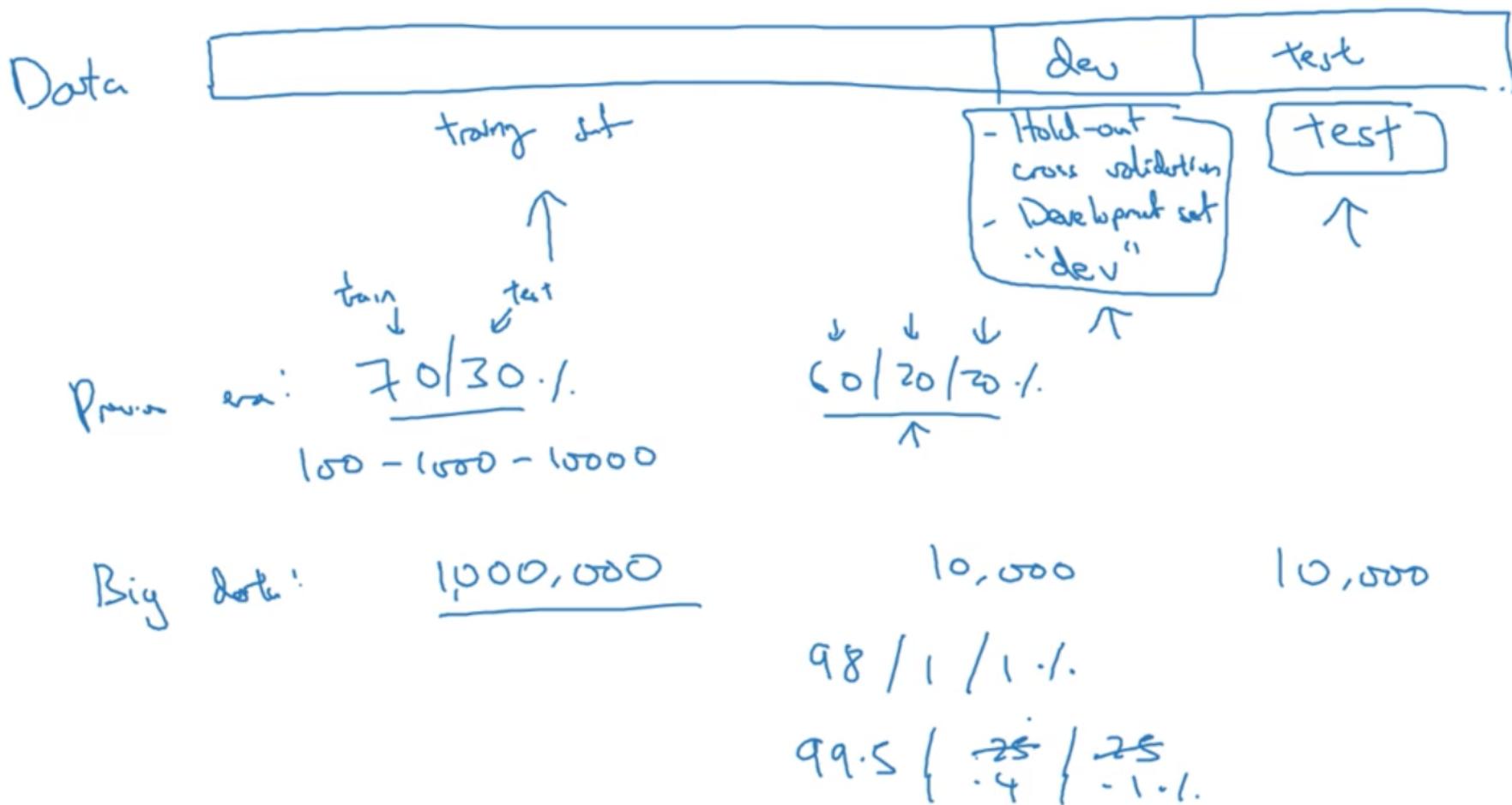


# Train/dev/test sets



Andrew Ng

# Mismatched train/test distribution

Corts

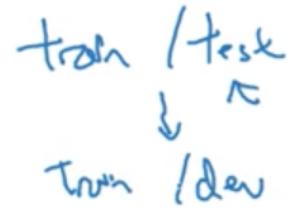
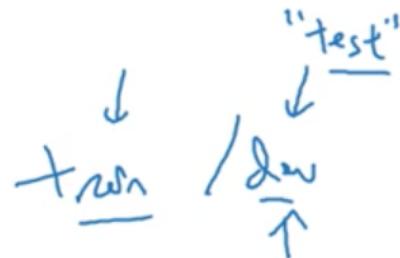
Training set:

Cat pictures from  
webpages

Dev/test sets:

Cat pictures from  
users using your app

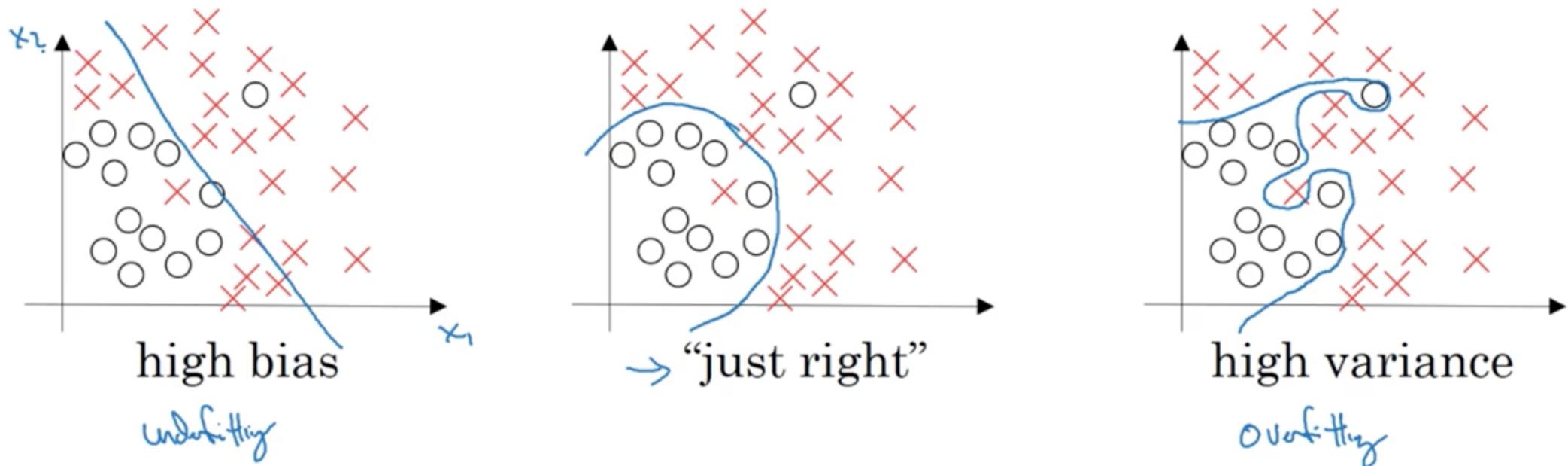
→ Make sure dev and test come from same distribution.



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

Andrew Ng

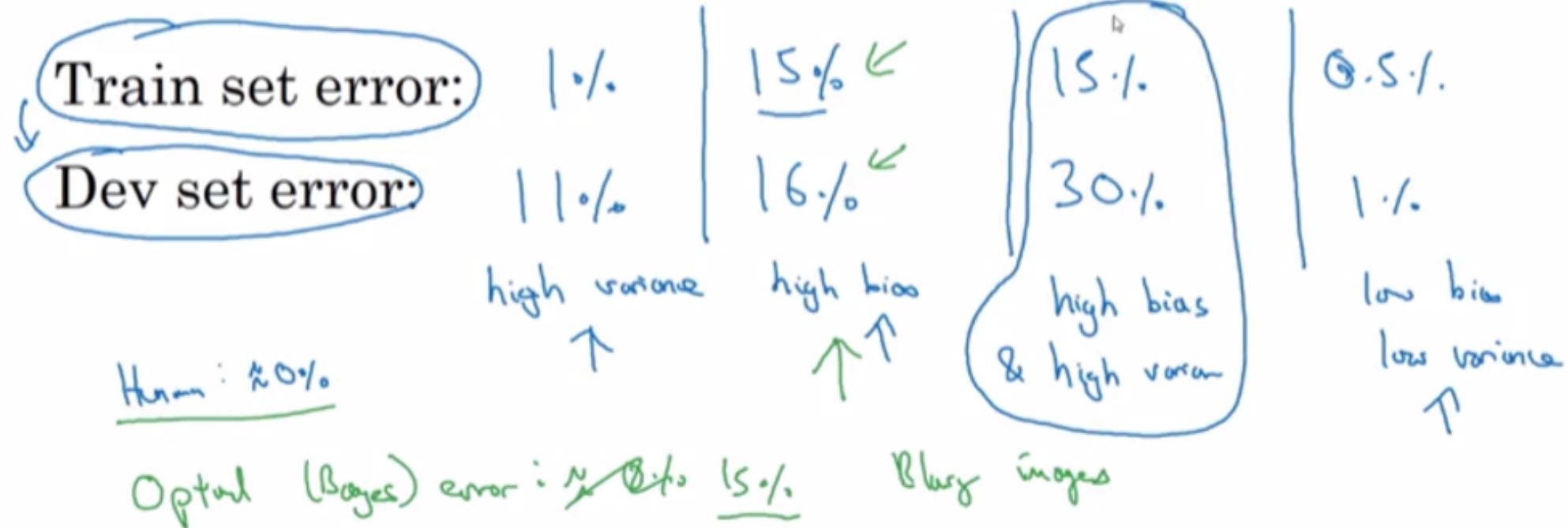
# Bias and Variance



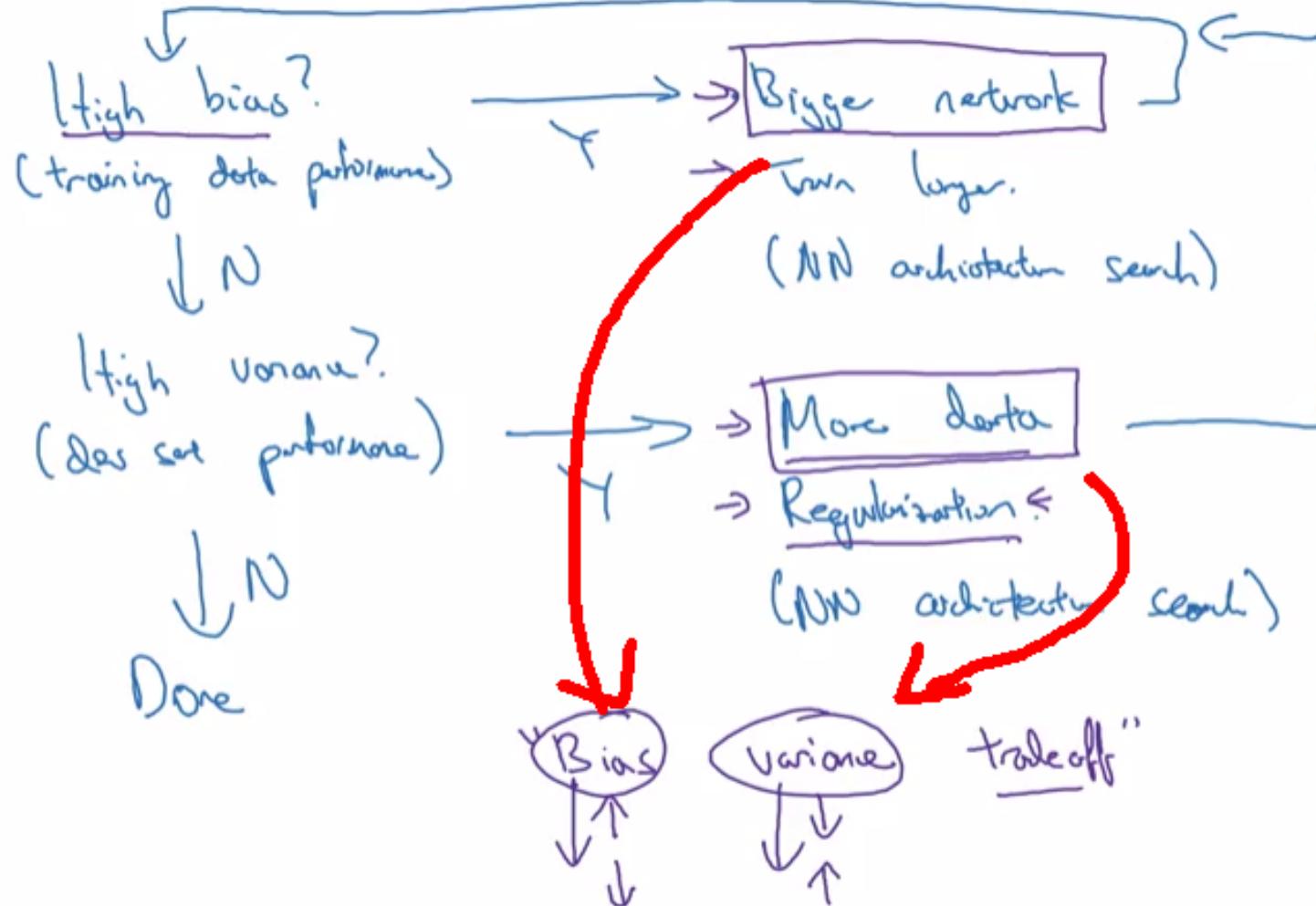
Andrew Ng

# Bias and Variance

Cat classification



# Basic recipe for machine learning



# Logistic regression

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

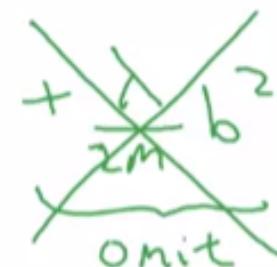
$$w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$$

$\lambda$  = regularization parameter

lambda

lambd

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



Popular!

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

L<sub>1</sub> regularization

$$\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \leftarrow w \text{ will be sparse}$$

# Neural network

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

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

"Frobenius norm"

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

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

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

?

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

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

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

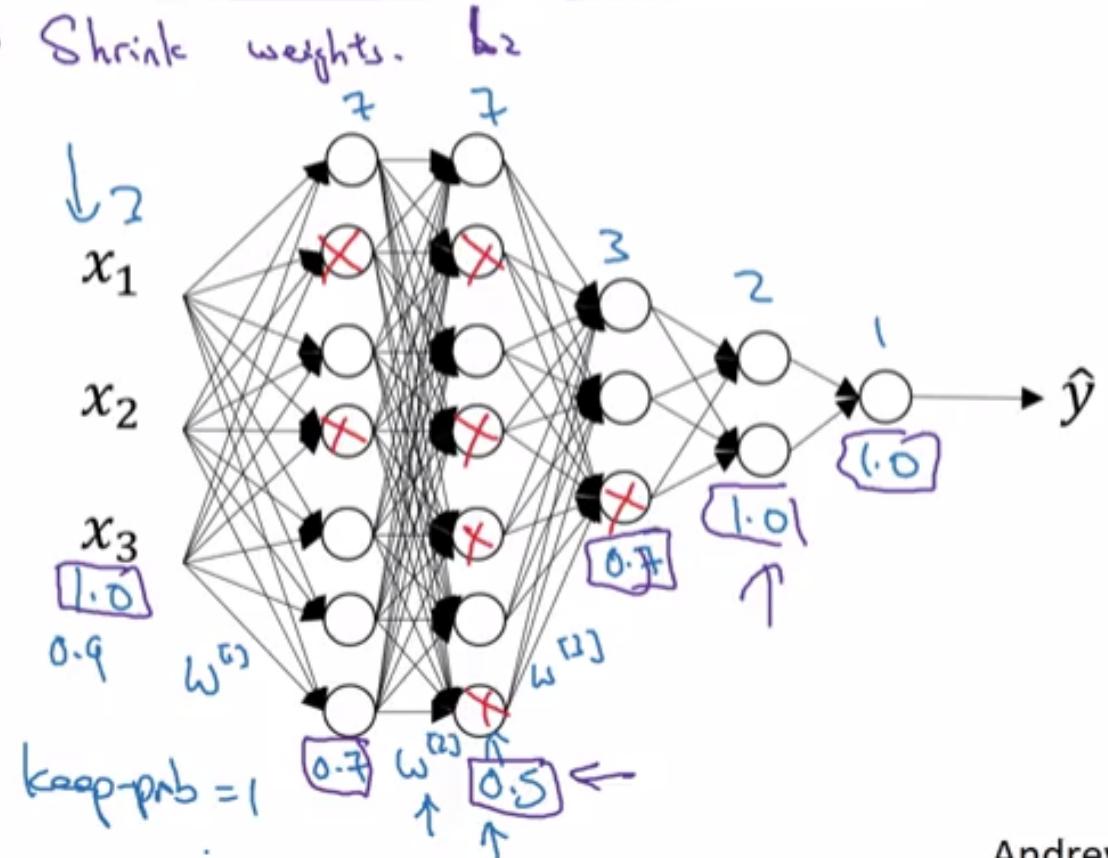
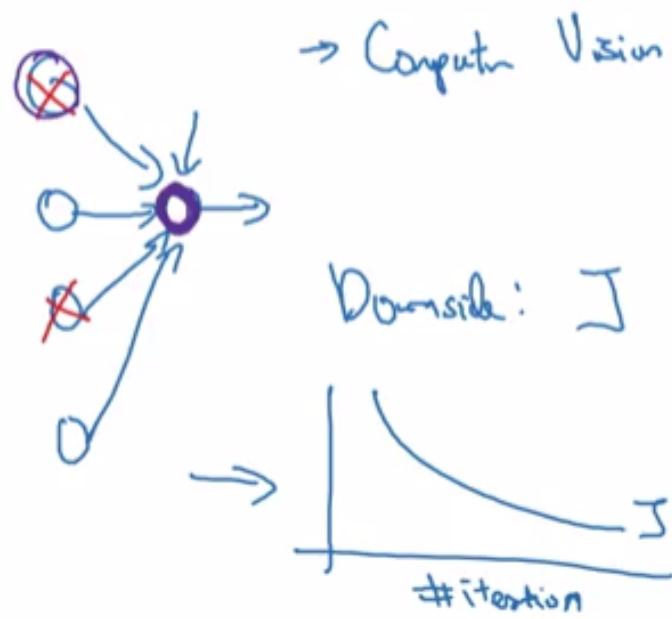
"Weight decay"

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

Andrew Ng

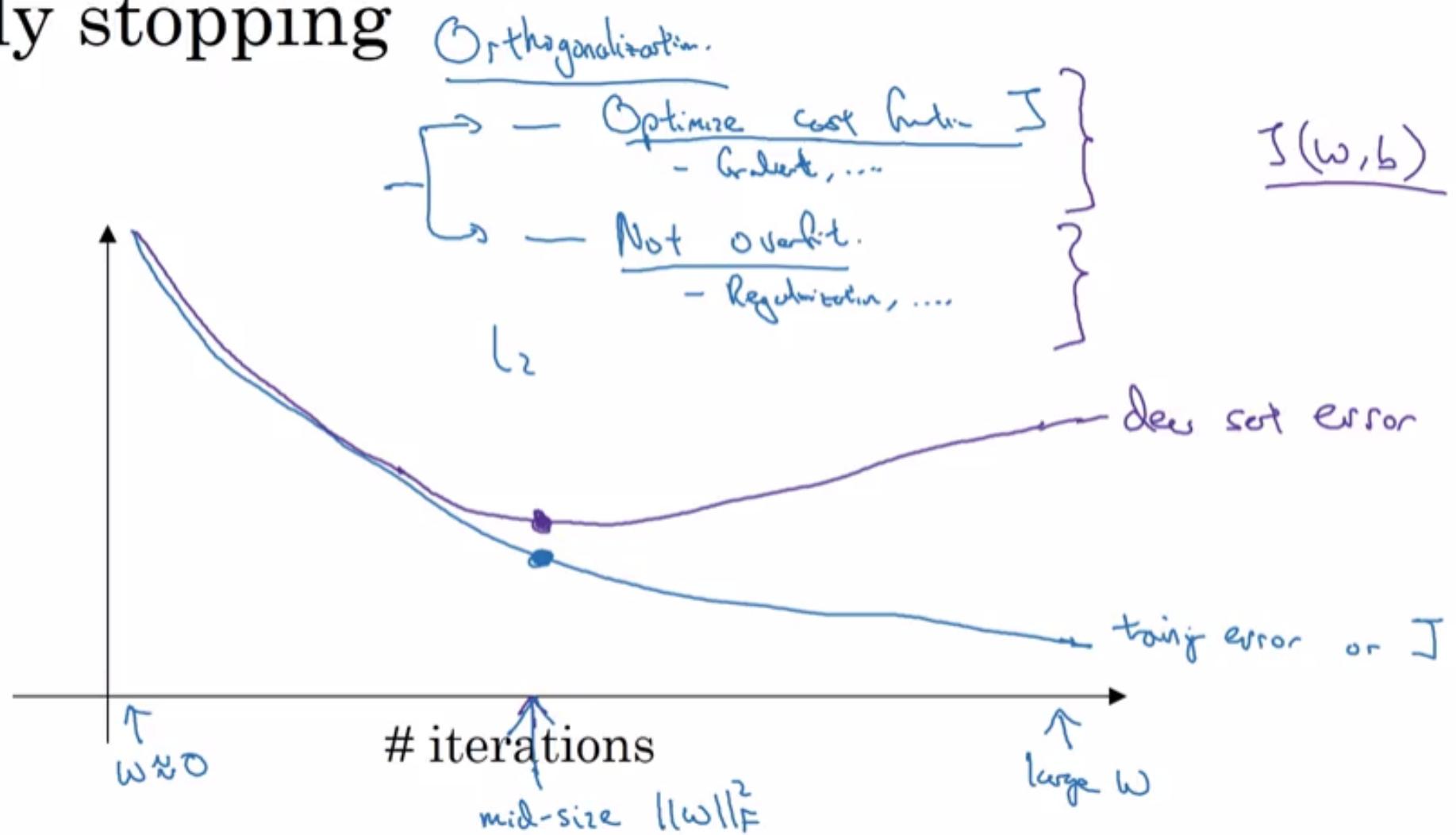
# Why does drop-out work?

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



Andrew Ng

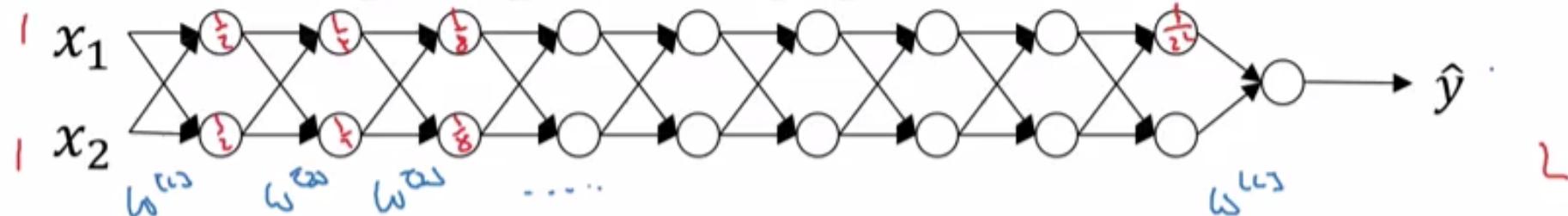
# Early stopping



Andrew Ng

# Vanishing/exploding gradients

$L=150$



$$\hat{y} = \underbrace{\omega^{[L]} \left( \underbrace{\omega^{[L-1]} \left( \dots \left( \underbrace{\omega^{[1]} x}_{z^{[1]}} \right) a^{[2]} \right) a^{[3]} \right) \dots a^{[L-1]} g(z^{[L]})}_{\hat{y}}$$

Diagram illustrating the forward pass of a neural network. The input  $x$  is multiplied by weight matrix  $\omega^{[1]}$  to produce  $z^{[1]}$ . This is followed by activation  $a^{[2]} = g(z^{[1]})$ , another weight matrix  $\omega^{[2]}$ , and so on. The final output is  $\hat{y} = \omega^{[L]} a^{[L]}$ .

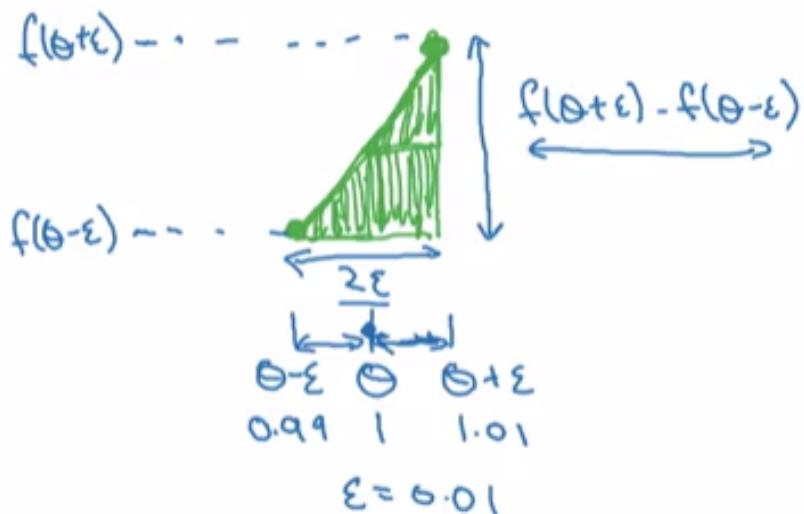
Explaining the terms:

- $\omega^{[1]} > I$
- $\omega^{[2]} < I$  [0.9 0.9]
- $\omega^{[L]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 1.5 \end{bmatrix}$
- $z^{[1]} = \omega^{[1]} x$
- $a^{[2]} = g(z^{[1]}) = z^{[1]}$
- $a^{[3]} = g(z^{[2]}) = g(\omega^{[2]} a^{[2]})$
- $\hat{y} = \underbrace{\omega^{[L]} \left[ \begin{bmatrix} 0.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L-1} x \right]}_{\hat{y}} \times 1.5^{L-1} \times 0.5^{L-1} x$

Andrew Ng

# Checking your derivative computation

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



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

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

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

approx error: 0.0001

(prev slide: 3.0301, error: 0.03)

$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon}$	$\frac{\mathcal{O}(\epsilon^2)}{\epsilon}$ $\frac{0.01}{0.01}$ $\underline{0.0001}$	$\frac{f(\theta+\epsilon) - f(\theta)}{\epsilon}$ $\uparrow$ $\text{error: } \mathcal{O}(\frac{1}{\epsilon})$ $0.01$
--	---	---

Andrew Ng

# Mini-batch gradient descent

repeat  $\{$   
for  $t = 1, \dots, 5000 \}$   $\}$

Forward prop on  $X^{[t]}$ .

$$z^{[t]} = w^{(1)} X^{[t]} + b^{(1)}$$

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

:

$$A^{[t]} = g^{(L)}(z^{[t]})$$

$$\text{Compute cost } J^{[t]} = \frac{1}{m} \sum_{i=1}^m l(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot m} \sum_{l=1}^L \|w^{(l)}\|_F^2.$$

Backprop to compute gradients w.r.t  $J^{[t]}$  (using  $(X^{[t]}, Y^{[t]})$ )

$$w^{(l)} := w^{(l)} - \alpha \nabla_w J^{[t]}, \quad b^{(l)} := b^{(l)} - \alpha \nabla_b J^{[t]}$$

3  
3

"1 epoch"  
└ pass through training set.

1 step of gradient descent  
using  $\frac{X^{[t]}, Y^{[t]}}{(m=1000)}$

$X, Y$

Andrew Ng

# Exponentially weighted averages

$$\underline{V_t} = \beta \underline{V_{t-1}} + (1-\beta) \underline{\theta_t} \leftarrow$$

$\beta = 0.9$  :  $\approx 10$  days' temper.

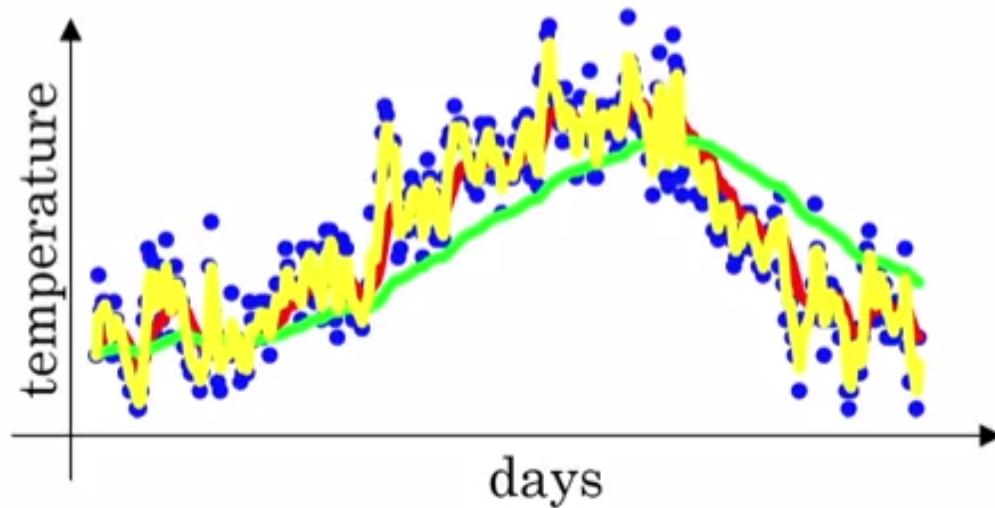
$\beta = 0.98$  :  $\approx 50$  days

$\beta = 0.5$  :  $\approx 2$  days

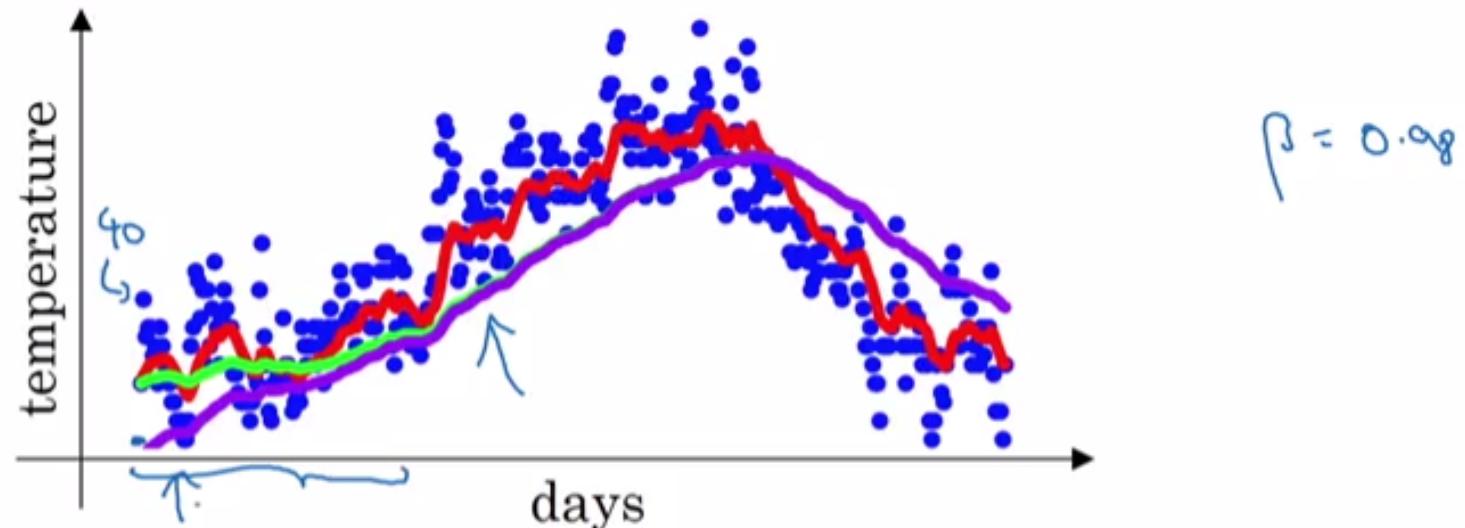
$V_t$  is approximately  
average over

$\rightarrow \approx \frac{1}{1-\beta}$  days'  
temperature.

$$\frac{1}{1-0.98} = 50$$



# Bias correction



$$\rightarrow v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

$$v_0 = 0$$

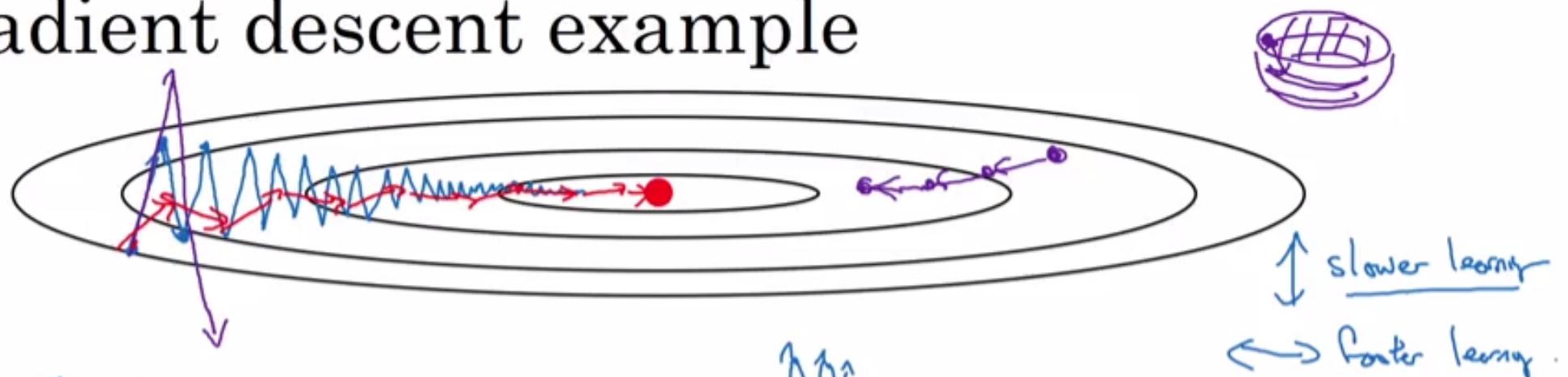
$$v_1 = \cancel{0.98 v_0} + 0.02 \theta_1$$

$$\begin{aligned} v_2 &= 0.98 v_1 + 0.02 \theta_2 \\ &= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2 \\ &= 0.0196 \theta_1 + 0.02 \theta_2 \end{aligned}$$

$$\left| \begin{array}{l} \frac{v_t}{1 - \beta^t} \\ t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396 \\ \frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396} \end{array} \right.$$

Andrew Ng

# Gradient descent example



Momentum:

On iteration  $t$ :

Compute  $\Delta w, \Delta b$  on current mini-batch.

$$v_{dw} = \beta v_{dw} + (1-\beta) \frac{\Delta w}{\text{velocity}}$$

$$v_{db} = \beta v_{db} + (1-\beta) \frac{\Delta b}{\text{acceleration}}$$

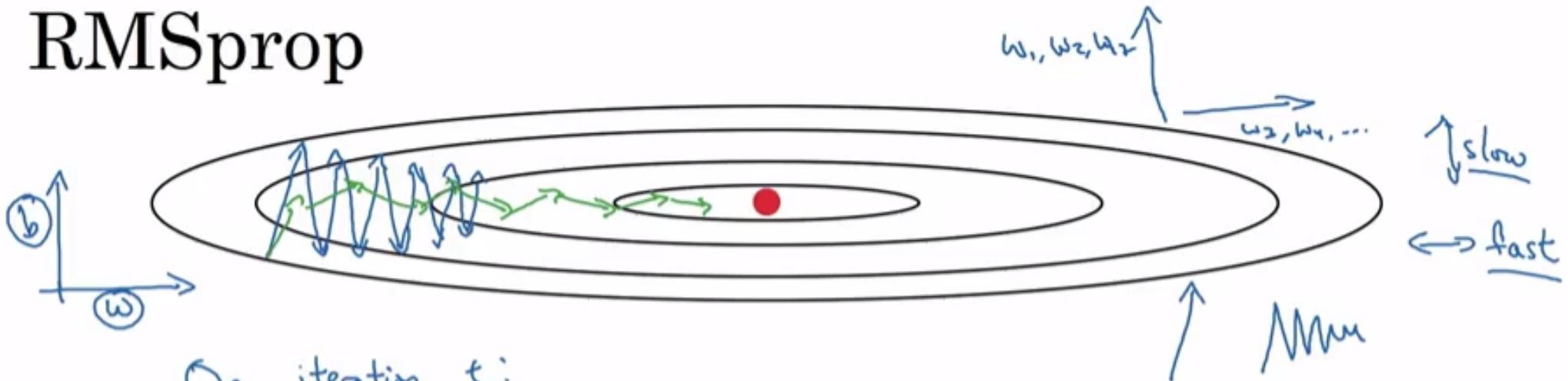
Friction  $\beta$   $\uparrow$  velocity  $\downarrow$  acceleration

$$w := w - \alpha v_{dw}, \quad b := b - \alpha v_{db}$$

$$v_w = \beta v_w + (1-\beta) \theta_w$$

Andrew Ng

# RMSprop



On iteration  $t$ :

Compute  $dW, db$  on current mini-batch

$$S_{dW} = \beta_2 S_{dW} + (1-\beta_2) dW^2 \quad \text{element-wise}$$

$$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \text{large}$$

$$w := w - \frac{\alpha}{\sqrt{S_{dW} + \epsilon}} dW \quad \leftarrow$$

$$b := b - \frac{\alpha}{\sqrt{S_{db} + \epsilon}} db \quad \begin{matrix} \checkmark \\ \text{can be relatively larger} \end{matrix} \leftarrow$$

$$\epsilon = 10^{-8}$$

# Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0. \quad V_{db} = 0, S_{db} = 0$$

On iteration  $t$ :

Compute  $\delta w, \delta b$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) \delta w, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) \delta b \leftarrow \text{"moment"} \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1 - \beta_2) \delta w^2, \quad S_{db} = \beta_2 S_{db} + (1 - \beta_2) \delta b \leftarrow \text{"RMSprop"} \beta_2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

$$w := w - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}$$

$$b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

# Hyperparameters choice:

- $\alpha$ : needs to be tune
- $\beta_1$ : 0.9  $\rightarrow (\underline{dw})$
- $\beta_2$ : 0.999  $\rightarrow (\underline{dw^2})$
- $\epsilon$ :  $10^{-8}$

Adam: Adaptive moment estimation



Adam Coates

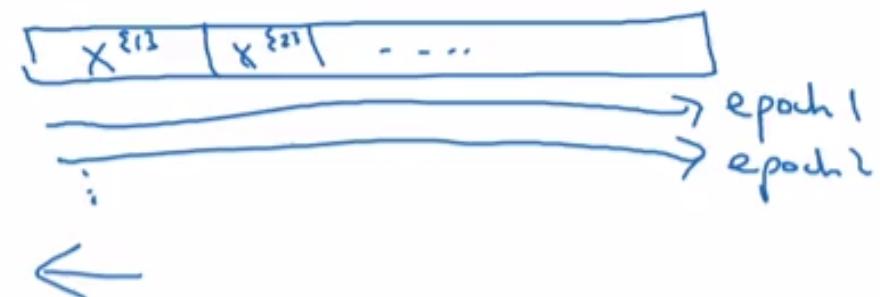
Andrew Ng

# Learning rate decay

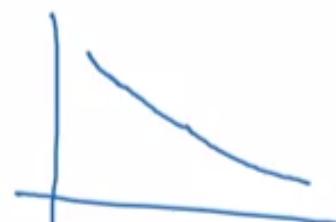
1 epoch = 1 pass through data.

$$\alpha = \frac{\alpha_0}{1 + \underbrace{\text{decay-rate} * \text{epoch-num}}_{\text{Curly bracket}}}$$

Epoch	$\alpha$
1	0.1
2	0.67
3	0.5
4	0.4
:	i

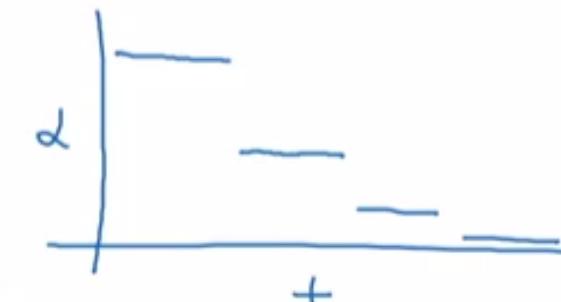


$$\alpha_0 = 0.2$$
$$\text{decay\_rate} = 1$$



# Other learning rate decay methods

Formulas

$$\alpha = \begin{cases} 0.95^{\text{epoch\_num}} \cdot \alpha_0 & - \text{exponentially decay} \\ \frac{k}{\sqrt{\text{epoch\_num}}} \cdot \alpha_0 & \text{or } \frac{k}{\sqrt{t}} \cdot \alpha_0 \end{cases}$$


Manual decay.

# Implementing Batch Norm

Given some intermediate values in NN

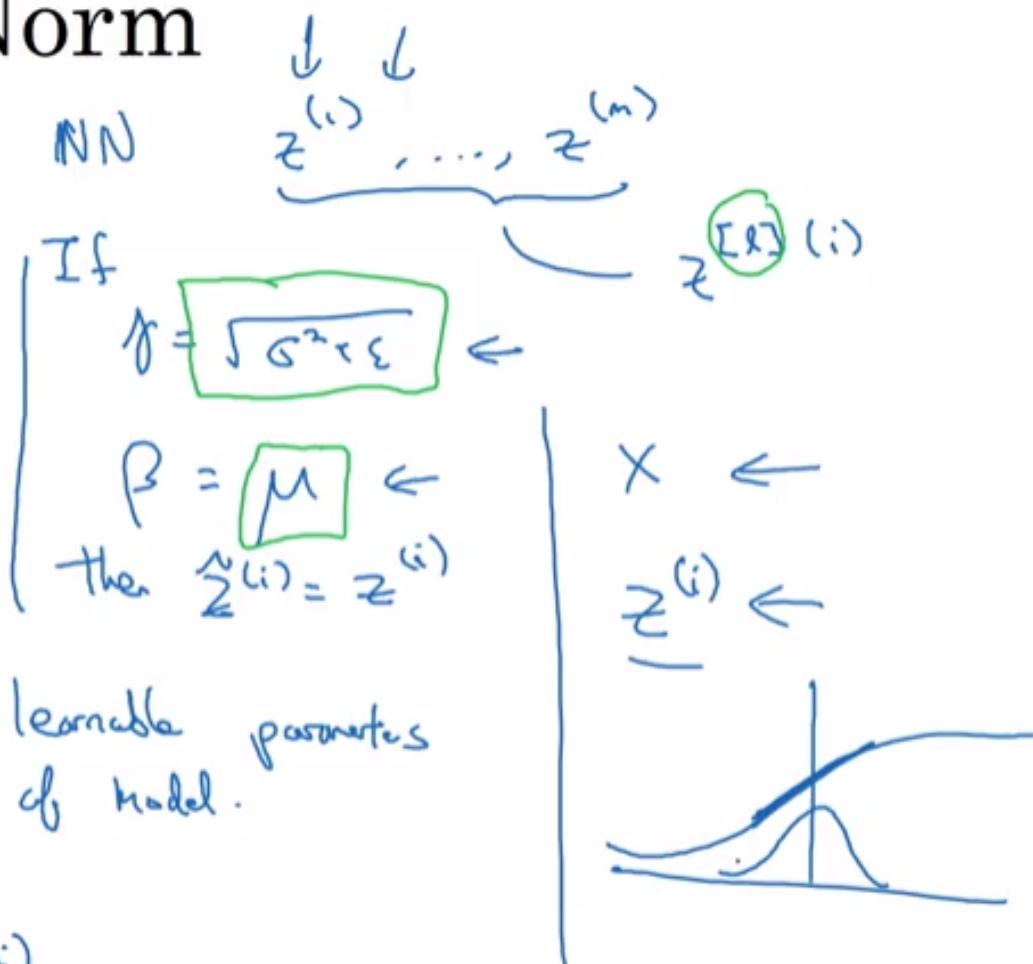
$$\mu = \frac{1}{m} \sum z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum (z_i - \mu)^2$$

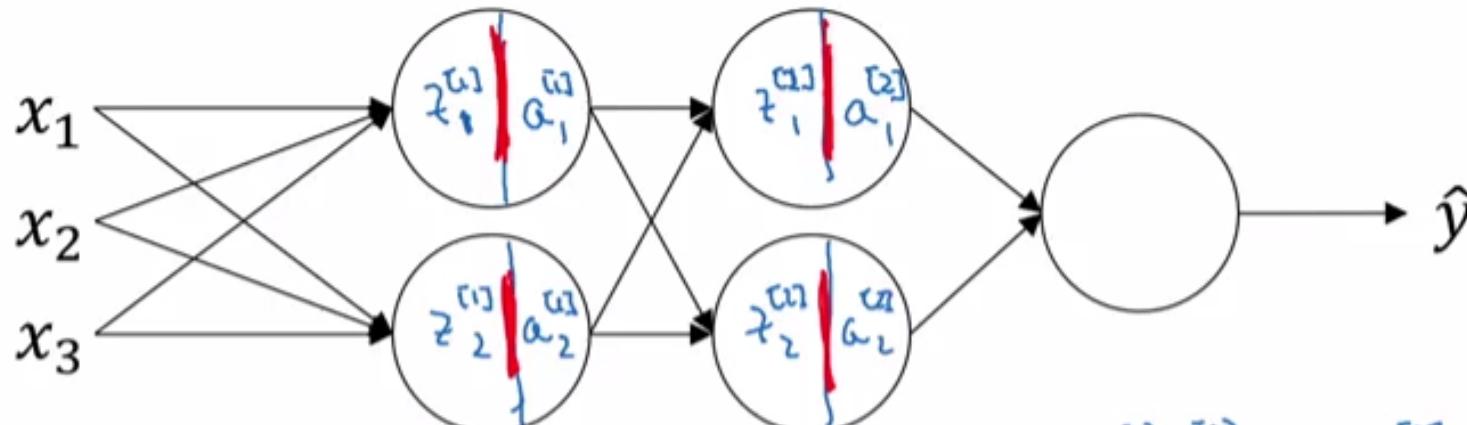
$$z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\hat{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta$$

Use  $\hat{z}^{(i)}$  instead of  $z^{(i)}$ .



# Adding Batch Norm to a network



$$x \xrightarrow{\omega^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\beta^{(1)}, \gamma^{(1)}} \underbrace{z^{(1)} \xrightarrow{\text{BatchNorm (BN)}} a^{(1)} = g(z^{(1)})}_{\text{BN}} \xrightarrow{\omega^{(2)}, b^{(2)}} z^{(2)} \xrightarrow{\beta^{(2)}, \gamma^{(2)}} \underbrace{z^{(2)} \xrightarrow{\text{BN}} a^{(2)} = g(z^{(2)})}_{\text{BN}}$$

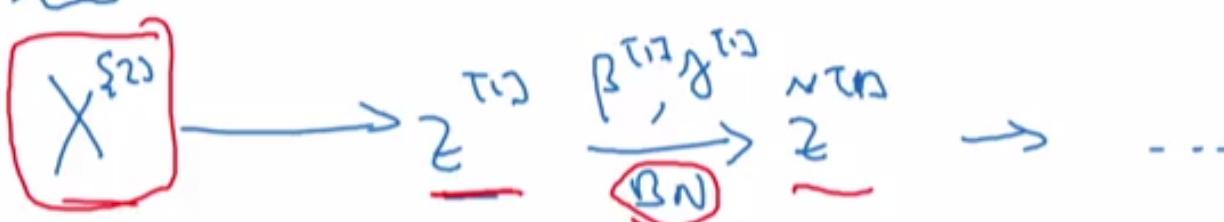
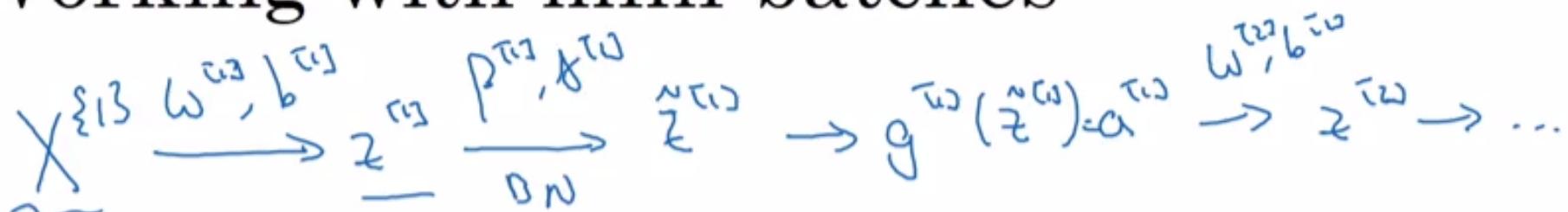
Parameters:

$$\left. \begin{array}{c} \omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}, \dots, \omega^{(L)}, b^{(L)}, \\ \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)} \end{array} \right\} \quad d\beta^{(l)} \quad \beta = \beta - \alpha d\beta^{(l)}$$

`tf.nn.batch_normalization`

$\rightarrow \beta$

# Working with mini-batches



$X^{i_2}$   $\rightarrow \dots$

Parameters:  $\cancel{W^{[l]}}$ ,  $\cancel{b^{[l]}}$ ,  $\beta^{[l]}$ ,  $\gamma^{[l]}$ .

$|$                    $|$                    $|$   
 $(n^{[l]}, 1)$        $(n^{[l]}, 1)$        $(n^{[l]}, 1)$

$z^{[l]}$   
 $(n^{[l]}, 1)$

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

$$z^{[l]} = W^{[l]} a^{[l-1]}$$

$$\underline{z^{[l]}} = \gamma^{[l]} z_{\text{norm}} + \beta^{[l]}$$

Andrew Ng

# Implementing gradient descent

for  $t = 1 \dots \text{num MiniBatches}$   
Compute forward pass on  $X^{[t]}$ .

In each hidden layer, use BN to replace  $\underline{z}^{[l]}$  with  $\hat{\underline{z}}^{[l]}$ .

Use backprop to compute  $\underline{dw}^{[l]}$ ,  ~~$\underline{db}^{[l]}$~~ ,  $\underline{dp}^{[l]}$ ,  $\underline{df}^{[l]}$

Update parameters  $\left. \begin{array}{l} w^{[l]} := w^{[l]} - \alpha \underline{dw}^{[l]} \\ b^{[l]} := b^{[l]} - \alpha \underline{db}^{[l]} \\ g^{[l]} := \dots \end{array} \right\} \leftarrow$

Works w/ momentum, RMSprop, Adam.