



deeplearning.ai

# Setting up your ML application

---

## Train/dev/test sets

# Applied ML is a highly iterative process

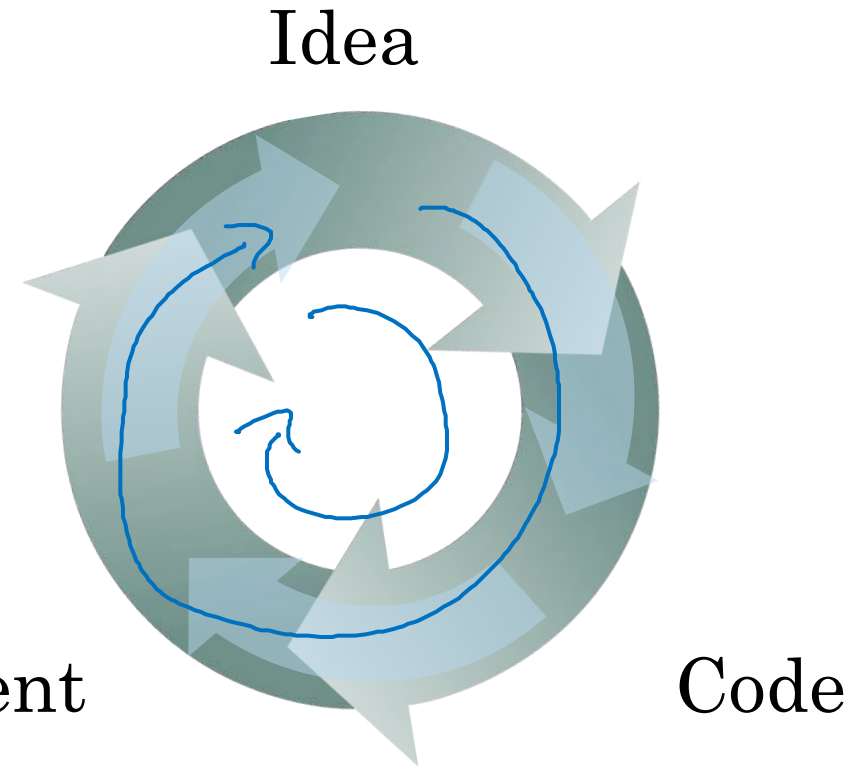
# layers

## # hidden units

# learning rates

## activation functions

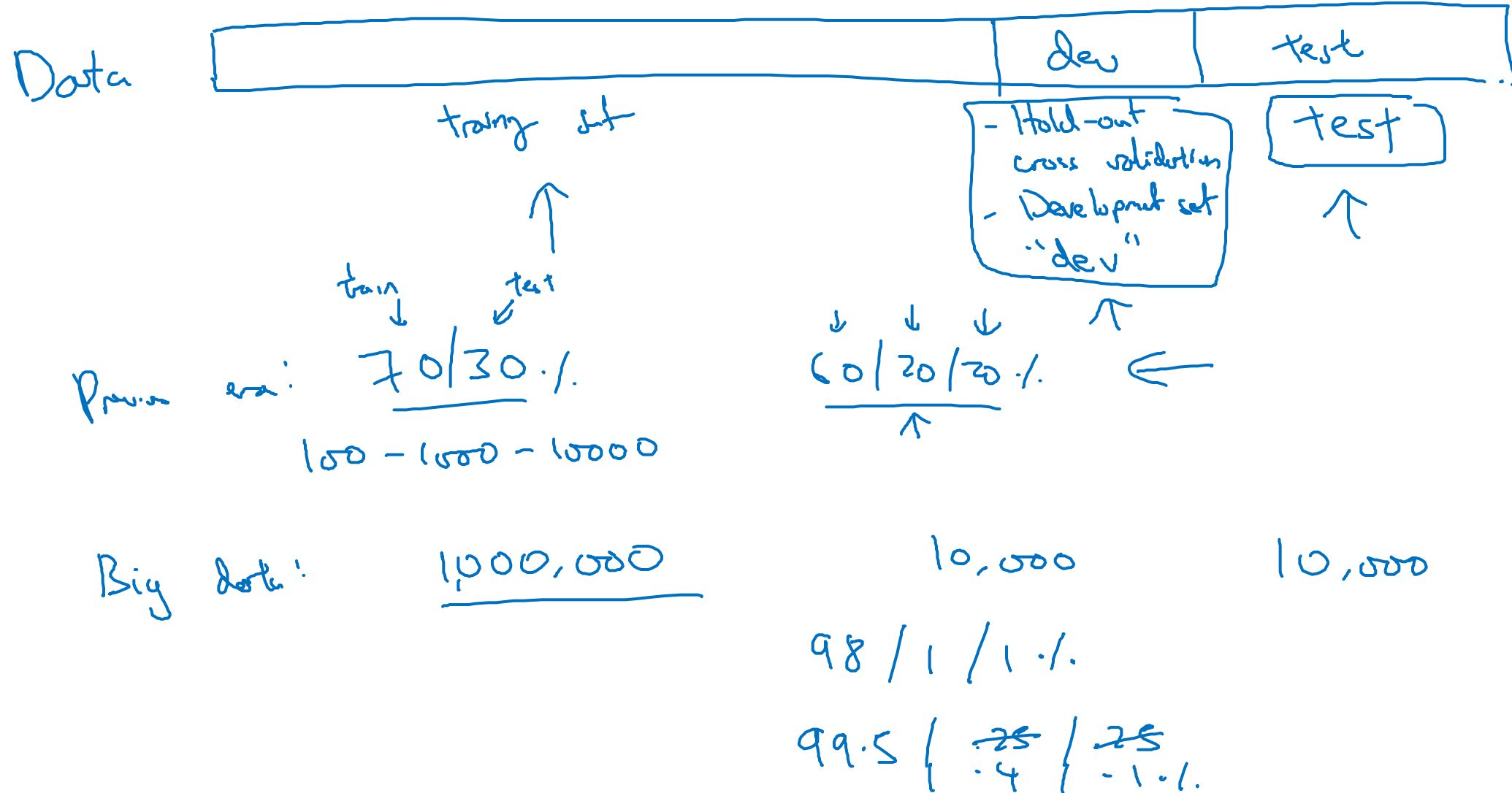
...



NLP, Vision, Speech, Structured Data

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

# Train/dev/test sets



# Mismatched train/test distribution

Certs

↙  
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 / dev      "test"

train / test  
↓      ↖  
→ train / dev

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



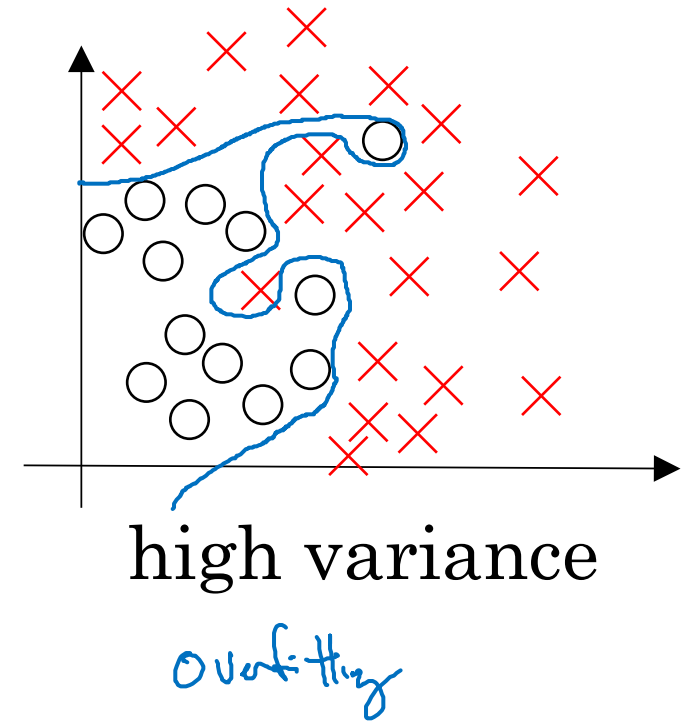
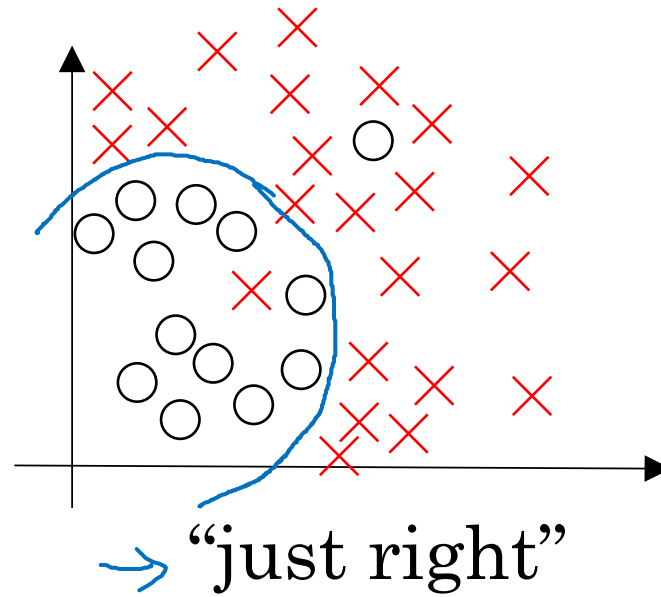
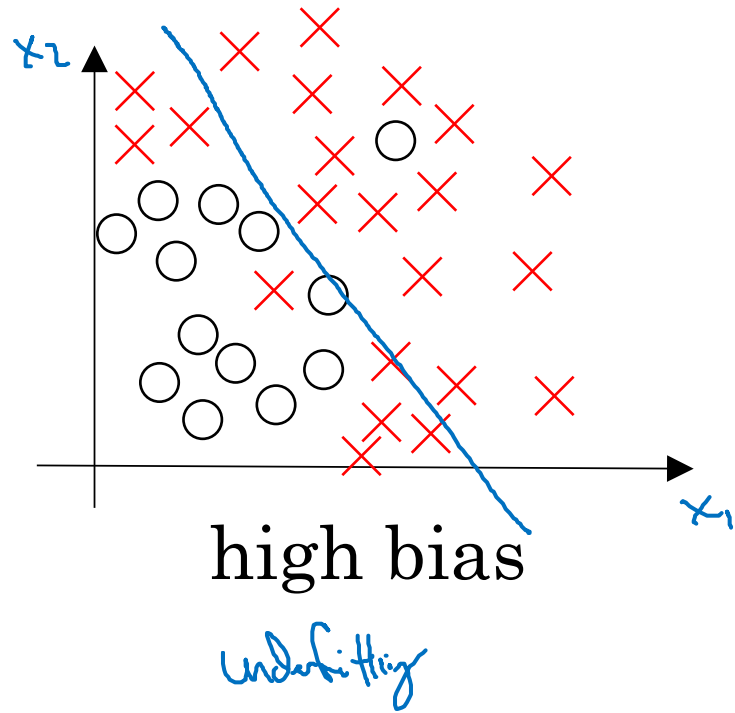
deeplearning.ai

# Setting up your ML application

---

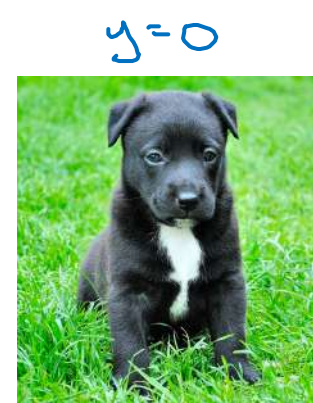
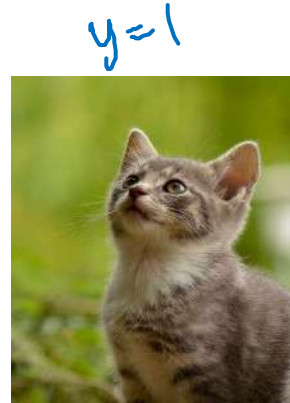
## Bias/Variance

# Bias and Variance



# Bias and Variance

Cat classification



Train set error:

Dev set error:

1%

11%

high variance  
↑

15% ←

16% ←

high bias  
↑ ↑

15%

30%

high bias  
& high variance

0.5%

1%

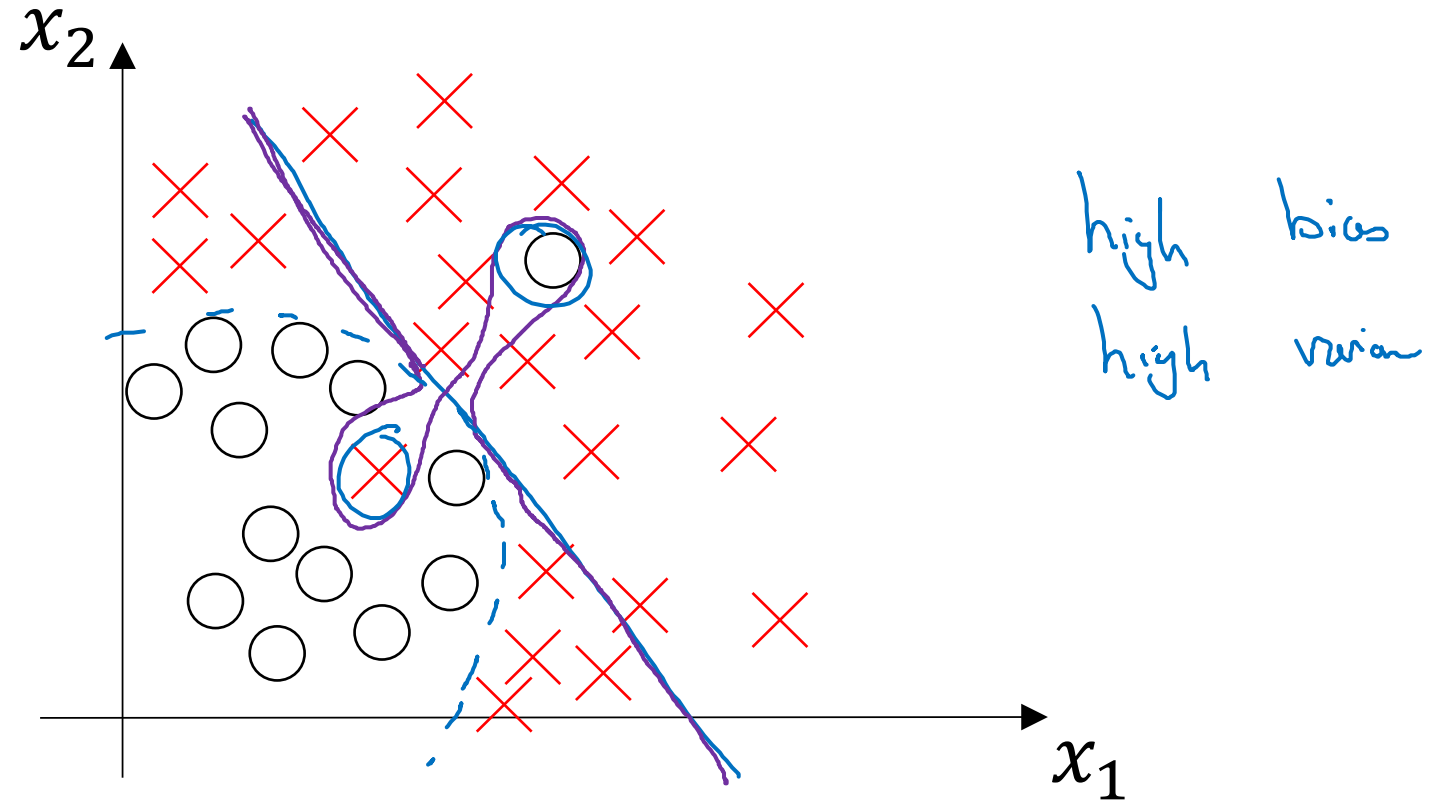
low bias  
low variance  
↑

Human: ~0%

Optimal (Bayes) error: ~~~0%~~ 15%

Blurry images

# High bias and high variance







deeplearning.ai

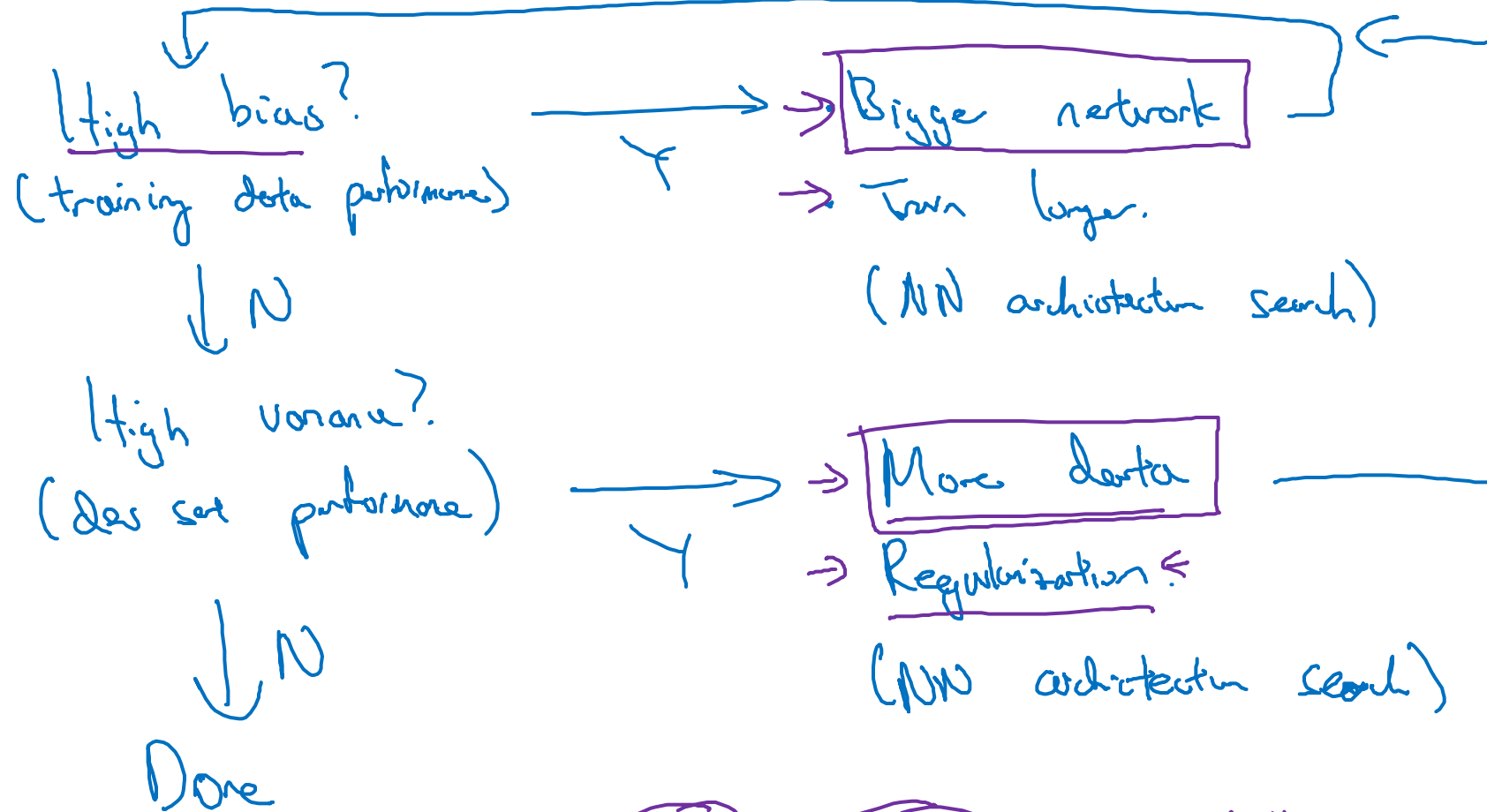
# Setting up your ML application

---

## Basic “recipe” for machine learning

# Basic “recipe” for machine learning

# Basic recipe for machine learning





deeplearning.ai

# Regularizing your neural network

---

## Regularization

# Logistic regression

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

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

$\lambda$  = regularization parameter  
lambda lambda

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

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

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

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

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

"Frobenius norm"

$\|\cdot\|_2^2$

$\|\cdot\|_F^2$

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

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

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

"Weight decay"

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

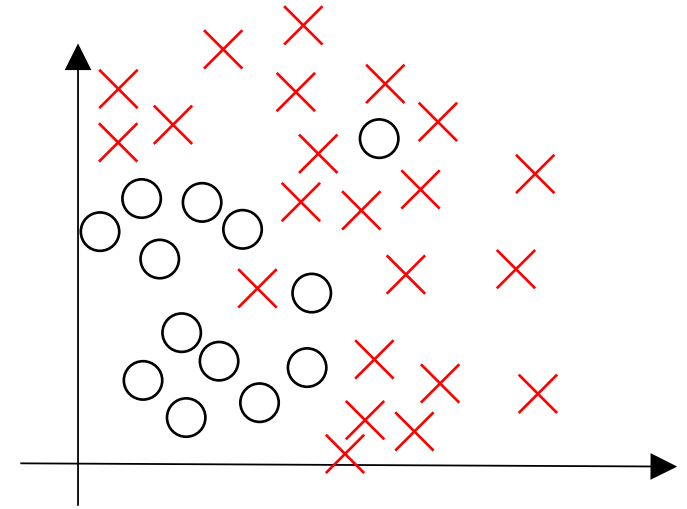
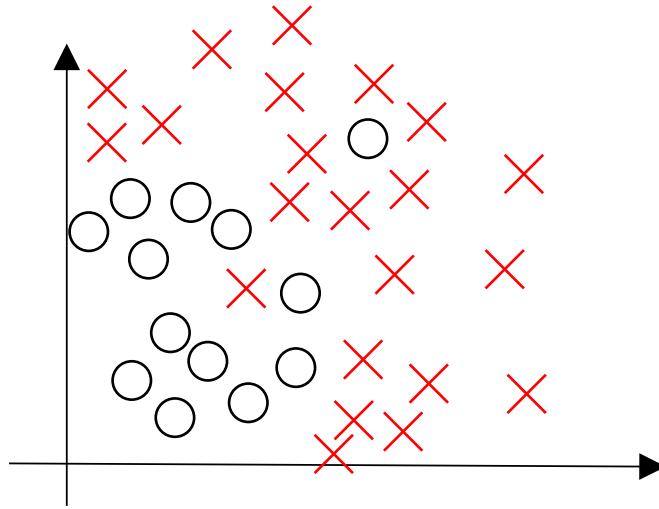
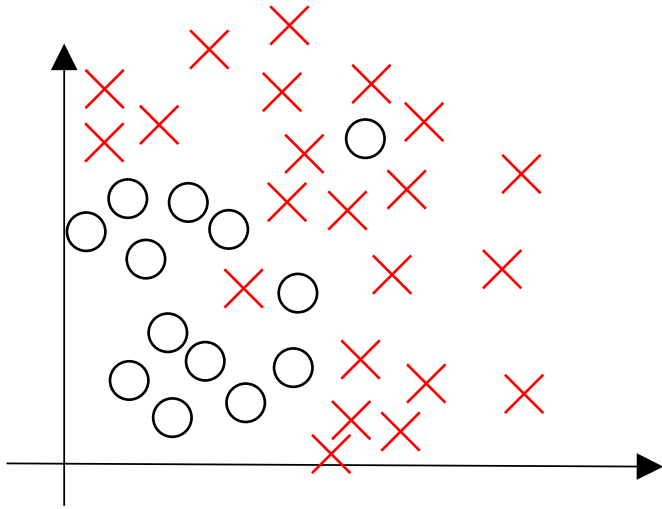
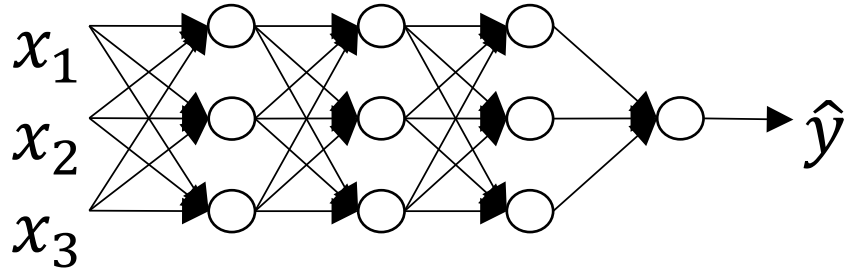
$$= w^{[l]} - \frac{\alpha \lambda}{n} w^{[l]} - \alpha \text{(from backprop)}$$

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

# Neural network

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

# How does regularization prevent overfitting?





# How does regularization prevent overfitting?



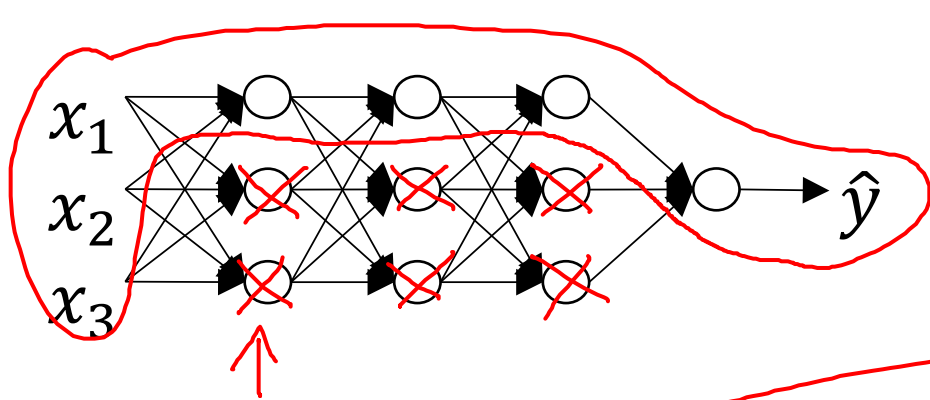
deeplearning.ai

# Regularizing your neural network

---

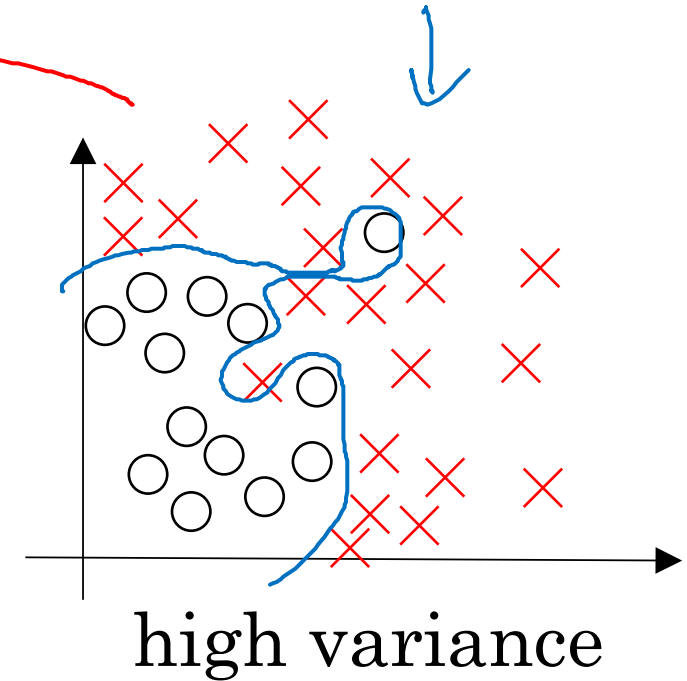
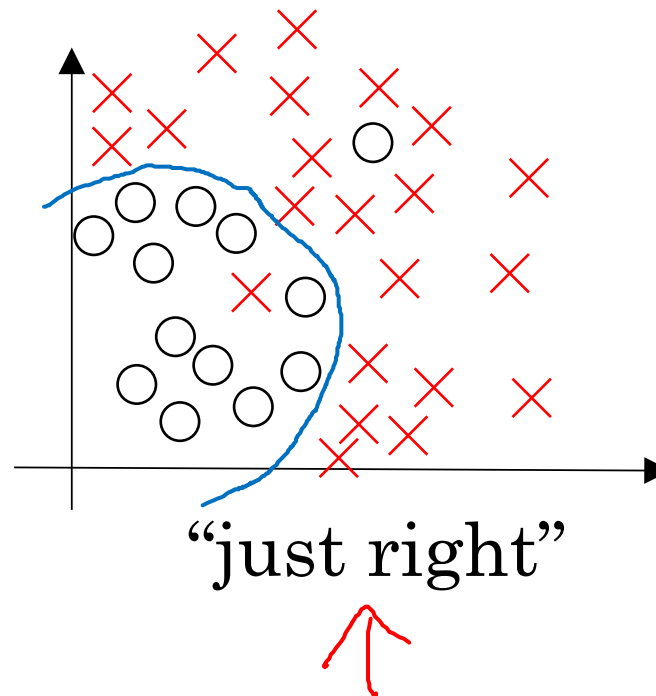
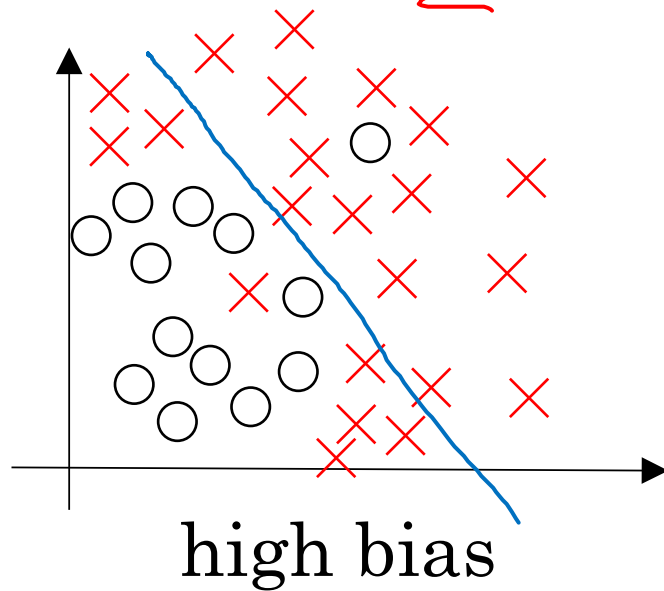
## Why regularization reduces overfitting

# How does regularization prevent overfitting?

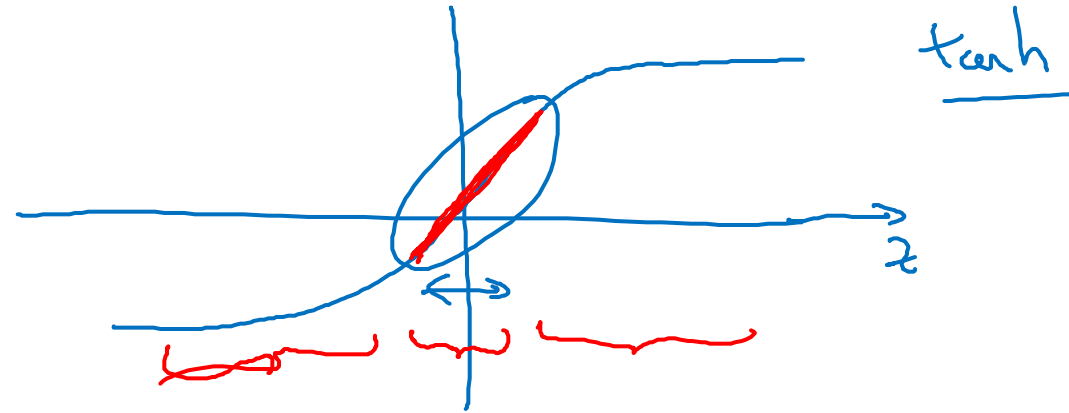


$$J(w^{(1)}, b^{(1)}) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

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



# How does regularization prevent overfitting?



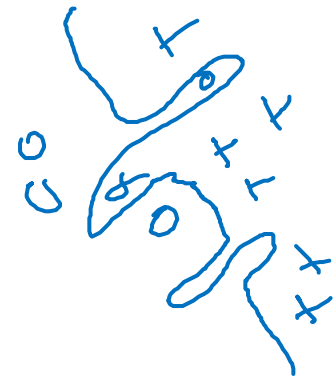
$$g(z) = \tanh(z)$$

$$\lambda \uparrow$$

$$W^{[L]} \downarrow$$

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

Every layer  $\approx$  linear.



$$J(\dots) = \underbrace{\sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)})}_{\text{training loss}} + \underbrace{\frac{\lambda}{2m} \sum_L \|W^{[L]}\|_F^2}_{\text{regularization term}}$$





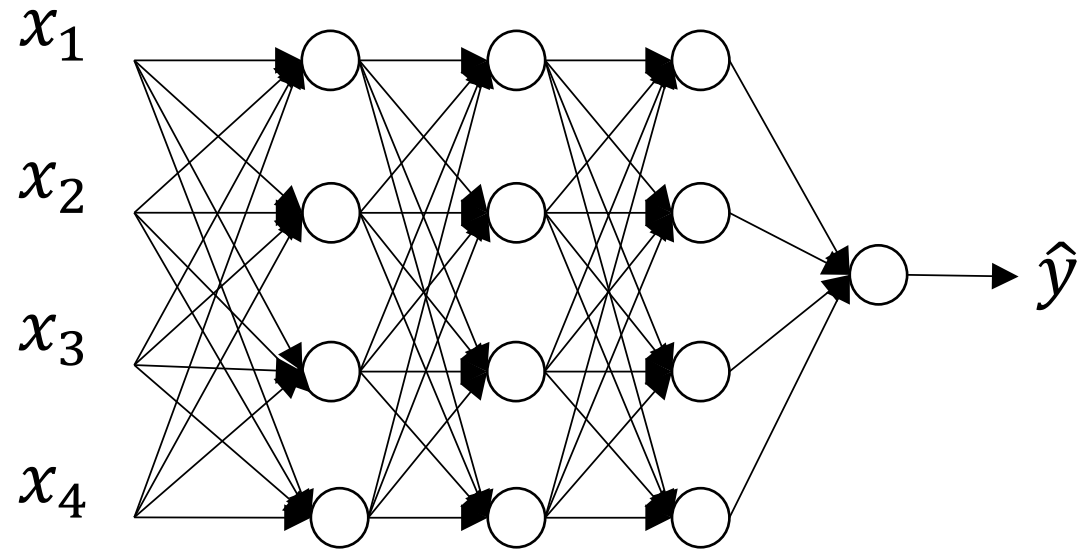
deeplearning.ai

# Regularizing your neural network

---

## Dropout regularization

# Dropout regularization



↑  
0.5    ↑  
0.5    ↑  
0.5

# Implementing dropout ("Inverted dropout")

Illustrate with layer  $l=3$ . keep-prob = 0.8 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

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

$\uparrow$

$\nwarrow$  reduced by 20%.

$\nwarrow$   $= 0.8$

Test

# Making predictions at test time

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

No drop out.

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

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

$$z^{[2]} = W^{[2]} \frac{a^{[1]}}{\text{keep-prob}} + b^{[2]}$$

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

↓  
↑  
y

/ = keep-prob





deeplearning.ai

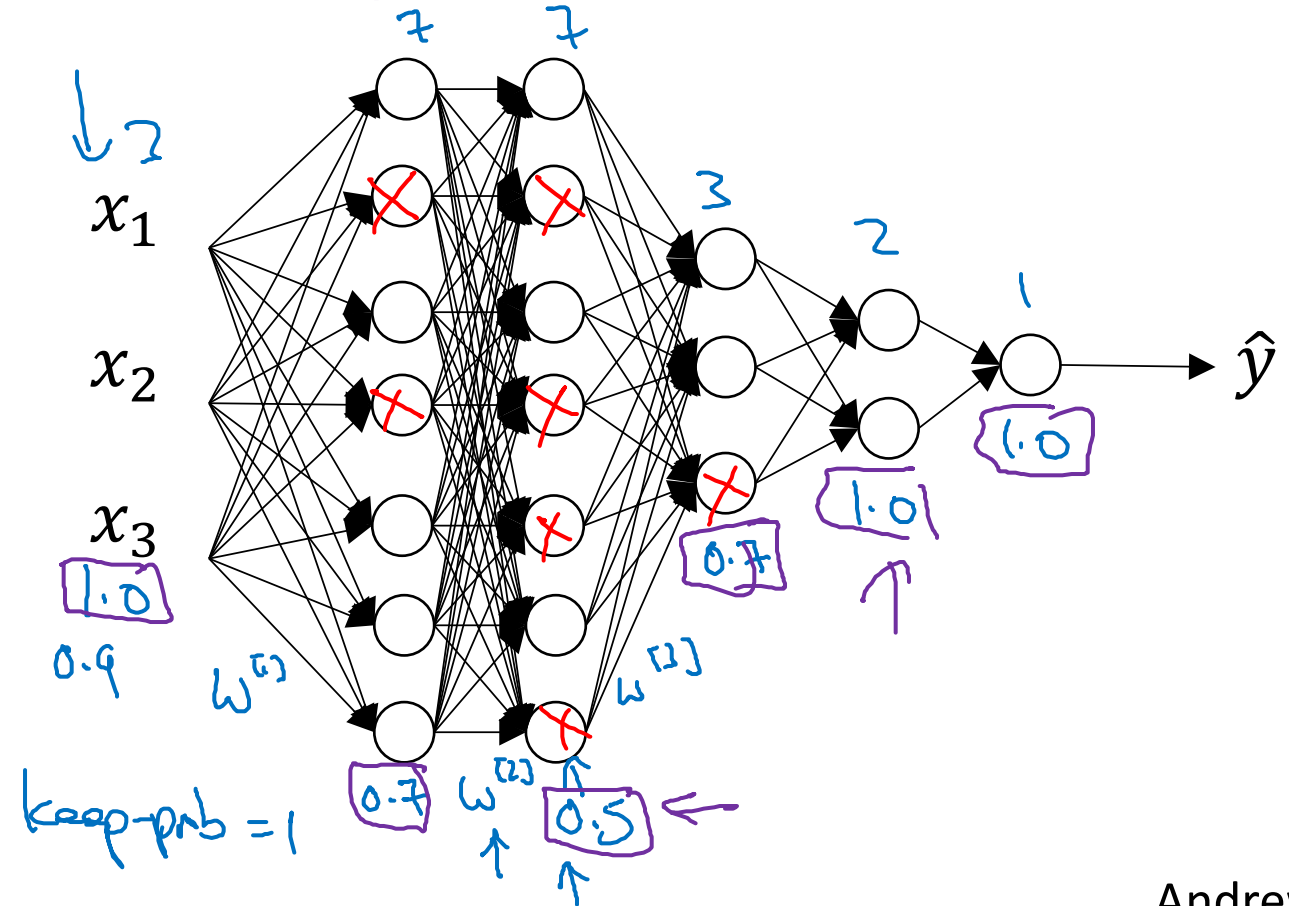
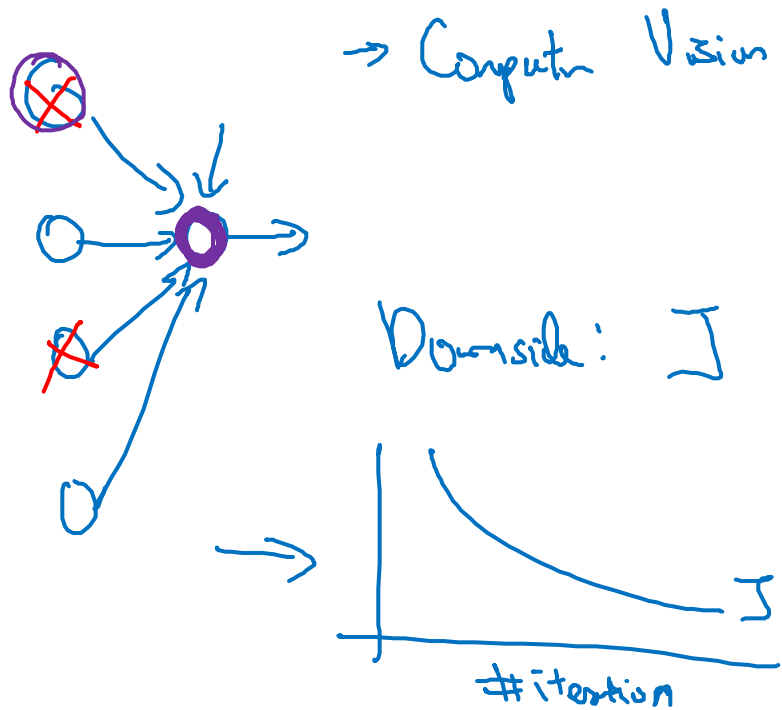
# Regularizing your neural network

---

## Understanding dropout

# Why does drop-out work?

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





deeplearning.ai

# Regularizing your neural network

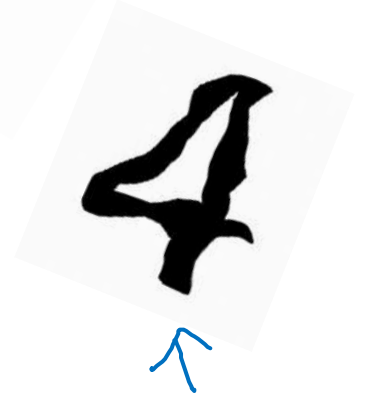
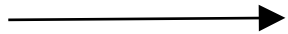
---

## Other regularization methods

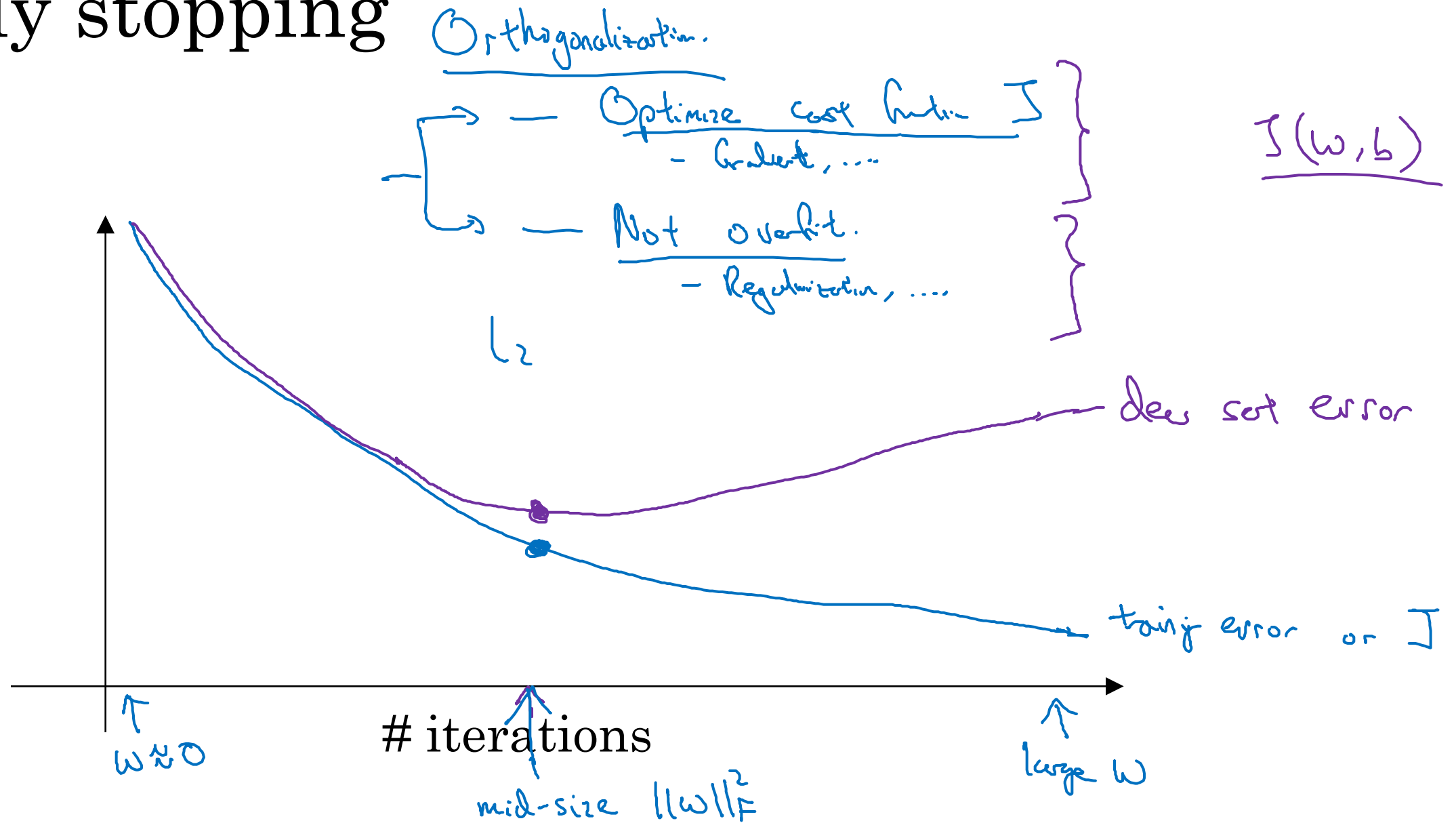
# Data augmentation



4



# Early stopping





deeplearning.ai

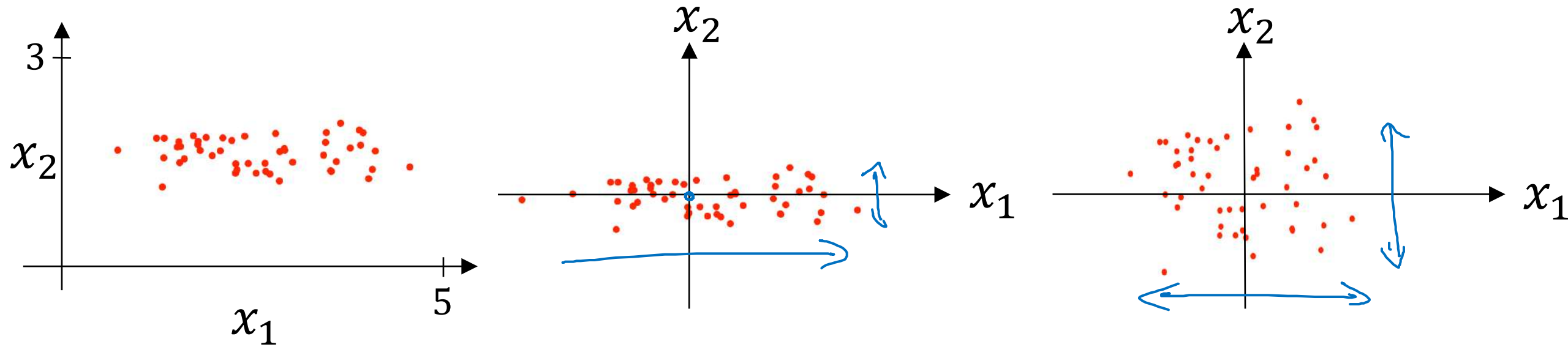
Setting up your  
optimization problem

---

Normalizing inputs

# Normalizing training sets

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



Subtract mean:

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

$$x := x - \mu$$

Normalize variance

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n x^{(i)} * x^{(i)T}$$

$\hookrightarrow$  element-wise

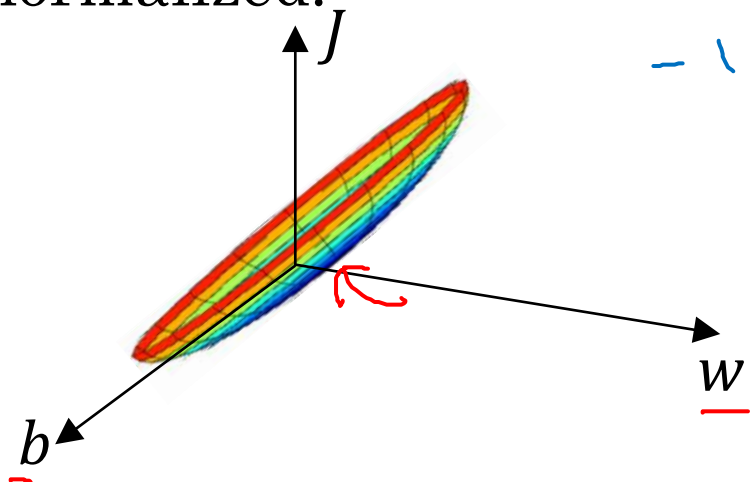
$$x /= \sigma^2$$

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

# Why normalize inputs?

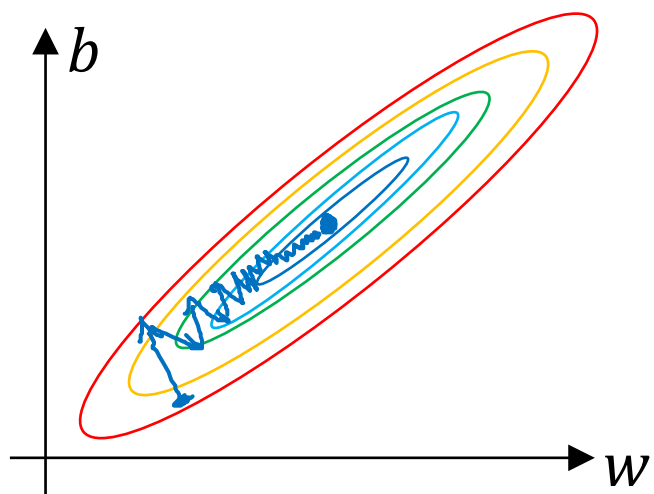
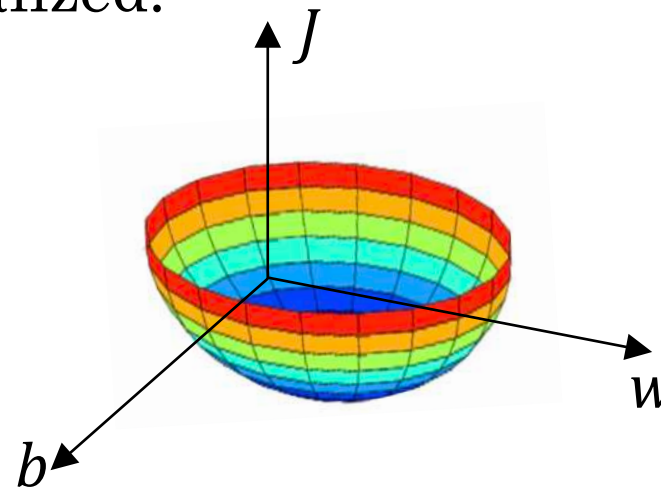
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Unnormalized:

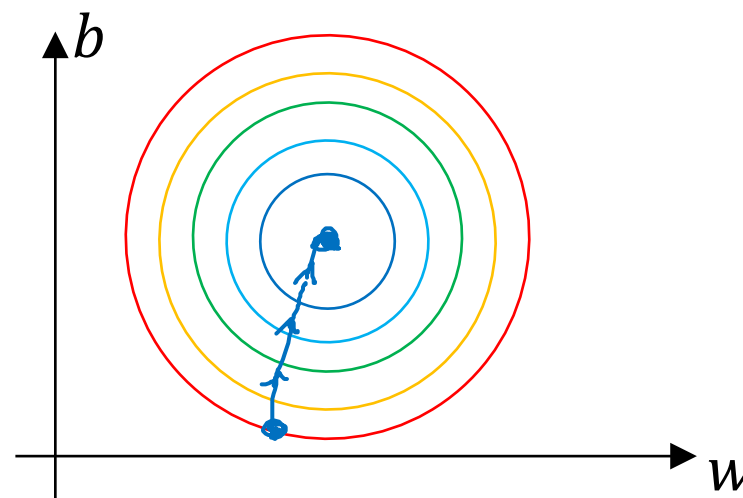


$w_1: x_1: \underline{1 \dots 1000} \leftarrow$   
 $w_2: x_2: \underline{0 \dots 1} \leftarrow$   
 $\quad \quad \quad -1 \dots 1$

Normalized:



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







deeplearning.ai

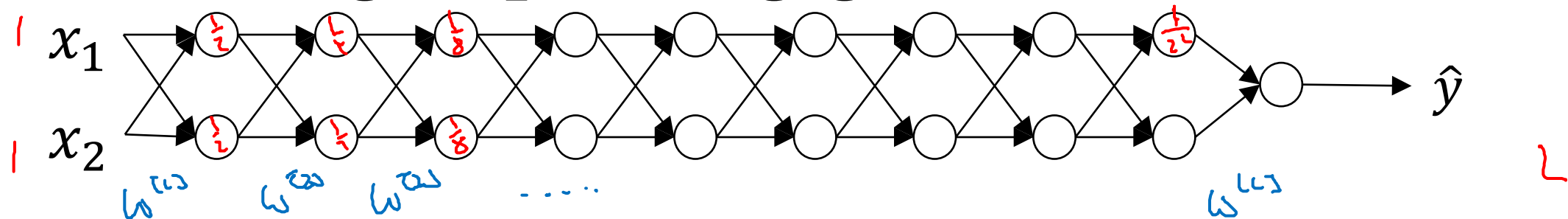
Setting up your  
optimization problem

---

Vanishing/exploding  
gradients

# Vanishing/exploding gradients

$L=150$



$g(z) = z$        $b^{(L)} = 0$

$\hat{y} = w^{(L,1)} \left( w^{(L-1,2)} w^{(L-2,2)} \dots \left( w^{(1,2)} w^{(1,1)} x \right) \right)$

$1.5^L$   
 $0.5^L$

$w^{(1,2)} > I$   
 $w^{(2,1)} < I$        $\begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$

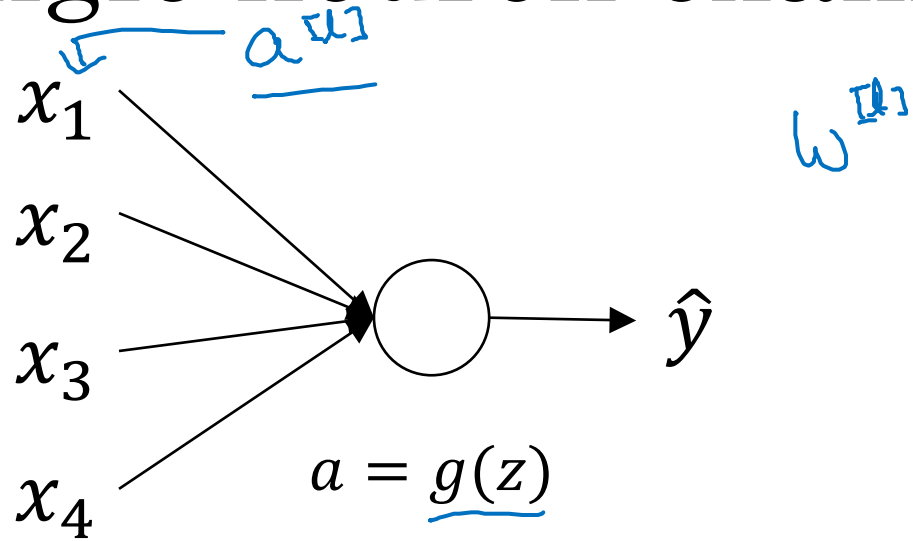
$w^{(2,1)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$

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

$\hat{y} = w^{(L,1)} \left[ \begin{matrix} 1.5 & 0 \\ 0 & 1.5 \end{matrix} \right]^{L-1} x$

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

# Single neuron example



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

large  $n \rightarrow$  Smaller  $w_i$

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

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

ReLU  $g^{[2]}(z) = \text{ReLU}(z)$

Other variants:

tanh

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

Xavier initialization ↑

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

↑



deeplearning.ai

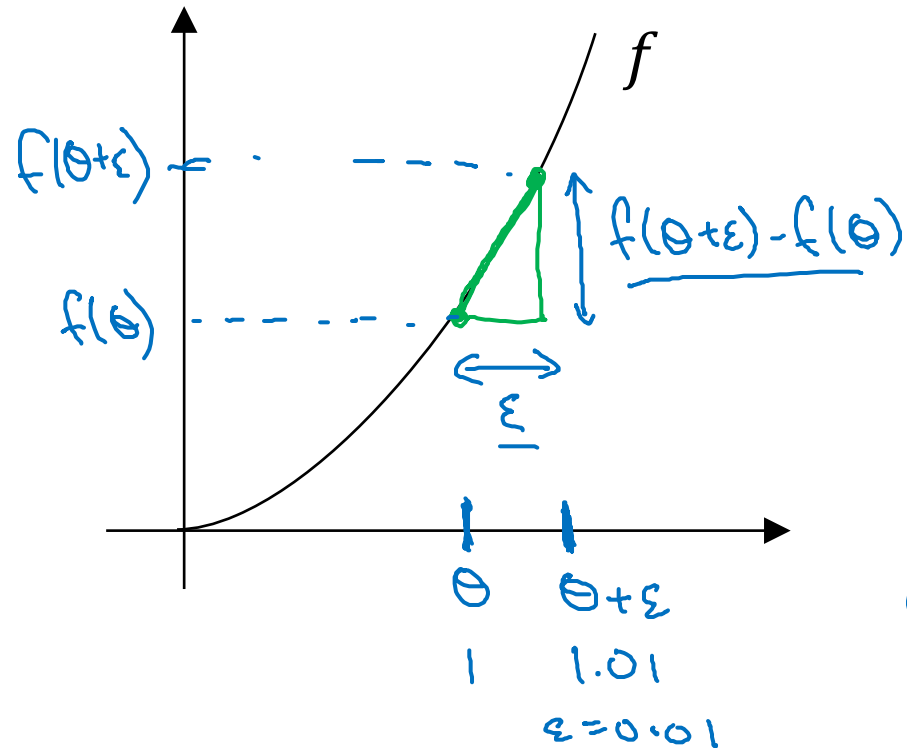
Setting up your  
optimization problem

---

Numerical approximation  
of gradients

# Checking your derivative computation

I  $f(\theta) = \theta^3$   
 $\theta \in \mathbb{R}.$



$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

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

$\frac{dw}{db}$

$g(\theta) = 3 \cdot (1)^2 = 3$   
 when  $\theta = 1$

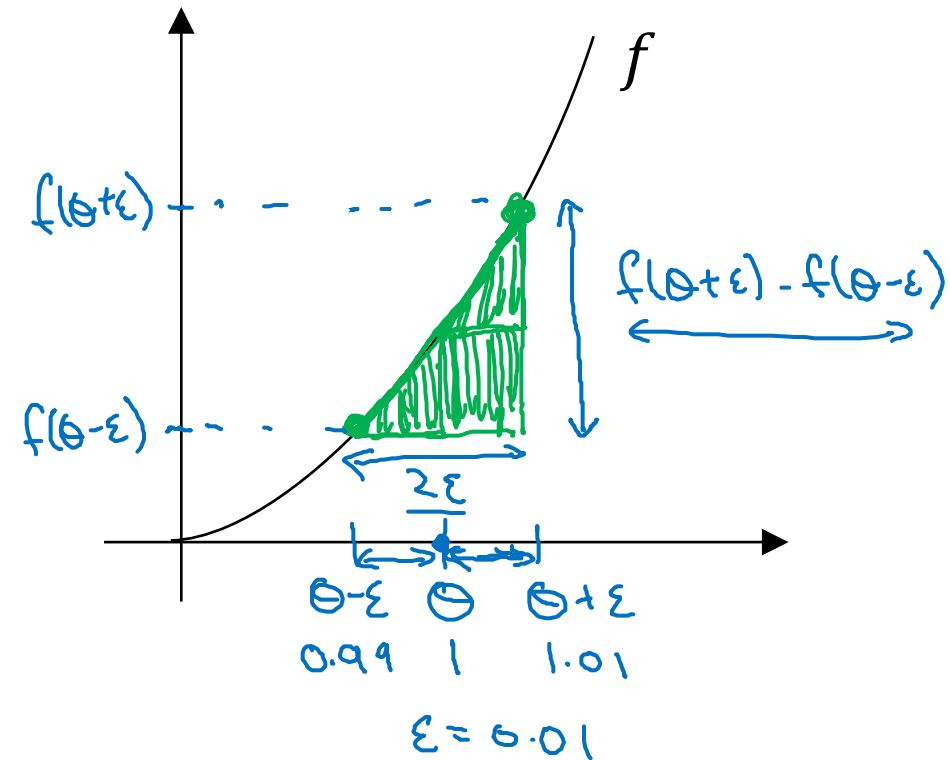
$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

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

Annotations:  $\theta = 1$ ,  $\theta + \epsilon = 1.01$ ,  $\epsilon = 0.01$ . The calculation shows  $3.1$  and  $3.2$  as intermediate steps in the division.

# Checking your derivative computation

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



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

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

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

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

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



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  $\theta$ .

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

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  $i$ :

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

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

Checks

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$
$$\underline{\epsilon = 10^{-7}}$$

$$\approx \frac{10^{-7}}{10^{-5}} - \text{great!} \leftarrow$$
$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$



deeplearning.ai

Setting up your  
optimization problem

---

Gradient Checking  
implementation notes

# Gradient checking implementation notes

- Don't use in training – only to debug

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

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

$$\frac{db^{[L]}}{\uparrow} \quad \frac{dW^{[L]}}{\uparrow}$$

- Remember regularization.

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

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

- Doesn't work with dropout.

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

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

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



deeplearning.ai

# Optimization Algorithms

---

Mini-batch  
gradient descent

# Batch vs. mini-batch gradient descent

$x, y$

$x^{\{t\}}, y^{\{t\}}$

Vectorization allows you to efficiently compute on  $m$  examples.

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_x, m)$        $\underbrace{\hspace{10em}}_{X^{\{1\}} \quad (n_x, 1000)}$        $\underbrace{\hspace{10em}}_{X^{\{2\}} \quad (n_x, 1000)}$        $\dots$        $\underbrace{\hspace{10em}}_{X^{\{5,000\}} \quad (n_x, 1000)}$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} & | & y^{(1001)} & \dots & y^{(2000)} & | & \dots & | & \dots & y^{(m)} \end{bmatrix}$$

$(1, m)$        $\underbrace{\hspace{10em}}_{Y^{\{1\}} \quad (1, 1000)}$        $\underbrace{\hspace{10em}}_{Y^{\{2\}} \quad (1, 1000)}$        $\dots$        $\underbrace{\hspace{10em}}_{Y^{\{5,000\}} \quad (1, 1000)}$

What if  $m = \underline{5,000,000}$ ?

5,000 mini-batches of 1,000 each

Mini-batch  $t$ :  $x^{\{t\}}, y^{\{t\}}$

$$\left| \begin{array}{l} x^{(i)} \\ z^{[l]} \\ x^{\{t\}}, y^{\{t\}} \end{array} \right.$$

# Mini-batch gradient descent

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

Forward prop on  $X^{\{t\}}$ .

$$Z^{(1)} = W^{(1)} X^{\{t\}} + b^{(1)}$$

$$A^{(1)} = g^{(1)}(Z^{(1)})$$

$\vdots$

$$A^{(L)} = g^{(L)}(Z^{(L)})$$

Vectorized implementation  
(1000 examples)

Compute cost  $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^L \ell(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{\ell} \|W^{(\ell)}\|_F^2$ .

↙ ↘ from  $X^{\{t\}}, Y^{\{t\}}$

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

$$W^{(1)} := W^{(1)} - \alpha dW^{(1)}, \quad b^{(1)} := b^{(1)} - \alpha db^{(1)}$$

}

"1 epoch"

pass through training set.

1 step of grad desc  
using  $X^{\{t+1\}}, Y^{\{t+1\}}$ .  
(as if  $m=1000$ )

$X, Y$



deeplearning.ai

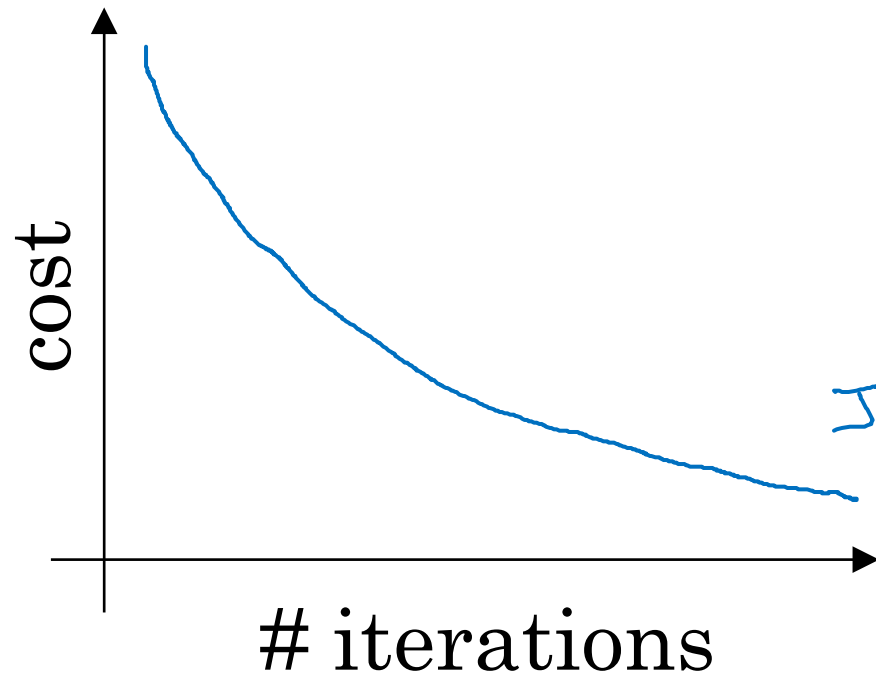
# Optimization Algorithms

---

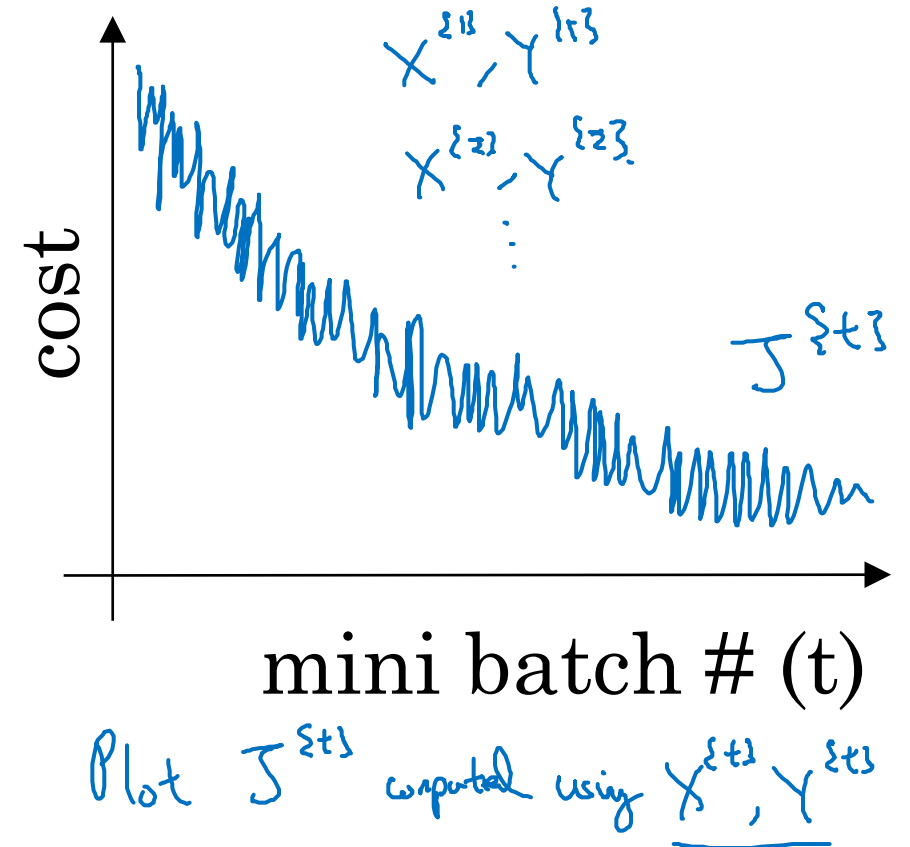
Understanding  
mini-batch  
gradient descent

# Training with mini batch gradient descent

Batch gradient descent



Mini-batch gradient descent





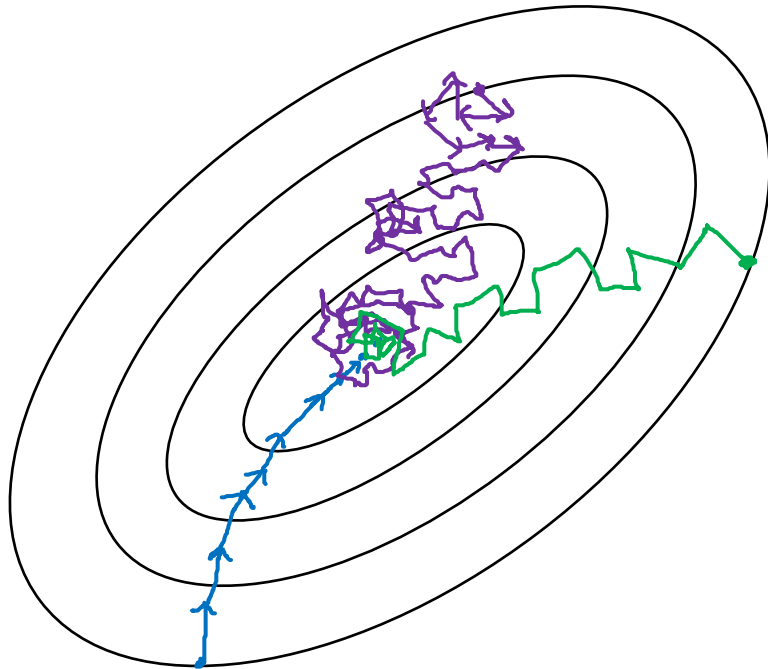
# Choosing your mini-batch size

→ If mini-batch size =  $m$  : Batch gradient descent.

$$(X^{(13)}, Y^{(13)}) = (X, Y).$$

→ If mini-batch size = 1 : Stochastic gradient descent. Every example is its own mini-batch.  
 $(X^{(13)}, Y^{(13)}) = (x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$  mini-batch.

In practice: Somewhere in-between 1 and  $m$



Stochastic  
gradient  
descent

Loss spikes  
from vectorization

In-between  
(mini-batch size  
not too big/small)

Fastest learning.

- Vectorization.  
( $n=1000$ )
- Make passes without  
processing entire training set.

Batch  
gradient descent  
(mini-batch size =  $m$ )

Too long  
per iteration

# Choosing your mini-batch size

If small toy set : Use batch gradient descent.  
( $m \leq 2000$ )

Typical mini-batch sizes:

→ 64, 128, 256, 512, 1024  
 $2^6, 2^7, 2^8, 2^9, 2^{10}$

Make sure mini-batch fit in CPU/GPU memory.  
 $X^{(t)}, Y^{(t)}$



deeplearning.ai

# Optimization Algorithms

---

## Exponentially weighted averages

# Temperature in London

$$\theta_1 = 40^\circ\text{F} \quad 4^\circ\text{C} \quad \leftarrow$$

$$\theta_2 = 49^\circ\text{F} \quad 9^\circ\text{C}$$

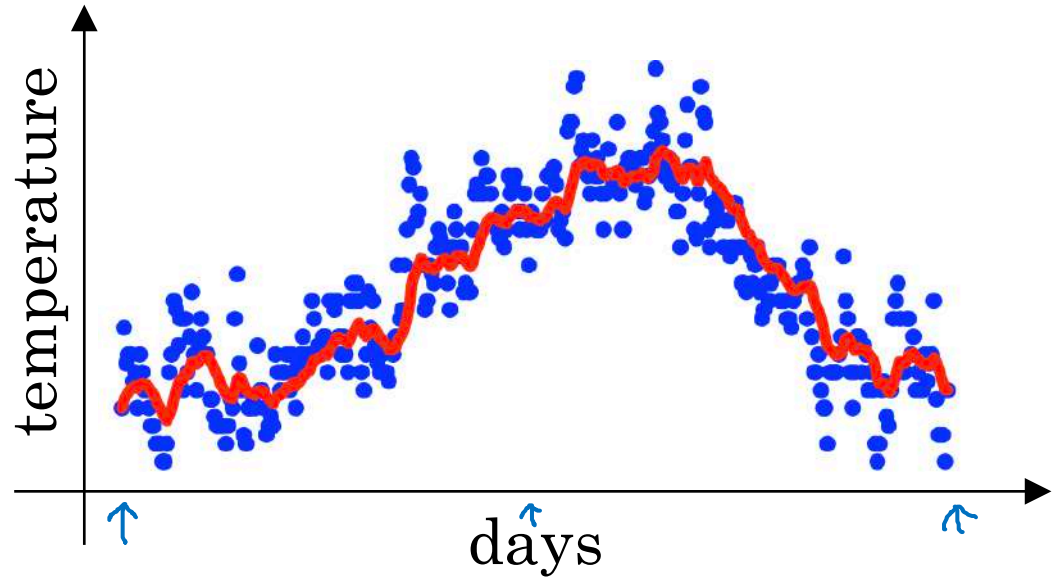
$$\theta_3 = 45^\circ\text{F} \quad \vdots$$

$\vdots$

$$\theta_{180} = 60^\circ\text{F} \quad 15^\circ\text{C}$$

$$\theta_{181} = 56^\circ\text{F} \quad \vdots$$

$\vdots$



$$V_0 = 0$$

$$V_1 = 0.9 V_0 + 0.1 \theta_1$$

$$V_2 = 0.9 V_1 + 0.1 \theta_2$$

$$V_3 = 0.9 V_2 + 0.1 \theta_3$$

$\vdots$

$$V_t = 0.9 V_{t-1} + 0.1 \theta_t$$

# Exponentially weighted averages <sup>moving</sup>

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

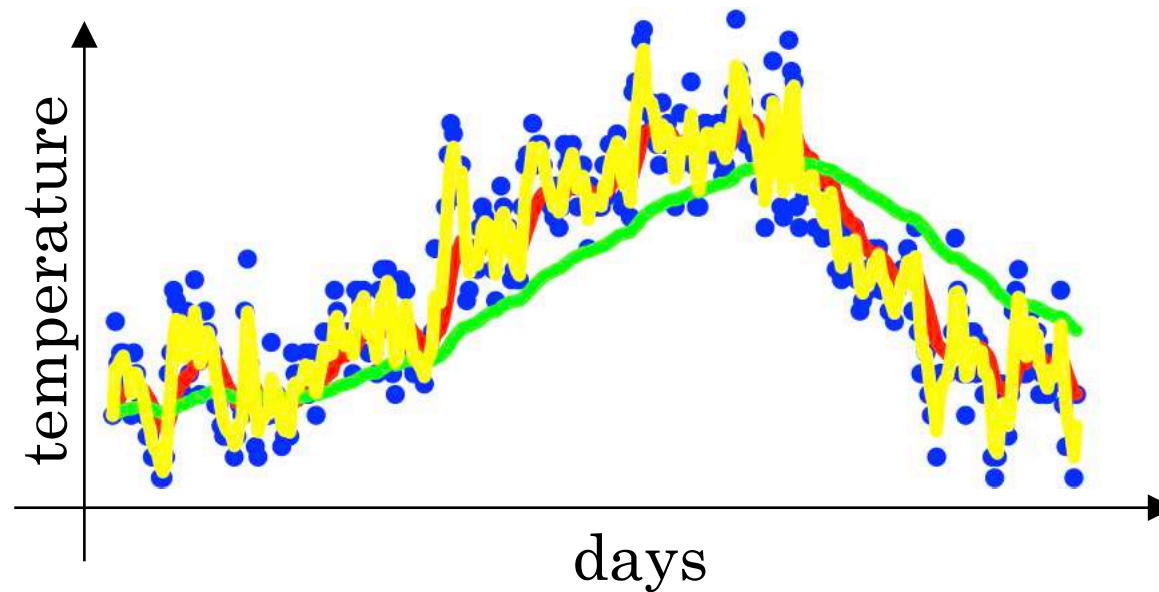
$\beta = 0.9$  :  $\approx 10$  days' temperature.

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

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

$V_t$  is approximately  
average over  
 $\rightarrow \approx \frac{1}{1-\beta}$  days' temperature.

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





deeplearning.ai

# Optimization Algorithms

---

Understanding  
exponentially  
weighted averages

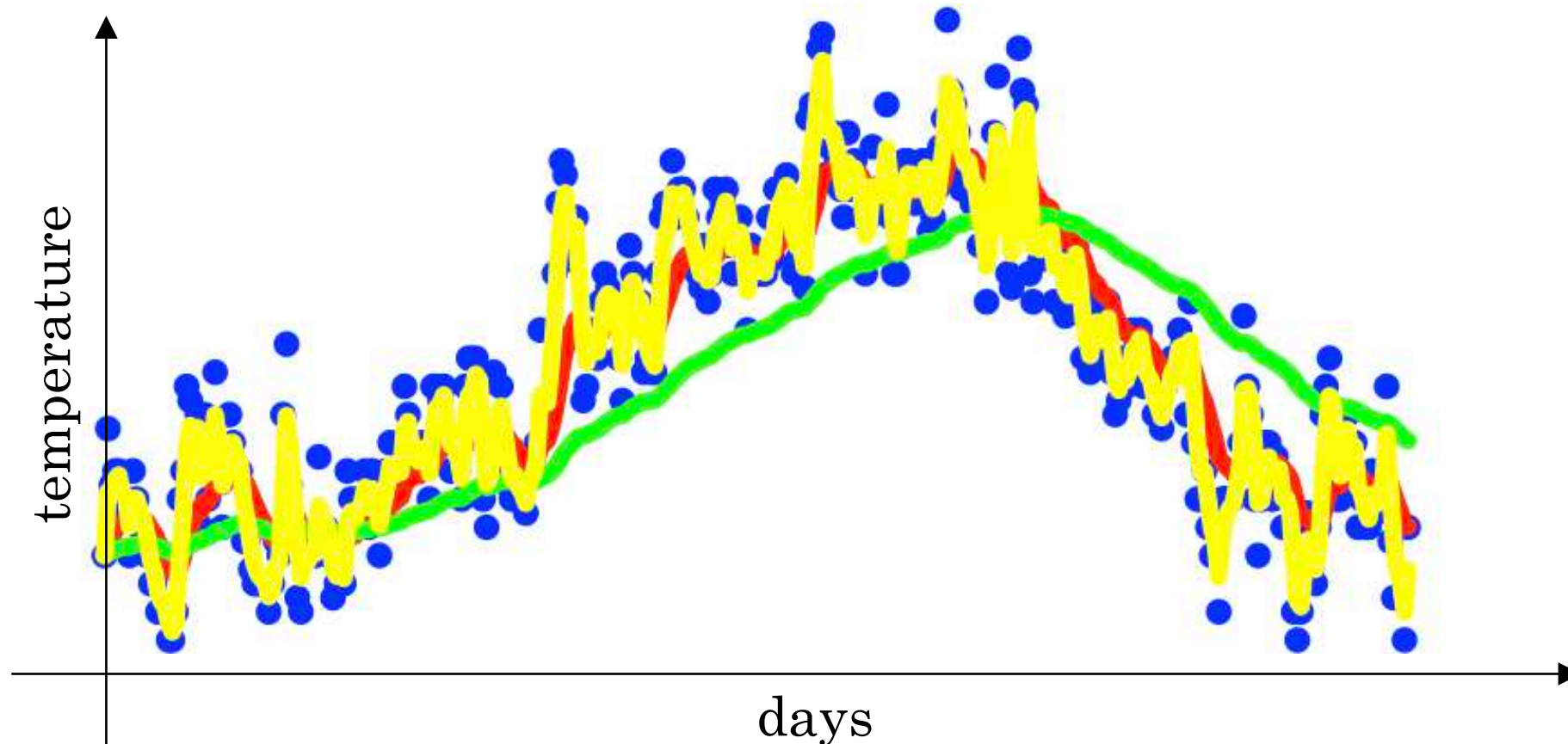
# Exponentially weighted averages

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

$$\beta = 0.9$$

$$0.98$$

$$0.5$$



# Exponentially weighted averages

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

$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

$$v_{99} = 0.9v_{98} + 0.1\theta_{99}$$

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

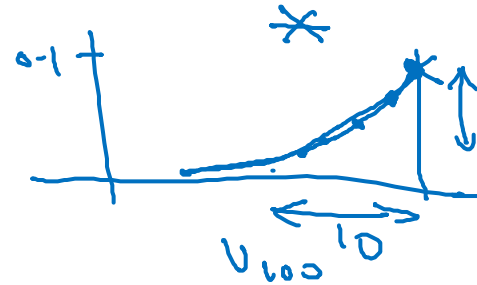
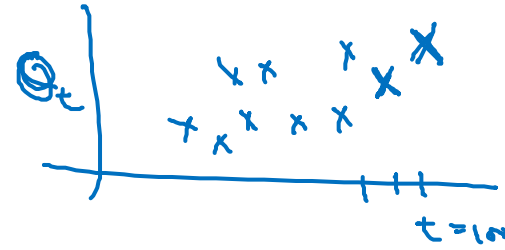
...

$$\begin{aligned} \rightarrow v_{100} &= 0.1\theta_{100} + 0.9 \cancel{v_{99}} (0.1\theta_{99} + 0.9 \cancel{v_{98}}) \\ &= \underbrace{0.1\theta_{100}} + \underbrace{0.1 \times 0.9 \cdot \theta_{99}} + \underbrace{0.1 (0.9)^2 \theta_{98}} + \underbrace{0.1 (0.9)^3 \theta_{97}} + \underbrace{0.1 (0.9)^4 \theta_{96}} + \dots \end{aligned}$$

$$\underbrace{0.9^{10}} \approx \underbrace{0.35} \approx \frac{1}{e}$$

$$\frac{(1-\epsilon)^{1/\epsilon}}{0.9} \approx \frac{1}{e}$$

$$\epsilon = 0.02 \rightarrow \underbrace{0.98^{50}} \approx \frac{1}{e}$$



$$\approx \frac{1}{1-\beta}$$

$$\epsilon = 1 - \beta$$

$$0.1\theta_{99} + 0.9v_{99}$$



# Implementing exponentially weighted averages

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

...

$$V_\theta := 0$$

$$V_\theta := \beta v + (1 - \beta) \theta_1$$

$$V_\theta := \beta v + (1 - \beta) \theta_2$$

⋮

---

$$\rightarrow V_\theta = 0$$

Repeat {

Get next  $\theta_t$

$$V_\theta := \beta V_\theta + (1 - \beta) \theta_t \leftarrow$$

}



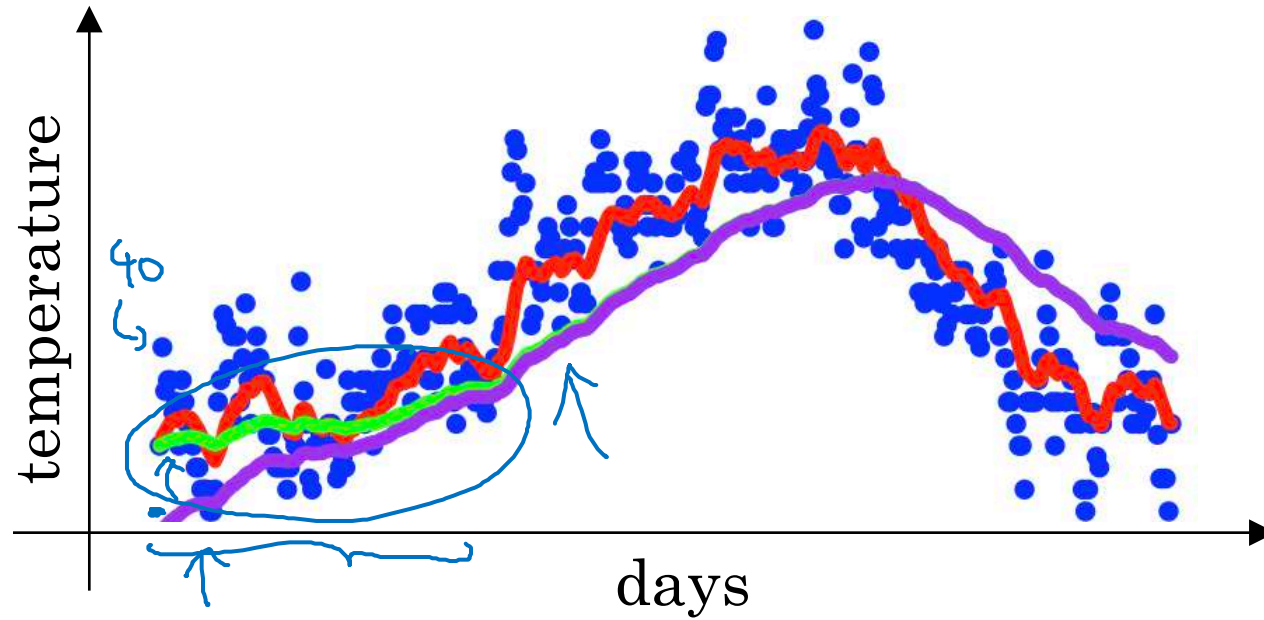
deeplearning.ai

# Optimization Algorithms

---

Bias correction  
in exponentially  
weighted average

# Bias correction



$$\beta = 0.98$$

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

$$v_0 = 0$$

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

$$v_2 = 0.98 v_1 + 0.02 \theta_2$$

$$= 0.98 \times 0.02 \times \theta_1 + 0.02 \theta_2$$

$$= \underline{0.0196} \theta_1 + \underline{0.02} \theta_2$$

$$\frac{v_t}{1 - \beta^t}$$

$$t=2: 1 - \beta^t = 1 - (0.98)^2 = 0.0396$$

$$\frac{v_2}{0.0396} = \frac{0.0196 \theta_1 + 0.02 \theta_2}{0.0396}$$



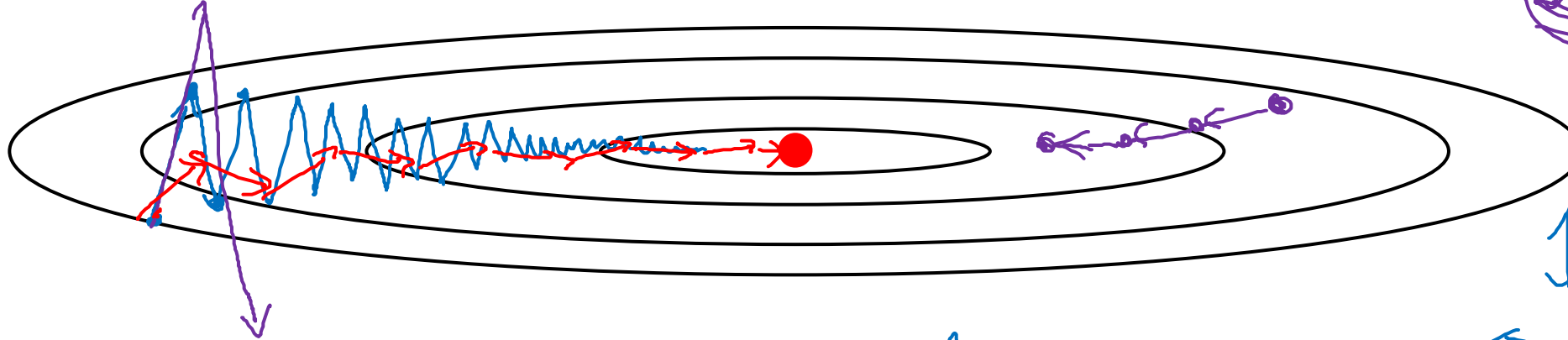
deeplearning.ai

# Optimization Algorithms

---

## Gradient descent with momentum

# Gradient descent example



↑ slower learning  
↔ faster learning

Momentum:

On iteration  $t$ :

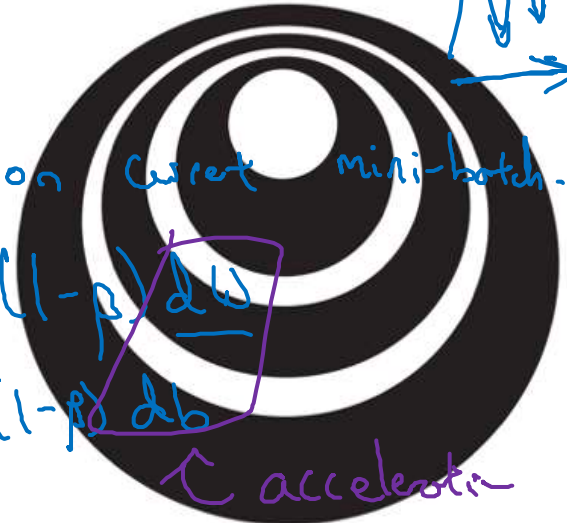
Compute  $\Delta W, \Delta b$  on current mini-batch.

$$V_{\Delta W} = \beta V_{\Delta W} + (1-\beta) \Delta W$$

$$V_{\Delta b} = \beta V_{\Delta b} + (1-\beta) \Delta b$$

friction — ↑ velocity

$$W := W - \alpha V_{\Delta W}, \quad b := b - \alpha V_{\Delta b}$$



$$V_{\theta} = \beta V_{\theta} + (1-\beta) \theta_t$$

deeplearning.ai

# Implementation details

$$v_{dW} = 0, \quad v_{db} = 0$$

On iteration  $t$ :

Compute  $dW, db$  on the current mini-batch

$$\left. \begin{aligned} \rightarrow v_{dW} &= \beta v_{dW} + (1 - \beta) dW \\ \rightarrow v_{db} &= \beta v_{db} + (1 - \beta) db \end{aligned} \right\} \quad \left| \quad \underbrace{v_{dW} = \beta v_{dW} + dW}_{\leftarrow}$$

$$W = W - \underbrace{\alpha v_{dW}}, \quad b = \underline{b} - \underbrace{\alpha v_{db}}$$

$$\frac{v_{dW}}{1 - \beta^t}$$

Hyperparameters:  $\alpha, \beta$

$$\underline{\beta = 0.9}$$

average over last  $\approx 10$  gradients



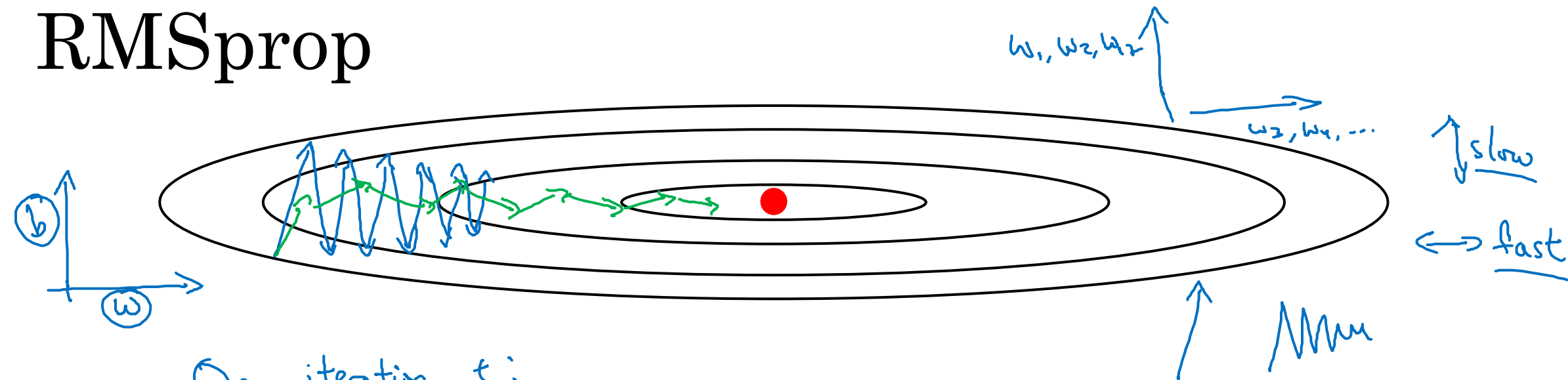
deeplearning.ai

# Optimization Algorithms

---

## RMSprop

# RMSprop



On iteration  $t$ :

Compute  $dw, db$  on current mini-batch

$$\underline{S_{dw}} = \beta_2 \underline{S_{dw}} + (1 - \beta_2) \underline{dw^2} \leftarrow \text{small}$$

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

$$w := w - \alpha \frac{dw}{\sqrt{S_{dw} + \epsilon}}$$

$$b := b - \alpha \frac{db}{\sqrt{S_{db} + \epsilon}}$$

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





deeplearning.ai

# Optimization Algorithms

---

## Adam optimization algorithm

# Adam optimization algorithm

$$V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$$

On iteration  $t$ :

Compute  $dw, db$  using current mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1 - \beta_1) dw, \quad V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \quad \leftarrow \text{"momentum"} \beta_1$$

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

yhat = np.array([.9, 0.2, 0.1, .4, .9])

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

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

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

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

# Hyperparameters choice:

→  $\alpha$  : needs to be tune

→  $\beta_1$  : 0.9 → ( $dw$ )

→  $\beta_2$  : 0.999 → ( $dw^2$ )

→  $\epsilon$  :  $10^{-8}$

Adam : Adaptive moment estimation



Adam Coates



deeplearning.ai

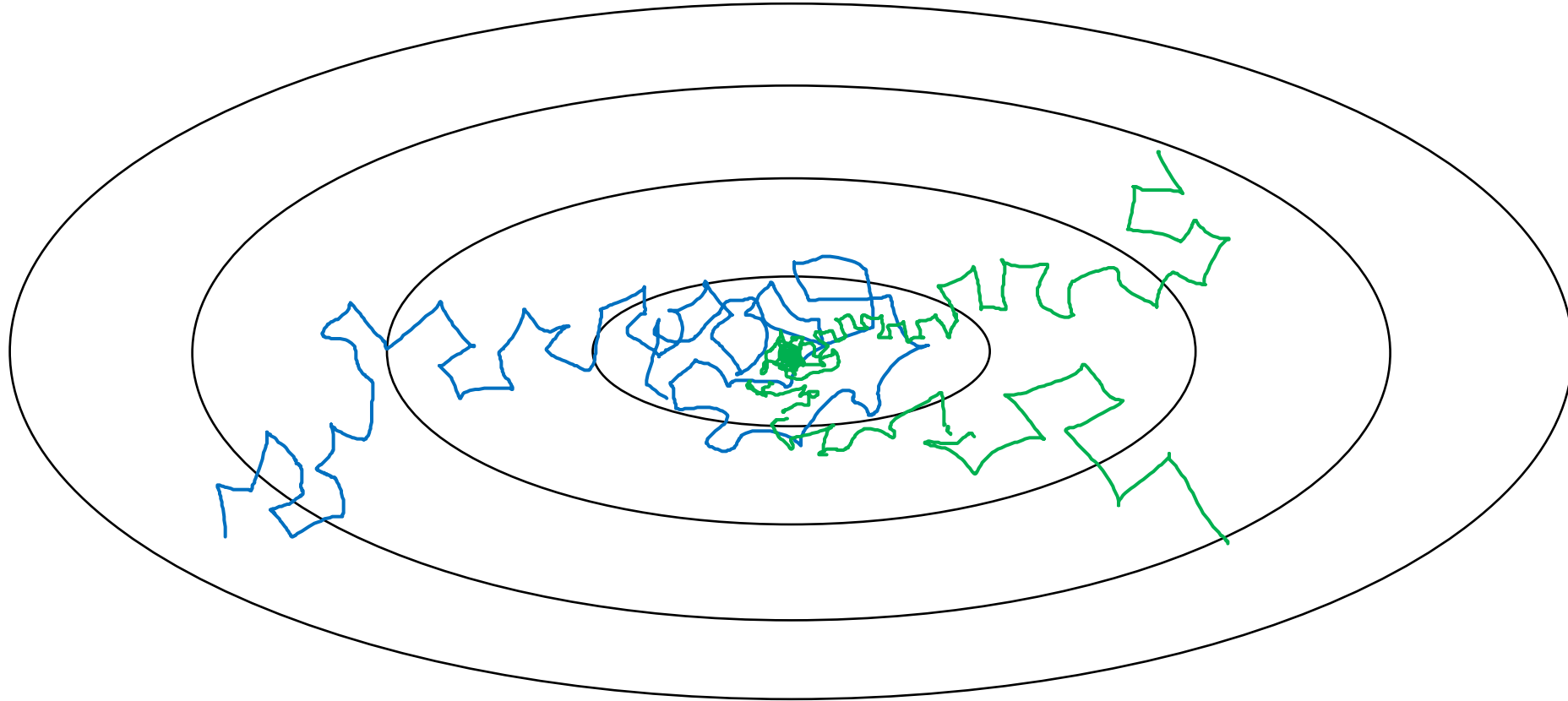
# Optimization Algorithms

---

## Learning rate decay

# Learning rate decay

Slowly reduce  $\alpha$

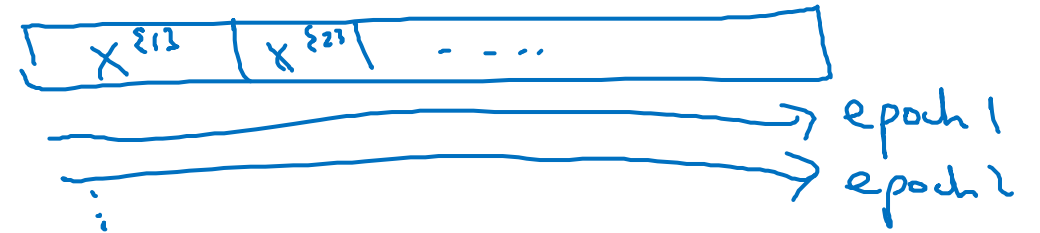


# Learning rate decay

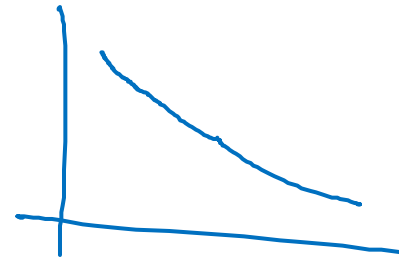
1 epoch = 1 pass through data.

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

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



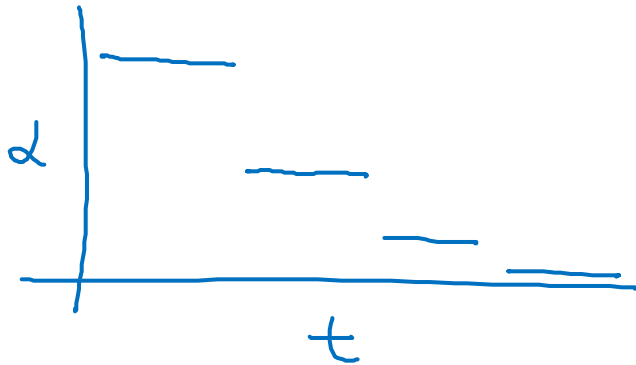
$$\alpha_0 = 0.2$$
$$\text{decay-rate} = 1$$



# Other learning rate decay methods

formula {  $\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0$  — exponentially decay.

$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0$  or  $\frac{k}{\sqrt{t}} \cdot \alpha_0$



discrete staircase

Manual decay.



deeplearning.ai

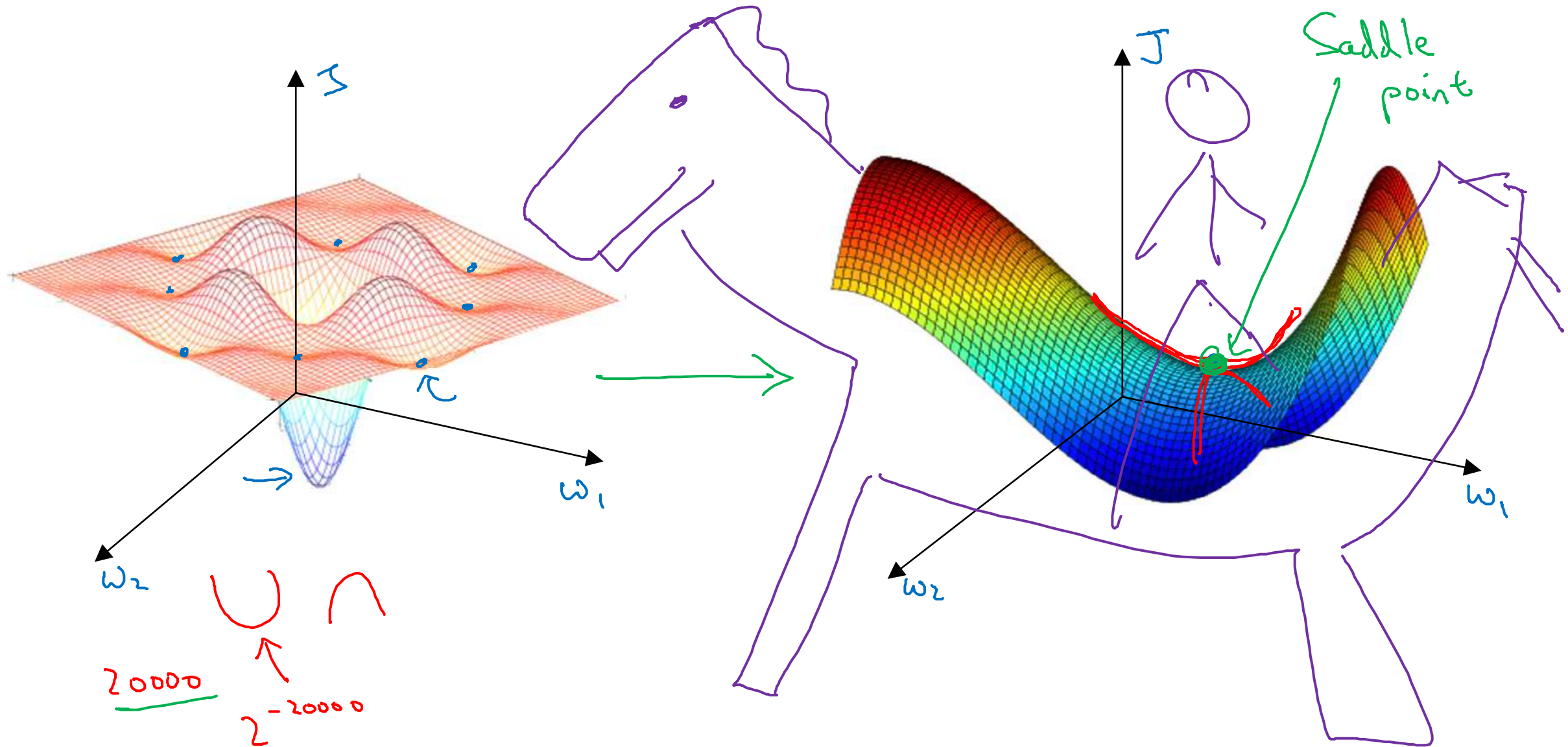
# Optimization Algorithms

---

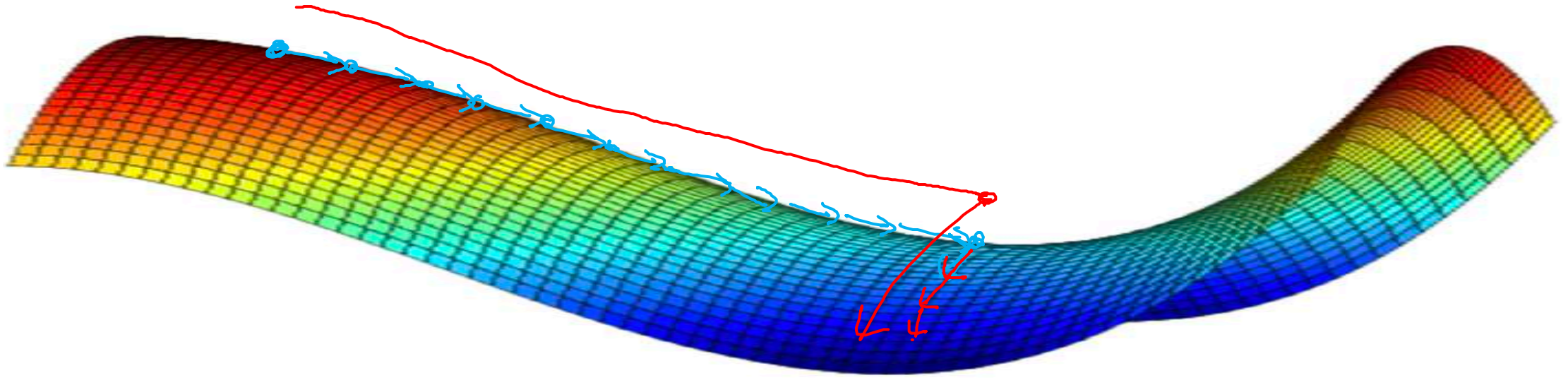
## The problem of local optima



# Local optima in neural networks



# Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow



deeplearning.ai

# Hyperparameter tuning

---

## Tuning process

# Hyperparameters

→  $\alpha$

$\beta$  0.9

$\beta_1, \beta_2, \epsilon$   
0.9 0.999  $10^{-8}$

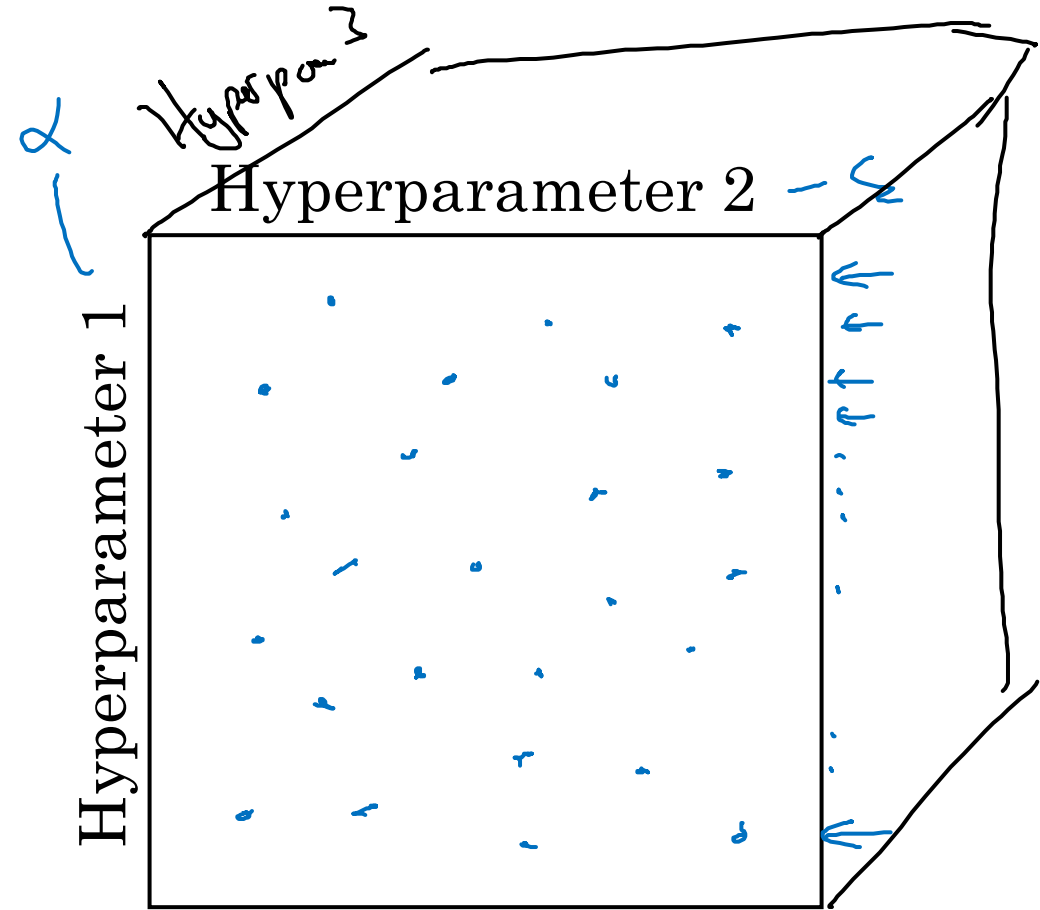
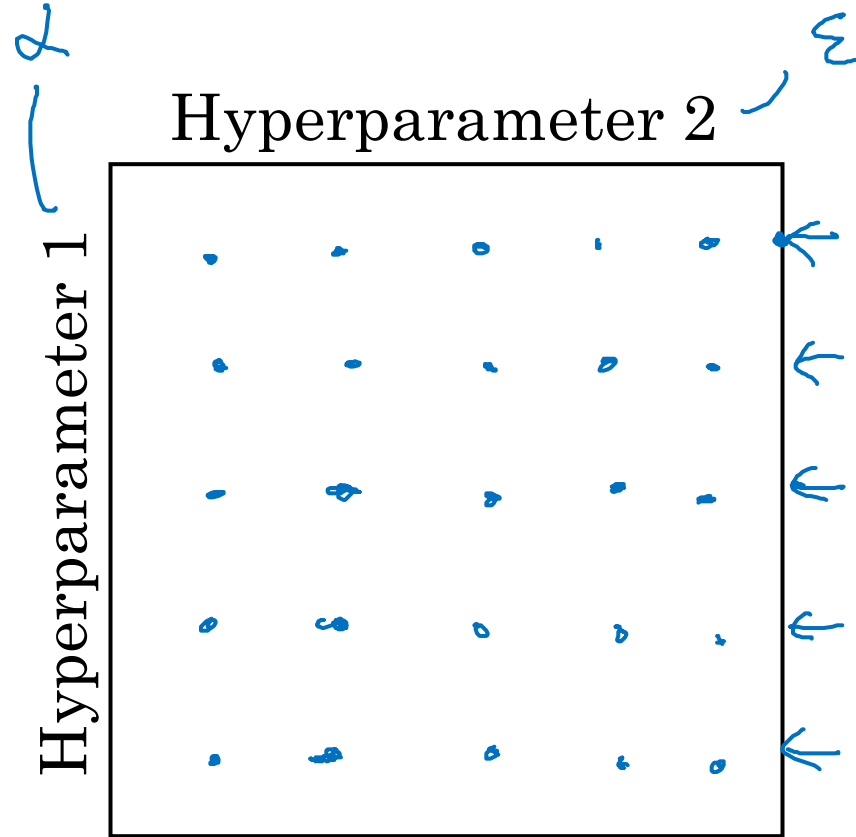
# layers

# hidden units

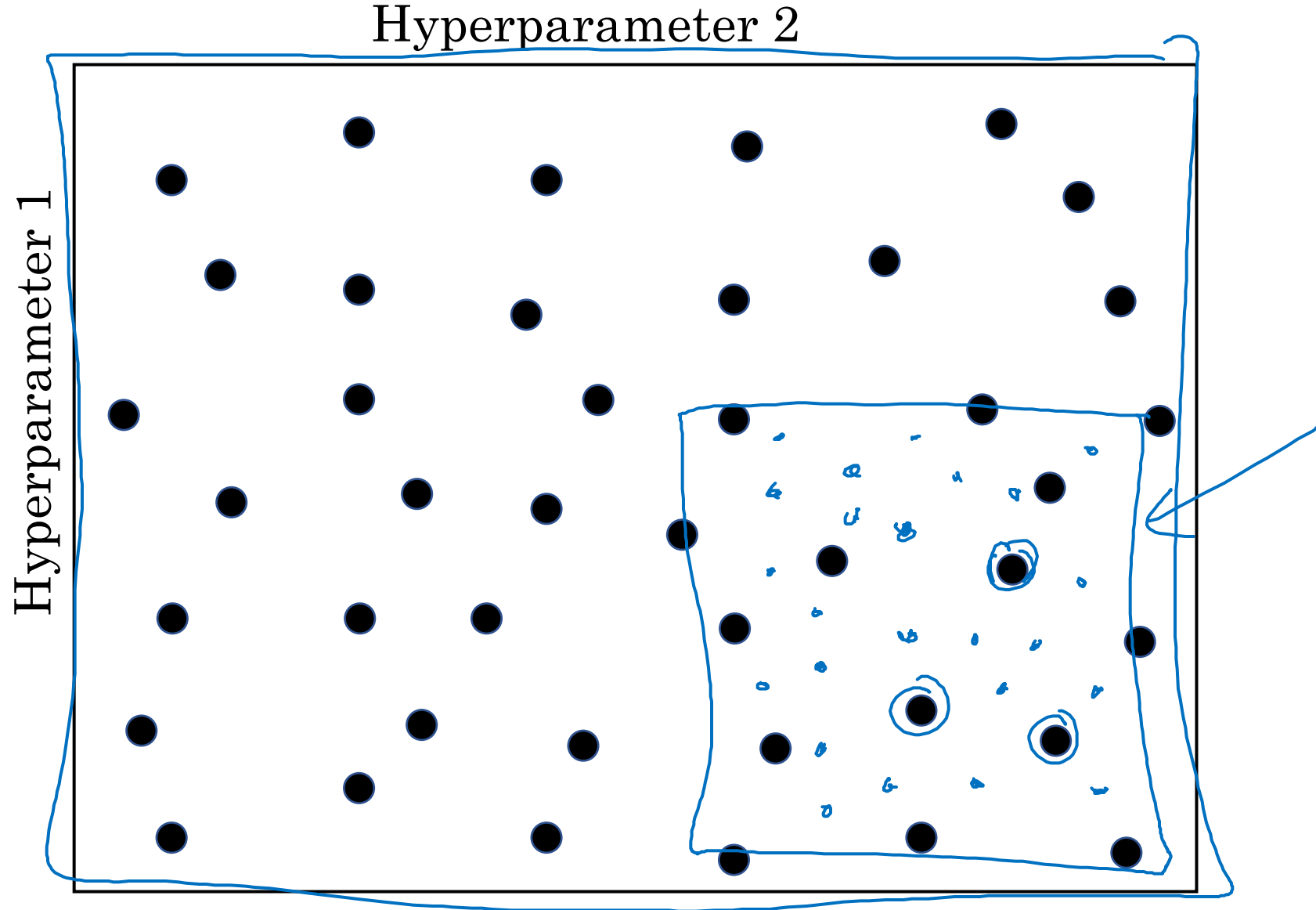
learning rate decay

mini-batch size

# Try random values: Don't use a grid



# Coarse to fine





deeplearning.ai

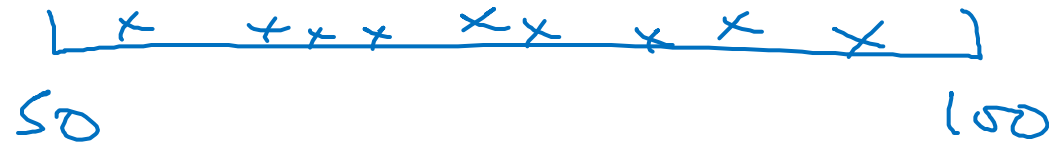
# Hyperparameter tuning

---

Using an appropriate  
scale to pick  
hyperparameters

# Picking hyperparameters at random

→  $n^{\text{test}} = 50, \dots, 100$



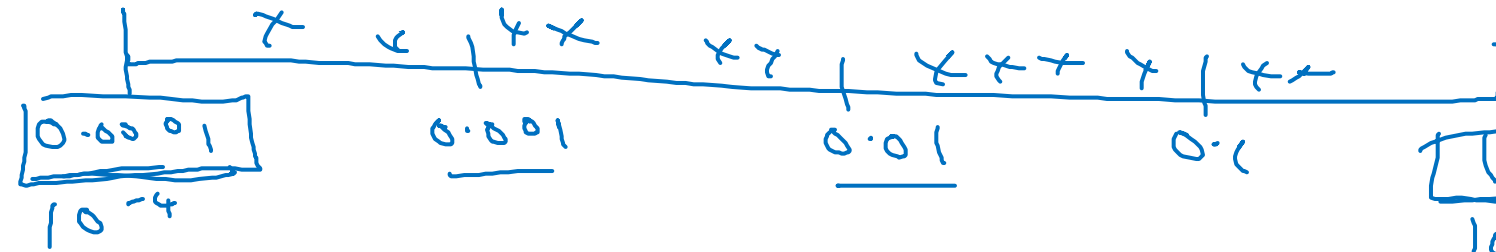
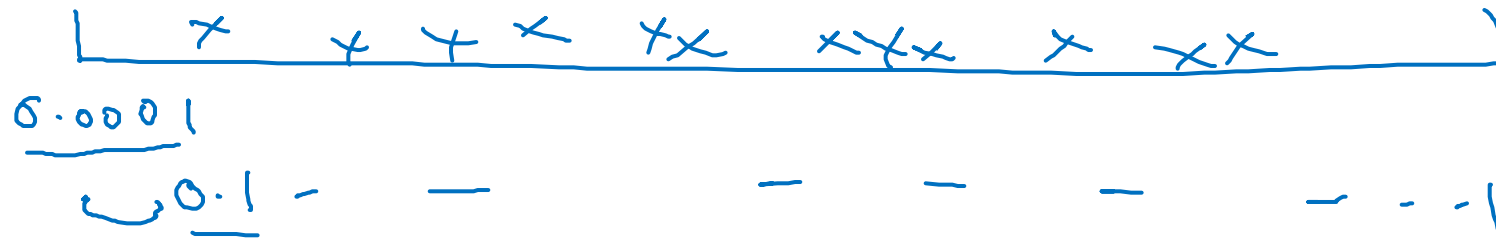
→ #layers      L :    2 - 4

2 , 3 , 4



# Appropriate scale for hyperparameters

$$\alpha = 0.0001, \dots, 1$$



$$a = \log_{10} 0.0001 = -4$$

$$r = -4 * \text{np.random.rand}()$$

$$\alpha = 10^r$$

$$r \in [-4, 0]$$

$$\alpha = 10^{-4} \dots 10^0$$

$$b = \log_{10} 1 = 0$$

$$\underline{10^{-4} \dots 10^0}$$

$$\underline{\frac{r \in [a, b]}{[-4, 0]}}$$

$$\underline{\alpha = 10^r}$$

# Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \quad \dots \quad 0.999$$

$\downarrow$   
10

$\downarrow$   
1000

$$1 - \beta = 0.1 \quad \dots \quad 0.001$$

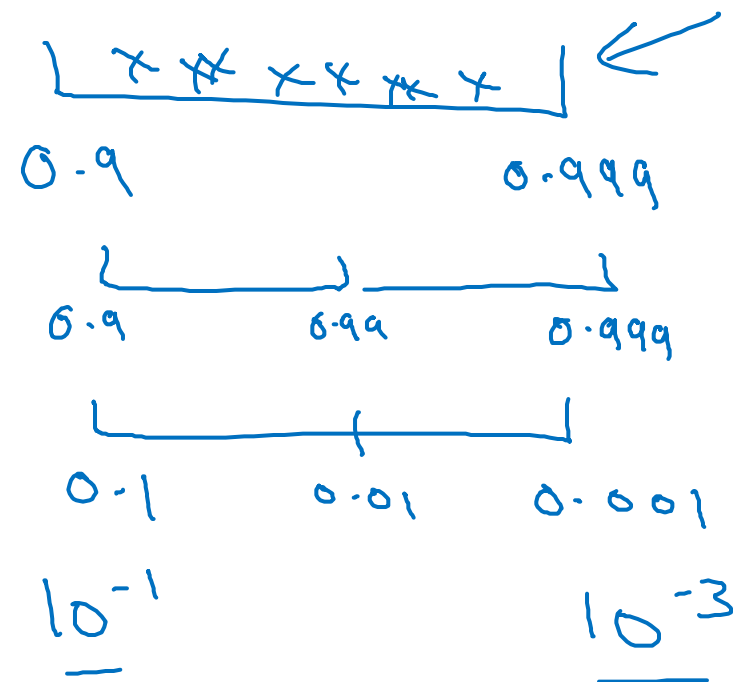

---

$$\beta: 0.999 \rightarrow 0.9995 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

$\sim 1000$ 
 $\sim 2000$

$$\frac{1}{1 - \beta_K}$$



$$r \in [-3, -1]$$

$$1 - \beta = 10^r$$

$$\beta = 1 - 10^r$$



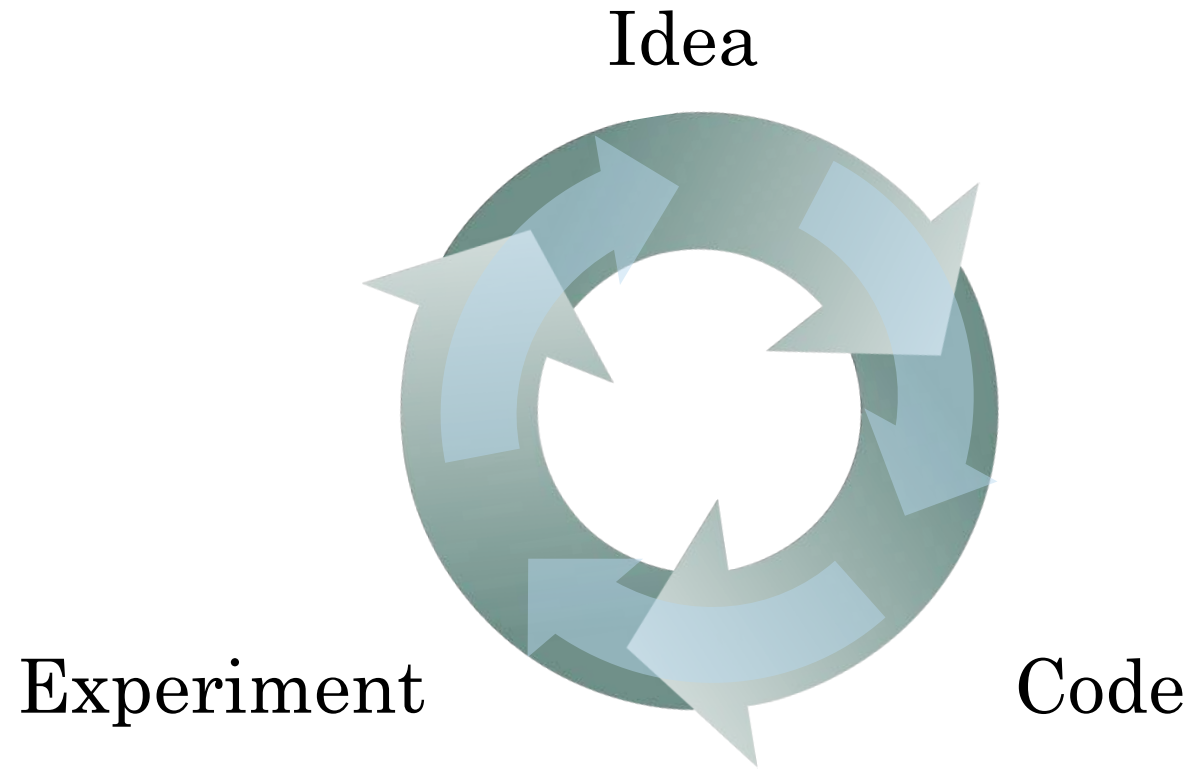
deeplearning.ai

# Hyperparameters tuning

---

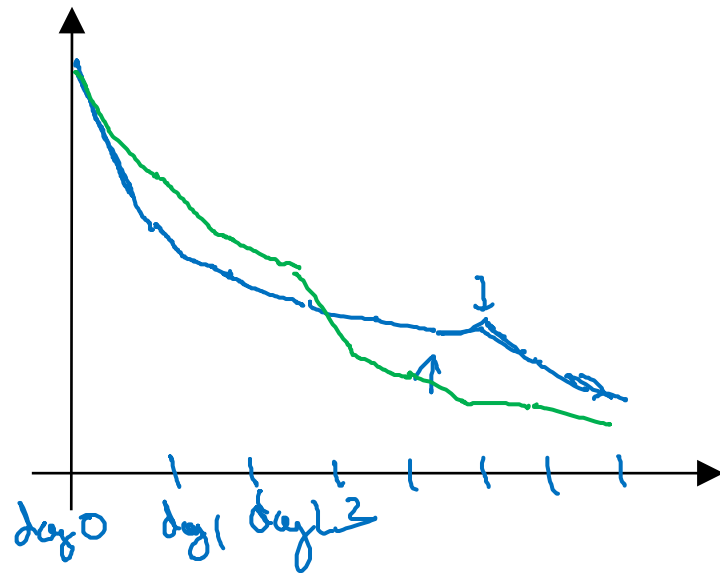
Hyperparameters  
tuning in practice:  
Pandas vs. Caviar

# Re-test hyperparameters occasionally



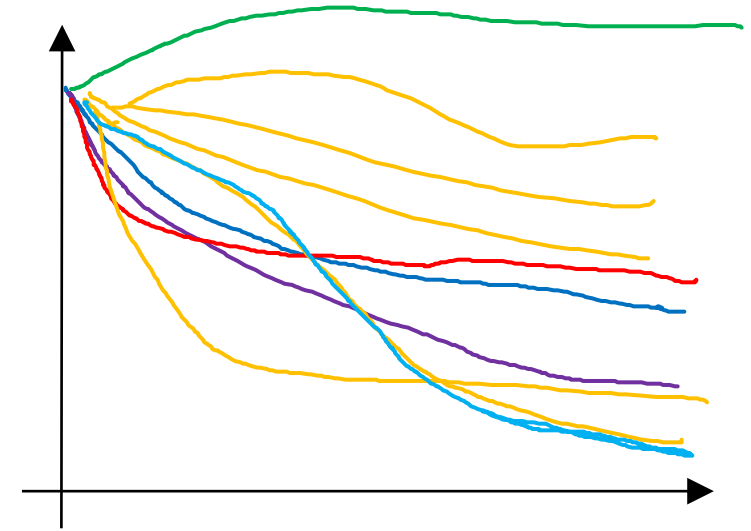
- NLP, Vision, Speech,  
Ads, logistics, ....
- Intuitions do get stale.  
Re-evaluate occasionally.

# Babysitting one model



Panda ←

# Training many models in parallel



Caviar ←



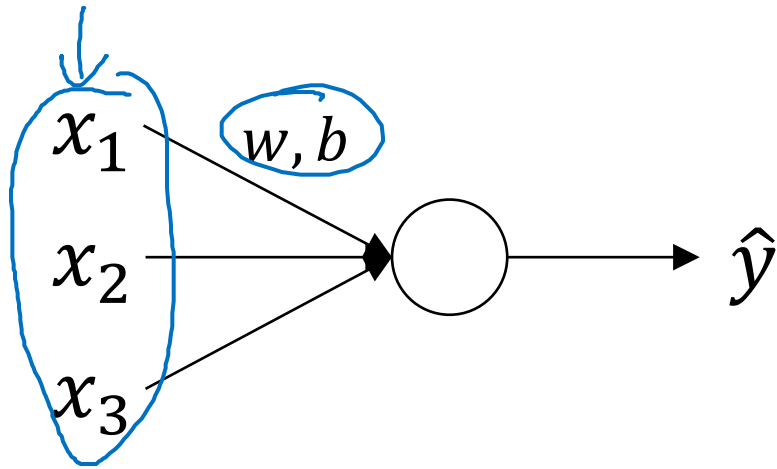
deeplearning.ai

# Batch Normalization

---

Normalizing activations  
in a network

# Normalizing inputs to speed up learning

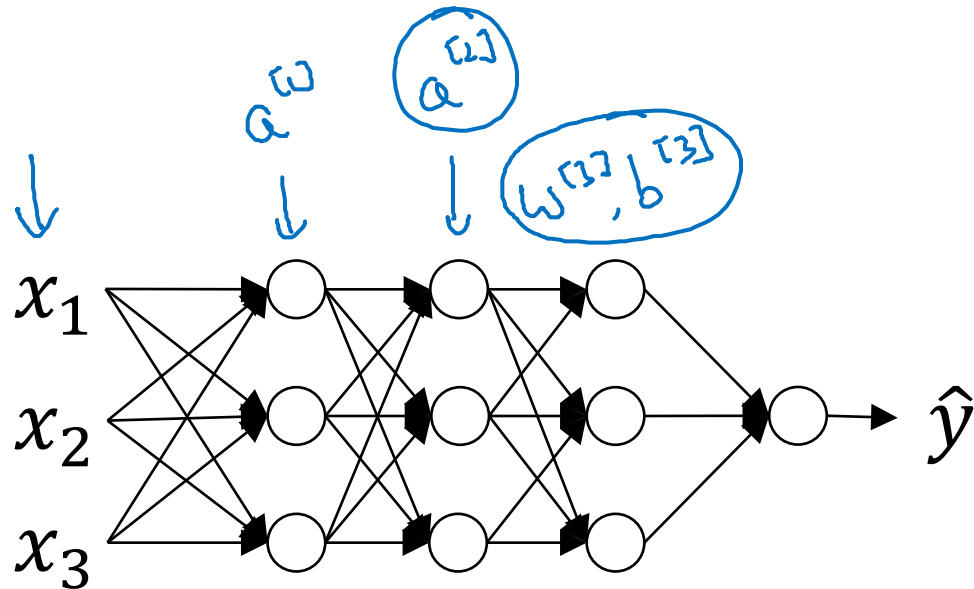
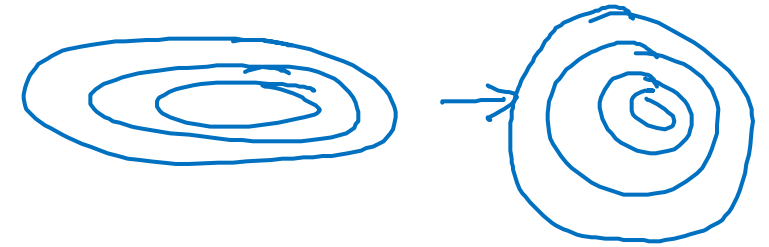


$$\mu = \frac{1}{n} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

$$X = X / \sigma^2$$



Can we normalize  $\frac{a^{[2]}}{w^{[2]}, b^{[2]}}$  so as to train faster

Normalize  $\frac{z^{[2]}}{\uparrow}$

# Implementing Batch Norm

Given some intermediate values in NN

$z^{(1)}, \dots, z^{(m)}$   
 $z^{[l]}(i)$

$$\begin{aligned} \mu &= \frac{1}{m} \sum_i z^{(i)} \\ \sigma^2 &= \frac{1}{m} \sum_i (z^{(i)} - \mu)^2 \\ z_{\text{norm}}^{(i)} &= \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}} \\ \hat{z}^{(i)} &= \gamma z_{\text{norm}}^{(i)} + \beta \end{aligned}$$

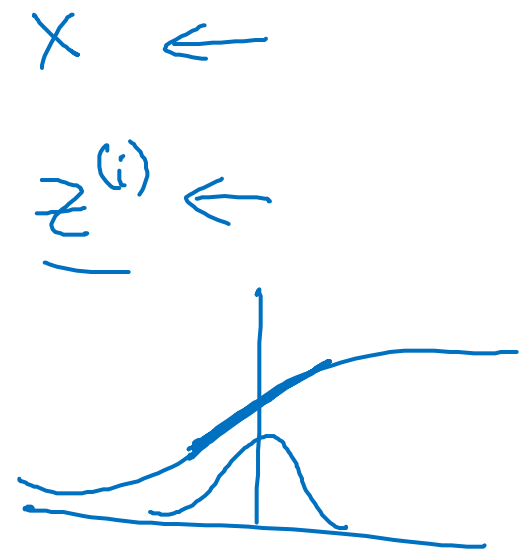
If

$$\gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

then  $\hat{z}^{(i)} = z^{(i)}$

learnable parameters of model.



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





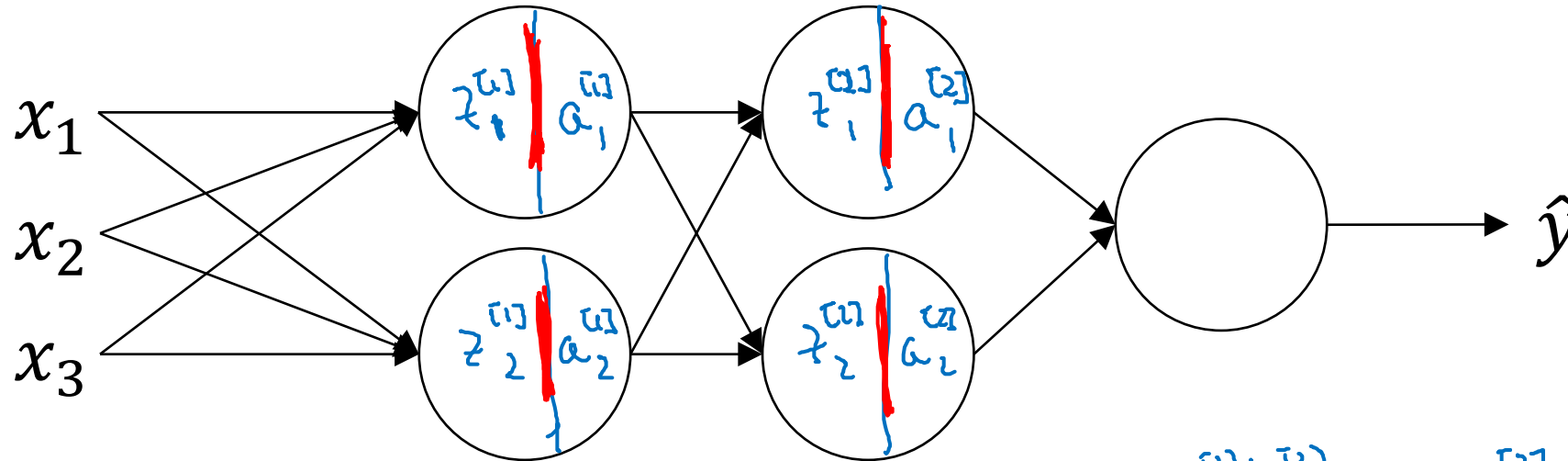
deeplearning.ai

# Batch Normalization

---

Fitting Batch Norm  
into a neural network

# Adding Batch Norm to a network



$$X \xrightarrow{W^{(1)}, b^{(1)}} \underline{z^{(1)}} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{(1)}, \gamma^{(1)}} \underline{z^{(1)}} \xrightarrow{W^{(2)}, b^{(2)}} \underline{z^{(2)}} \xrightarrow[\text{BN}]{\beta^{(2)}, \gamma^{(2)}} \underline{z^{(2)}} \rightarrow a^{(2)} \rightarrow \dots$$

$a = g(z)$

Parameters:  $\left\{ W^{(1)}, b^{(1)}, W^{(2)}, b^{(2)}, \dots, W^{(L)}, b^{(L)} \right\}$

$\rightarrow \underline{\beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)}}$

$\rightarrow \underline{\beta}$

$$d\beta^{(2)} \quad \beta = \beta - \alpha d\beta^{(2)}$$

tf.nn.batch-normalization ←

# Working with mini-batches

$$\underline{X^{[1]}} \xrightarrow{W^{[1]}, b^{[1]}} \underline{z^{[1]}} \xrightarrow[\text{BN}]{\beta^{[1]}, \gamma^{[1]}} \underline{\tilde{z}^{[1]}} \rightarrow g^{[1]}(\tilde{z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} \underline{z^{[2]}} \rightarrow \dots$$

$$\boxed{X^{[2]}} \rightarrow \underline{z^{[2]}} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \underline{\tilde{z}^{[2]}} \rightarrow \dots$$

$$X^{[3]} \rightarrow \dots$$

Parameters:  $W^{[1]}, \cancel{b^{[1]}}, \beta^{[1]}, \gamma^{[1]}$

$\uparrow$   $(n^{[1]}, 1)$      $\uparrow$   $(n^{[1]}, 1)$      $\uparrow$   $(n^{[1]}, 1)$

$\tilde{z}^{[1]}_{(n^{[1]}, 1)}$

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

$$z^{[2]} = W^{[2]} a^{[1]}$$

$$z^{[2]}_{\text{norm}}$$

$$\rightarrow \tilde{z}^{[2]} = \gamma^{[2]} z^{[2]}_{\text{norm}} + \boxed{\beta^{[2]}}$$

# Implementing gradient descent

for  $t = 1 \dots \text{num Mini Batches}$

Compute forward pass on  $X^{\{t\}}$ .

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

Use backprop to compute  $\underline{dw}^{\{t\}}$ ,  ~~$\underline{db}^{\{t\}}$~~ ,  $\underline{dp}^{\{t\}}$ ,  $\underline{d\delta}^{\{t\}}$

Update params 
$$\left. \begin{aligned} w^{\{t\}} &:= w^{\{t-1\}} - \alpha dw^{\{t\}} \\ \beta^{\{t\}} &:= \beta^{\{t-1\}} - \alpha dp^{\{t\}} \\ \gamma^{\{t\}} &:= \dots \end{aligned} \right\} \leftarrow$$

Works w/ momentum, RMSprop, Adam.



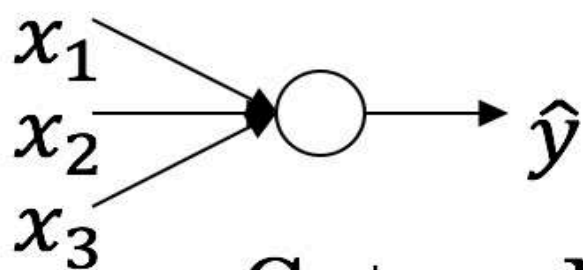
deeplearning.ai

# Batch Normalization

---

Why does  
Batch Norm work?

# Learning on shifting input distribution

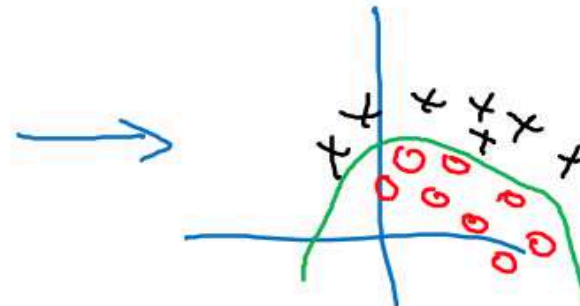
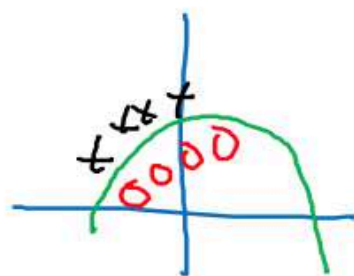
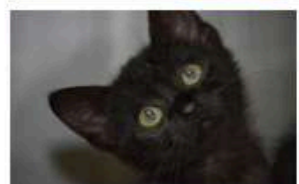
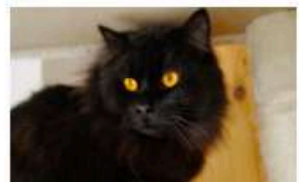
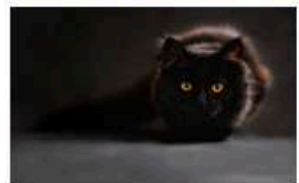


Cat

Non-Cat

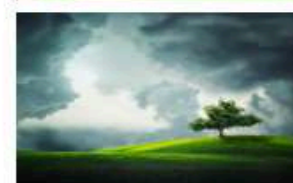
$y = 1$  ✓

$y = 0$



$y = 1$  ✓

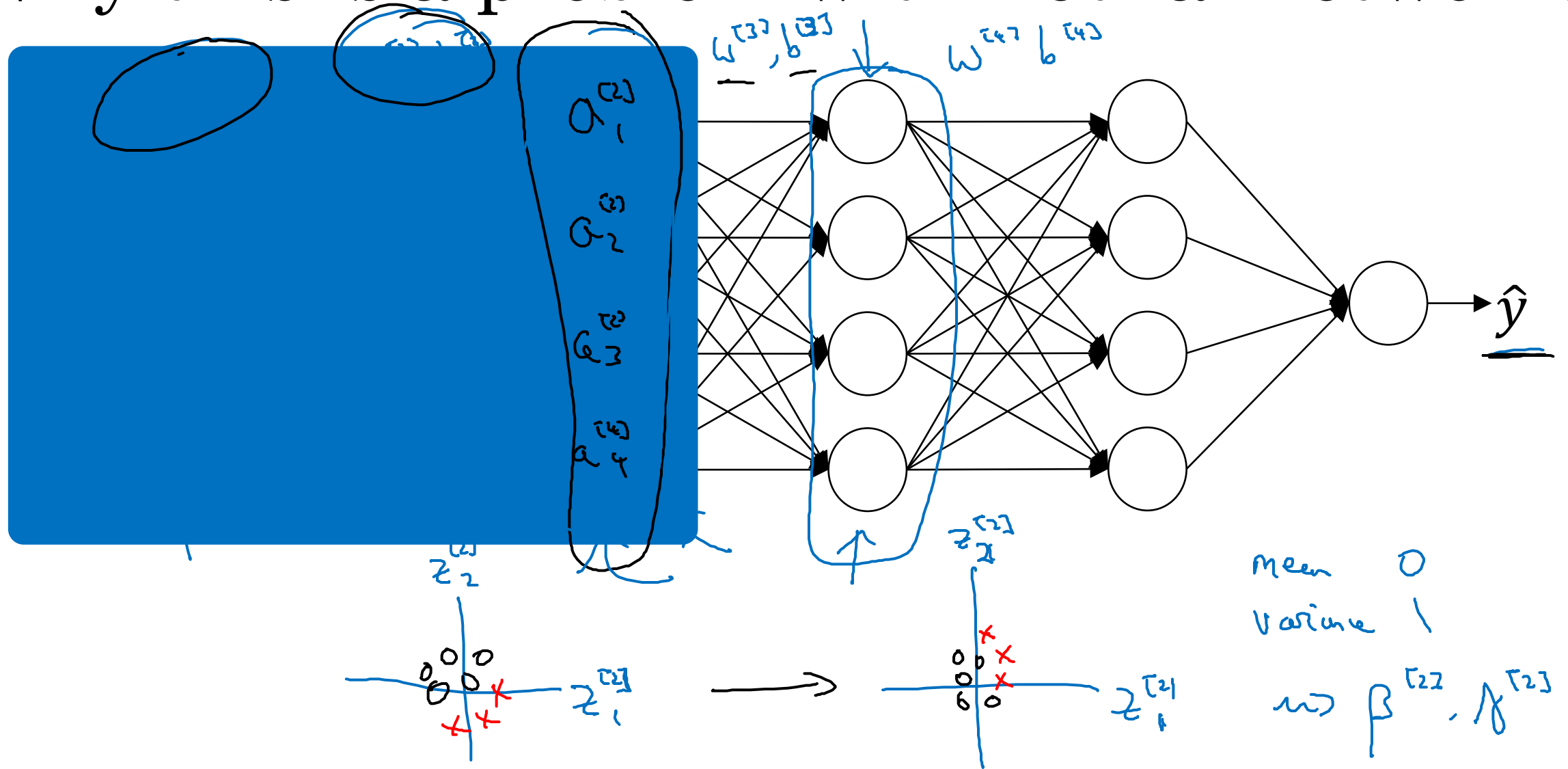
$y = 0$



"Covariate shift"

$\underline{x} \rightarrow y$

# Why this is a problem with neural networks?



# Batch Norm as regularization

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.
- This adds some noise to the values  $z^{[l]}$  within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.
- This has a slight regularization effect.

mini-batch : 64  $\longrightarrow$  512





deeplearning.ai

# Batch Normalization

---

## Batch Norm at test time

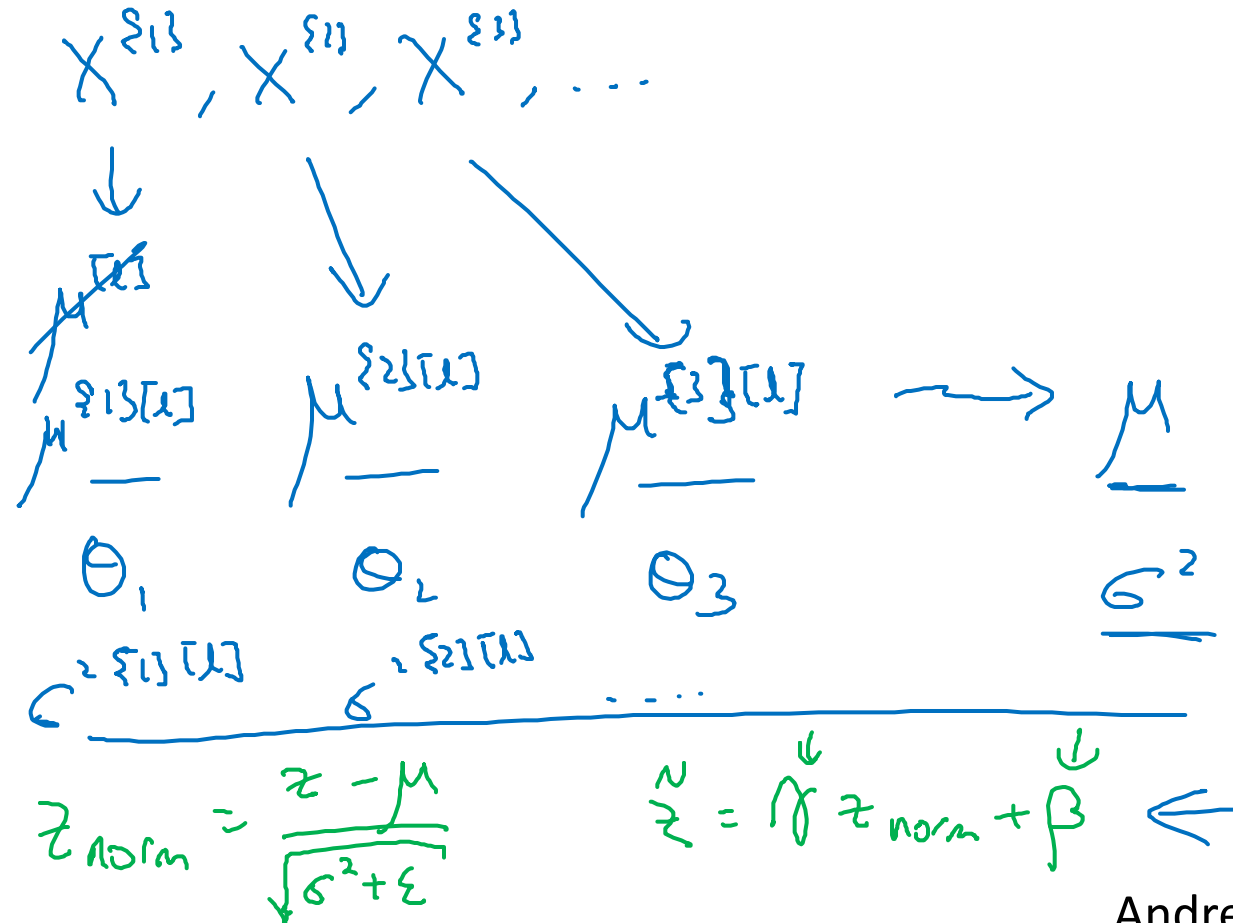
# Batch Norm at test time

$$\begin{aligned} \rightarrow \underline{\mu} &= \frac{1}{\underline{m}} \sum_i \underline{z^{(i)}} \\ \rightarrow \underline{\sigma^2} &= \frac{1}{\underline{m}} \sum_i (\underline{z^{(i)}} - \underline{\mu})^2 \end{aligned}$$

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

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

$\underline{\mu}, \underline{\sigma^2}$ : estimate using exponentially weighted average (across mini-batches).





deeplearning.ai

# Multi-class classification

---

## Softmax regression

# Recognizing cats, dogs, and baby chicks



3



1



2



0



3



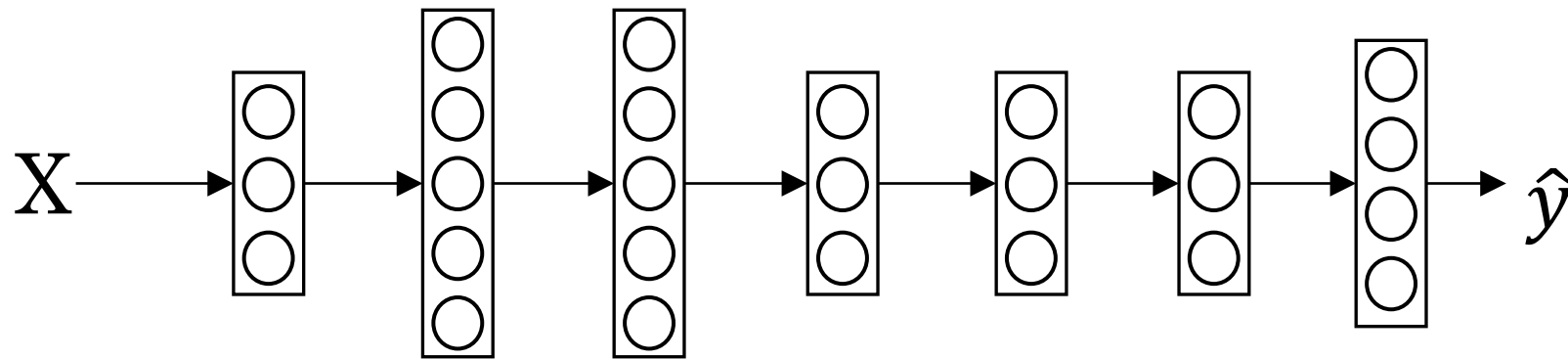
2



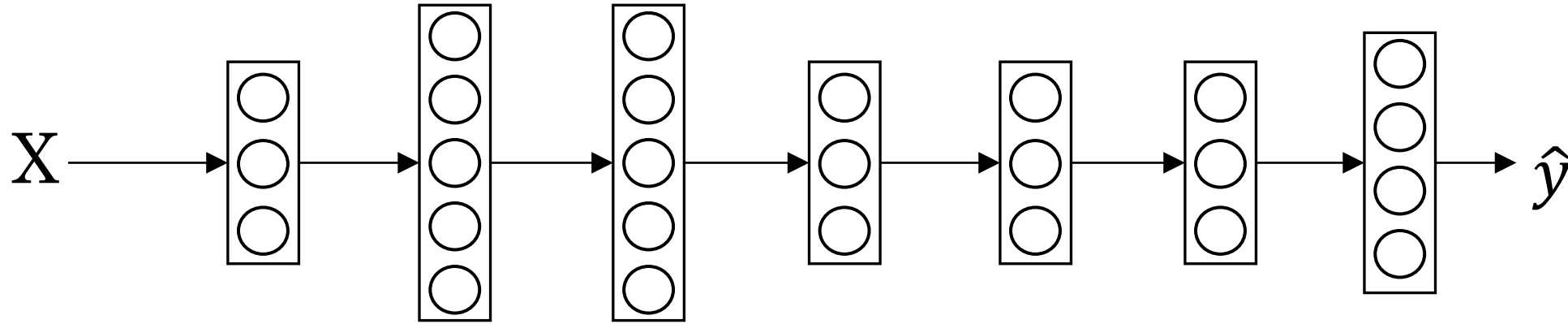
0



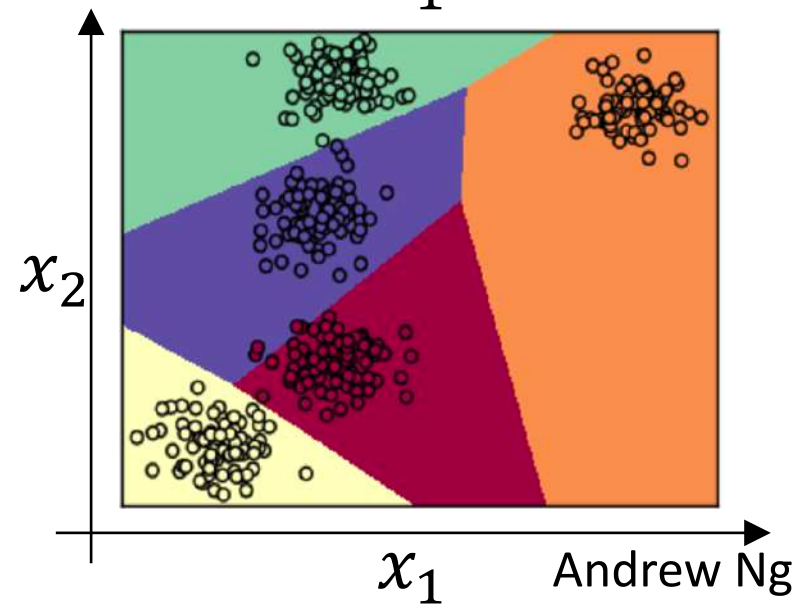
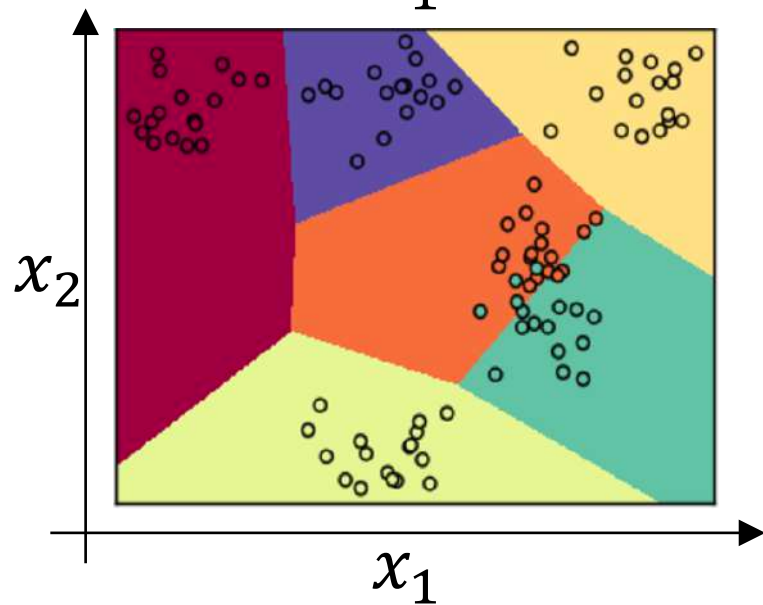
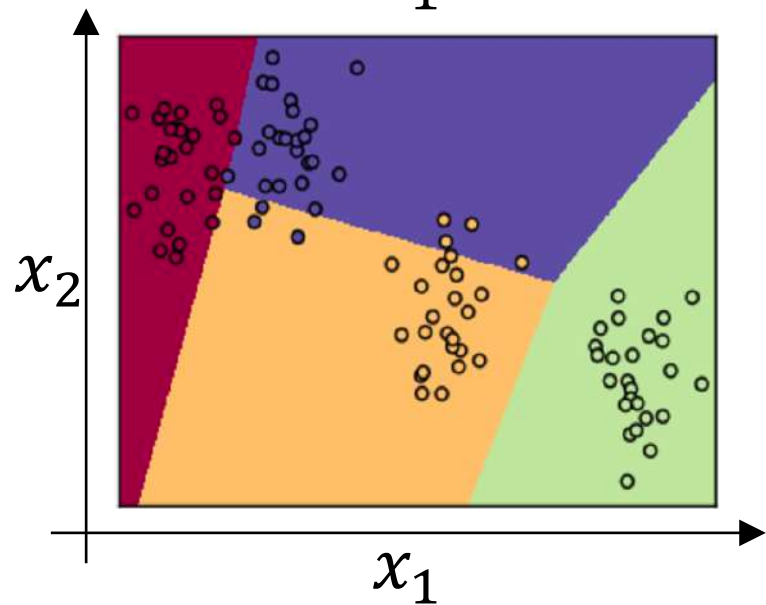
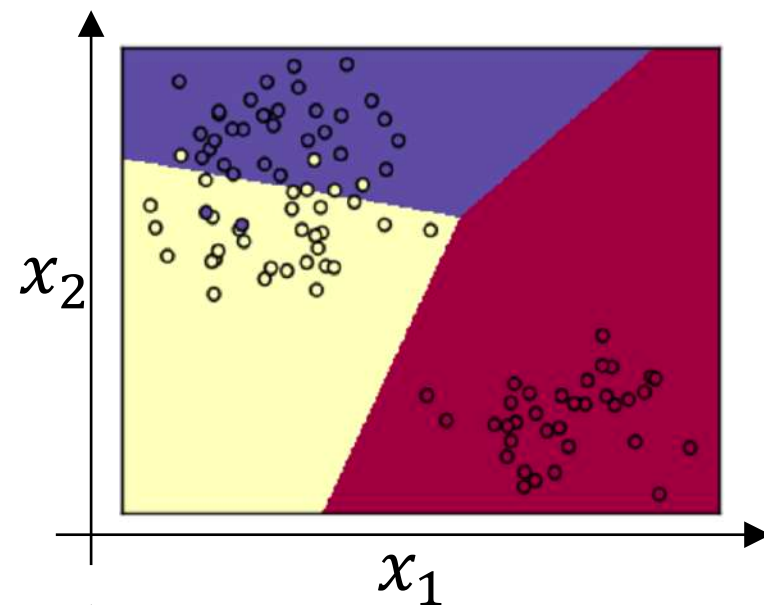
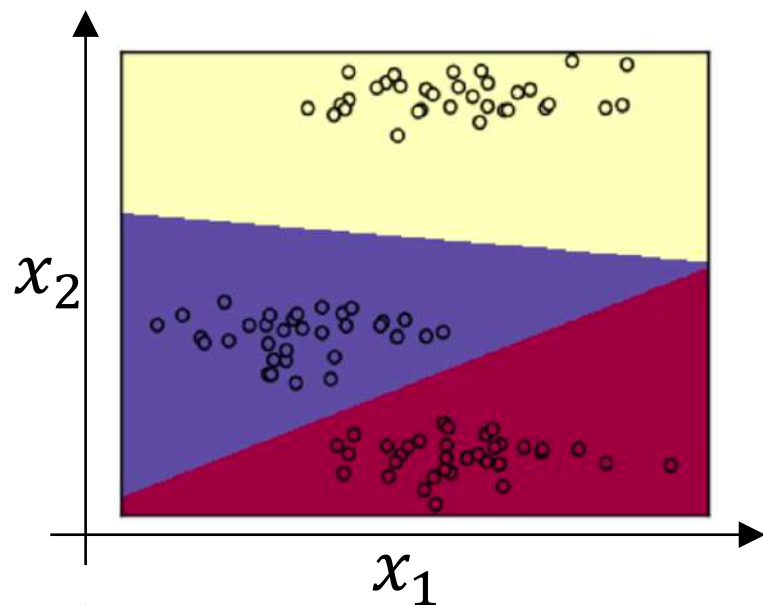
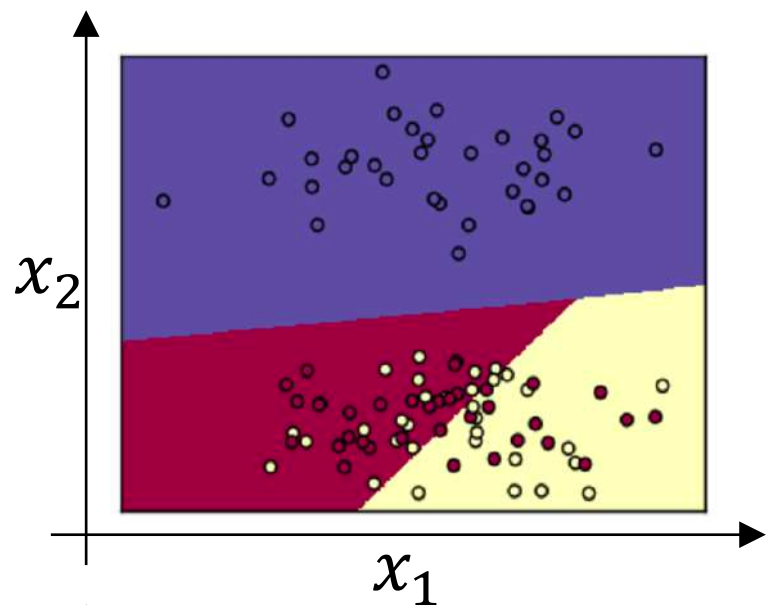
1



# Softmax layer



# Softmax examples





deeplearning.ai

# Multi-class classification

---

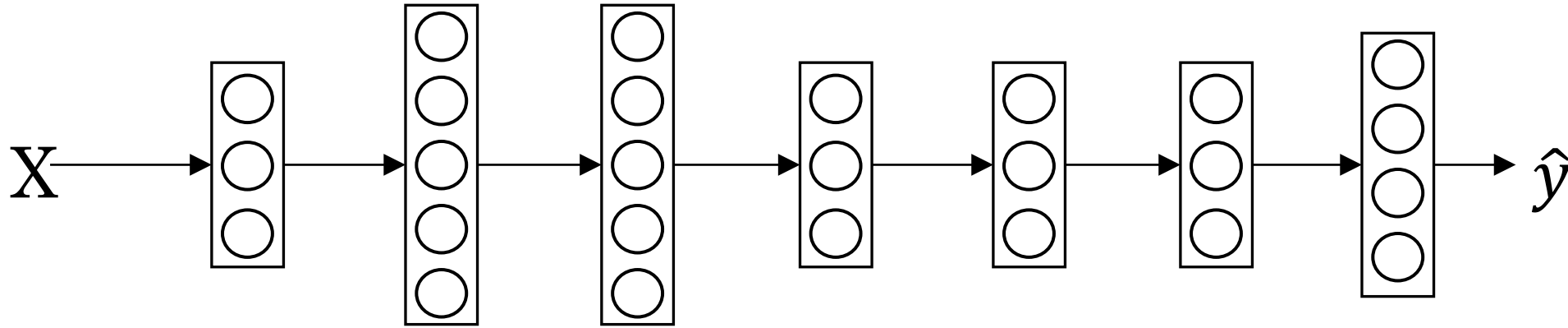
## Trying a softmax classifier

# Understanding softmax



# Loss function

# Summary of softmax classifier





deeplearning.ai

# Programming Frameworks

---

# Deep Learning frameworks

# Deep learning frameworks

- Caffe/Caffe2
- CNTK
- DL4J
- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

## Choosing deep learning frameworks

- Ease of programming (development and deployment)
- Running speed
- - Truly open (open source with good governance)



deeplearning.ai

# Programming Frameworks

---

## TensorFlow

# Motivating problem

$$\begin{aligned} J(w) &= \boxed{w^2 - 10w + 25} \\ &\quad \swarrow \\ &\quad (w-5)^2 \\ &\quad w=5 \end{aligned}$$

$$\begin{aligned} J(w, b) \\ \uparrow \quad \uparrow \end{aligned}$$

# Code example

```
import numpy as np
import tensorflow as tf
```

```
coefficients = np.array([[1], [-20], [25]])
```

```
w = tf.Variable([0], dtype=tf.float32)
```

```
x = tf.placeholder(tf.float32, [3, 1])
```

```
cost = x[0][0]*w**2 + x[1][0]*w + x[2][0] # (w-5)**2
```

```
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
init = tf.global_variables_initializer()
```

```
session = tf.Session()
```

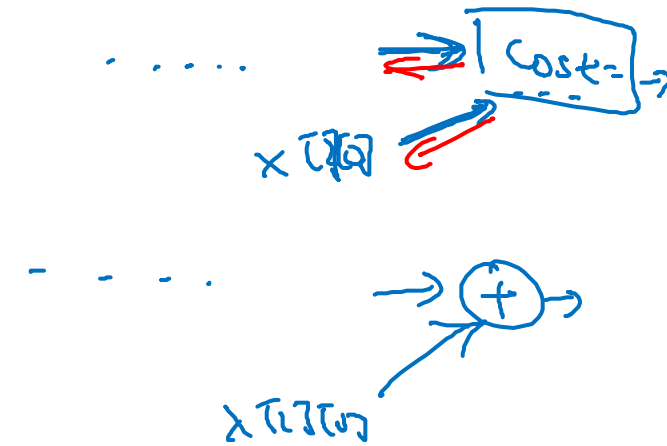
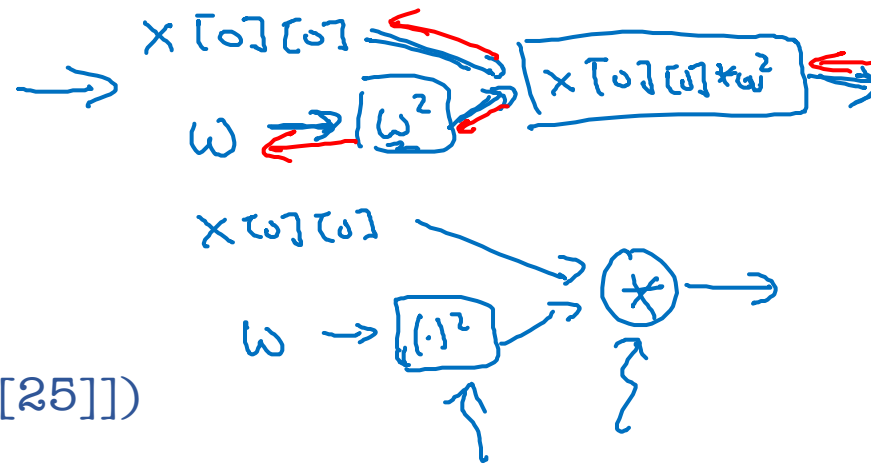
```
session.run(init)
```

```
print(session.run(w))
```

```
for i in range(1000):
```

```
    session.run(train, feed_dict={x:coefficients})
```

```
print(session.run(w))
```



```
with tf.Session() as session:
```

```
    session.run(init)
```

```
    print(session.run(w))
```