



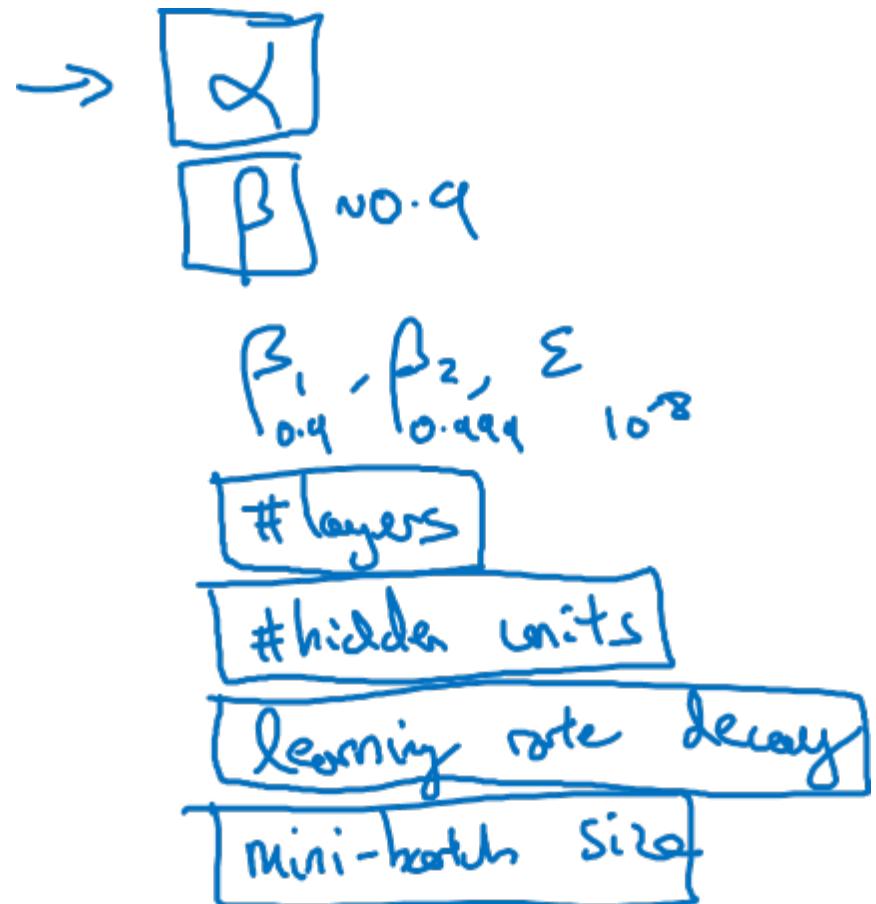
deeplearning.ai

# Hyperparameter tuning

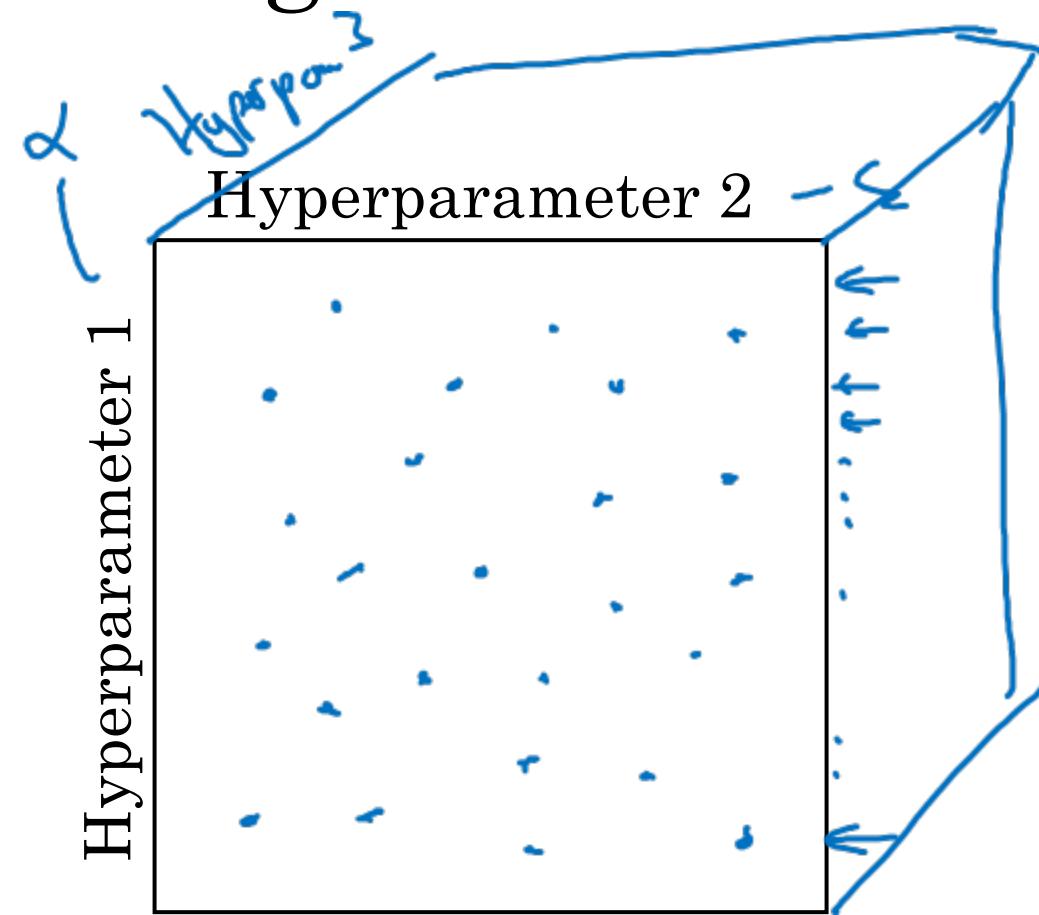
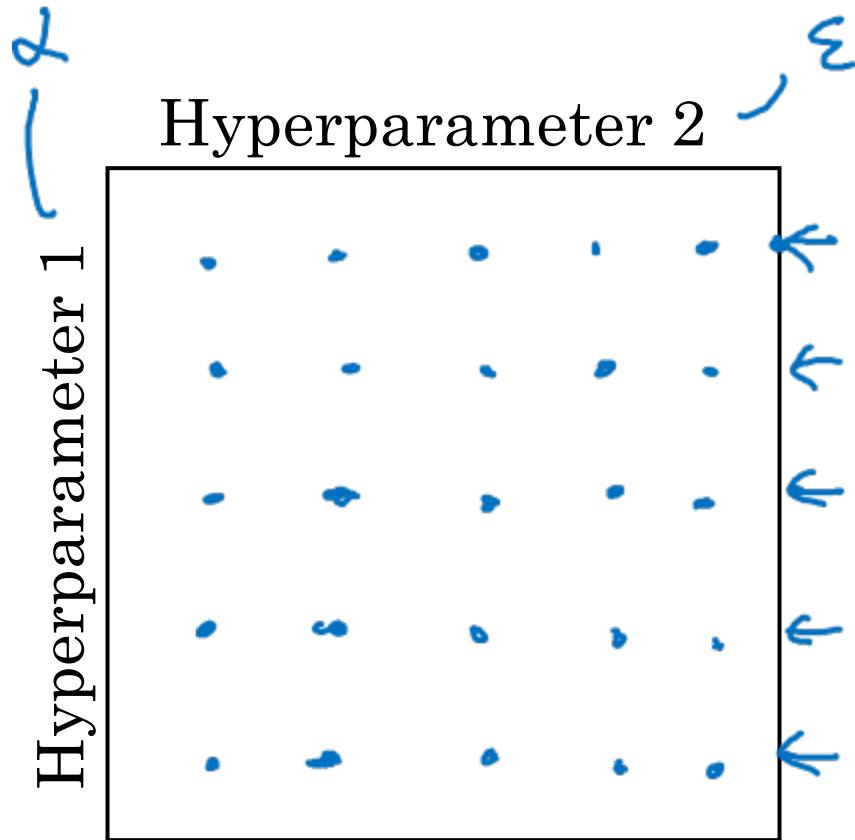
---

## Tuning process

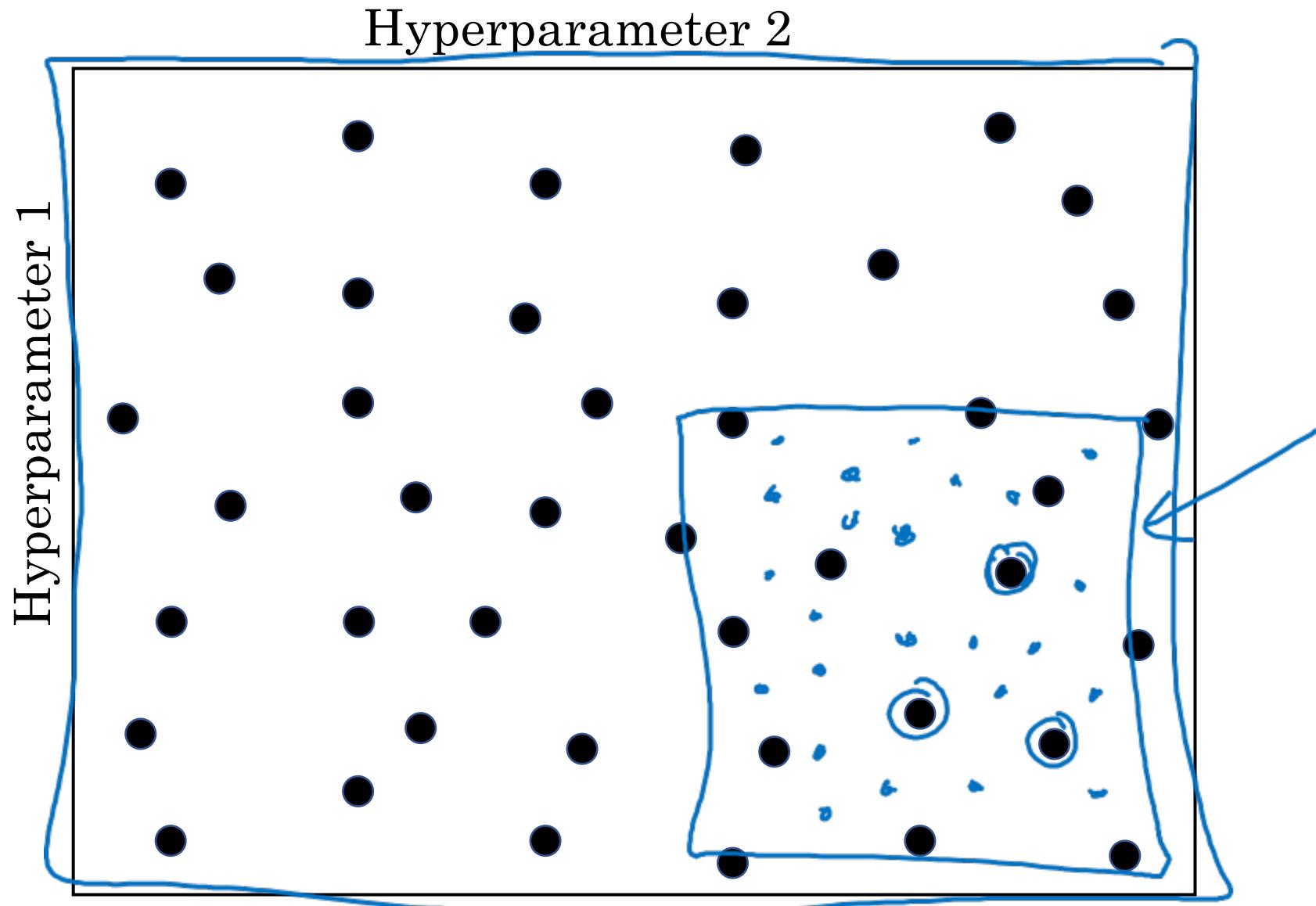
# Hyperparameters



# Try random values: Don't use a grid



# Coarse to fine





deeplearning.ai

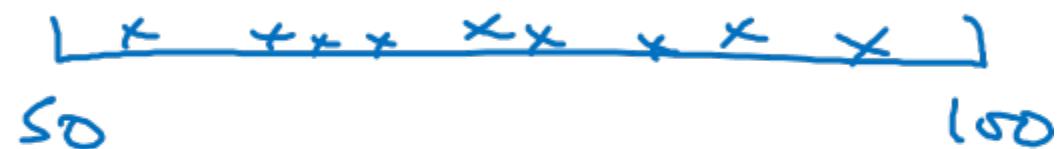
# Hyperparameter tuning

---

Using an appropriate  
scale to pick  
hyperparameters

# Picking hyperparameters at random

→  $\eta^{[l]} = 50, \dots, 100$

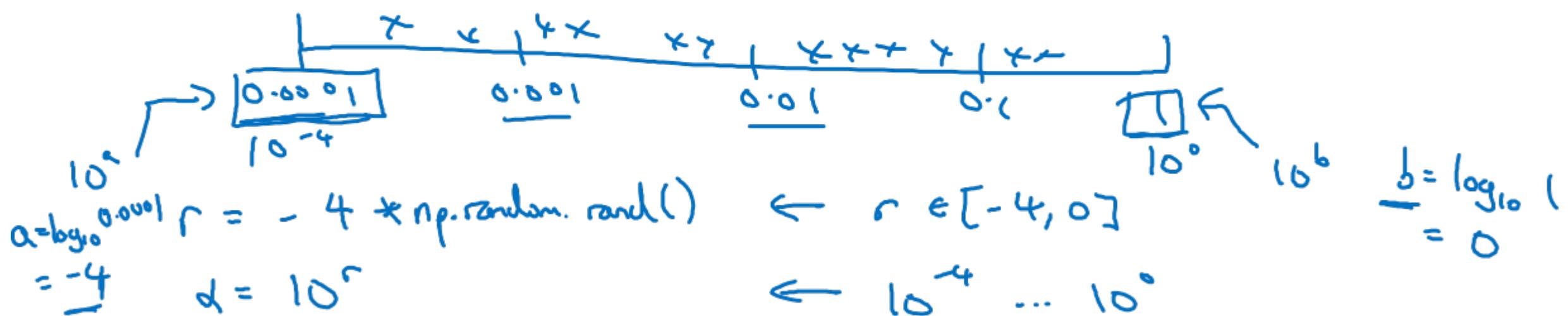
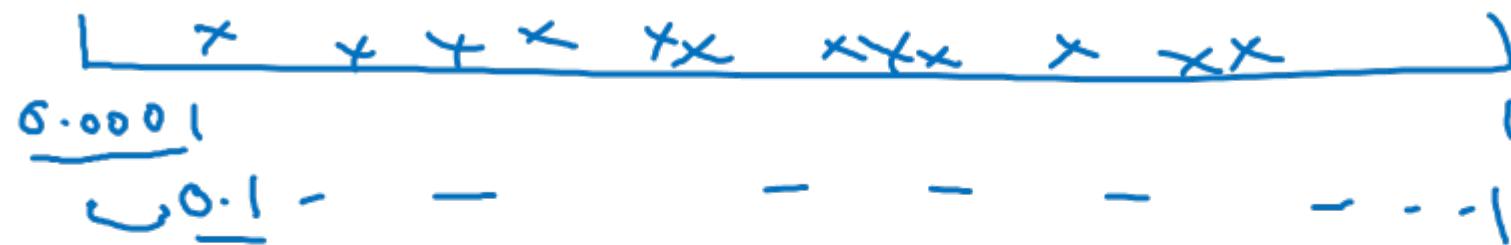


→ #layers    L : 2 - 4

2, 3, 4

# Appropriate scale for hyperparameters

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



$$10^a \dots 10^b$$

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

$$\alpha = 10^r$$

# Hyperparameters for exponentially weighted averages

$$\beta = 0.9 \dots 0.999$$

$\downarrow$                      $\downarrow$   
 $10$                      $1000$

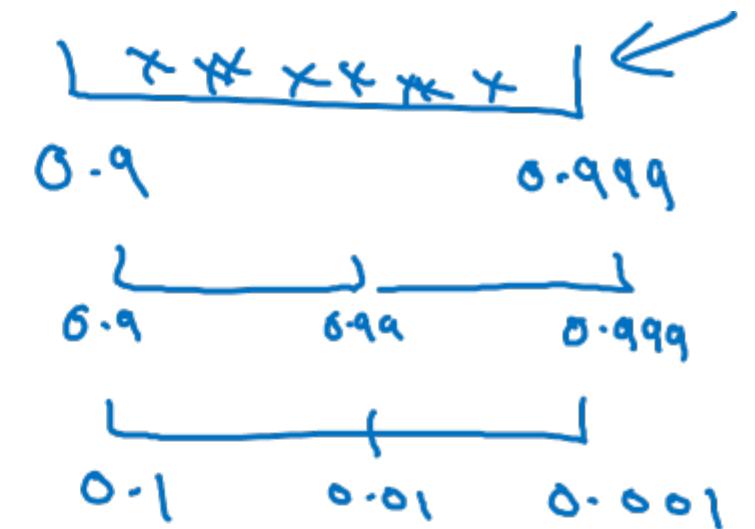
$$1-\beta = 0.1 \dots 0.001$$

$$\beta: 0.900 \rightarrow 0.9005 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

$\sim 1000$                      $\sim 2000$

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



$$\frac{10^{-1}}{10^{-3}}$$

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

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



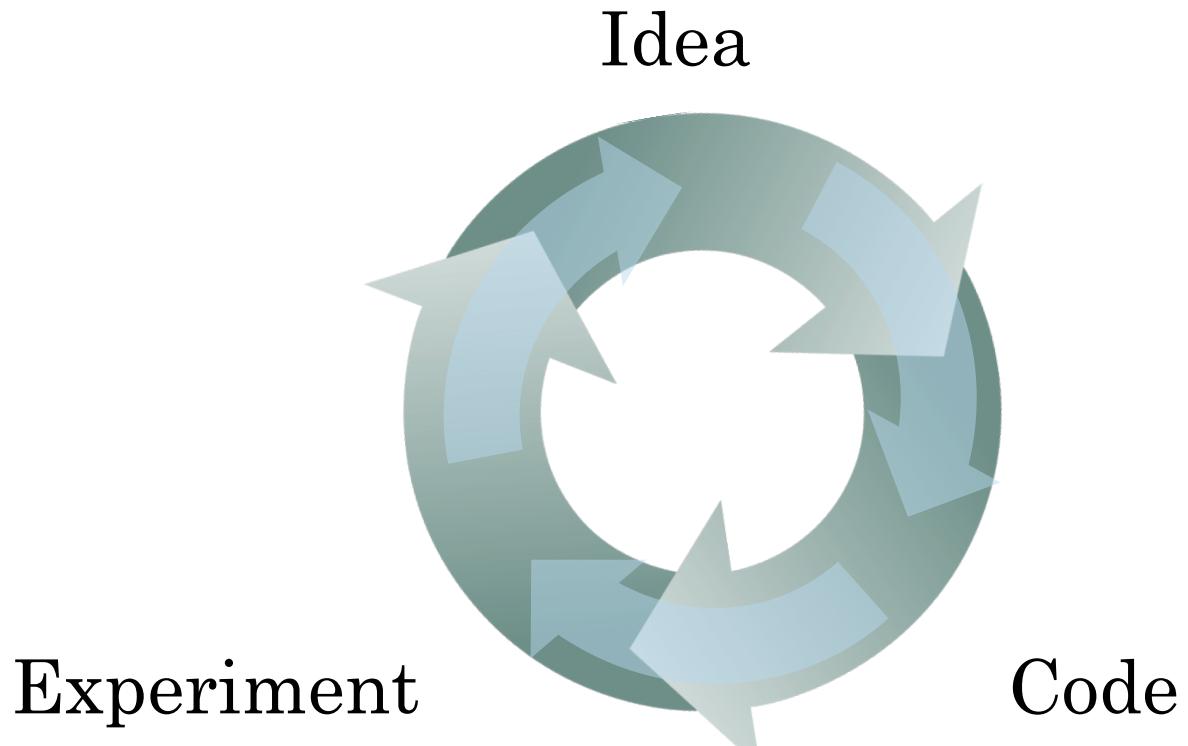
deeplearning.ai

# Hyperparameters tuning

---

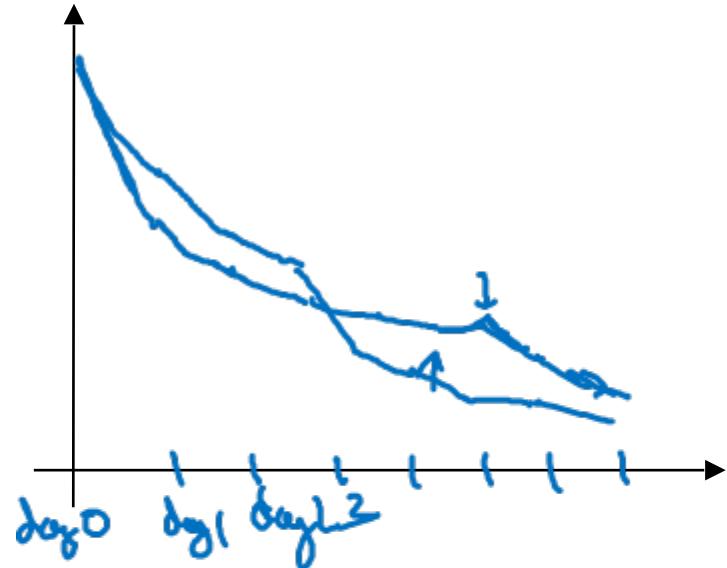
## Hyperparameters tuning in practice: Pandas vs. Caviar

# Re-test hyperparameters occasionally



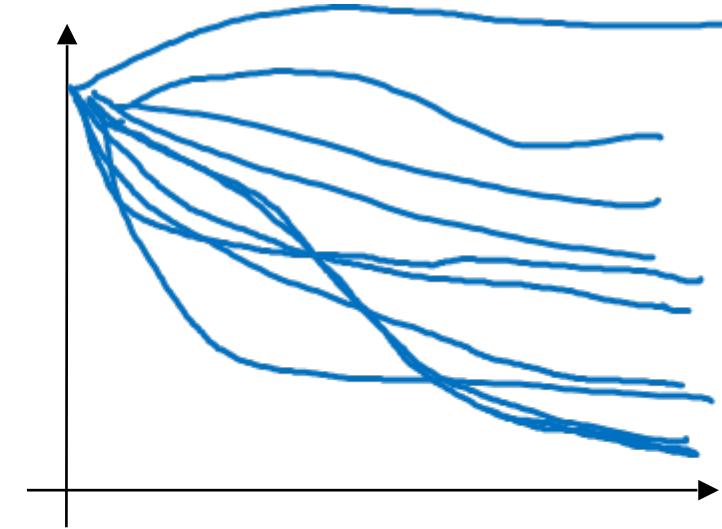
- NLP, Vision, Speech,  
Ads, logistics, ....
- Intuitions do get stale.  
Re-evaluate occasionally.

# Babysitting one model



Panda ↪

# Training many models in parallel



Caviar ↪

Andrew Ng



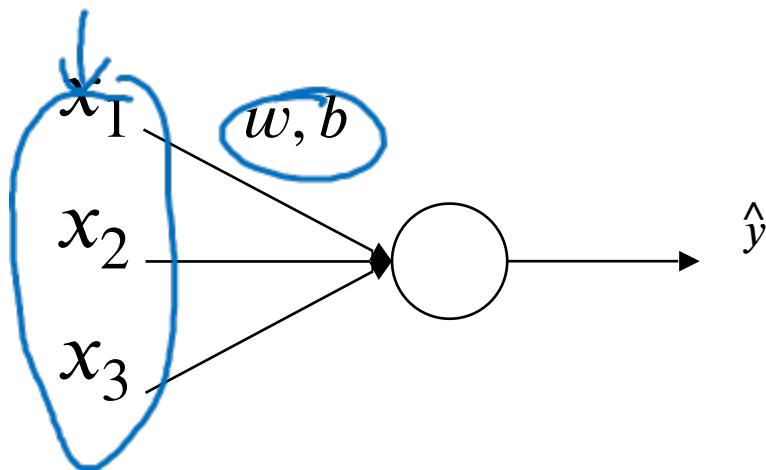
deeplearning.ai

# Batch Normalization

---

Normalizing  
activations  
in a network

# Normalizing inputs to speed up learning

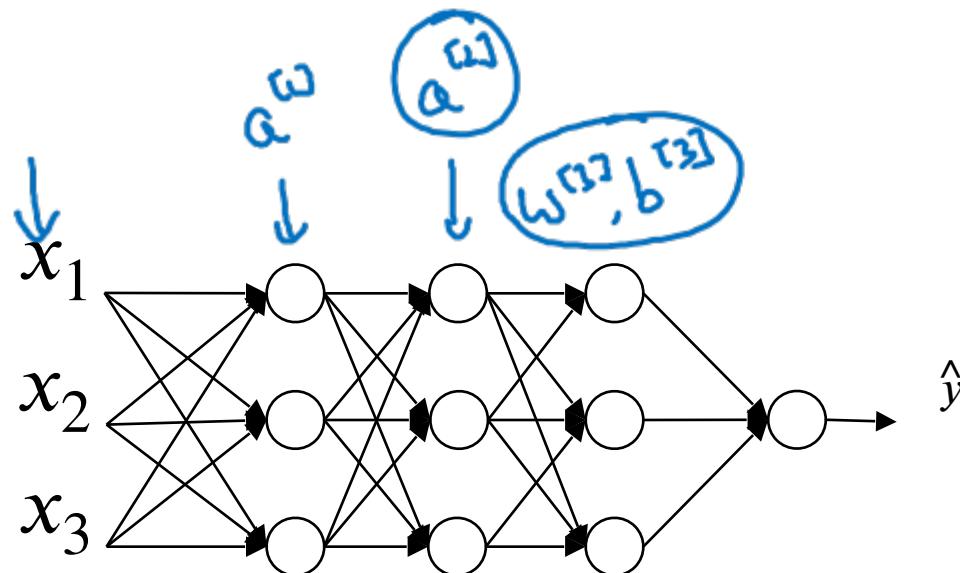
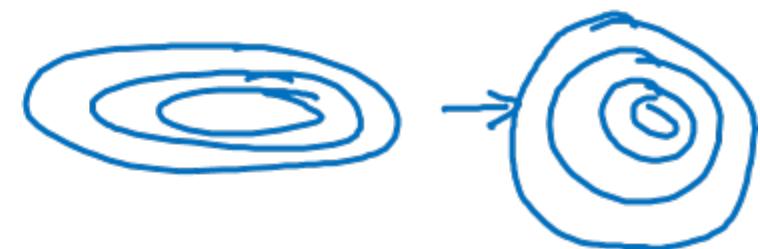


$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$X = X - \mu \quad \leftarrow \text{element-wise}$$

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

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



Can we normalize  $\frac{a^{(2)}}{w^{(2)}, b^{(2)}}$  so  
as to tan  $w^{(2)}, b^{(2)}$  fast?

Normalize  $\frac{z^{(2)}}{\uparrow}$

# Implementing Batch Norm

Given some intermediate values in NN

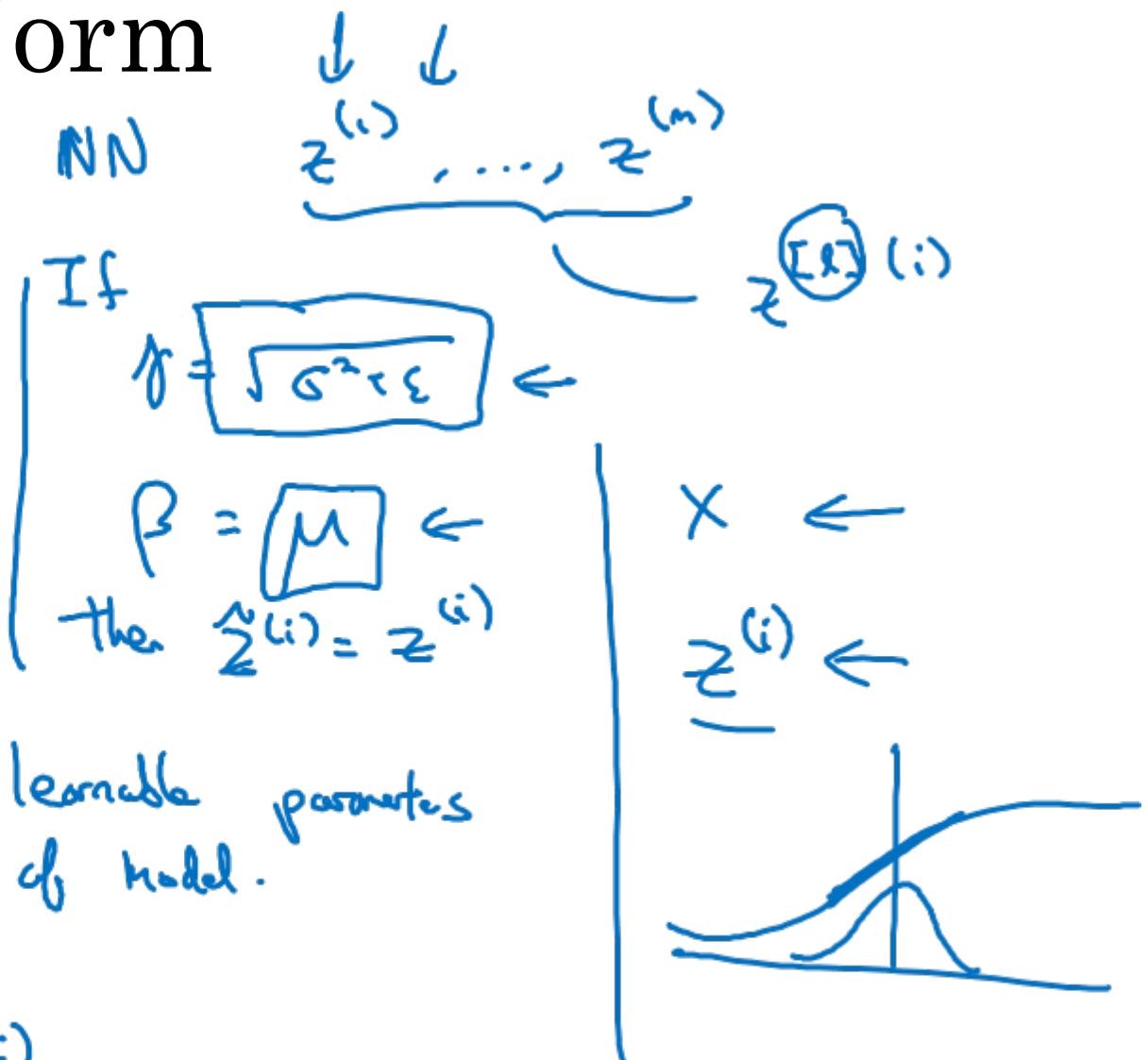
$$\mu = \frac{1}{m} \sum z^{(i)}$$

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

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

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

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





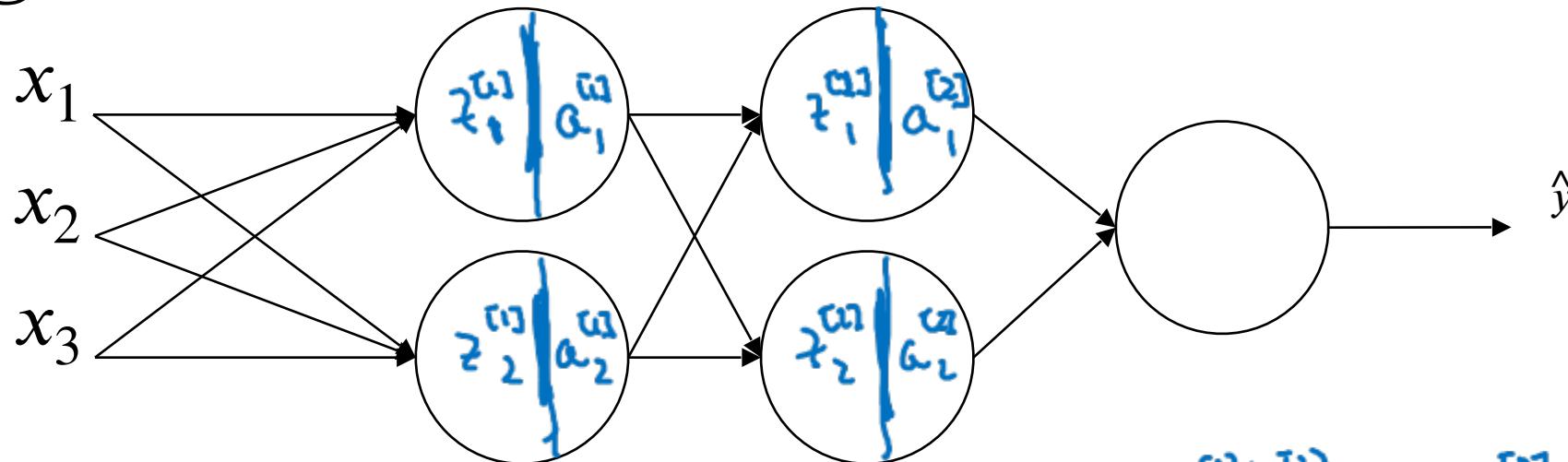
deeplearning.ai

# Batch Normalization

---

## Fitting Batch Norm into a neural network

# Adding Batch Norm to a network



$$x \xrightarrow{\substack{w^{(l)}, b^{(l)} \\ z^{(l)}}} \xrightarrow[\text{BatchNorm (BN)}]{\substack{\beta^{(l)}, \gamma^{(l)} \\ \bar{z}}} \xrightarrow{\substack{N^{(l)} \\ z^{(l)}}} a^{(l)} = g(\bar{z}^{(l)}) \xrightarrow{\substack{w^{(l)}, b^{(l)} \\ z^{(l)}}} \xrightarrow[\text{BN}]{\substack{\beta^{(l)}, \gamma^{(l)} \\ \bar{z}}} \xrightarrow{\substack{N^{(l)} \\ z^{(l)}}} a^{(l)} \rightarrow \dots$$

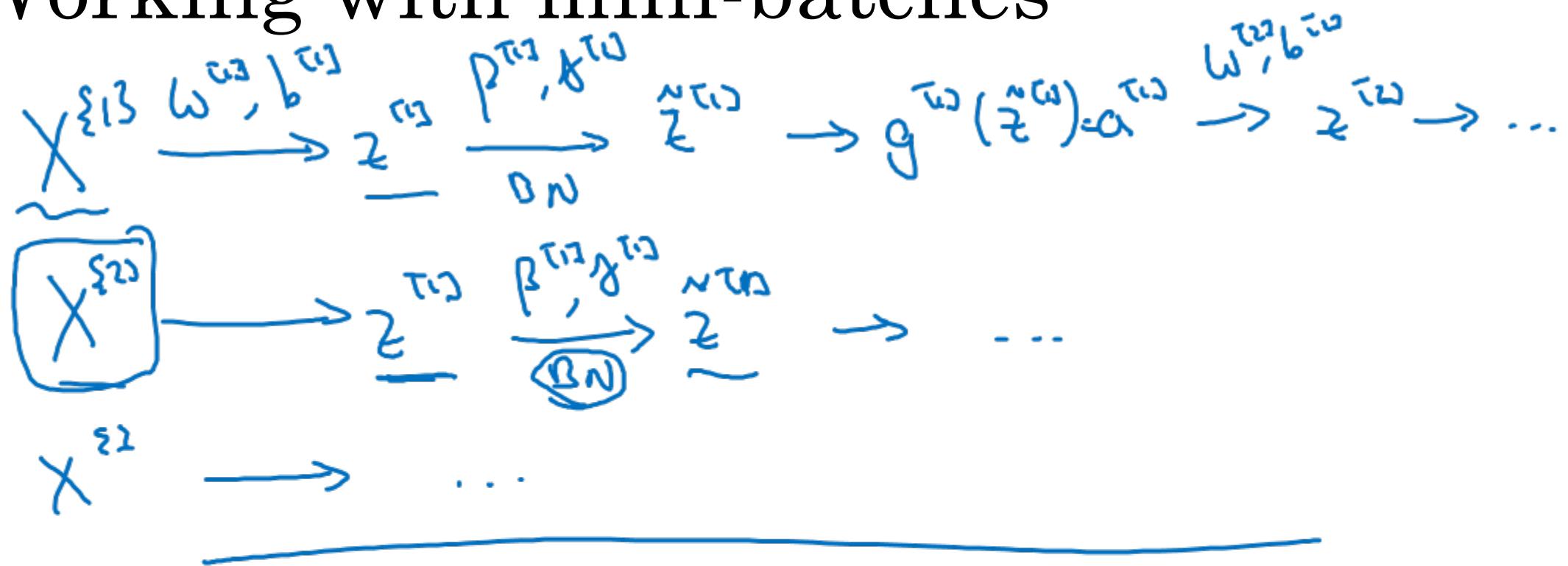
Parameters:

$$\left. \begin{array}{c} w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}, \\ \rightarrow \beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)} \end{array} \right\} \quad \rightarrow \beta$$

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

`tf.nn.batch_normalization` ←

# Working with mini-batches



Parameters:  $\omega^{(1)}, \gamma^{(1)}, \beta^{(1)}, \gamma^{(2)}, \beta^{(2)}, \dots, \omega^{(L)}, \gamma^{(L)}, \beta^{(L)}$ .

$$z^{(1)} \\ (n^{(1)}, 1)$$

$$\begin{matrix} \omega^{(1)} \\ (n^{(1)}, n^{(2)}) \end{matrix}, b^{(1)} \\ \begin{matrix} \beta^{(1)} \\ (n^{(2)}, 1) \end{matrix}, \gamma^{(1)} \\ \begin{matrix} \omega^{(2)} \\ (n^{(2)}, n^{(3)}) \end{matrix}, b^{(2)} \\ \vdots \\ \begin{matrix} \omega^{(L)} \\ (n^{(L)}, n^{(L+1)}) \end{matrix}, b^{(L)}$$

$$\rightarrow z^{(1)} = \omega^{(1)} a^{(1)} + \beta^{(1)}$$
$$z^{(2)} = \omega^{(2)} a^{(2)}$$
$$z_{\text{norm}}^{(2)} = \frac{z^{(2)}}{\|\omega^{(2)}\|}$$
$$\rightarrow z^{(L)} = \omega^{(L)} z_{\text{norm}}^{(L)} + \beta^{(L)}$$

# Implementing gradient descent

for  $t = 1 \dots \text{numMiniBatches}$   
Compute forward prop on  $X^{[t]}$ .

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

Use backprop to compute  $\underline{dw}^{[t]}$ ,  ~~$\underline{dx}^{[t]}$~~ ,  $\underline{d\beta}^{[t]}$ ,  $\underline{dg}^{[t]}$

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

Works w/ momentum, RMSprop, Adam.



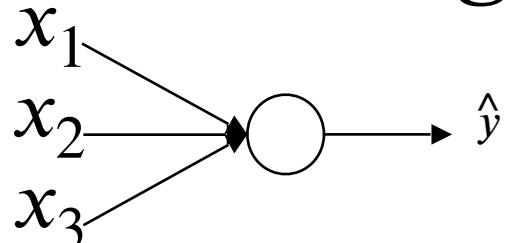
deeplearning.ai

# Batch Normalization

---

Why does  
Batch Norm work?

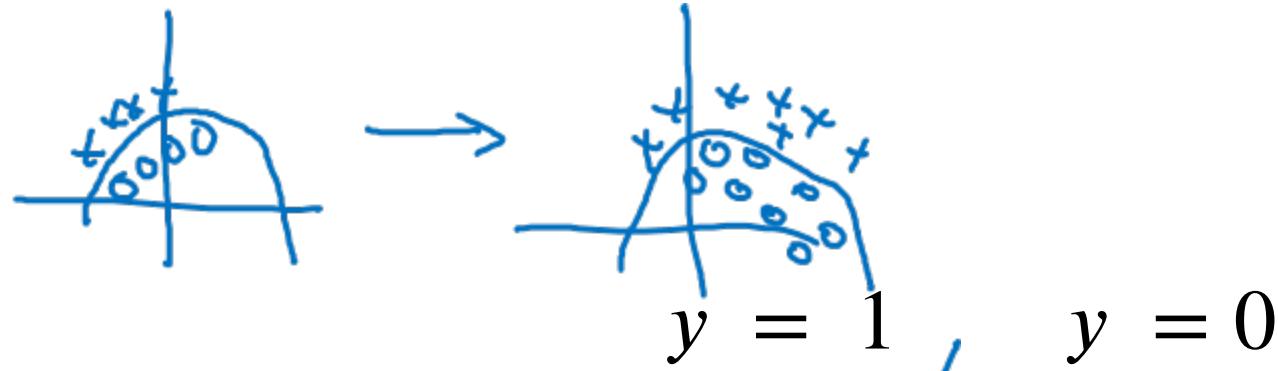
# Learning on shifting input distribution



Cat



Non-Cat

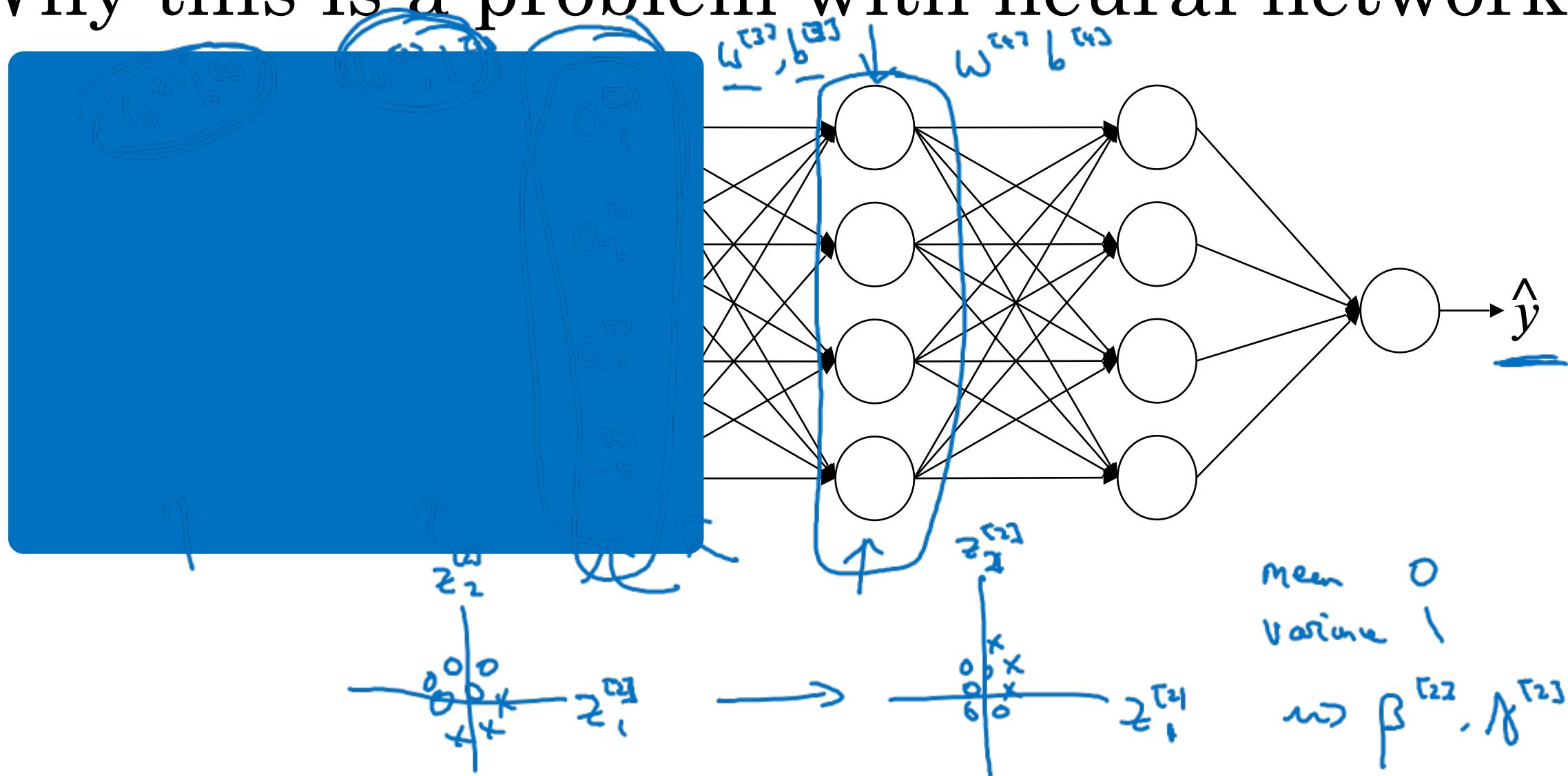


"Covariate shift"

$$\underline{x} \rightarrow y$$



# Why this is a problem with neural networks?



# Batch Norm as regularization

X

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.  
 $\hat{z}^{(i)}$      $\mu, \sigma^2$   
 $\{x_i\}$
- This adds some noise to the values within that minibatch. So similar to dropout, it adds some noise to each hidden layer's activations.  
 $\mu, \sigma^2$
- 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

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$
$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2$$

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

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

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

$$x^{e_1}, x^{e_2}, x^{e_3}, \dots$$

$$\mu^{e_1} = \frac{x^{e_1}}{N}$$

$$\mu^{e_2} = \frac{x^{e_2}}{N}$$

$$\mu^{e_3} = \frac{x^{e_3}}{N}$$

$$\underline{\mu}$$

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

$$\theta_1 = \frac{\epsilon^{e_1}}{N}, \theta_2 = \frac{\epsilon^{e_2}}{N}, \theta_3 = \frac{\epsilon^{e_3}}{N}, \dots$$

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

$$\tilde{z} = \gamma z_{\text{norm}} + \beta$$



deeplearning.ai

Multi-class  
classification

---

Softmax regression

# Recognizing cats, dogs, and baby chicks



3

1

2

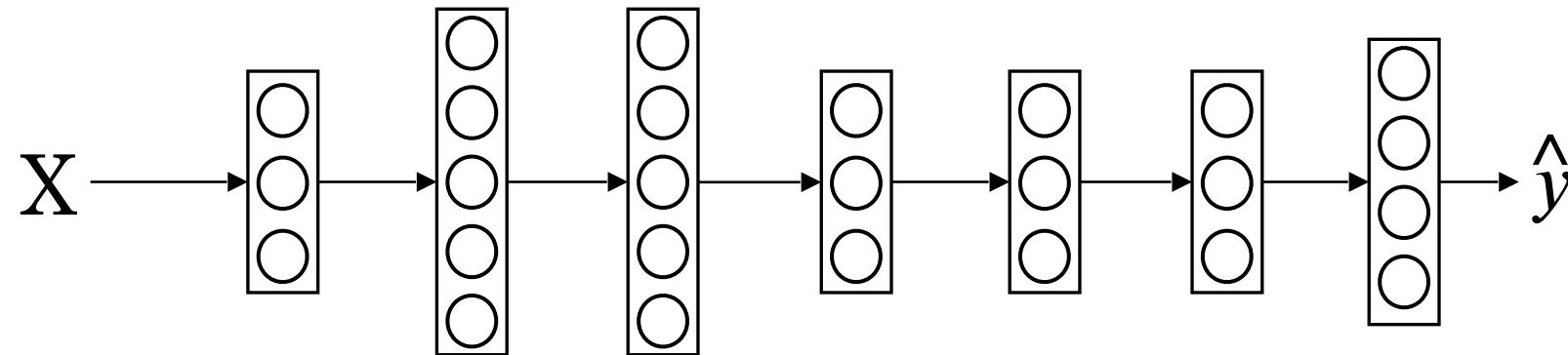
0

3

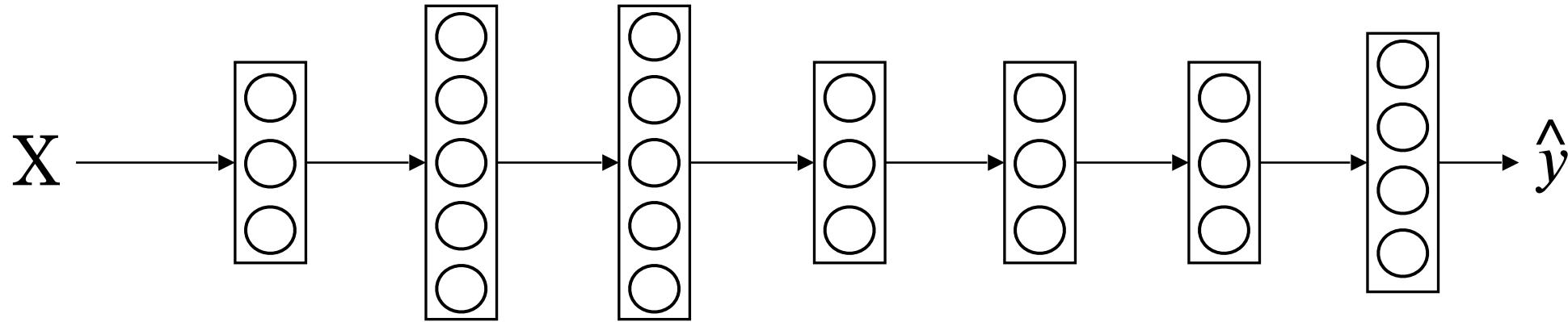
2

0

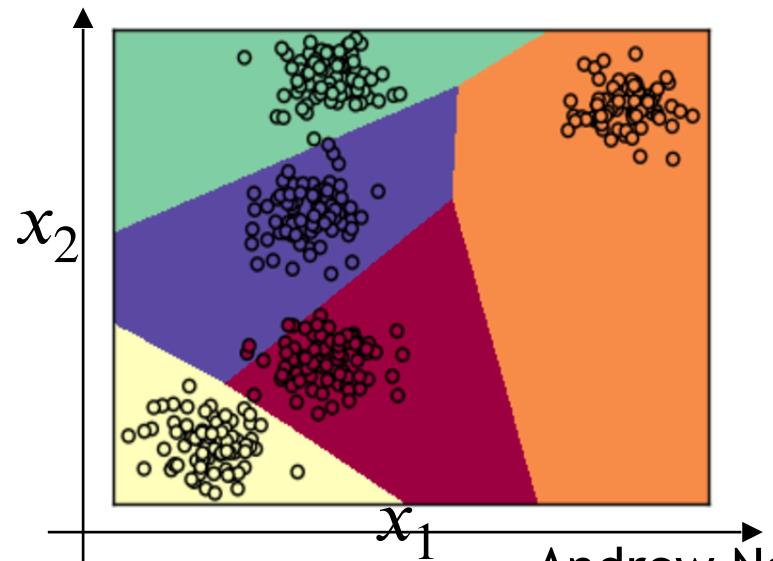
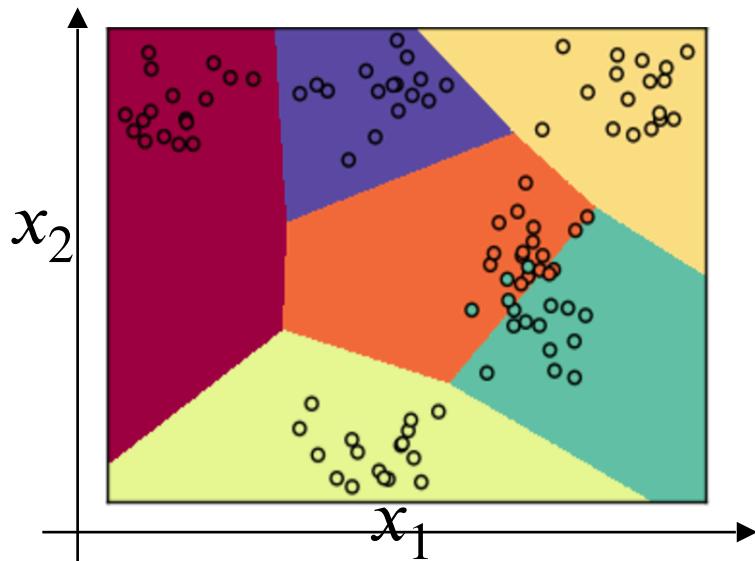
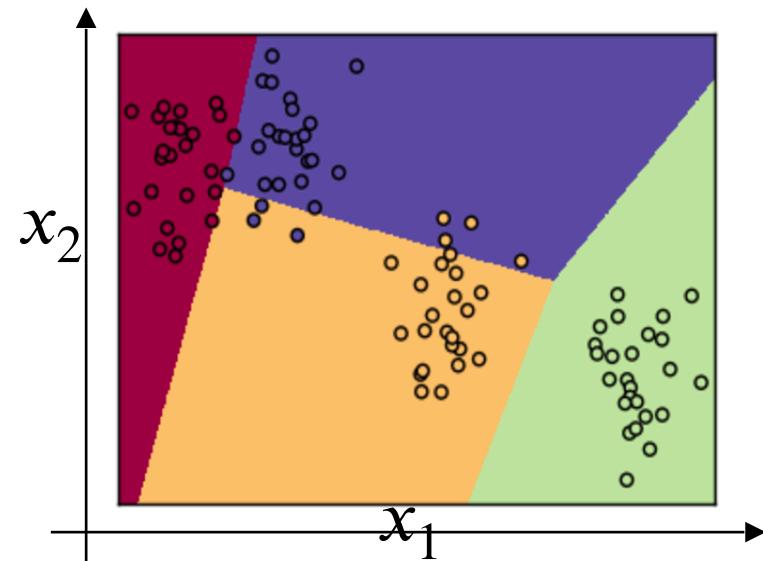
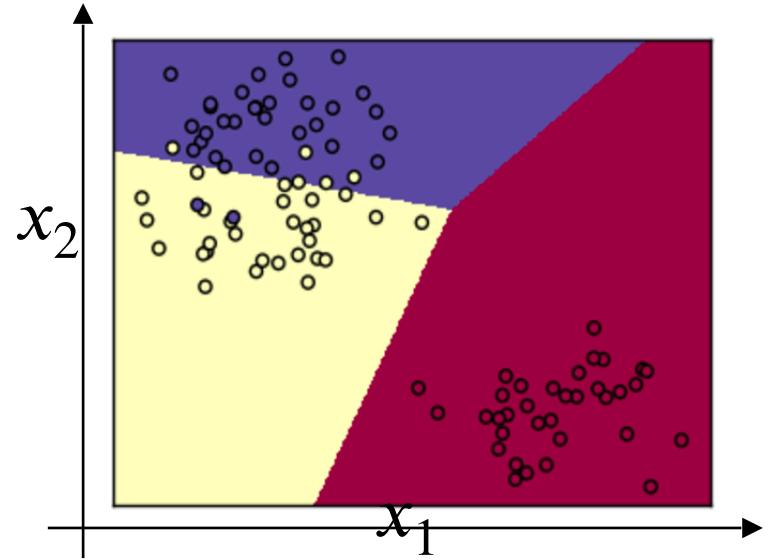
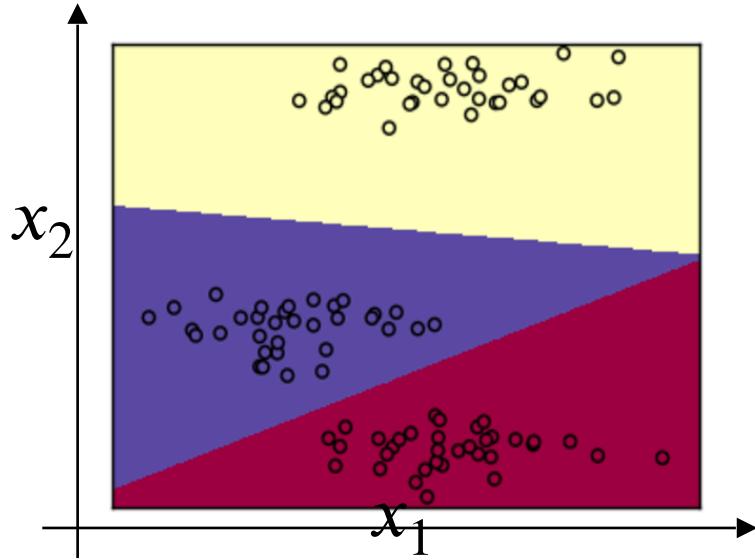
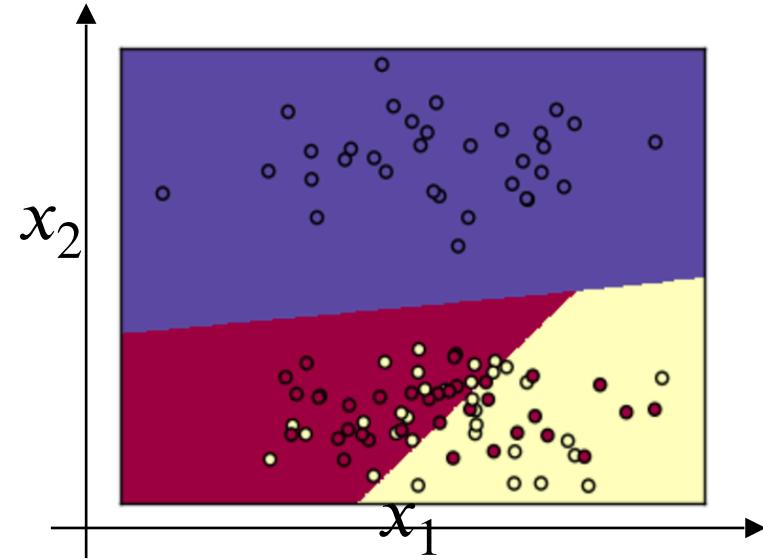
1



# Softmax layer



# Softmax examples



Andrew Ng



deeplearning.ai

# Multi-class classification

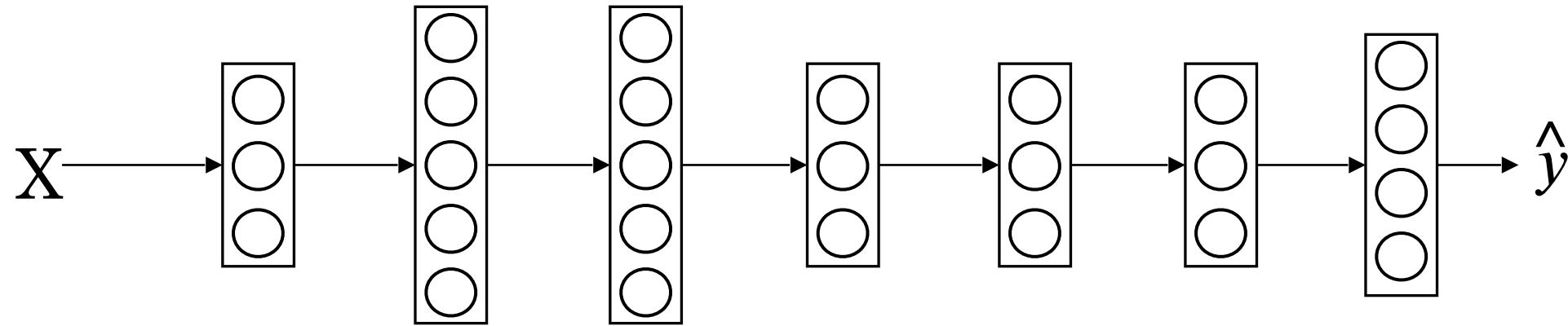
---

## 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

$$J(\omega) = \frac{\omega^2 - 10\omega + 25}{(\omega - 5)^2}$$

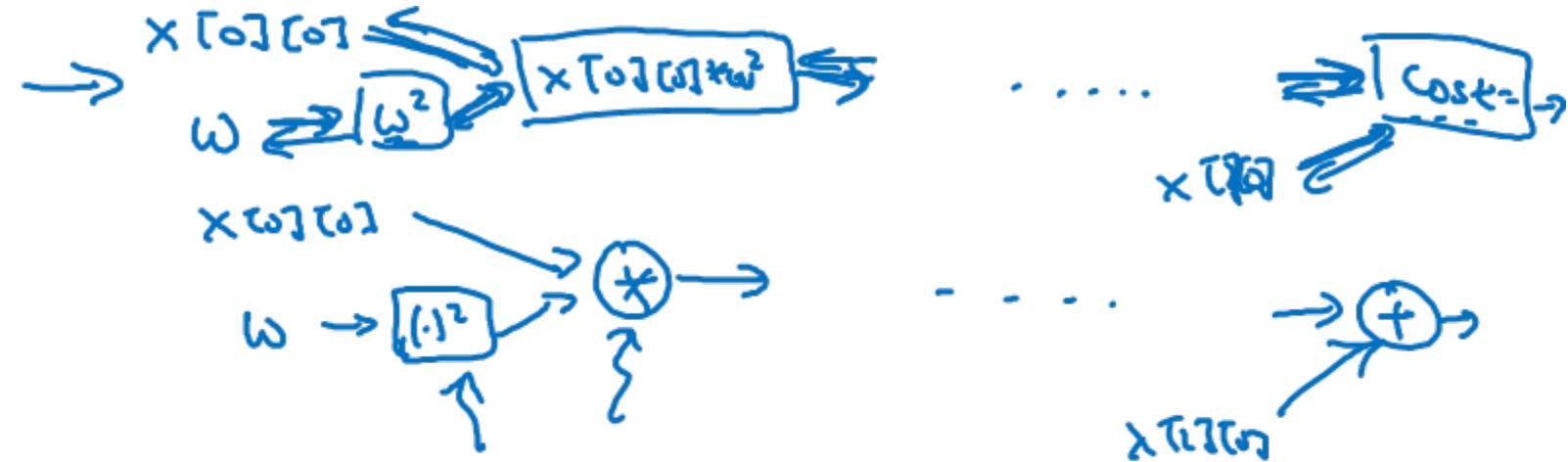
$\omega = 5$

$$J(\omega, b)$$

↑↑

# 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))
```