



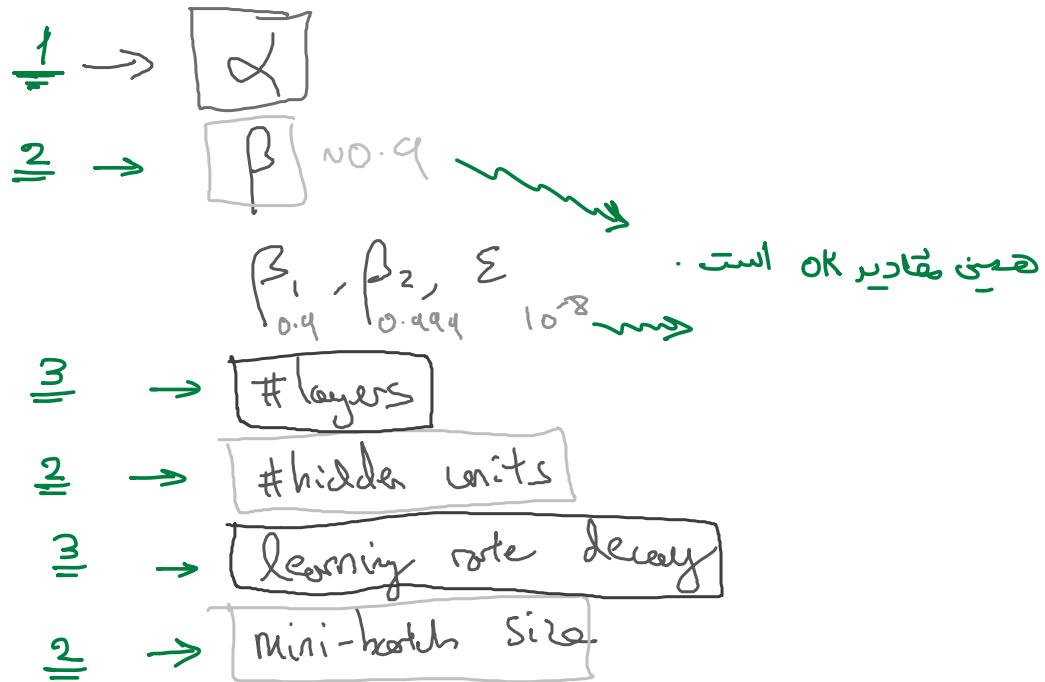
deeplearning.ai

# Hyperparameter tuning

---

## Tuning process

# Hyperparameters

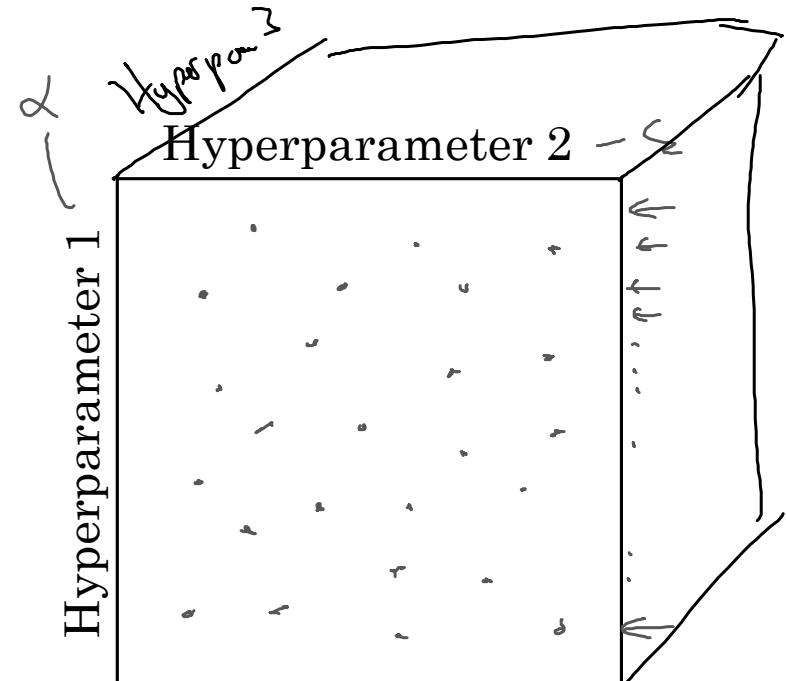
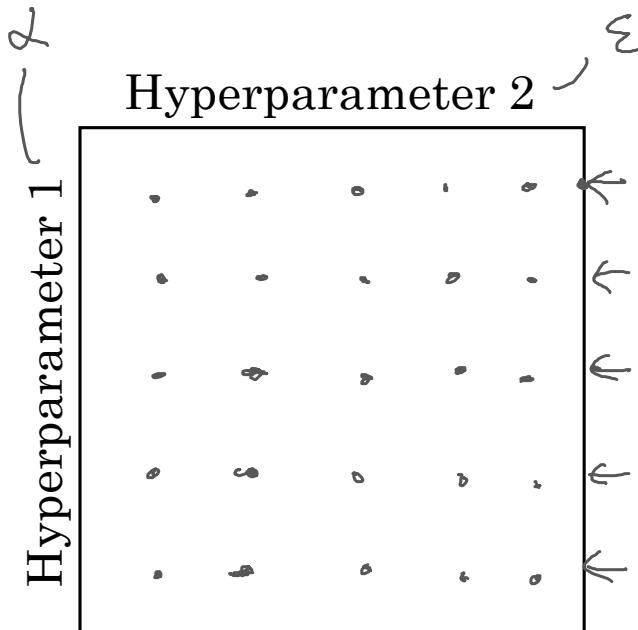


importance

When working with Big Data to tune the hyperparameters

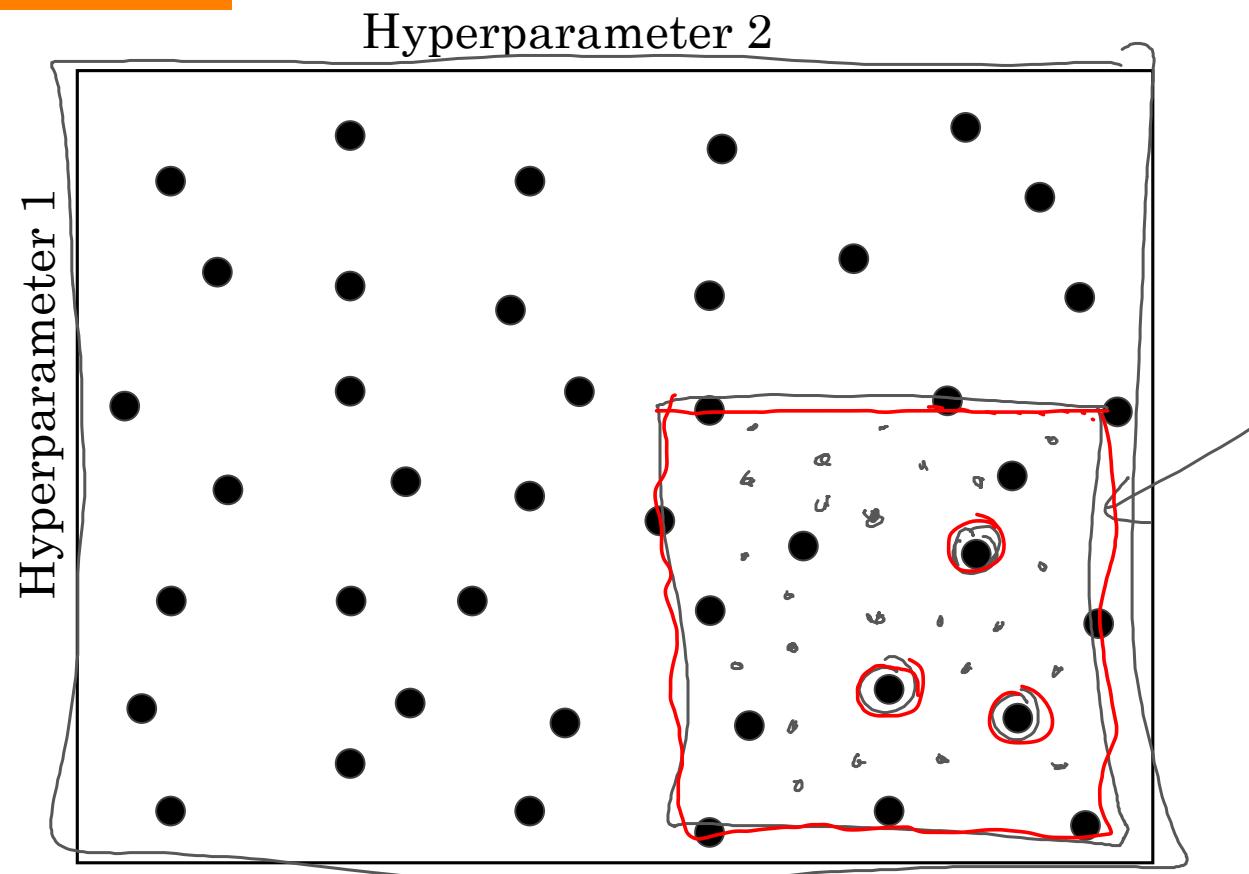
⇒ (not necessarily uniform) Random Search.

Try random values: Don't use a grid



لذه آن سه نکله خوبی پردازشی، ادماز search را در حوالی همان انجام می‌دهی!

## Coarse to fine





deeplearning.ai

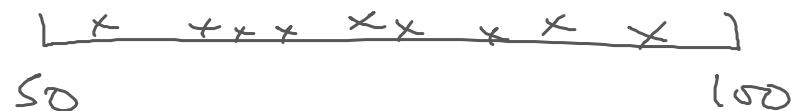
## Hyperparameter tuning

---

Using an appropriate  
scale to pick  
hyperparameters

# Picking hyperparameters at random

→  $n^{[l]} = 50, \dots, 100$



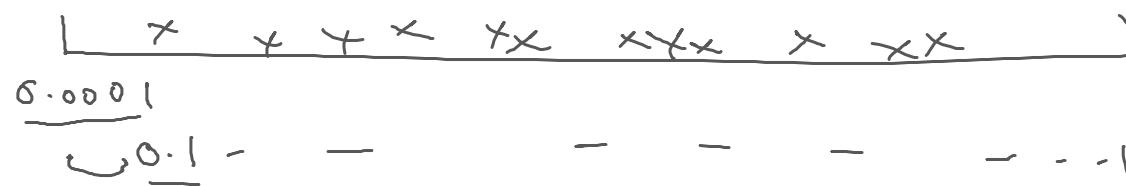
→ #layers    L : 2 - 4

2 , 3 , 4

مناسب برای random search می باشد از تبدیل آن  
که uniform search فضای جستجو را می بیند و می باشد  $\alpha \in [10^{-3}, 1]$ .  
•  $r = \log_{10} \alpha \in [-4, 0]$

## Appropriate scale for hyperparameters

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



$$\begin{aligned}
& \text{Range: } 10^{-4} \dots 10^0 \\
& \alpha = \log_{10} 0.0001 \quad r = -4 * \text{np.random.rand()} \quad \leftarrow r \in [-4, 0] \quad b = \log_{10} 1 \\
& = -4 \quad \alpha = 10^r \quad \leftarrow 10^{-4} \dots 10^0 \quad = 0
\end{aligned}$$

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

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

$$\alpha = 10^r$$

تبدیل مناسب برای چهارمین  $\beta$  دوست  $\log(1 - \beta)$  uniform search است.

## Hyperparameters for exponentially weighted averages

$$\beta = 0.9, \dots, 0.999$$

↓                            ↓

1.0                        1.000

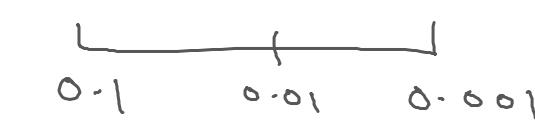
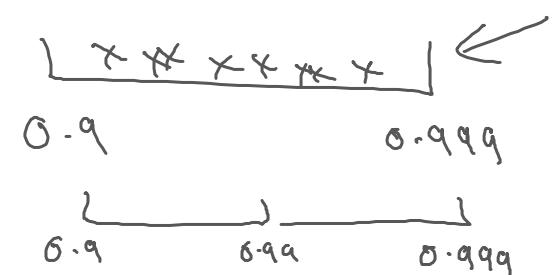
$$1 - \beta = 0.1, \dots, 0.001$$

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

$$\beta: 0.999 \rightarrow 0.9995$$

$\sim 1000$                        $\sim 2000$

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



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

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

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

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



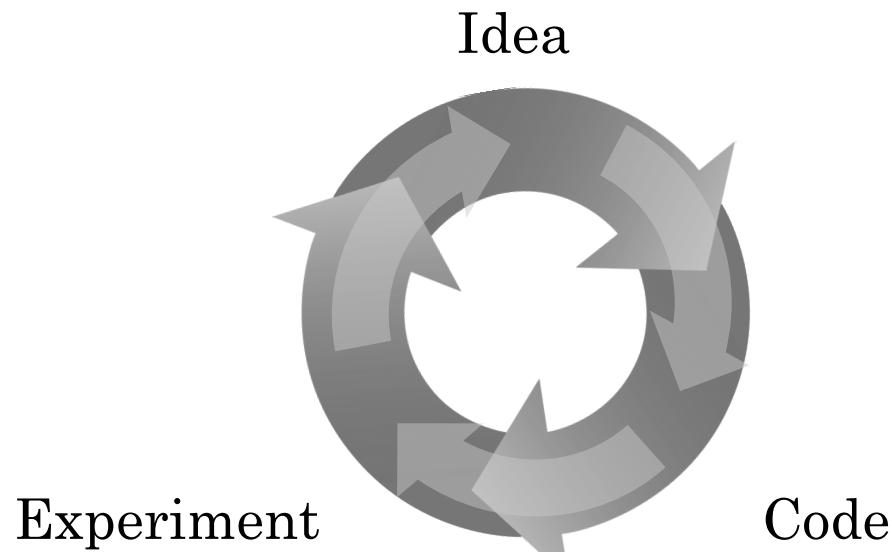
deeplearning.ai

# Hyperparameters tuning

---

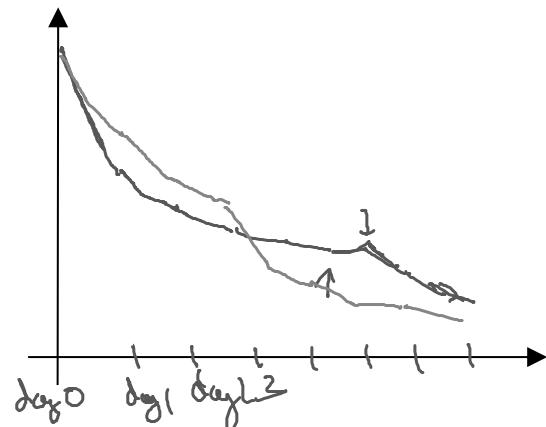
## Hyperparameters tuning in practice: Pandas vs. Caviar

# Re-test hyperparameters occasionally



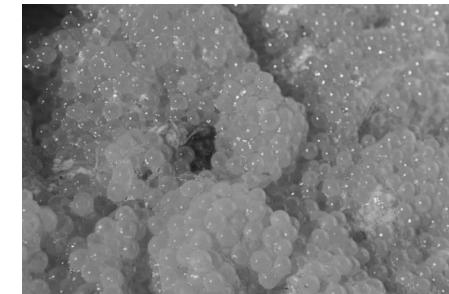
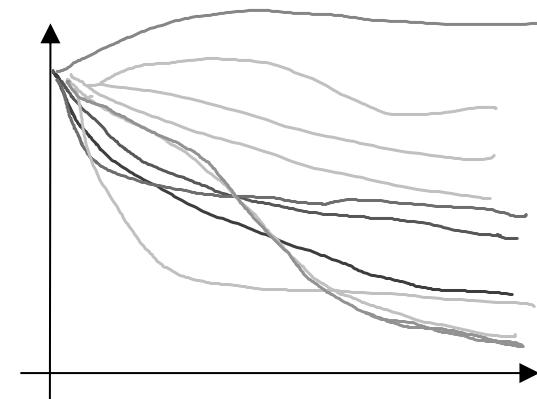
- NLP, Vision, Speech,  
Ads, logistics, ....
- Intuitions do get stale.  
Re-evaluate occasionally.

Babysitting one  
model



Panda ↵

Training many  
models in parallel



Caviar ↵



deeplearning.ai

# Batch Normalization

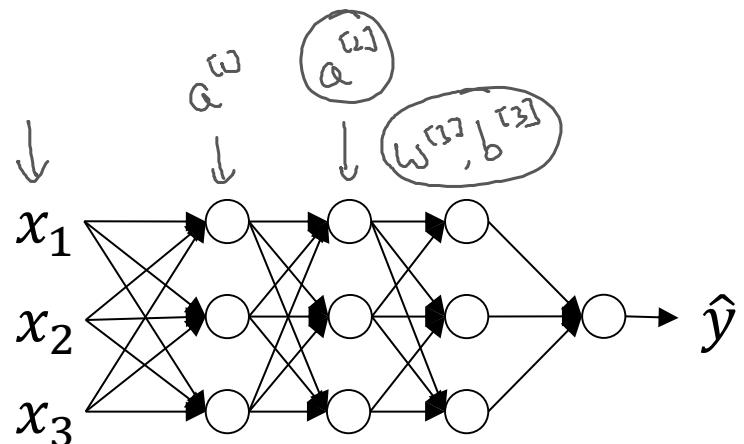
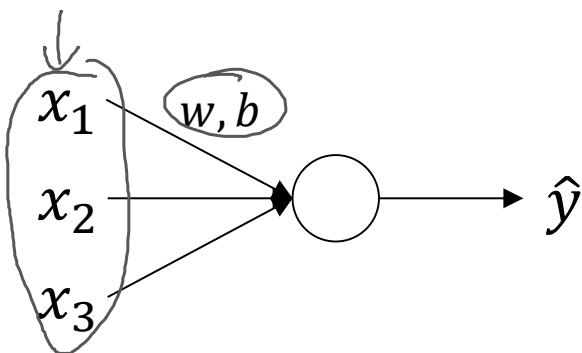
---

## Normalizing activations in a network

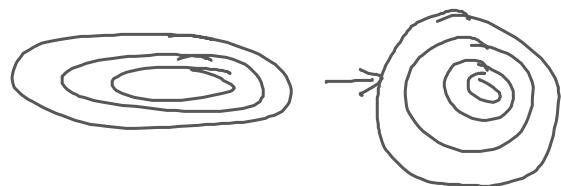
Batch normalization makes your hyperparameter search easier and makes your NN more robust. The choice of hyperparameters is a bigger range of hyperparameters that work well, and enable you to easily train deep networks.

Hidden layers  $\rightarrow (\alpha^{[l]} \leftarrow \gamma^{[l]} z^{[l]})$  خود ری باعث سرعت تمرین نرم افزار می شوند همان طوره normalize می شوند.

## Normalizing inputs to speed up learning



$$\begin{aligned}\mu &= \frac{1}{m} \sum_i x^{(i)} \\ X &= X - \mu \quad \text{element-wise} \\ \sigma^2 &= \frac{1}{m} \sum_i (x^{(i)} - \mu)^2 \\ X &= X / \sigma^2\end{aligned}$$



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

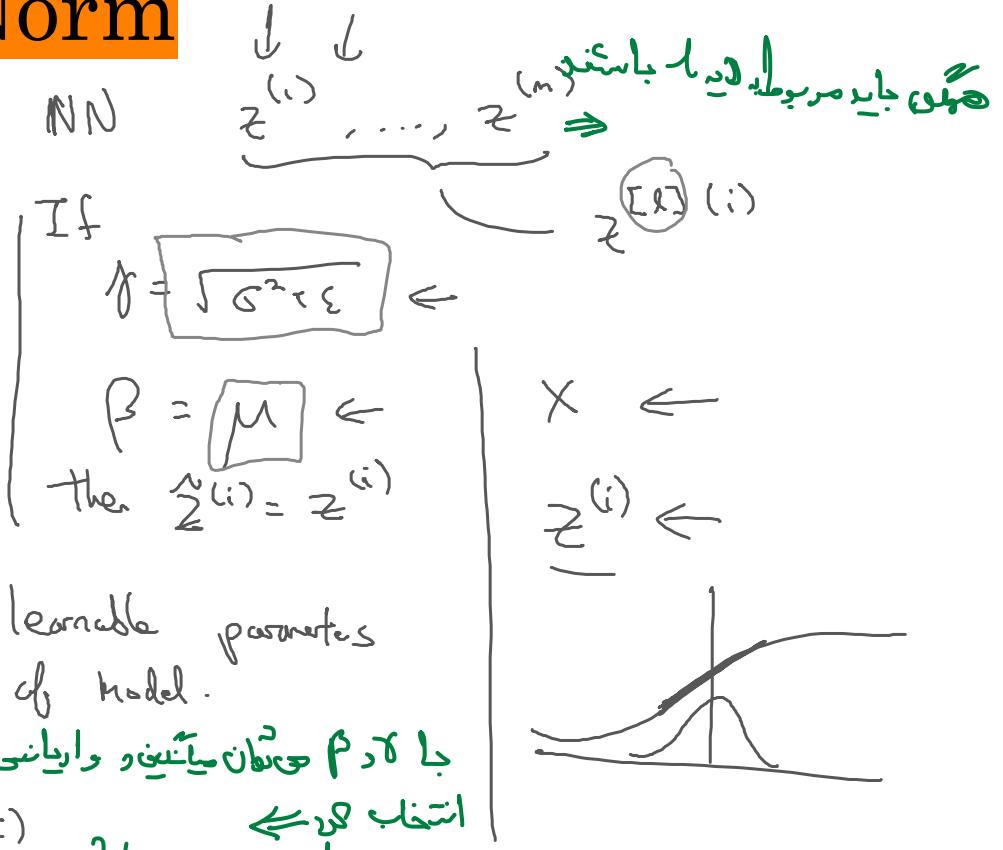
$$\text{Normalize } \frac{z^{[2]}}{\uparrow}$$



## Implementing Batch Norm

Given some intermediate values in NN

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



Use  $\hat{z}^{(i)}$  insted of  $z^{(i)}$ . استخاب کری Hyperparameter



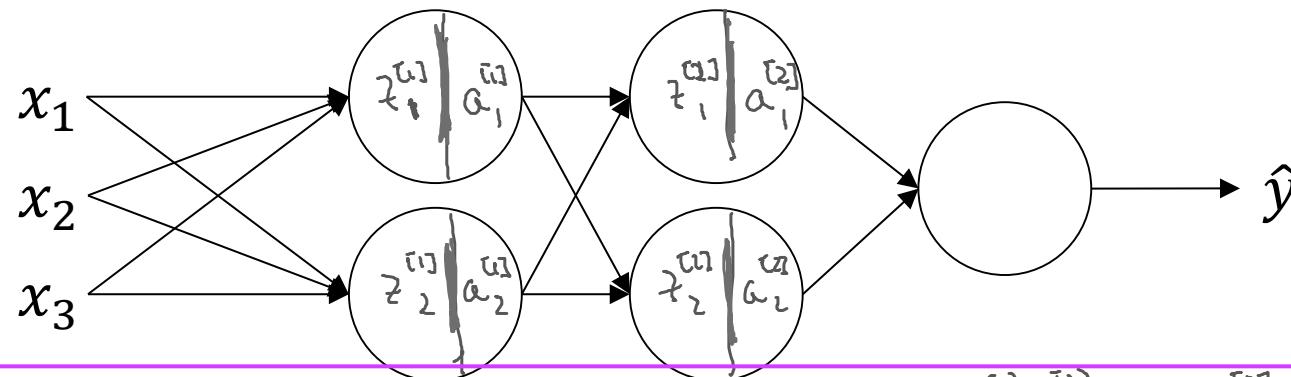
deeplearning.ai

# Batch Normalization

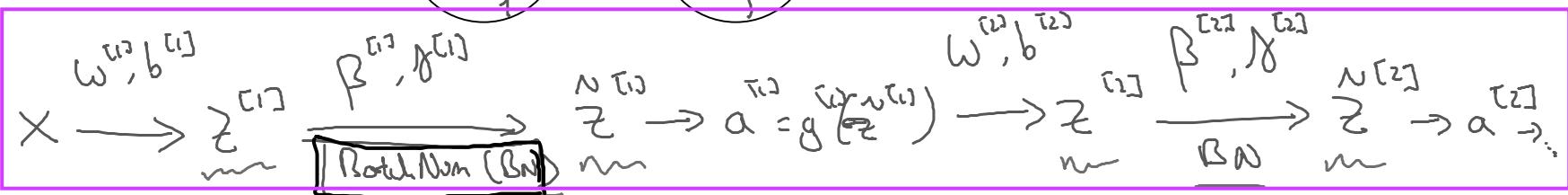
---

## Fitting Batch Norm into a neural network

# Adding Batch Norm to a network



Batch Norm  
in DL



Parameters:  $w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}$   
 $\beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)}$   
 $\rightarrow \beta$

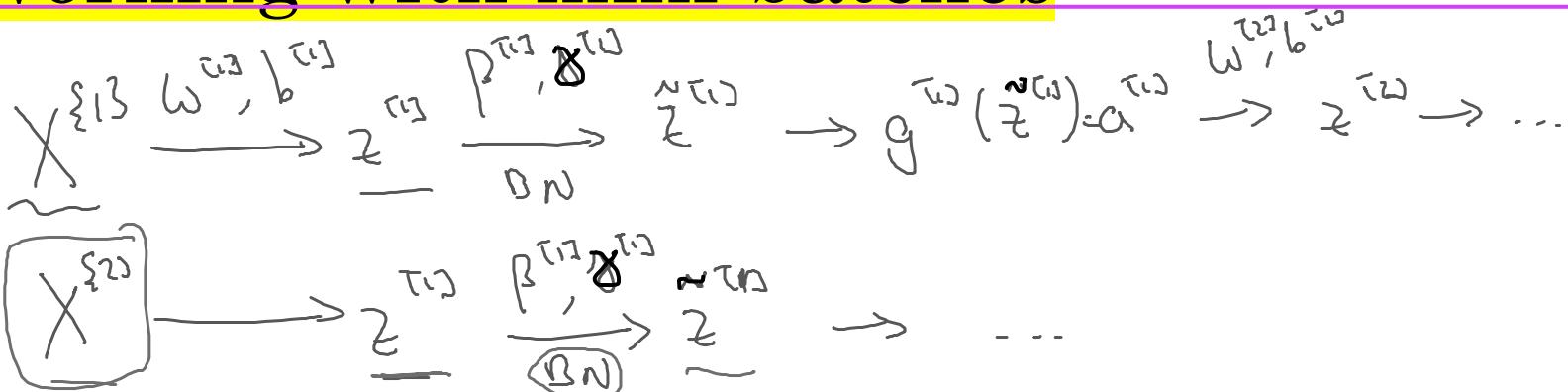
$$d\beta^{(l)} \quad \beta = \beta - d\beta^{(l)}$$

tf.nn.batch\_normalization ←  
 ↗ tf ۰۵۰ جاوزه نیست.

Remark: Here,  $\beta_{BN}$  is not  $\beta_{Adam}$  ⇒ just bad notation

In practice, BN should be combined with mini-batches

## Working with mini-batches



$X^{(1)}$   $\rightarrow \dots$

Parameters:  $w^{(1)}, \cancel{b^{(1)}}, \beta^{(1)}, \gamma^{(1)}$ .

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

$\uparrow$   
 $\uparrow$   
 $\uparrow$

$\text{dimension}$

$$\begin{aligned} \underline{z}^{(1)} &= w^{(1)} a^{(1)} + \cancel{b^{(1)}} \\ z^{(1)} &= w^{(1)} a^{(1)} \\ z^{(1)}_{\text{norm}} &= \gamma^{(1)} z^{(1)}_{\text{norm}} + \sqrt{\beta^{(1)}} \end{aligned}$$

وچن BN داری می توانی پارامتر  $b = 0$   $\Rightarrow$  چون BN میانی و دارای میانی است نتیج حاصلد، اگرچه، از بین حاصلد

## Final Implementation

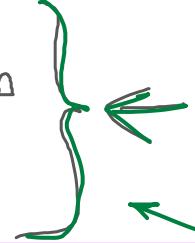
### Implementing gradient descent

for  $t = 1 \dots \text{num MiniBatches}$   
Compute forward prop on  $X^{(t)}$ .

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

Use backprop to compute  $\underline{dw}^{(l)}, \underline{dx}^{(l)}, \underline{d\beta}^{(l)}, \underline{dg}^{(l)}$

Update parameters

$$\begin{aligned} w^{(l)} &:= w^{(l)} - \alpha \underline{dw}^{(l)} \\ \beta^{(l)} &:= \beta^{(l)} - \alpha \underline{d\beta}^{(l)} \\ g^{(l)} &:= \dots \end{aligned}$$


کافی است این بخوبی Works form / momentum, RMSprop, Adam.  $\Leftrightarrow$  updating کنیم.



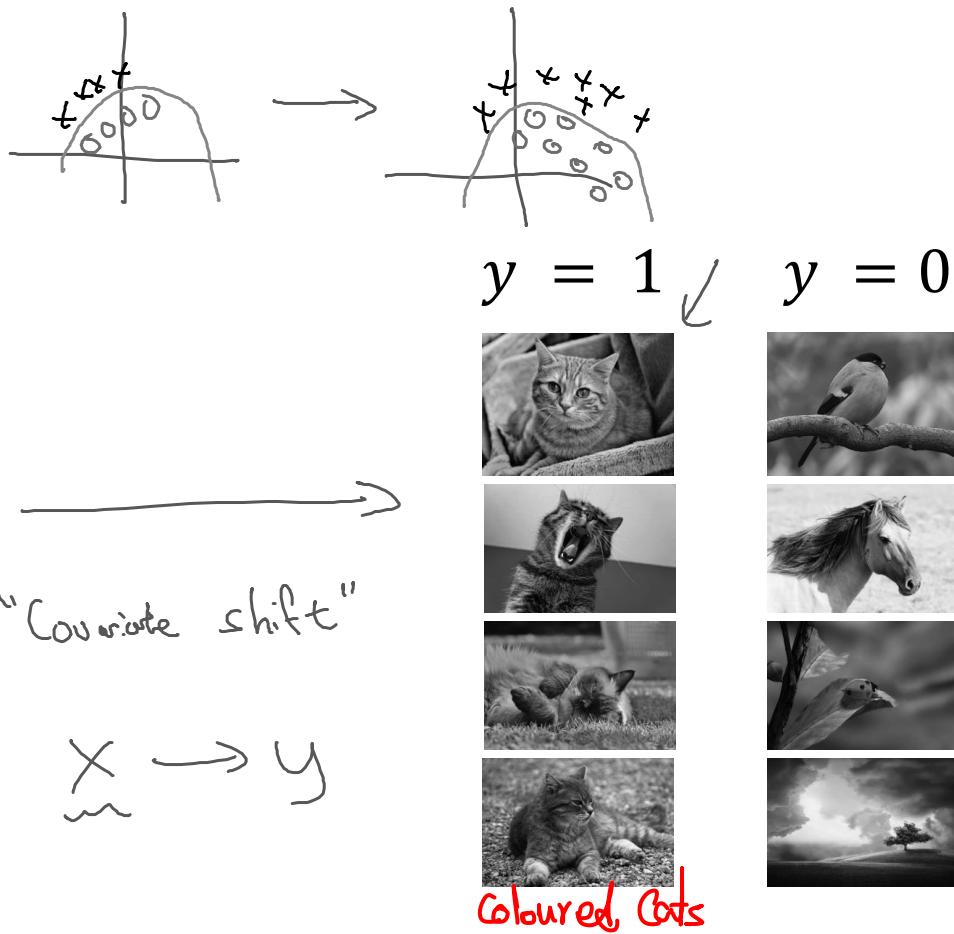
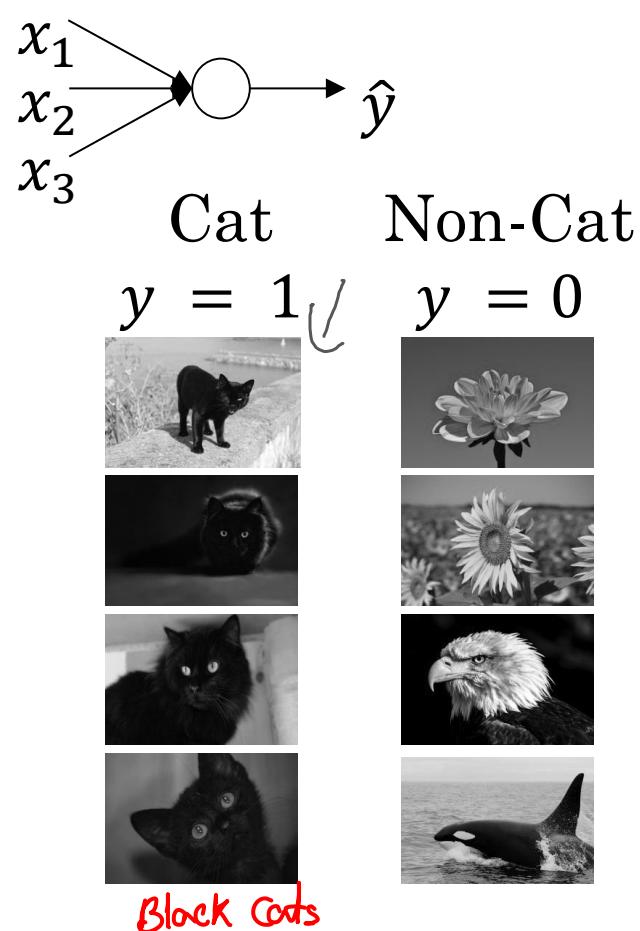
deeplearning.ai

# Batch Normalization

---

## Why does Batch Norm work?

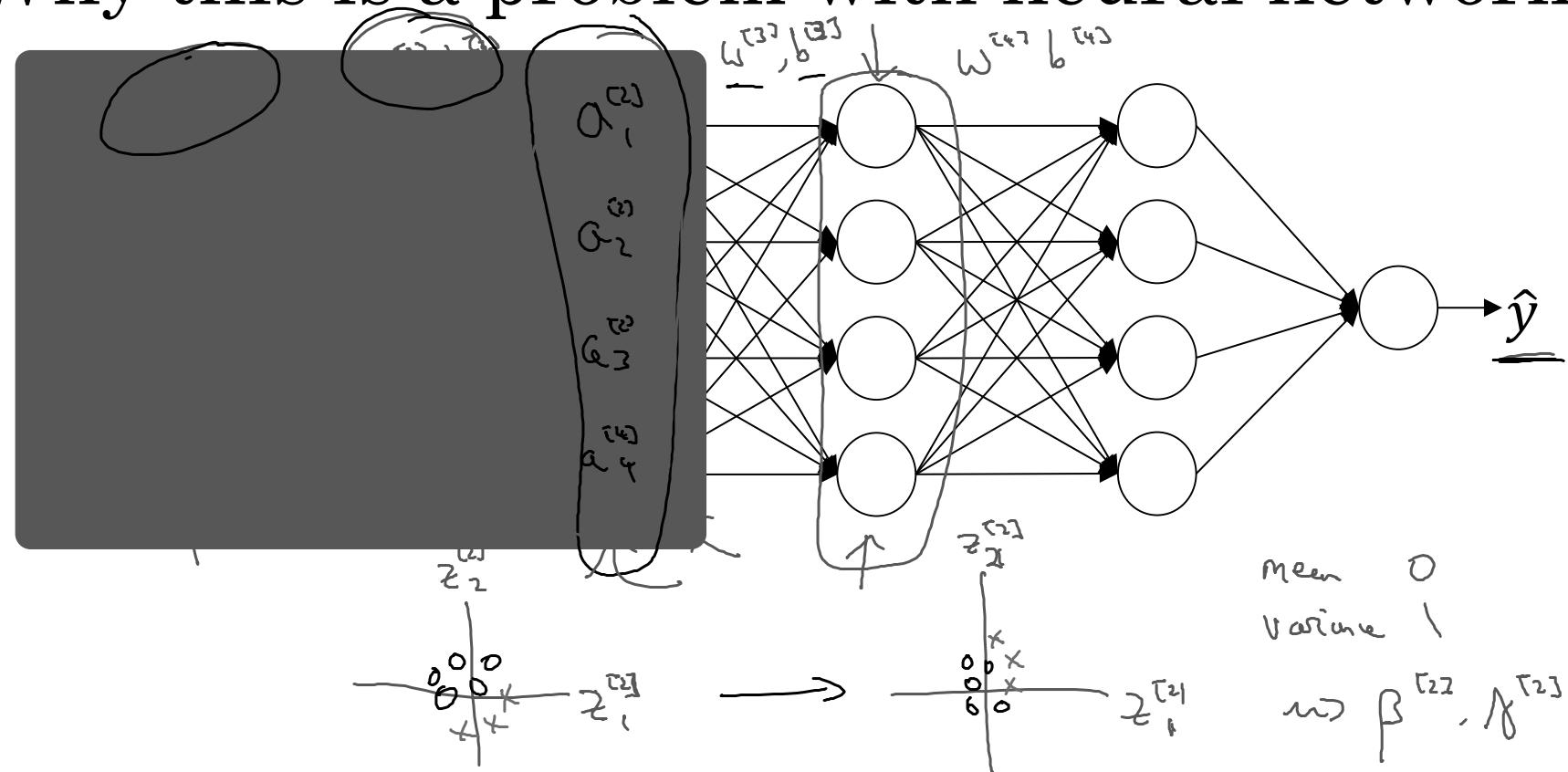
# Learning on shifting input distribution



آندر classifier بای کویهای سیله بسازی، الزاماً بای کویهای دئن خوب چارمه کند چون covariate shift داشته باش.

BN limits the amount to which updating the parameters in the earlier layers can affect the distribution of values that the third layer sees. Hence, it reduces the problem of the input values changing: it makes these values more stable so that the later layers have more firm ground to stand on.

# 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.
- 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.

$\tilde{z}^{[l]}$

$\underline{64}, \underline{128}$

$X^{\{+1\}}$

$\underline{2}^{(0)}$

$\mu, \sigma^2$

mini-batch :  $\underline{64} \longrightarrow \underline{512}$



deeplearning.ai

# Batch Normalization

---

## Batch Norm at test time

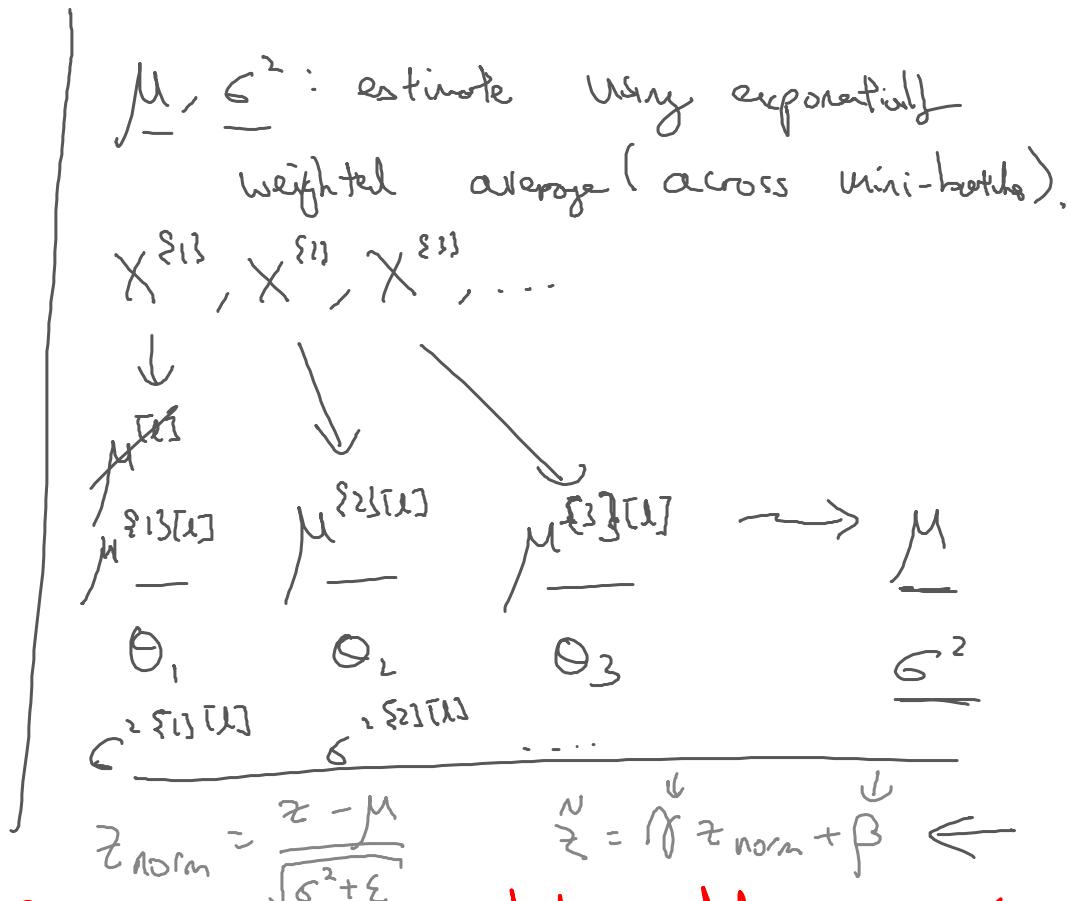
BN handles data one mini-batch at a time: it computes mean and variances on each mini-batch.

At test time, you might not have a mini-batch of examples: you might be processing one single example at the time. So, you need to do something different at test time.

## Standard BN (for training)

# Batch Norm \at test time

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



training minibatches over  $\sigma^2, \mu$  جزوی از داده ها، ممکن است جاید exponentially weighted ave با کم و سریع BN (for test) بدل شود و سپسی z را ایجاد کند.



deeplearning.ai

Multi-class  
classification

---

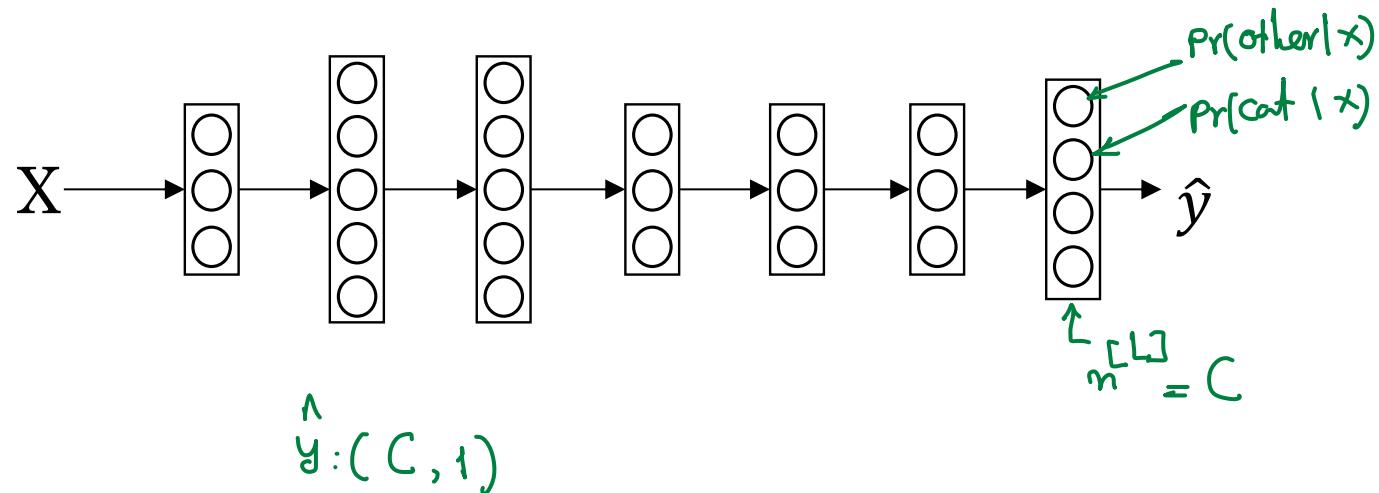
Softmax regression

Four classes  $\{1, 2, 3, 0\}$   
 Recognizing cats, dogs, and baby chicks other



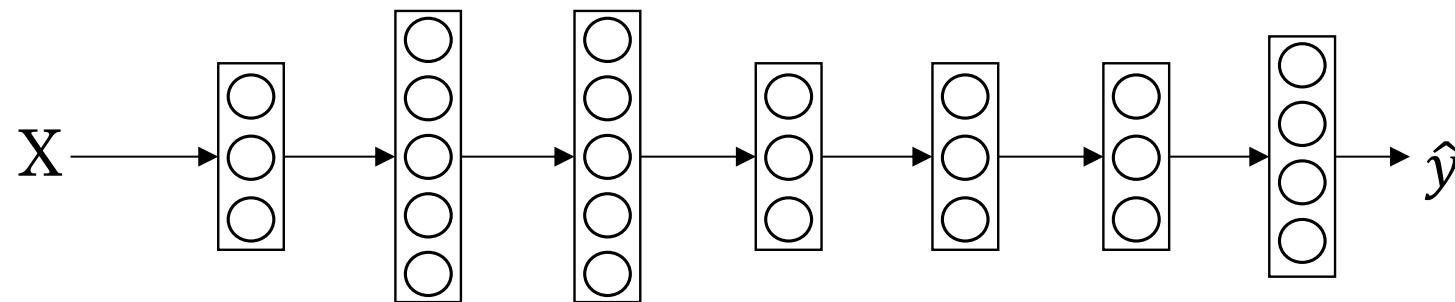
3      1      2      0      3      2      0      1

$$C = \# \text{classes} = 4 \quad (0, \dots, 3)$$



$$a^{[L]} = \frac{e^{z_i^L}}{\sum_{i=1}^n t_i}$$

## Softmax layer



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

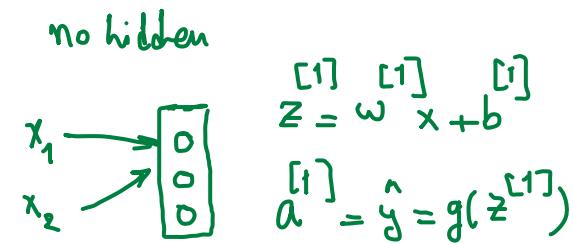
Activation Function

$$a^{[L]} = \frac{\exp(z_i^{[L]})}{\sum_{i=1}^n \exp(z_i^{[L]})} = \frac{t_i}{\sum t_i} \text{ where } t_i = \exp(z_i^{[L]})$$

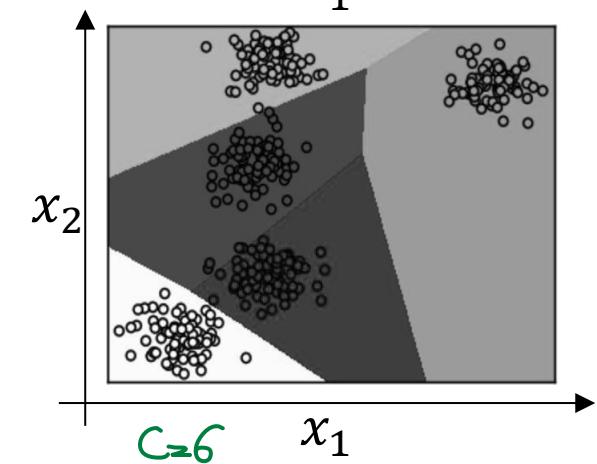
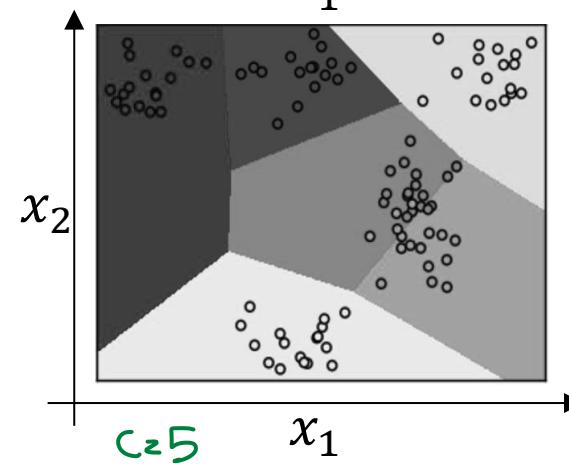
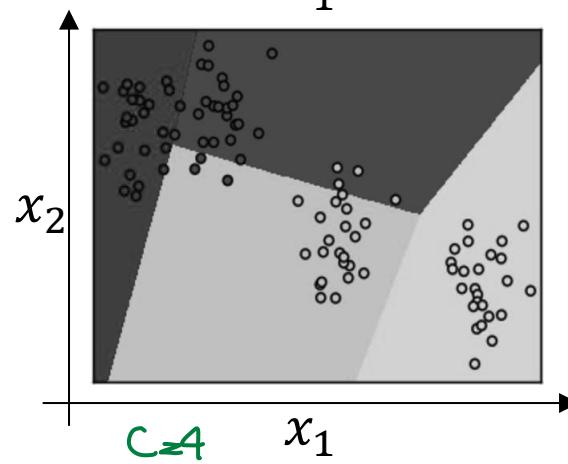
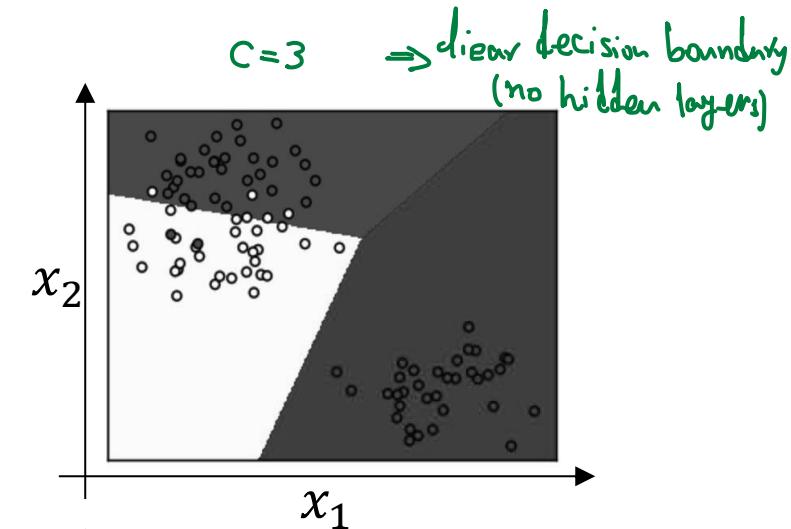
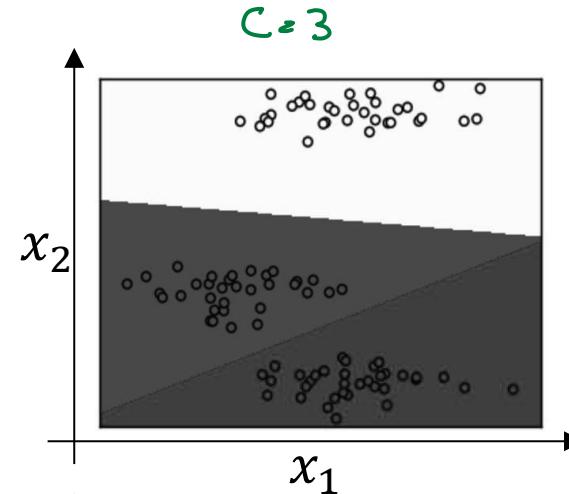
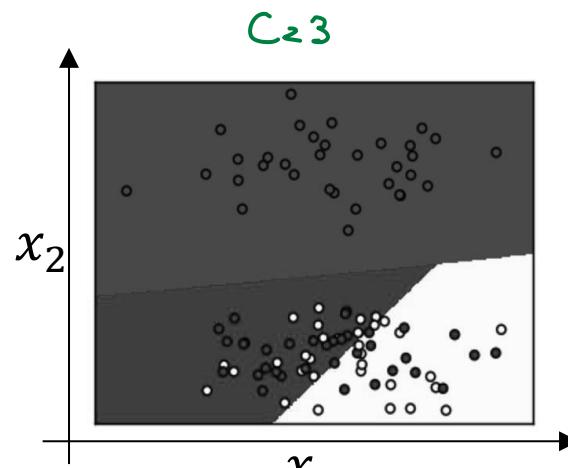
$\Rightarrow$

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

softmax activation function



## Softmax examples





deeplearning.ai

Multi-class  
classification

---

Trying a softmax  
classifier

# Understanding softmax

$$\begin{array}{l} \text{Softmax} \\ \text{Input } z = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \Rightarrow t_z = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} \Rightarrow \alpha = \begin{bmatrix} 0,84 \\ 0,04 \\ 0,002 \\ 0,11 \end{bmatrix} \\ \text{if } C=4 \end{array}$$

$$\text{Hardmax} \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Softmax regression generalizes logistic regression to C classes.  $\Rightarrow$  If  $C=2$ ,  $\text{Softmax} = \text{Logistic}$

# Loss function

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j \quad \text{and} \quad J(w^{[l]}, b^{[l]}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

With DL frameworks, you just need to focus on the forward prop: the primary will figure out how to do back prop.



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

# Code example

$$J(w) = w^2 - 10w + 25 \Rightarrow w_{\min} = 5$$

```
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))
```

equiv  
≡

```
{ with tf.Session() as session:  
    session.run(init)  
    print(session.run(w))
```

```
for i in range(1000):  
    session.run(train, feed_dict={x:coefficients})  
    print(session.run(w))
```