

# Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>



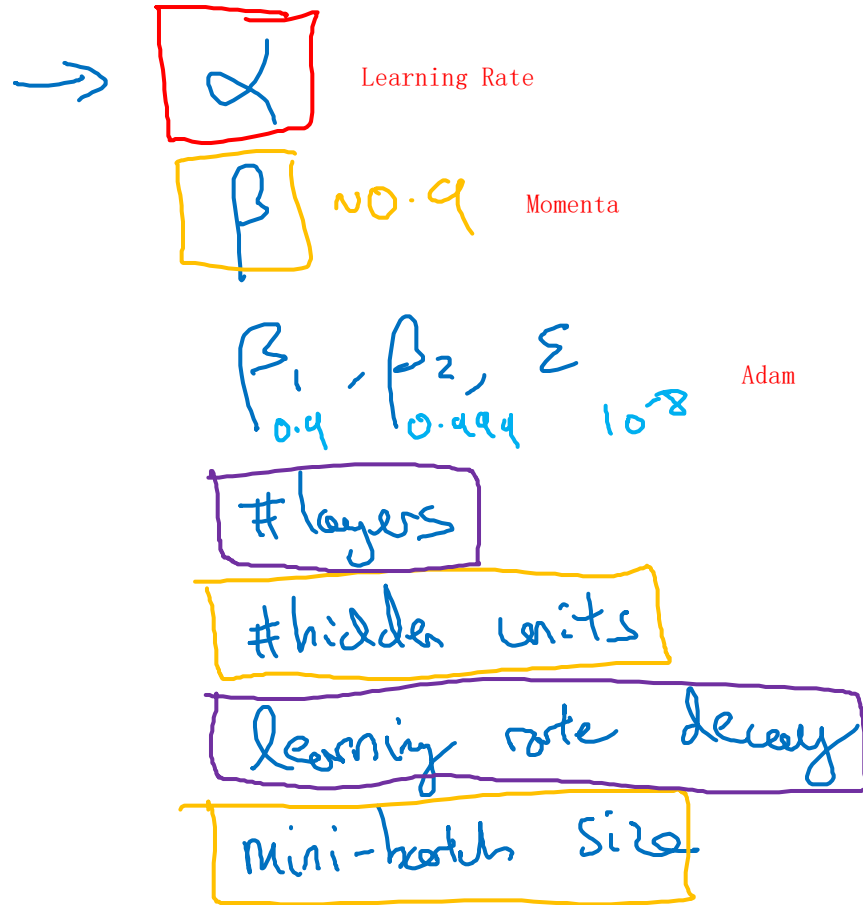
deeplearning.ai

# Hyperparameter tuning

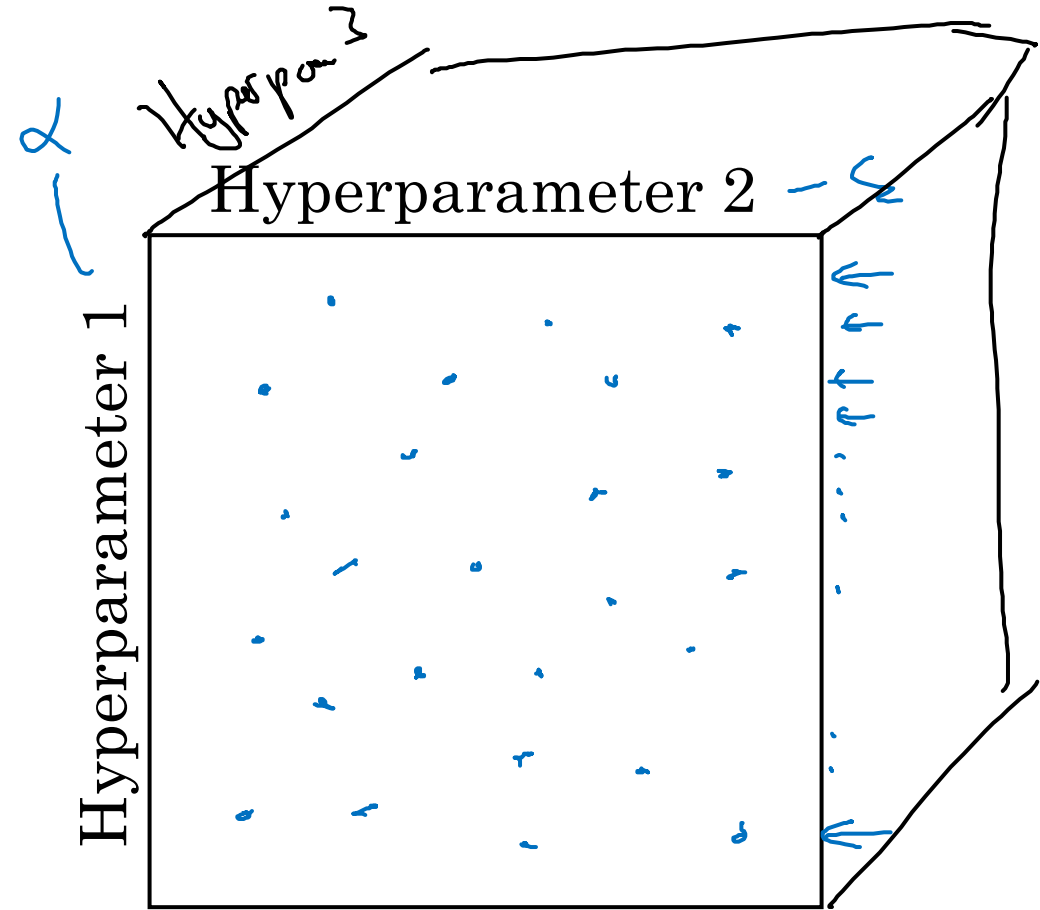
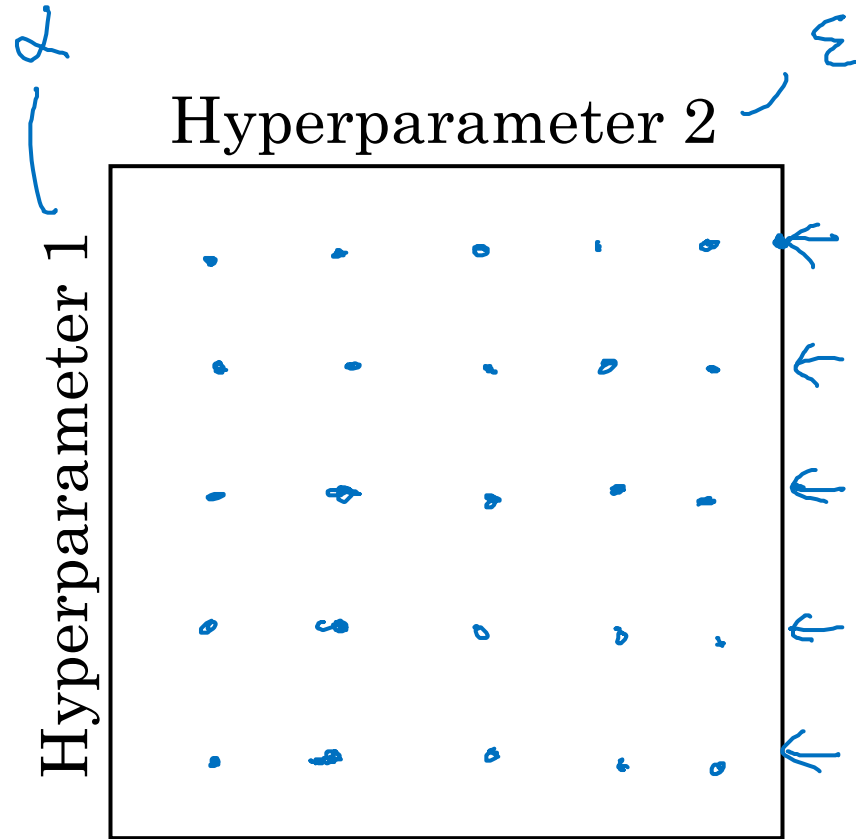
---

## Tuning process

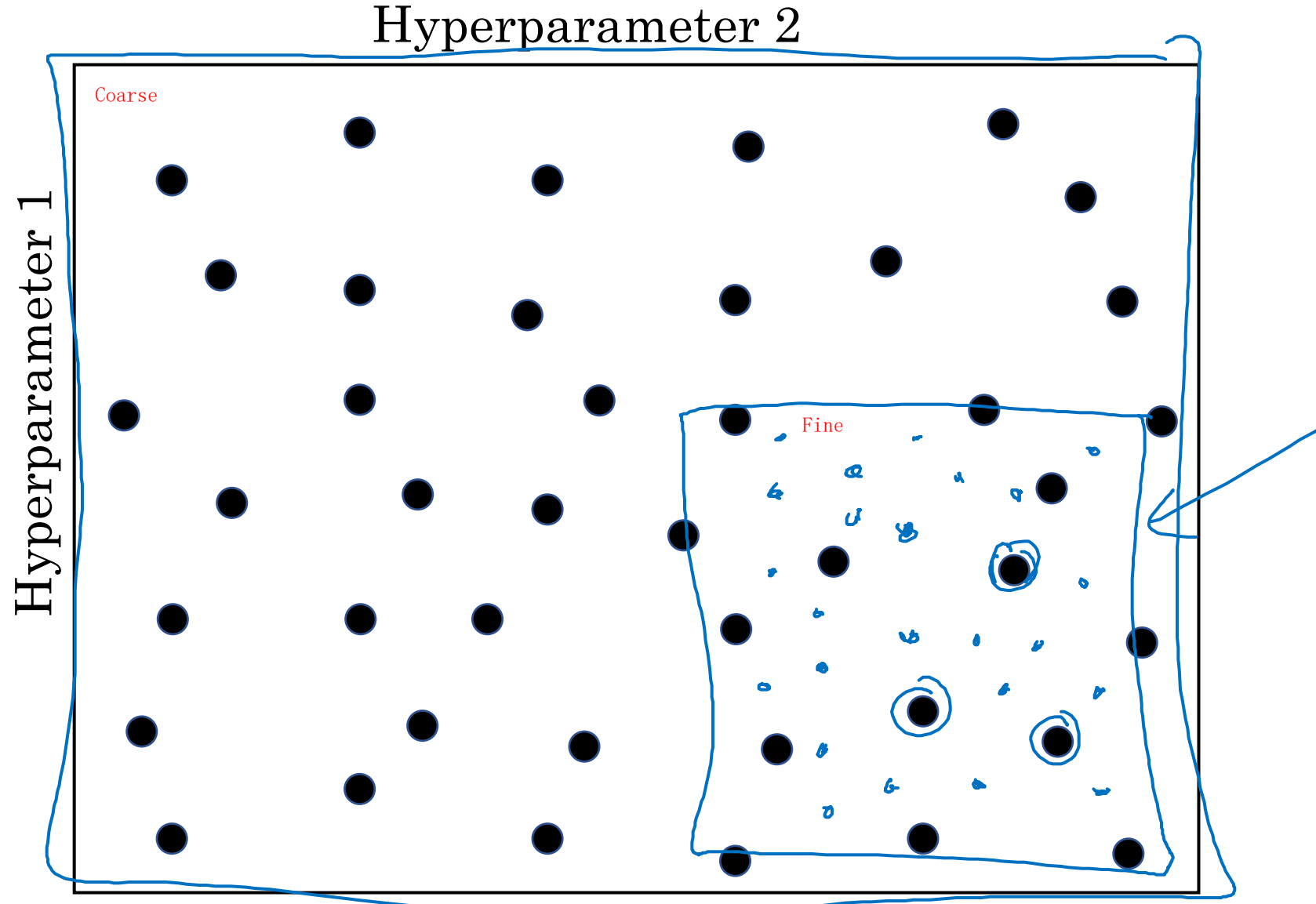
# Hyperparameters



# Try random values: Don't use a grid



# Coarse to fine





deeplearning.ai

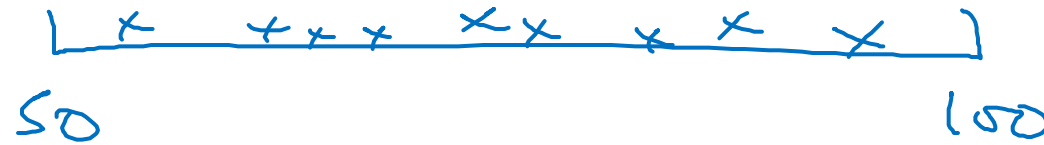
# Hyperparameter tuning

---

Using an appropriate  
scale to pick  
hyperparameters

# Picking hyperparameters at random

→  $n^{\text{test}} = 50, \dots, 100$

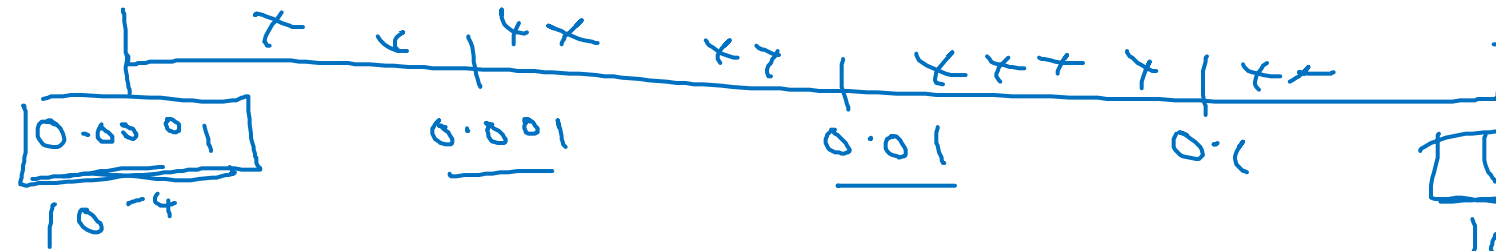
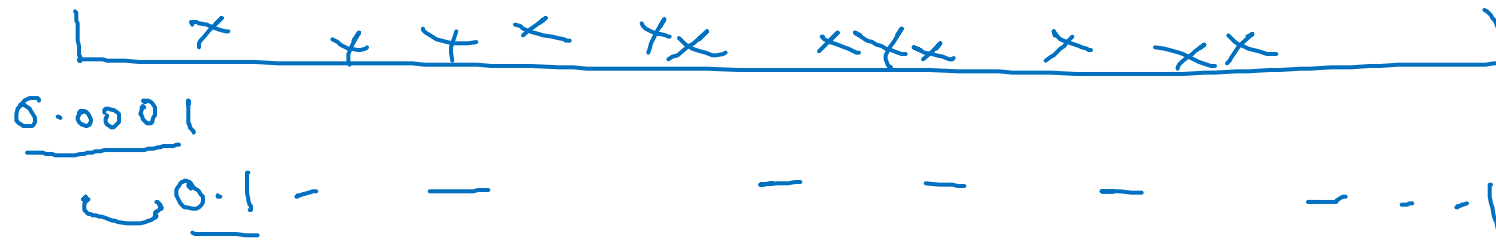


→ #layers     $L : 2 - 4$     sometimes sampling uniformly at random might be reasonable

2, 3, 4

# Appropriate scale for hyperparameters

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



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

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

$$\alpha = 10^r$$

$$\underline{10^a \dots 10^b}$$

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

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

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

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

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



# Hyperparameters for exponentially weighted averages

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

$\downarrow$                        $\downarrow$   
 average over last 10                      1000

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

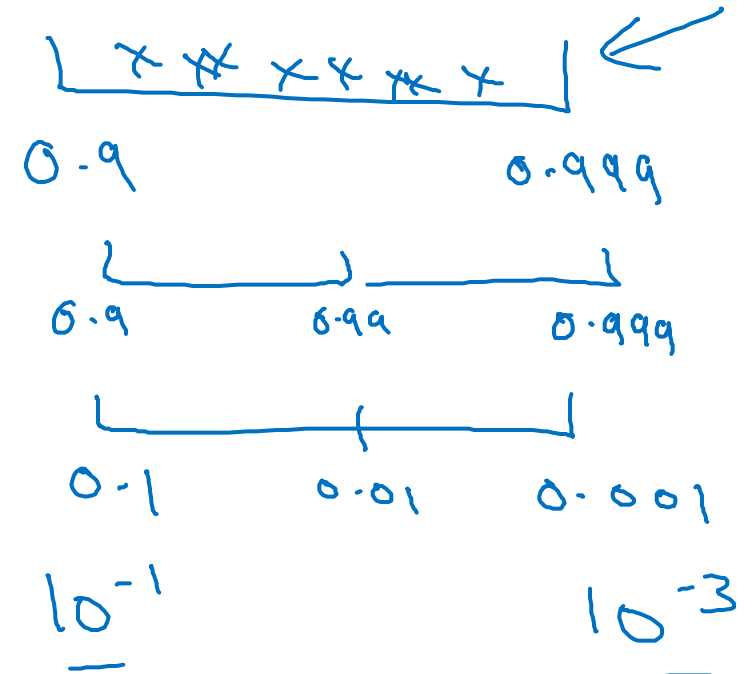
$$\beta: 0.999 \rightarrow 0.9995 \quad \} \sim 10$$

$$\beta: 0.999 \rightarrow 0.9995$$

$\sim 1000$                        $\sim 2000$

become away  
more sensitive

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



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

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

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



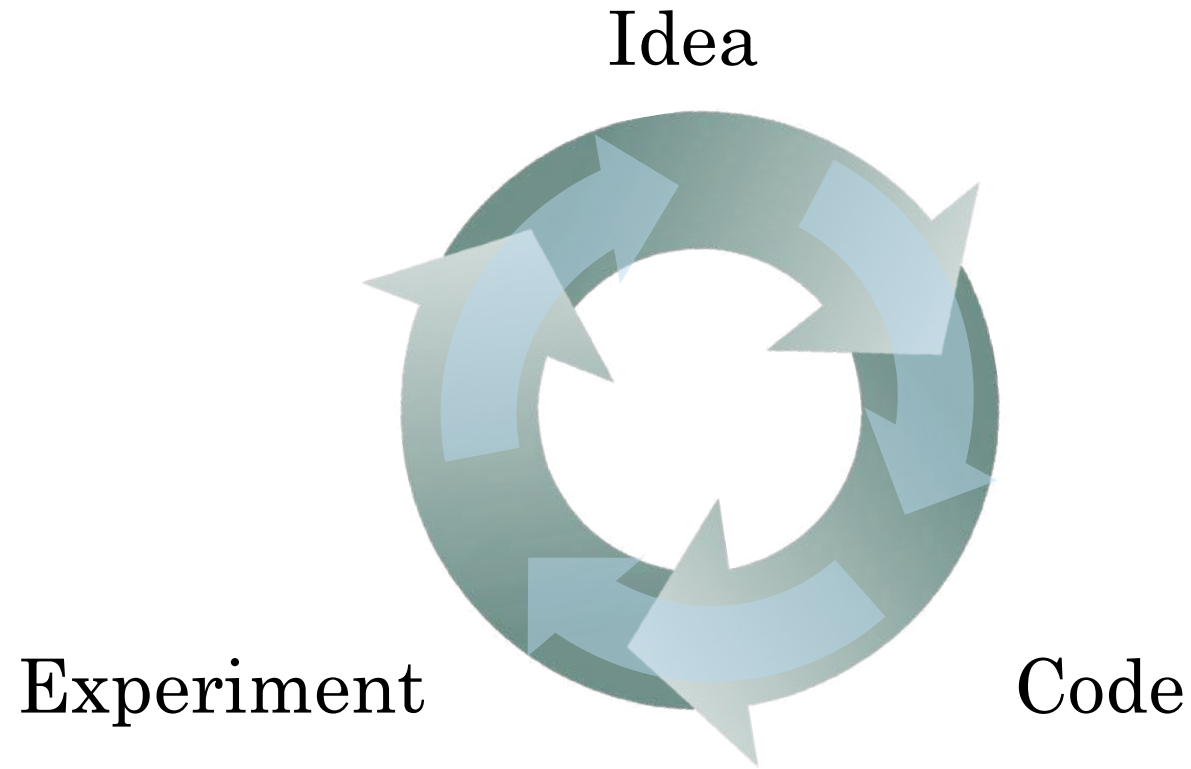
deeplearning.ai

# Hyperparameters tuning

---

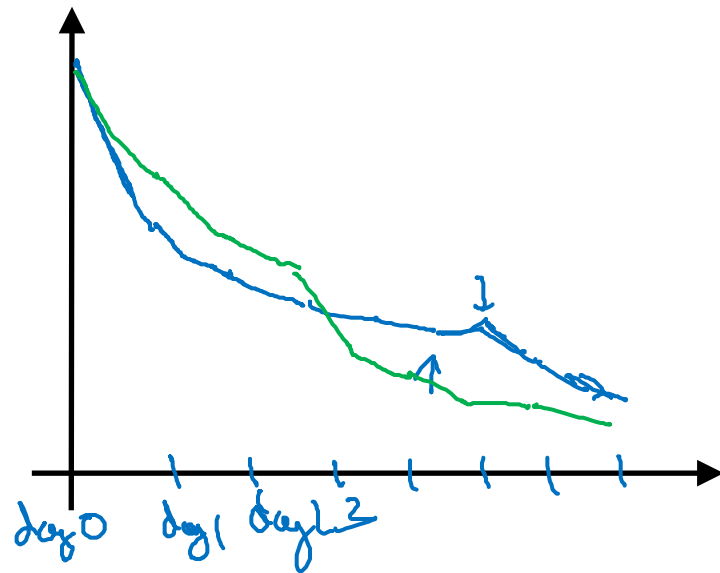
Hyperparameters  
tuning in practice:  
Pandas vs. Caviar

# Re-test hyperparameters occasionally



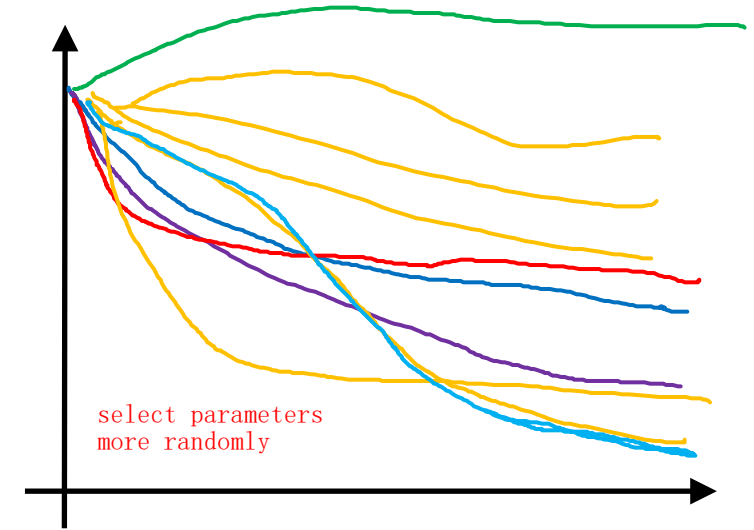
- NLP, Vision, Speech,  
Ads, logistics, ....
- Intuitions do get stale.  
Re-evaluate occasionally.

# Babysitting one model



Panda ←

# Training many models in parallel



Caviar ←



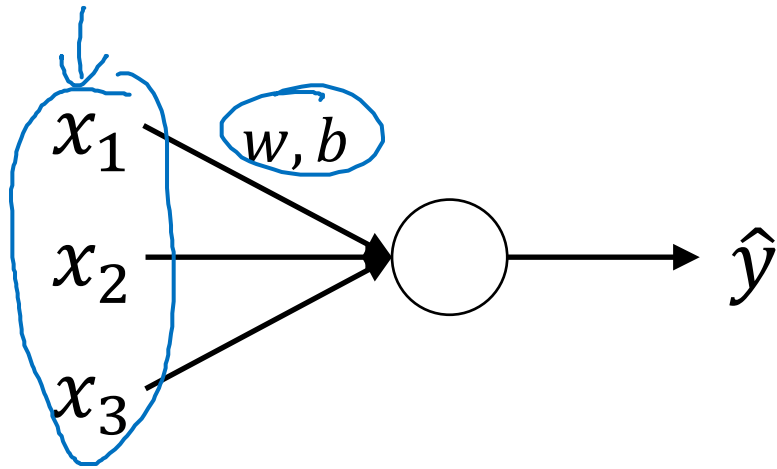
deeplearning.ai

# Batch Normalization

---

Normalizing activations  
in a network

# Normalizing inputs to speed up learning

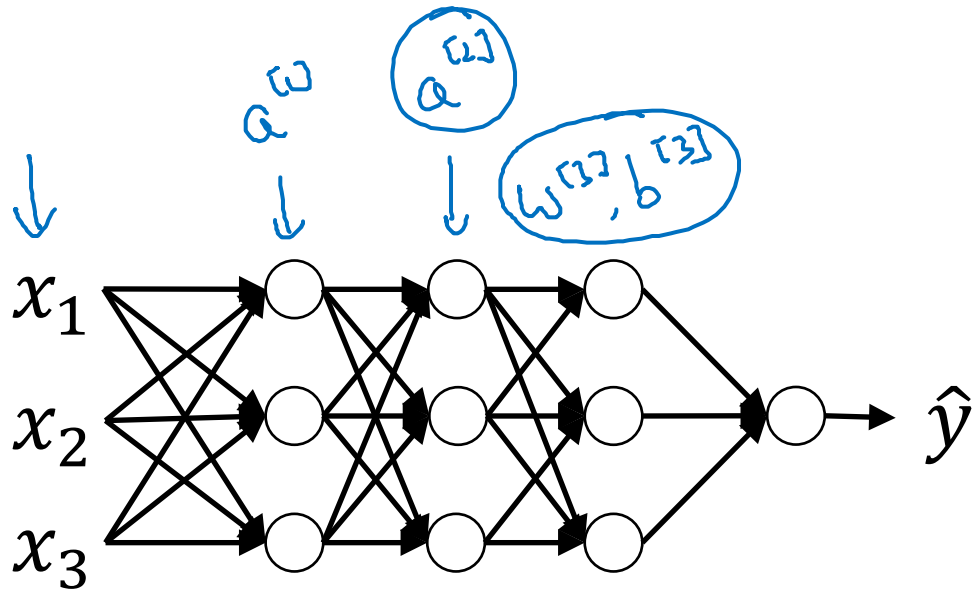
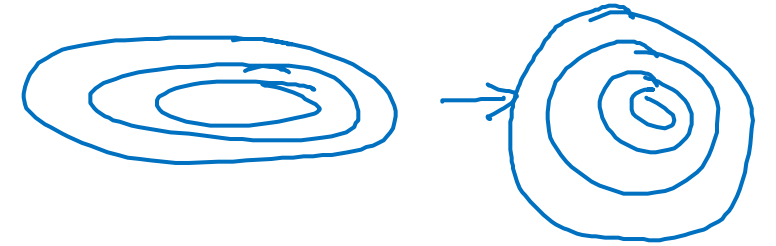


$$\mu = \frac{1}{n} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

$$X = X / \sigma$$



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

Normalize  $\frac{z^{[2]}}{\sigma}$   
↑

# Implementing Batch Norm

Given some intermediate values in NN

$$\begin{matrix} \downarrow & \downarrow \\ z^{(1)} & \dots, z^{(m)} \end{matrix}$$

$z^{[l]}(i)$

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

If

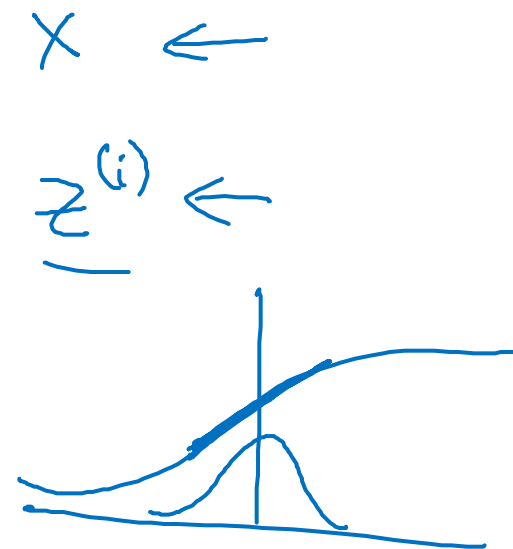
$$\gamma = \sqrt{\sigma^2 + \epsilon}$$

$$\beta = \mu$$

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

learnable parameters of model.

not necessarily mean 0, variance 1



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

apply normalization not only in the input layer but also in hidden layers



deeplearning.ai

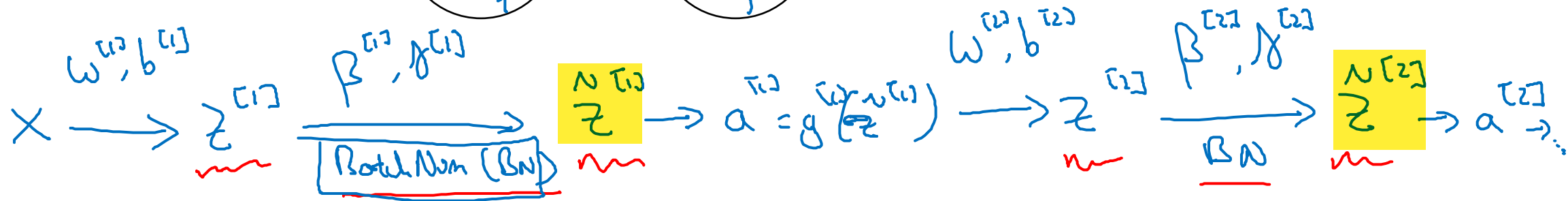
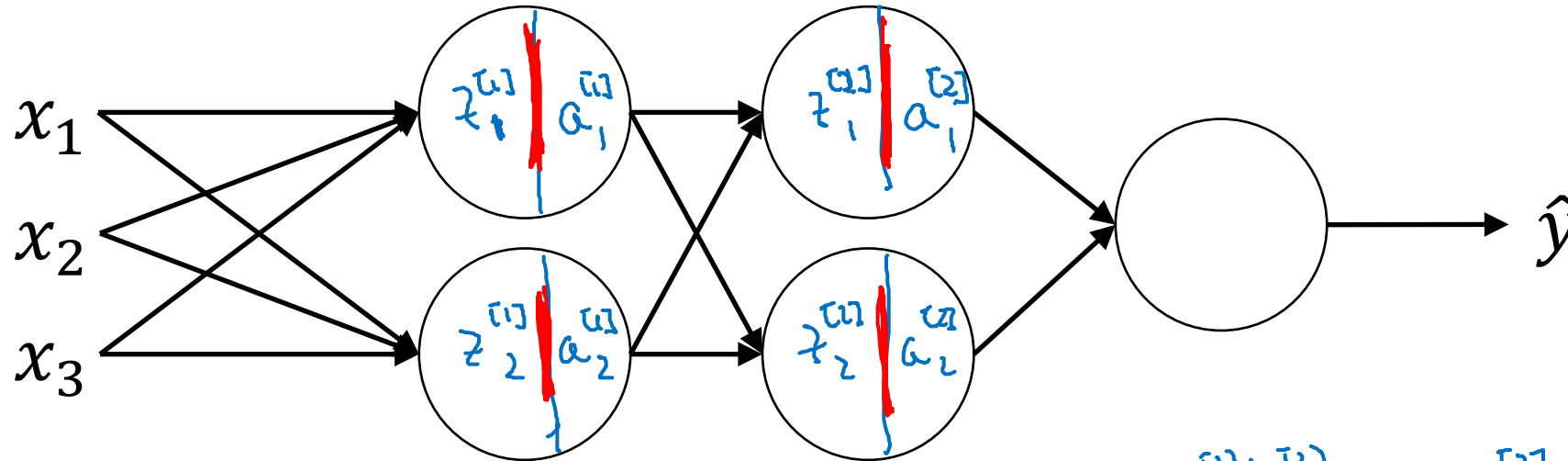
# Batch Normalization

---

Fitting Batch Norm  
into a neural network



# Adding Batch Norm to a network



Parameters:  $\left\{ W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]} \right\}$   
 $\rightarrow \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$   
 $\rightarrow \beta$

$d\beta^{[2]}$   
 $\beta = \beta - \alpha d\beta^{[2]}$

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

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

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

subtracted out

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

$\uparrow$

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

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

# 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  $\underline{dw}^{\{t\}}$ ,  ~~$\underline{db}^{\{t\}}$~~ ,  $\underline{dp}^{\{t\}}$ ,  $\underline{df}^{\{t\}}$

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

Works w/ momentum, RMSprop, Adam.



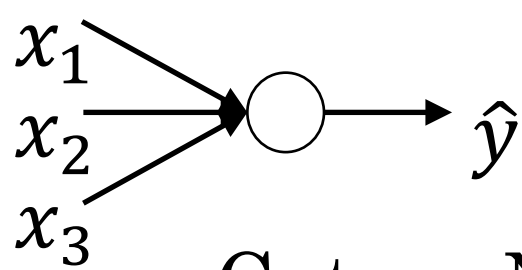
deeplearning.ai

# Batch Normalization

---

Why does  
Batch Norm work?

# Learning on shifting input distribution

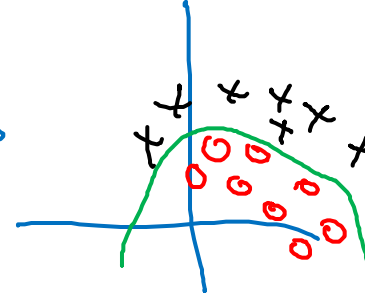
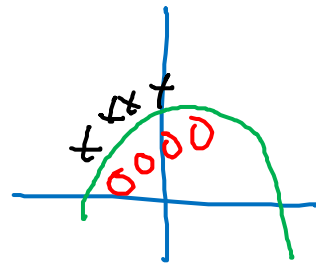


Cat

Non-Cat

$y = 1$

$y = 0$



$y = 1$

$y = 0$

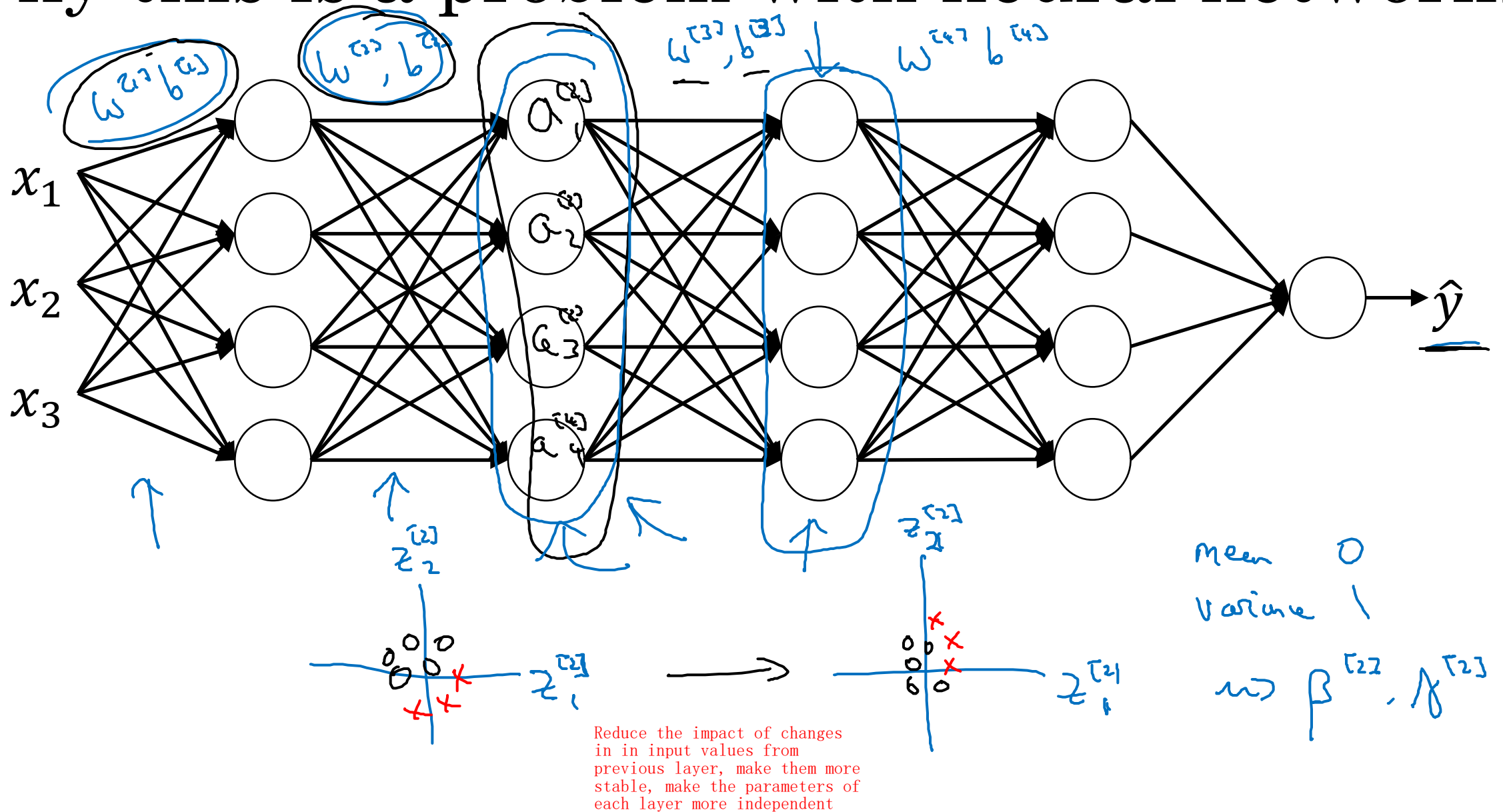


"Covariate shift"

if the distribution  
of  $x$  changes, then  
we need to retrain  
the model

$x \rightarrow y$

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



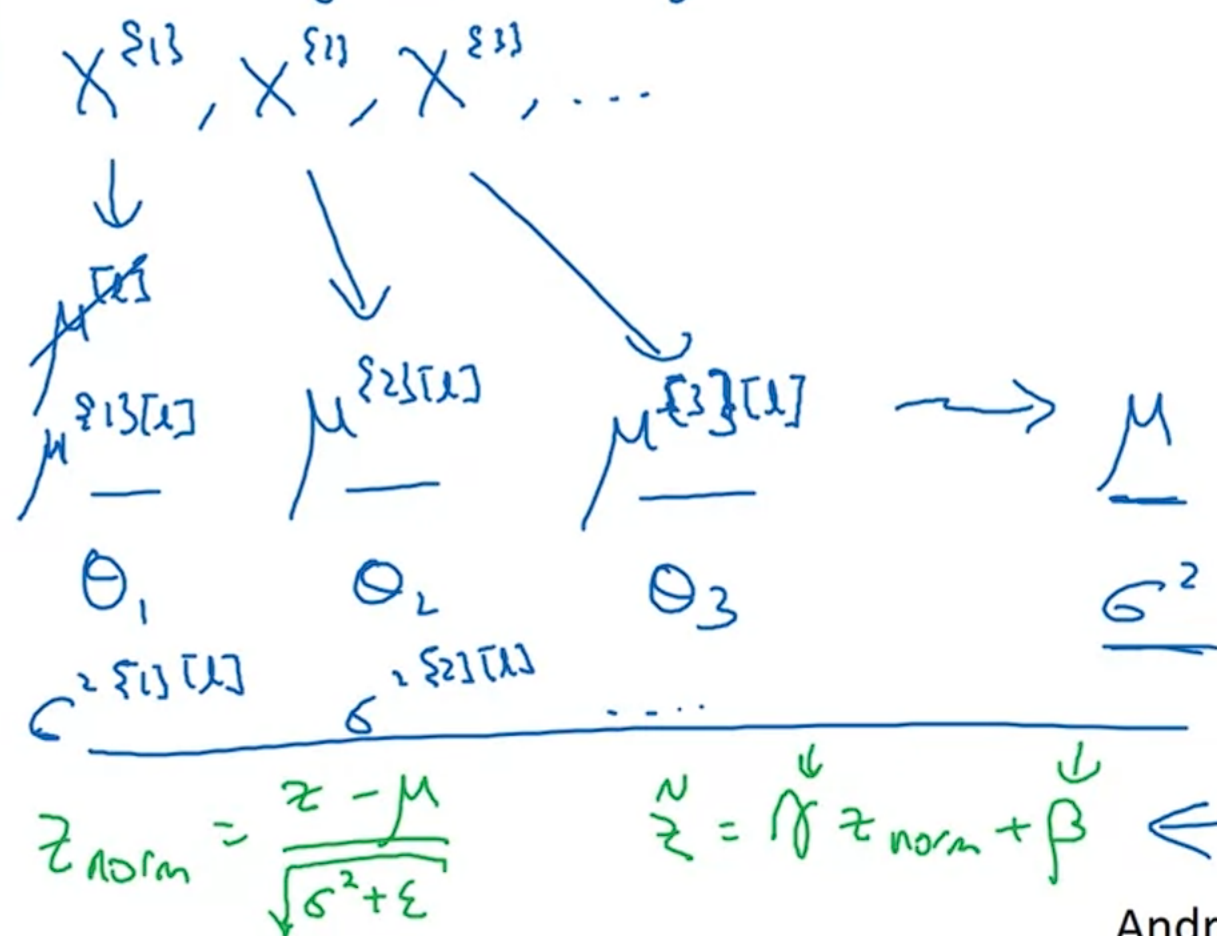
# Batch Norm at test time

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

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

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

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







deeplearning.ai

# Multi-class classification

---

## Softmax regression

# Recognizing cats, dogs, and baby chicks



3



1



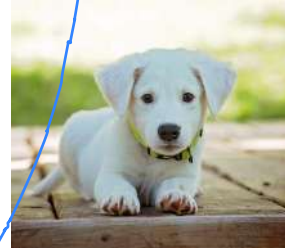
2



0



3



2

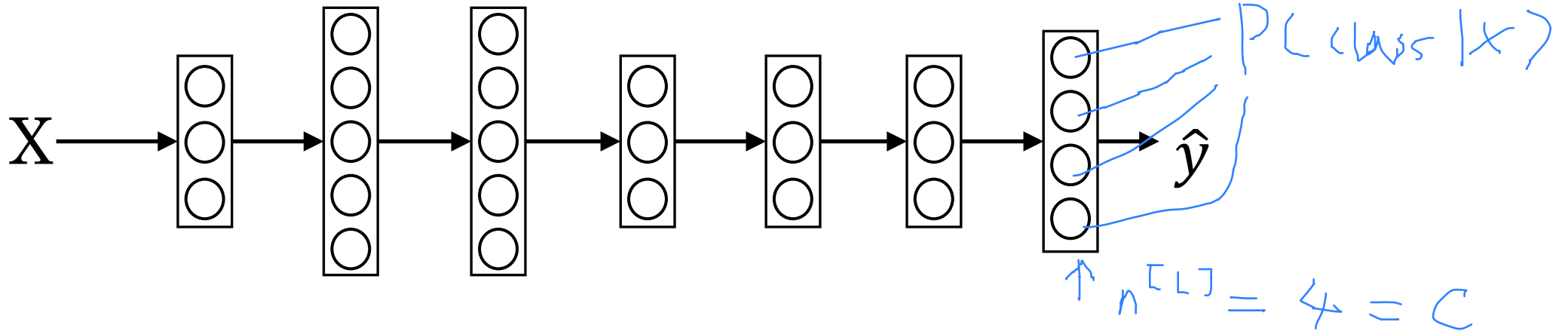


0

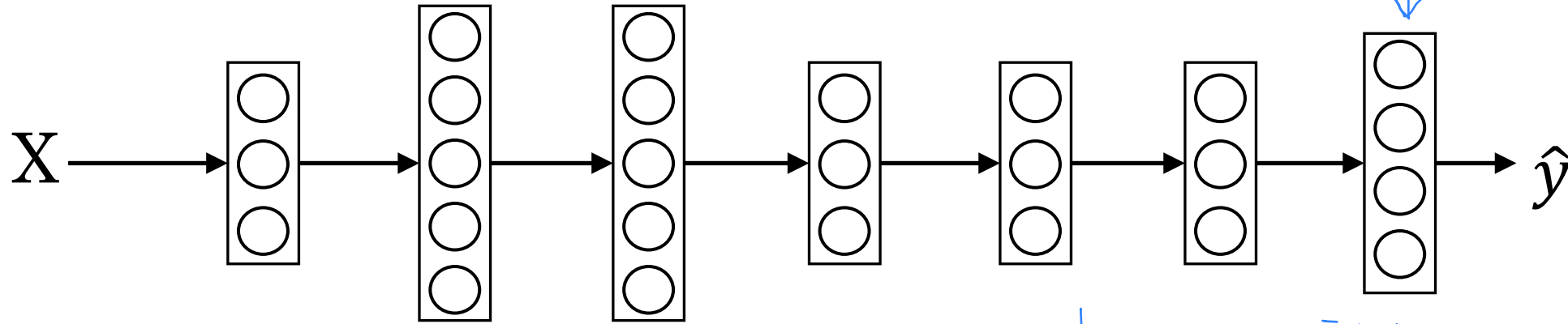


1

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



# Softmax layer



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

$(4, 1)$

Activation Function:

$$\rightarrow t = e^{z^{[L]}}$$

$$\rightarrow a^{[L]} = \frac{e^{z^{[L]}}}{\sum_{j=1}^4 t_j}, \quad a_2^{[L]} = \frac{t_2}{\sum_{j=1}^4 t_j}$$

$(4, 1)$

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \leftarrow$$

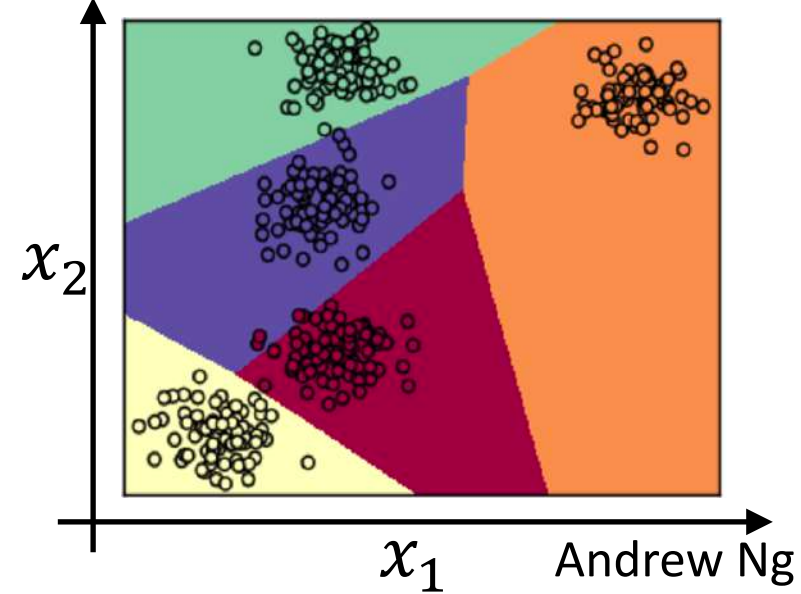
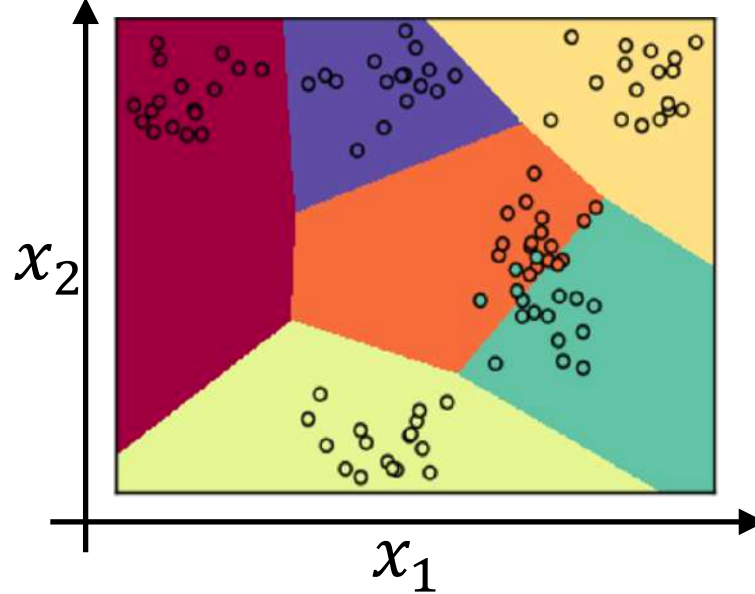
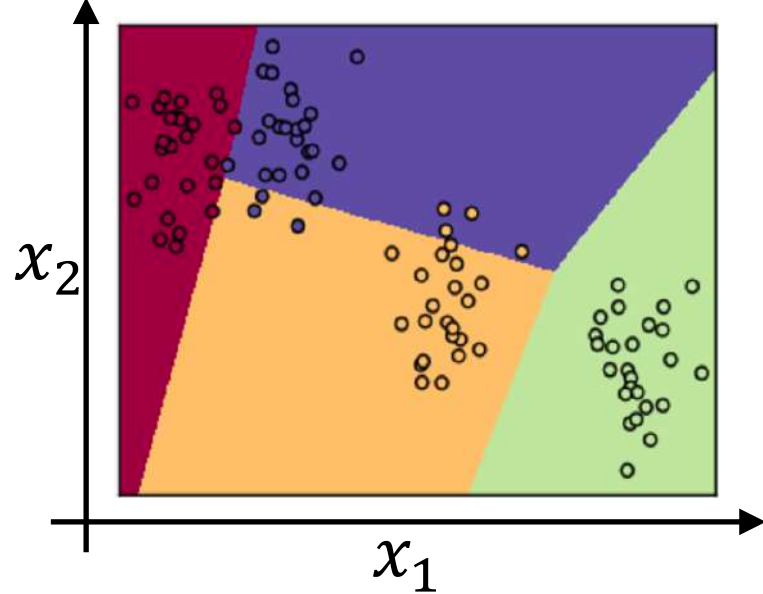
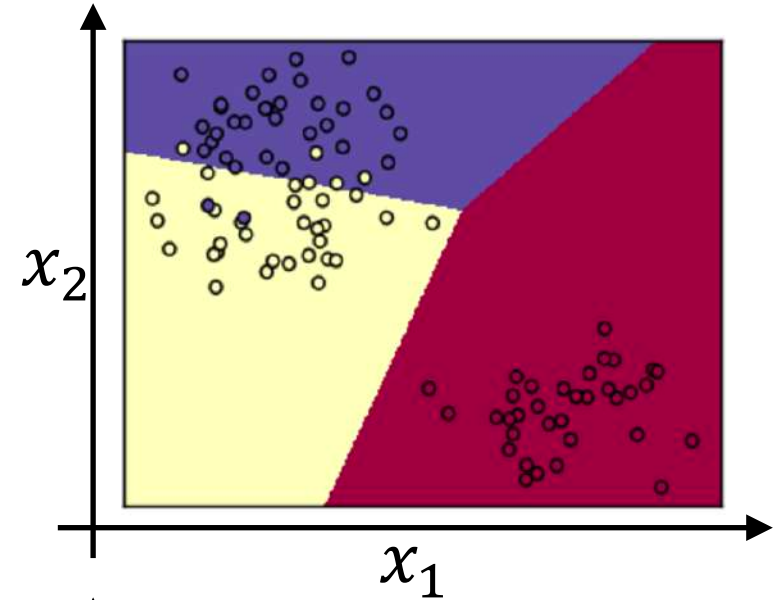
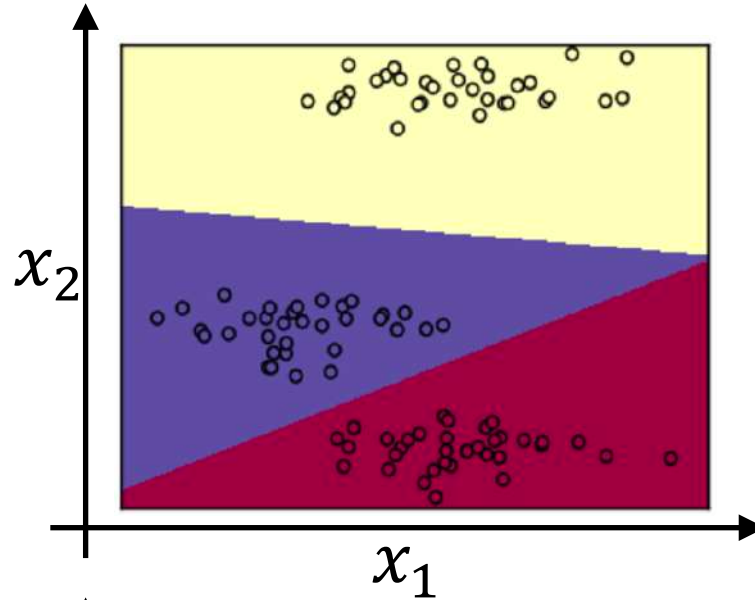
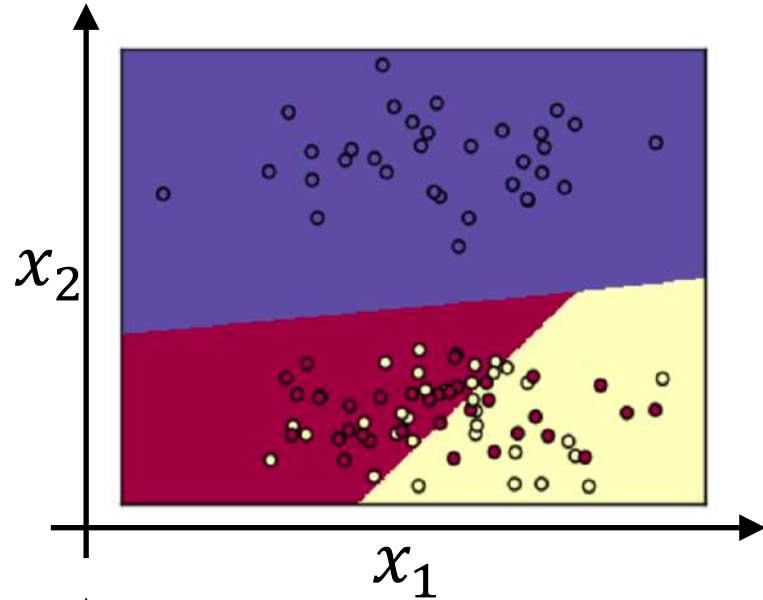
$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.7 \end{bmatrix} \cdot \sum_{j=1}^4 t_j = 176.3$$

$$a^{[L]} = t / 176.3$$

# Softmax examples

$$\begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \rightarrow \hat{y}$$

$$z^{[L]} = W^{[L]}x + b^{[L]}$$
$$a^{[L]} = \hat{y} = \sigma(z^{[L]})$$



# Understanding softmax

(4,1)

$$z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix} \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix}$$

$C=4$

$g^{[L]}(\cdot)$

"soft max"

$$a^{[L]} = g^{[L]}(z^{[L]}) = \begin{bmatrix} e^5 / (e^5 + e^2 + e^{-1} + e^3) \\ e^2 / (e^5 + e^2 + e^{-1} + e^3) \\ e^{-1} / (e^5 + e^2 + e^{-1} + e^3) \\ e^3 / (e^5 + e^2 + e^{-1} + e^3) \end{bmatrix} = \begin{bmatrix} 0.842 \\ 0.042 \\ 0.002 \\ 0.114 \end{bmatrix}$$

"hard max"

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

Softmax regression generalizes logistic regression to  $C$  classes.

If  $C=2$ , softmax reduces to logistic regression.  $a^{[L]} = \begin{bmatrix} 0.842 \\ 0.158 \end{bmatrix}$



# Loss function

$(4,1)$

$y^{(1)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$  - cat  
 $y_2 = 1$   
 $y_1 = y_3 = y_4 = 0$

small

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^C y_j \log \hat{y}_j$$

$$- y_2 \log \hat{y}_2 = - \log \hat{y}_2$$

Make  $\hat{y}_2$  big.

$(4,1)$

$a^{(1)} \approx \hat{y}^{(1)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$

$C=4$

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

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}$$

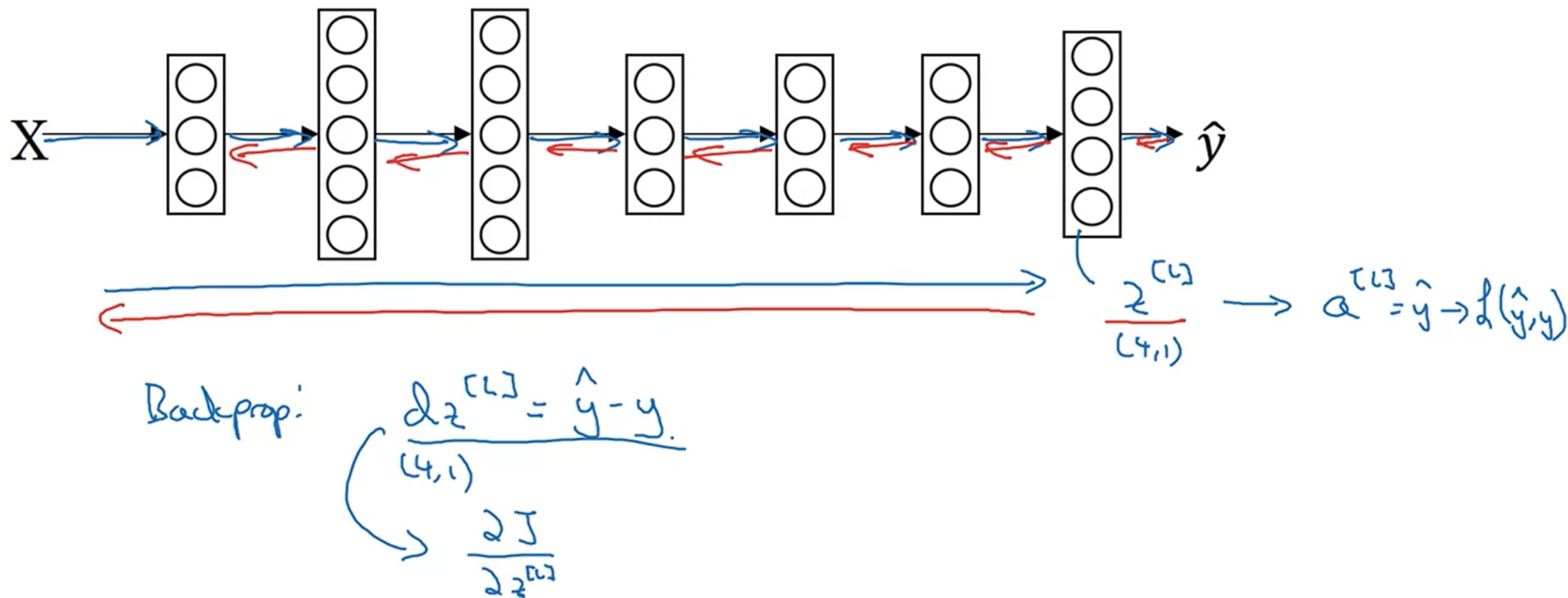
$(4,m)$

$$\hat{Y} = [\hat{y}^{(1)} \ \dots \ \hat{y}^{(m)}]$$

$$= \begin{bmatrix} 0.3 & \dots \\ 0.2 & \dots \\ 0.1 & \dots \\ 0.4 & \dots \end{bmatrix}$$

$(4,m)$

# Gradient descent with softmax





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))
```

```
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))
```

