



deeplearning.ai

# Hyperparameter tuning

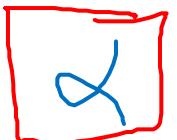
---

## Tuning process

# Hyperparameters

\*\*\* most important to tune

\*\*\*



(learning rate)

\*\*



$\beta_{0.9}$

(momentum term)

\*

least important to tune

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

(adam)

\*



\*\*



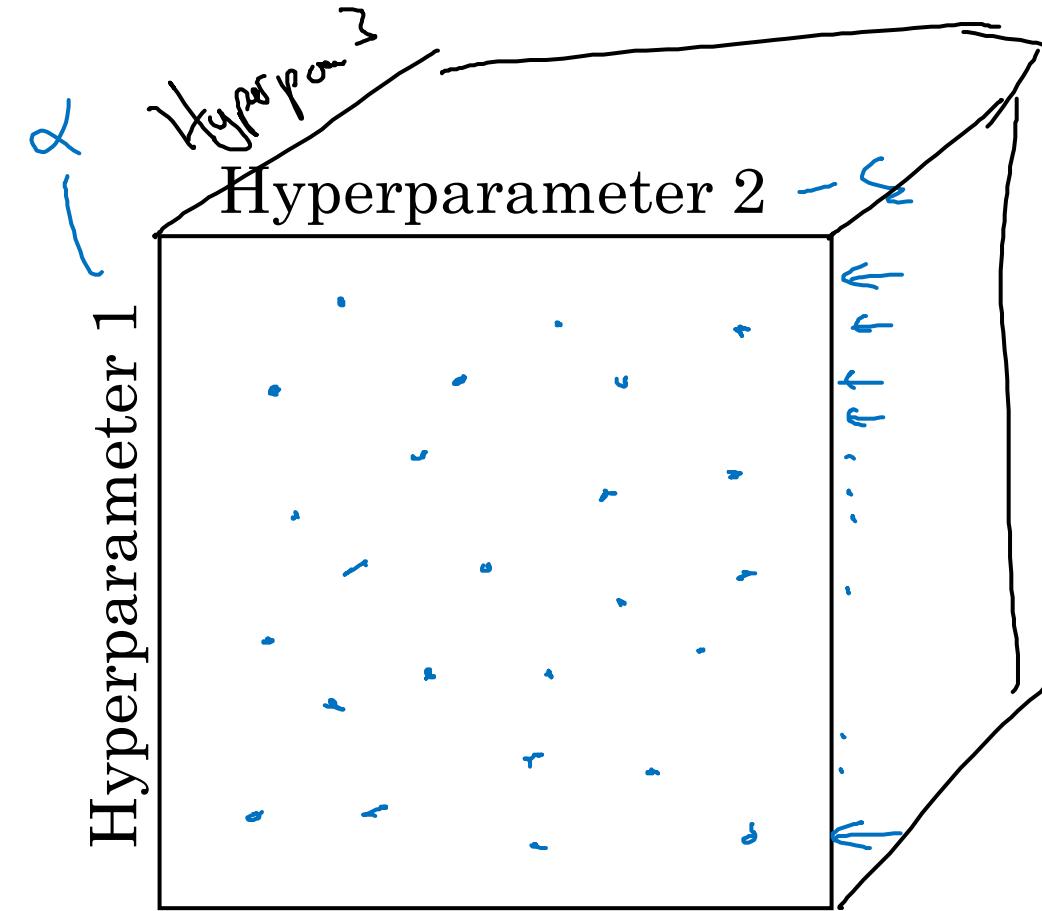
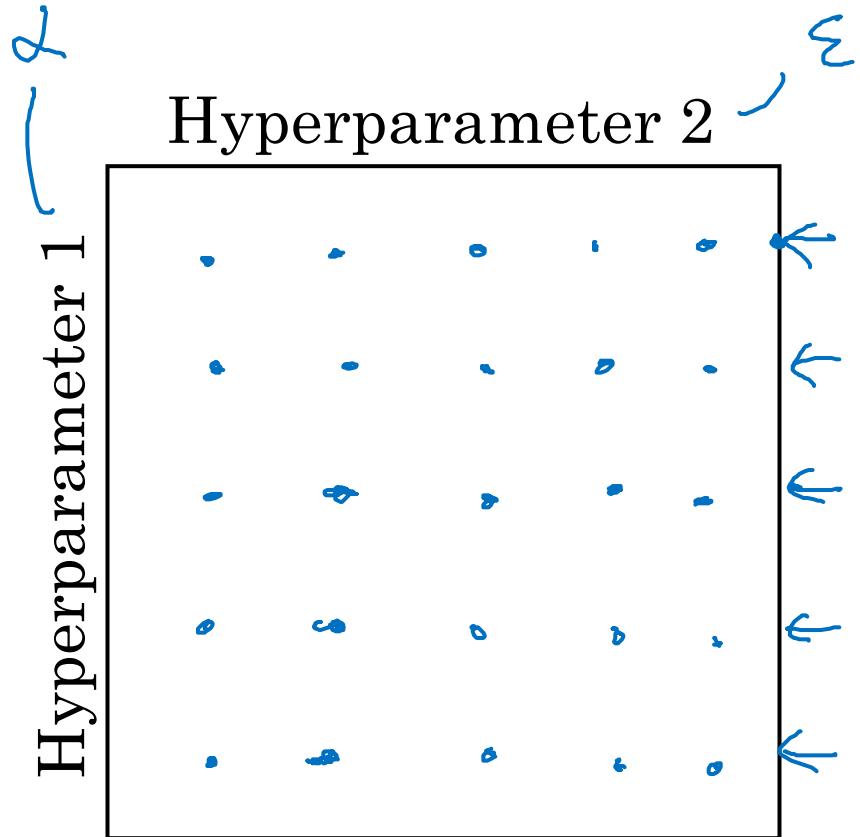
\*



\*\*



# Try random values: Don't use a grid

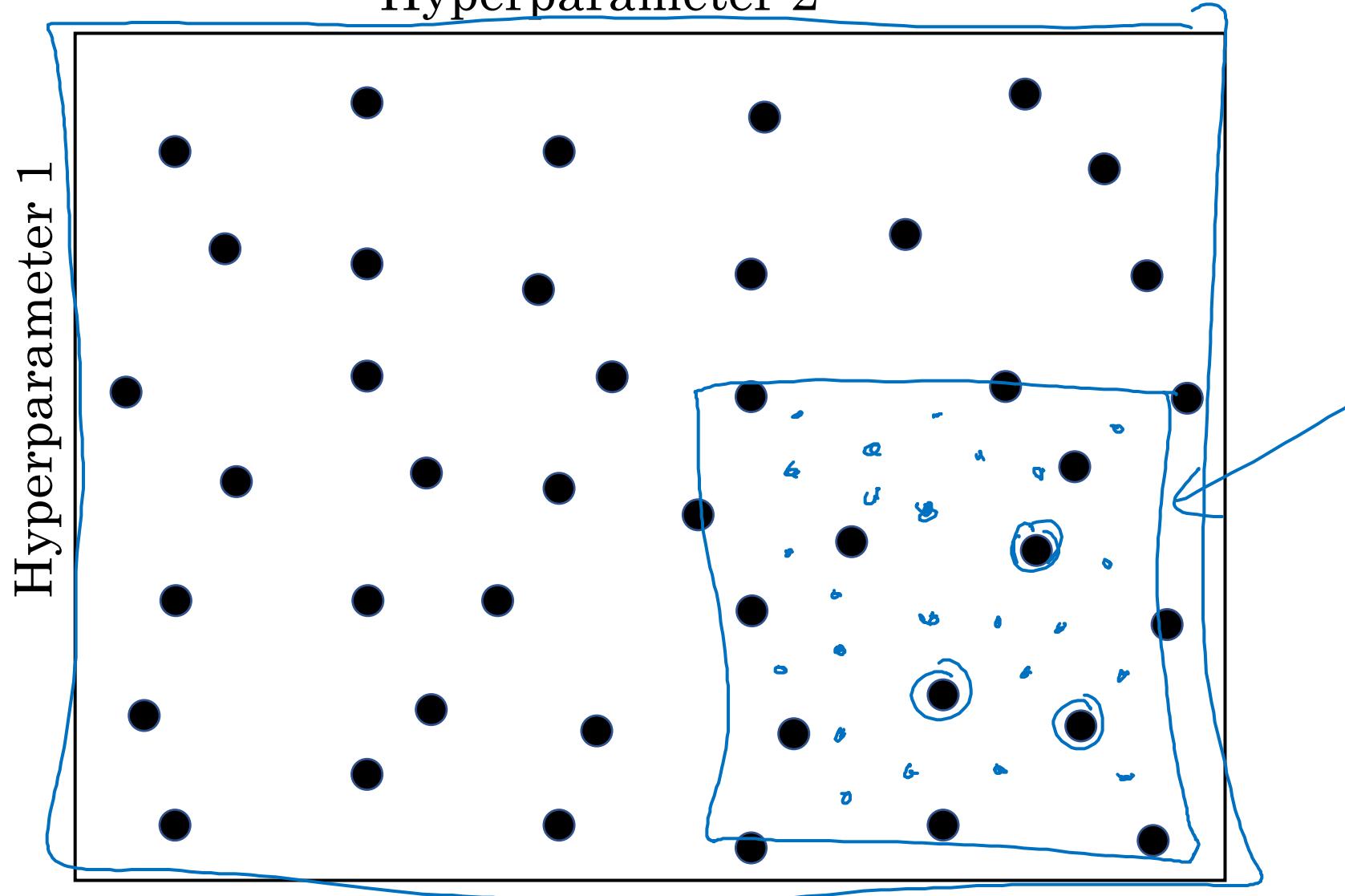


intuition about why random search is useful in practice:  
you tested 25 values of each variable and not only 5.

# Coarse to fine

→ Special search scheme where you progressively repeat the search sampling more densely in the areas of the search space that performed best in the previous search.

Hyperparameter 2





deeplearning.ai

# Hyperparameter tuning

---

Using an appropriate  
scale to pick  
hyperparameters

# Picking hyperparameters at random

$$\rightarrow n^{[l]} = 50, \dots, 100$$



$$\rightarrow \# \text{ layers} \quad L : 2 - 4$$

2, 3, 4

- Random search for hyp-opt should be performed at specific scales depending on the nature of the hyperparameter.
- discrete/integer hyperparameters should be sampled from a uniform distribution on a linear scale  
 $\Theta \sim U([\Theta_{\min}, \Theta_{\max}])$

- Continuous hyperparameters should be sampled from a uniform distribution on a logarithmic scale.

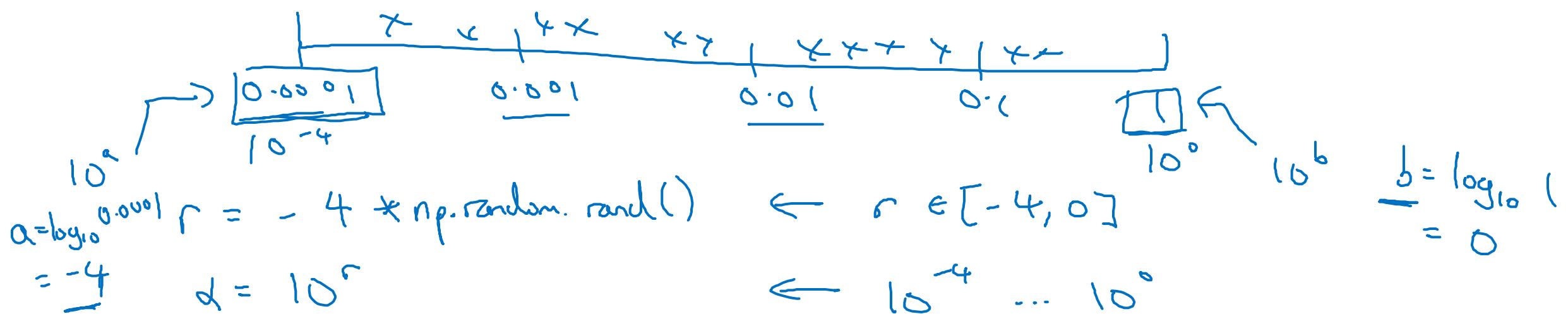
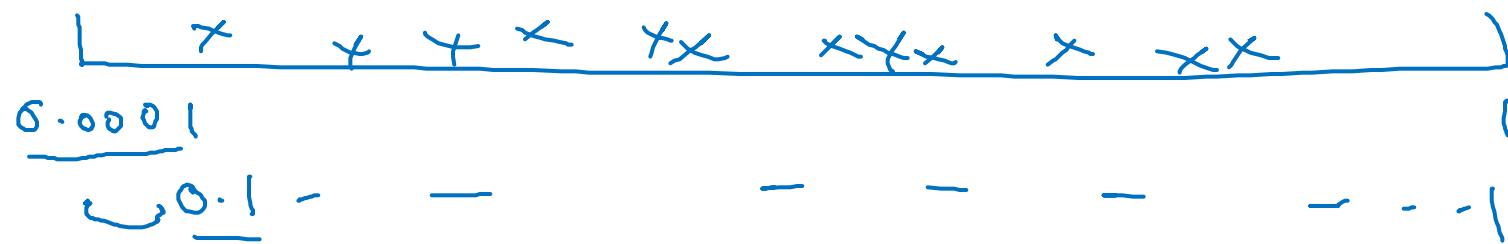
$$\Theta = 10^x \sim U([a, b])$$

This can be done by:

- seeing  $\Theta$  as  $10^x$  and sampling the exponent  $x$  uniformly.  
 $\Theta$  in  $[10^a, 10^b] \Rightarrow$  sample  $\Theta = 10^{x \sim U([a, b])}$
- sampling  $\Theta \sim \log U([a, b])$

# Appropriate scale for hyperparameters

$$\lambda = 0.0001, \dots, 1$$



$$\boxed{10^a \dots 10^b} \quad \boxed{\begin{matrix} r \in [a, b] \\ [-4, 0] \end{matrix}} \quad \boxed{\lambda = 10^r}$$

# Hyperparameters for exponentially weighted averages

② Example: if points in  $[0.9, 0.99]$ , write it as  $\log[10^{-2}, 10^{-1}]$  and then sample exponent in  $[-2, -1]$ .

$$\beta = 0.9 \dots 0.999$$

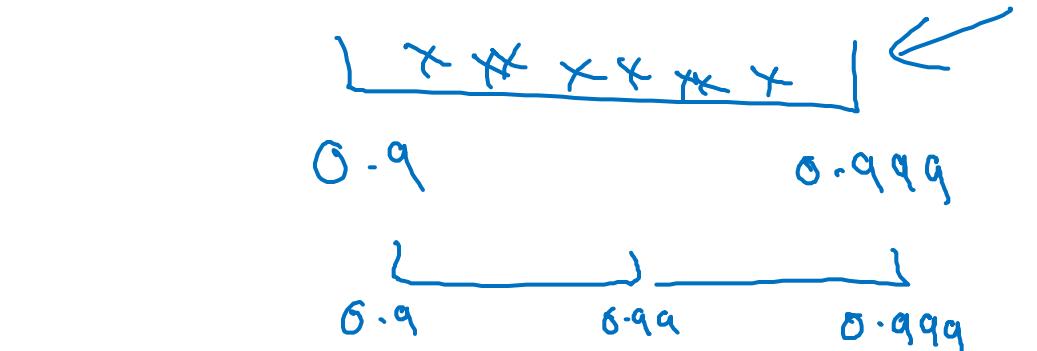
$\downarrow$                      $\downarrow$   
 $10$                      $1000$

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

$$\begin{aligned} \beta: 0.900 &\rightarrow 0.9005 \\ \beta: 0.999 &\rightarrow 0.9995 \end{aligned}$$

$\sim 1000$                      $\sim 2000$

① intuition: an example in which logarithmic sampling is very relevant is when sampling the momentum parameter  $\beta$ . As  $\beta$  moves toward 1, smaller intervals become much denser of relevant points. Varying  $\beta$  in  $[0.9, 0.9995]$  is much more sensitive than varying it in  $[0.9, 0.9005]$ .



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

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

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

$\Rightarrow$  Loguniform sampling is crucial

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

average over  $n$  samples



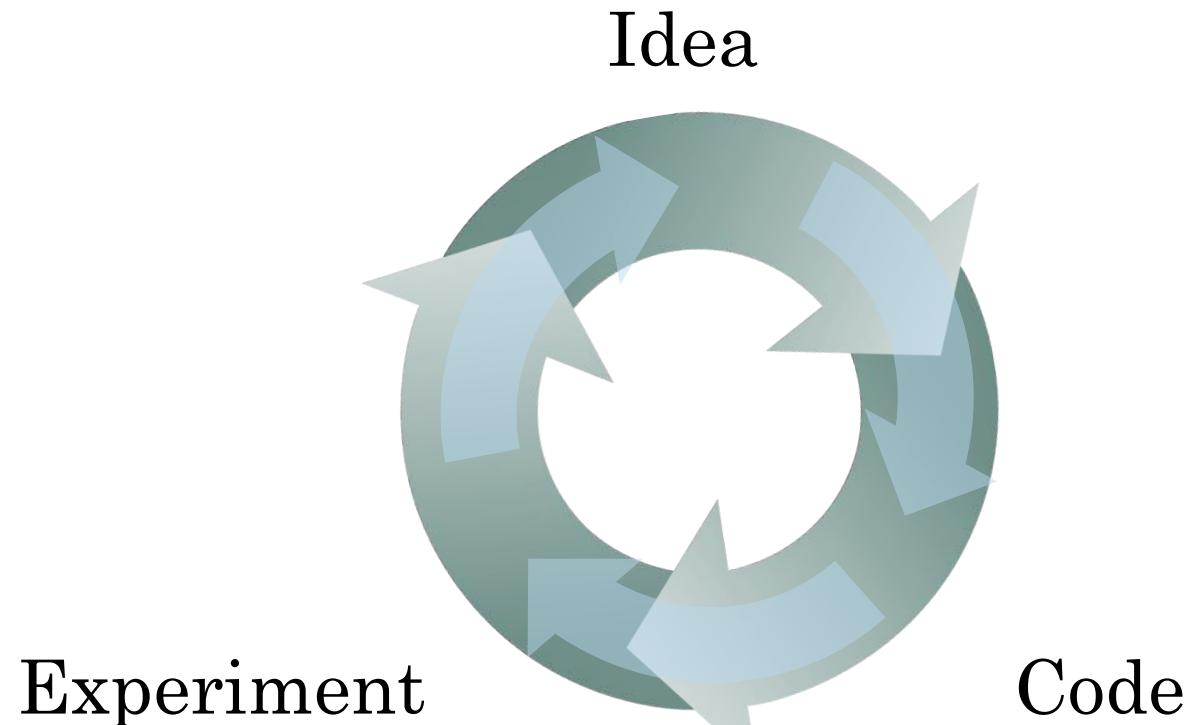
deeplearning.ai

# Hyperparameters tuning

---

## Hyperparameters tuning in practice: Pandas vs. Caviar

# Re-test hyperparameters occasionally

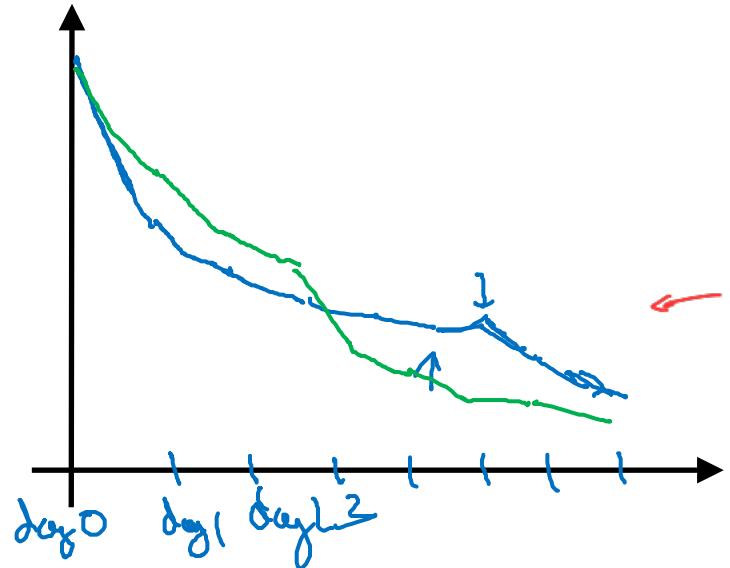


- NLP, Vision, Speech, Ads, logistics, ....
- Intuitions do get stale.  
Re-evaluate occasionally.



Real world data (in real applications) may be non-stationary, therefore it may be useful to periodically re-check the model's hyperparameters.

# Babysitting one model

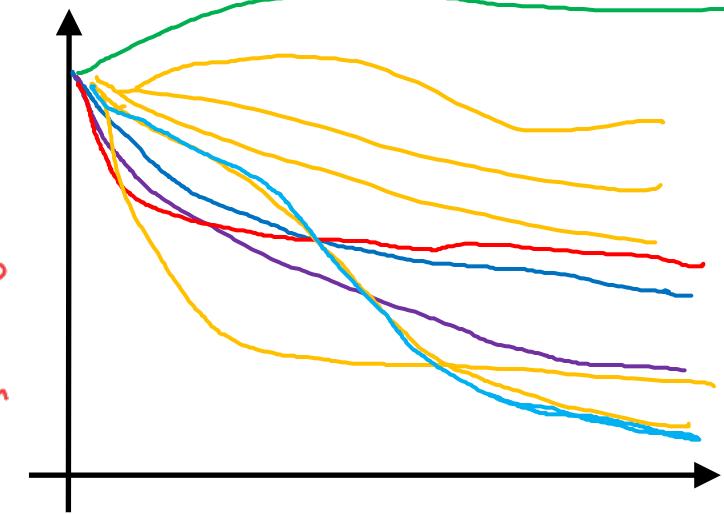


Panda ↵

How to test for different hyperparameters depends on the computational resources.

limited resources:  
test 1 set of hyperparameters (1 model),  
as it trains (maybe over days of training)  
keep adjusting the hyperparameters based  
on the convergence trend.

# Training many models in parallel



Lots of resources: →  
test many hyperparameters in parallel.



Caviar ↵



deeplearning.ai

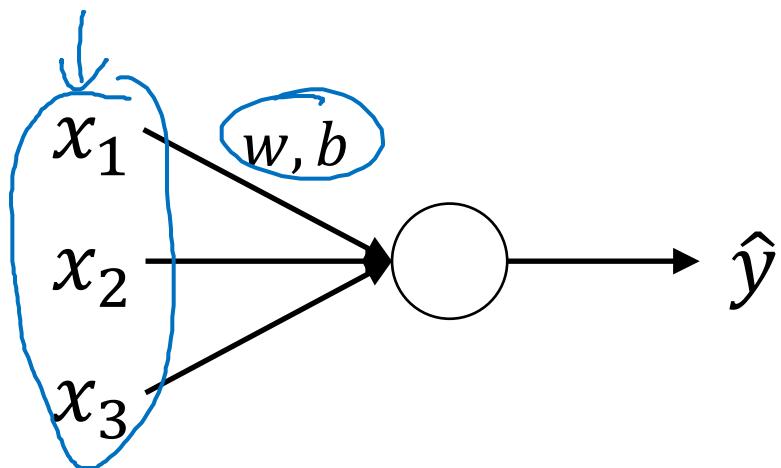
# Batch Normalization

---

## Normalizing activations in a network

Batch normalization simplifies the hyperparameter search  
and speeds up convergence of the network.

# Normalizing inputs to speed up learning



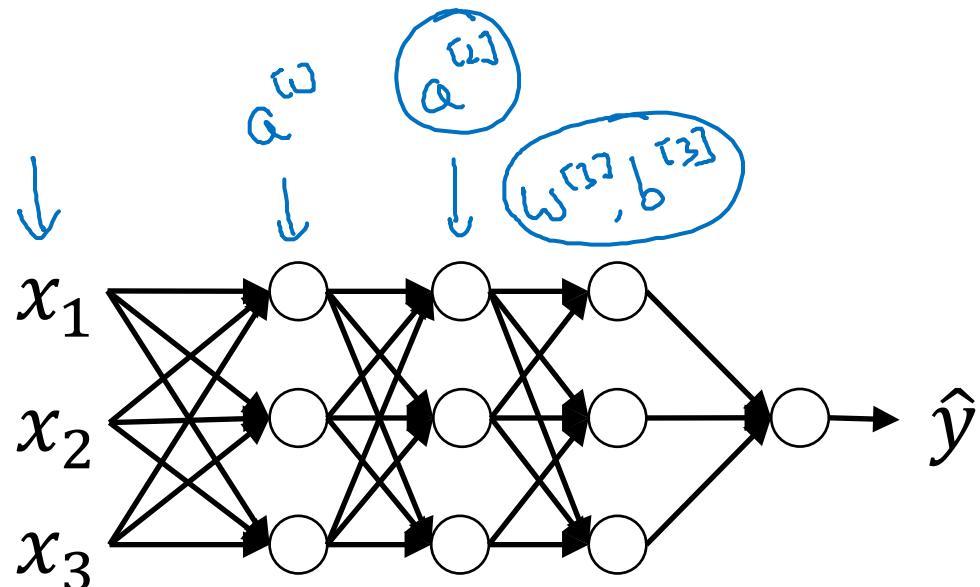
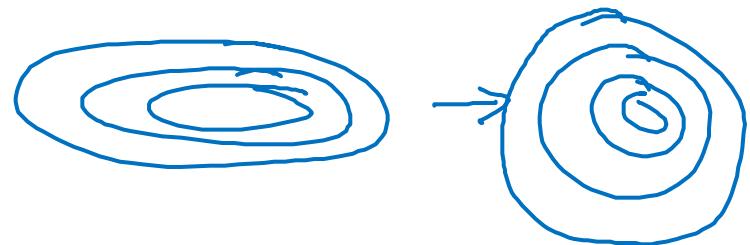
$$\mu = \frac{1}{m} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

$$Z = \frac{X}{\sigma}$$

Annotations:  $\mu$  is labeled "mean" and  $X - \mu$  is labeled "element-wise".  $\sigma^2$  is labeled "variance" and  $Z = \frac{X}{\sigma}$  is labeled "standard deviation".



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

Normalize  $\frac{z^{[2]}}{}$   
so to train neurons of layer  $n+1$  faster.

( $z^{[n]}$  → actually  $z^{[n]}$  is most often normalized.)

# Implementing Batch Norm

Given some intermediate values in NN

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

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

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

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

modifies variance

modifies mean

Use

$\gamma$  initial  $\beta$

- NOTES:
  - $\gamma, \beta$  are params and not hyperparams ( $\Rightarrow$  they are learned)
  - $\gamma^{(l)}, \beta^{(l)}$  defined for each layer  $l$  and hidden unit  $i$ .
  - The param  $\beta^{(l)}$  has nothing to do with the momentum hyperparameter

IDEA of BATCHNORM: normalizing the output of layer  $n$  ( $z^{(n)}$ , or more commonly  $\tilde{z}^{(n)}$ ) so to train the parameters of layer  $n+1$  faster.

$z^{(1)}, \dots, z^{(n)}$  of a certain layer  $l$

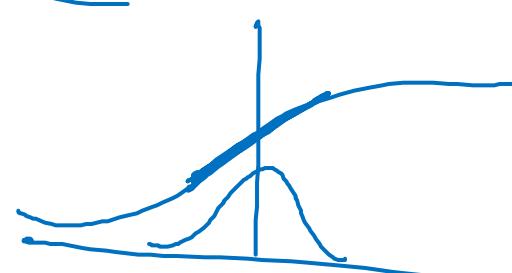
If  $\gamma = \sqrt{\sigma^2 + \epsilon}$

$\beta = \mu$

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

$x$

$\tilde{z}^{(i)}$



learnable parameters of model.

- For layer  $l$  and all the samples in the minibatch, compute the mean and variance of the activations  $z^{(l)}(i)$  of the units in that layer. Then normalize the activation  $z^{(l)}(i)$  of each unit  $i$  using  $\mu$  and  $\sigma$ .
- Differently from the normalization of the inputs (input layer) we may want  $z^{(l)}(i)$  to have  $\mu$  and  $\sigma$  different from 0 and 1 respectively. For this reason the params  $\gamma, \beta$  are introduced and tuned.  $\tilde{z}^{(l)}(i) = \gamma z^{(l)}(i) + \beta^{(l)}$  modifies the variance of  $\gamma$  and  $\beta$  the mean.

Andrew Ng



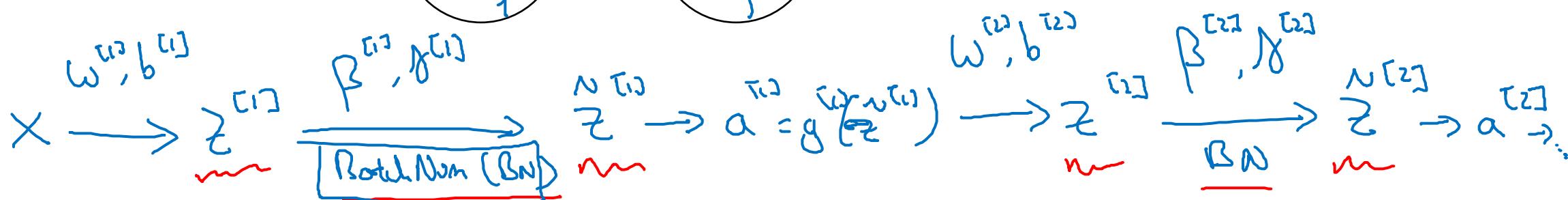
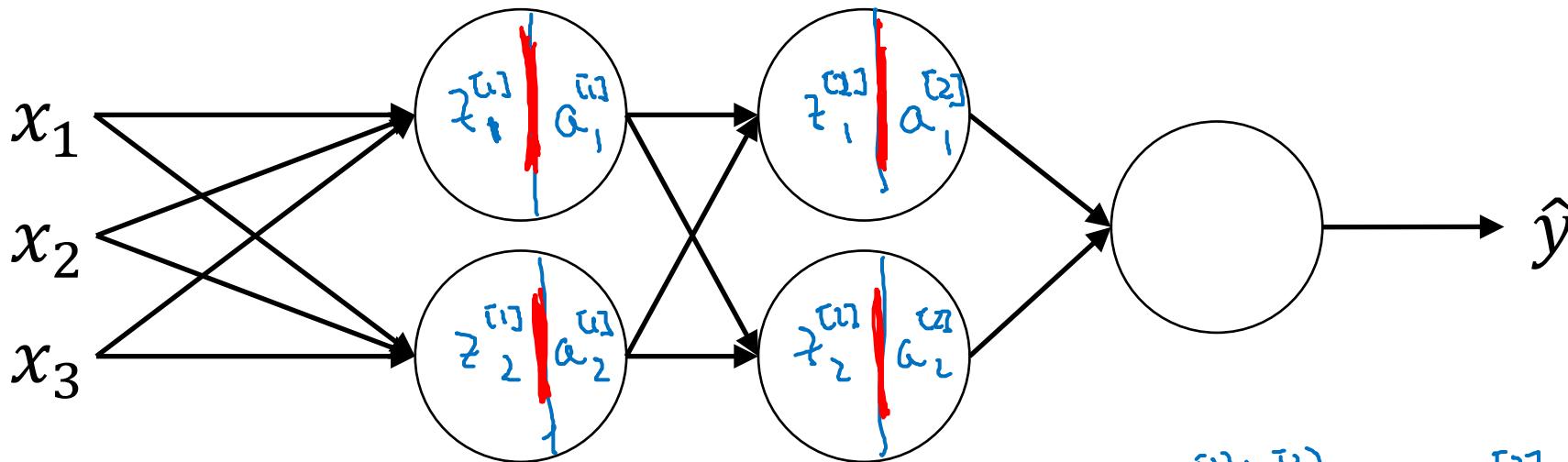
deeplearning.ai

# Batch Normalization

---

Fitting Batch Norm  
into a neural network

# Adding Batch Norm to a network



Parameters:

$$W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]},$$

$$\beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]}$$

$$\frac{d\beta^{[L]}}{d\beta} = \frac{\partial \beta^{[L]}}{\partial \beta} = 1$$

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

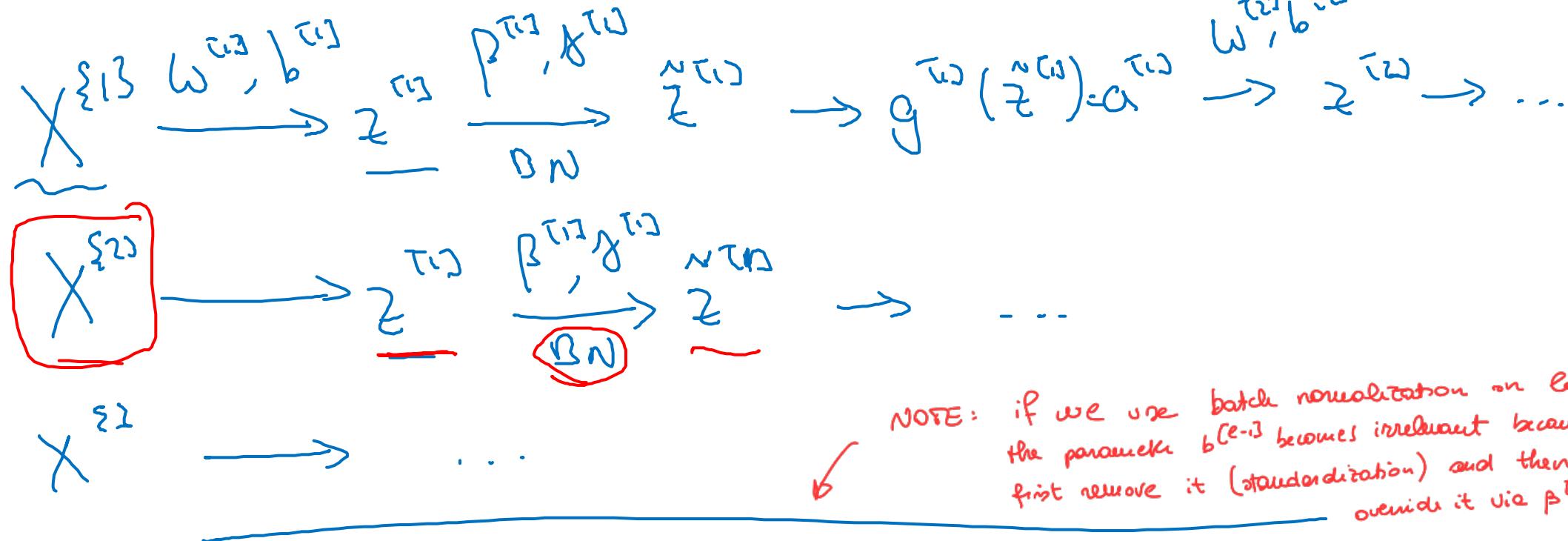
`tf.nn.batch_normalization` ←

-  $\beta$  and  $\gamma$  are learned with backpropagation (like  $W$  and  $b$ ).



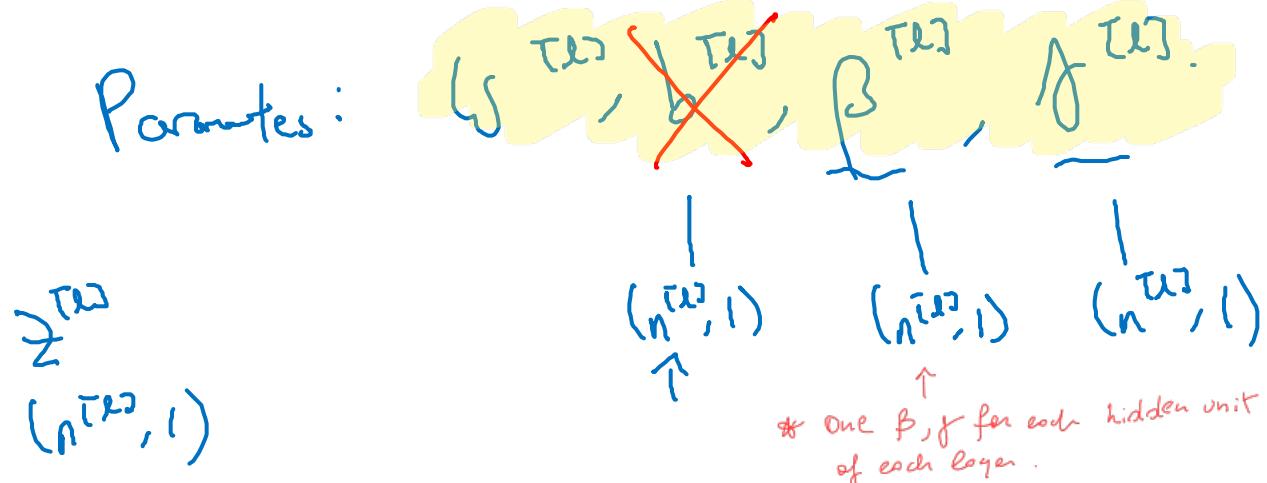
# Working with mini-batches

- For minibatch GD, the batch norm is repeated for each minibatch (i.e. using only the data from that minibatch).



NOTE: if we use batch normalization on layer  $l$ , then the parameter  $b^{[l-1]}$  becomes irrelevant because the BN $^{[l]}$  will first remove it (standardization) and then potentially override it via  $\beta^{[l]}$ .

Parameters:



$$\begin{aligned} \underline{z}^{[l]} &= W^{[l]} \underline{a}^{[l-1]} + \cancel{b^{[l]}} \\ z^{[l]} &= W^{[l]} a^{[l-1]} \\ \underline{z}^{[l]}_{\text{norm}} &= \gamma^{[l]} \underline{z}^{[l]}_{\text{norm}} + \beta^{[l]} \end{aligned}$$

Andrew Ng

# Implementing gradient descent with Batch Norm

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  $\underline{dw}^{[l]}$ ,  ~~$\underline{db}^{[l]}$~~ ,  $\underline{d\beta}^{[l]}$ ,  $\underline{dg}^{[l]}$

Update parameters

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

Works w/ momentum, RMSprop, Adam.

Same as without BN. Works also with other optimization algorithms such as Adam.

Motivation 1: (As said before) because it standardizes the inputs of intermediate layers therefore making GD faster at learning those layers' parameters  
(similar to what happens for the input normalization).

M2: because it makes the weights of later layers more robust to changes to weights in earlier layers. It does so by mitigating the distribution shift of the inputs  $\mathbf{x}^{l-1}$  to the layer  $l$  due to changes in the weights of earlier layers ( $1, \dots, l-1$ ). Reducing the distribution shift "decouples" earlier and later layers, that is, makes later layers a bit more independent from changes of the previous ones therefore simplifying/speeding up learning.



deeplearning.ai

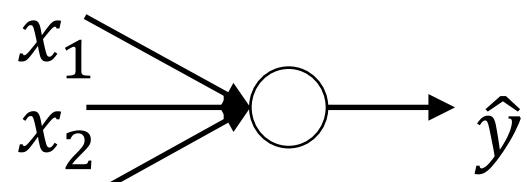
# Batch Normalization

---

## Why does Batch Norm work?

NOTE: Batch normalization has a slight normalization side-effect.

# Learning on shifting input distribution



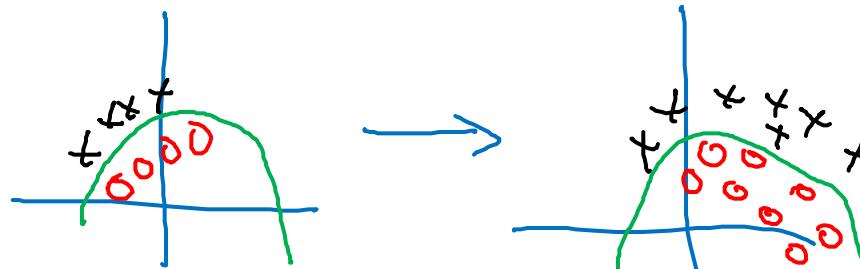
Cat



Non-Cat



$$y = 1$$



$$y = 1$$



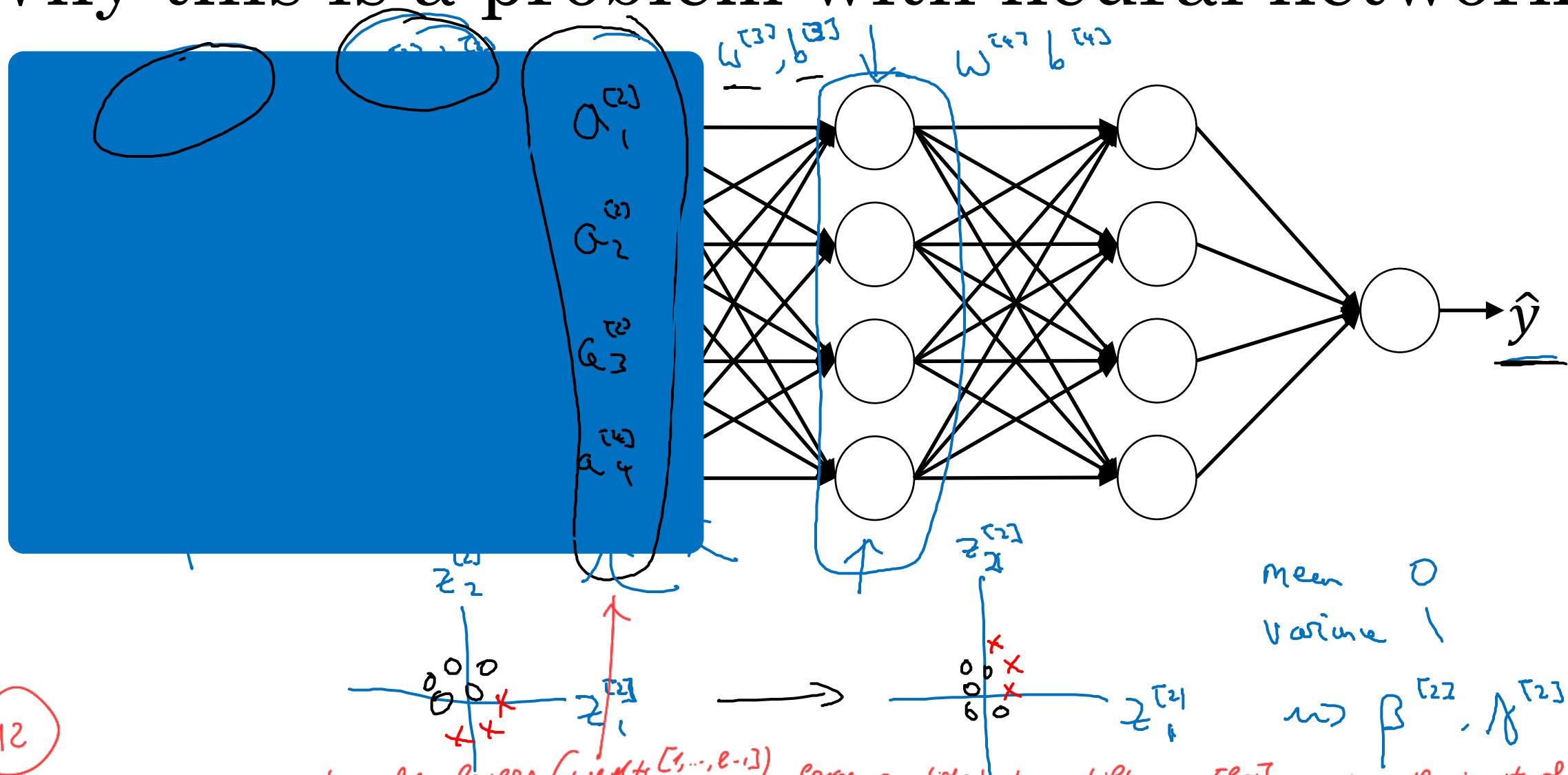
$$y = 0$$



"Covariate shift"

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

# Why this is a problem with neural networks?



- M2
- changes in earlier layers (weights  $[1, \dots, l-1]$ ) cause a distribution shift in  $a^{(l-1)}$ , which is the input of layer  $l$ .
  - BN applied to layer  $N$  reduces the distribution shift by ensuring a constant  $\mu = 0$  of the input  $a$ .
  - This makes learning a bit easier and also decouples the layers (forces the network to learn layers independently).

Andrew Ng

# Batch Norm as regularization

Secondary / side - effect  
of batch norm.

- Each mini-batch is scaled by the mean/variance computed on just that mini-batch.  
 $\hat{z}^{[l]}$        $\underline{54}, \underline{128}$        $\underline{z}^{[l]}$
- 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.  
 $\mu, \sigma^2$
- This has a slight regularization effect.

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

- BN in combination with minibatch learning has a slight regularization effect because layer activations are scaled by quantities computed on a portion of the data and are, therefore, not exactly right. This has the effect of adding some noise to each layer, therefore forcing the NN to "spread" the weights. Similar to what dropout does.  
↳ => regularization effect.    => SMALLER MINIBATCH = more noise = more regularization.

Andrew Ng

## BATCH NORM AT TEST TIME (missing slide)

Problem: at test time, examples are potentially fed into the NN one at a time and is therefore impossible to estimate  $\mu^{[l]}, \sigma^{[l]}$  to standardize the activation of the layer  $l$ .

$\Rightarrow$  At training time: compute an exponentially weighted average of  $\mu^{[l]}, \sigma^{[l]}$  over consecutive minibatches (and keep learning  $\beta, \gamma$ ).

- then, at test time, use that  $\mu^{[l]}, \sigma^{[l]}$  to perform layer standardization.

### Batch Norm at test time

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

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

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

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

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

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

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

Andrew Ng



deeplearning.ai

Multi-class  
classification

---

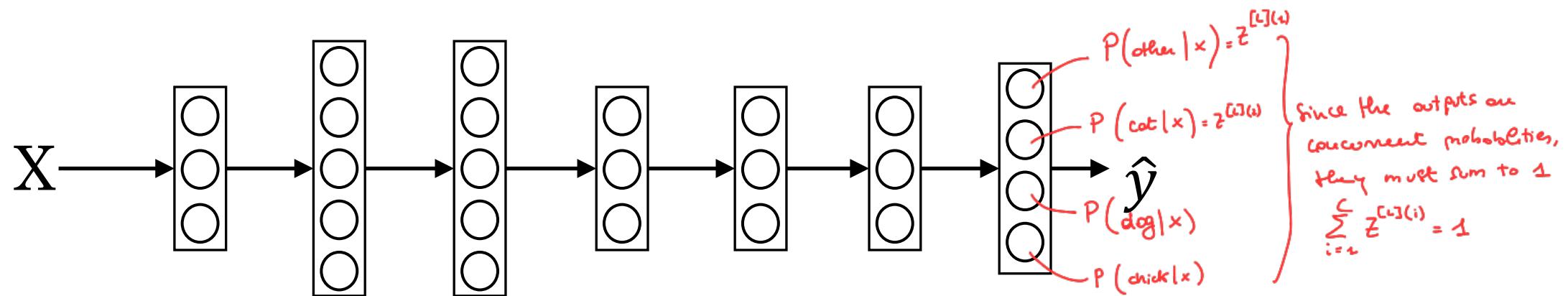
Softmax regression

# Recognizing cats, dogs, and baby chicks

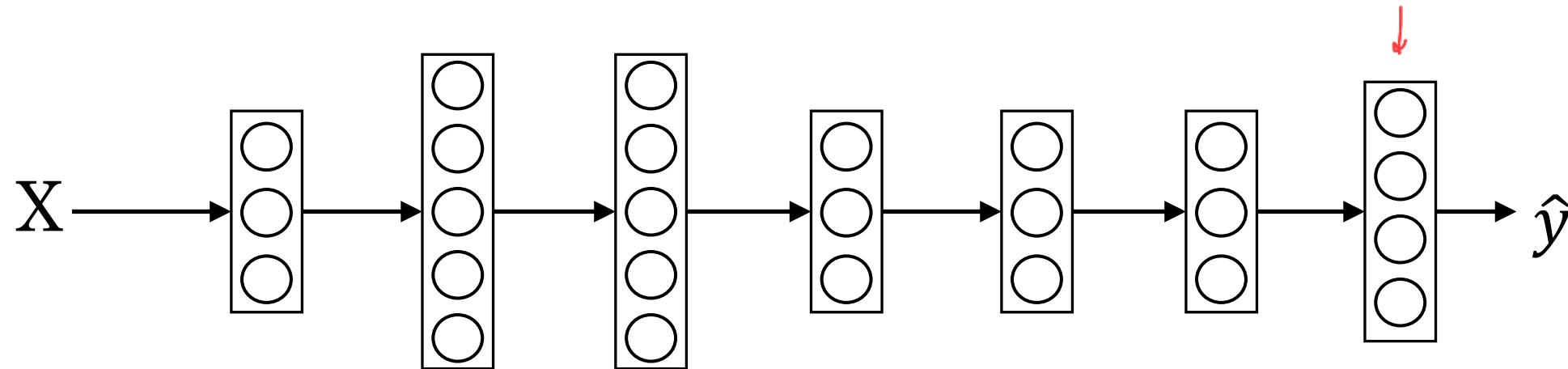


3            1            2            0            3            2            0            1

$C = \# \text{ classes}$  (in this case  $C=4$ )  $\rightarrow$  the  $n$ -units in the output layer must be  $= C$



# Softmax layer

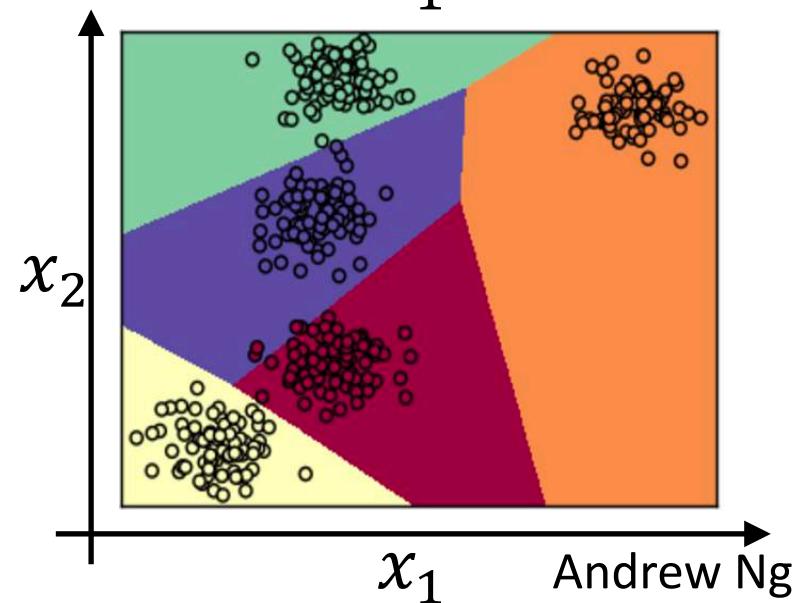
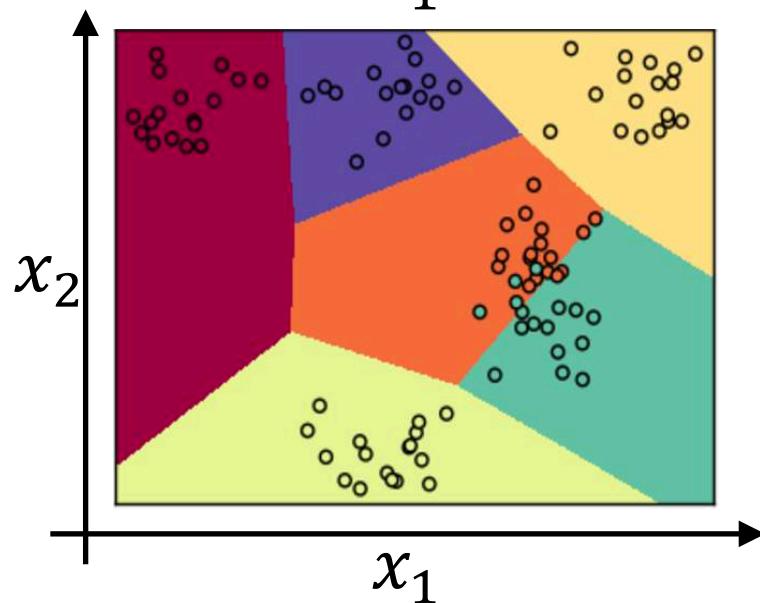
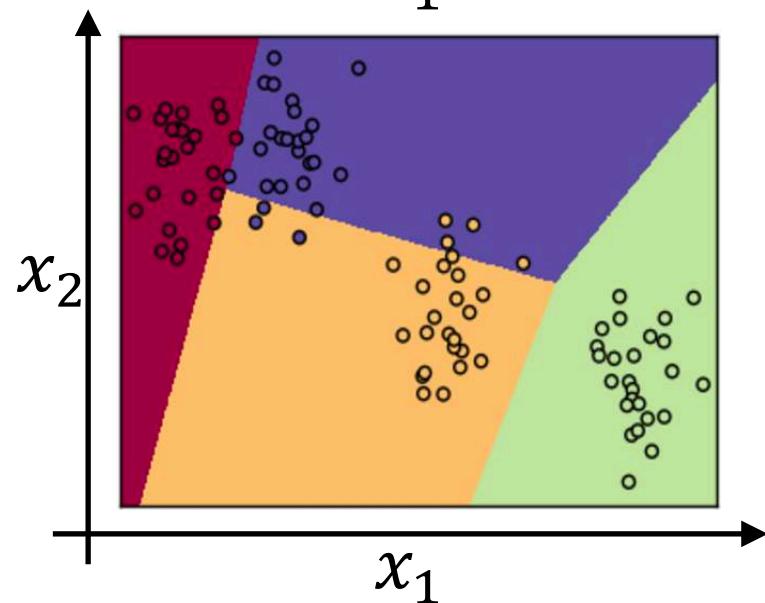
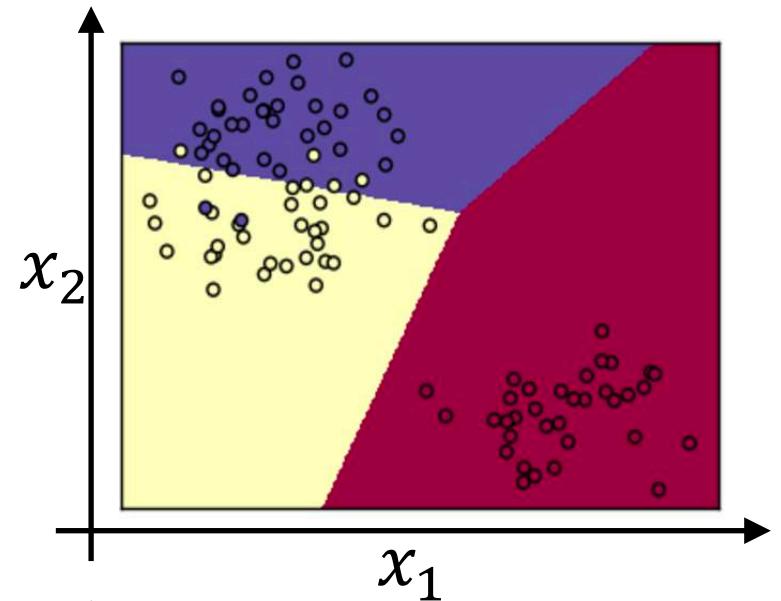
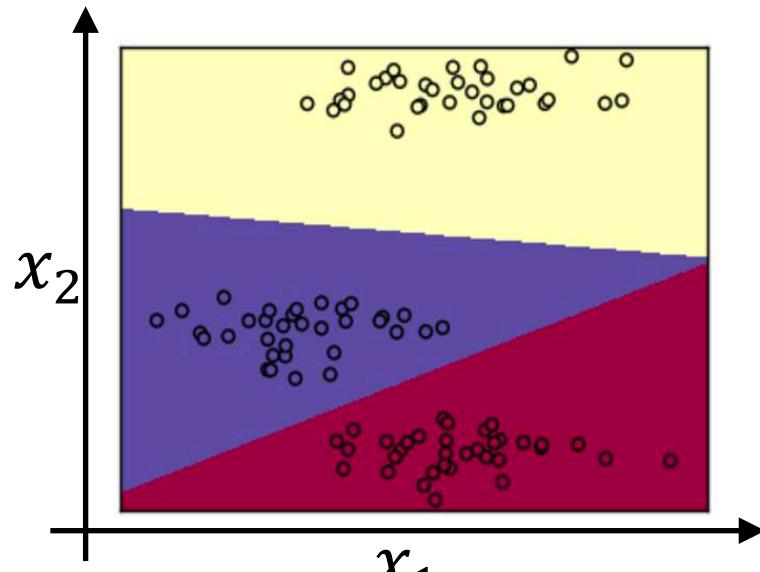
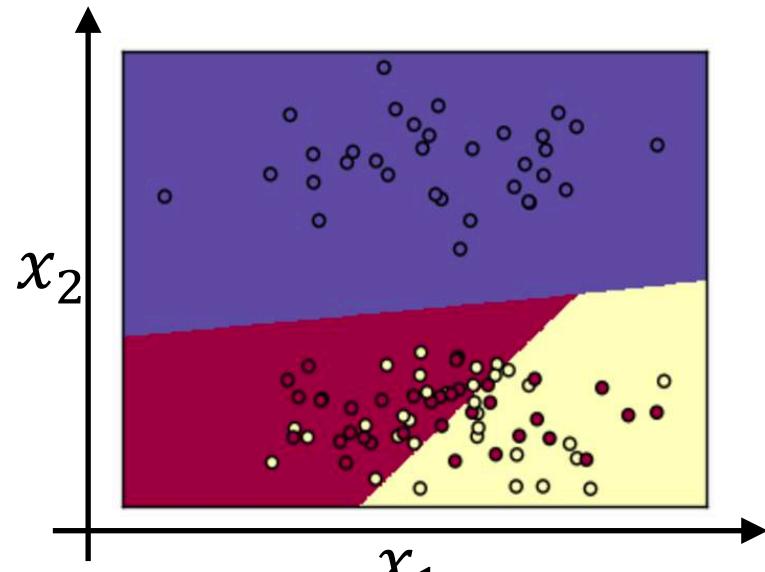


- to make the outputs of the last layer be probabilities that sum up to one, the last layer is made implement a softmax function. The activation of each hidden unit becomes

$$a^{[L](i)} = \frac{e^{z^{[L]}(i)}}{\sum_{j=1}^n e^{z^{[L]}(j)}}$$

- Softmax regression generalizes logistic regression to  $C > 2$  classes. For  $C=2$  the two activations correspond.

# Softmax examples



## Loss function for softmax regression

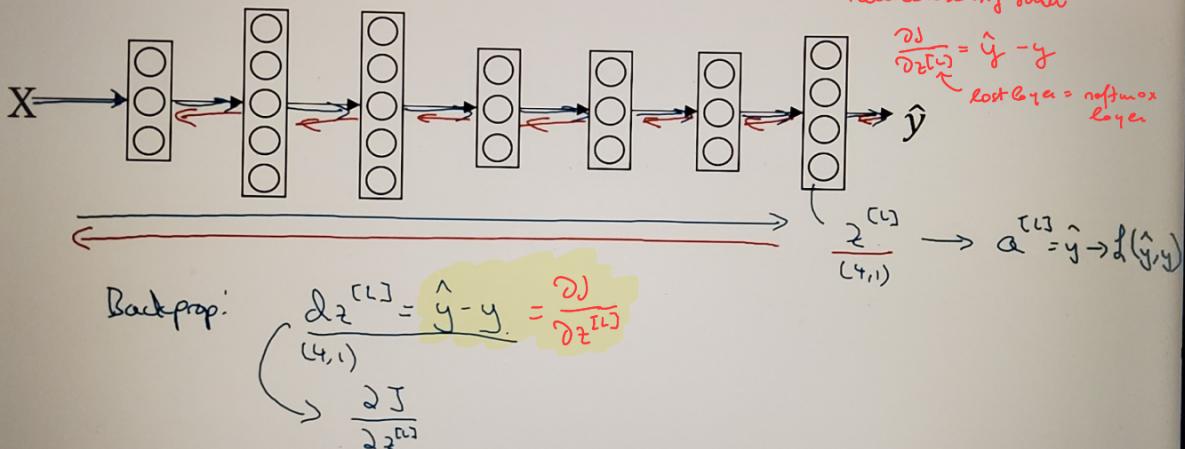
### Loss function

$$\begin{aligned}
 & \text{Diagram showing a 4-class softmax output vector } \hat{y} = \begin{bmatrix} 0.1 \\ 0.6 \\ 0.1 \\ 0.2 \end{bmatrix} \text{ where } \hat{y}_1 = \hat{y}_3 = \hat{y}_4 = 0 \\
 & \text{and } \hat{y}_2 = 1. \\
 & \text{Equation: } L(\hat{y}, y) = -\sum_{j=1}^C y_j \log \hat{y}_j \\
 & \quad \text{with a note: "small" under the summation sign.} \\
 & \text{Equation: } J(w^{(i)}, b^{(i)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\
 & \quad \text{minimizing this loss requires maximizing the probability } \hat{y}_j \text{ of the true class } y_i. \\
 & \quad \text{Note: } -y_2 \log \hat{y}_2 = -\log \hat{y}_2.
 \end{aligned}$$

$$\begin{aligned}
 Y &= [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \\
 &= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \dots \\
 &\quad (4, m)
 \end{aligned}
 \qquad
 \begin{aligned}
 \hat{Y} &= [\hat{y}^{(1)} \ \dots \ \hat{y}^{(m)}] \\
 &= \begin{bmatrix} 0.1 \\ 0.6 \\ 0.1 \\ 0.2 \end{bmatrix} \dots \\
 &\quad (4, m)
 \end{aligned}$$

Andrew Ng

### Gradient descent with softmax



Andrew Ng



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)

(↳ how likely is it that the company writing the software keeps it open over the years.



deeplearning.ai

# Programming Frameworks

---

## TensorFlow

- Cool feature: you only have to implement the forward step and it will compute the backward automatically.

# Motivating problem

$$\xrightarrow{w} ((w-5)^2) \rightarrow \hat{y}$$
$$J(w) = \boxed{w^2 - 10w + 25}$$

(cost)

$$(w-5)^2$$
$$\left. \begin{array}{c} J(w, b) \\ \uparrow \uparrow \end{array} \right\}$$

$\uparrow$   $w = 5$

The boss only wants to minimize the parabola. No error wrt inputs is considered. We know that the optimum will be  $w=5$ . Let's find it using TensorFlow.

*solution:*

```
In [1]: import numpy as np
import tensorflow as tf

In [2]: w = tf.Variable(0, dtype=tf.float32)
optimizer = tf.keras.optimizers.Adam(0.1)

def train_step():
    with tf.GradientTape() as tape:
        cost = w ** 2 - 10 * w + 25
    trainable_variables = [w]
    grads = tape.gradient(cost, trainable_variables)
    optimizer.apply_gradients(zip(grads, trainable_variables))

print(w)

<tf.Variable 'Variable:0' shape=() dtype=float32, numpy=0.0>

In [3]: train_step()
print(w)

<tf.Variable 'Variable:0' shape=() dtype=float32, numpy=0.09999997>

In [4]: for i in range(1000):
    train_step()
print(w)

<tf.Variable 'Variable:0' shape=() dtype=float32, numpy=5.000001>

In [ ]: |
```



*Other problem: the cost function includes the training data*

$$J(w, b) = w_1^2 x_1 + w_2 x_2 + x_3$$

```
In [7]: w = tf.Variable(0, dtype=tf.float32)
x = np.array([1.0, -10.0, 25.0], dtype=np.float32)
optimizer = tf.keras.optimizers.Adam(0.1)

def training(x, w, optimizer):
    def cost_fn():
        return x[0] * w ** 2 + x[1] * w + x[2]
    for i in range(1000):
        optimizer.minimize(cost_fn, [w])
    return w

w = training(x, w, optimizer)
print(w)

<tf.Variable 'Variable:0' shape=() dtype=float32, numpy=5.000001>
```

# 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()
```

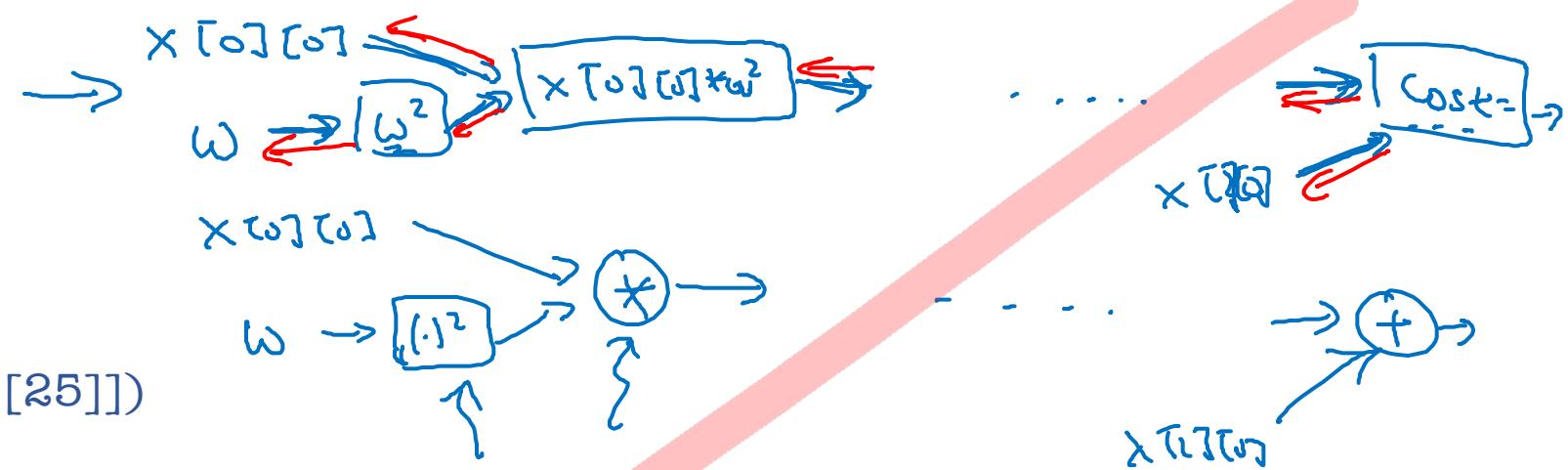
```
session.run(init)
```

```
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))
```