



deeplearning.ai

Hyperparameter
tuning

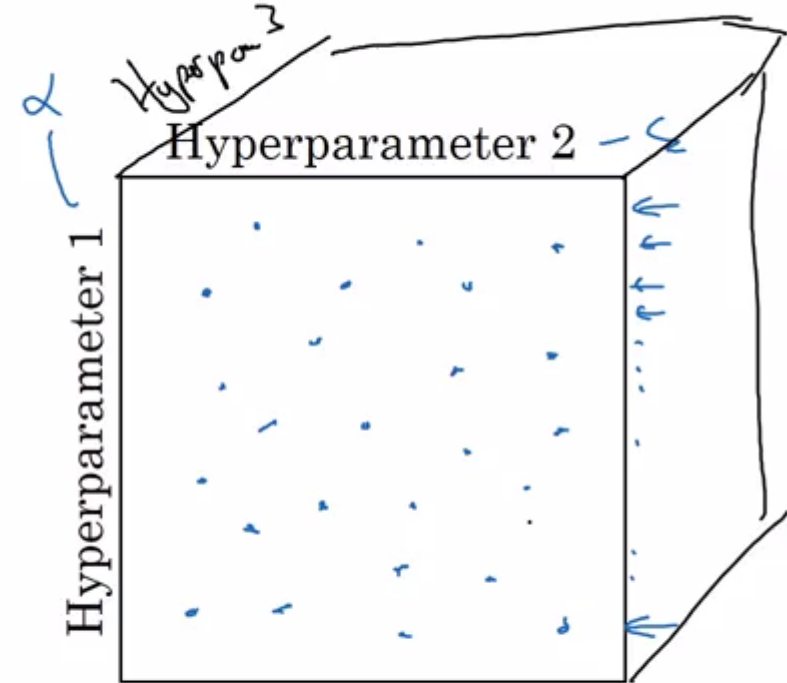
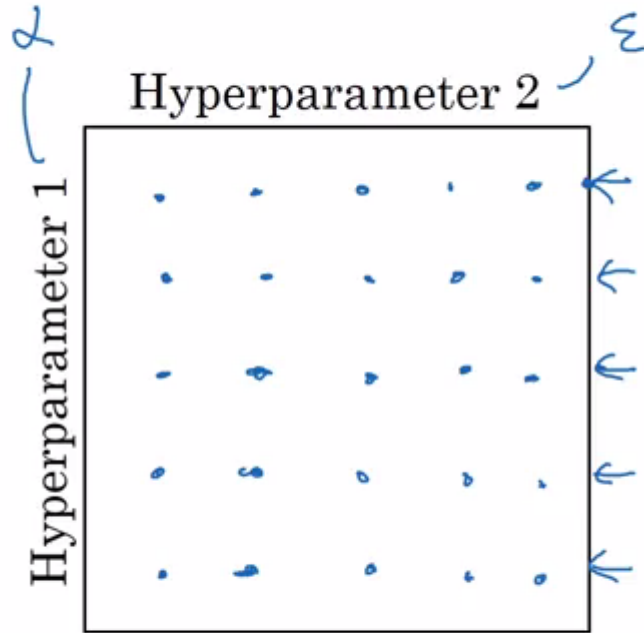
Tuning process



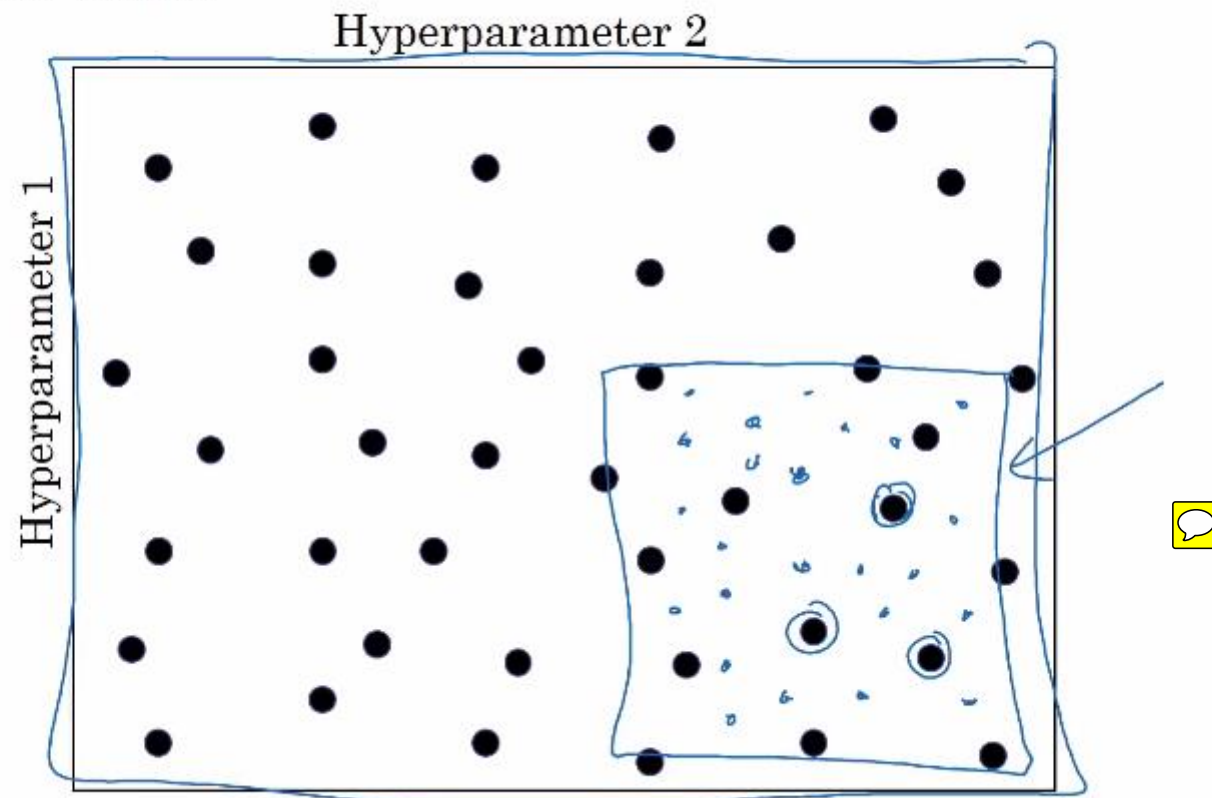
Hyperparameters

→ α
 β 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



Andrew Ng



deeplearning.ai

Hyperparameter tuning

Using an appropriate
scale to pick
hyperparameters



Picking hyperparameters at random

→ $n^{\text{test}} = 50, \dots, 100$

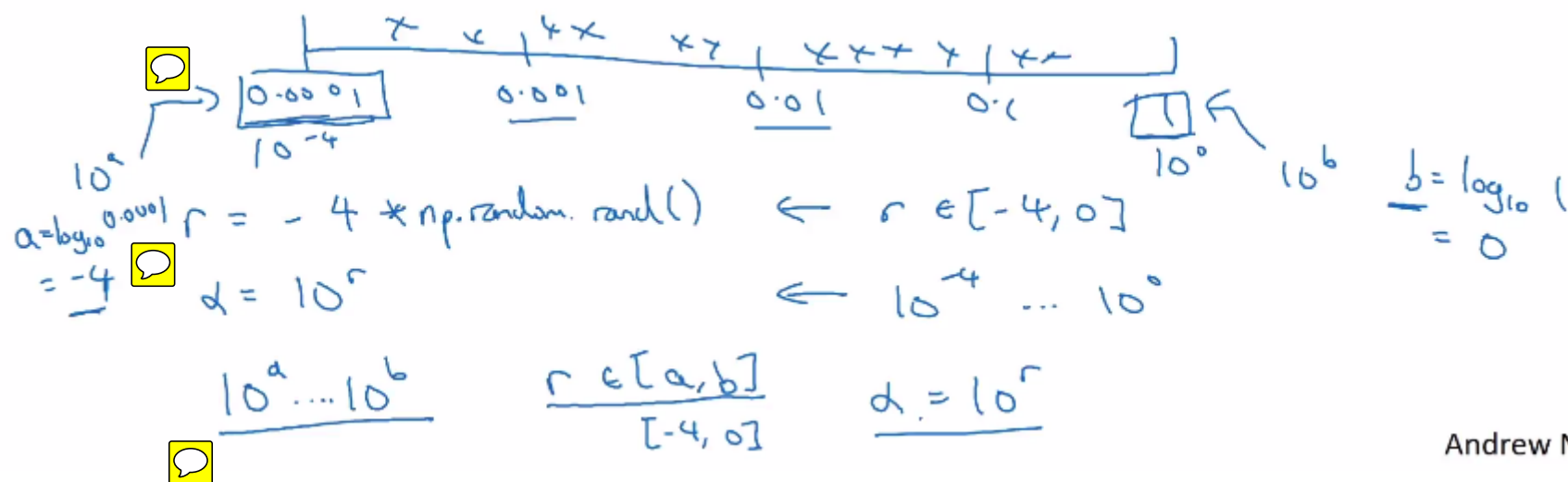
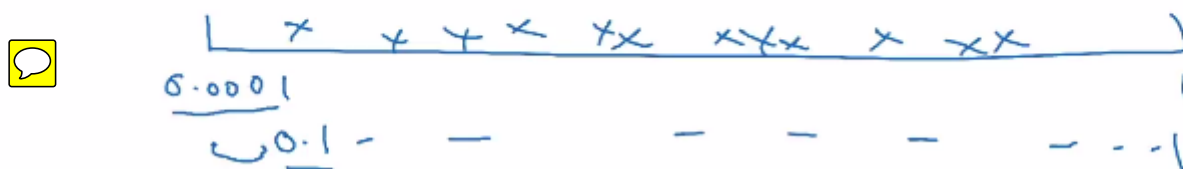


→ #layers L : 2 - 4

2, 3, 4

Appropriate scale for hyperparameters

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



Hyperparameters for exponentially weighted averages

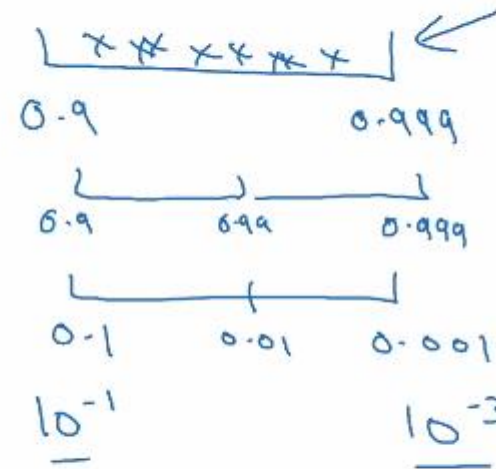
$\beta = 0.9 \dots 0.999$
 $\downarrow \qquad \qquad \downarrow$
 $10 \qquad \qquad 1000$

$1 - \beta = 0.1 \dots 0.001$

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

$\beta: 0.999 \rightarrow 0.9995$
 $\sim 1000 \qquad \sim 2000$

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



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

$1 - \beta = 10^r$

$\beta = 1 - 10^r$



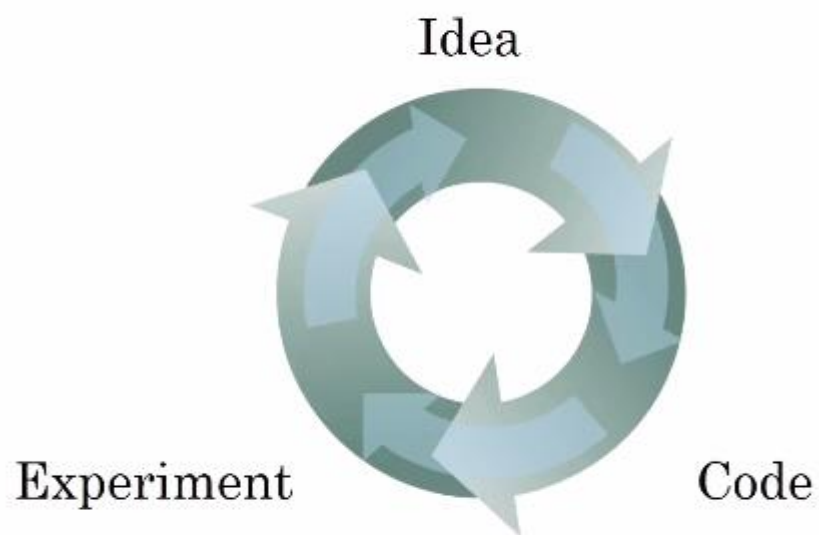
deeplearning.ai

Hyperparameters tuning

Hyperparameters tuning in practice: Pandas vs. Caviar



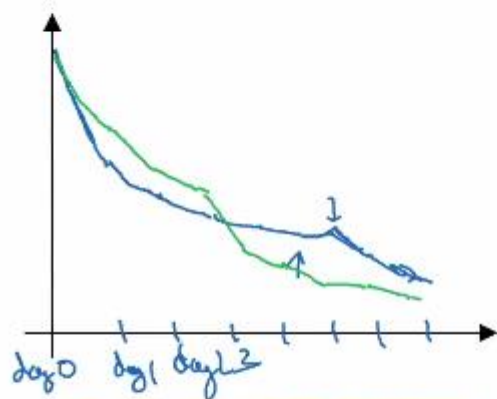
Re-test hyperparameters occasionally



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



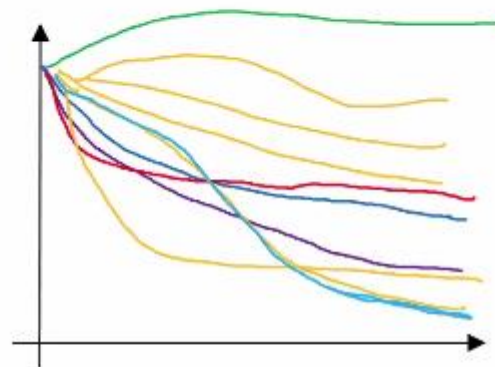
Babysitting one model



Panda ←



Training many models in parallel



Caviar ←

Andrew Ng



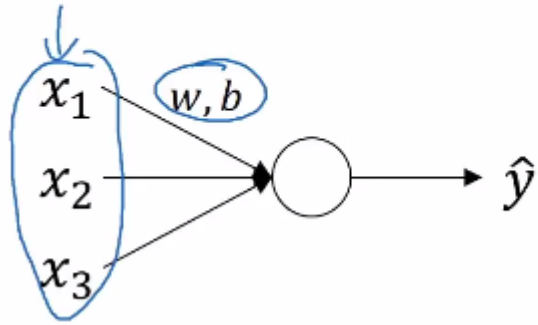
deeplearning.ai

Batch Normalization

Normalizing activations
in a network



Normalizing inputs to speed up learning

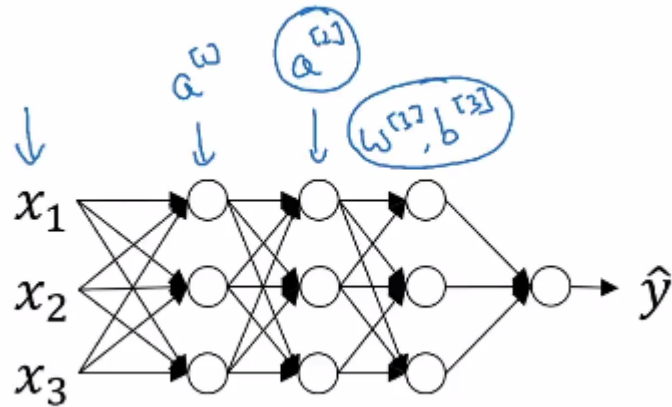
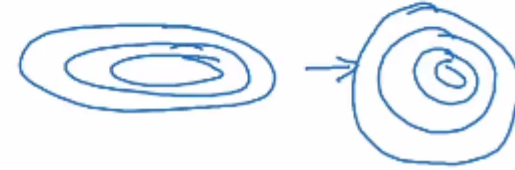


$$\mu = \frac{1}{n} \sum_i x^{(i)}$$

$$X = X - \mu$$

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

$$X = X / \sigma$$



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

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

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



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

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

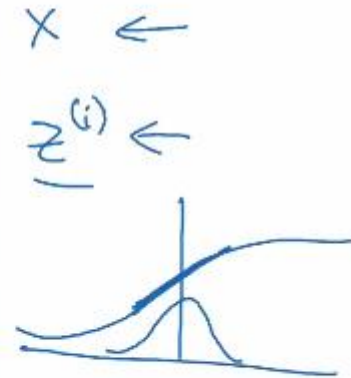
If

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

$$\beta = \mu$$

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

learnable parameters of model.





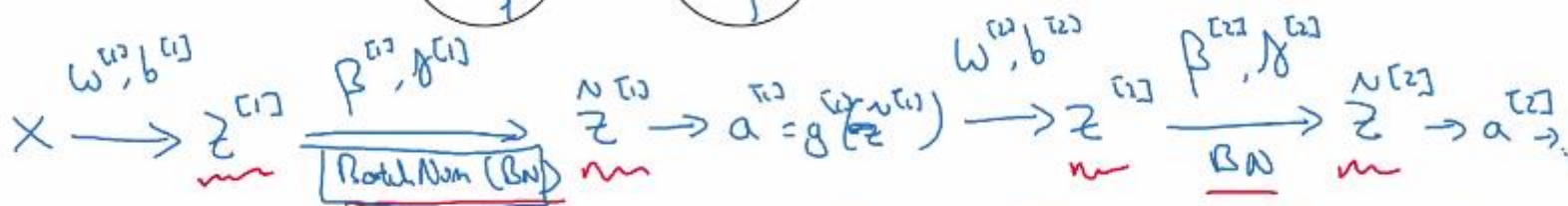
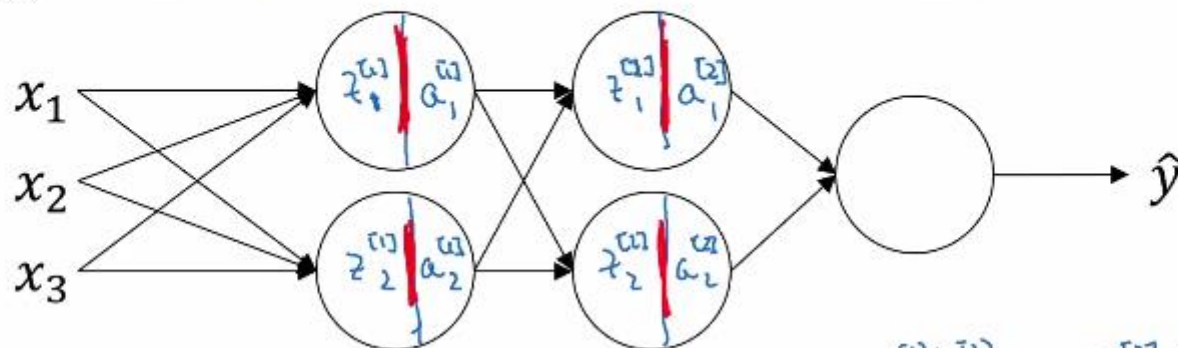
deeplearning.ai

Batch Normalization

Fitting Batch Norm into a neural network



Adding Batch Norm to a network

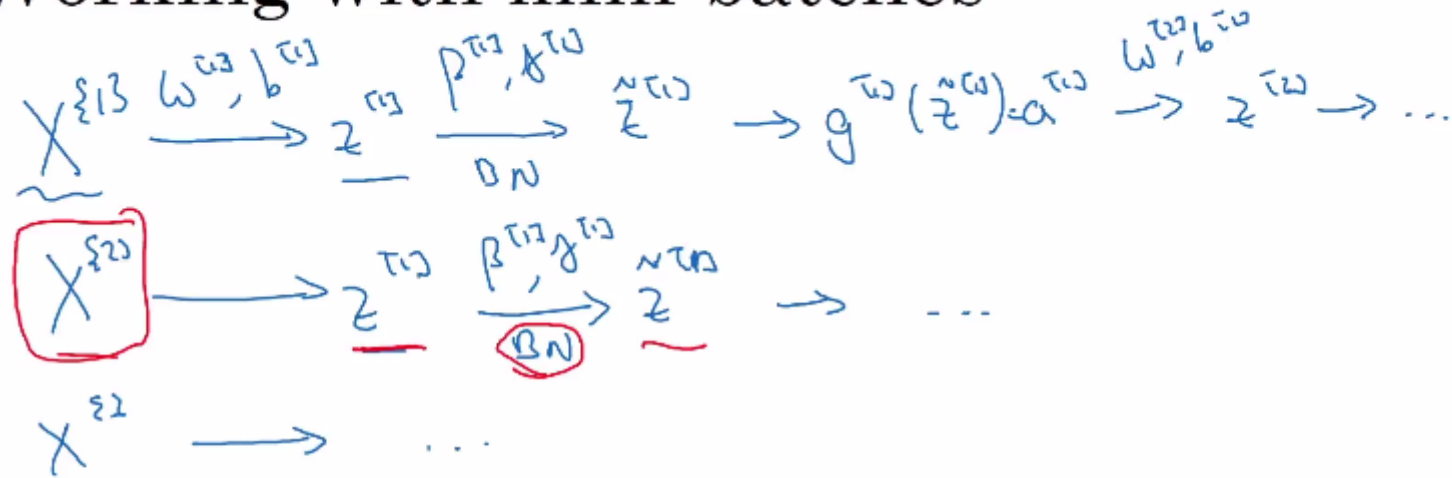


Parameters: $\{w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}\}$
 $\rightarrow \underline{\beta}^{[1]}, \underline{\gamma}^{[1]}, \underline{\beta}^{[2]}, \underline{\gamma}^{[2]}, \dots, \underline{\beta}^{[L]}, \underline{\gamma}^{[L]}$
 $\rightarrow \underline{\beta}$

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

tf.nn.batch-normalization ←

Working with mini-batches



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

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

$\rightarrow \underline{z}^{[1]} = W^{[1]} a^{[1-1]} + \cancel{b^{[1]}}$
 $\tilde{z}^{[1]} = W^{[1]} a^{[1-1]}$
 $\tilde{z}^{[1]}_{\text{norm}} = \gamma^{[1]} \tilde{z}^{[1]} + \beta^{[1]}$

Andrew Ng

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 $\hat{\underline{z}}^{\{t\}}$.

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

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

Works w/ momentum, RMSprop, Adam.



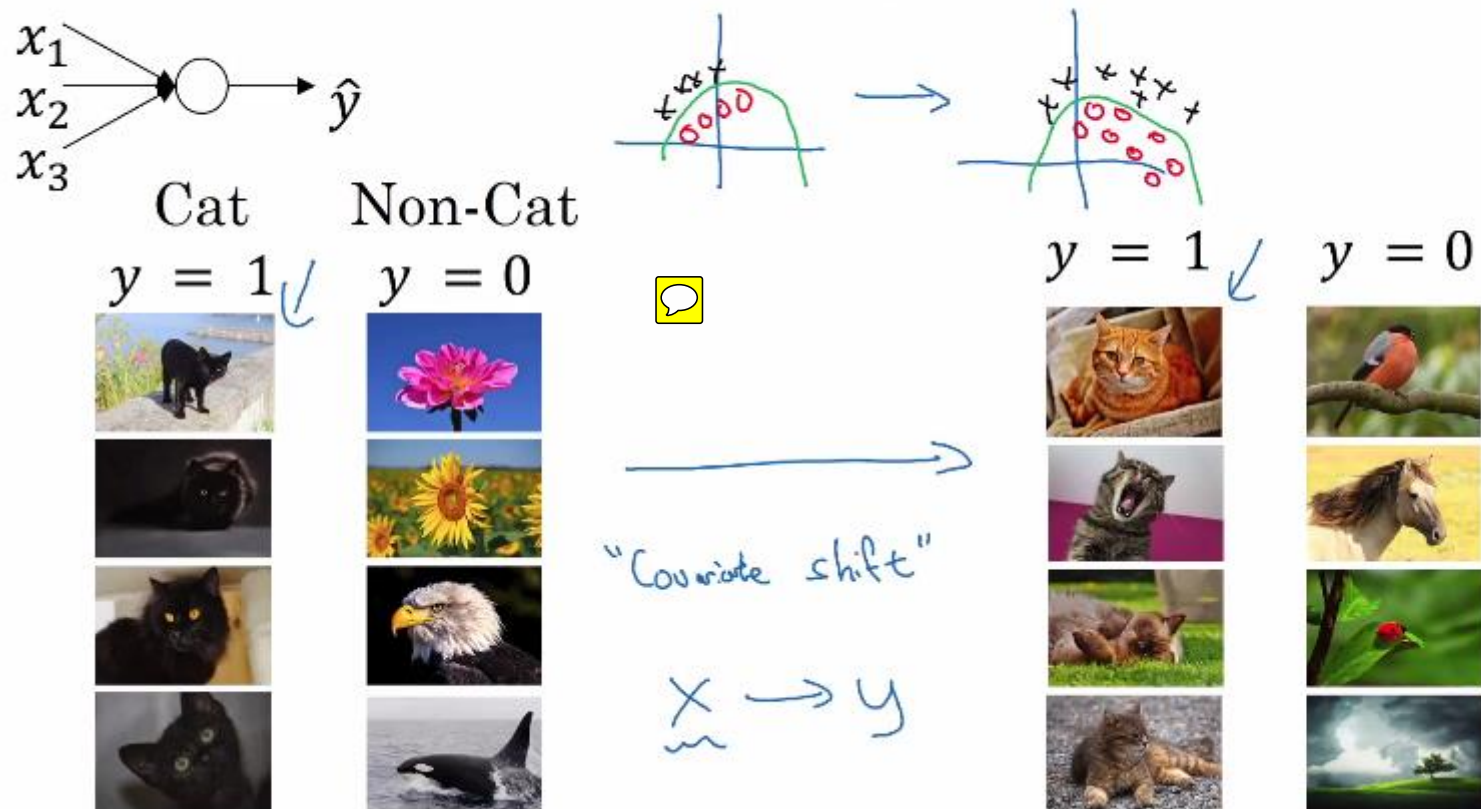
deeplearning.ai

Batch Normalization

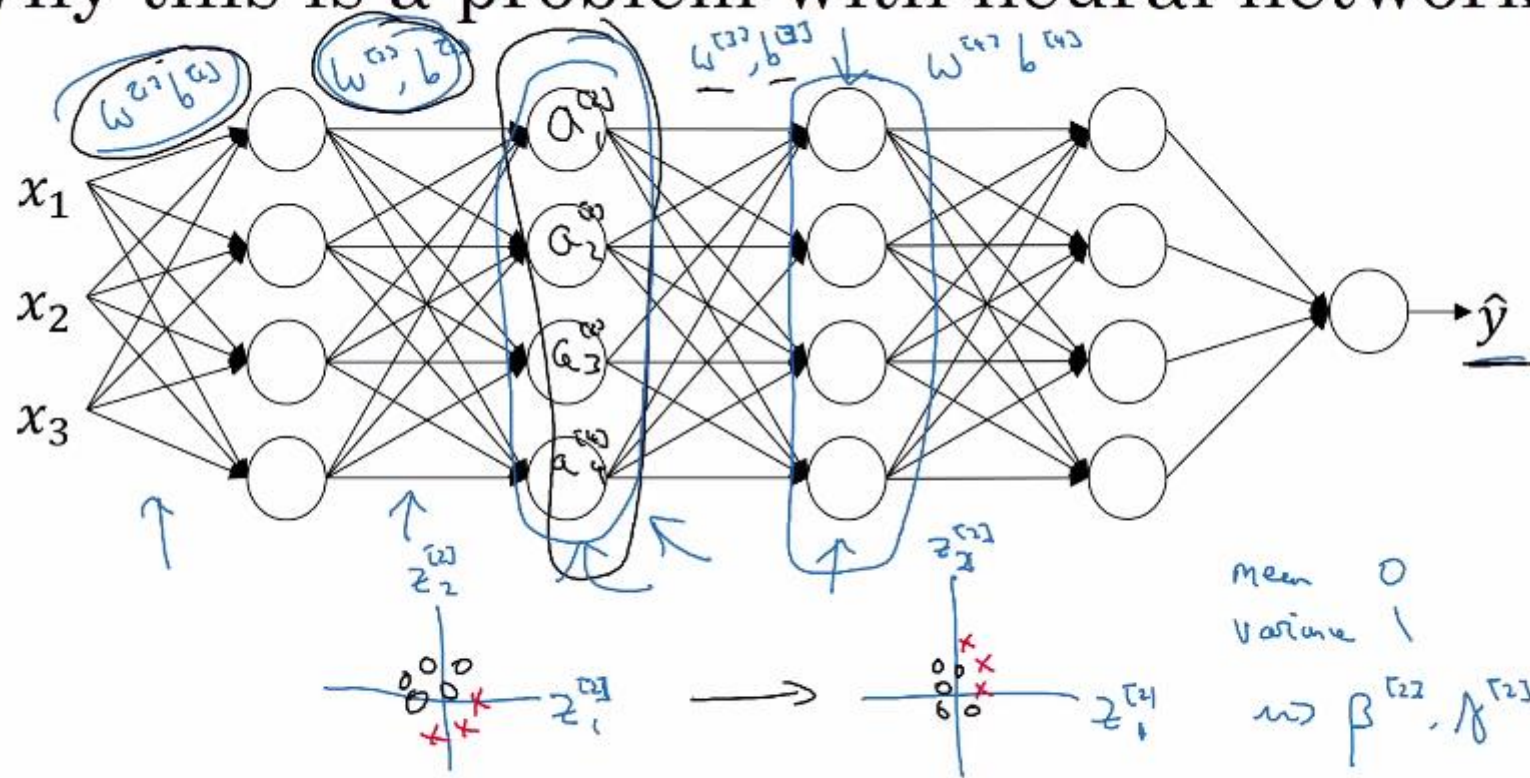
Why does
Batch Norm work?



Learning on shifting input distribution



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

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

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

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

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

$$X^{[1]}, X^{[2]}, X^{[3]}, \dots$$

$$\begin{array}{ccc} \downarrow & \downarrow & \downarrow \\ \mu_{[1]}^{[1]} & \mu_{[2]}^{[2]} & \mu_{[3]}^{[3]} \rightarrow \underline{\mu} \\ \theta_1 & \theta_2 & \theta_3 \\ \sigma_{[1]}^{[1]} & \sigma_{[2]}^{[2]} & \dots \rightarrow \underline{\sigma}^2 \end{array}$$

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

Andrew Ng



deeplearning.ai

Multi-class
classification

Softmax regression

Recognizing cats, dogs, and baby chicks, *other*



3

1

2

0

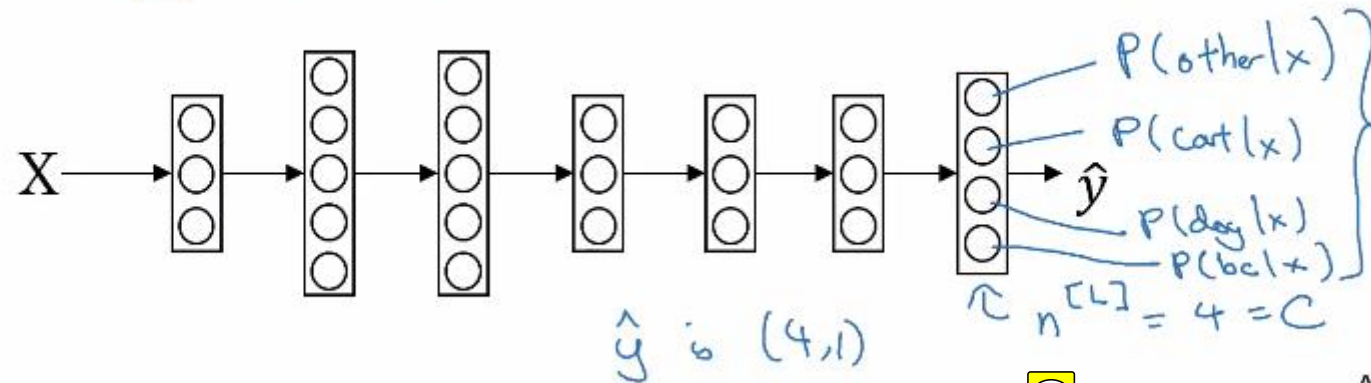
3

2

0

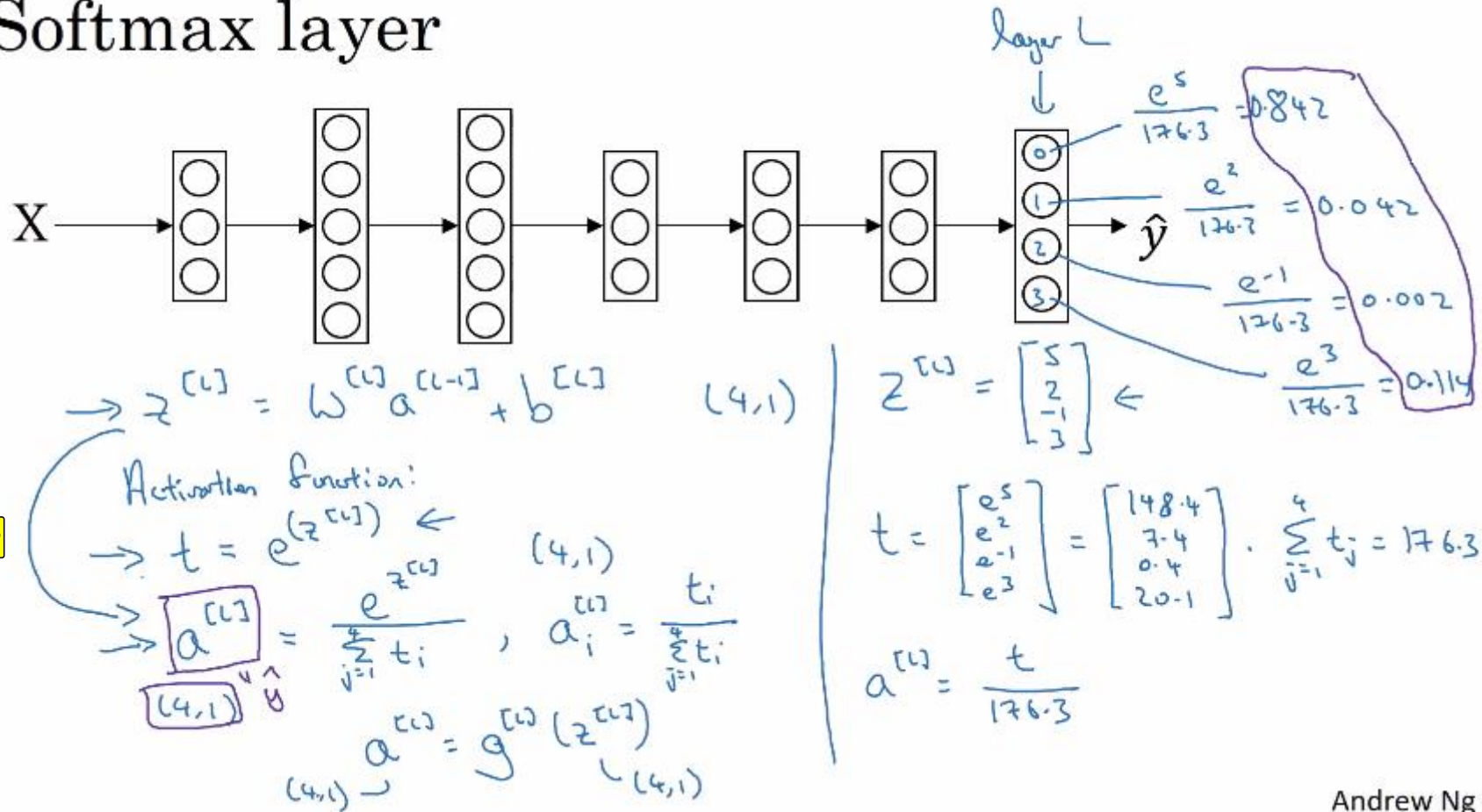
1

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



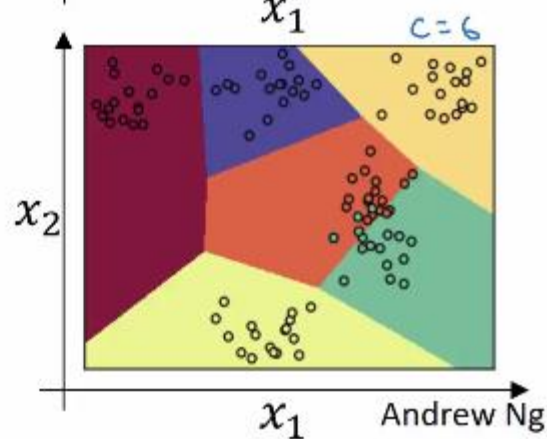
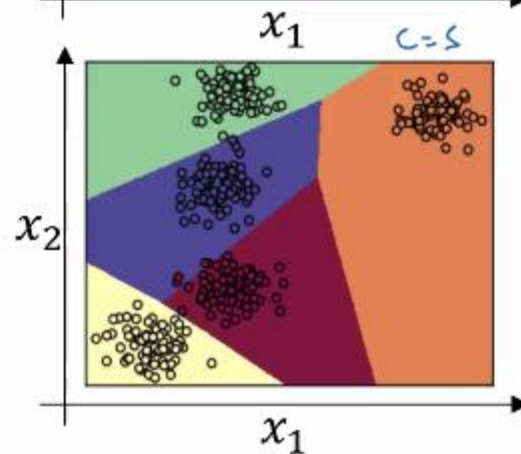
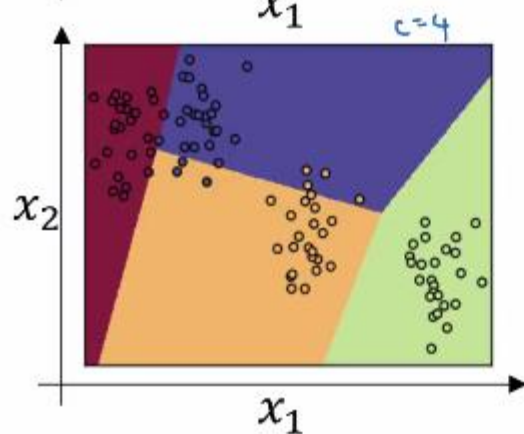
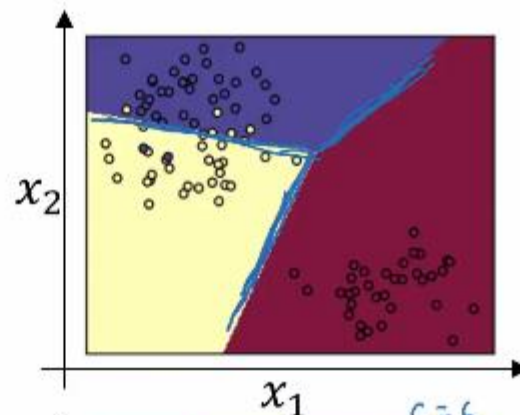
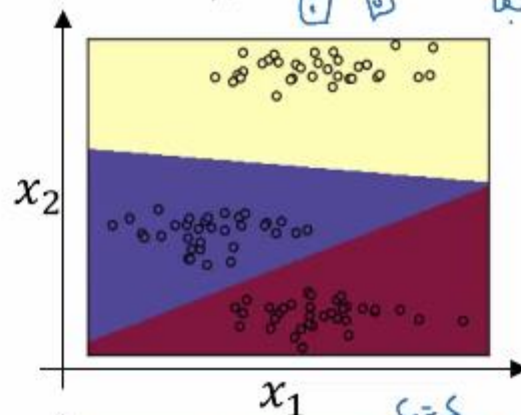
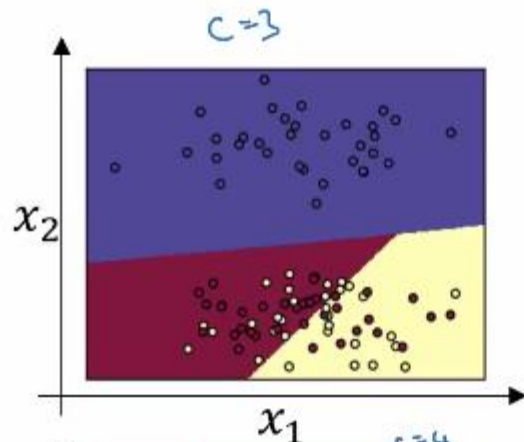
Andrew Ng

Softmax layer



Softmax examples

$$\begin{aligned}
 & \begin{matrix} x_1 \\ x_2 \end{matrix} \rightarrow \begin{bmatrix} 0.00 \\ 0.00 \end{bmatrix} \rightarrow \hat{y} \\
 & x \rightarrow \begin{bmatrix} 0.00 \\ 0.00 \end{bmatrix} \rightarrow \dots \rightarrow \begin{bmatrix} 0.00 \\ 0.00 \end{bmatrix} \rightarrow \hat{y} \\
 & z^{(1)} = W^{(1)}x + b^{(1)} \\
 & a^{(1)} = \hat{y} = g(z^{(1)})
 \end{aligned}$$



Andrew Ng



deeplearning.ai

Multi-class
classification

Training a softmax
classifier

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

$(4,1)$
 $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ - cat
 $y_2 = 1$
 $y_1 = y_3 = y_4 = 0$

$(4,1)$
 $\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$

$C=4$

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

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

$-y_2 \log \hat{y}_2 = -\log \hat{y}_2$ Make \hat{y}_2 big.

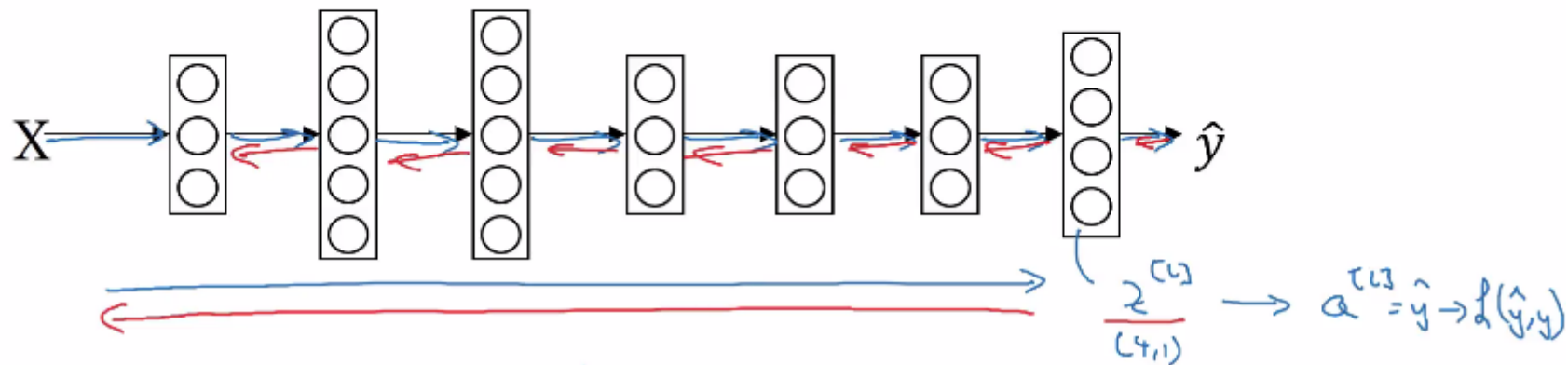
$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(n)}]$

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

$= \begin{bmatrix} 0 & 0 & 1 & \dots \\ 1 & 0 & 0 & \dots \\ 0 & 1 & 0 & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}$
 $(4, m)$

$= \begin{bmatrix} 0.3 & \dots \\ 0.2 & \dots \\ 0.1 & \dots \\ 0.4 & \dots \end{bmatrix}$
 $(4, m)$

Gradient descent with softmax



Backprop:

$$\frac{\partial z^{[L]}_{(4,1)}}{\partial z^{[L]}_{(4,1)}} = \hat{y} - y$$

$$\rightarrow \frac{\partial J}{\partial z^{[L]}_{(4,1)}}$$



deeplearning.ai

Programming
Frameworks

Deep Learning
frameworks

4:15



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

$$J(w) = \boxed{w^2 - 10w + 25}$$

\nwarrow
 $(w-5)^2$
 $w=5$

$$J(w, b)$$

$\uparrow \uparrow$
.

TensorFlow example

localhost:8891/notebooks/TensorFlow%20example%20.ipynb

jupyter TensorFlow example Last Checkpoint: 4 minutes ago (unsaved changes) Logout

File Edit View Insert Cell Kernel Help Trusted Python 3

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

In [2]: `w = tf.Variable(0, dtype=tf.float32)
cost = tf.add(tf.add(w**2, tf.multiply(-10., w)), 25)
train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)

init = tf.global_variables_initializer()
session = tf.Session()
session.run(init)
print(session.run(w))`
0.0

In [3]: `session.run(train)
print(session.run(w))`
0.1

In [4]: `for i in range(1000):
 session.run(train)
 print(session.run(w))`
4.99999

TensorFlow example Last Checkpoint: 10 minutes ago (autosaved) Logout

File Edit View Insert Cell Kernel Help Notebook saved Trusted Python 3

import tensorflow as tf

In [8]: `coefficients = np.array([[1.], [-10.], [25.]])`

`w = tf.Variable(0, dtype=tf.float32)`
`x = tf.placeholder(tf.float32, [3, 1])`
`#cost = tf.add(tf.add(w**2, tf.multiply(-10., w)), 25)`
`#cost = w**2 - 10*w + 25`
`cost = x[0][0]*w**2 + x[1][0]*w + x[2][0]`
`train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)`

`init = tf.global_variables_initializer()`
`session = tf.Session()`
`session.run(init)`
`print(session.run(w))`

0.0

In [9]: `session.run(train, feed_dict={x:coefficients})`
`print(session.run(w))`

0.1

In [10]: `for i in range(1000):`
 `session.run(train, feed_dict={x:coefficients})`
`print(session.run(w))`

0.999999



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))
```

