# Copyright Notice

These slides are distributed under the Creative Commons License.

Hyperparameter tuning

Tuning process

deeplearning.ai

# Hyperparameters

$\rightarrow$ $\boxed{\alpha}$

$\boxed{\beta}$ $\approx 0.9$

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

$\boxed{\text{\# layers}}$

$\boxed{\text{\#hidden units}}$

$\boxed{\text{learning rate decay}}$

$\boxed{\text{Mini-batch size}}$

Andrew Ng

# Try random values: Don't use a grid

Andrew Ng

# Coarse to fine

Hyperparameter 2

Hyperparameter 1

deeplearning.ai
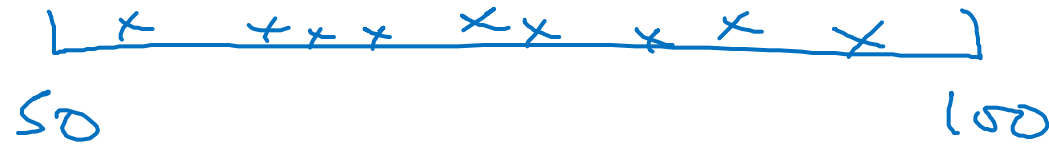
Hyperparameter tuning

Using an appropriate scale to pick hyperparameters

# Picking hyperparameters at random

$$\rightarrow n^{[\ell]} = 50, \ldots, 100$$
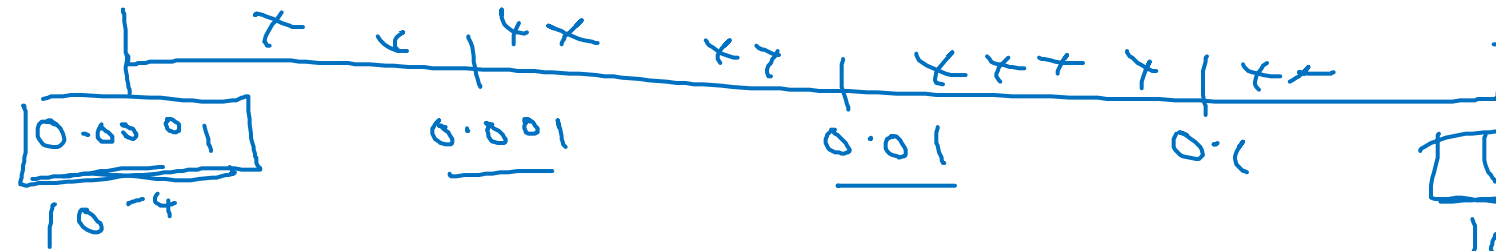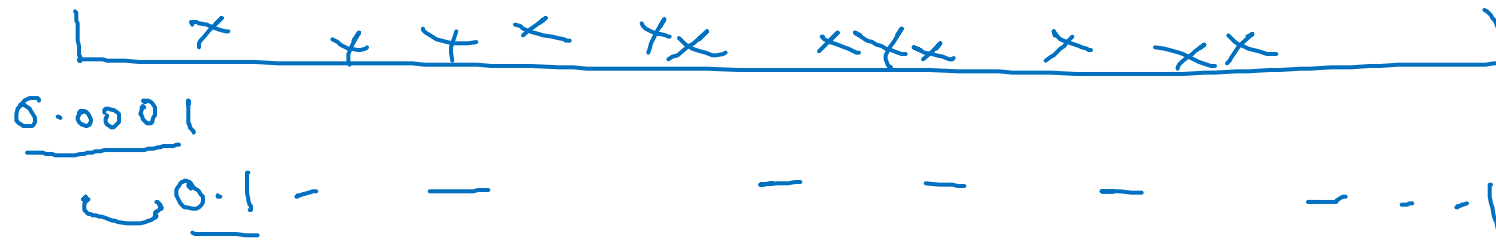
$$\underbrace{[\times \quad \times\times\times \quad \times\times \quad \times \quad \times \quad \times]}_{}$$
50                                                            100

$$\rightarrow \#layers \quad L: \quad 2 - 4$$

$$2, 3, 4$$

# Appropriate scale for hyperparameters

$\alpha = 0.0001 \ldots\ldots, 1$



$0.0001$            $1$

$\sim, 0.1 - \quad - \quad\quad - \quad - \quad - \quad -\cdots 1$

$0.0001$     $0.001$     $0.01$     $0.1$     $10^0$   $10^6$

$10^{-4}$          $10^0$

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

$a = \log_{10} 0.0001 \qquad r = -4 * np.random.rand() \qquad \leftarrow \quad r \in [-4, 0]$

$= -4 \qquad\qquad \alpha = 10^r \qquad\qquad\qquad \leftarrow \quad 10^{-4} \ldots 10^0$

$10^a \ldots 10^b \qquad\qquad \underset{[-4, 0]}{r \in [a, b]} \qquad\qquad \alpha = 10^r$

Andrew Ng

# Hyperparameters for exponentially weighted averages

$\beta = 0.9 \quad \cdots \quad 0.999$

$\downarrow \qquad\qquad \downarrow$

$10 \qquad\qquad\quad 1000$

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

_____

$\beta: \quad 0.9000 \to 0.9005 \quad \big\} \sim 10$

$\beta: \quad 0.999 \to 0.9995$

$\sim 1000 \qquad\qquad \sim 2000$

$\dfrac{1}{1-\beta}$

$0.9 \qquad\qquad\qquad 0.999$

$0.9 \qquad 0.99 \qquad 0.999$

$0.1 \qquad 0.01 \qquad 0.001$

$10^{-1} \qquad\qquad\qquad 10^{-3}$
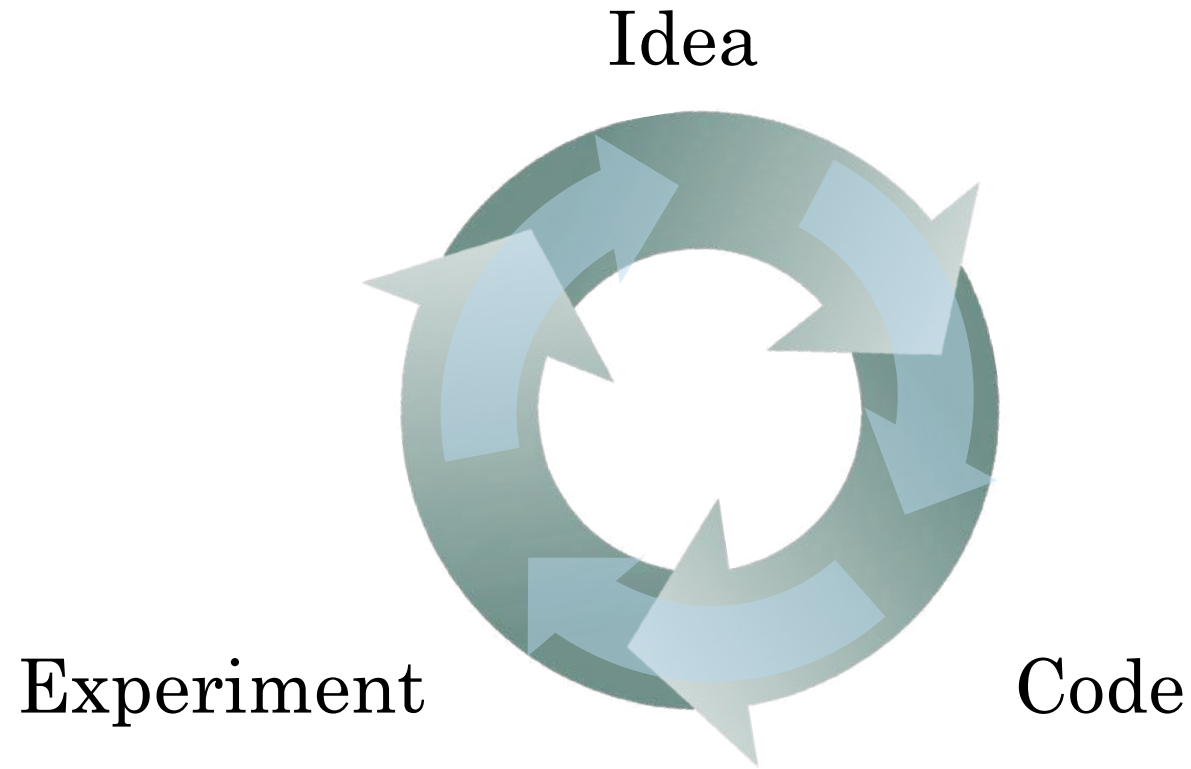
$r \in [-3, -1]$

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

$\beta = 1 - 10^{r}$
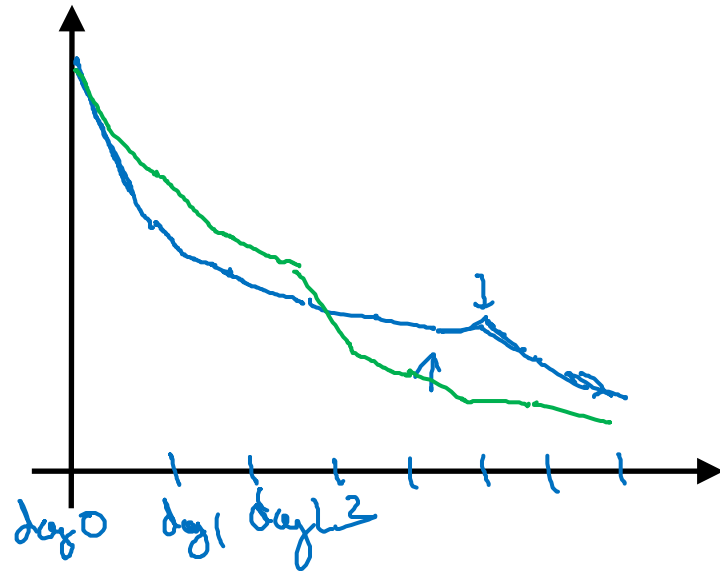
deeplearning.ai

Hyperparameters tuning

Hyperparameters tuning in practice: Pandas vs. Caviar

# Re-test hyperparameters occasionally
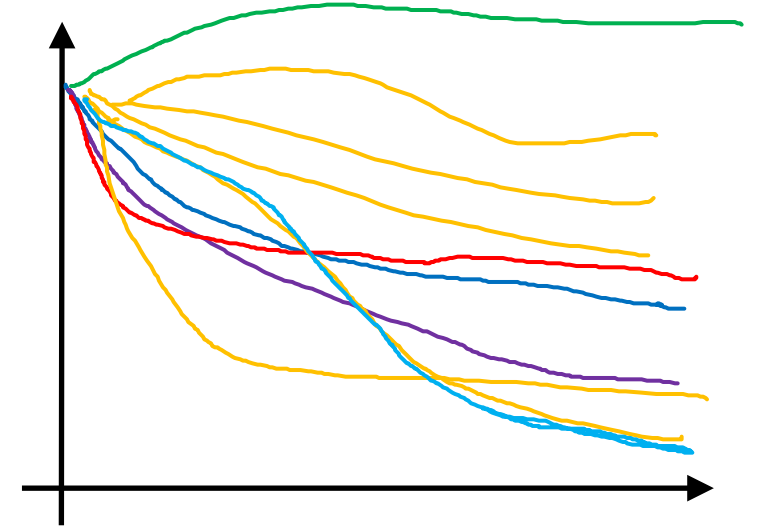


Idea

Experiment

Code

- NLP, Vision, Speech, Ads, logistics, ….

- Intuitions do get stale. Re-evaluate occasionally.

# Babysitting one model



day 0    day 1   day 2  3

Panda ←

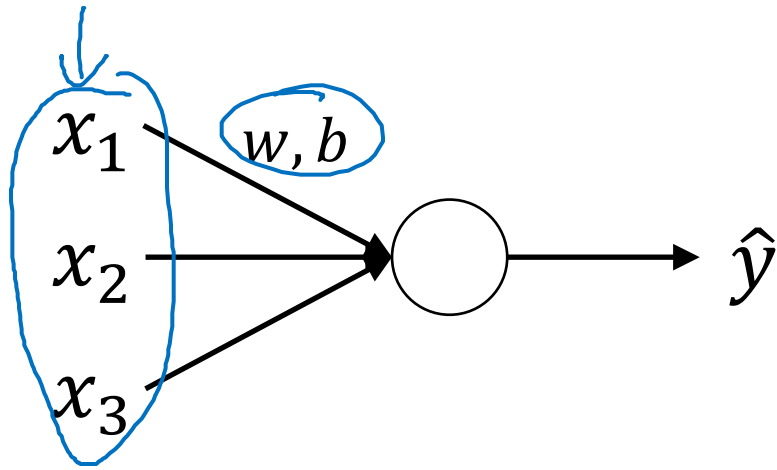# Training many models in parallel



Caviar ←

Andrew Ng

Batch
Normalization

Normalizing activations
in a network

deeplearning.ai

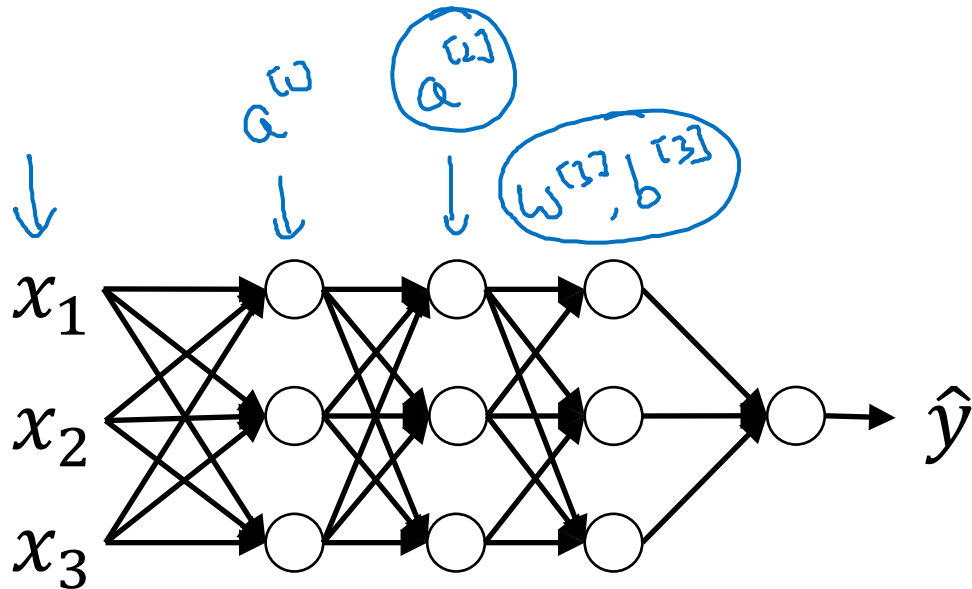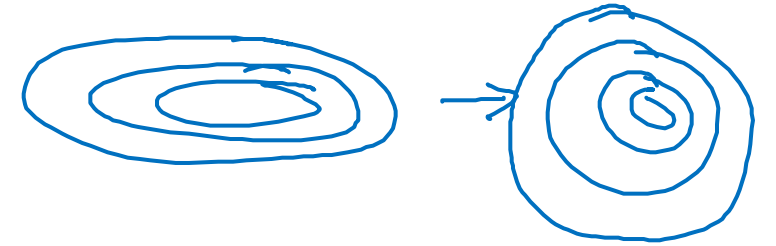# Normalizing inputs to speed up learning



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

$$X = X - \mu$$

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

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

Can we normalize $a^{[2]}$ so as to train $W^{[3]}, b^{[3]}$ faster

Normalize $z^{[2]}$

# Implementing Batch Norm

Given some intermediate values in NN $\quad Z^{(i)}, \ldots, Z^{(m)}$

$\downarrow \quad \downarrow$

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

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

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

$Z_{norm}^{(i)} = \dfrac{Z^{(i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}} \quad \Leftarrow$

$\tilde{Z}^{(i)} = \gamma Z_{norm}^{(i)} + \beta$

learnable parameters of model.

If

$\gamma = \sqrt{\sigma^2 + \varepsilon} \quad \Leftarrow$

$\beta = \mu \quad \Leftarrow$

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

$X \quad \Leftarrow$

$Z^{(i)} \quad \Leftarrow$

Use $\tilde{Z}^{[l](i)}$ instead of $Z^{[l](i)}$.

Batch
Normalization

Fitting Batch Norm
into a neural network

deeplearning.ai

# Adding Batch Norm to a network



$X \xrightarrow{W^{[1]}, b^{[1]}} z^{[1]} \xrightarrow{\boxed{\text{Batch Norm (BN)}}} \tilde{z}^{[1]} \xrightarrow{\beta^{[1]}, \gamma^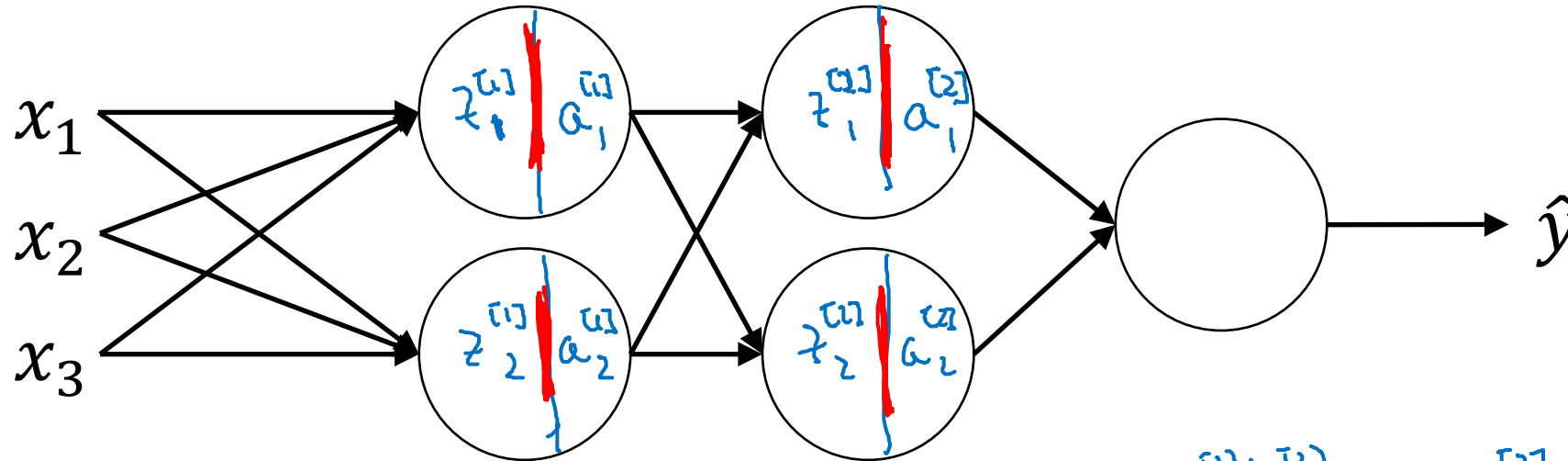{[1]}} a^{[1]} = g^{[1]}(\tilde{z}^{[1]}) \xrightarrow{W^{[2]}, b^{[2]}} z^{[2]} \xrightarrow[\text{BN}]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \to a^{[2]} \to \dots$

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, \dots, W^{[L]}, b^{[L]},$
$\to \beta^{[1]}, \gamma^{[1]}, \beta^{[2]}, \gamma^{[2]}, \dots, \beta^{[L]}, \gamma^{[L]} \Big\}$

$d\beta^{[l]}$

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

$\to \beta$

tf.nn.batch-normalization

# Working with mini-batches

$$X^{\{1\}} \xrightarrow{W^{[1]}, b^{[1]}} Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow g^{[1]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[2]}, b^{[2]}} Z^{[2]} \rightarrow \dots$$

$$X^{\{2\}} \longrightarrow Z^{[1]} \xrightarrow[BN]{\beta^{[1]}, \gamma^{[1]}} \tilde{Z}^{[1]} \rightarrow \dots$$

$$X^{\{3\}} \longrightarrow \dots$$

Parameters: $w^{[l]}, \cancel{b^{[l]}}, \beta^{[l]}, \gamma^{[l]}$.

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

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

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

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

$Z^{[l]}_{norm}$

$\rightarrow \tilde{Z}^{[l]} = \gamma^{[l]} Z^{[l]}_{norm} + \beta^{[l]} \leftarrow$

Andrew Ng

# Implementing gradient descent

for $t = 1 \ldots$ num MiniBatches

Compute forward prop on $X^{\{t\}}$.

In each hidden layer, use BN to replace $Z^{[l]}$ with $\tilde{Z}^{[l]}$.

Use backprop to compute $dW^{[l]}$, $\cancel{db^{[l]}}$, $d\beta^{[l]}$, $d\gamma^{[l]}$

Update parameters

$$W^{[l]} := W^{[l]} - \alpha \, dW^{[l]}$$
$$\beta^{[l]} := \beta^{[l]} - \alpha \, d\beta^{[l]}$$
$$\gamma^{[l]} := \ldots$$

$\longleftarrow$

Works w/ momentum, RMSprop, Adam.

Andrew Ng

Batch
Normalization

Why does
Batch Norm work?

deeplearning.ai

# Learning on shifting input distribution



$x_1$
$x_2$
$x_3$
$\hat{y}$
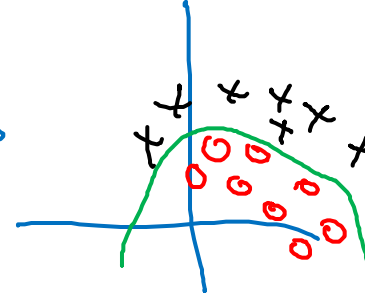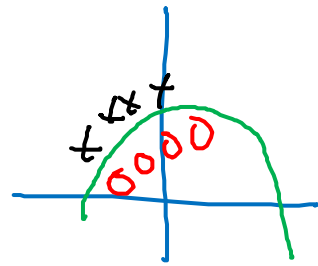
Cat    Non-Cat
$y = 1$    $y = 0$

$y = 1$    $y = 0$

"Covariate shift"

$x \longrightarrow y$

Andrew Ng

# Why this is a problem with neural networks?



$W^{[3]}, b^{[3]}$

$W^{[4]}, b^{[4]}$

$a_1^{[2]}$

$a_2^{[2]}$

$a_3^{[2]}$

$a_4^{[2]}$

$\hat{y}$

$z_2^{[2]}$

$z_1^{[2]}$

$z_1^{[2]}$

$z_1^{[2]}$

Mean 0

Variance 1

$\Rightarrow \beta^{[2]}, \gamma^{[2]}$

Andrew Ng

# 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.

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

$\{64, 128\}$

$z^{[l]}$

$\mu, \sigma^2$

Mini-batch : $64 \longrightarrow 512$

# Batch Norm at test time

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

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

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

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

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

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

$\mu^{\{1\}[l]}$  $\mu^{\{2\}[l]}$  $\mu^{\{3\}[l]}$ $\longrightarrow \mu$

$\theta_1$   $\theta_2$   $\theta_3$   $\sigma^2$

$\sigma^{2\{1\}[l]}$ , $\sigma^{2\{2\}[l]}$ , $\dots$

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

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

Andrew Ng

# Batch Norm at test time

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

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

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

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

Multi-class classification

Softmax regression

deeplearning.ai

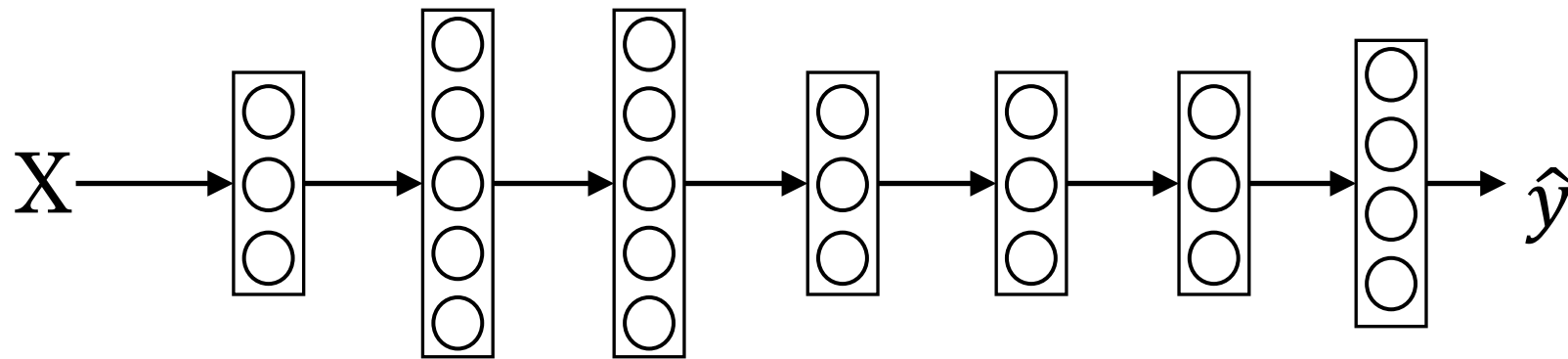# Recognizing cats, dogs, and baby chicks
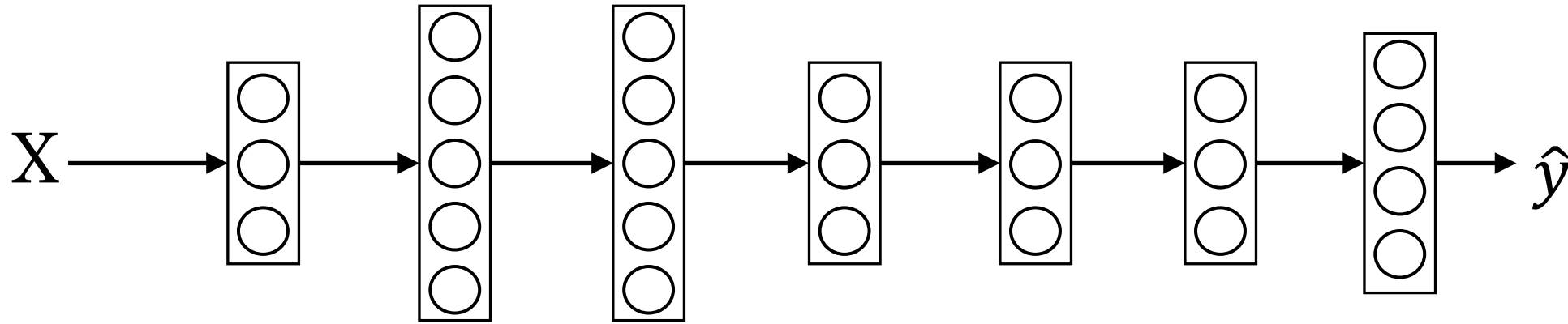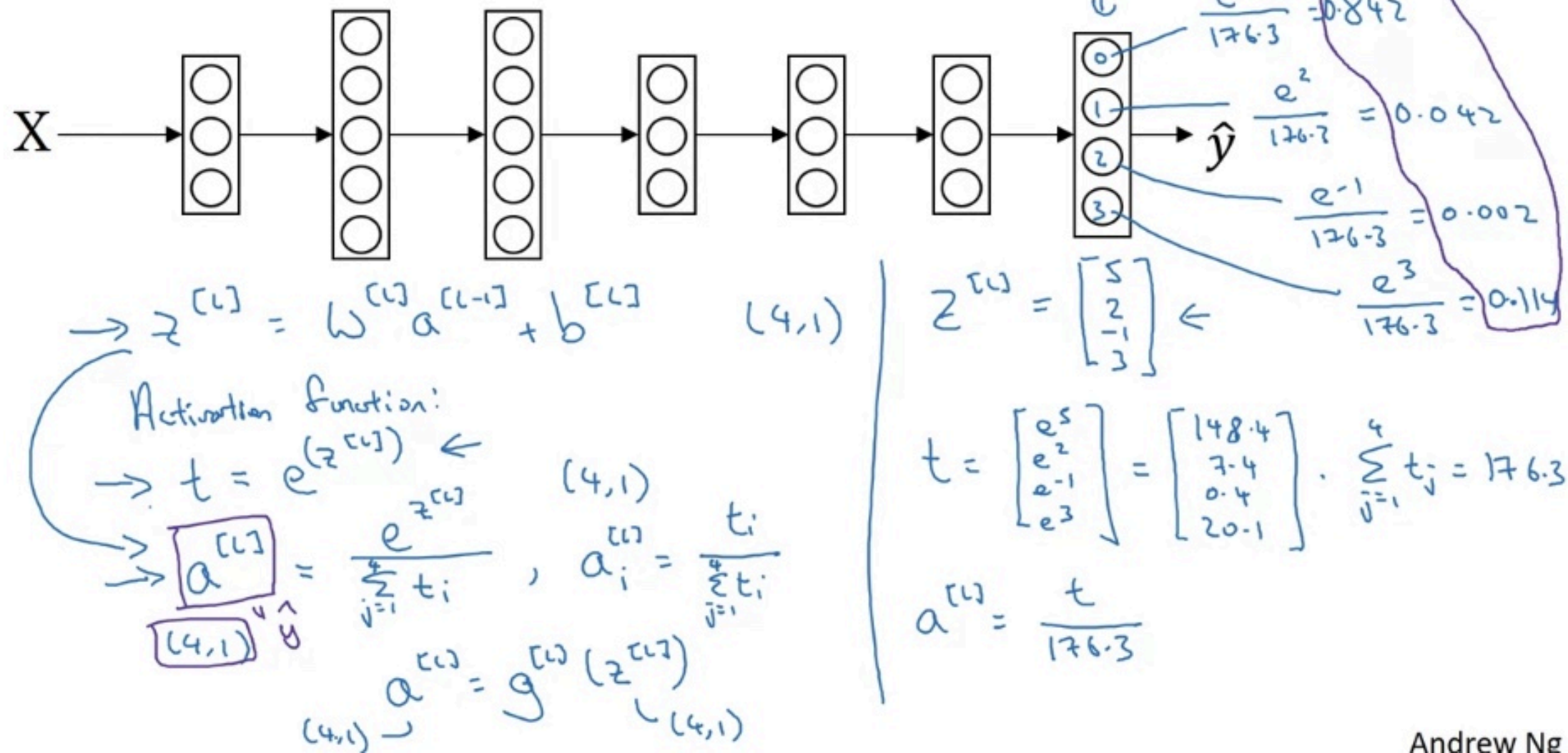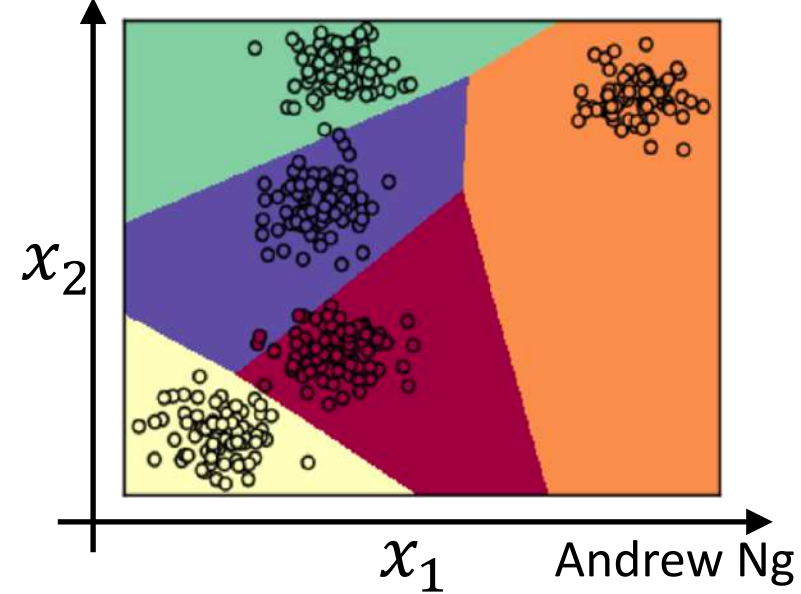


3      1      2      0      3      2      0      1

$$X \rightarrow \cdots \rightarrow \hat{y}$$

# Softmax layer



X → $\hat{y}$

# Softmax layer



layer L

$$\frac{e^5}{176.3} = 0.842$$

$$\frac{e^2}{176.3} = 0.042$$

$$\frac{e^{-1}}{176.3} = 0.002$$

$$\frac{e^3}{176.3} = 0.114$$

$$\rightarrow z^{[L]} = w^{[L]} a^{[L-1]} + b^{[L]} \qquad (4,1)$$

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

Activation function:

$$\rightarrow t = e^{(z^{[L]})} \leftarrow \qquad (4,1)$$

$$\rightarrow \boxed{a^{[L]}} = \frac{e^{z^{[L]}}}{\sum\limits_{j=1}^{4} t_i} \quad , \quad a_i^{[L]} = \frac{t_i}{\sum\limits_{j=1}^{4} t_i}$$

$$\boxed{(4,1)} \quad \hat{y}$$

$$\underset{(4,1)}{a^{[L]}} = g^{[L]} \underset{(4,1)}{(z^{[L]})}$$

$$t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix} \cdot \sum\limits_{j=1}^{4} t_j = 176.3$$

$$a^{[L]} = \frac{t}{176.3}$$

Andrew Ng
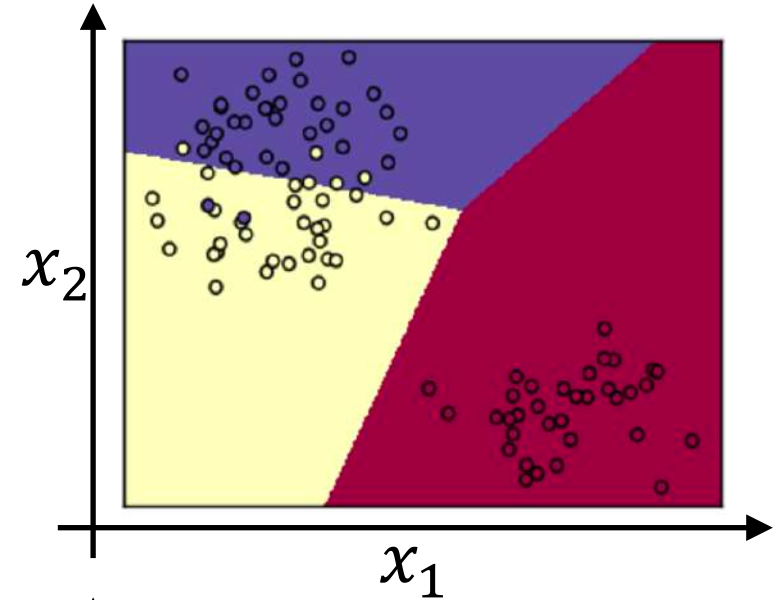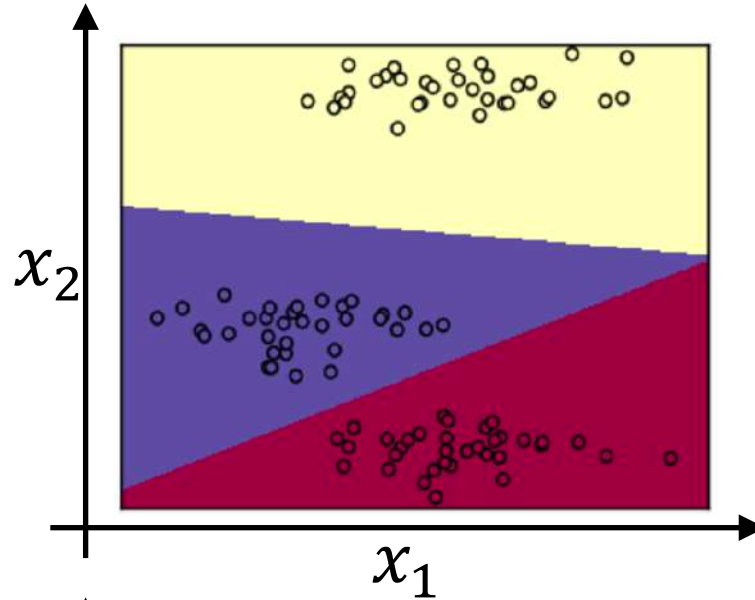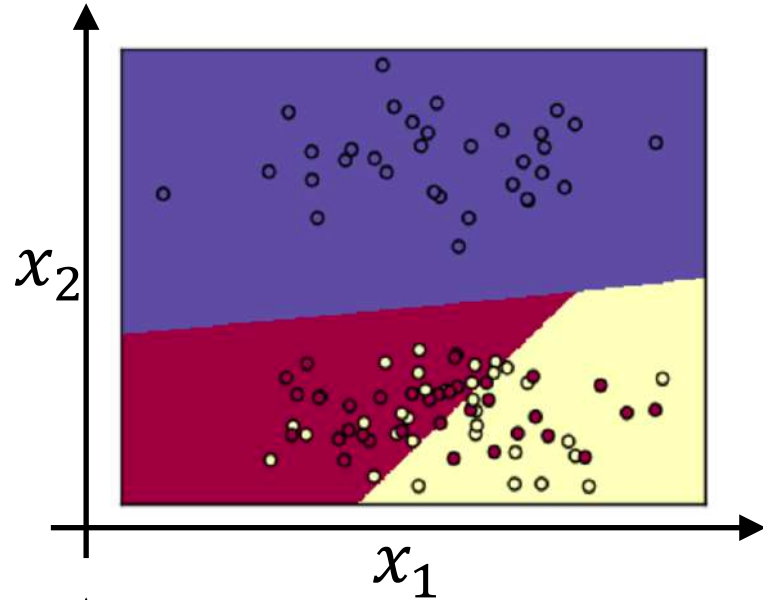
# Softmax examples

# Understanding softmax

$(4,1)$

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

$C = 4$ $\qquad g^{[L]}(\cdot)$

"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

$(4,1)$

$$y^{(i)} = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} \text{— cat} \\ y_2 = 1 \end{matrix}$$

$y_1 = y_3 = y_4 = 0$

$(4,1)$

$$a^{[L](i)} \approx \hat{y}^{(i)} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$

$C = 4$

$$\mathcal{L}(\hat{y}, y) = -\sum_{j=1}^{4} y_j \log \hat{y}_j$$

small

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

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

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

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

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

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

$(4, m)$

$(4, m)$

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)

Programming
Frameworks

TensorFlow

deeplearning.ai

# Motivating problem

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

(cost)

$(\omega - 5)^2$

$\omega = 5$

$J(W, b)$

# Code example

```python
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))
```

```python
with tf.Session() as session:

    session.run(init)

    print(session.run(w))
```

Andrew Ng