Branch: **master** ▾

Find file   Copy path

**deep-learning-specialization-coursera** / 02-Improving-Deep-Neural-Networks / week3 /

## hyperparameter-tuning-and-programming-frameworks.ipynb

**andersy005** add datasets

64a8031   on Oct 15, 2018

**1** contributor

‹› 📄      Raw   Blame   History      🖥 ✏ 🗑
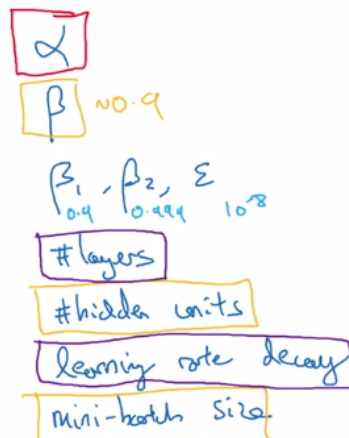
455 lines (455 sloc)    12.7 KB

# Table of Contents
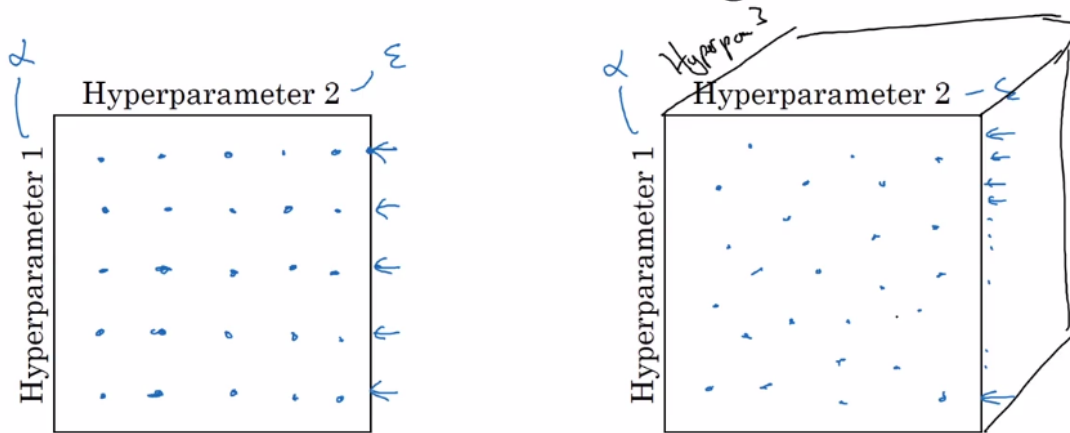
# Hyperparameter tuning, Batch Normalization and Programming Frameworks
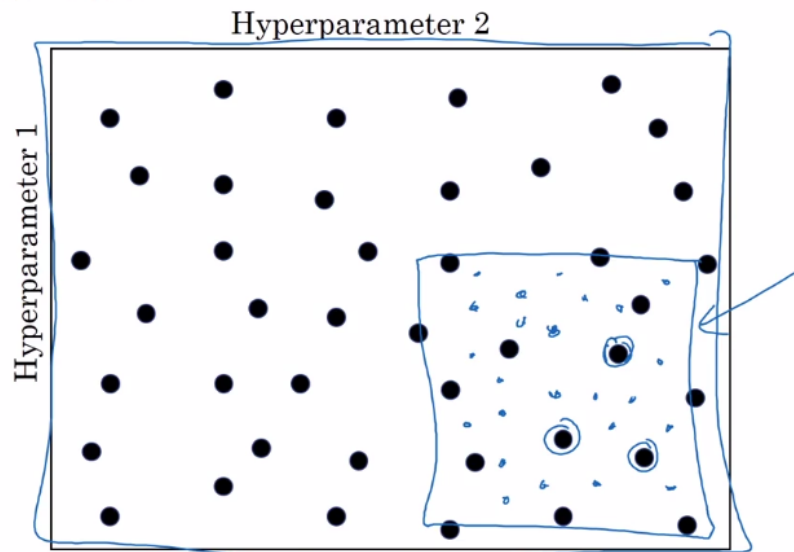
## Hyperparameter Tuning

### Tuning Process
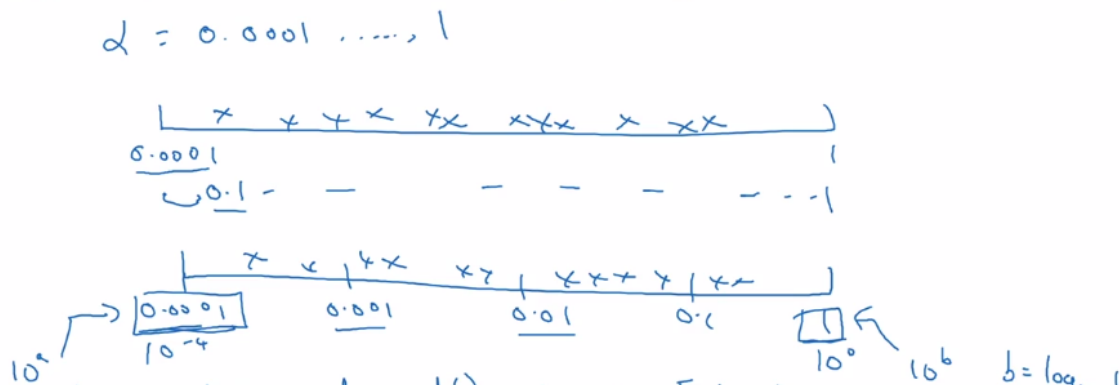
# Try random values: Don't use a grid



Andrew Ng

# Coarse to fine



Andrew Ng

**Using an appropriate scale to pick hyperparameters**

# Appropriate scale for hyperparameters

$a = by_0 0.0001$   $r = -4 * np.random.rand()$   $\leftarrow$   $r \in [-4, 0]$     $\frac{}{} = \sqrt{10}\ ^1$
$= -4$    $\alpha = 10^r$      $\leftarrow$   $10^{-4} \ldots 10^0$      $= 0$

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

# Hyperparameters for exponentially weighted averages

$\beta = 0.9 \quad \ldots \quad 0.999$
     $\downarrow$           $\downarrow$
     $10$          $1000$

$1-\beta = 0.1 \quad \ldots \quad 0.001$
_____

$\beta: \ 0.9000 \rightarrow 0.9005 \ \} \sim 10$
$\beta: \ 0.999 \rightarrow 0.9995$
    $\sim 1000$      $\sim 2000$

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

| $\times$ $\times$ $\times$ $\times$ $\times$ $\times$ | $\leftarrow$
0.9           0.999

0.9     0.99     0.999

0.1    0.01    0.001
$10^{-1}$           $10^{-3}$

$r \in [-3, -1]$
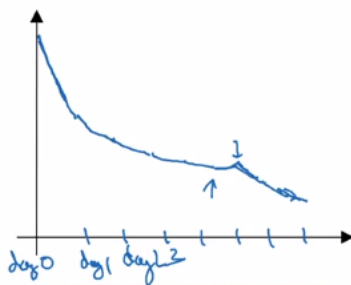$1-\beta = 10^r$
$\beta = 1 - 10^r$

## Hyperparameters tuning in practice: Pandas vs. Caviar

# Babysitting one model



day 0   day 1   day 2



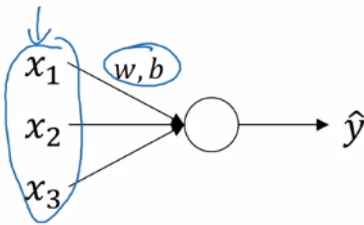Panda $\leftarrow$

# Training many models in parallel





Caviar $\leftarrow$

# Batch Normalization

## Normalizing activations in a network
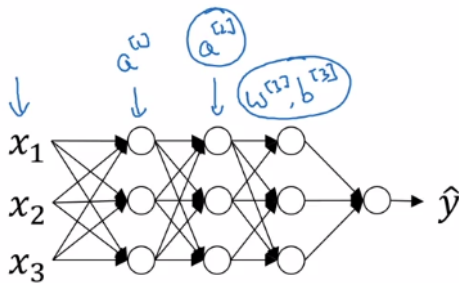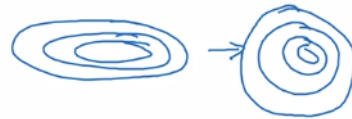
# Normalizing inputs to speed up learning



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

$$X = X - \mu$$

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

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

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

Normalize $z^{[2]}$

Andrew Ng

# Implementing Batch Norm

Given some intermediate values in NN $z^{(1)}, \ldots, z^{(m)}$

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

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

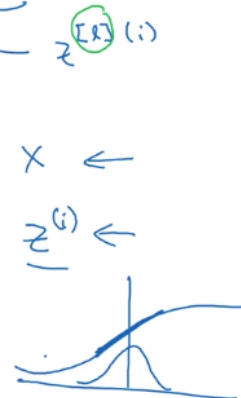$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}}$$

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

$\gamma, \beta$ — learnable parameters of model.

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

$\beta = \mu$ $\leftarrow$

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

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

$X \leftarrow$

$z^{(i)} \leftarrow$

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

Andrew Ng

## Fitting Batch Norm into a neural network

# Adding Batch Norm to a network

$$X \longrightarrow z^{[1]} \xrightarrow[\text{Batch Norm (BN)}]{\beta^{[1]}, \gamma^{[1]}} \tilde{z}^{[1]} \rightarrow a^{[1]} = g(\tilde{z}^{[1]}) \longrightarrow z^{[2]} \xrightarrow[BN]{\beta^{[2]}, \gamma^{[2]}} \tilde{z}^{[2]} \rightarrow a^{[2]} \rightarrow$$

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

$\beta$

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

tf.nn.batch-normalization $\Leftarrow$

Andrew Ng

# Working with mini-batches

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

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

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

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

$z^{[L]} \quad (n^{[L]}, 1)$

$\quad\quad (n^{[L]},1) \quad (n^{[L]},1) \quad (n^{[L]},1)$

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

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

$z^{[L]}_{norm}$

$\tilde{z}^{[L]} = \gamma^{[L]} z^{[L]}_{norm} + \beta^{[L]}$

Andrew Ng

# Implementing gradient descent

for $t = 1 \ldots$ num MiniBatches

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

     In each hidden layer, use BN to repace $z^{[L]}$ with $\tilde{z}^{[L]}$.

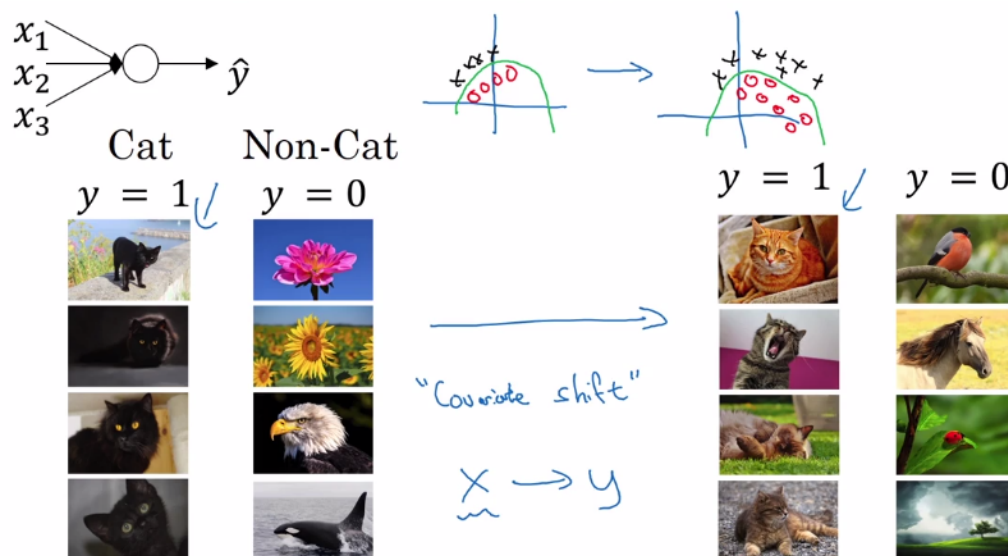   Use backprop to compute $dW^{[L]}, \cancel{db^{[L]}}, d\beta^{[L]}, d\gamma^{[L]}$

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

Works w/ momentum, RMSprop, Adam.

**Why does Batch Norm work?**

# Learning on shifting input distribution

$x_1$
$x_2$
$x_3$
$\rightarrow \hat{y}$

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?



$x_1$
$x_2$
$x_3$
$\rightarrow \hat{y}$

Mean   0
Variance   1

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

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.

$\{4, 128\}$   $z^{[l]}$

$\mu, \sigma^2$

Mini-batch : $64 \longrightarrow 512$

## Batch Norm at test time



# Multi-class classificiation

## Softmax Regression

$$z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]} \quad (4,1)$$

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

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

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

Activation function:

$$t = e^{(z^{[l]})} \quad (4,1)$$

$$a^{[l]} = \frac{e^{z^{[l]}}}{\sum_{j=1}^{4} t_i} \,, \quad a_i^{[l]} = \frac{t_i}{\sum_{j}^{4} t_i}$$

$$a^{[l]} = g^{[l]}(z^{[l]})$$
(4,1)  (4,1)

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

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

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

Andrew Ng

# Softmax examples

$$x_1 \\ x_2$$ → $\hat{y}$

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

C=3



Andrew Ng

## Training a softmax classifier

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

"soft max"

"hard max"

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

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

Softmax regression generalizes logistic regression to $C$ classes.

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

# Introduction to programming frameworks

## Deep Learning Frameworks

- Keras
- Lasagne
- mxnet
- PaddlePaddle
- TensorFlow
- Theano
- Torch

- Running speed
→ - Truly open (open source with good governance)

Andrew

## TensorFlow

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

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

```
In [3]:  train = tf.train.GradientDescentOptimizer(0.01).minimize(cost)
```

```
In [4]:  init = tf.global_variables_initializer()
```

```
In [5]:  session = tf.Session()
```

```
In [6]:  %time session.run(init)
```

```
CPU times: user 6.07 ms, sys: 2.41 ms, total: 8.48 ms
Wall time: 5.82 ms
```

```
In [7]:  print(session.run(w))
```

```
0.0
```

```
In [8]:  %time session.run(train)
```

```
CPU times: user 16.1 ms, sys: 2.34 ms, total: 18.4 ms
Wall time: 15.9 ms
```

```
In [9]:  %time print(session.run(w))
```

```
0.099999994
CPU times: user 815 $\mu$s, sys: 243 $\mu$s, total: 1.06 ms
Wall time: 807 $\mu$s
```

```
In [10]:  %load_ext version_information
          %version_information tensorflow, numpy
```

Out[10]:

| Software | Version |
|---|---|
| Python | 3.6.6 64bit [GCC 4.2.1 Compatible Apple LLVM 6.1.0 (clang- |

| Python | 602.0.53)] |
|---|---|
| IPython | 7.0.1 |
| OS | Darwin 17.7.0 x86_64 i386 64bit |
| tensorflow | 1.10.0 |
| numpy | 1.15.1 |
| Sun Oct 14 21:48:53 2018 MDT | |