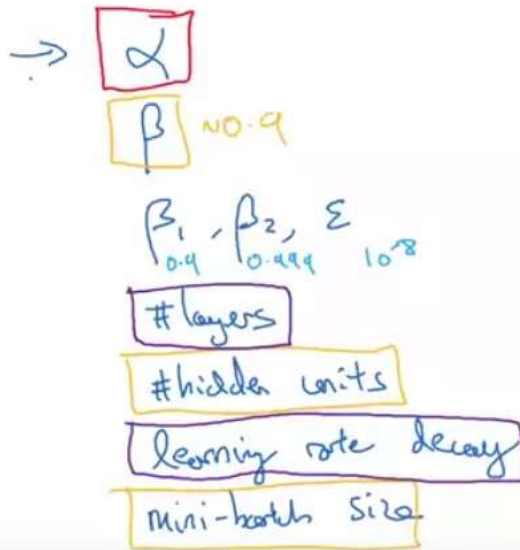


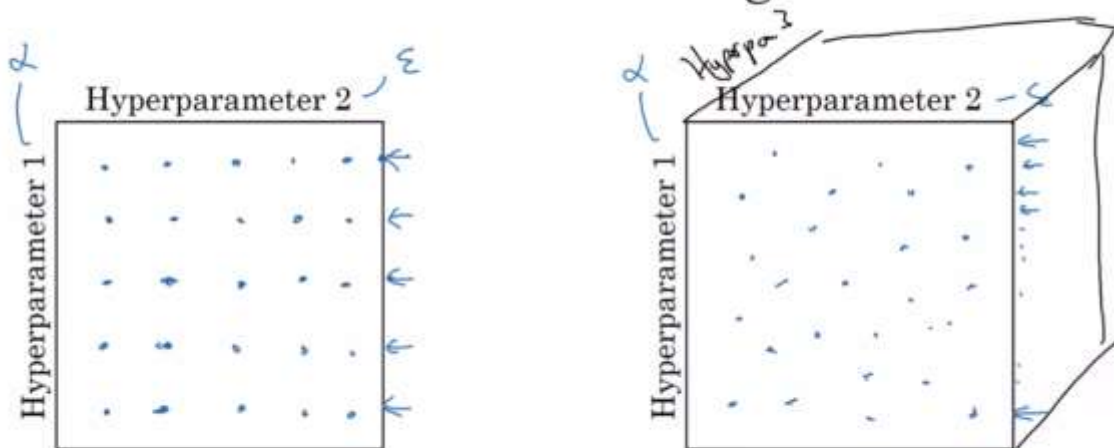
Hyperparameters Tuning

Week 3

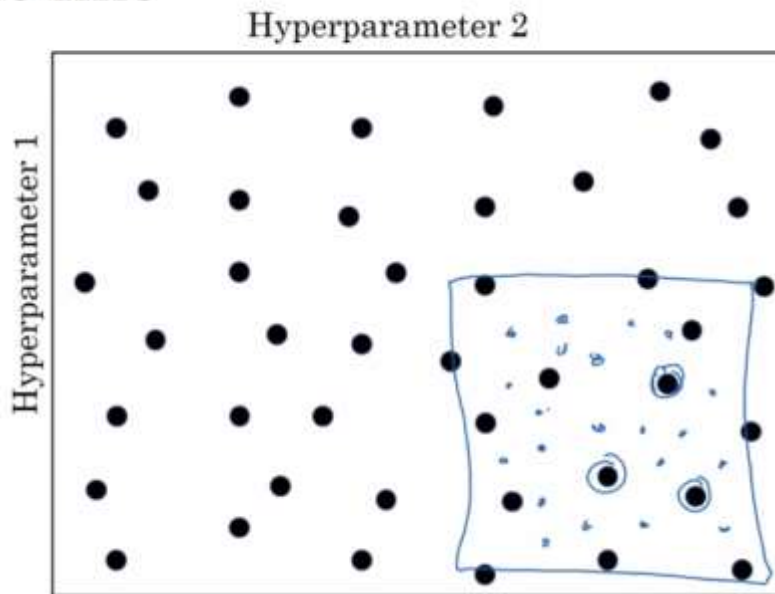
Hyperparameters



Try random values: Don't use a grid

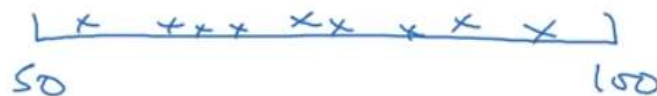


Coarse to fine



Picking hyperparameters at random

$$\rightarrow n^{\text{test}} = 50, \dots, 100$$



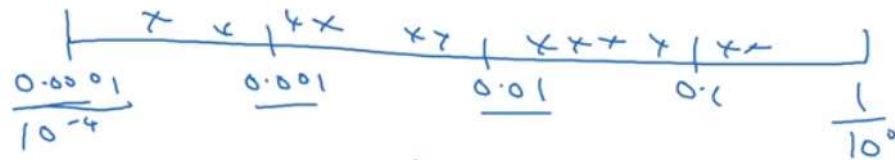
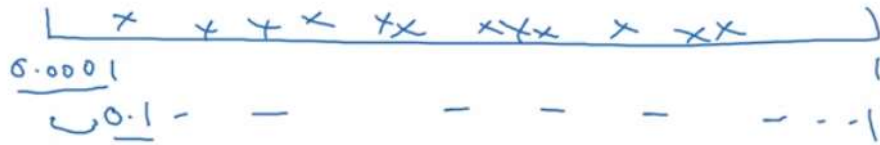
$$\rightarrow \# \text{layers} \quad L: \quad 2 - 4$$

$$2, 3, 4$$

We should not randomly search for hyperparameters, but randomly search in a specific range of values.

Appropriate scale for hyperparameters

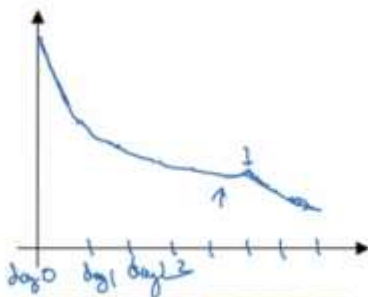
$$\alpha = 0.0001, \dots, 1$$



$$r = -4 * \text{np.random.randn}() \quad \leftarrow r \in [-4, 0]$$

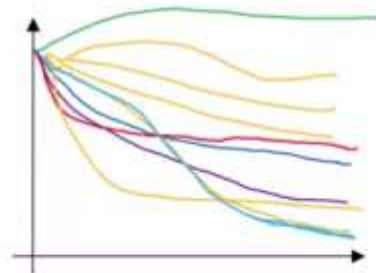
$$\alpha = 10^r \quad \leftarrow 10^{-4} \dots 10^0$$

Babysitting one model



Panda \leftarrow

Training many models in parallel



Caviar \leftarrow

Andrew Ni

Batch Normalization

Implementing Batch Norm

Given some intermediate values in NN

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

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

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

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

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

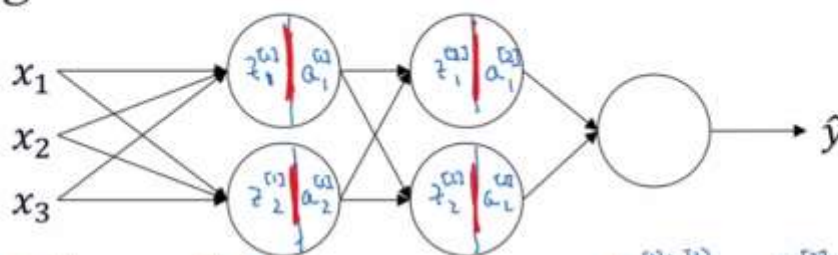
learnable parameters of model.

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

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

Use $\hat{z}^{(i)}$ instad of $z^{(i)}$.

Adding Batch Norm to a network



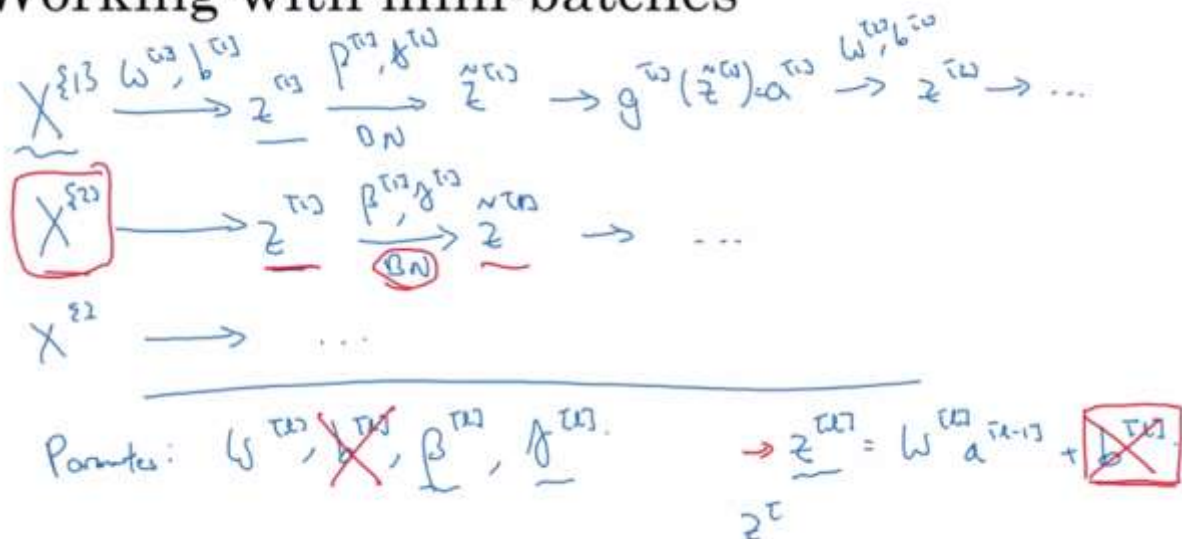
$$x \xrightarrow{w^{(1)}, b^{(1)}} z^{(1)} \xrightarrow{\beta^{(1)}, \gamma^{(1)}} \hat{z}^{(1)} \rightarrow a^{(1)} = g(\hat{z}^{(1)}) \xrightarrow{w^{(2)}, b^{(2)}} z^{(2)} \xrightarrow{\beta^{(2)}, \gamma^{(2)}} \hat{z}^{(2)} \rightarrow a^{(2)}$$

Batch Norm (BN)

Parameters: $\{w^{(1)}, b^{(1)}, w^{(2)}, b^{(2)}, \dots, w^{(L)}, b^{(L)}\}$
 $\rightarrow \{\beta^{(1)}, \gamma^{(1)}, \beta^{(2)}, \gamma^{(2)}, \dots, \beta^{(L)}, \gamma^{(L)}\}$

$d\beta^{(1)} \quad \beta^{(1)} = \beta^{(1)} - \alpha d\beta^{(1)}$

Working with mini-batches



Implementing gradient descent

for $t = 1 \dots \text{num Mini Batches}$

Compute forward pass on $\underline{X}^{[t]}$.

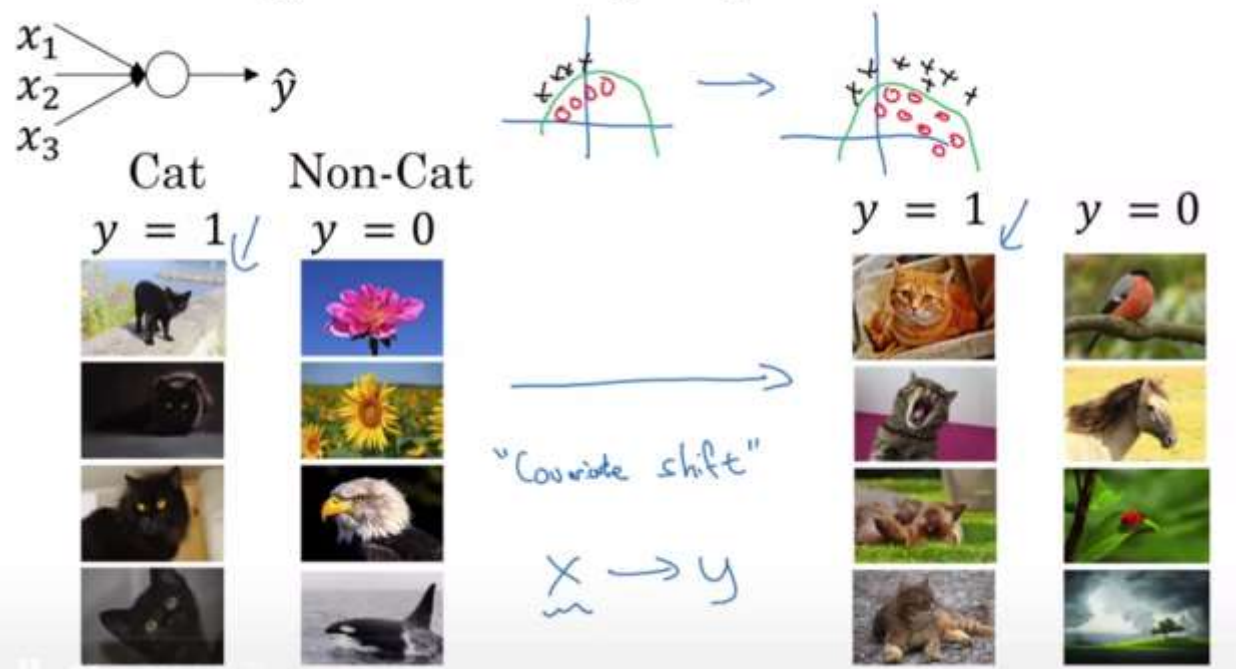
In each hidden layer, use BN to replace $\underline{z}^{[l]}$ with $\underline{\tilde{z}}^{[l]}$.

Use backprop to compute $\underline{dW}^{[1]}, \cancel{d\beta^{[1]}}, \underline{d\beta}^{[1]}, \underline{d\gamma}^{[1]}$

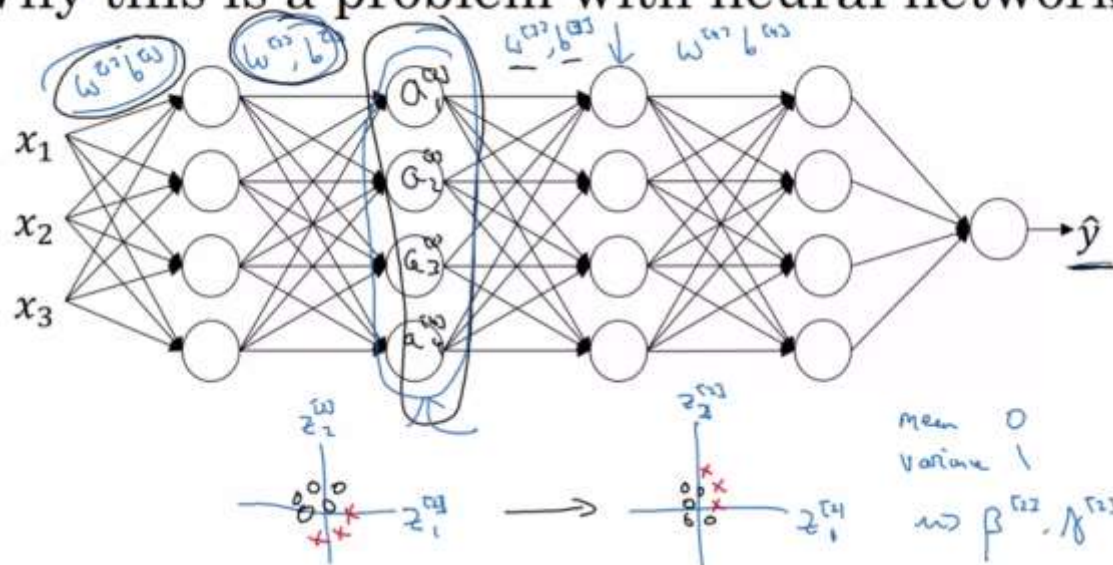
Update parameters $\left. \begin{aligned} W^{[1]} &:= W^{[1]} - \alpha dW^{[1]} \\ \beta^{[1]} &:= \beta^{[1]} - \alpha d\beta^{[1]} \\ \gamma^{[1]} &:= \dots \end{aligned} \right\}$

Works w/ momentum, RMSprop, Adam.

Learning on shifting input distribution



Why this is a problem with neural networks?



Batch Norm at test time

→ $\mu = \frac{1}{m} \sum_i z^{(i)}$

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

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

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

Handwritten notes:

μ, σ^2 : estimate using exponentially weighted average (across mini-batches).

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

↓

$\mu^{(1)}, \mu^{(2)}, \mu^{(3)}, \dots \rightarrow \mu$

$\sigma^{(1)}, \sigma^{(2)}, \sigma^{(3)}, \dots \rightarrow \sigma^2$

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

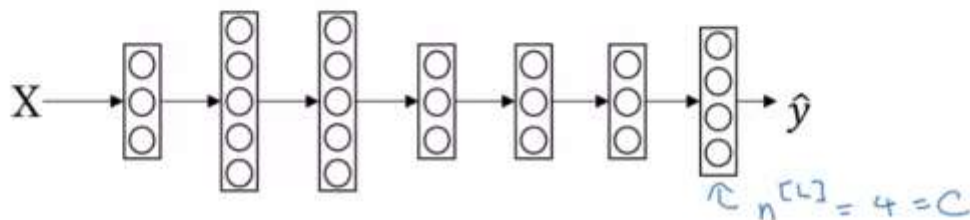
At test time, we will not have a mini-batch to compute the mean and std; during training our network, we will keep an exponentially weighted average for the mean and std., for each layer.

Recognizing cats, dogs, and baby chicks, ¹ ² ³ other ⁰

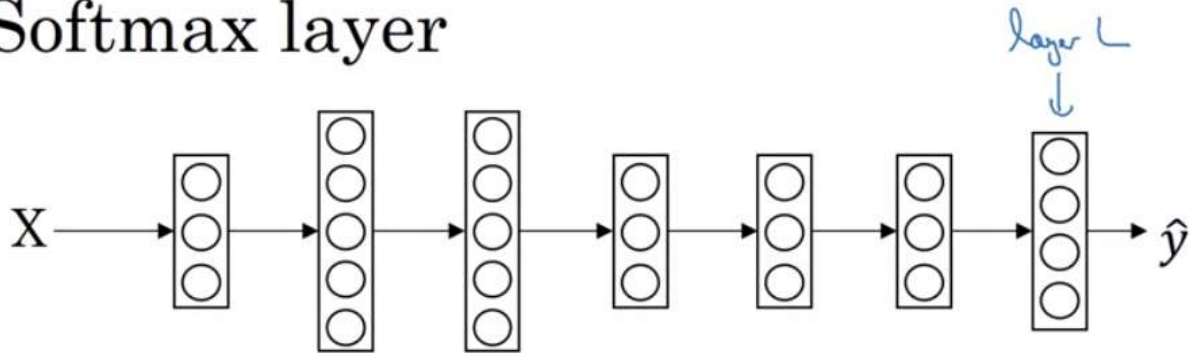


3 1 2 0 3 2 0 1

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



Softmax layer



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

Activation function:

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

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

(4,1)

$$\underbrace{\mathcal{L}(\hat{y}, y)}_{\text{small}} = - \sum_{j=1}^4 \underbrace{y_j \log \hat{y}_j}_{\uparrow} \quad \left| \quad \mathcal{J}(W^{[L]}, b^{[L]}, \dots) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \right.$$

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

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)

Tensorflow

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

```
In [5]: w = tf.Variable(0, dtype=tf.float32)
#cost = tf.add(tf.add(w**2, tf.multiply(-10., w)), 25)
cost = w**2 - 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 [6]: session.run(train)
print(session.run(w))

0.1
```

```
In [7]: for i in range(1000):
        session.run(train)
        print(session.run(w))
```

4.99999



6:14 / 16:07

```
: 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
```

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

As of 2019, Google launched TensorFlow 2. TensorFlow 2 borrowed its syntax from Keras. Keras is a high level library that can operate on top of TensorFlow 1 as well as other deep learning libraries. In later courses of this specialization, you will learn Keras.

- When you code in tensorflow you have to take the following steps:
 - Create a graph containing Tensors (Variables, Placeholders ...) and Operations (tf.matmul, tf.add, ...)
 - Create a session
 - Initialize the session
 - Run the session to execute the graph
- You can execute the graph multiple times as you've seen in model()
- The backpropagation and optimization is automatically done when running the session on the "optimizer" object.