

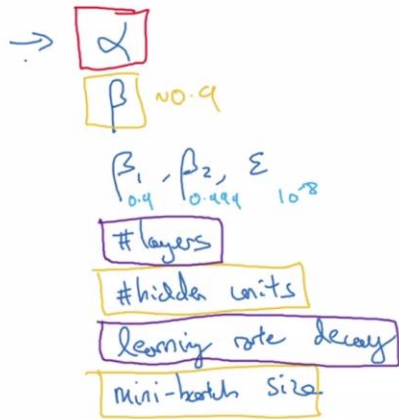
# Coursera Deep Learning Specialization Week 3

collaged by  
Kivanc Babacan

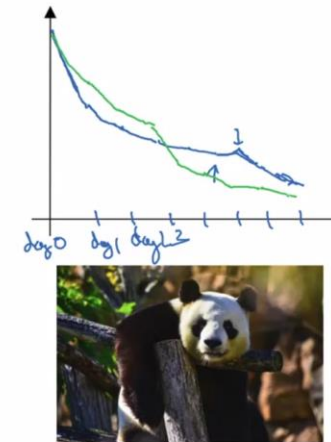
# Hyperparameter tuning, Batch Normalization and Programming Frameworks

1. Tuning process
2. Using an appropriate scale to pick hyperparameters
3. Hyperparameters tuning in practice: Pandas vs. Caviar
4. Normalizing activations in a network
5. Fitting Batch Norm into a neural network
6. Why does Batch Norm work?
7. Batch Norm at test time
8. Softmax Regression
9. Training a softmax classifier
10. Deep learning frameworks
11. TensorFlow

# Hyperparameters

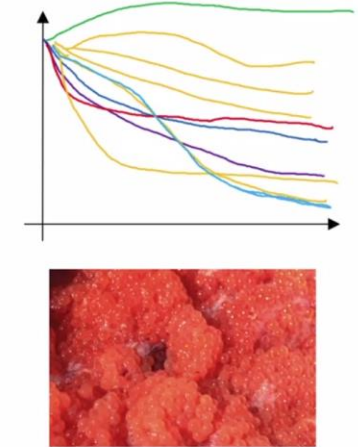


## Babysitting one model



Panda ←

## Training many models in parallel



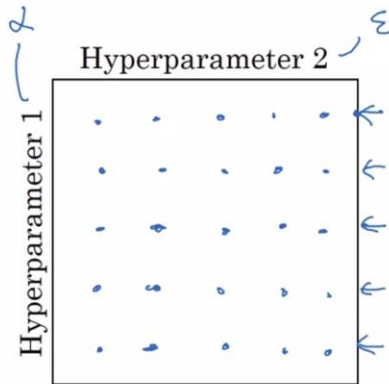
Caviar ←

Andrew Ng

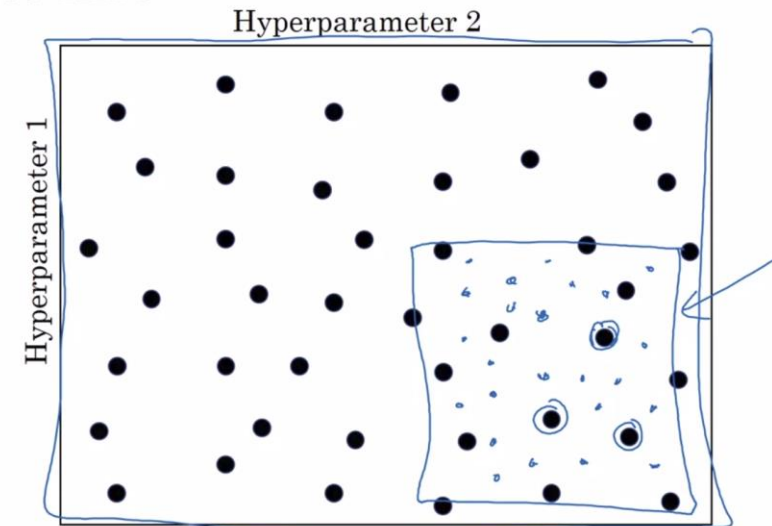
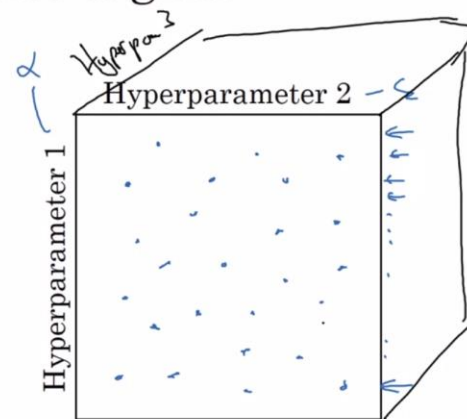
Andrew Ng

# HYPERPARAMETERS

Try random values: Don't use a grid



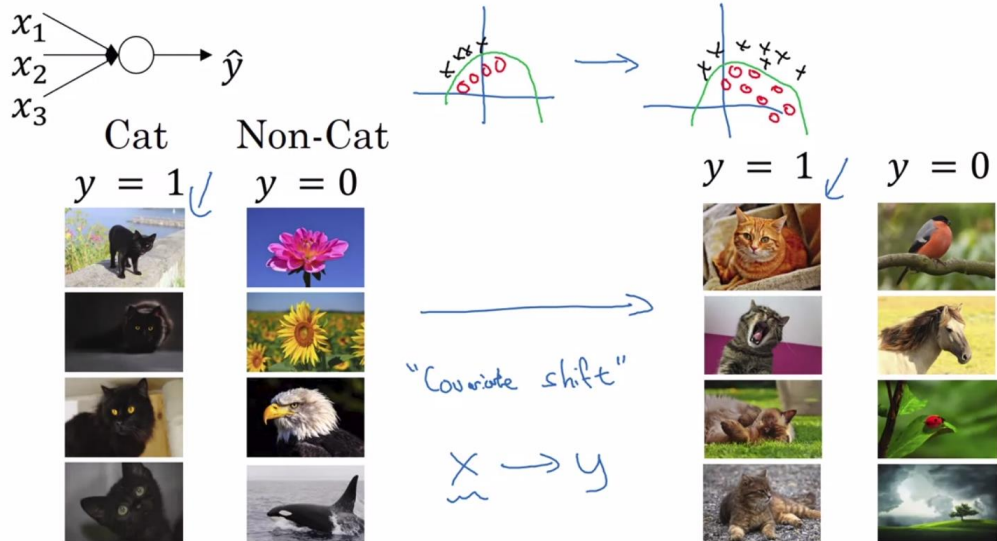
Coarse to fine



Andrew Ng

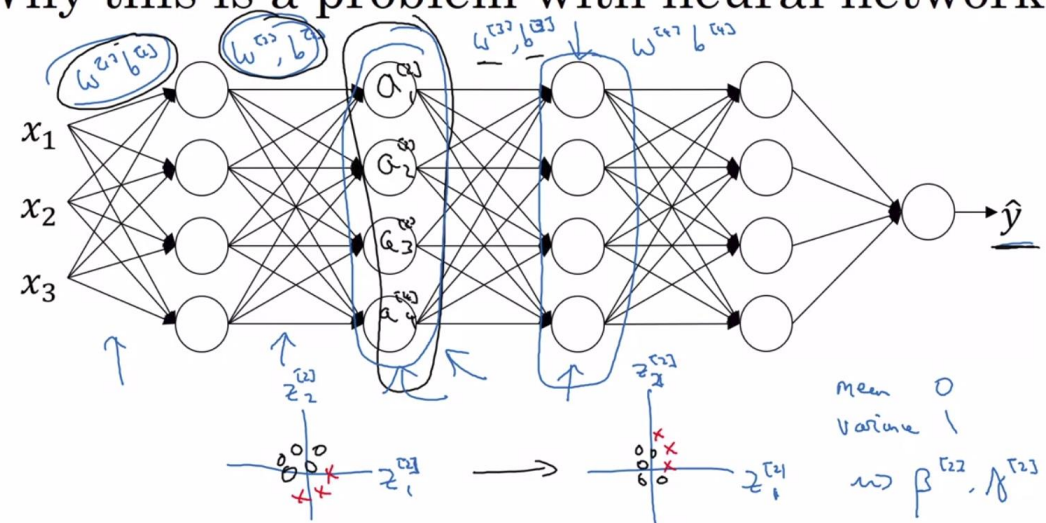
Andrew Ng

# Learning on shifting input distribution



Andrew Ng

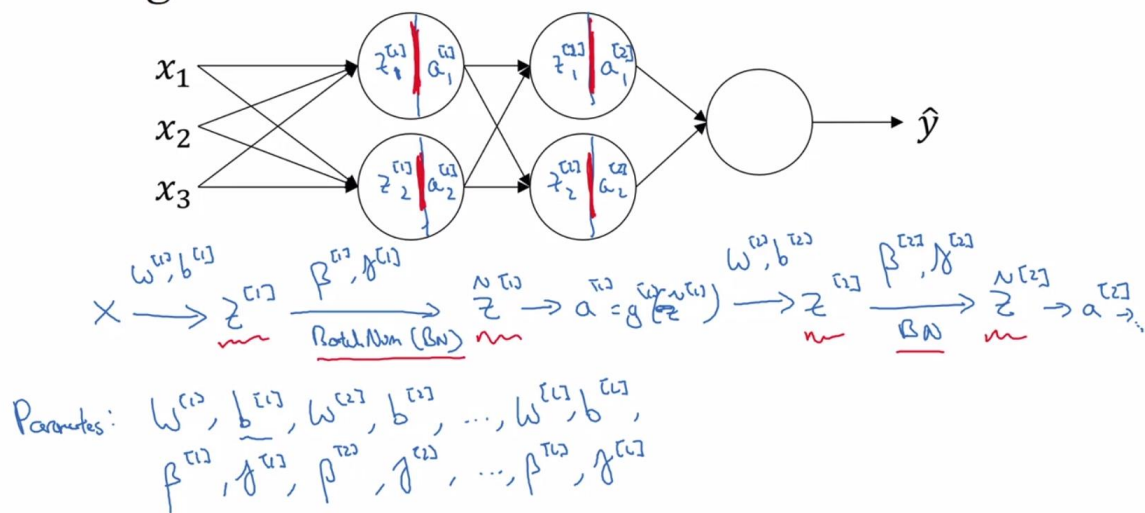
# Why this is a problem with neural networks?



Andrew Ng

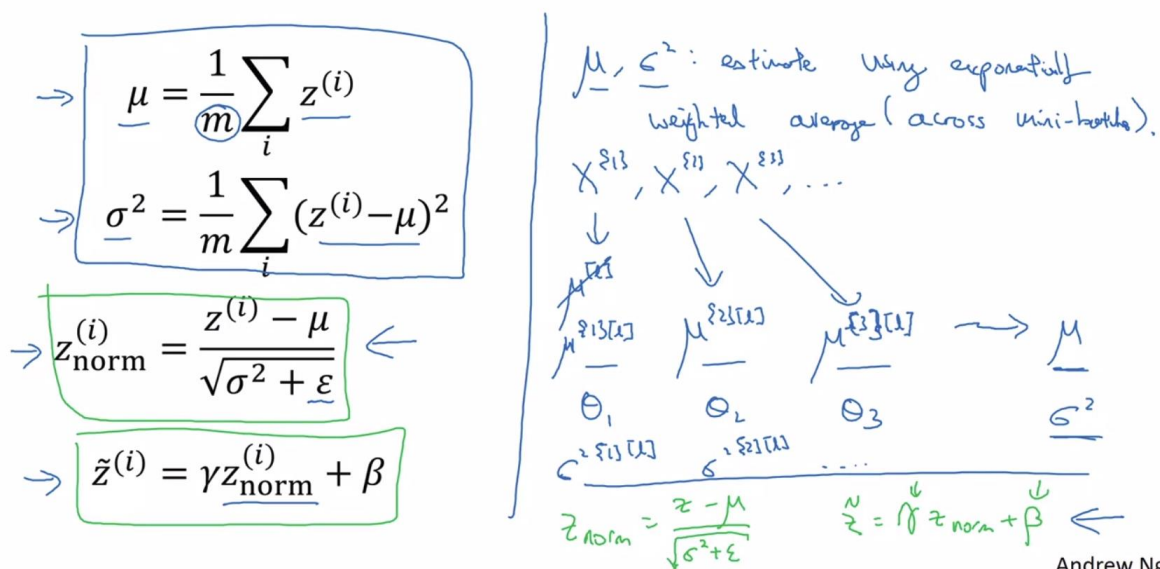
## BATCH NORMALIZATION

### Adding Batch Norm to a network



Andrew Ng

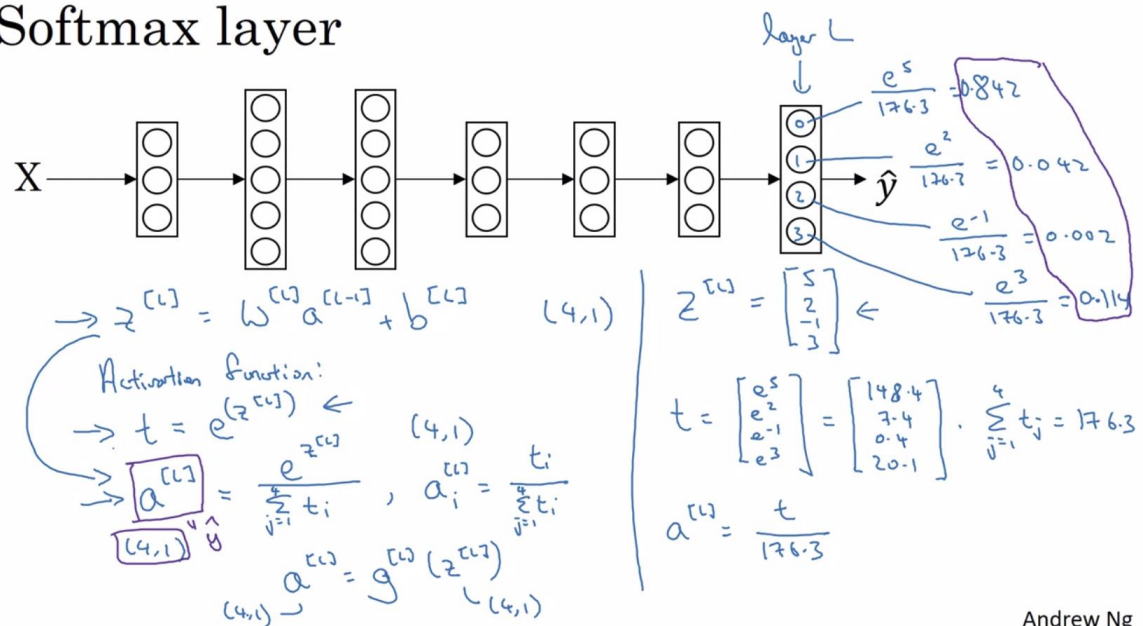
### Batch Norm at test time



Andrew Ng



## Softmax layer

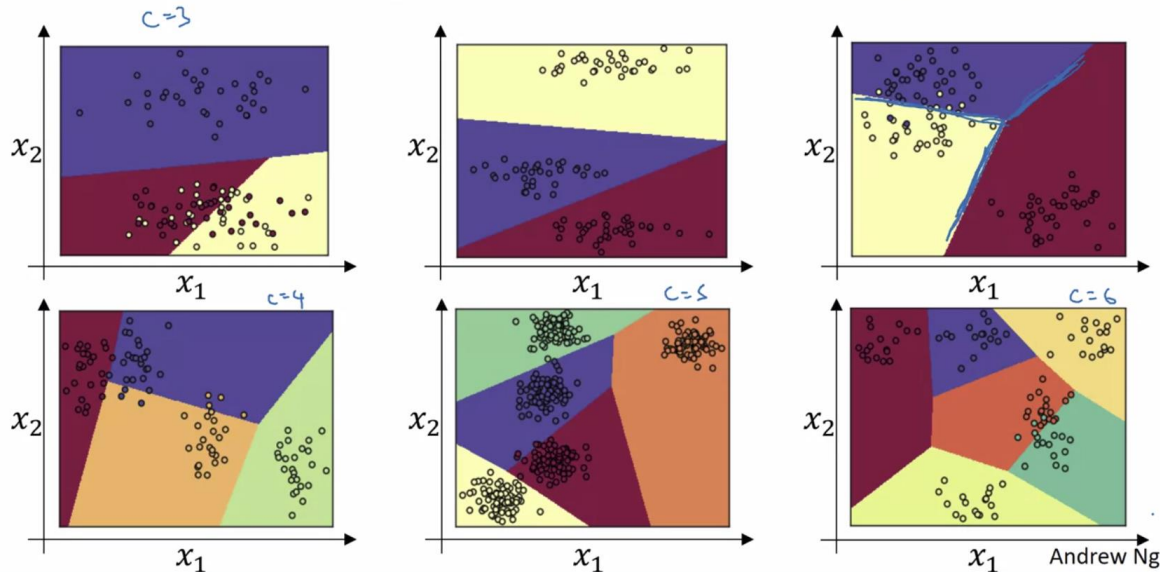


## Softmax examples

$$x_1 > x_2 \Rightarrow \hat{y} = 1$$

$$z^{(L)} = W^{(L)}x + b^{(L)}$$

$$a^{(L)} = \hat{y} = g(z^{(L)})$$



# SOFTMAX

## Loss function

Diagram of a Softmax layer with  $C=4$  classes. The input  $y$  is a vector of size  $(4,1)$ . The output  $\hat{y}$  is a vector of size  $(4,1)$ . The loss function is calculated as:

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

Example calculation for  $y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$  and  $\hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$ :

$$\mathcal{L}(\hat{y}, y) = -1 \log 0.3 = -\log 0.3$$

Make  $\hat{y}_2$  big.

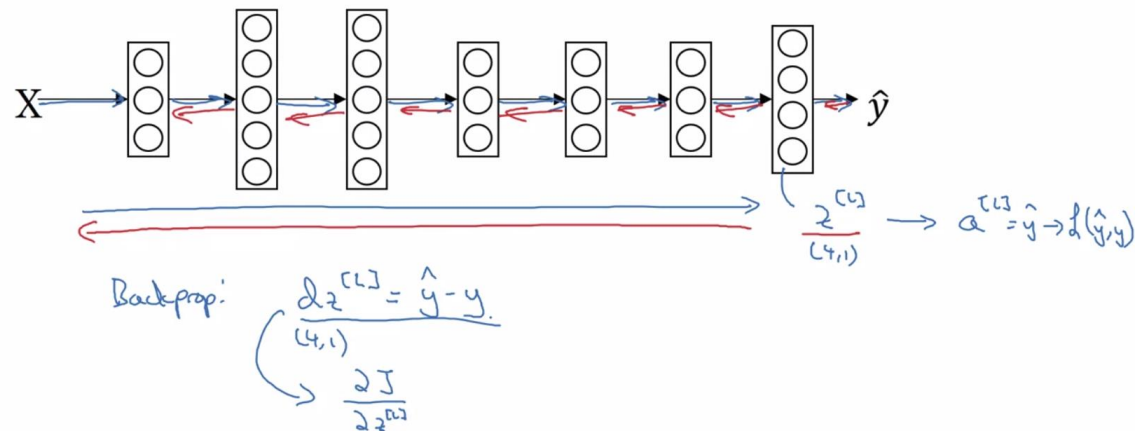
Overall loss function:

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

Input matrix  $Y = [y^{(1)} y^{(2)} \dots y^{(m)}]$  of size  $(4, m)$ .

Output matrix  $\hat{Y} = [\hat{y}^{(1)} \dots \hat{y}^{(m)}]$  of size  $(4, m)$ .

## Gradient descent with softmax



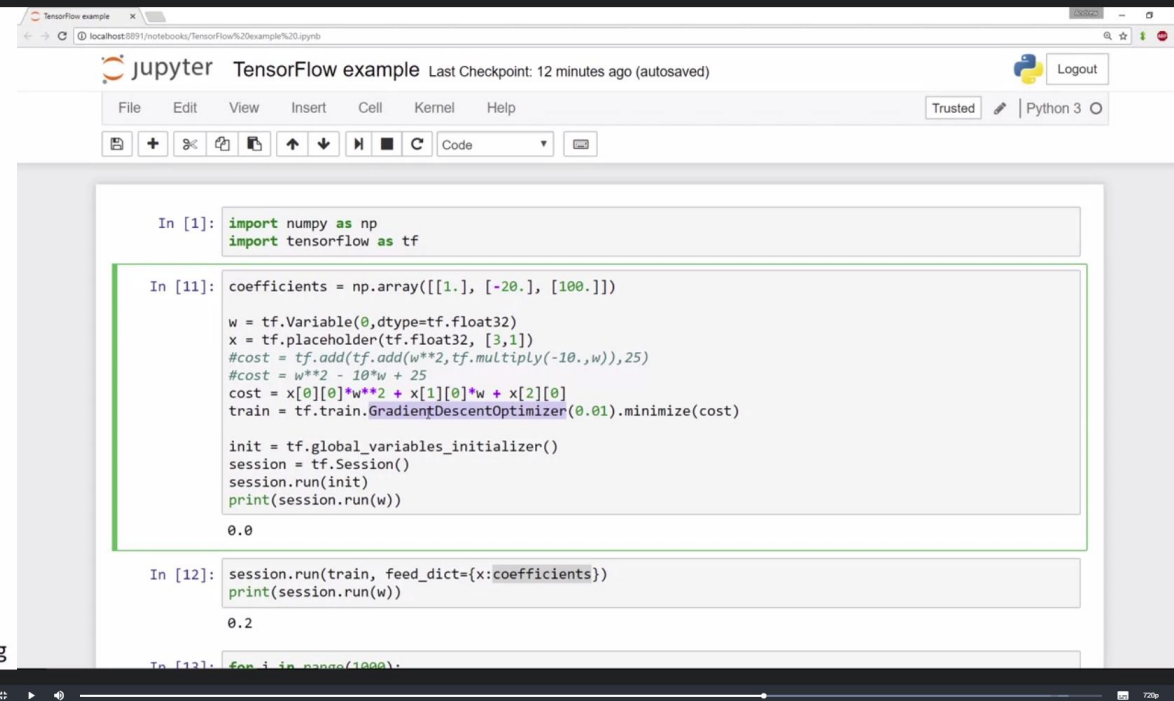
# Deep learning frameworks

- 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)

Andrew Ng



The screenshot shows a Jupyter Notebook interface with the title 'TensorFlow example'. The notebook contains three code cells. The first cell imports numpy and tensorflow. The second cell defines coefficients, creates a variable w, a placeholder x, and sets up the cost function and training process. The third cell runs the training process and prints the result.

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

In [11]: coefficients = np.array([[1.], [-20.], [100.]])
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 [12]: session.run(train, feed_dict={x: coefficients})
print(session.run(w))

0.2

In [13]: for i in range(1000):
```

# Take aways

- Random hyperparameter search instead of using grids, to cover more variant values for the more important parameter. Then refine coarse to fine.
- Just as normalizing your input, normalize every layer input to alleviate the covariance shift. Also adds small noise in the minibatch – free regularizer.
- Softmax ,a generalization of logistic regression to multiple classes linearizes decision boundaries.
- MatConvNet, the deep learning framework employed in GEOICT lab for the first is not counted. Use Tensorflow to match with fancy stuff.