

## **Course 1 Neural Networks and Deep Learning:**

### **Week 1**

#### What is a neural network?

Neural Network takes the input features, and we let the neural network decide what the hidden layers should be.

The input layer is densely connected to the hidden layer, every feature is connect to the hidden layer neurons.

Giving Neural network input  $x$  and output  $y$ , neural networks are really good at finding functions from  $x$  to  $y$

#### Types of Neural Networks

- Standard
- Image Recognition (Convolution Neural Networks)
- Sequence Data (Example Audio, translation, time series) (Recurrent NN )

Structured Data: Like A DataBase (Table)

UnStructured Data: Images, Audio, Text

More Data, Larger Neural Network -> better performance

Small training data, can't tell what size of neural net to use

For larger training data, bigger the size of neural net the better!

/

Algorithms: (FASTER) better to use a RELU function, stays at zero than a straight line \_\_\_\_\_/

RELU:  $h = \max(0, a)$  where  $a = Wx + b$

Reason: One major benefit is the reduced likelihood of the gradient to vanish. This arises when  $a > 0$ . In this regime the gradient has a constant value. In contrast, the gradient of

sigmoids becomes increasingly small as the absolute value of  $x$  increases. The constant gradient of ReLUs results in faster learning.

## Week 2

### Binary Classification (USE Logistic Regression)

To train neural nets the matrix  $X$  has each sample data down the column, where the number of columns represent the number of samples in the training set and the number of rows represent the number of features.

So if training set is  $(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$ , where each  $x_i$  is a vector of features for the input and  $y_i$  is the result (0,1)

$$X = [x_1 \ x_2 \ \dots \ x_n]$$

$$Y = [y_1 \ y_2 \ \dots \ y_n] \leftarrow \text{result matrix}$$

Error function the function to get loss between what the estimated probability and the true result. (Estimated probability is  $\hat{y}$ ) (Remember  $y$  can only be 0 or 1)(its a cat or not a cat)

- should only have one minimum so gradient descent doesn't get stuck on a local minimum, and should be convex.
- Logistic regression (loss) error function is  $L(\hat{y}, y) = -(y \log(\hat{y}) + (1-y)\log(1 - \hat{y}))$
- The equation explains that the closer  $\hat{y}$  is to  $y$  the smaller the error function.
- Remember this done on one sample, a specific training example.

The cost function is just the sum of the loss(error) function values over the number of samples in the training set  $J(w, b) = \frac{1}{m} (\sum_{i=1}^m L(\hat{y}_i, y_i))$

- Basically trying to find  $w$  and  $b$  to minimize the cost function

Gradient Descent (Take steps in the steepest path downward, by subtracting the slope)

- Fill parameters w and b with random values, (for logistic regression usually 0) and gradient descent, figures out the parameter values for where the cost function is minimized.
- Algorithm: Repeat {  $w = w - \alpha(\text{derivative of the cost function respect to } w)$
- Where alpha is the learning rate

### Computation Graph

Basically the cost function is made up of many variables  
 forward propagation step is just one operation in the calculation of the cost function. (Like a binary tree of operations and going down one node)  
 Backward propagation is used to find the derivatives respect to one of the variables in the cost function. (Chain Rule) (Need this derivative for the gradient descent algorithm)

### Logistic Regression Gradient Descent

To get the derivative value in the gradient descent formula:

Backpropagation, we get

$dL/dz = a-y$  where L is the loss functions and  $z = Wx + b$

To get the derivatives for all parameters  $w_1 \dots w_n$  since  $z = w_1x_1 + \dots + w_nx_n + b$

$dL/dw_i = (dL/dz)x_i$  which is just  $(a - y)x_i$  and  $dL/db = dL/dz = a - y$

So the gradient descent algorithm would just be:

(let  $d_{wi}$  represent  $dL/dw_i$  since we always finding derivatives regarding the loss function

$w_1 := w_1 - \alpha(d_{w1})$

$w_2 := w_2 - \alpha(d_{w2})$

,

where alpha is just the learning rate

,

$w_n := w_n - \alpha(d_{wn})$

$b := b - \alpha(db)$

## SUMMARY OF GRADIENT DESCENT ALGORITHM for LOGISTIC REGRESSION

For the cost function, since we have m samples in training set, just need to sum up all the derivative values and divide by m

For i = 1 to m

$$z_i = W^T(x_i) + b$$

$$a_i = \sigma(z_i)$$

$$J += y_i(\log(a_i)) + (1 - y_i)$$

$$d_{zi} = a_i - y_i$$

$$d_{w1} += (x_{1i})(d_{zi})$$

.

.

$$d_{wn} += (x_{ni})(d_{zi})$$

$$d_b += (d_{zi})$$

Done

$$J /= m$$

$$d_{w1} /= m; \dots d_{wn} /= m; d_b /= m$$

Then just adjust all  $w_i = w_i - \alpha(d_{wi})$

THAT IS ONE STEP IN GRADIENT DESCENT

# Logistic regression on $m$ examples

$$J=0; \underline{\Delta w_1}=0; \underline{\Delta w_2}=0; \underline{\Delta b}=0$$

For  $i = 1$  to  $m$

$$z^{(i)} = \omega^T x^{(i)} + b$$

$$\alpha^{(i)} = \sigma(z^{(i)})$$

$$J_t = -[y^{(i)} \log \alpha^{(i)} + (1-y^{(i)}) \log(1-\alpha^{(i)})]$$

$$\underline{\Delta z^{(i)}} = \alpha^{(i)} - y^{(i)}$$

$$\begin{aligned} \Delta w_1 &+= x_1^{(i)} \Delta z^{(i)} \\ \Delta w_2 &+= x_2^{(i)} \Delta z^{(i)} \end{aligned}$$

$$\Delta b += \Delta z^{(i)}$$

$$J/m \leftarrow$$

$$\Delta w_1/m; \Delta w_2/m; \Delta b/m. \leftarrow$$

$$\Delta w_1 = \frac{\partial J}{\partial w_1}$$

$$w_1 := w_1 - \alpha \underline{\Delta w_1}$$

$$w_2 := w_2 - \alpha \underline{\Delta w_2}$$

$$b := b - \alpha \underline{\Delta b}$$

Vectorization

## Vectorization

Basically getting rid of for loops for faster execution

Basically to do a dot product, (Example Av, or AB )

Numpy:

Every Time You need a for loop, see if there is a numpy build in function to the same thing .

# Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

$\rightarrow u = np.zeros((n, 1))$   
 $\boxed{\text{for } i \text{ in range}(n):} \leftarrow$   
 $\rightarrow u[i] = \text{math.exp}(v[i])$

import numpy as np  
u = np.exp(v) ←  
np.log(u)  
np.abs(u)  
np.maximum(v, 0)  
v\*\*2      1/u

---

Basically, can cut for loops by making a vector and then doing vectorized

## Vectorization on Logistic Regression

Broadcasting -> when a  $1 \times 1$  matrix like b acts like a  $1 \times m$  matrix when doing ...

$$Z = np.dot(w.T, X) + b$$

## Vectorized Implementation of Gradient Descent

- 1)  $Z = np.dot(w.T, X) + b$
- 2)  $A = \text{sigmoi}(Z)$
- 3)  $dZ = A - Y$     <<--- dZ is a  $1 \times m$  matrix for the derivative with respect to z for each Training example
- 4)  $db = 1/m \cdot \text{np.sum}(dZ)$
- 5)  $dw = 1/m \cdot (\text{np.dot}(X, dZ.T))$     <<--- shown above in a for loop, dw is a  $1 \times m$  matrix for all optimized parameters

## SUMMARY OF ONE STEP OF GRADIENT DESCENT

# Implementing Logistic Regression

```

 $J = 0, dw_1 = 0, dw_2 = 0, db = 0$ 
for i = 1 to m:
     $z^{(i)} = w^T x^{(i)} + b \leftarrow$ 
     $a^{(i)} = \sigma(z^{(i)}) \leftarrow$ 
     $J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$ 
     $dz^{(i)} = a^{(i)} - y^{(i)} \leftarrow$ 
     $\begin{cases} dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{cases}$ 
 $\left. \frac{\partial J}{\partial w} = \frac{1}{m} X^T dz^{(i)} \right\}$ 
 $J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$ 
 $db = db/m$ 

```

```

for iter in range(1000):
     $Z = w^T X + b$ 
     $= np.dot(w.T, X) + b$ 
     $A = \sigma(Z)$ 
     $dZ = A - Y$ 
     $dw = \frac{1}{m} X^T dZ$ 
     $db = \frac{1}{m} np.sum(dZ)$ 
     $w := w - \alpha dw$ 
     $b := b - \alpha db$ 

```

## Broadcasting

- 1) | 1 |  
| 2 | + 100    python automatically converts it too a 4 x 1 matrix of 100's  
| 3 |            add 100 to all elements  
| 4 |

## General Principle of BroadCasting

## General Principle

$$\begin{array}{ccc} (m, n) & \xrightarrow{\text{+}} & (1, n) \rightsquigarrow (m, n) \\ \underline{\text{matrix}} & \diagdown \times & (m, 1) \rightsquigarrow (m, n) \end{array}$$

$$\begin{array}{ccc} (m, 1) & + & \mathbb{R} \\ \left[ \begin{smallmatrix} 1 \\ 2 \\ 3 \end{smallmatrix} \right] & + & 100 = \left[ \begin{smallmatrix} 101 \\ 102 \\ 103 \end{smallmatrix} \right] \\ [1 \ 2 \ 3] & + & 100 = [101 \ 102 \ 103] \end{array}$$

Python: A.sum(axis = 0) sums axis vertically , if axis = 1 the horizontally

Python TIPS -> AVOID RANK 1 Array's and use col/row vectors

## Python/numpy vectors

```
a = np.random.randn(5)
      a.shape = (5, )
      "rank 1 array"
      } Don't use
```

```
a = np.random.randn(5, 1) → a.shape = (5, 1) column vector ✓
```

```
a = np.random.randn(1, 5) → a.shape = (1, 5) row vector. ✓
```

```
assert(a.shape == (5, 1)) ←
```

```
a = a.reshape((5, 1))
```

TO MAKE a IMAGE MATRIX Flatter -> ( So that every column vector repersents one image

[image1, image2 ..... Imagem]

A trick when you want to flatten a matrix X of shape (a,b,c,d) to a matrix X\_flatten of shape (b\*c\*d, a) is to use:

```
X_flatten = X.reshape(X.shape[0], -1).T      # X.T is the transpose of X
```

**Feature scaling** is a method used to standardize the range of independent variables or features of data

#### Motivation:

Since the range of values of raw data varies widely, in some machine learning algorithms, objective functions will not work properly without normalization. For example, the majority of classifiers calculate the distance between two points by the Euclidean distance. If one of the features has a broad range of values, the distance will be governed by this particular feature. Therefore, the range of all features should be normalized so that each feature contributes approximately proportionately to the final distance.

Another reason why feature scaling is applied is that gradient descent converges much faster with feature scaling than without it.[1]

## We use the standardization technique, **but for images better to just divide by 255**

To represent color images, the red, green and blue channels (RGB) must be specified for each pixel, and so the pixel value is actually a vector of three numbers ranging from 0 to 255.

One common preprocessing step in machine learning is to center and standardize your dataset, meaning that you subtract the mean of the whole numpy array from each example, and then divide each example by the standard deviation of the whole numpy array. But for picture datasets, it is simpler and more convenient and works almost as well to just divide every row of the dataset by 255 (the maximum value of a pixel channel).

Let's standardize our dataset.

**Overfitting:** Is extremely affected by the number of iteration for gradient descent, since the more you optimize the function based on the training data, the more noise you can pick up within the parameters.

- Also affected by the learning rate

## More complex Neural Networks

The superscripts [ ] represent a layer

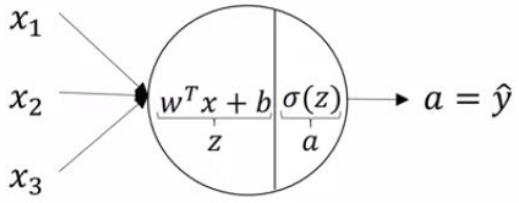
Taking logistic regression and repeating it

Hidden Layer: Don't see it in the training set

Number of Layers: Don't count the input layer

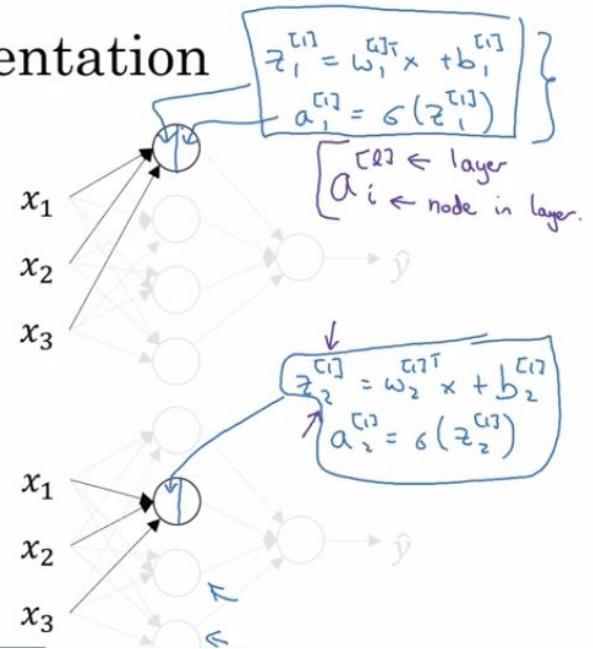
Hidden Layers and Output Layers have parameters associated with it

# Neural Network Representation



$$z = w^T x + b$$

$$a = \sigma(z)$$

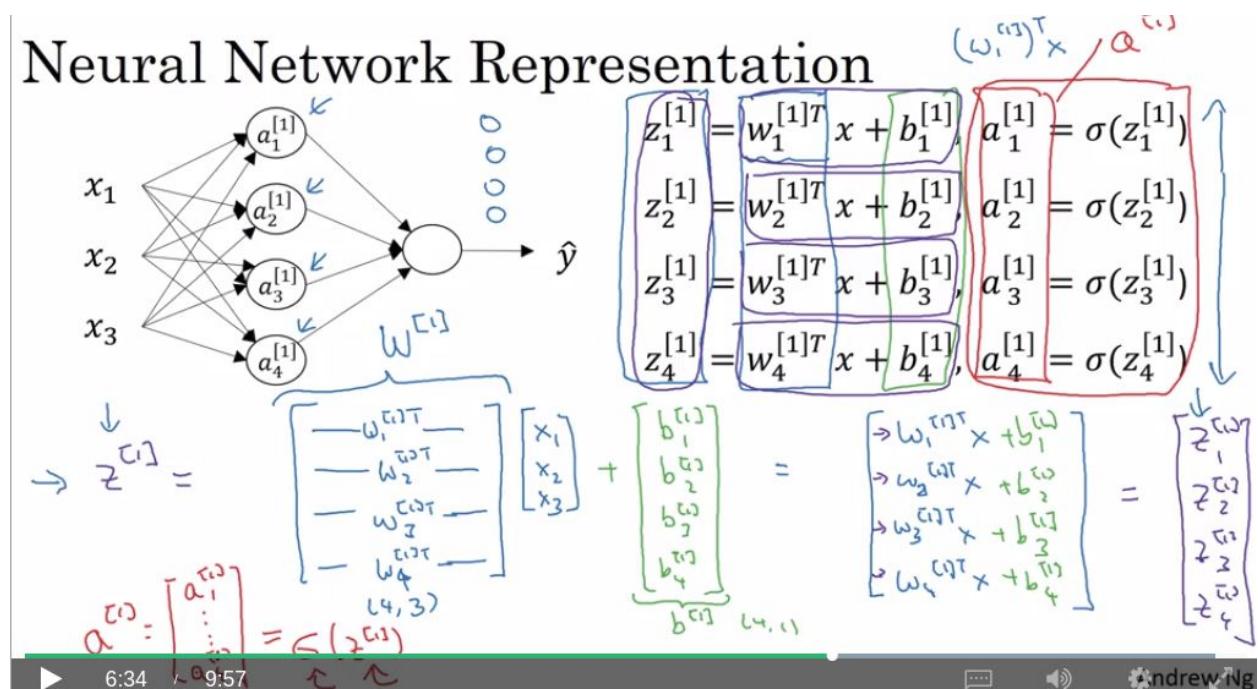


## Neural Networks dimension

Output Layer Dimension  $W$ : (hidden layer nodes, 1)

Hidden Layer Dimension  $W[1]$ : (hidden layer nodes, # of inputs)

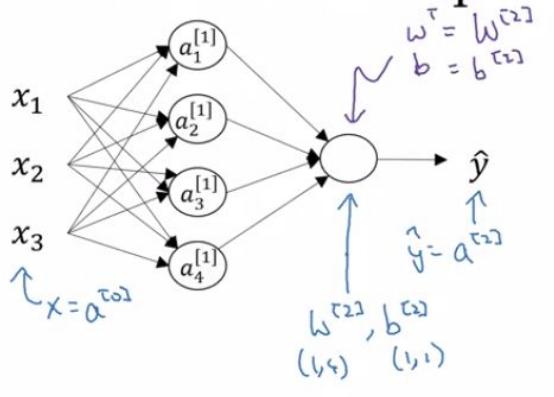
# Neural Network Representation



Summarizing the Notation:

$W^{[1]}$  basically now have multiple pairs of parameters where the number of pairs equals the nodes in the hidden layer and the number of parameters always equals the number of inputs. Thats how we get the  $W^{[1]}$  matrix

# Neural Network Representation learning



Given input  $x$ :

$$\rightarrow z^{[1]} = W^{[1]} \underset{(4,1)}{\alpha^{[0]}} + b^{[1]} \underset{(4,1)}{}$$

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

$$\rightarrow z^{[2]} = W^{[2]} a^{[1]} + b^{[2]} \underset{(1,1)}{(1,4)} \underset{(4,1)}{(4,1)} \underset{(1,1)}{}$$

$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \underset{(1,1)}{(1,1)}$$

## W[2] is a little bit confusing

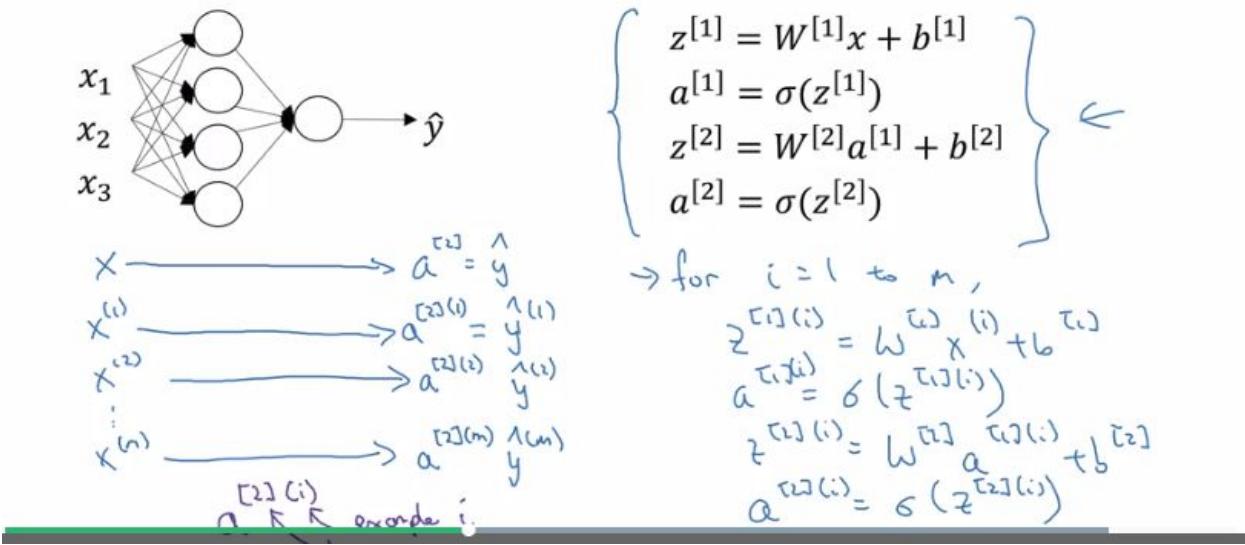
Basically  $W[2]$  is  $(1,4)$  since we only have one node in this layer and 4 inputs to this layer.

In the example  $z[2]$  and  $a[2]$  represent the Output layer

The four instructions above represent code for the above neural network

## Vectorization of Layered Neural Networks

# Vectorizing across multiple examples



All the capital letter Matrices -> each column represents a different training example

-> each row represents a different hidden unit / node / feature / parameter

## Justification for vectorized implementation

$$\begin{aligned}
 z^{(1)(1)} &= w^{(1)} x^{(1)} + b^{(1)}, & z^{(1)(2)} &= w^{(1)} x^{(2)} + b^{(1)}, & z^{(1)(3)} &= w^{(1)} x^{(3)} + b^{(1)} \\
 w^{(1)} &= \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}, & w^{(1)} x^{(1)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(2)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}, & w^{(1)} x^{(3)} &= \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \\
 z^{(1)} &= \begin{bmatrix} 1 & 1 & 1 \\ x^{(1)} & x^{(2)} & x^{(3)} \\ 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & z^{(1)(3)} \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = z^{(1)}
 \end{aligned}$$

### Activation functions

Example: Sigmoid

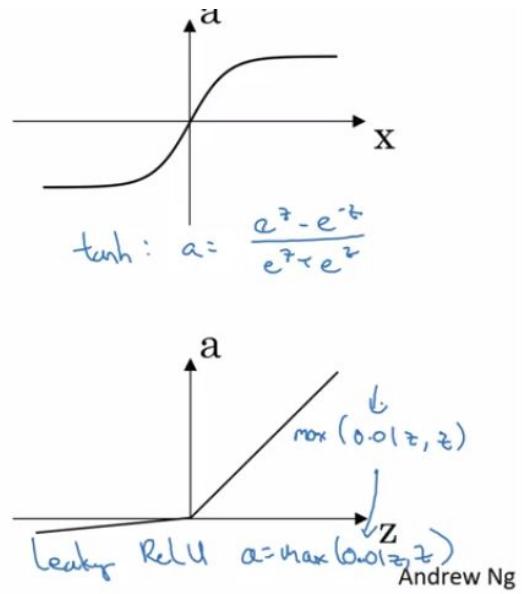
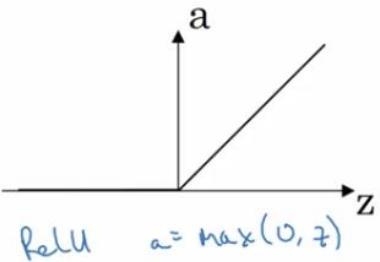
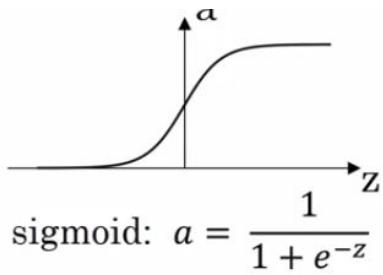
Tanh function almost always works better than the sigmoid function

Makes sure that the data is centered and the mean is 0 makes it easier for the next layer

Exception: Use the sigmoid function if you're doing a binary classification problem only for the output layer, since  $y \rightarrow \{0, 1\}$  so you want  $0 \leq y_{\text{hat}} \leq 1$  ( $y_{\text{hat}}$  is basically col of A) If

If  $Z$  is very large or small, in both tanh and sigmoid function the gradient becomes very slow since the slope decreases (look at visual below)

USE RELU ACTIVATION!!! (Neural Network work will work faster, Slope is greater)



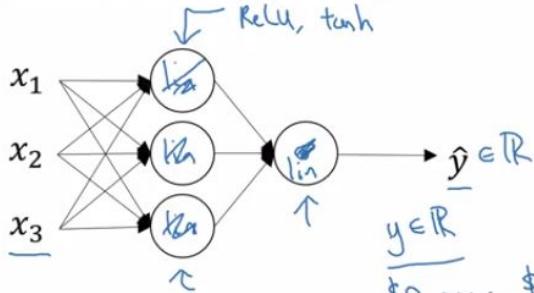
### Why do we need a activation function

Because if we don't use a activation function we get a linear function.

A Linear hidden layer is useless (a layer that has no activation function)

Because compositions of linear functions are linear, not really doing anything interesting

# Activation function



Given  $x$ :

$$\begin{aligned}\rightarrow z^{[1]} &= W^{[1]}x + b^{[1]} \\ \rightarrow a^{[1]} &= g^{[1]}(z^{[1]}) \quad z^{[1]} \\ \rightarrow z^{[2]} &= W^{[2]}a^{[1]} + b^{[2]} \\ \rightarrow a^{[2]} &= g^{[2]}(z^{[2]}) \quad z^{[2]}\end{aligned}$$

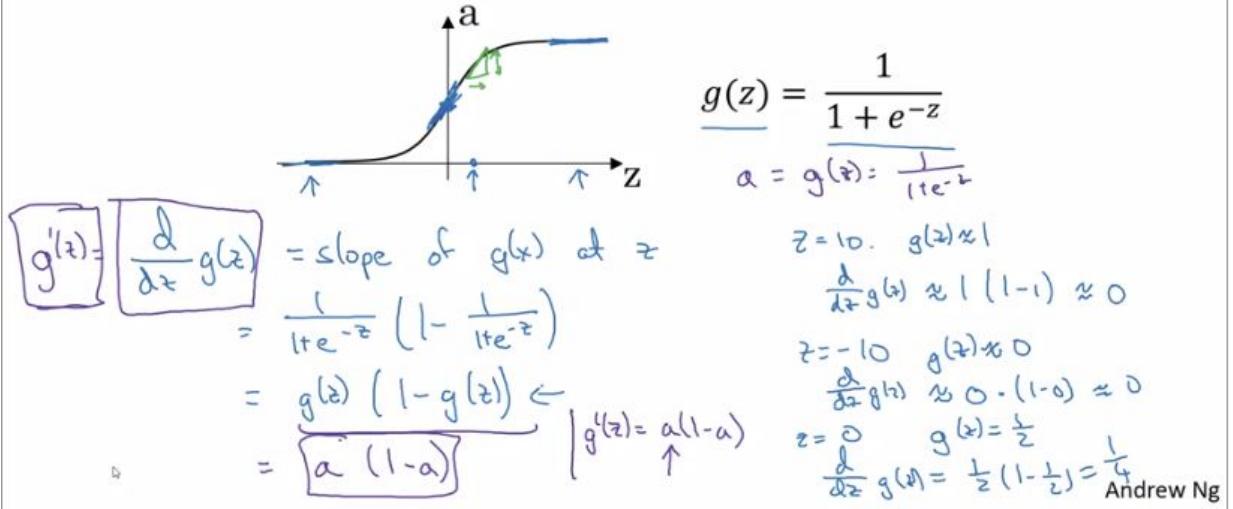
$g(z) = z$   
"linear activation  
function"

$$\begin{aligned}a^{[1]} &= z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[2]} &= z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} &= W^{[2]}(\underbrace{W^{[1]}x + b^{[1]}}_{a^{[1]}}) + b^{[2]} \\ &= (\underbrace{W^{[2]}W^{[1]}}_{W'})x + (\underbrace{W^{[2]}b^{[1]} + b^{[2]}}_{b'}) \\ &= W'x + b'\end{aligned}$$

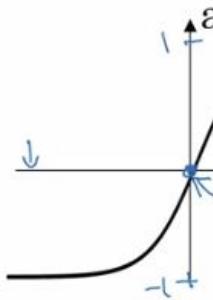
Andrew Ng

A place where you could use linear activation functions (aka no activation function) is in the output layer if the output you want for example is a real number from negative .... to a million dollars

# Sigmoid activation function



## Tanh activation function



$$g(z) = \tanh(z)$$

$$= \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = \frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z$$

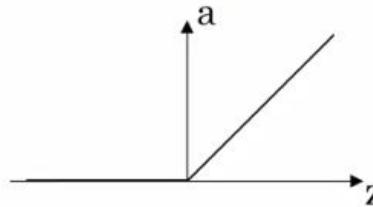
$$= 1 - (\tanh(z))^2 \leftarrow$$

$a = g(z), \quad g'(z) = 1 - a^2$

$$\begin{cases} z = 10 & \tanh(z) \approx 1 \\ g'(z) \approx 0 \\ z = -10 & \tanh(z) \approx -1 \\ g'(z) \approx 0 \\ z = 0 & \tanh(z) = 0 \\ g'(z) = 1 \end{cases}$$

Andrew Ng

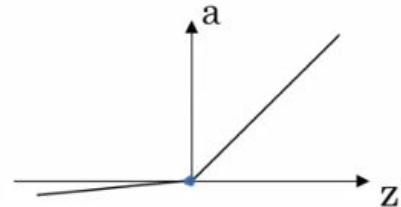
# ReLU and Leaky ReLU



ReLU

$$g(z) = \max(0, z)$$
$$\rightarrow g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

~~unlike if  $z=0$~~   
 $z=0.0000\cdots 0$



Leaky ReLU

$$g(z) = \max(0.01z, z)$$
$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

Andrew Ng

## Gradient Descent for Neural Networks

dimension :

$N[0]$  is input

$N[1]$  is hidden layer

$N[2]$  is output

Summary if how gradient descent will work

# Gradient descent for neural networks

Parameters:  $w^{[0]}, b^{[0]}$ ,  $w^{[1]}, b^{[1]}$ ,  $\dots$ ,  $w^{[L]}, b^{[L]}$   
 $n_x = n^{[0]}, n^{[1]}, \dots, n^{[L]} = 1$

Cost function:  $J(w^{[0]}, b^{[0]}, w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$

Gradient Descent:

→ Repeat {  
    → Compute predicts  $(\hat{y}^{(i)}, i=1 \dots m)$   
     $\frac{\partial J}{\partial w^{[0]}} = \frac{\partial J}{\partial w^{[0]}}$ ,  $\frac{\partial J}{\partial b^{[0]}} = \frac{\partial J}{\partial b^{[0]}}$ , ...  
     $w^{[0]} := w^{[0]} - \alpha \frac{\partial J}{\partial w^{[0]}}$   
     $b^{[0]} := b^{[0]} - \alpha \frac{\partial J}{\partial b^{[0]}}$   
    ...  
}

Derivatives For Gradient Descent

# Formulas for computing derivatives

Forward propagation:

$$z^{(1)} = w^{(1)}X + b^{(1)}$$

$$A^{(1)} = g^{(1)}(z^{(1)}) \leftarrow$$

$$z^{(2)} = w^{(2)}A^{(1)} + b^{(2)}$$

$$A^{(2)} = g^{(2)}(z^{(2)}) = \underline{c}(z^{(2)})$$

Back propagation:

$$dz^{(2)} = A^{(2)} - Y \leftarrow$$

$$dW^{(2)} = \frac{1}{m} dz^{(2)} A^{(1)T}$$

$$db^{(2)} = \frac{1}{m} \text{np.sum}(dz^{(2)}, \text{axis}=1, \text{keepdims=True})$$

$$dz^{(1)} = \underbrace{(w^{(1)T} dz^{(2)} \times g^{(1)'}(z^{(1)}))}_{\text{element-wise product}} \quad (n^{(2)}, m)$$

$$dW^{(1)} = \frac{1}{m} dz^{(1)} X^T$$

$$db^{(1)} = \underbrace{\frac{1}{m} \text{np.sum}(dz^{(1)}, \text{axis}=1, \text{keepdims=True})}_{(n^{(1)}, 1)} \quad \text{reshape} \uparrow$$

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

$$(n^{(0)}) \leftarrow$$

$$\downarrow (n^{(1)}, 1) \leftarrow$$

Andrew Ng

## Random Initialization

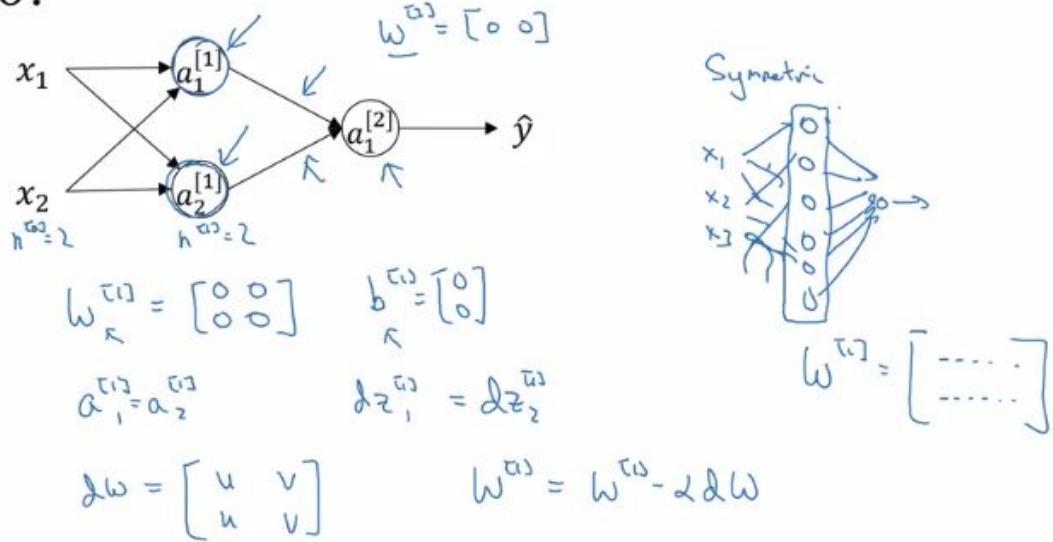
For Logistic Regression okay to start with values of zero, but for Neural Networks must randomize the w and b

**Explanation For Diagram Below:** Both Units in the Hidden Layer compute the same function, so by symmetry  $dz[1]1 = dz[1]2$ , and each row of dw (aka  $dL/dw$ ) end up being identical, so when you do

$W[1] = W[1] - \text{learning\_rate}(dw) W[1]$  gets identical rows

YOU WANT DIFFERENT HIDDEN UNITS TO COMPUTE DIFFERENT FUNCTIONS THUS DIFFERENT PARAMETERS

# What happens if you initialize weights to zero?

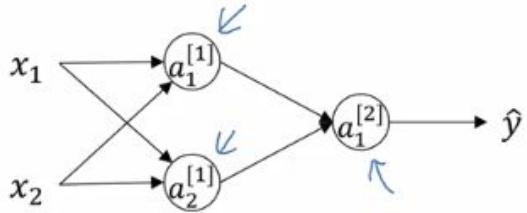


Andrew Ng

SO MUST INITIALIZE W RANDOMLY, b doesn't matter since W is already random, so we get random number finally

WE LIKE THE RANDOM INITIALIZATION of numbers to be very small so that's why we divide by 100. This is because if it's very large then for tanh and sigmoid function you end up at the edges where slope is very little

# Random initialization

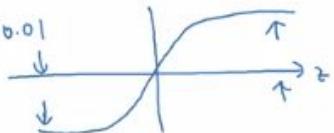


$$\rightarrow \omega^{[1]} = \text{np.random.randn}(1, 2) * \frac{0.01}{100?} \quad \rightarrow z^{[1]} = \omega^{[1]} x + b^{[1]}$$

$$b^{[1]} = \text{np.zeros}(1, 1)$$

$$\omega^{[2]} = \text{np.random.randn}(1, 2) * 0.01 \quad \downarrow \quad \uparrow$$

$$b^{[2]} = 0 \quad \quad \quad z \quad \quad \quad a^{[2]} = g^{[2]}(z^{[2]})$$



Andrew Ng

FOR DEEP NEURAL NETWORKS might want to multiply initialization of W with different constant.

## SUMMARY OF DERIVATIVES NEURAL NETWORK

## Summary of gradient descent

$dz^{[2]} = a^{[2]} - y$	$dZ^{[2]} = A^{[2]} - Y$
$dW^{[2]} = dz^{[2]} a^{[1]T}$	$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$
$db^{[2]} = dz^{[2]}$	$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis = 1, keepdims = True)$
$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$	$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$
$dW^{[1]} = dz^{[1]} X^T$	$dW^{[1]} = \frac{1}{m} dZ^{[1]} X^T$
$db^{[1]} = dz^{[1]}$	$db^{[1]} = \frac{1}{m} np.sum(dZ^{[1]}, axis = 1, keepdims = True)$

Andrew Ng

Accuracy is really high compared to Logistic Regression. The model has learnt the leaf patterns of the flower! Neural networks are able to learn even highly non-linear decision boundaries, unlike logistic regression.

## Deep Neural Networks

The output layer is denoted as L, and the activation if layer L

$a[L]$  is  $y_{\text{hat}}$ , it is the predicted output

## Forward Propagation

$$\boxed{\begin{aligned} Z^{[l]} &= W^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(Z^{[l]}) \end{aligned}}$$

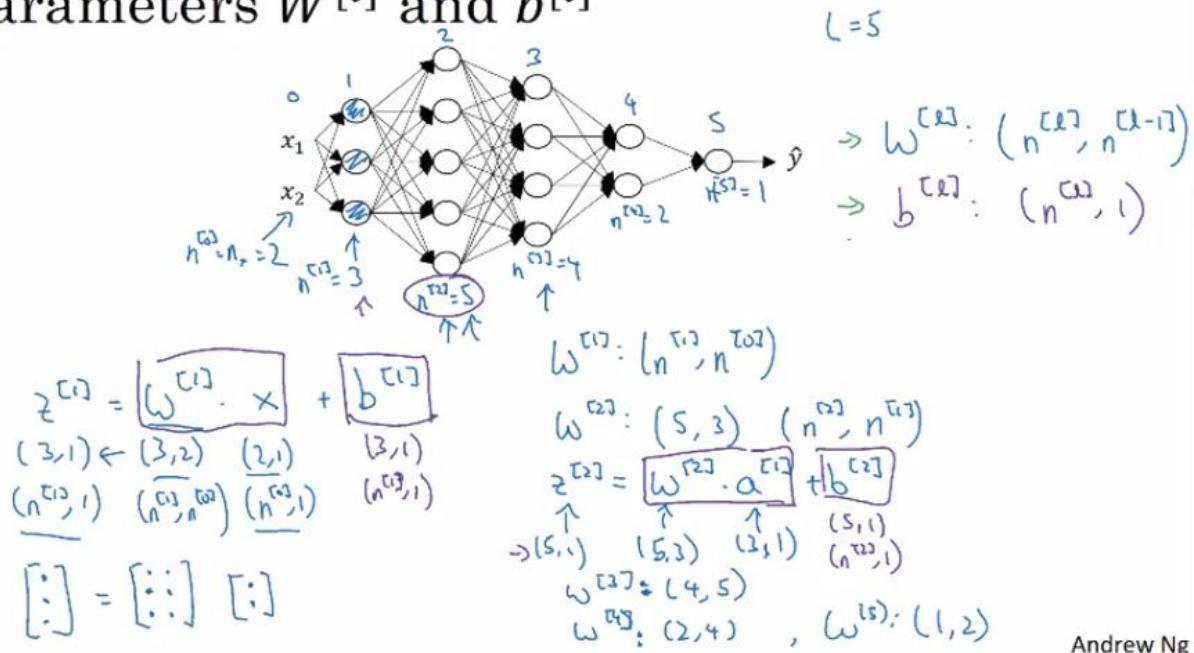
Okay to have a for loop

Sizes of W parameters:

$W[L]: \{n[L] \times n[L-1]\}$

$b[L]: \{n[L] \times 1\}$

## Parameters $W^{[l]}$ and $b^{[l]}$



Andrew Ng

If making a neural network more shallow, you might need to add more nodes in the hidden layer to make up for the logical thought process.

Cache in implementation is used to store the parameters  $W_i$ ,  $Z_i$ ,  $b_i$  so you can use it in backpropagation to calculate derivatives and calculate new  $W_i$  and  $b_i$ 's.

Back Propagation Step Summary:

## Backward propagation for layer $l$

→ Input  $\underline{da^{[l]}}$

→ Output  $\underline{da^{[l-1]}}, \underline{dW^{[l]}}, \underline{db^{[l]}}$

$$\begin{cases} \underline{dz^{[l]}} = \underline{da^{[l]}} * g^{[l]}(z^{[l]}) \\ \underline{dW^{[l]}} = \underline{dz^{[l]}} \cdot \underline{a^{[l-1]}} \\ \underline{db^{[l]}} = \underline{dz^{[l]}} \\ \underline{da^{[l-1]}} = \underline{w^{[l]T} \cdot dz^{[l]}} \\ \underline{dz^{[l-1]}} = \underline{w^{[l-1]T} \cdot dz^{[l]}} + g^{[l-1]}(z^{[l]}) \end{cases}$$

$$dz^{[l]} = dA^{[l]} * g^{[l]}(z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dZ^{[l]}, \text{axis}=1, \text{keepdim=True})$$

$$dA^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

Andrew Ng

Hyperparameters

Are things that affect the parameters

## What are hyperparameters?

Parameters:  $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}, W^{[3]}, b^{[3]} \dots$

Hyperparameters: learning rate  $\alpha$   
#iterations  
#hidden layers  $L$   
#hidden units  $n^{[1]}, n^{[2]}, \dots$   
choice of activation function

Others: Momentum, mini-batch size, regularizations, ...

Try multiple values for hyper parameters

Back Propagation In Practise CODE

```

def linear_backward(dZ, cache):
    """
    Implement the linear portion of backward propagation for a single layer (layer l)

    Arguments:
    dZ -- Gradient of the cost with respect to the linear output (of current layer l)
    cache -- tuple of values (A_prev, W, b) coming from the forward propagation in the current layer

    Returns:
    dA_prev -- Gradient of the cost with respect to the activation (of the previous layer l-1), same shape as A_prev
    dW -- Gradient of the cost with respect to W (current layer l), same shape as W
    db -- Gradient of the cost with respect to b (current layer l), same shape as b
    """
    A_prev, W, b = cache
    m = A_prev.shape[1]

    ### START CODE HERE ### (≈ 3 lines of code)
    dW = 1/m * np.dot(dZ, A_prev.T)
    db = 1/m * np.sum(dZ, axis = 1, keepdims = True)
    dA_prev = np.dot(W.T,dZ)
    ### END CODE HERE ###

```

## Course 2: Improving Deep Neural Networks, Hyper Parameter Tuning, Regularization and Optimization

### Hyperparameters

# Layers

# Hidden Units

# learning Rates

# Activation Functions

Very iterative project idea, experiment, code

### Data training/dev/test

Dev set is to try and tune the model

Test set is only required for unbiased estimator

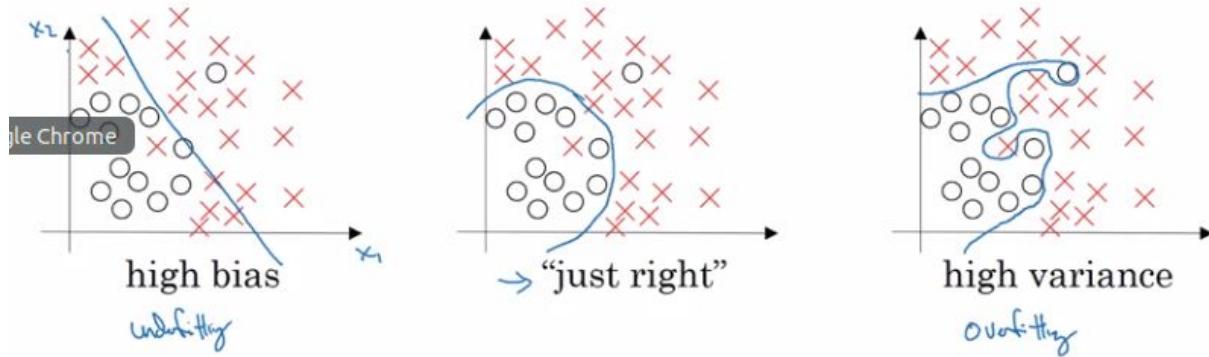
1000, 10000 -> should do 70/30 60/20/20

Big Data -> millions just pick 99/1 98/1/1 (like 10000 tests examples)

Make sure dev and test sets come from the same distribution

### Bias/Variance

## Bias and Variance



High Variance -> Overfitting

High Bias -> Underfitting

# Bias and Variance

Cat classification



<u>Train set error:</u>	1%	15% ↗	15%	0.5%
<u>Dev set error:</u>	11%	16% ↗	30%	1%
	high variance	high bias	high bias & high variance	low bias low variance
<u>Human:</u>	20%			
Optimal (Bayes) error:	15%	Blurry images		

Optimal is how much error in terms of human prediction, the bias and variance above is for optimal of 0%, if it was 15% #2 col would be low bias and variance

## Recipe To Improve Error Bias/Variance

### Training Data

- High Bias      -> Bigger Network (More Layers) (more computational Time)
- > Train Longer
- > Different Architecture

### Dev Data

- High Variance   -> More Data

-> Regularization

-> Architecture

Improvement has been shown with just adding more layers and including more dev data, reduces both bias and variance

### Regularization (To reduce overfitting and variance)

Use L2 Norm for logistic regression

## Logistic regression

$$\min_{w,b} J(w,b) \quad w \in \mathbb{R}^{n_x}, b \in \mathbb{R} \quad \lambda = \text{regularization parameter}$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m l(y_i, \hat{y}_i) + \frac{\lambda}{2m} \|w\|_2^2$$



$$\text{L}_2 \text{ regularization} \quad \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$$

$$\text{L}_1 \text{ regularization} \quad \frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1 \quad w \text{ will be sparse}$$

Using L2 / Weight Decay Regularization Technique to make values in W a bit smaller

# Neural network

$$J(w^{(0)}, b^{(0)}, \dots, w^{(L)}, b^{(L)}) = \underbrace{\frac{1}{m} \sum_{i=1}^m f(\hat{y}^{(i)}, y^{(i)})}_{\text{"Frobenius norm"}} + \underbrace{\frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2}_{\text{regularization}}$$

$$\|w^{(l)}\|_F^2 = \sum_{i=1}^{n^{(l+1)}} \sum_{j=1}^{n^{(l)}} (w_{ij}^{(l)})^2$$

$w: (n^{(L+1)} \times n^{(L)})$

"Frobenius norm"       $\|\cdot\|_2^2$        $\|\cdot\|_F^2$

$$\begin{aligned} d\omega^{(l)} &= (\text{from backprop}) + \frac{\lambda}{m} w^{(l)} \\ \rightarrow \omega^{(l)} &:= \omega^{(l)} - \alpha d\omega^{(l)} \end{aligned}$$

$$\frac{\partial J}{\partial \omega^{(l)}} = \alpha d\omega^{(l)}$$

"Weight decay"

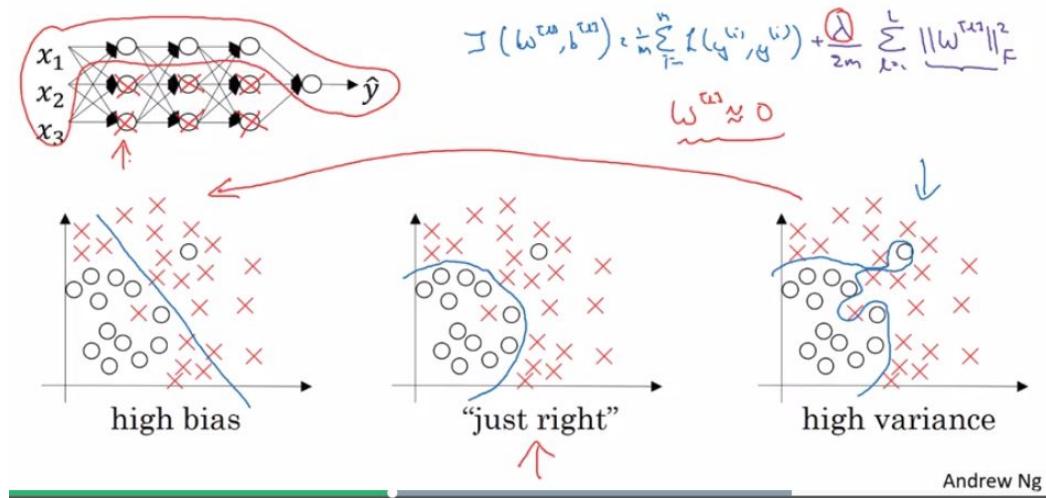
$$\begin{aligned} \underline{\omega^{(l)}} &:= \omega^{(l)} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} w^{(l)} \right] \\ &= \underline{\omega^{(l)}} - \underline{\left( \frac{\alpha \lambda}{m} \right) \omega^{(l)}} - \underline{\alpha (\text{from Backprop})} \end{aligned}$$

Andrew Ng

## How does regularization minimize overfitting?

Hidden Units become very close to zero, so it becomes a simple neural network and follows the arrow to being less variance , but increasing bias, there must be an intermediate values for lambda thats just right.

## How does regularization prevent overfitting?

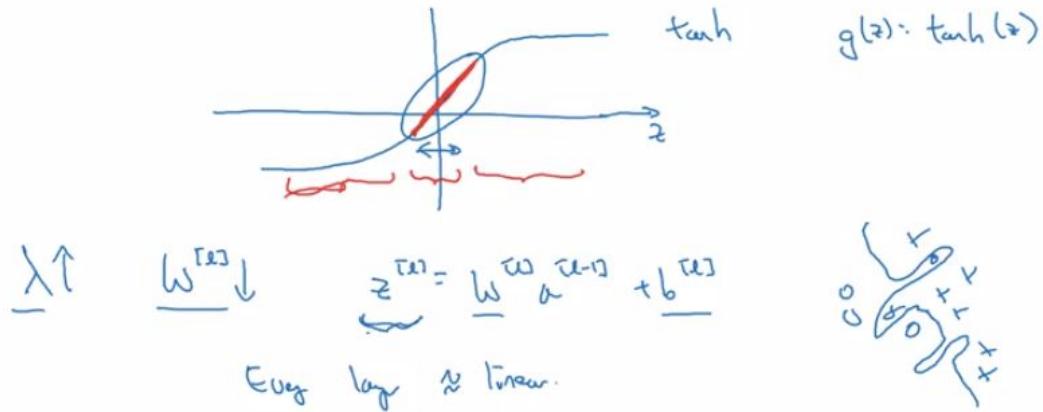


A more intuitive reason to why regularization, prevents overfitting

Lambda is bigger  $\rightarrow$  W is smaller  $\rightarrow$  Z has smaller values so in the graph below shows that  $g(Z)$  can be only linear. So if every layer is linear, then your whole network is only linear.

Can't get some crazy function like the one shown below to fit the data.

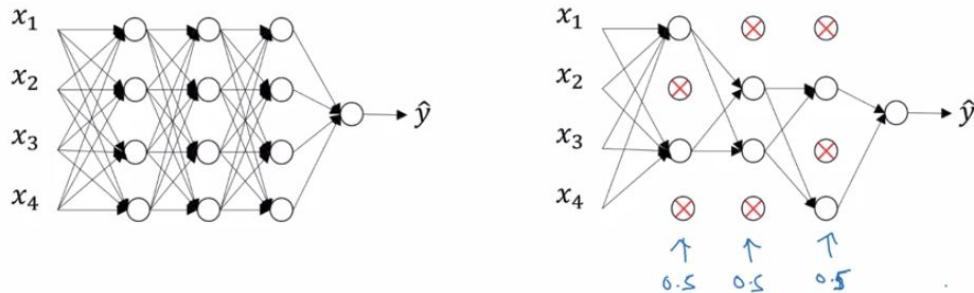
# How does regularization prevent overfitting?



## Dropout Regularization

Basically each layer has a probability of dropping some of its hidden units, (not using them), taking the links away

## Dropout regularization



## Implementing Dropout

### Implementing dropout ("Inverted dropout")

Illustrate with layer  $l=3$ .  $\text{keep-prob} = \frac{0.8}{x}$   $\underline{0.2}$

$$\rightarrow d_3 = \underbrace{\text{np.random.rand}(a_3.shape[0], a_3.shape[1]) < \text{keep-prob}}$$

$$a_3 = \underbrace{\text{np.multiply}(a_3, d_3)}_{\# a_3 \neq d_3}$$

$$\rightarrow a_3 /= \cancel{\text{keep-prob}} \leftarrow$$

50 units  $\rightsquigarrow$  10 units shut off

$$z^{(4)} = w^{(4)} \cdot \underbrace{\frac{a^{(3)}}{x}}_{\text{reduced by } 20\%} + b^{(4)}$$

$$/ = \underline{0.8}$$

Test

Basically On Every Iteration of Gradient Descent You Drop Different Hidden Units Randomly.

During Testing you shouldn't drop any hidden units!!!!

Can vary Keep\_Prob for each hidden layer, Depending on which layer you feel there is overfitting

Make Sure gradient descent is working without dropout, decreasing monotonically, then try dropout

## Why Dropout Actually Works

Can't rely on any single feature, this is how overfitting happens, Spread out the weights for each features

Dropout is used greatly in COMPUTER VISION, since you don't have a lot of data and way to many pixels

## Other ways of decreasing overfitting

Increasing Data:

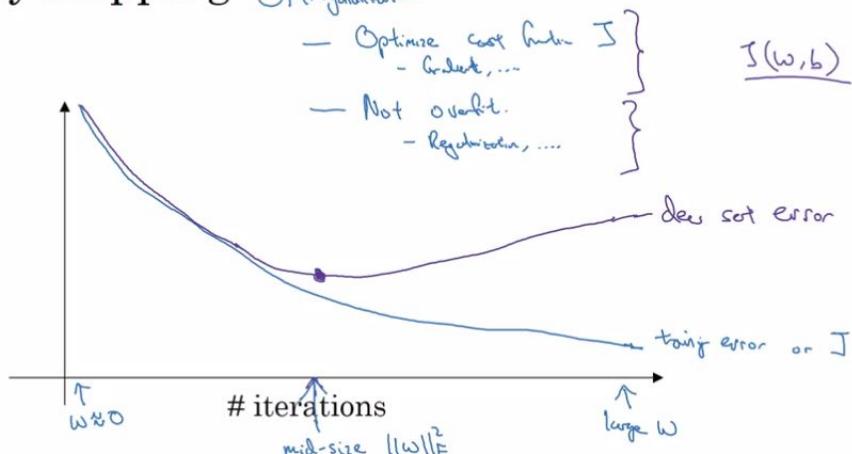
Can increase data for computer vision by rotating cat images, and flipping them, , zooming double, triple your data.

Early Stopping:

Stopping gradient descent early since W is smaller usually at the beginning, when W is smaller it minimizes overfitting

Disadvantage: Cannot Optimize the cost function J, should deal with optimization and overfitting separately

## Early stopping (Orthogonalization)



Disadvantage to L2 is the amount of values you have to test for lamda

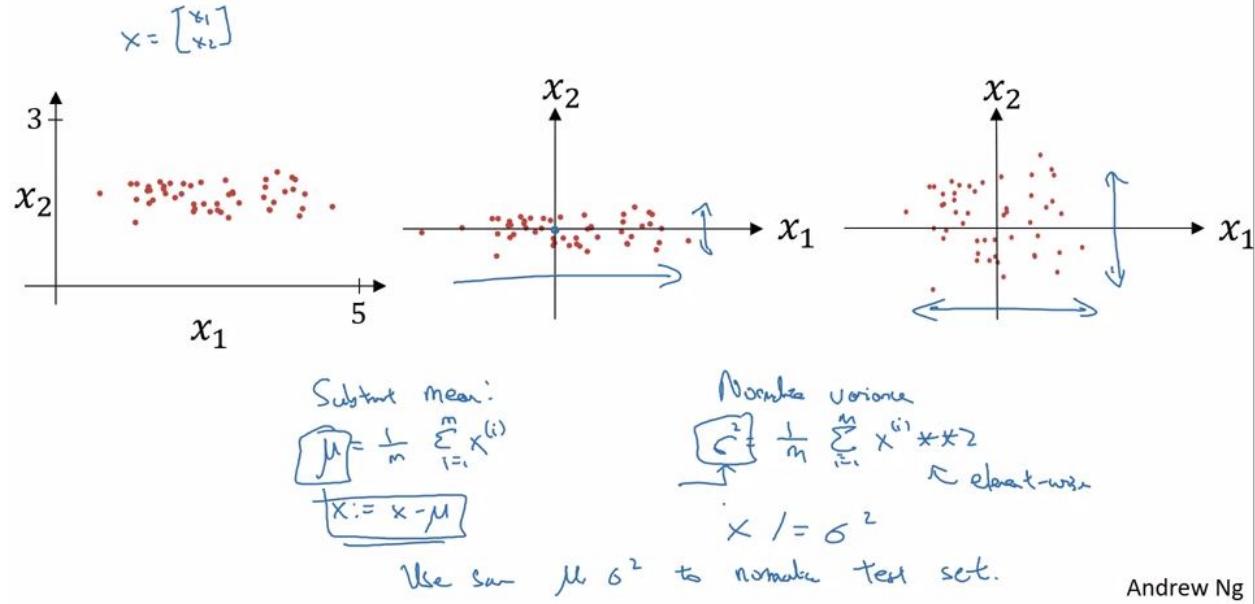
## Normalization

Subtract The Mean to have 0 mean

Normalize the variance

Don't normalize the test and training set differently, us the same mean and variance

## Normalizing training sets



### Why to Normalize?

Since features are in different scales, w1 and w2 will have different scale of values, elongated convex bowl gradient descent.

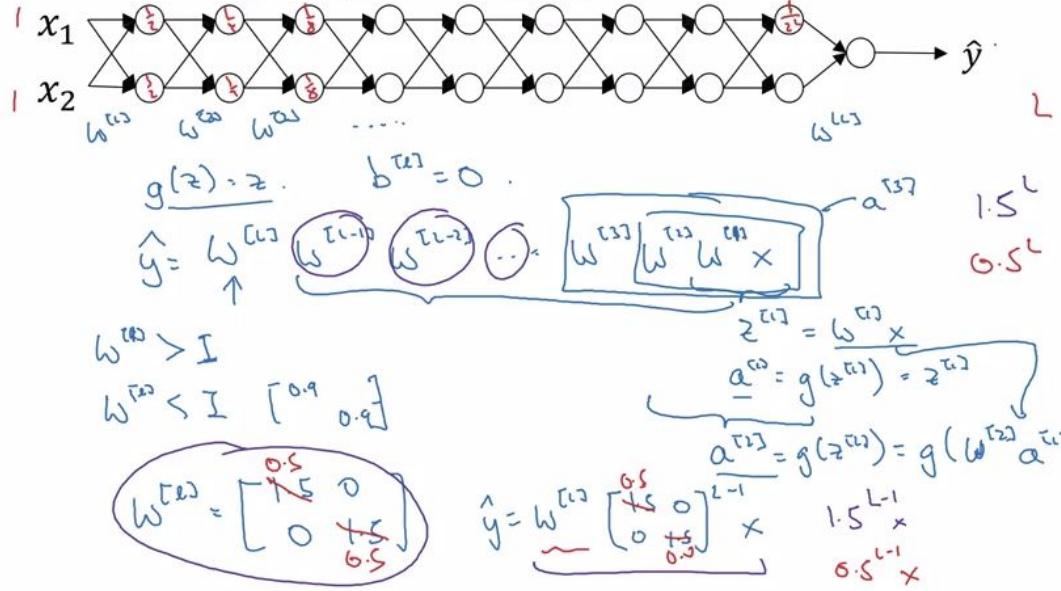
For Unnormalized data the learning needs to be small, since you have to be careful to take small steps to get to the minimum, once normalized, learning rate can be large, since it goes straight to the minimum, less spread in data.

### Vanishing / Exploding Gradients

For Deep Networks, The values of the activation functions can increase or decrease exponential, as L gets very large, therefore

# Vanishing/exploding gradients

$L=150$



Andrew Ng

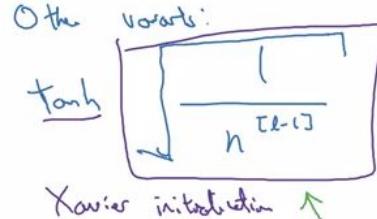
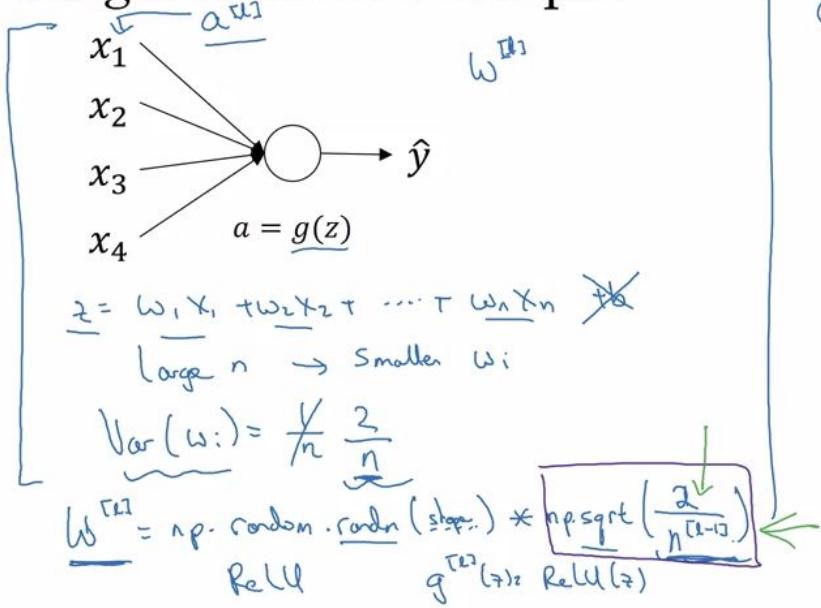
Derivative Values can get very Big Or Very Small, so learning can happen very slowly or too large steps to never reach minimum

## Partial Solution

Weight Initialization

Not Too Important to Tune

# Single neuron example



Xavier initialization ↑

$$\frac{2}{n^{(l-1)} + n^{(l)}}$$

Andrew Ng

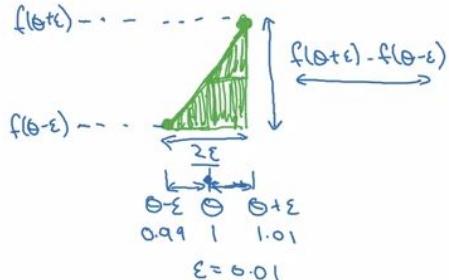
## Numerical Approximation Of Gradients

2 sided derivative calculation  $(f(\theta + \epsilon) - f(\theta - \epsilon)) / 2\epsilon$  has a lower approximation error of the derivative compared to the one sided derivative of  $(f(\theta + \epsilon) - f(\theta)) / \epsilon$

Where  $\epsilon = \text{epsilon}$

## Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(prev slide: 3.0301, error: 0.03)

USE THE TWO SIDED DIFFERENCE

Gradient Checking Debugging

## Gradient check for a neural network

Take  $\underbrace{W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}}$  and reshape into a big vector  $\theta$ .

$$\underbrace{\text{concatenate}}_{\mathcal{J}(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]})} = \underline{\mathcal{J}(\theta)}$$

Take  $\underbrace{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}}$  and reshape into a big vector  $d\theta$ .

## Gradient checking (Grad check)

$$J(\theta) = J(\theta_0, \theta_1, \theta_2, \dots)$$

for each  $i$ :

$$\rightarrow \underline{d\theta_{approx}[i]} = \frac{J(\theta_0, \theta_1, \dots, \theta_i + \varepsilon, \dots) - J(\theta_0, \theta_1, \dots, \theta_i - \varepsilon, \dots)}{2\varepsilon}$$

$$\approx \underline{d\theta[i]} = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{approx} \stackrel{?}{\approx} d\theta$$

Check

$$\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2}$$

$$\varepsilon = 10^{-7}$$

$$\approx \boxed{\frac{10^{-7}}{10^{-5}} - \text{great!}} \leftarrow$$

$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$

## Gradient checking implementation notes

- Don't use in training – only to debug
- If algorithm fails grad check, look at components to try to identify bug  
 $\underline{db}^{[i]}$      $\underline{dw}^{[i]}$
- Remember regularization.
- Doesn't work with dropout.     $J$      $\underline{\text{keep\_prob} = 1.0}$
- Run at random initialization; perhaps again after some training.

$$\underline{w, b \approx 0}$$

Andrew Ng

## Initialization of Parameters W

### Random Initialization

The model is predicting 0 for every example.

In general, initializing all the weights to zero results in the network failing to break symmetry. This means that every neuron in each layer will learn the same thing, and you might as well be training a neural network with  $n_{ij}=1$  for every layer, and the network is no more powerful than a linear classifier such as logistic regression.

#### **What you should remember:**

- The weights  $W_{ij}$  should be initialized randomly to break symmetry.
- It is however okay to initialize the biases  $b_i$  to zeros. Symmetry is still broken so long as  $W_{ij}b_i$  is initialized randomly.

### The Values of Initialization

- The cost starts very high. This is because with large random-valued weights, the last activation (sigmoid) outputs results that are very close to 0 or 1 (mathematically) for some examples, and when it gets that example wrong it incurs a very high loss for that example (because you're comparing  $y_{\text{hat}}$  and  $y$ ).
- Poor initialization can lead to vanishing/exploding gradients, which also slows down the optimization algorithm.
- If you train this network longer you will see better results, but initializing with overly large random numbers slows down the optimization.

#### **In summary:**

- Initializing weights to very large random values does not work well.
- Hopefully initializing with small random values does better. The important question is: how small should be these random values be? Lets find out in the next part!

## SUMMARY PARAMETER INITIALIZATION

- [Different initializations lead to different results](#)
- [Random initialization is used to break symmetry and make sure different hidden units can learn different things](#)
- [Don't initialize to values that are too large](#)

- [He initialization works well for networks with ReLU activations.](#)
- [He = sqrt\(1./layers\\_dims\[l-1\]\)](#)

## Regularization Applied

### L2

Forward Propagation:

$$J_{regularized} = \underbrace{-\frac{1}{m} \sum_{i=1}^m (y^{(i)} \log(a^{[L](i)}) + (1 - y^{(i)}) \log(1 - a^{[L](i)}))}_{\text{cross-entropy cost}} + \underbrace{\frac{1}{m} \frac{\lambda}{2} \sum_l \sum_k \sum_j W_{k,j}^{[l]2}}_{\text{L2 regularization cost}}$$

Backward propagation:

**Exercise:** Implement the changes needed in backward propagation to take into account regularization. The changes only concern  $dW1$ ,  $dW2$  and  $dW3$ . For each, you have to add the regularization term's gradient ( $\frac{d}{dW}(\frac{1}{2} \frac{\lambda}{m} W^2) = \frac{\lambda}{m} W$ ).

## SUMMARY L2

**Observations:**

- The value of  $\lambda$  is a hyperparameter that you can tune using a dev set.
- L2 regularization makes your decision boundary smoother. If  $\lambda$  is too large, it is also possible to "oversmooth", resulting in a model with high bias.

**What is L2-regularization actually doing?:**

L2-regularization relies on the assumption that a model with small weights is simpler than a model with large weights. Thus, by penalizing the square values of the weights in the cost function you drive all the weights to smaller values.

This happens because greater the lambda the higher the value of the derivative dW and smaller the W.

Smoother model because as described earlier, when W is small  $Z_i = WA_i + b$ ,  $Z_i$  is also small, then the activation function  $g(Z_i)$  gets value closer to zero, where the slope is linear, so  $dA$  tends to be linear, and you get a lot linear functions being multiplied with matrices.

$$dW_i = 1/m * np.dot(dZ_i, A[i-1]) + (\lambda/m)*W_i$$

So you end up with a linear prediction

It becomes too costly for the cost to have large weights! This leads to a smoother model in which the output changes more slowly as the input changes.

#### **What you should remember -- the implications of L2-regularization on:**

- The cost computation:
  - A regularization term is added to the cost
- The backpropagation function:
  - There are extra terms in the gradients with respect to weight matrices
- Weights end up smaller ("weight decay"):
  - Weights are pushed to smaller values.

#### DROPOUT SUMMARY:

- A **common mistake** when using dropout is to use it both in training and testing. You should use dropout (randomly eliminate nodes) only in training.
- Deep learning frameworks like [tensorflow](#), [PaddlePaddle](#), [keras](#) or [caffe](#) come with a dropout layer implementation. Don't stress - you will soon learn some of these frameworks.

#### **What you should remember about dropout:**

- Dropout is a regularization technique.
- You only use dropout during training. Don't use dropout (randomly eliminate nodes) during test time.
- Apply dropout both during forward and backward propagation.

- During training time, divide each dropout layer by keep\_prob to keep the same expected value for the activations. For example, if keep\_prob is 0.5, then we will on average shut down half the nodes, so the output will be scaled by 0.5 since only the remaining half are contributing to the solution. Dividing by 0.5 is equivalent to multiplying by 2. Hence, the output now has the same expected value. You can check that this works even when keep\_prob is other values than 0.5.

## HOW TO IMPLEMENT GRADIENT CHECK

**Exercise:** Implement gradient\_check\_n().

**Instructions:** Here is pseudo-code that will help you implement the gradient check.

For each i in num\_parameters:

- To compute  $J_{\text{plus}}[i]$ :
  1. Set  $\theta^+$  to `np.copy(parameters_values)`
  2. Set  $\theta_i^+$  to  $\theta_i^+ + \epsilon$
  3. Calculate  $J_i^+$  using `to_forward_propagation_n(x, y, vector_to_dictionary( $\theta^+$ ))`.
- To compute  $J_{\text{minus}}[i]$ : do the same thing with  $\theta^-$
- Compute  $\text{gradapprox}[i] = \frac{J_i^+ - J_i^-}{2\epsilon}$

Thus, you get a vector gradapprox, where  $\text{gradapprox}[i]$  is an approximation of the gradient with respect to `parameter_values[i]`. You can now compare this gradapprox vector to the gradients vector from backpropagation. Just like for the 1D case (Steps 1', 2', 3'), compute:

$$\text{difference} = \frac{\|\text{grad} - \text{gradapprox}\|_2}{\|\text{grad}\|_2 + \|\text{gradapprox}\|_2} \quad (3)$$

## Moving Averages

Number of days to average:

1/1- Beta = Number of Days

Calculate Beta and then plus below to graph Vt

Momentum for gradient descent with dW tho

$V_t = \text{Beta}(V_{t-1}) + (1-\text{Beta})(\theta_t)$

{ $\theta_t$  is basically 1 piece of current data} and { $V_{t-1}$  is summed up past data}

Bias Stoppage at the beginning

The formula is  $Vt/1 - B^t$

RMS is basically squaring the derivative terms since you want to reduce the vertical distance in the step. So  $dWi$  or  $db$  (any derivative term represents the vertical distance) so when you square it.

Bigger Number -> Get Very Big

Small Number -> Gets Very Small

Adam = Momentum Optimization + RMS (Root Mean Square) Optimization

ONLY PARAMETER YOU NEED TO TUNE IS the learning\_rate

Alpha =  $\alpha_0 / (1 + \text{decay\_rate} * \text{epoch})$  (number of mini batches completed)

Just need to tune decay rate and alpha0

The interview question on why gradient descent is bad for low amounts of data ->

Answer: This is a huge misconception

Based on the number of features

High-> then you could go toward saddle points, not a local optima. To get stuck at a local minimum there is  $2^n$ -features probability of that happening.

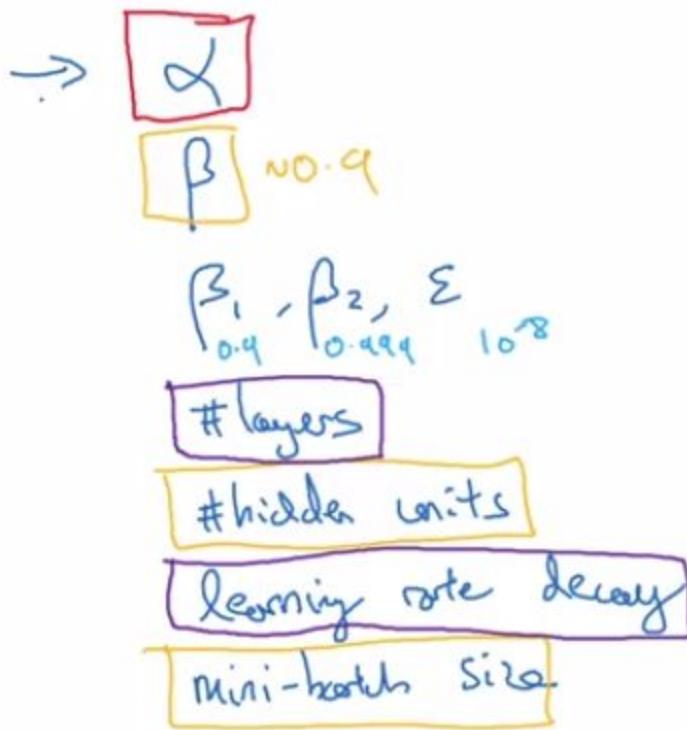
When you low amount of features then probbility increases, but the number of training and test data sets does not affect this

Since your solving problems in such high dimensional spaces, people don't have the intution of what the looks like, so YOU ARE VER VERY UNLUKILY TO GET STUCK IN ON LOCAL MINIMUM

However when you man features, you get alot saddle points, which slow down the learning and make gradient descent really slow

#### Hyperparameters Tuning Rules

# Hyperparameters



Red -> 1st

Yellow -> second

Purple -> 3rd

Not Circled -> like beta1 and beta2 and epsilon are left usually constant

USE RANDOM Sampling, however some variables might need to be scaled like alpha

Alpha -> logarithmic

Beta is also logarithmic

Layers -> Linear

Training HyperParameter Values can

Babysitting the variable or just having many computers testing variables in parallel

#### Normalizing (Batch Norm) (Quicker Traing and takes away Input Distribution)

Not just normalizing input but also all the layers

Normalize  $Z[i]$  before the activation function (there are variables gamma and beta)

Beta and Gamma are parameters

{beta[1], gamma[1], ..., beta[n], gamma[n]}

Each beta[i], gamma[i] are used for separate layer

The parameter b[i] is not needed anymore, since we are subtracting the mean and changing it up by using beta[i] (zero's out the mean)

#### Input Values Changing Distribution

Normalizing speeds up learning

Batch norm works the same way for each layer but additionally:

Covariate Change: After training the model, there is a shift in the input distribution:

Coloured Cats vs Non-Coloured

**Batch Norm keeps this distribution consistent (Mean and Variance Will stay the same)**

Earlier Layers can't shift in distribution as much, since mean and variance at a later layer is always consistent, so the each layer applies its independent logic

Has a slight regularization effect:

Since: Each Mini Batch is being scales by the mean and variance of just that mini-batch

So it adds noise just like dropout to each hidden layer -> slight regularization effect

Test Time -> just one mean and variance so used a exponentially weighted average over your training mini-batches

### Muti-Class Classification

Softmax regression:

Just Has a Different Activation Function at the end

$$1) t = e^Z[i]$$

$$2) t_i / \text{Sum of all } t_i's$$

$t$  is same dimension as  $Z$  and represents all the classes shape(Number of classes, 1)

### MultiClass Loss Function:

# Loss function

$$\begin{array}{c} \text{Diagram: } y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \xrightarrow{\text{- cat}} \begin{bmatrix} y_1 \\ y_2=1 \\ y_3 \\ y_4=0 \end{bmatrix} \\ \text{Equation: } L(\hat{y}, y) = \underbrace{-\sum_{j=1}^4 y_j \log \hat{y}_j}_{\text{small}} \\ \text{Value: } \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \quad C=4 \\ \text{Cost Function: } J(w^{(1)}, b^{(1)}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) \\ -y_2 \log \hat{y}_2 = -\log \hat{y}_2. \quad \text{make } \hat{y}_2 \text{ big.} \end{array}$$

Tensorflow does backprop but i have done a lot of the backprop myself using numpy

## Tensorflow

Minimize Cost

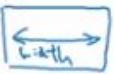
## Structuring Machine Learning Projects

Orthogonalization:

Design the model so that each control, controls one variable (Car, TV Model)

ML Controls General:

# Chain of assumptions in ML

- Fit training set well on cost function   *(≈ human-level performance)*    
*bigger network  
Adam  
...*
- Fit dev set well on cost function      
*Regularization  
Digger tiny set*
- Fit test set well on cost function     
*Digger dev set*
- Performs well in real world  *(Happy cat pic app users.)*    
*Change dev set or cost function*

---

Precision: Of the examples classified as cats, what percentage are cats

Recall: What percentage of the actual cats are recognized

Don't use both separately

get the harmonic mean together (F1 Score) take the best one, done in the dev set

## Optimizing and Satisfying

Optimizing: Maximize variable

Satisfying: one or more ( Having a threshold ) less than 100

## Setting Up Your Training / Dev / Test Sets

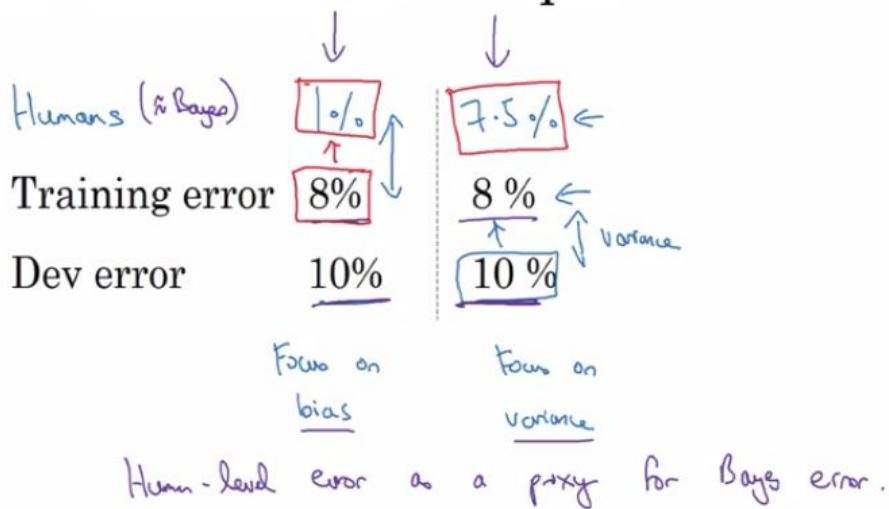
- Make sure your dev and test sets are from the same distribution
- Lots of Data today: So choose like a 1% dev/test set (millions of data)
- Tuning Set is the dev set
- Should have a test set -> no overfitting

## How to improve algorithm? (Based on Human Error)

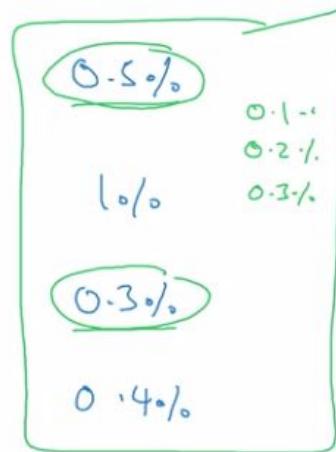
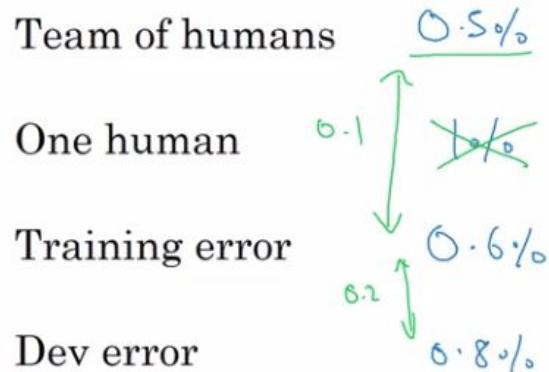
If your training set error is much larger than human error, clearly there is a bias, not training well, need to run gradient descent longer, Bigger Network.

If huge gap between dev and training -> regularization

## Cat classification example



# Surpassing human-level performance



What is avoidable bias?

Harder to tell if there is any bias on the right example cause we never know bayes error

## Error Analysis

Is it worth reducing the number of dogs being detected as cats

Just go through all images and calculate the percentage of images that were misrecognized by category and get the error for each category

If the error is huge, then focus on the category that produces the most error

## Incorrect Labelled Data

Deep Learning Quite robust to random errors in the training set

Not systematic errors

For the test and dev sets add another column to error analysis for incorrect labels

% t

## Error analysis

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeler missed cat in background
99		✓			
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	

% total states the the % error in that category

Overall dev set error ..... 100%  
Errors due incorrect labels ..... 0.6% ←  
Errors due to other causes ..... 9.4% ← | 2%  
  0.6% ←  
  1.4%

0.6% is nothing not worth fixing up

Correct dev/test set same time -> continue to come from same distribution

It is okay to have different distributions in train and dev/test data

Build new system quickly then iterate

## Training and testing on different distributions

### Cat app example

Data from webpages



$\rightarrow \approx 200,000$

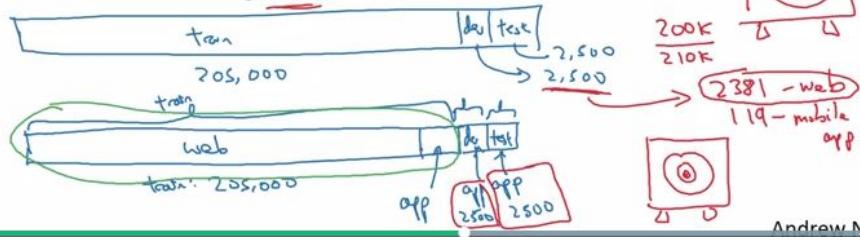
care about this

Data from mobile app

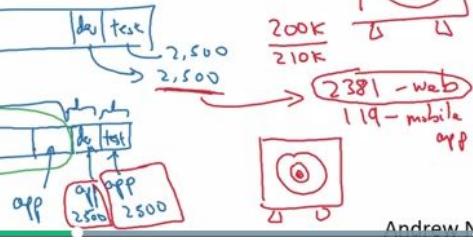


$\rightarrow \approx 10,000$

X Option 1:



Option 2:



Andrew Ng

Should add the target dev/test data into the training set, which also has data from a different distribution

## Bias and Variance with mismatched data distributions

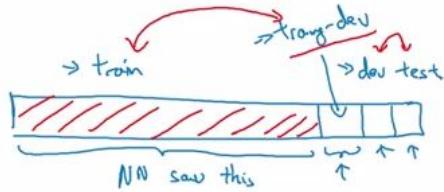
Training set and training-dev set have the same distribution

# Cat classifier example

Assume humans get  $\approx 0\%$  error.

Training error ..... 1%  $\downarrow$  9%  
 Dev error ..... 10%

Training-dev set: Same distribution as training set, but not used for training



Training error	1%	Variance	1%
$\rightarrow$ Training-dev error	9%	9%	1.5%
$\rightarrow$ Dev error	10%	10%	10% mismatch
		Variance	
Human error	0%	$\uparrow$ Avoidable bias	$\uparrow$ Avoidable bias
Training error	10%	$\uparrow$ Bias	$\uparrow$ Variance
Training-dev error	11%		$\uparrow$ Variance
Dev error	12%	$\uparrow$ Data mismatch	$\uparrow$ Data mismatch
Bias		$\uparrow$ Bias	$\uparrow$ Data mismatch
		Andrew Ng	

Human level

Training set error

Training-dev set error

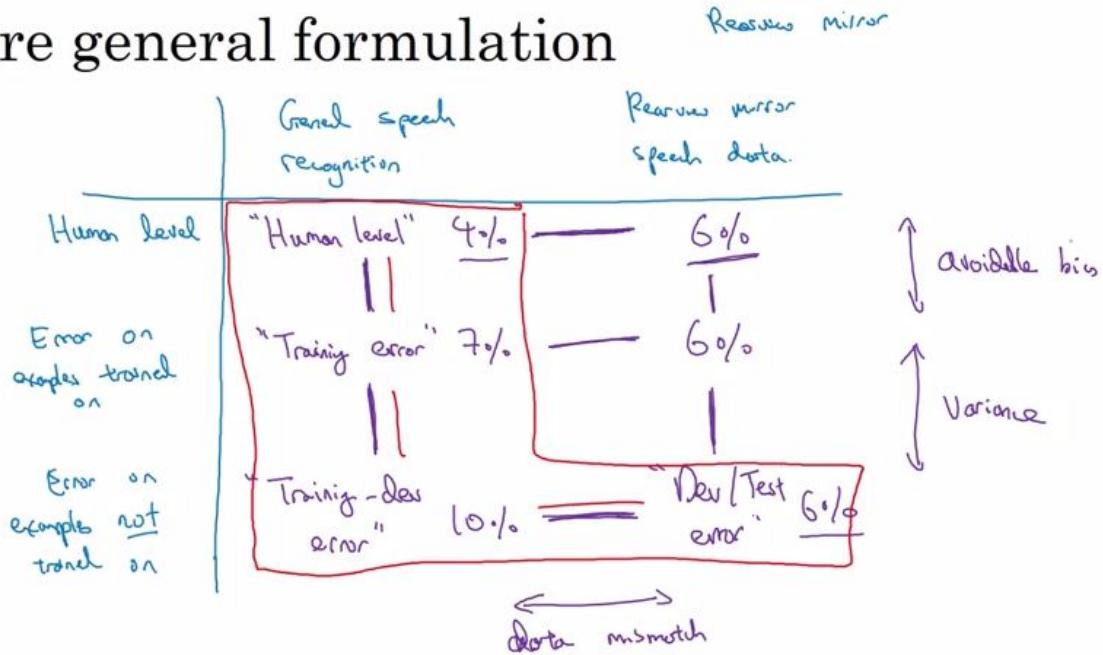
$\rightarrow$  Dev error

$\Rightarrow$  Test error

4%  $\uparrow$  avoidable bias  
 7%  $\uparrow$  variance  
 10%  $\uparrow$  data mismatch  
 12%  $\uparrow$  degree of overfit to dev set.

Purple overfit the dev-set

# More general formulation



How to handle data mismatch??

**Just collect more data relevant to the dev/test sets**

Artificial Data Synthesis (Could Overfit)

Transfer learning used a lot!!!!

PreTraining -> Image Recognition

Fine Tuning -> the last layer and calculating cost

**Used when you have a lot of data from the problem you are transferring from and less data for the problem your transferring too.**

Image recognition to radiology system

Speech recognition and drive thru communication

### Multi-task learning

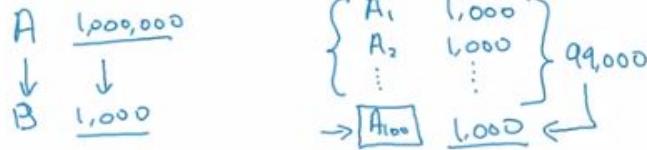
Multiple labels to one image

Solving multiple problems at the same time

Using one neural network to do multiple tasks (labels) is effective by sharing low level features.

## When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.



Creates more data

Multitasking hurts performance if neural network is small, otherwise is very helpful

Used more in computer vision

### End to End Deep Learning

Need more Data

Worked well for machine translation english to french.

Not always the best idea

face recognition

2 step approach is better -> a lot more data for individual problems

Estimation child's age through bone structure

Image -> bones -> age is better than image -> age

Excludes potentially useful hand-designed components, giving it manual knowledge by breaking it into smaller pieces, but sometimes those manual pieces of knowledge can be flawed like phonemes before words. (where computer is better of learning its own way)

#### COURSE 4: CONVOLUTION NEURAL NETWORKS

### Vertical edge detection

$$3 \times 1 + 1 \times 1 + 2 \times 1 + 0 \times 0 + 3 \times 0 + 7 \times 0 + 1 \times -1 + 8 \times -1 + 2 \times -1 = -5$$

Diagram illustrating the convolution process for vertical edge detection:

The input image is a  $6 \times 6$  matrix:

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

The filter is a  $3 \times 3$  matrix:

1	0	-1
1	0	-1
1	0	-1

The result is a  $4 \times 4$  output matrix:

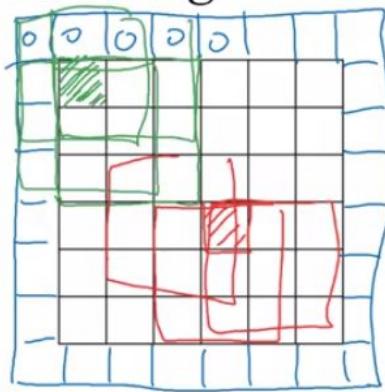
-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

Annotations include:  
"convolution" above the filter matrix.  
A multiplication symbol (\*) between the input and filter matrices.  
Dimensions  $6 \times 6$  for the input,  $3 \times 3$  for the filter, and  $4 \times 4$  for the output.

In CNNs, filters are not defined. The value of each filter is learned during the training process.

By being able to learn the values of different filters, CNNs can find more meaning from images that humans and human designed filters might not be able to find.

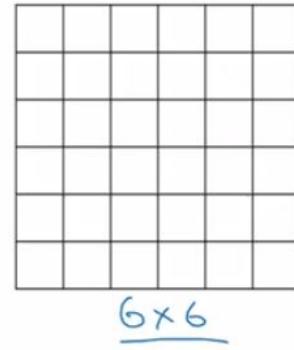
## Padding



- shrinks output  
- throws away info from edge

$$* \quad \begin{array}{|c|c|c|} \hline & & \\ \hline & & \\ \hline & & \\ \hline \end{array} =$$

$3 \times 3$   
 $f \times f$



$$\frac{6 \times L}{n \times n} \rightarrow 8 \times 8$$

$$n-f+1 \times n-f+1$$

$6-3+1=4$

$$P = \text{padding} = 1$$

$$n+2p-f+1 \times n+2p-f-1$$

$$6+2-3+1 \times \underline{\underline{ }} = 6 \times 6$$

Andrew Ng

## Padding

Another name Same Convolution (Same size output)

No padding => valid convolutions  $P = 0$

Good for two reasons

- 1) Just so the corner pixels count the same
- 2) And the output after the convolution step stays the same size

Done: Adding zeros around the initial image pixels (1 layer on each side)

Padding Size:

$$P = (f - 1) / 2$$

$f$  is odd (it is the size of the filter)

$P$  is the padding

### Stride Convolution

Skipping the filter more than once

### Strided convolution

The diagram illustrates a strided convolution operation. On the left, an input matrix of size  $7 \times 7$  is shown with values ranging from 0 to 9. A blue circle highlights the top-left element (2). A blue bracket labeled  $3 \times 3$  indicates the receptive field of this element, which covers the central 3x3 area of the input. To the right of the input is a multiplication symbol (\*). Next to it is a 3x3 kernel matrix with values 3, 4, 4; 1, 0, 2; and -1, 0, 3. Below the input is the label "stride = 2". To the right of the multiplication symbol is an equals sign (=). To the right of the equals sign is the output matrix, which is a  $3 \times 3$  matrix with values 91, 100, 83; 69, 91, 127; and 44, 72, 74. Above the output matrix is a blue arrow pointing downwards. Below the input matrix, the dimensions  $7 \times 7$  are written above the input. Below the input matrix, the labels "n x n", "p", "f x f", "strides s", and "s=2" are written. Below the output matrix, there is a large bracketed formula:  $\left[ \frac{n+2p-f}{s} + 1 \right] \times \left[ \frac{n+2p-f}{s} + 1 \right]$ . Below this formula is the calculation  $\frac{7+0-3}{2} + 1 = \frac{4}{2} + 1 = 3$ . In the bottom right corner of the box, the name "Andrew Ng" is written.

Equation below

## Summary of convolutions

$n \times n$  image       $f \times f$  filter

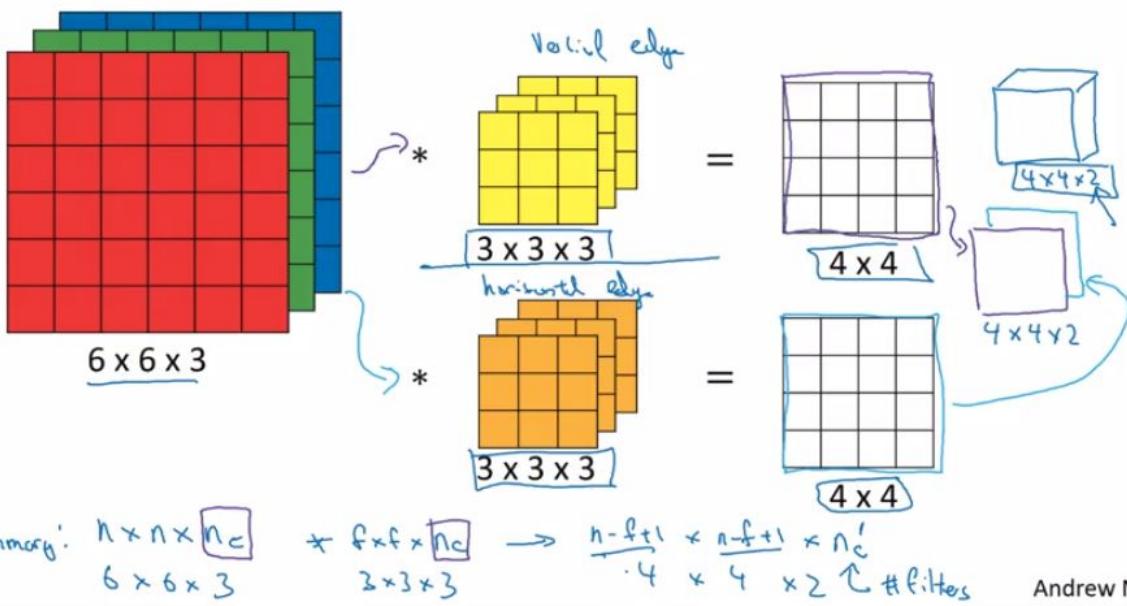
padding  $p$       stride  $s$

$$\left\lfloor \frac{n+2p-f}{s} + 1 \right\rfloor \quad \times \quad \left\lfloor \underbrace{\frac{n+2p-f}{s}}_{s} + 1 \right\rfloor$$

Don't need to flip the filter: just avoid in neural networks, only used in signal processing

Convolution on Volume:

## Multiple filters



$n'_c$  above is the number of filters

### Convolution Networks Connected to Forward Propagation in Deep learning

X: The image pixels

W: Filter

B: constant

$Z = ((X * W) + B)$  ->> the \* is not multiplication but the convolution operation

$A = \text{Relu}(Z)$

## Notation Convolution Neural Networks

### Summary of notation

If layer  $l$  is a convolution layer:

$f^{[l]}$  = filter size

$p^{[l]}$  = padding

$s^{[l]}$  = stride

$n_c^{[l]}$  = number of filters

→ Each filter is:  $f^{[l]} \times f^{[l]} \times n_c^{[l]}$

Activations:  $A^{[l-1]} \rightarrow n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Weights:  $f^{[l]} \times f^{[l]} \times n_c^{[l-1]} \times n_c^{[l]}$

bias:  $n_c^{[l]} - (1, 1, 1, n_c^{[l]})$  ← #f: bias in layer  $l$ .

Input:  $n_H^{[l-1]} \times n_W^{[l-1]} \times n_C^{[l-1]}$

Output:  $n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$

$$n_H^{[l]} = \left\lfloor \frac{n_H^{[l-1]} + 2p^{[l]} - f^{[l]}}{s^{[l]}} + 1 \right\rfloor$$

$$A^{[l]} \rightarrow m \times n_H^{[l]} \times n_W^{[l]} \times n_C^{[l]}$$

Andrew Ng

## Pooling / Average Pooling

Usually using  $f = 2$  and  $s = 2$

Just finding max number in the filter size  $f$

Same formula

$$((N + 2p - f) / s) + 1$$

# Summary of pooling

Hyperparameters:

$f$ : filter size       $f=2, s=2$   
 $s$ : stride             $f=3, s=2$   
Max or average pooling  
→ p: padding.

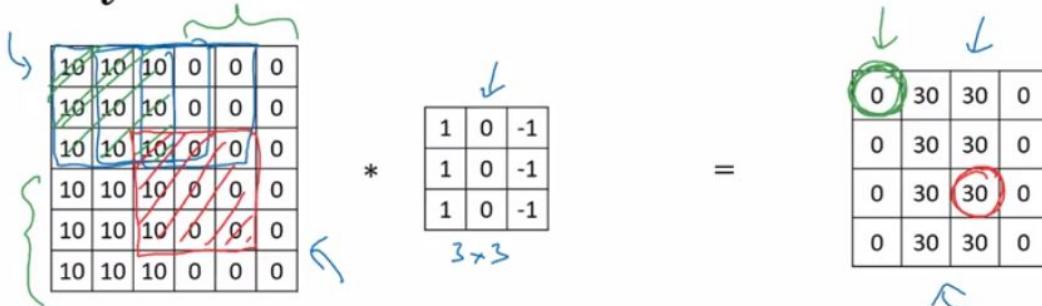
No parameters to learn!

$$\begin{aligned} & n_H \times n_w \times n_c \\ & \downarrow \\ & \left\lfloor \frac{n_H-f+1}{s} \right\rfloor \times \left\lfloor \frac{n_w-f}{s} + 1 \right\rfloor \\ & \quad \times n_c \end{aligned}$$

No parameters to learn in pooling!

Why less parameter

## Why convolutions



**Parameter sharing:** A feature detector (such as a vertical edge detector) that's useful in one part of the image is probably useful in another part of the image.

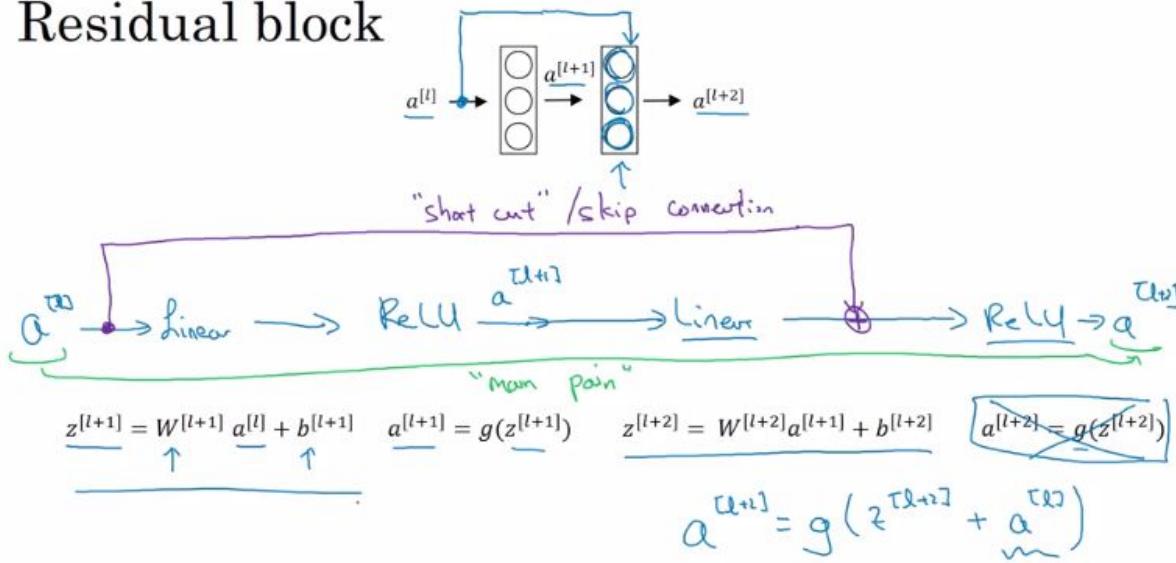
→ **Sparsity of connections:** In each layer, each output value depends only on a small number of inputs.

Andrew Ng

ResNets

Residual Block:

## Residual block



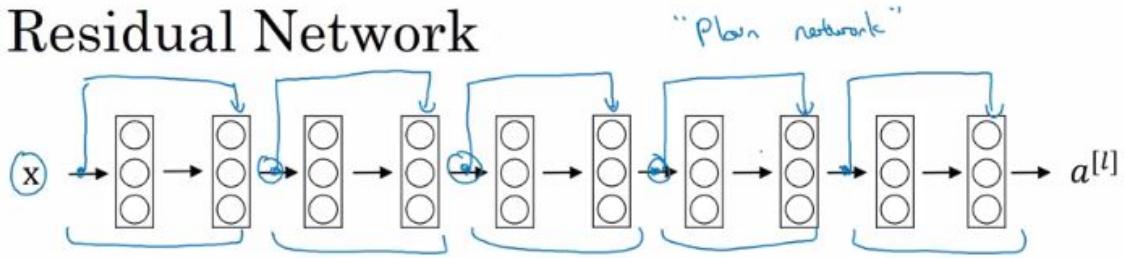
[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

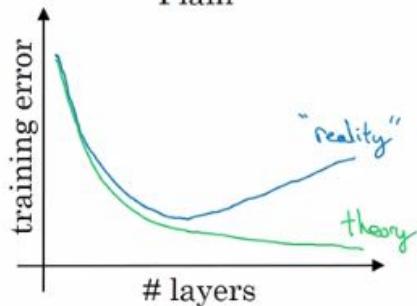
Residual network:

In reality when your network gets too deep your training error follows the blue curve, that's why we add residual blocks. Training acts like theory

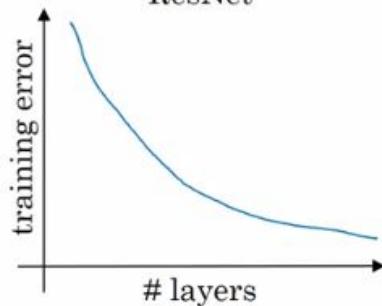
# Residual Network



Plain



ResNet



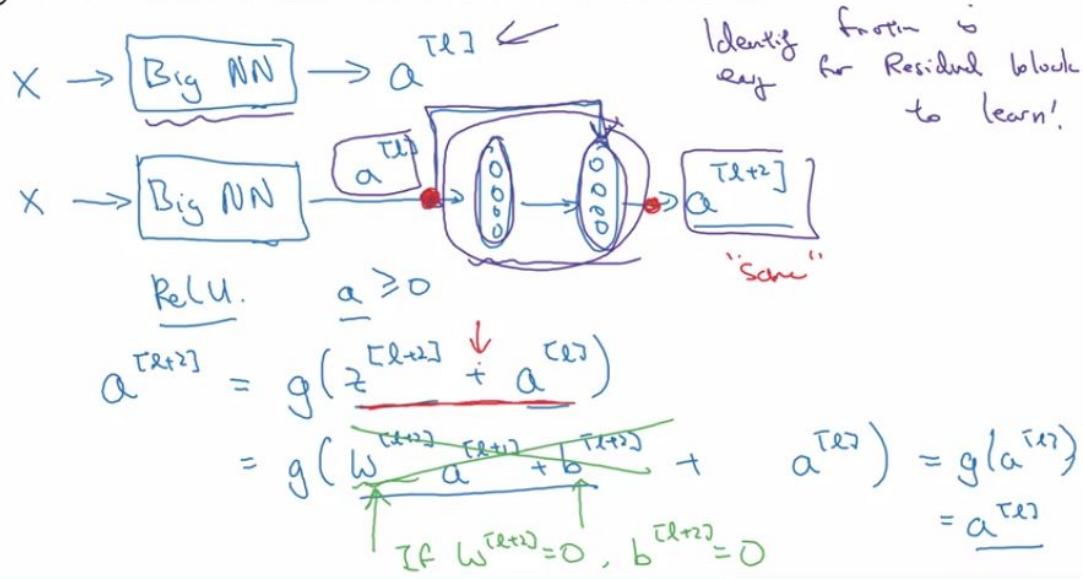
[He et al., 2015. Deep residual networks for image recognition]

Andrew Ng

## Why residual networks work?

Because of L2 regularization/ weight decay  $W[l+1]$  and  $b[l+1]$  become smaller close to zero which means you just have  $a[l]$  left and  $g(a[l])$  is just  $a[l]$  so the extra layers don't affect anything in this case.

# Why do residual networks work?

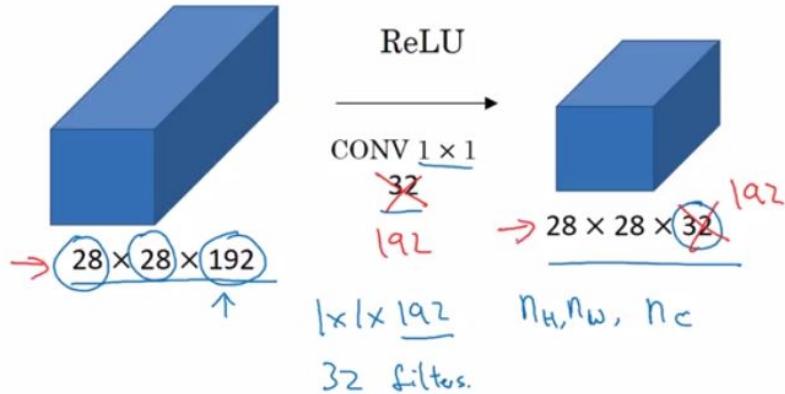


A lot of same convolutions adding padding so  $g(z[l+2] + a[l])$  is possible that you can  $z$  and  $a$ 's have the same matrix shape

## 1 x 1 convolutions

1 x 1 convolutions help you shrink the number of channels

## Using $1 \times 1$ convolutions

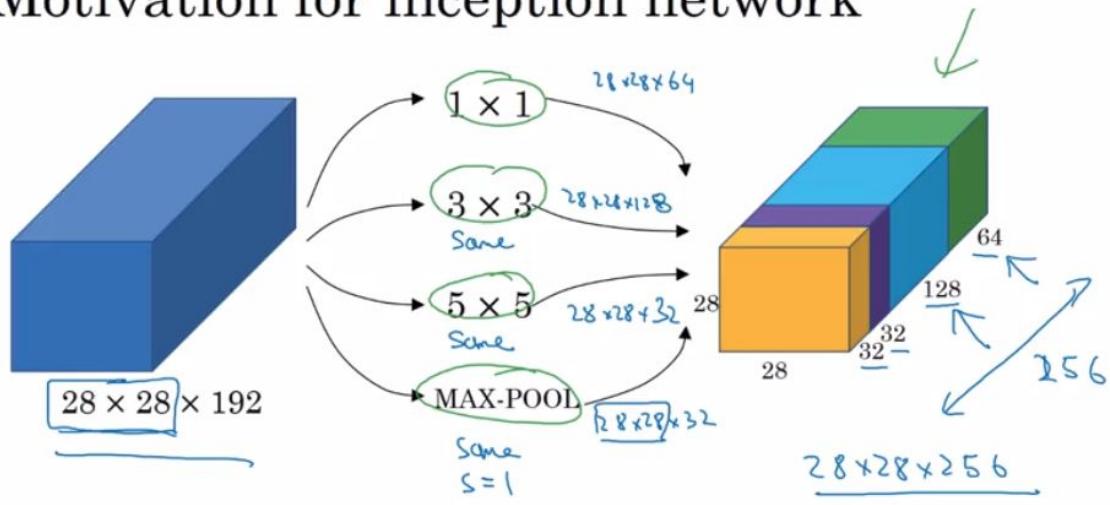


[Lin et al., 2013. Network in network]

Andrew Ng

Inception Network

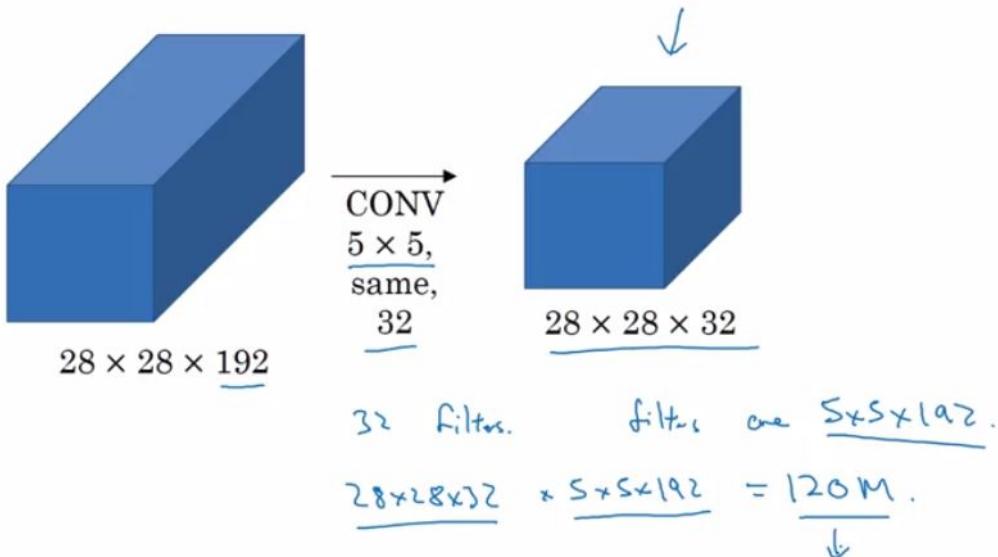
## Motivation for inception network



[Szegedy et al. 2014. Going deeper with convolutions]

Andrew Ng

## The problem of computational cost

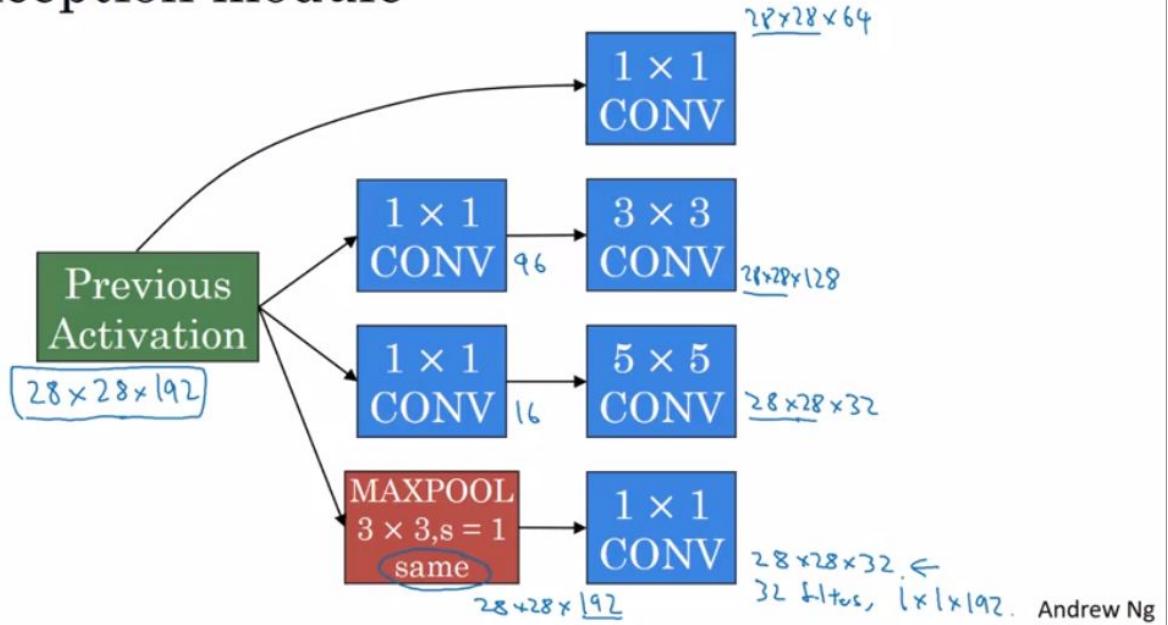


Andrei

$1 \times 1$  filter can resolve the crazy amount of calculations we need to do

So that for inception model you can do a  $1 \times 1$  conv between (bottleneck) the actual  $f \times f$  calculation the the  $1 \times 1 \times$  channel helps you reduce the channels in the original volume

## Inception module



Inception network / googleLeNet

Use Open Source Projects

- to get the architecture for your project
- And to get the pretrained weights (transfer learning)

For Computer Vision Use Transfer Learning, already pre trained models on huge data sets

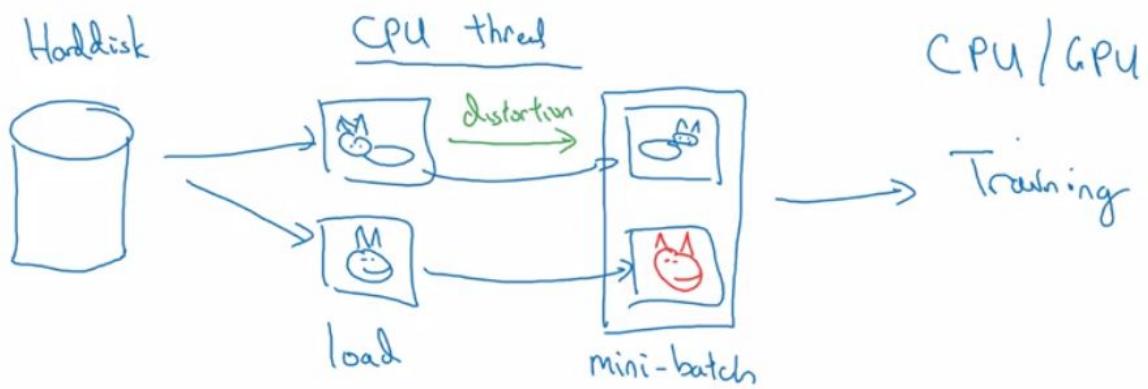
Data Augmentation

Very useful for computer vision

- Mirroring
- Random Cropping
- Shearing
- Rotation

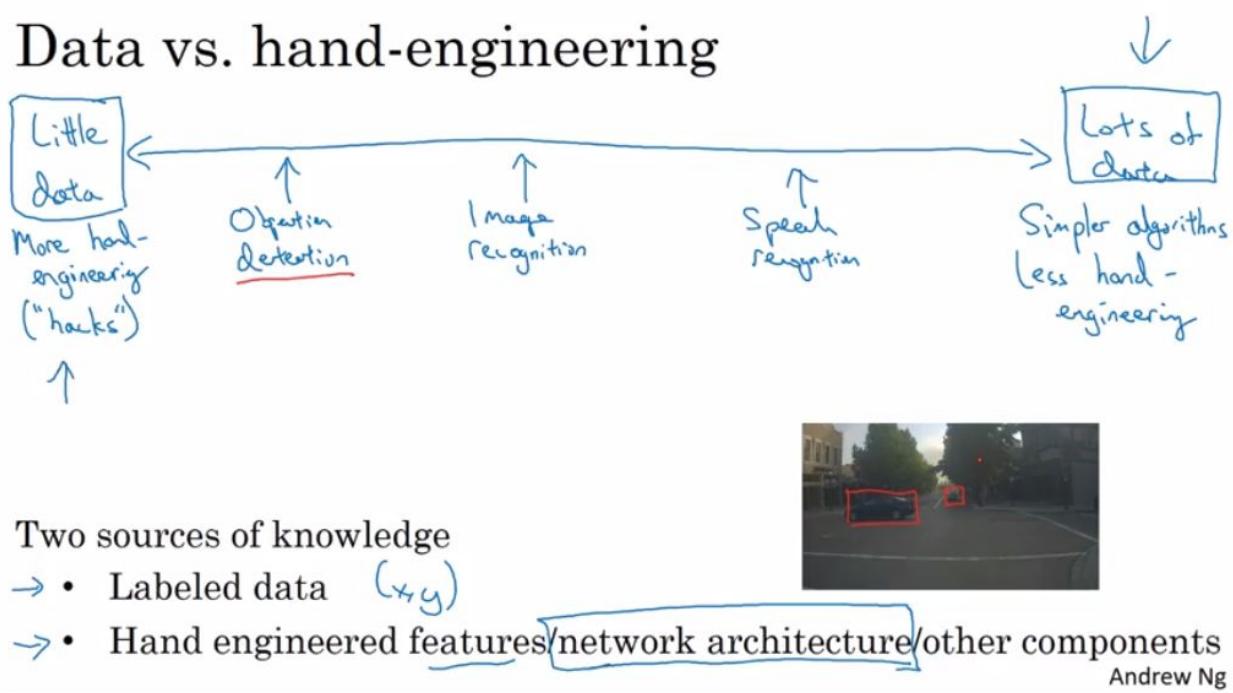
- Colour shifting (Changing the RGB with constants)

## Implementing distortions during training



CPU thread continually brings in images and then the cpu/gpu and the cpu thread run in parallel

## Data vs. hand-engineering



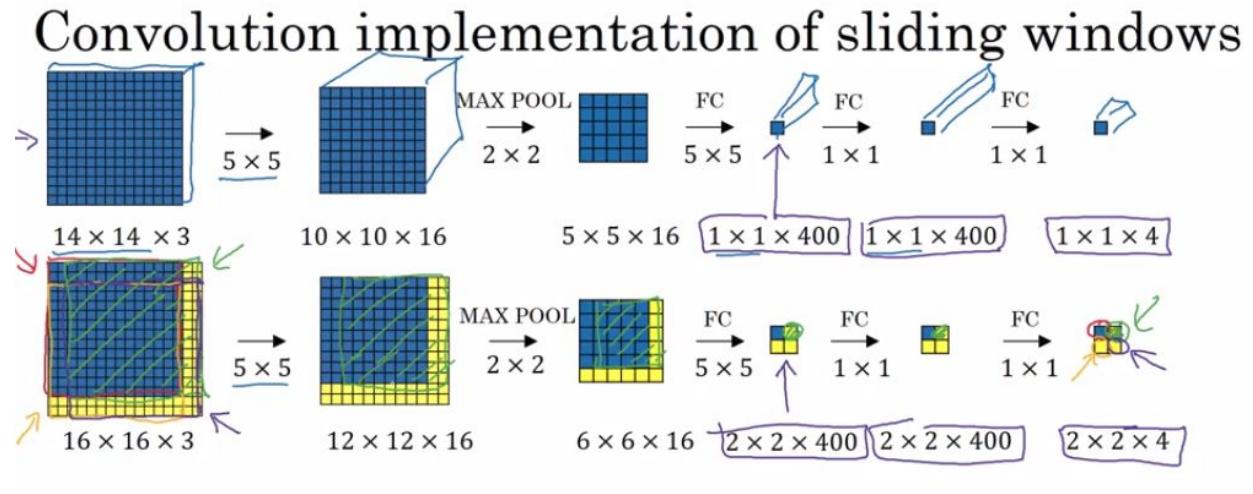
## Use open source code

- Use architectures of networks published in the literature
- Use open source implementations if possible
- Use pretrained models and fine-tune on your dataset

Localization -> Detecting where object is in the picture

BEST ALGORITHM IS THE YOLO algorithm for detection

### Sliding Windows



The above shows how to make computational cost lower, can just run the convolution layers on the image, and the final result (2D) will tell you where the object is

Since the  $14 \times 14 \times 3$  input gives you a  $1 \times 1$  output in the original convolution layers, any image you put has to be bigger than  $14 \times 14 \times 3$  and the sliding window size now is  $14 \times 14 \times 3$

Disadvantage the sliding window is fixed size

# Evaluating object localization



Intersection over Union (IoU)

$$= \frac{\text{size of } \begin{array}{c} \text{yellow} \\ \text{box} \end{array}}{\text{size of } \begin{array}{c} \text{green} \\ \text{box} \end{array}}$$

"Correct" if  $\text{IoU} \geq \underline{0.5}$  ←

0.6 ←

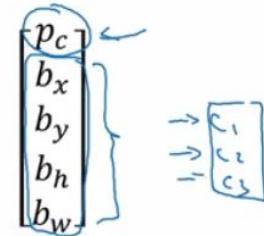
More generally, IoU is a measure of the overlap between two bounding boxes.

non-Max Suppression Algo

## Non-max suppression algorithm



Each output prediction is:



Discard all boxes with  $p_c \leq 0.6$

→ While there are any remaining boxes:

- Pick the box with the largest  $p_c$   
Output that as a prediction.
- Discard any remaining box with  
 $\text{IoU} \geq 0.5$  with the box output  
in the previous step

Andrew Ng

Anchor Box For Multiple Object Detection

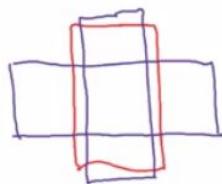
# Anchor box algorithm

Previously:

Each object in training image is assigned to grid cell that contains that object's midpoint.

Output y:

$3 \times 3 \times 8$



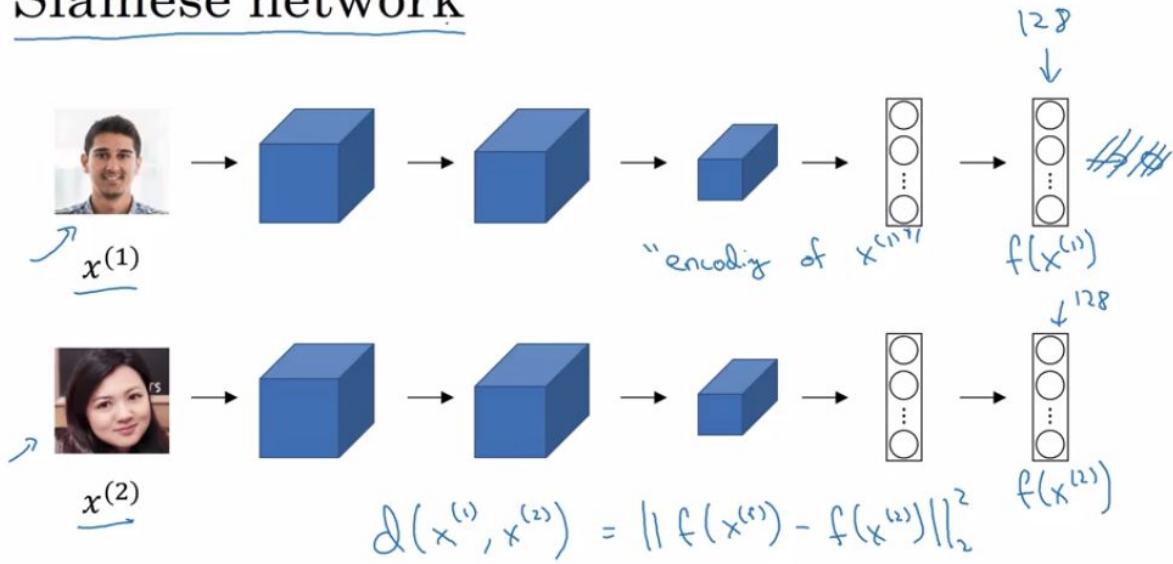
With two anchor boxes:

Each object in training image is assigned to grid cell that contains object's midpoint and anchor box for the grid cell with highest IoU.

(grid cell, anchor box.)

FACE RECOGNITION

## Siamese network



A way determining how similar the faces are using a one shot approach. (only one picture available)

The 128 dimension output are the encodings

Triplet Loss Function:

Basically comparing a positive match and a negative match.

Getting a difference function for both positive and negative match.

Then applying alpha to make sure there is a certain margin of difference between them

Need to make sure that the negative match also has a small difference function, so you can train harder pictures and get a better result

## Face Verification and Binary Classification

Just compare two cnn output from two different images and label the images a 0 or 1. (Logistic Regression)

### What is neural style transfer?

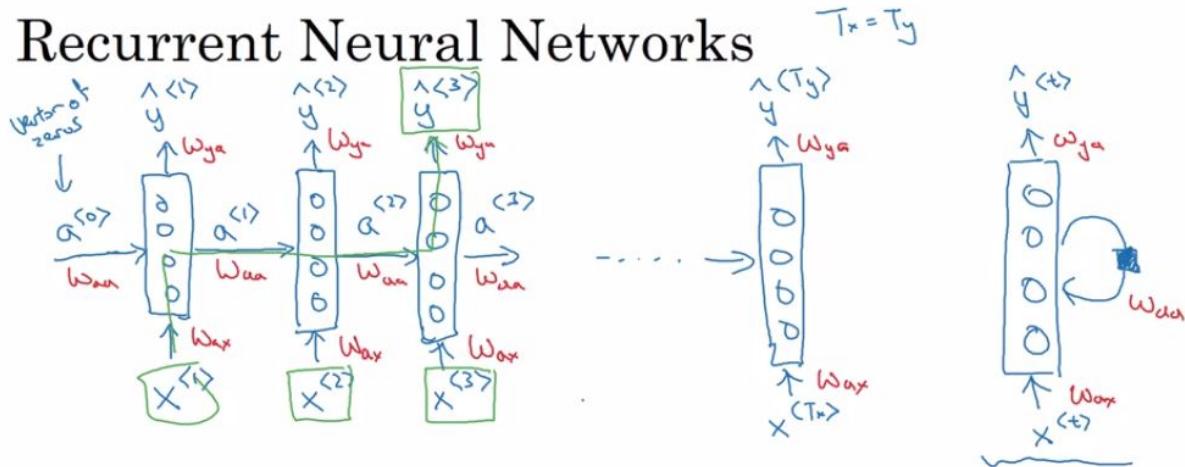
Combining one picture's content and another picture style

Content is just the difference between the activations for a layer in content picture and the final picture

Style has a G style matrix for picture S the style picture and G the final picture and figures out the their difference for the cost function

### Sequence Models

## Recurrent Neural Networks



Can have many relationships depending on the size of input and output

One to One

One to Many (Music Generation)

Many to One (Movie Rating)

Many to Many (2 situations equal length and different size lengths)

### **Exploding Gradients:**

Gradient Clipping: Just rescale

### **Vanishing Gradients Problem**

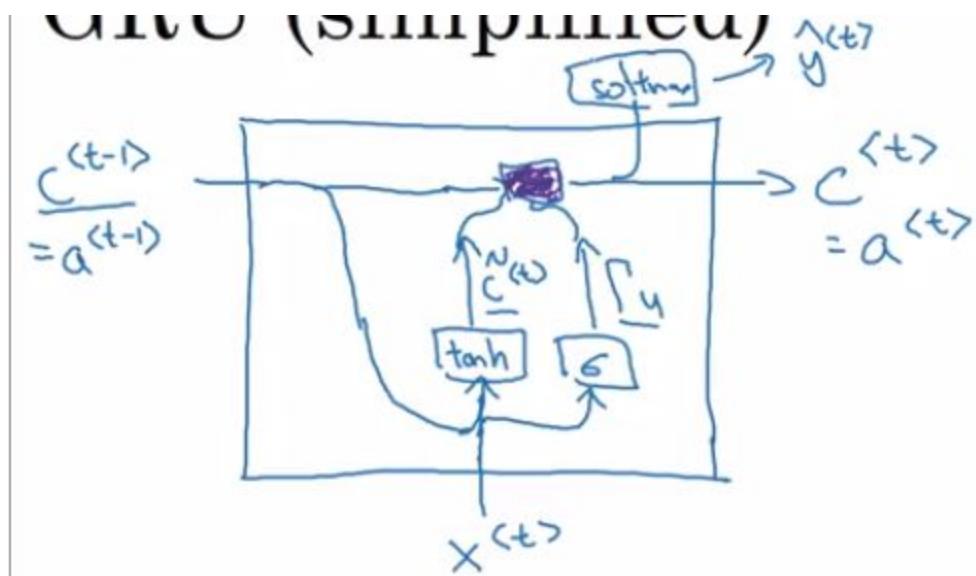
RNN is not very good at catching long term dependencies

Cats ..... Were

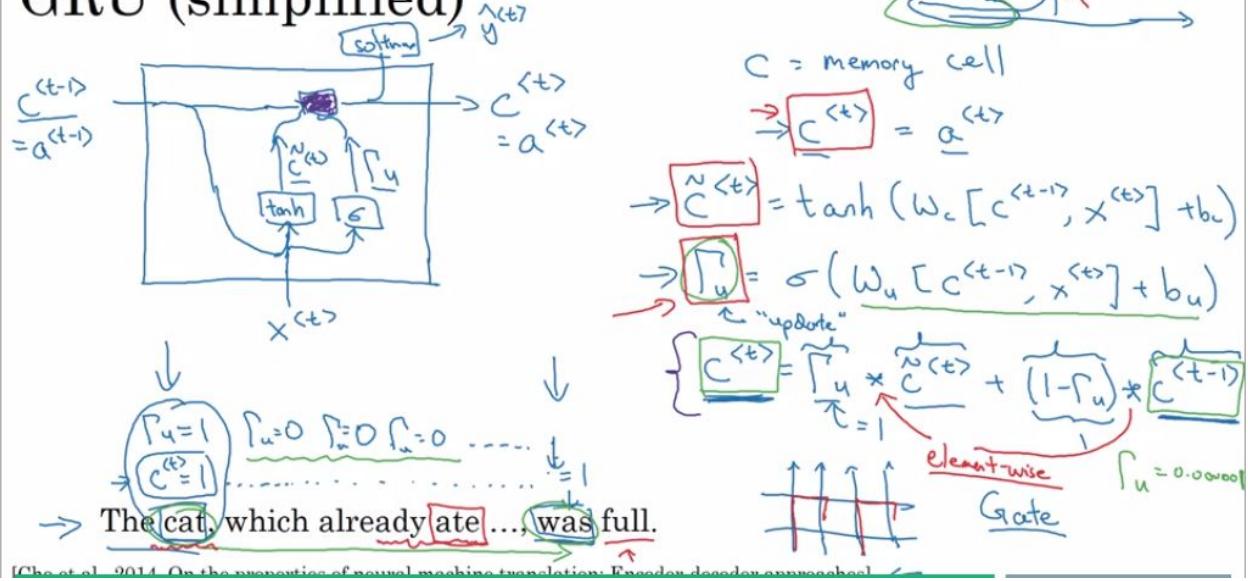
Cats ..... Was

Many local layers highly effect the output values (Because its very hard for an error to backpropagate long distances away. (long distances mean lots of layers)

### Solution: Gated Recurrent Unit (GRU)



## GRU (simplified)



$C_t$ ,  $\Gamma_u$  and  $C_t^{\sim}$  is actually a vector with multiple bits representing different facts about a sentence

$C_t^{\sim}$  is for what facts the current words adds to the sentence, the gate filters out what info is important what is not okay to update

## GRU

$$\tilde{c}^{} = \tanh(W_c[\Gamma_r * \underline{c}^{}, x^{}] + b_c)$$

$$\underline{\Gamma_u} = \sigma(W_u[\underline{c}^{}, x^{}] + b_u)$$

$$\underline{\Gamma_r} = \sigma(W_r[\underline{c}^{}, x^{}] + b_r)$$

$$\underline{\underline{c}^{}} = \underline{\Gamma_u} * \underline{\tilde{c}^{}} + (1 - \underline{\Gamma_u}) * \underline{\underline{c}^{}}$$

$$\underline{\underline{a}^{}} = \underline{\underline{c}^{}}$$

Another Solution: LSTM

# LSTM in pictures

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

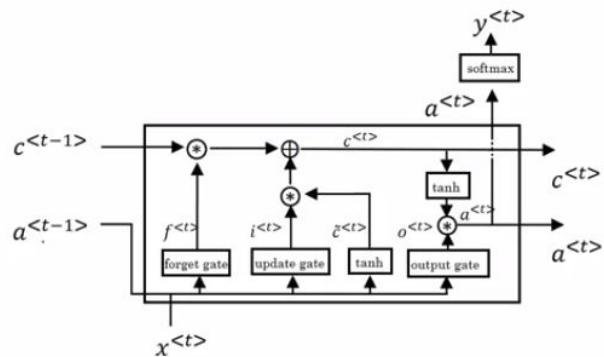
$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$



## Bidirectional RNN

NLP: Uses Bidirectional RNN + LSTM

However for real time conversation you would have to wait until the person stops talking to start the algo, not very effective!

## Deep RNN

Not many layers as CNN

## Word Embedding and Transfer Learning

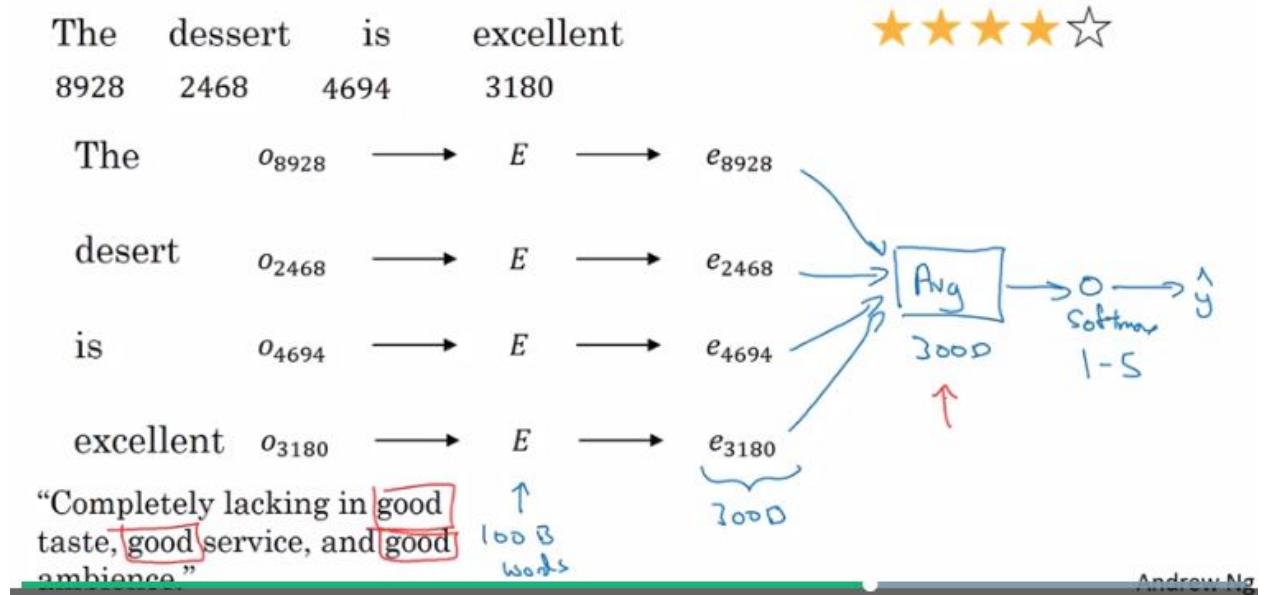
Can be useful for transfer learning when you have a small data set

## Application of Sentiment Classification Model

Word embeddings

"Can't just take the average of all your word embeddings",

# Simple sentiment classification model



Apply RNN "not good" does not mean "good" a sequence comes to affect here

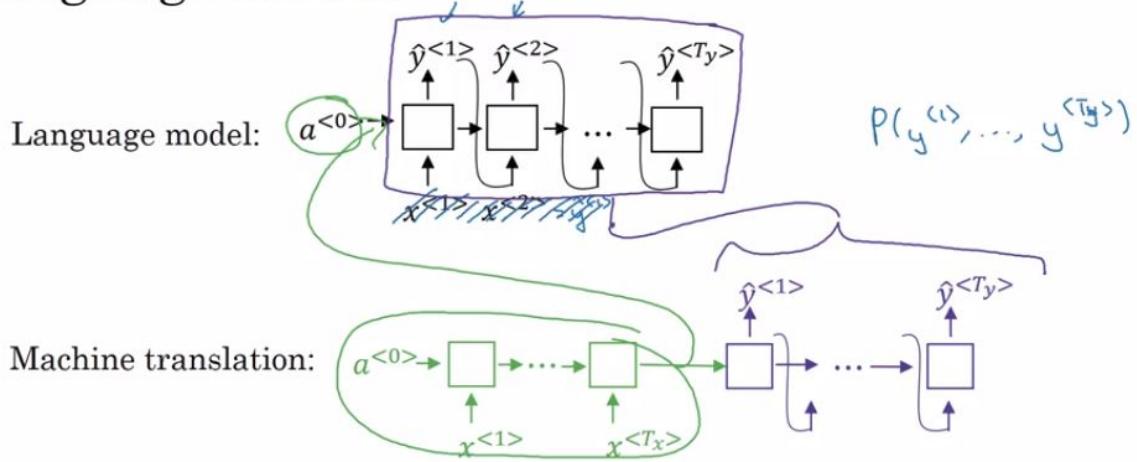
## Emojify/Sentiment Analysis

All about using Embedding vectors, training them if large data set, otherwise use transfer learning and applying it to a sentiment analysis. Using LSTM to keep a sequence words actual meaning.

## Sequence to Sequence Models

Image to caption -> google images

# Machine translation as building a conditional language model



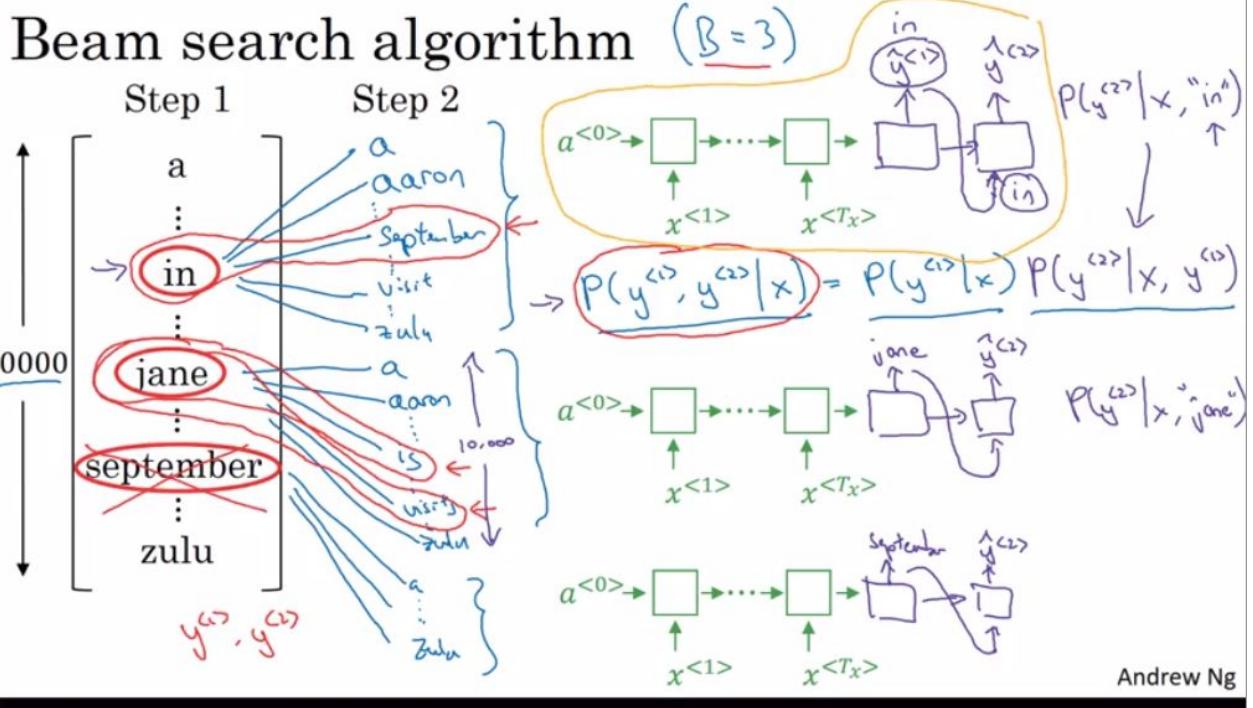
Machine translation is a conditional language model  $P(y_1 \dots y_{T_y} | x_1 \dots x_{T_x})$

Greedy Approach doesn't work, since some words might be used more often, but another word makes sense in the sequence

## Beam Search Algorithm

Beam = # of groups you can select being mostly likely

1. Select beam number of words most likely for first word
2. Then for each beam # of words look at the next possible word options
3. select beam # of groups of 2 words that end up with greatest probability
4. Then repeat again for the beam # of groups to look at the next possible word options



If beam # = 1 -> that's just greedy search algo

#### Fixing Issues in Beam Algo Underflow

Since multiplying very small numbers various times for the probabilities could result in an underflow, we must log each probability

Beam Algo likes shorter sentences since probability is higher

Just normalize it, divide by the number of words in the sequence

## Length normalization

$$\begin{aligned}
 & p(y^{(1)} \dots y^{(T_y)} | x) = \frac{p(y^{(1)} | x) p(y^{(2)} | x, y^{(1)}) \dots}{p(y^{(T_y)} | x, y^{(1)}, \dots, y^{(T_y-1)})} \\
 & \arg \max_y \prod_{t=1}^{T_y} P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \\
 & \log \left( \arg \max_y \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \right) \leftarrow \overline{T_y} \\
 & \frac{1}{\overline{T_y}^\alpha} \sum_{t=1}^{T_y} \log P(y^{(t)} | x, y^{(1)}, \dots, y^{(t-1)}) \quad \alpha = 0.7 \quad \frac{\alpha = 1}{\alpha = 0}
 \end{aligned}$$

Andrew Ng

Larger Beam  $\rightarrow$  Slower  $\rightarrow$  better outcome tho

Small Beam  $\rightarrow$  Faster  $\rightarrow$  worst outcome

### Error analysis in beam search

1 beam is at fault  $\rightarrow$  increase beam number

If RNN is at fault  $\rightarrow$  add regularization, change architect, more data

$y^*$  being the human output

$y^\wedge$  being the machine output

## Error analysis process

Human	Algorithm	$P(y^* x)$	$P(\hat{y} x)$	At fault?
Jane visits Africa in September.	Jane visited Africa last September.	$2 \times 10^{-10}$	$1 \times 10^{-10}$	(B)
...	...	—	—	(R)
...	...	—	—	R R R ...

Figures out what fraction of errors are “due to” beam search vs. RNN model

An's new Na

Evaluating how good the machine output is compared to the human output

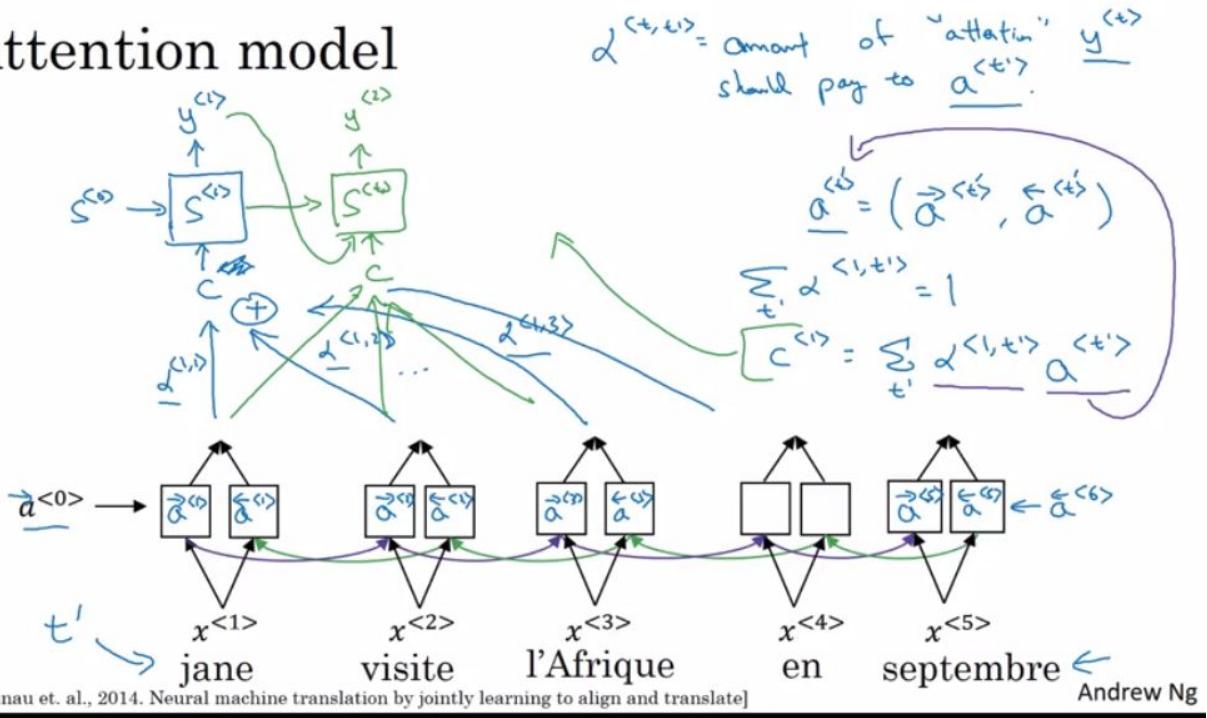
Called the Blue Score

BP penalizes shorter length sentences, cause easier to get a high blue score

Networks have a hard time memorizing long sentences so we use the attention model

Attention Model: Audio Data

# Attention model

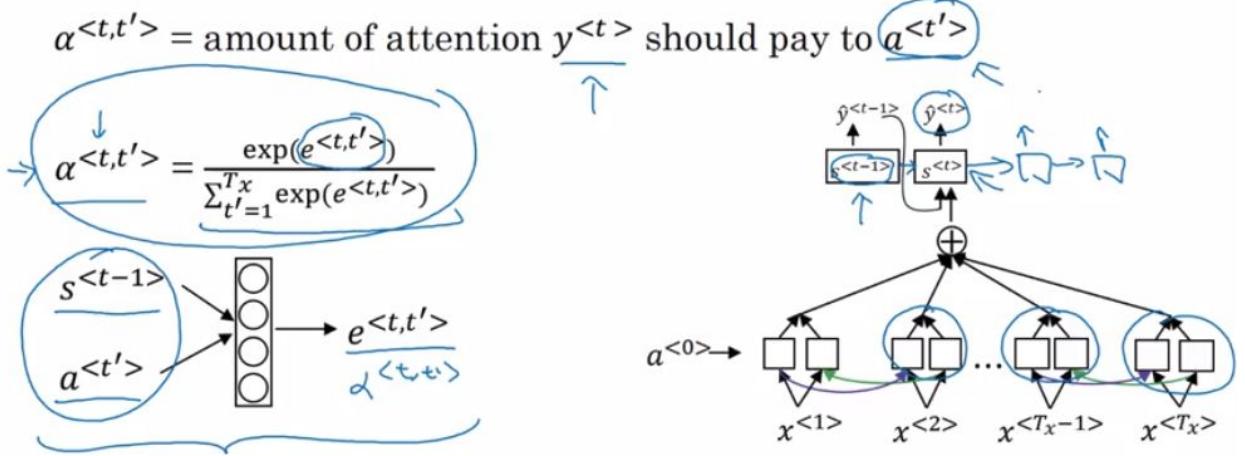


[ahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]

Andrew Ng

- Bidirectional RNN's at the bottom
- The alpha represents how much attention each output should give to each word.
- Alpha's sum to 1 for each output

# Computing attention $\alpha^{<t,t'>}$



[Bahdanau et. al., 2014. Neural machine translation by jointly learning to align and translate]  
[Xu et. al., 2015. Show, attend and tell: Neural image caption generation with visual attention]

Andrew Ng

Quadratic Cost

Speech Recognition

- Trigger Words are just a long RNN that output 0's and 1's