

Course summary

Here are some course summary as its given on the course link:

If you want to break into the cutting edge AI, this course will help you do so. Deep learning engineers are highly sought after and mastering deep learning will give you numerous new career opportunities. Deep learning is also a new "super power" that will let you build AI systems that just weren't possible a few years ago.

In this course, you will learn the foundations of deep learning. When you finish this class, you will:

- Understand the major technology trends driving Deep Learning
- Be able to build, train and apply fully connected deep neural networks
- Know how to implement efficient (vectorized) neural networks
- Understand the key parameters in a neural network's architecture

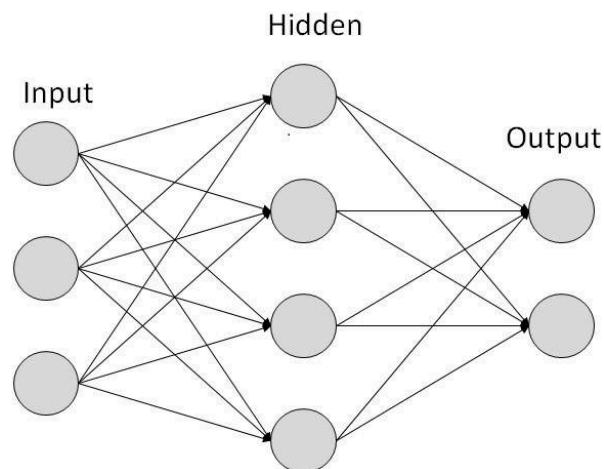
This course also teaches you how Deep Learning actually works, rather than presenting only a cursory or surface level description. So after completing it, you will be able to apply deep learning to your own applications. If you are looking for a job in AI, after this course you will also be able to answer basic interview questions.

Introduction to Deep Learning

"Be able to explain the major trends driving the rise of deep learning, and understand where and how it is applied today."

What is a (Neural Network) NN?

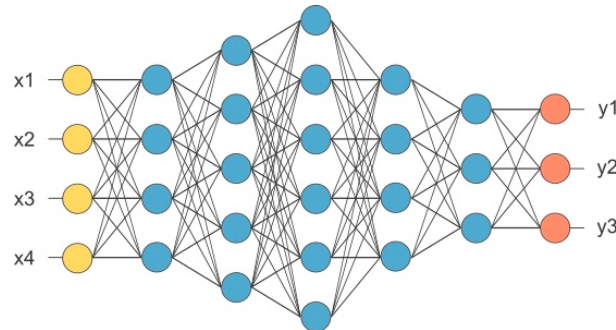
- Single neuron == linear regression
- Simple NN graph



- **ReLU** stands for **Rectified Linear Unit** - is the most popular activation function right now that

makes deep NNs train faster now.

- Hidden layers predicts connection between inputs automatically, that's what deep learning is good at.
- Deep NN consists of more hidden layers (deep layers).



- Each input will be connected to the hidden layer and the NN will decide the connections.
- Supervised learning means we have the (X, Y) and we need to get the function that maps X to Y.

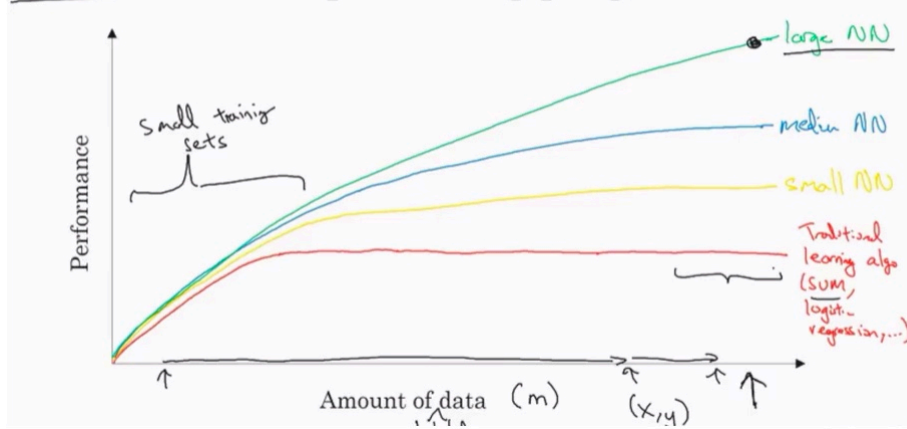
Supervised learning with neural networks

- Different types of neural networks for supervised learning which includes:
 - **CNN** or Convolutional Neural Networks (Mạng Neuron tích chập) (Useful for Computer Vision)
 - **RNN** or Recurrent Neural Networks (Mạng neuron hồi quy) (Useful in Speech Recognition or Natural Language Processing)
 - **Standard NN** (Useful for Structured data)
 - Hybrid/ custom NN or a collection of NNs types
- Structured data is like the databases and tables.
- Unstructured data is like videos, audios, texts.
- Structured data gives more money because companies relies on prediction on its big data.

Why is deep learning taking off?

- Deep learning is taking off for 3 reasons:
 - Data:
 - Using this image we can conclude:

Scale drives deep learning progress



- For small data, NN can perform as Linear regression or SVM
- For big data, a small NN is better than SVM
- For big data, a big NN is better than a medium NN, and better than a small NN
- Hopefully we have a lot of data because the world is using the computer a little bit more.
 - Mobiles
 - IOT
- Computation:
 - GPUs
 - Powerful CPUs
 - Distributed Computing
 - ASICs
- Algorithms:
 - Creative algorithms has appeared that changed the way NN works.
 - For example using ReLu function is so much better than using Sigmoid function in training a NN because it helps with the vanishing gradient problem.

Neural network basics

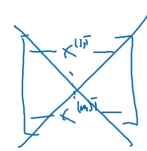
Learn to set up a machine learning problem with a neural network mindset. Learn to use vectorization to speed up your models.

Binary classification

- Mainly he (Andrew Ng) is talking about how to do a logistic regression to make a binary classifier.
- He talked about an example of knowing if the current image contains a cat or not.
- Here are some notations:
 - M is the number of examples in the dataset
 - N_x is the size of the input vector
 - N_y is the size of the output vector
 - $X(1)$ is the first input vector
 - $Y(1)$ is the first output vector
 - $X = [x(1) \ x(2) \ x(3) \ \dots \ x(M)]$
 - $Y = [y(1) \ y(2) \ y(3) \ \dots \ y(M)]$
- We will use python in this course

- In Numpy we can make matrices and make operations on them in a fast and reliable time.

$(x, y) \quad x \in \mathbb{R}^{n_x}, y \in \{0, 1\}$
 m training examples: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
 $M = M_{\text{train}} \quad M_{\text{test}} = \# \text{test examples.}$

$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$

 $X \in \mathbb{R}^{n_x \times m} \quad X.\text{shape} = (n_x, m)$

$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$
 $Y \in \mathbb{R}^{1 \times m}$
 $Y.\text{shape} = (1, m)$

Logistic Regression

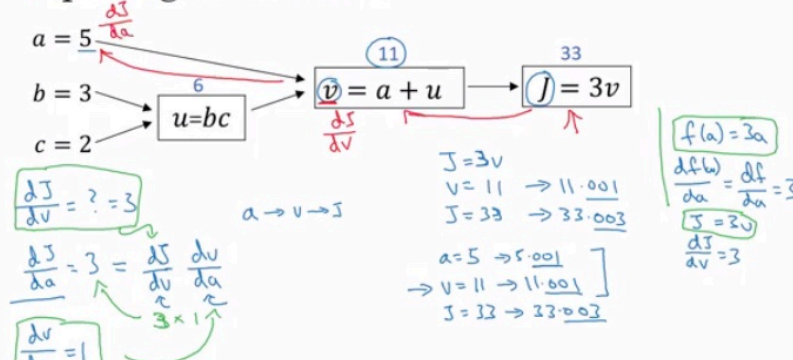
- Algorithm is used for classification algorithm of 2 classes.
- Equations:
 - Simple equation: $y = wT.x + b$
 - If we need y to be in between 0 and 1: $y = \text{sigmoid}(wT.x + b)$
 - In some notations this might be used: $y = \text{sigmoid}(w(\text{transpose})x)$
 - While b is w_0 of w and we add $x_0 = 1$. But we won't use this notation in this course.
- In binary classification Y has to be between 0 and 1
- In the last equation w is a vector of N_x and b is a real number

Logistic Regression Cost Function

- First lost function would be the square root error: $L(y', y) = 1/2 (y' - y)^2$
- But we won't use this notation because it leads us to optimization problem which is non convex, means it contains local optimum points.
- This is the function that we will use: $L(y', y) = -(y' \log(y') + (1 - y') \log(1 - y'))$
- To explain the last function lets see:
 - if $y = 1 \implies L(y', 1) = -\log(y') \implies$ we want y' to be the largest $\implies y'$ biggest value is 1
 - if $y = 0 \implies L(y', 0) = -\log(1 - y') \implies$ we want $1 - y'$ to be the largest $\implies y'$ to be smaller as possible because it can only has 1 value.
- Then the Cost function will be: $J(w, b) = (1/m) * \text{Sum}(L(y'[i], y[i]))$
- The loss function computes the error for a single training example; the cost function is the average of the loss functions of the entire training set.

Gradient Descent

- We want to predict w and b that minimize the cost function
- Our cost function is convex.
- First we initialize w and b to 0.0 or initialize them to a random value in the convex function and then try to improve the values of the cost function to reach global optimum.
- In logistic regression people always use 0.0 instead of random.
- The gradient descent algorithm repeats: $w = w - \alpha * dw$



Gradient Descent on m Examples

- Lets say we have these variables:

X1	Feature
X2	Feature
W1	Weight of the first feature.
W2	Weight of the second feature.
B	Logistic Regression parameter.
M	Number of training examples
Y(i)	Expected output of i

- So we have:

```
X1 \
W1 \
X2 ==> z(i) = X1W1 + X2W2 + B ==> a(i) = Sigmoid(z(i)) ==> l(a(i),Y(i)) = - (Y(i)*log(a(i)) + (1-Y(i))*log(1-a(i)))
Y2 /
B /
```

- Then from right to left we will compute the derivations compared to the result:
 - $d(a) = d(l)/d(a) = -(y/a) + ((1-y)/(1-a))$
 - $d(z) = d(l)/d(z) = a - y$
 - $d(W1) = X1 * d(z)$
 - $d(W2) = X2 * d(z)$
 - $d(B) = d(z)$
- From the above we can conclude the logistic regression pseudo code:

```
J = 0; dw1 = 0; dw2 = 0; db = 0
w1 = 0; w2 = 0; b = 0
```

```
for i = 1 to m:
```

```
  # forward pass
```

```
  z(i) = w1 * x1(i) + w2 * x2(i) + b
```

```
  a(i) = Sigmoid(z(i))
```

```
  J += (Y(i) * log(a(i)) + (1 - Y(i)) * log(1 - a(i)))
```

```
  # Backward pass
```

```
  dz(i) = a(i) - Y(i)
```

```
  dw1 += dz(i) * x1(i)
```

```
  dw2 += dz(i) * x2(i)
```

```
  db += dz(i)
```

```
J /= m
```

```
dw1 /= m
```

```
dw2 /= m
```

```
db /= m
```

```
# Gradient descent
w1 = w1 - alpha * dw1
w2 = w2 - alpha * dw2
b = b - alpha * db
```

- The above code should run for some iterations to minimize error.
- So there will be two linear loops to implement the logistic regression
- Vectorization is so important on deep learning to reduce loops. In the last code we can make the whole loop in one step using vectorization!

Vectorization

- Deep learning shines when the dataset are big. However for loops will make you wait a lot for a result. That's why we need vectorization to get rid of some of our for loops.
- Numpy library (dot) function is using vectorization by default.
- The vectorization can be done on CPU or GPU through the SIMD operation. But it's faster on GPU.
- Whenever possible avoid explicit for loops.
- Most of the Numpy library methods are vectorized version.

Vectorizing Logistic Regression

- We will implement Logistic Regression using one for loop.
- As an input we have a matrix X and its shape is [Nx, m] and a matrix Y which is [Ny, m]
- We will then compute at instance $[z_1, z_2, z_3, \dots, z_m] = W' * X + [b, b, \dots, b]$.
- This can be written in Python as:
 - $Z = \text{np.dot}(W.T, X) + b$ # Vectorization, then broadcasting, Z.shape = (1, m)
 - $A = 1 / (1 + \text{np.exp}(-Z))$ # Vectorization, A.shape = (1, m)
- Vectorizing Logistic Regression's Gradient Output:
 - $dz = A - Y$ # Vectorization, dz shape is (1, m)
 - $dw = \text{np.dot}(X, dz.T) / m$ # Vectorization, dw shape is (Nx, 1)
 - $db = \text{np.sum}(dz) / m$ # Vectorization, db shape is (1, 1)

Notes on Python and NumPy

- In NumPy, `obj.sum(axis = 0)` sums the columns while `obj.sum(axis = 1)` sums the rows.
- In NumPy, `obj.reshape(1,4)` changes the shape of the matrix by broadcasting the values.
- Reshape is cheap in calculations so put it everywhere you're not sure about the calculations.
- Broadcasting works when you do a matrix operation with matrices that doesn't match for the operation, in this case NumPy automatically makes the shapes ready for the operation by broadcasting the values.
- In general principle of broadcasting. If you have an (m,n) matrix and you add(+) or subtract(-) or multiply(*) or divide(/) with a (1,n) matrix, then this will copy it m times into an (m,n) matrix. The same with if you use those operations with a (m, 1) matrix, then this will copy it n times into (m, n) matrix. And then apply the addition, subtraction, and multiplication or division element wise.

- Some tricks to eliminate all the strange bugs in the code:
 - If you didn't specify the shape of a vector, it will take a shape of (m,) and the transpose operation won't work. You have to reshape it to (m, 1)
 - Try to not use the rank one matrix in ANN
 - Don't hesitate to use `assert(a.shape == (5,1))` to check if your matrix shape is the required one.
 - If you've found a rank one matrix try to run reshape on it.
 - Jupyter / IPython notebooks are so useful library in python that makes it easy to integrate code and document at the same time. It runs in the browser and doesn't need an IDE to run.
- To open Jupyter Notebook, open the command line and call: `jupyter-notebook` It should be installed to work.
- To Compute the derivative of Sigmoid:


```
s = sigmoid(x)
ds = s * (1 - s)    # derivative using calculus
```
- To make an image of (width,height,depth) be a vector, use this:
 - `v = image.reshape(image.shape[0]*image.shape[1]*image.shape[2],1)` #reshapes the image.
- Gradient descent converges faster after normalization of the input matrices.

General Notes

- The main steps for building a Neural Network are:
 - Define the model structure (such as number of input features and outputs)
 - Initialize the model's parameters.
 - Loop.
 - Calculate current loss (forward propagation)
 - Calculate current gradient (backward propagation)
 - Update parameters (gradient descent)
- Preprocessing the dataset is important.
- Tuning the learning rate (which is an example of a "hyperparameter") can make a big difference to the algorithm.
- kaggle.com is a good place for datasets and competitions.
- Pieter Abbeel is one of the best in deep reinforcement learning.

Shallow Neural Network

"Learn to build a neural network with one hidden layer, using forward propagation and back propagation"

Neural Network Overview

- In logistic regression we had:

$$\begin{array}{l} X1 \setminus \\ X2 \implies z = XW + B \implies a = \text{Sigmoid}(z) \implies l(a,Y) \\ X3 / \end{array}$$

- In neural networks with one layer we will have:

X1 \

X2 => $z_1 = XW_1 + B_1 \Rightarrow a_1 = \text{Sigmoid}(z_1) \Rightarrow z_2 = a_1W_2 + B_2 \Rightarrow a_2 = \text{Sigmoid}(z_2) \Rightarrow l(a_2, Y)$

X3 /

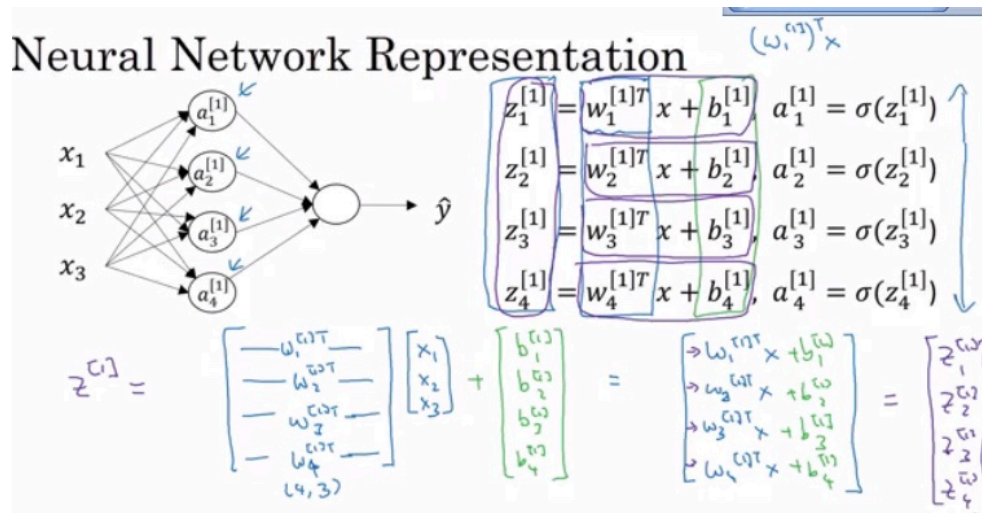
- X is the input vector (X1, X2, X3), and Y is the output variable (1 x 1)
- NN is stack of logistic regression objects

Neural Network Representation

- We will define the neural networks that has one hidden layer
- NN contains of input layers, hidden layers, output layers.
- Hidden layer means we can't see that layers' output in the training phase
- $a_0 = x$ (the input layer)
- a will represent the activation of the hidden neurons
- a_2 will represent the output layer
- we are talking about 2 layers NN. The input layer isn't counted.

Computing a Neural Network's Output

- Equations of Hidden Layers:



- Here are some informations about the last image:

- noOfHiddenNeuron = 4
- $N_x = 3$
- Shapes of the variables:
 - W_1 is the matrix of the first hidden layer, it has a shape of (noOfHiddenNeurons, N_x)
 - b_1 is the matrix of the first hidden layer, it has the shape of (noOfHiddenNeurons, 1)
 - z_1 is the result of the equation $z_1 = W_1 * X + b$, it has the shape of (noOfHiddenNeurons, 1)
 - a_1 is the result of the equation $a_1 = \text{sigmoid}(z_1)$, it has the shape of (noOfHiddenNeurons, 1)
 - W_2 is the matrix of the second hidden layer, it has a shape of (1, noOfHiddenNeurons)
 - b_2 is the matrix of the second hidden layer, it has the shape of (1, 1)

- z_2 is the result of the equation $z_2 = W_2 * a_1 + b$, it has the shape of (1, 1)
- a_2 is the result of the equation $a_2 = \text{sigmoid}(z_2)$, it has the shape of (1, 1)

Vectorizing across multiple examples

- Pseudo code for forward propagation for the 2 layers NN:

```

for i = 1 to m
    z[1, i] = W1*x[i] + b1      # shape of z[1, i] is (noOfHiddenNeurons,1)
    a[1, i] = sigmoid(z[1, i])  # shape of a[1, i] is (noOfHiddenNeurons,1)
    z[2, i] = W2*a[1, i] + b2   # shape of z[2, i] is (1,1)
    a[2, i] = sigmoid(z[2, i])  # shape of a[2, i] is (1,1)

```

- Lets say we have X on shape (Nx, m). So the new pseudo code:

```

Z1 = W1X + b1      # shape of Z1 (noOfHiddenNeurons,m)
A1 = sigmoid(Z1)    # shape of A1 (noOfHiddenNeurons,m)
Z2 = W2A1 + b2     # shape of Z2 is (1,m)
A2 = sigmoid(Z2)    # shape of A2 is (1,m)

```

- If you notice always m is the number of columns
- In the last example we can call $X = A_0$. So the previous step can be rewritten as:

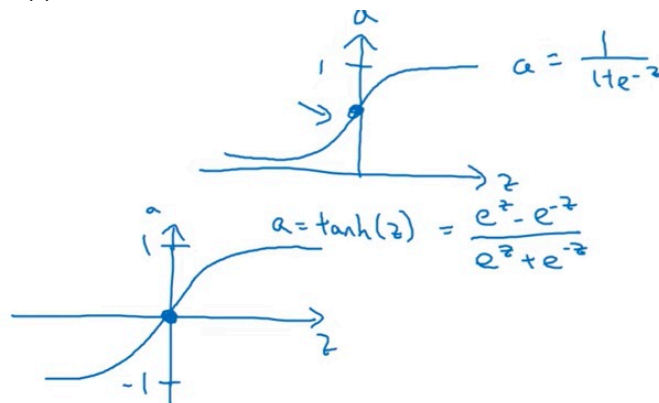
```

Z1 = W1A0 + b1      # shape of Z1 (noOfHiddenNeurons,m)
A1 = sigmoid(Z1)    # shape of A1 (noOfHiddenNeurons,m)
Z2 = W2A1 + b2     # shape of Z2 is (1,m)
A2 = sigmoid(Z2)    # shape of A2 is (1,m)

```

Activation Functions

- So far we are using sigmoid, but in some cases other functions can be a lot better.
- Sigmoid can lead us to gradient descent problem where the updates are so low.
- Sigmoid activation function range is [0, 1] $A = 1 / (1 + \text{np.exp}(-z))$ where z is the input matrix
- Tanh activation function range is [-1, 1] Shifted version of sigmoid function
 - In NumPy we can implement Tanh using one of these methods:
 - $A = (\text{np.exp}(z) - \text{np.exp}(-z)) / (\text{np.exp}(z) + \text{np.exp}(-z))$ where z is the input matrix
 - Or $A = \text{np.tanh}(z)$ where z is the input matrix



- It turns out the the tanh activation ususally works better than sigmoid activation function for hidden units because the mean of its output is closer to zero, and so it centers the data better for the next layer.
- Sigmoid or Tanh function disadvantage is that if the input is too small or too high, the slope will be near to zero which will cause us the gradient descent problem.
- One of the popular activation functions that solved the slow gradient descent is the RELU activation function.
 - $\text{ReLU} = \max(0, z)$ so if z is negative the slope is 0 and if z is positive the slope remains linear
- So here is some basic rule for choosing activation functions, if your classification is between 0 and 1, use the output activation as sigmoid and the others as ReLU
- Leaky ReLU activation function different from ReLU is that if the input is negative the slope will be so small. It works as ReLU but most people uses ReLU. $\text{Leaky_ReLU} = \max(0.01z, z)$ #the 0.01 can be a parameter for your algorithm.
- In NN you will decide a lot of choices like:
 - No of hidden layers.
 - No of neurons in each hidden layer.
 - Learning rate. (The most important parameter)
 - Activation functions.
 - And others.
- It turns out there are no guide lines for that. You should try all activation functions.

Why do you need non-linear activation functions?

- If we removed the activation function from our algorithm that can be called linear activation function.
- Linear activation function will output linear activations
 - Whatever hidden layers you add, the activation will always be linear like logistic regression (So its useless in alot of complex problems)
 - You might use linear activation function in one place - in the ouotput layer if the output is real numbers (regression problem). But even in this case if the output value is non-negative you could use ReLU instead.

Derivatives of activation functions

- Derivation of Sigmoid activation function:

$$\begin{aligned} g(z) &= 1 / (1 + \text{np.exp}(-z)) \\ g'(z) &= (1 / (1 + \text{np.exp}(-z))) * (1 - (1 / (1 + \text{np.exp}(-z)))) \\ g'(z) &= g(z) * (1 - g(z)) \end{aligned}$$

- Derivation of Tanh activation function:

$$\begin{aligned} g(z) &= (\text{e}^z - \text{e}^{-z}) / (\text{e}^z + \text{e}^{-z}) \\ g'(z) &= 1 - \text{np.tanh}(z)^2 = 1 - g(z)^2 \end{aligned}$$

- Derivation of ReLU activation function:

$$\begin{aligned} g(z) &= \text{np.maximum}(0, z) \\ g'(z) &= \begin{cases} 0 & \text{if } z < 0 \end{cases} \end{aligned}$$

$$1 \text{ if } z \geq 0 \}$$

- Derivation of leaky ReLU activation function:

$$g(z) = \text{np.maximum}(0.01 * z, z)$$

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

Gradient descent for Neural Works

- In this section we will have the full back propagation of the neural network (just the equations with no explanations).
- Gradient descent algorithm:

- NN parameters:
 - $n[0] = N_x$
 - $n[1] = \text{NoOfHiddenNeurons}$
 - $n[2] = \text{NoOfOutputNeurons} = 1$
 - $W1$ shape is $(n[1], n[0])$
 - $b1$ shape is $(n[1], 1)$
 - $W2$ shape is $(n[2], n[1])$
 - $b2$ shape is $(n[2], 1)$
- Cost function $I = I(W1, b1, W2, b2) = (1/m) * \text{Sum}(L(Y, A2))$
- Then Gradient descent:

Repeat:

Compute predictions ($y'[i], i = 0, \dots, m$)
 Get derivatives: $dW1, db1, dW2, db2$
 Update: $W1 = W1 - \text{LearningRate} * dW1$
 $b1 = b1 - \text{LearningRate} * db1$
 $W2 = W2 - \text{LearningRate} * dW2$
 $b2 = b2 - \text{LearningRate} * db2$

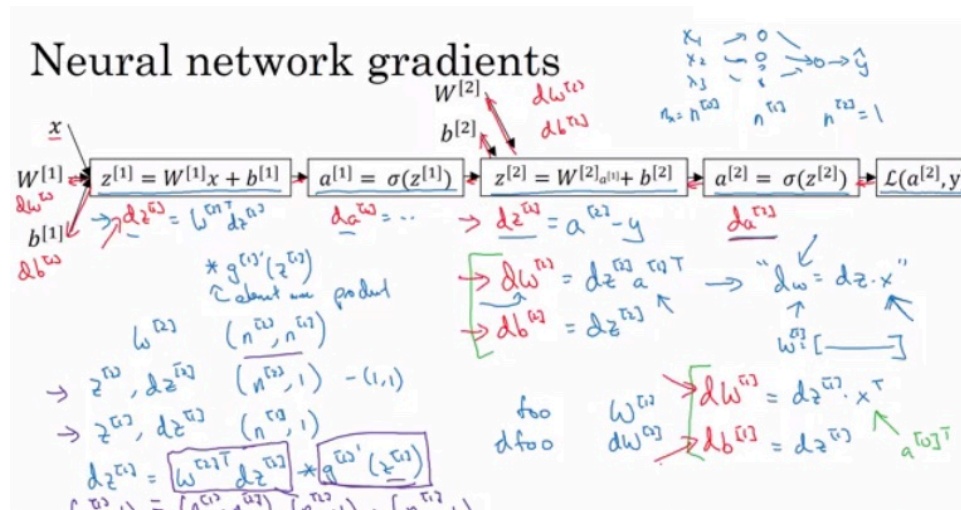
- Forward propagation:

$Z1 = W1A0 + b1$ # $A0$ is X
 $A1 = g1(Z1)$
 $Z2 = W2A1 + b2$
 $A2 = \text{Sigmoid}(Z2)$ # Sigmoid because the output is between 0 and 1

- Backpropagation (derivations):

$dZ2 = A2 - Y$ # derivative of cost function we used * derivative of the sigmoid function
 $dW2 = (dZ2 * A1.T) / m$
 $db2 = \text{Sum}(dZ2) / m$
 $dZ1 = (W2.T * dZ2) * g'1(Z1)$ # element wise product (*)
 $dW1 = (dZ1 * A0.T) / m$ # $A0 = X$
 $db1 = \text{Sum}(dZ1) / m$
 # Hint there are transposes with multiplication because to keep dimensions correct

- How we derived the 6 equations of the backpropagation:



Random Initialization

- In logistic regression it wasn't important to initialize the weights randomly, while in NN we have to initialize them randomly.
- If we initialize all the weights with zeros in NN it won't work (initializing bias with zero is OK):
 - all hidden units will be completely identical (symmetric) - compute exactly the same function
 - on each gradient descent iteration all the hidden units will always update the same
- To solve this we initialize the W 's with a small random numbers:

`W1 = np.random.randn((2,2)) * 0.01` # 0.01 to make it small enough

`b1 = np.zeros((2,1))` # its ok to have b as zero, it won't get us to the symmetry

breaking problem

- We need small values because in sigmoid (or tanh), for example, if the weight is too large you are more likely to end up even at the very start of training with very large values of Z . Which causes your tanh or your sigmoid activation function to be saturated, thus slowing down learning. If you don't have any sigmoid or tanh activation functions throughout your neural network, this is less of an issue.
- Constant 0.01 is alright for 1 hidden layer networks, but if the NN is deep this number can be changed but it will always be a small number.

Deep Neural Network

" Understand the key computations underlying deep learning, use them to build and train deep neural networks, and apply it to computer vision "

Deep L-Layer neural network

- Shallow NN is a NN with one or two layers.
- Deep NN is a NN with three or more layers.
- We will use the notation L to denote the number of layers in a NN
- $n[l]$ is the number of neurons in a specific layer l .
- $n[0]$ denotes the number of neurons input layer. $n[L]$ denotes the number of neurons in output layer.
- $g[l]$ is the activation function.
- $a[l] = g[l](z[l])$
- $w[l]$ weights is used for $z[l]$
- $x = a[0]$, $a[l] = y'$
- These were the notation we will use for deep neural network.
- So we have:
 - A vector n of shape $(1, \text{NoOfLayers}+1)$
 - A vector g of shape $(1, \text{NoOfLayers})$
 - A list of different shapes w based on the number of neurons on the previous and the current layer.
 - A list of different shapes b based on the number of neurons on the current layer.

Forward propagation in a Deep Network

- Forward propagation general rule for one input:

$$\begin{aligned} z[l] &= w[l].a[l-1] + b[l] \\ a[l] &= g[l].(z[l]) \end{aligned}$$

- Forward propagation general rule for m inputs:

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

- We can't compute the whole layers forward propagation without a for loop so its OK to have a for loop here.
- The dimensions of the matrices are so important you need to figure it out.

Getting your matrix dimension right

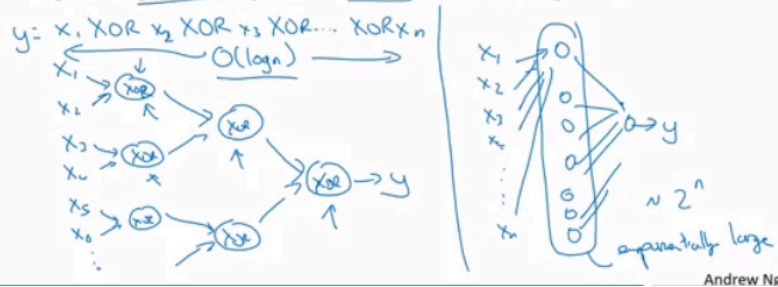
- The best way to debug your matrices demensions is by a pencil and paper.
- Dimension of W is $(n[l], n[l-1])$. Can be thought right to left.
- Dimension of b is $(n[l], 1)$
- dw has the same shape as W , while db is the same shape as b
- Dimension of $Z[l]$, $A[l]$, $dZ[l]$, and $dA[l:]$ is $(n[l], m)$

Why deep representations?

- Why deep NN works well, we will discuss this question in this section.
- Deep NN makes relations with data from simpler to complex. In each layer it tries to make a relation with the previous layer. E.g.:
 - Face recognition application:
 - Image ==> Edges ==> Face parts ==> Faces ==> desired face
 - Audio recognition application:
 - Audio ==> Low level sound features like (sss,bb) ==> Phonemes ==> Words ==> Sentences
- Neural Researchers think that deep neural networks "think" like brains (simple ==> complex)
- Circuit theory and deep learning:

Circuit theory and deep learning

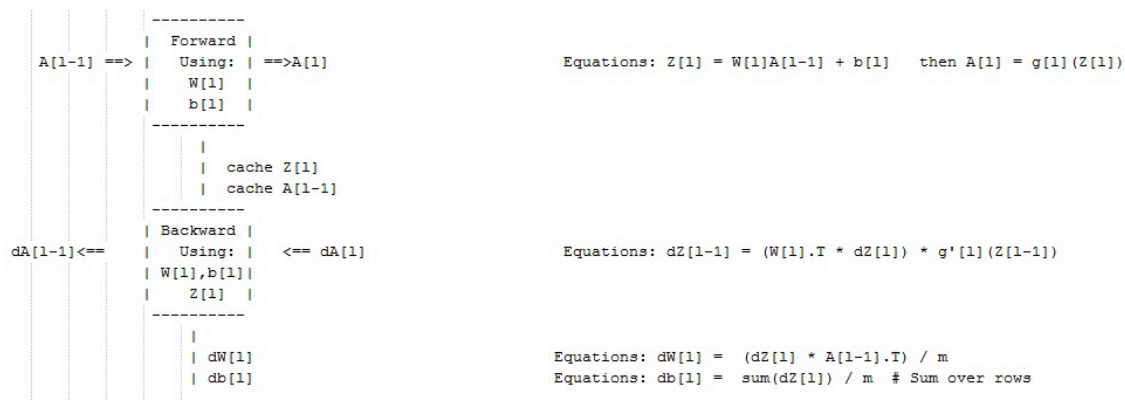
Informally: There are functions you can compute with a "small" L-layer deep neural network that shallower networks require exponentially more hidden units to compute.



- When starting on an application don't start directly by dozens of hidden layers. Try the simplest solutions (e.g. Logistic Regression), then try the shallow neural network and so on.

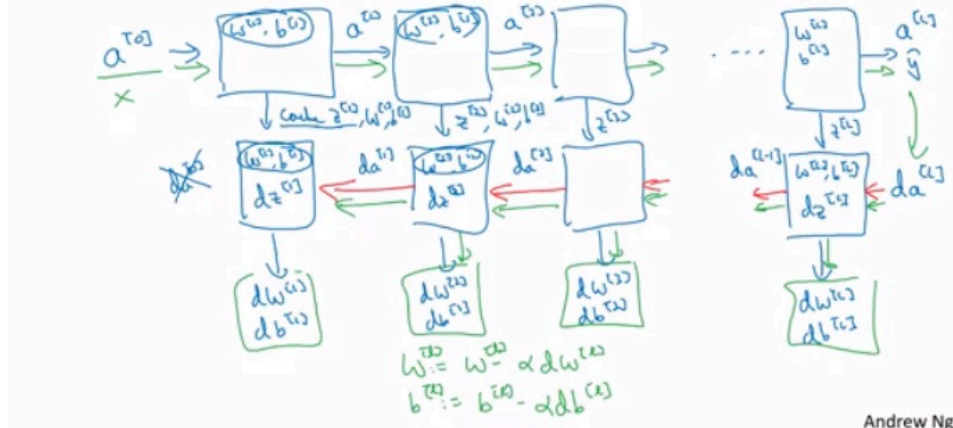
Building blocks of deep neural networks

- Forward and back propagation for layer l:



- Deep NN blocks:

Forward and backward functions



Forward and Backward Propagation

- Pseudo code for forward propagation for layer l:

Input $A[l-1]$
 $Z[l] = W[l]A[l-1] + b[l]$
 $A[l] = g[l](Z[l])$
 Output $A[l]$, cache($Z[l]$)

- Pseudo code for back propagation for layer l:

Input $da[l]$, Caches
 $dZ[l] = dA[l] * g'[l](Z[l])$
 $dW[l] = (dZ[l]A[l-1].T) / m$
 $db[l] = \text{sum}(dZ[l]) / m$
 $dA[l-1] = w[l].T * dZ[l]$
 Output $dA[l-1]$, $dW[l]$, $db[l]$

Dont forget axis=1, keepdims=True
 # The multiplication here are a dot product.

- If we have used our loss function then:

$$dA[L] = -(y/a) + ((1-y)/(1-a))$$

Parameters vs Hyperparameters

- Main parameters of the NN is W and b
- Hyper parameters (parameters that control the algorithm) are like:
 - Learning rate.
 - Number of iteration.
 - Number of hidden layers L .
 - Number of hidden units n .
 - Choice of activation functions.

- You have to try values yourself of hyper parameters.
- In the earlier days of DL and ML learning rate was often called a parameter, but it really is (and now everybody call it) a hyperparameter.
- On the next course we will see how to optimize hyperparameters.

What does this have to do with the brain

- The analogy that "It is like the brain" has become really an oversimplified explanation.
- There is a very simplistic analogy between a single logistic unit and a single neuron in the brain.
- No human today understand how a human brain neuron works.
- No human today know exactly how many neurons on the brain.
- Deep learning in Andrew's opinion is very good at learning very flexible, complex functions to learn X to Y mappings, to learn input-output mappings (supervised learning).
- The field of computer vision has taken a bit more inspiration from the human brains then other disciplines that also apply deep learning.
- NN is a small representation of how brain work. The most near model of human brain is in the computer vision (CNN)