

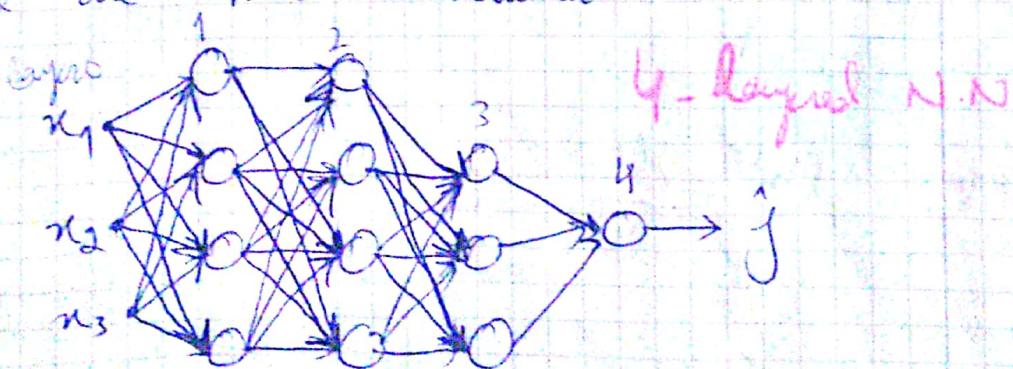
Week #4: Deep Neural Network

Deep L-layer Neural Network:-

Previously we saw 1 layered network where we implemented the logistic regression and later in week 3 ~ 2-layered network
→ Remember \Rightarrow no of layers = total hidden layers plus output layer

Let's work out some notation that we'll be using here.

e.g. here we have a network.



$L = 4 \Rightarrow$ (No. of layers)

$n^{[L]}$ = no. of units in layer
 $n^{[1]} = 4$, $n^{[2]} = 4$, $n^{[3]} = 3$, $n^{[4]} = n_x = 1$

$a^{[l]}$ = activation in layer l $n^{[l]} = n_x = 3$

$a^{[l]} = g^{[l]}(z^{[l]})$, $w^{[l]}$ = weights for $z^{[l]}$
 $b^{[l]}$ = bias for $z^{[l]}$.

Forward Propagation in Deep network :-

Considering the last example of network we designed.
let's go over how forward propagation will look like
for one training example.

Given $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$ for first layer we can compute
 $z^{[1]} = w^{[1]}x + b^{[1]}$

the activation for the same layer is given by:-

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

How about layer 2:-

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]}) \text{ and so on}$$

Until we get to output layer

$$z^{[L]} = w^{[L]}a^{[L-1]} + b^{[L]}$$

$$a^{[L]} = g^{[L]}(z^{[L]})$$

Now we have seen it for one training example
let's go through it for m training examples via
Vectorization.

Vectorization for m training examples:-

$$Z^{[1]} = W^{[1]} A^{[0]} + b^{[1]} \quad X = A^{[0]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

$$Z = \begin{bmatrix} Z^{1} & Z^{[1](2)} & \dots & Z^{[1](m)} \\ \vdots & \vdots & \ddots & \vdots \\ Z^{[2](1)} & Z^{2} & \dots & Z^{[2](m)} \\ \vdots & \vdots & \ddots & \vdots \\ Z^{[L](1)} & Z^{[L](2)} & \dots & Z^{[L](m)} \end{bmatrix}$$

1st training example mth training ex.

In this process we end up with

$$\hat{Y} = g(Z^{[L]}) = A^{[L]}$$

Getting your matrix dimensions right:-

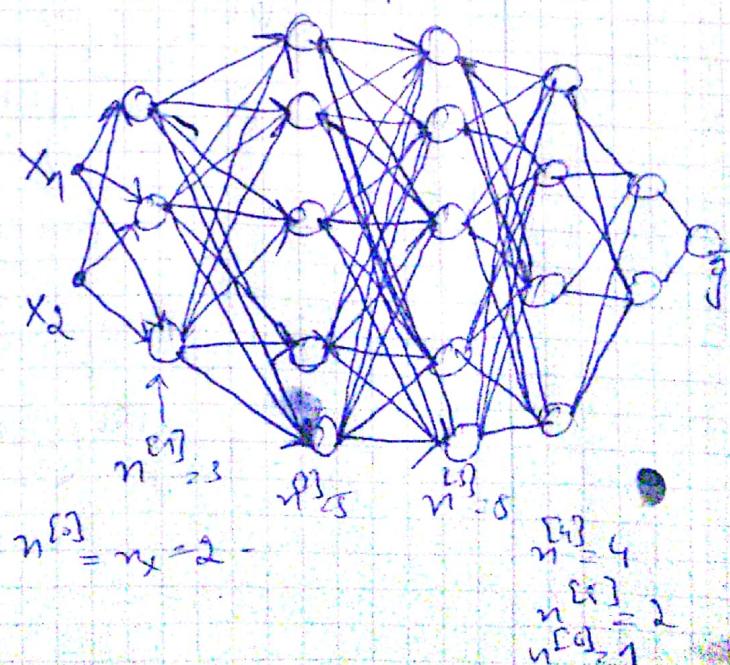
getting the matrix dimensions right is extremely important when implementing your neural network.

So working that out on a piece of paper can really save you a lot of headache. Let's see how to do that.

→ let's say we have a 5-layered neural network

$$L = 5$$

$$Z^{[1]} = W^{[1]} \cdot X + b^{[1]}$$



Let's forget the bias term for a while.

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$(3, 1) \quad (3, 2) \quad (2, 1) \quad (3, 1) \quad W^{[1]} = (n^{[1]}, n^{[2]})$$

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} \vdots & \vdots \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$W^{[2]} = (5, 3)$$

$$W^{[3]} \in \mathbb{R}^{3 \times 5}$$

$$Z^{[3]} = W^{[3]} a^{[3]} + b^{[3]}$$

$$(5, 1) \quad (5, 5), (5, 1) \quad (5, 1)$$

$$W^{[4]} \in \mathbb{R}^{2 \times 4}$$

$$Z^{[4]} = W^{[4]} a^{[4]} + b^{[4]}$$

$$(2, 1) \quad (2, 4) \quad (4, 2) \quad (2, 1)$$

$$W^{[4]} \in \mathbb{R}^{4 \times 5}$$

$$Z^{[4]} = W^{[4]} a^{[4]} + b^{[4]}$$

$$(4, 1) = (4, 5) (5, 1) (4, 1)$$

$$W^{[5]} \in \mathbb{R}^{1 \times 2}$$

$$Z^{[5]} = W^{[5]} a^{[5]} + b^{[5]}$$

$$(1, 1) = (1, 2) (2, 1) (1, 1)$$

In backward propagation $dW^{[l]}$'s dimension = W^l 's

$d b^l$'s ... = b^l 's

In Vectorized version:

$Z^{[1]} = W^{[1]} \cdot X + b^{[1]} \Rightarrow$ this is still $(n^{[1]}, 1)$, but
 $(n^{[0]}, n) \cdot (n^{[0]}, n^{[1]}) \cdot (n^{[1]}, n)$ it expands into $(n^{[1]}, n)$ by
 python broadcasting and
 added into the $W^{[1]} X$ matrix
 element wise.

$$Z^{[1]}, a^{[0]} \in (\mathbb{R}^{n^{[0]}, 1})$$

$$Z^{[0]}, A^{[0]} : (n^{[0]}, n)$$

$$dZ^{[0]}, dA^{[0]} : (n^{[0]}, n)$$

$$l=0 \quad A^{[0]} = X = (n^{[0]}, n)$$

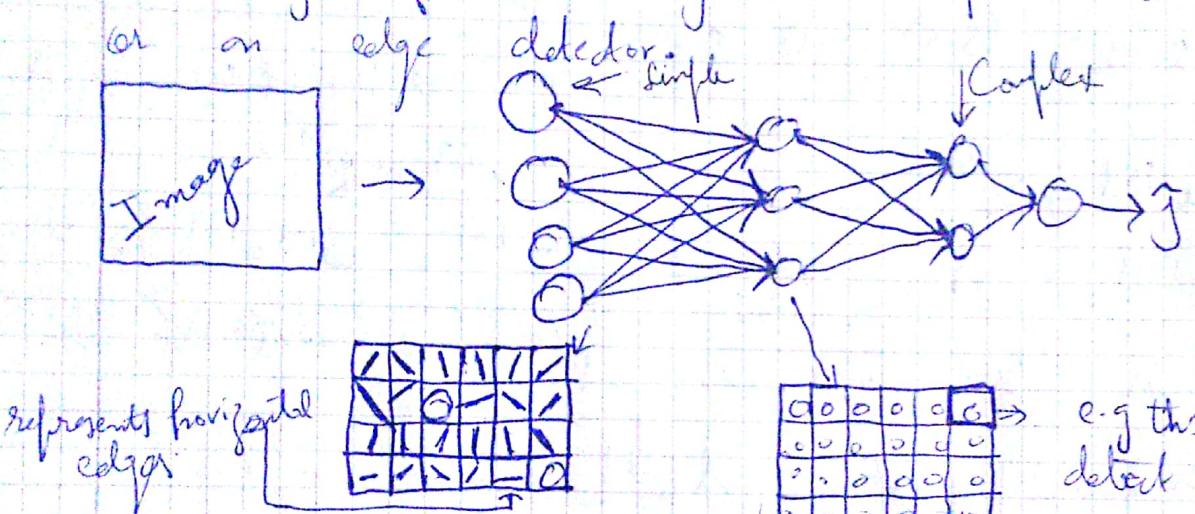
Why deep neural networks work better?

→ Let's go through a couple of examples and try to develop an intuition why they might work well.

So, first. What is a deep network computing?

→ If you are building a system of face recognition or face detection, here is what a deep neural network could be doing:-

→ You input a face's ^{image} and then the first layer could be thought of as may be a feature detector or an edge detector.



→ first layer of a neural network may be thought of as edge detector. Now we have the detected edges. Now In this example I am computing which the next layer can take the first 20 hidden units try to the detected edges and compute on this image. So this hidden unit is trying to figure out where are the edges of the orientation in the image. Now by taking together those computed parts of face you can take these to detect different types of faces.

→ So, you can think that first layer is being used to compute simple functions like detecting features and later layers are trying to compute more complex functions by taking the previously outputted simple functions.

Main Intuitive takeaway:-

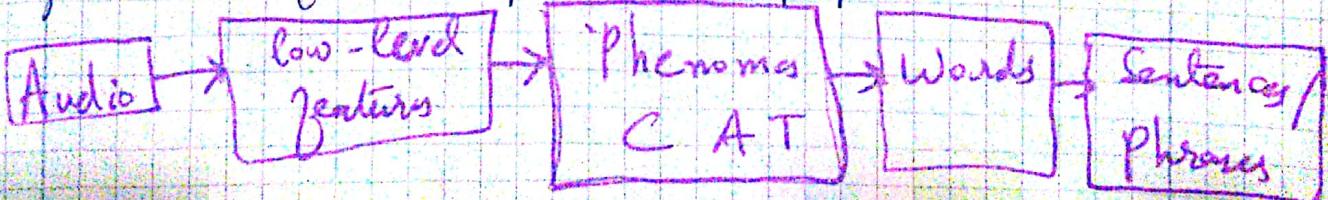
Finding simple things like edges then building them up, Composing them together to detect more complex things.

→ This type of simple to complex hierarchical representation applies in other types of data rather than just imagery data as well.

e.g. If you want to build speech recognition system

→ If you input an audio clip then may be first level of neural network is being used to compute low level audio features like audio waveform, like is it going up or down or is it white noise or pitch.

→ Then by those low level features you learn to develop some high level things like sounds. In linguistics they are called Phonemes. Composing these sounds together to recognize words and in turn composing words together to form complete sentences/phrases.



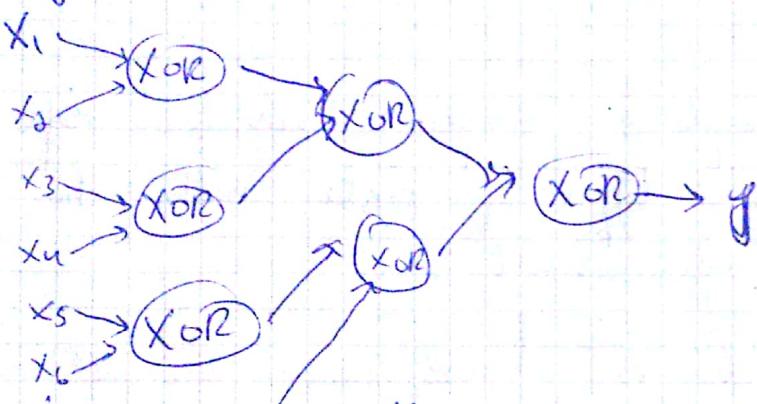
Circuit Theory and Deep Learning:-

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.

Let's motivate this by an example:-

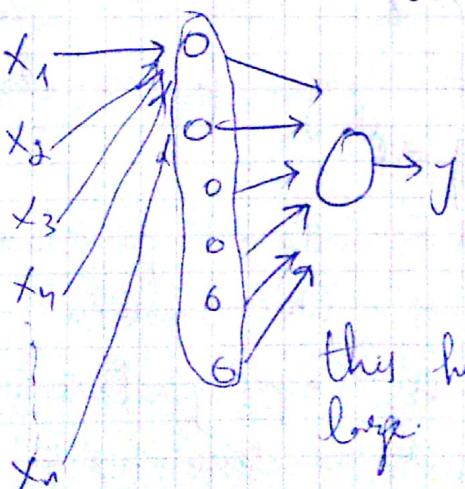
$$x_1 \text{ XOR } x_2 \text{ XOR } x_3 \text{ XOR } \dots \text{ XOR } x_n$$

One way to do that is to compute it like this:



So, to compute all this, the depth of the network will be of the order $O(\log(n))$

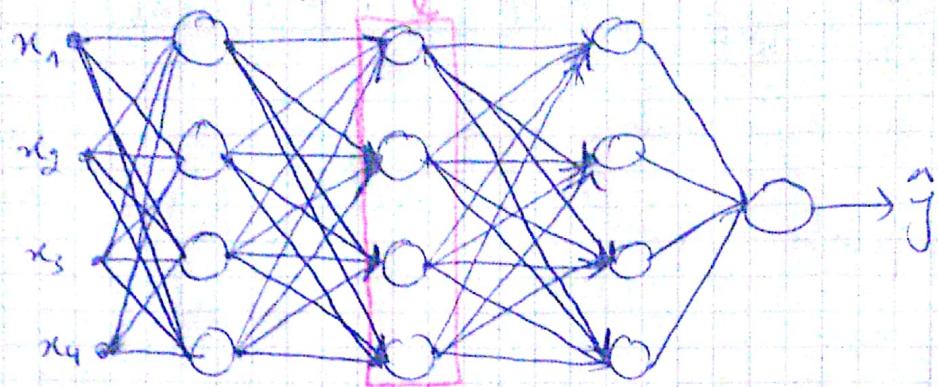
Now with shallower network if we do it like this:-



this hidden layer will be exponentially large: $2^{n-1} \Rightarrow O(2^n)$

Building blocks of deep neural Networks:

We have already seen basic building blocks for forward propagation and backward propagation. Let's see how you can take all those building blocks together and build a deep neural network.



Let's consider the above drawn neural network and focus on rectangular layer
layer l : $w^{[l]}, b^{[l]}$

Forward: I/p $a^{[l-1]}$, output $a^{[l]}$

$$\text{O/p } z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \text{ Cache } z^{[l]}$$

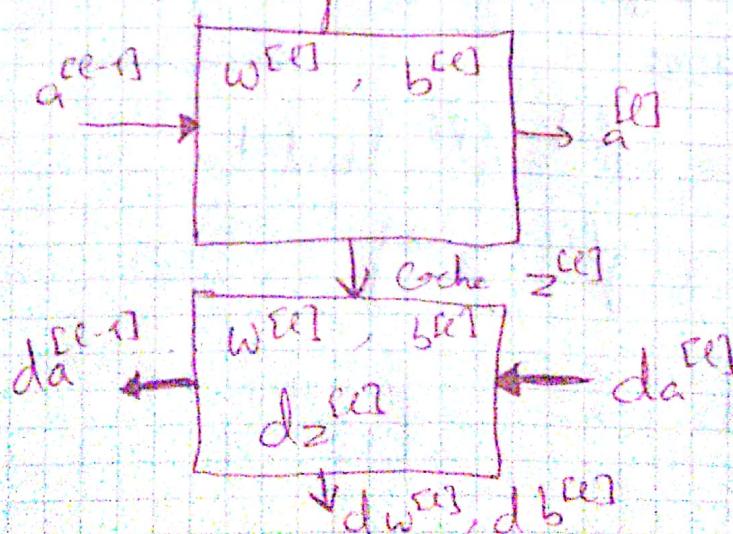
$$a^{[l]} = g^{[l]}(z^{[l]})$$

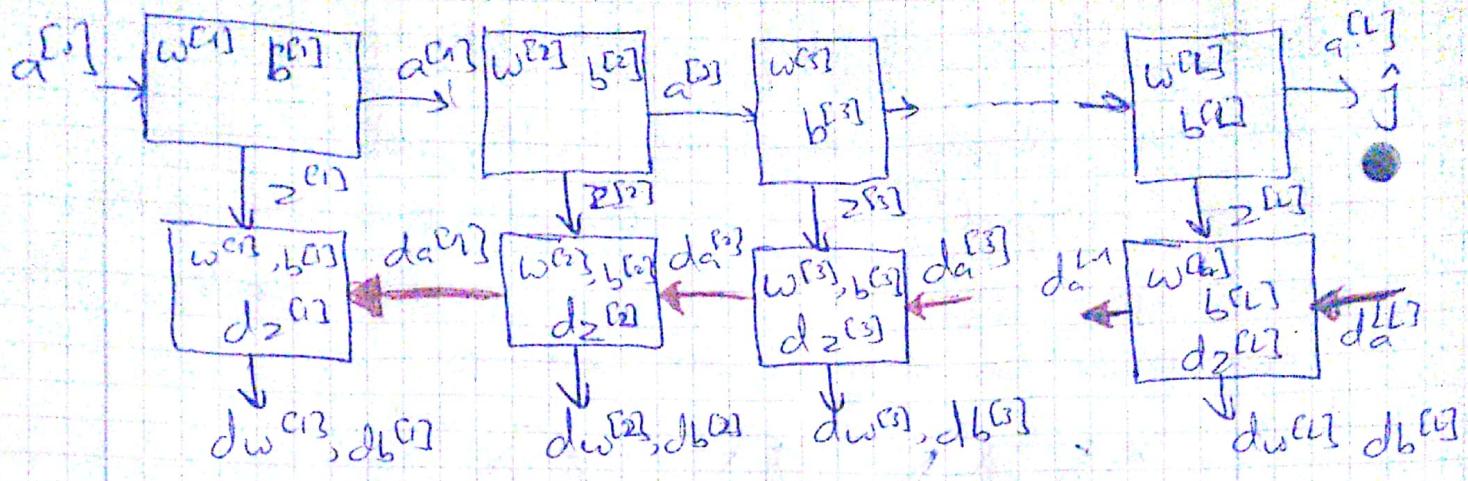
Backward propagation: I/p $d_a^{[l]}$, Cache ($z^{[l]}$)

$$\text{O/p } d_a^{[l-1]}, dw^{[l]}, db^{[l]}$$

Layer l

To summarize





Red denotes back propagation.

$$\text{update} \quad \begin{cases} w^{[l]} = w^{[l]} - \alpha dw^{[l]} \\ b^{[l]} = b^{[l]} - \alpha db^{[l]} \end{cases}$$

Forward propagation step and backward propagation step:-

Let's see how we can develop code to implement these algorithms. Let's start with forward propagation step.

→ Forward propagation for layer l :-

$$\begin{aligned} &\rightarrow \text{Input } a^{[l-1]} \\ &\rightarrow \text{Output } a^{[l]}, \text{ cache } (z^{[l]}) \end{aligned}$$

$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

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

Vectorized Version:

$$z^{[l]} = w^{[l]} \cdot A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$



→ Backward propagation for layer l:-

→ Input $d_a^{[l]}$

→ Output $d_z^{[l-1]}, d_w^{[l]}, d_b^{[l]}$

$$\cdot d_z^{[l]} = d_a^{[l]} \times g^{[l]}'(z^{[l]})$$

$$\cdot d_w^{[l]} = d_z^{[l-1]} \cdot a^{[l-1]}$$

$$\cdot d_b^{[l]} = d_z^{[l]}$$

$$\cdot d_a^{[l-1]} = W^{[l]T} \cdot d_z^{[l]}$$

$$\cdot d_z^{[l]} = W^{[l+1]T} \cdot d_z^{[l+1]} \times g^{[l]}'(z^{[l]})$$

Let's write the vectorized version:-

$$dZ^{[l]} = dA^{[l]} \times g^{[l]}'(z^{[l]})$$

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

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

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

Parameters Vs Hyperparameters :-

Being a good neural network developer means that, you not only manage your parameters well, but also your hyperparameters.

So what are hyperparameters? Let's discuss that.

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

Hyperparameters: Learning rate α
iterations

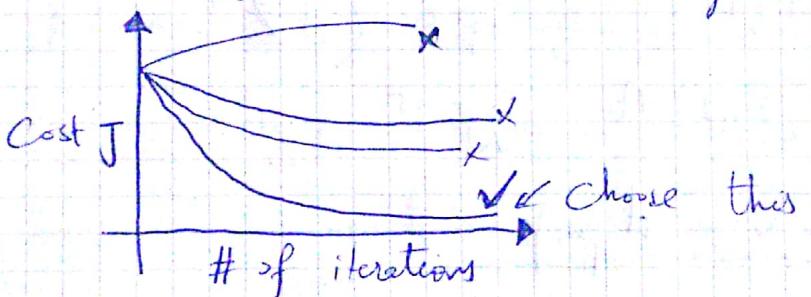
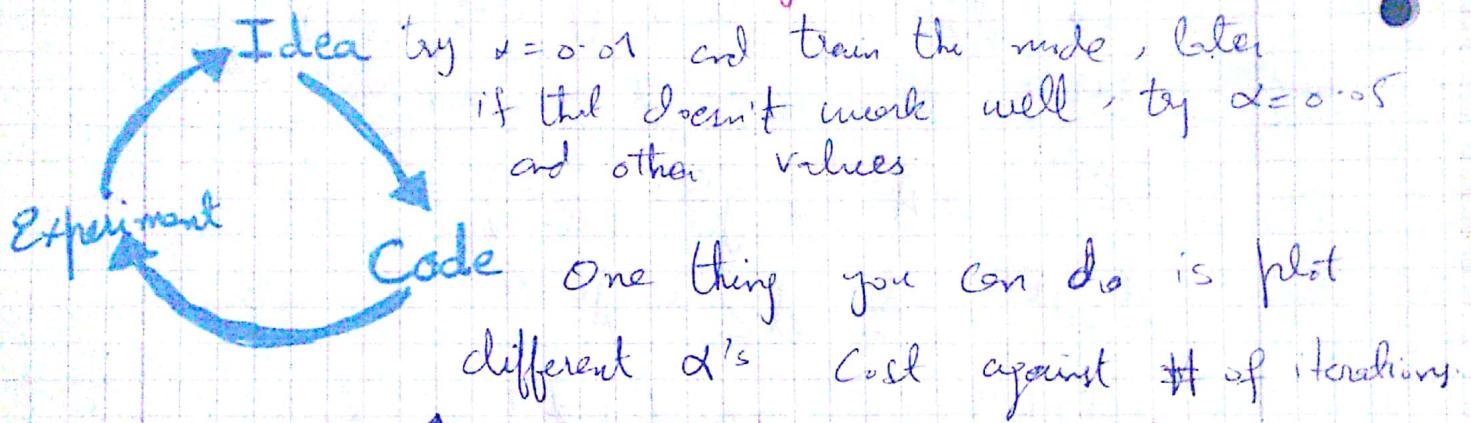
of hidden layers

hidden units $n[3], n[2], \dots$

Choice of activation function

- as all these ~~other~~ things control the values of your parameters w, b so we call them hyperparameters
- later in the course we will see other hyperparameters like momentum, minibatch, size, regularization.

Applied Deep Learning is very empirical process



- figuring out hyperparameter values at the start of training is very difficult, so you should try different values and go around the above blue cycle like hit and trial and see what works.

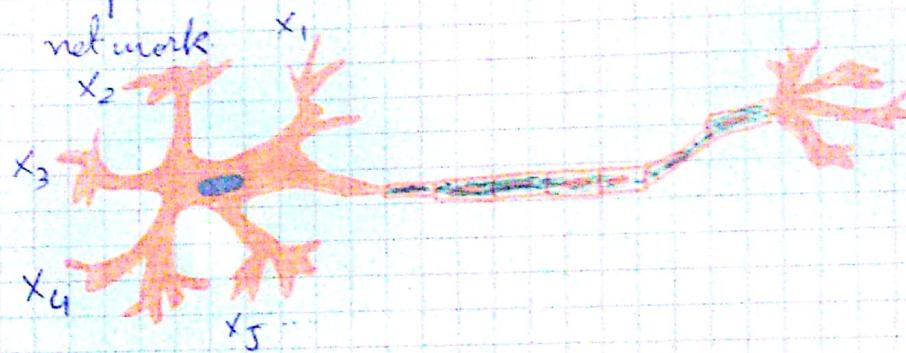
What does Deep Learning have to do with human brain?

In short \Rightarrow not a whole lot ;)

when you implement a neural network you implement forward propagation and backpropagation

the equations give very less intuition about the whole things relation with the brain

\rightarrow If we look at the neurons cell of human brain structure of that is more or less the same as neural network



There is very little information today that how one neuron operate or if they use an algorithm like backprop. and forward prop to operate or if there is totally different learning algorithm.