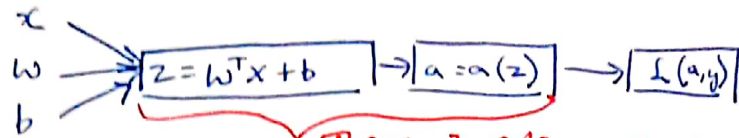
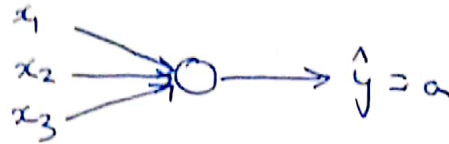


Week 3: One Hidden Layer NN

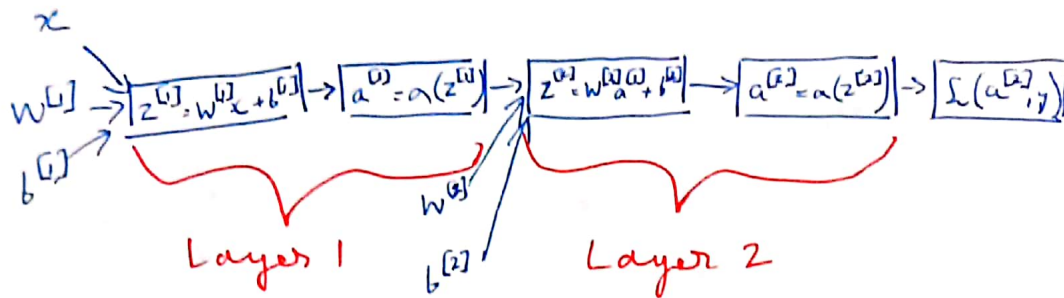
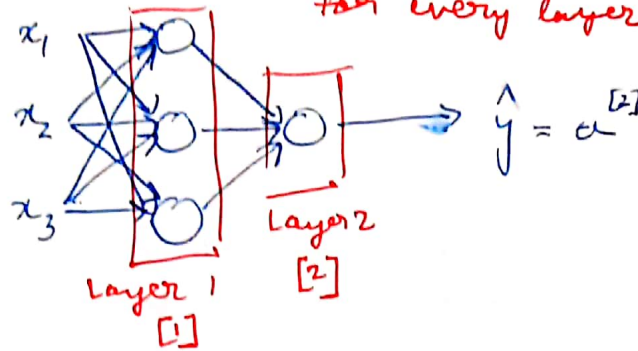
Overview

Single Neuron:

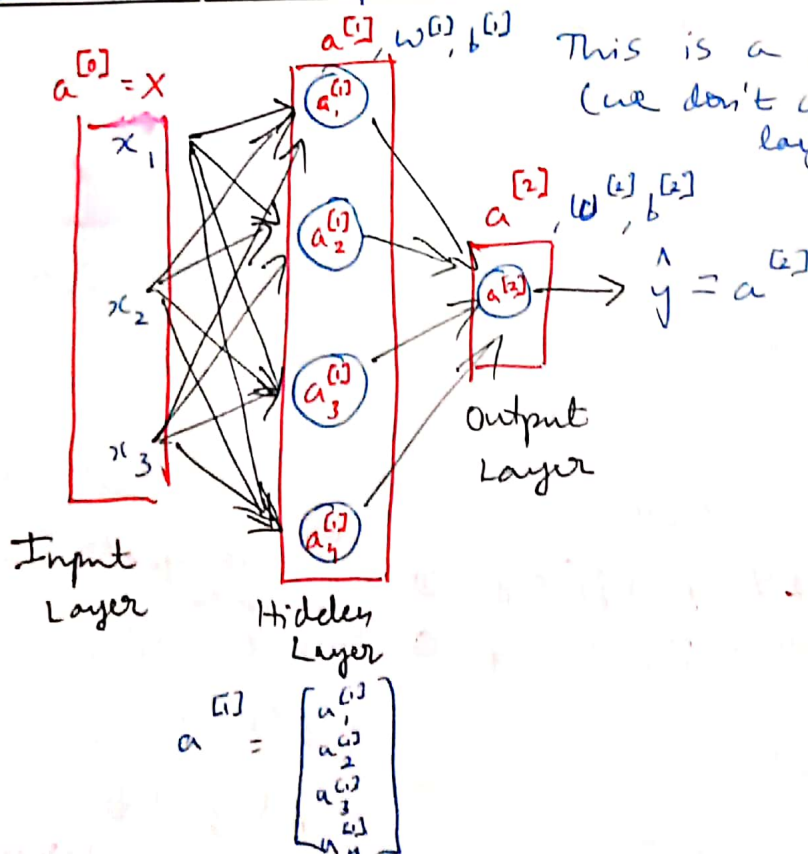


These 2 steps are required for every layer

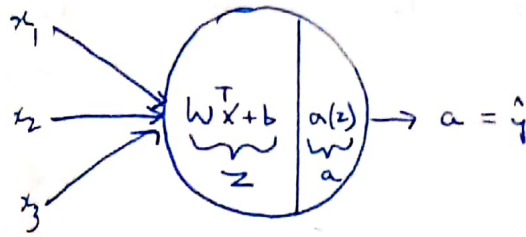
Simple NN:



Neural Network Representation



Computing a Neural Network's Output

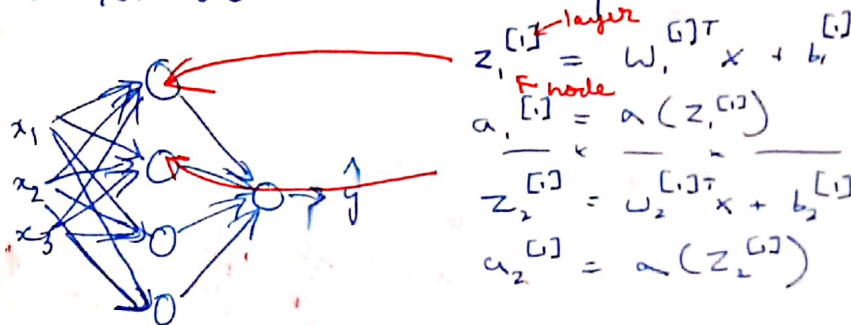


logistic regression represents 2 steps:

$$\rightarrow z = W^T x + b$$

$$\rightarrow a = \alpha(z)$$

So a NN does this a lot of times



$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]}$$

$$a_1^{[1]} = \alpha(z_1^{[1]})$$

$$z_2^{[1]} = W_2^{[1]T} x + b_2^{[1]}$$

$$a_2^{[1]} = \alpha(z_2^{[1]})$$

\therefore Finding equations for all the nodes,

$$z_1^{[1]} = W_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \alpha(z_1^{[1]})$$

$$z_2^{[1]} = W_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \alpha(z_2^{[1]})$$

$$z_3^{[1]} = W_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \alpha(z_3^{[1]})$$

$$z_n^{[1]} = W_n^{[1]T} x + b_n^{[1]}, \quad a_n^{[1]} = \alpha(z_n^{[1]})$$

$$Z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ \vdots \\ z_n^{[1]} \end{bmatrix}$$

$$W^{[1]} = \begin{bmatrix} -W_1^{[1]T} & - \\ -W_2^{[1]T} & - \\ -W_3^{[1]T} & - \\ \vdots & \vdots \\ -W_n^{[1]T} & - \end{bmatrix}$$

$$x = a^0 = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ \vdots \\ b_n^{[1]} \end{bmatrix}$$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ \vdots \\ a_n^{[1]} \end{bmatrix}$$

Given input x ,

$$z^{[1]} = W^{[1]} x + b^{[1]} \quad \leftarrow \text{since } x = a^{[0]}$$

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

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

$$a^{[2]} = \alpha(z^{[2]})$$

we are not doing W^T since in $W^{[2]}$ we have W^T

Vectorising across Multiple Examples

In the previous section, we saw how to compute the prediction on a neural network for a single training example. In this section, we will vectorize across multiple training examples.

In looping it would be as follows,

for $i = 1$ to m :

$$z^{[1]}(i) = W^{[1]}x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

$$z^{[2]}(i) = W^{[2]}a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

where $z^{[1]}(i)$ \uparrow training example i
layer

Consider we stack up these vectors,

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

For intuition

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix}$$

nodes
3

Training examples
1 \dots m

$$Z^{[1]} = \begin{bmatrix} | & | & \dots & | \\ z^{[1]}(1) & z^{[1]}(2) & \dots & z^{[1]}(m) \\ | & | & \dots & | \end{bmatrix}$$

$$A^{[1]} = \begin{bmatrix} | & | & \dots & | \\ a^{[1]}(1) & a^{[1]}(2) & \dots & a^{[1]}(m) \\ | & | & \dots & | \end{bmatrix}$$

In the previous section we only considered one training example so took $\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

For multiple examples,

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

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

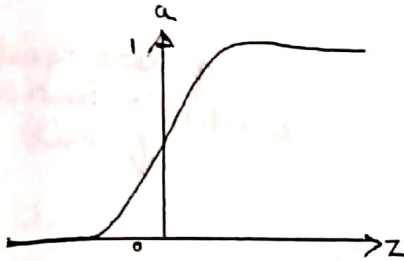
$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

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

Activation Functions

Till now we have been using the sigmoid function σ . However there are more superior options as well.

Sigmoid



$$a = \frac{1}{1 + e^{-z}}$$

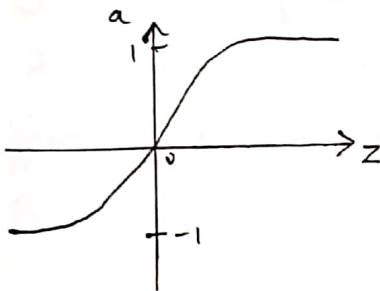
- Gives value between 0 to 1
so only use this for the output ~~function~~ layer

For Output Layer

(If ^{binary} classification $y = \{0, 1\}$)

- Worst option to use for hidden layers

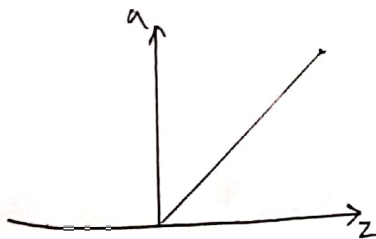
tanh



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

- A better option than Sigmoid, but still sucks
- Give values between -1 to 1
- 0 centred
- Don't use lol

ReLU (Rectified Linear Unit)



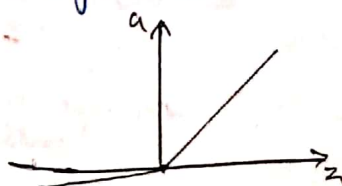
$$a = \max(0, z)$$

- The perfect function for hidden layers

For hidden layers

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

Leaky ReLU



$$a = \max(0.01z, z)$$

- If neurons start dying (giving 0 value for everything), then use this

Why can't we use linear activation functions?

Linear function - slope is constant ✗

Non-linear function - slope varies ✗

- The activation function has to be a non-linear function (Ex: σ , tanh, ReLU)

- If it is linear, let's say

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

$$a^{[1]} = z^{[1]} \text{ instead of } a^{[1]} = g^{[1]}(z^{[1]})$$

g here represents
an activation function like
 $\sigma(z)$

Then the output ~~will be~~ of the NN
will be $w'x + b' \Rightarrow (w^{[2]} \underbrace{w^{[1]}}_{w'})x + (w^{[2]} \underbrace{b^{[1]}}_{b'}) + b^{[2]}$

Therefore having hidden layers will
become pointless as they won't work

\therefore Only use Non-linear activation functions.

Derivatives of Activation Functions

① Sigmoid Function:

$$g(z) = \frac{1}{1+e^{-z}}$$

since $a = g(z)$

$$\therefore g'(z) = \frac{d}{dz} g(z) = \frac{1}{1+e^{-z}} \left(1 - \frac{1}{1+e^{-z}}\right) = g(z)(1-g(z)) = a(1-a)$$

② Tanh Function:

$$g(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

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

according to
maths

③ ReLU Function:

$$g(z) = \max(0, z)$$

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

undefined if $z=0$
But in coding,
1 if $z \geq 0$

④ Leaky ReLU Function:

$$g(z) = \max(0.01z, z)$$

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

in software it
works

if coding

Gradient Descent for NN

parameters: $w^{[1]}$, $b^{[1]}$, $w^{[2]}$, $b^{[2]}$

Cost Function: $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y} - y)$
 $\square a^{[2]}$

Gradient Descent:

Repeat Σ

Compute predictions ($\hat{y}^{(i)}$, $i=1 \dots m$)

$$dw^{[1]} = \frac{\partial J}{\partial w^{[1]}}, \quad db^{[1]} = \frac{\partial J}{\partial b^{[1]}}, \dots$$

$$w^{[1]} := w^{[1]} - \alpha dw^{[1]}$$

$$b^{[1]} := b^{[1]} - \alpha db^{[1]}$$

...

}

Formulas for computing derivatives:

Forward propagation:

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

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

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

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

Here for simplification
 $= \sigma(z^{[2]})$

Backward propagation:

$$dz^{[2]} = A^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]T}$$

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

$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

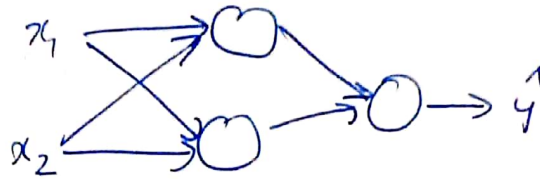
$$dw^{[1]} = \frac{1}{m} dz^{[1]} x^T$$

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

Match the matrix $(n^{[2]}, 1)$
instead of $(n^{[2]},)$
↓

Random Initialization:

If you set W to 0 i.e. ~~is~~



$$W_{\text{all}}^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad \text{for}$$

Then $dW = \begin{bmatrix} u & u \\ u & u \end{bmatrix}$ - The change of all the neurons in the layer will end up being same.

\therefore Therefore W values will always be the same and the NN will become symmetrical. Therefore having n no. of nodes won't matter since it'll be equivalent to having 1 node (since all values same).

We use ~~of~~ random initializing to break the symmetry:

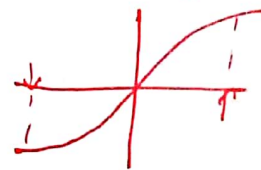
$$W^{[1]} = \text{np.random.randn}(2,2) * 0.01$$

$$b^{[1]} = \text{np.zeros}(2,1)$$

$$W^{[2]} = \dots$$

$$b^{[2]} = 0$$

\uparrow
This needs to be a small number so that
 $a^{[1]} = g^{[1]}(z^{[1]})$



The value of $a^{[1]}$ doesn't end up here