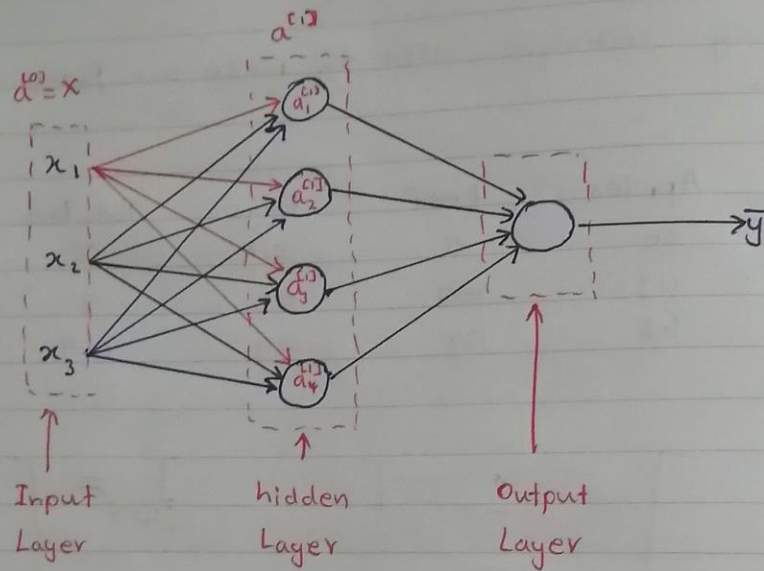


# Week 3

## Shallow Neural Networks

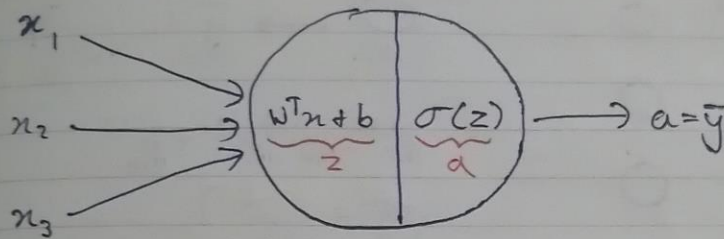
### Neural Networks Overview



- \* When we count number of layers in neural network, we don't count input layer. Therefore, above network has only 2 Layers (hidden layer and output layer)
- \* In the ~~train~~ training set, the true values ~~a~~ for nodes in hidden layer are not observed. It means you don't see what should be in the training set. that is the meaning of hidden layer.

## Computing a Neural Network's Output

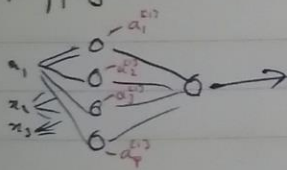
\* Let's look in to one node of hidden layer and see what happened in that node



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

$$a = \sigma(z)$$

\* Now, apply above theory into ~~to~~ our previous neural network



$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

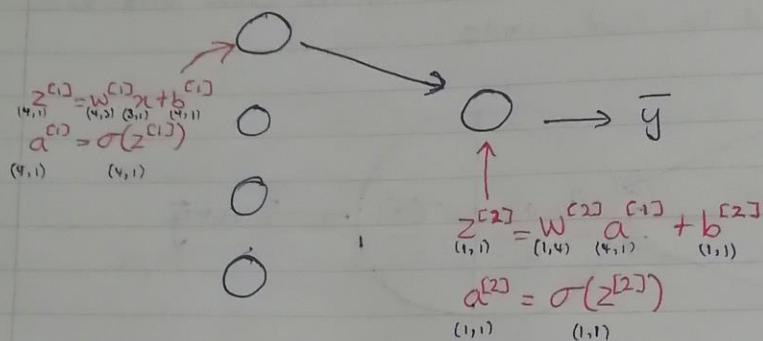
$$z_2^{(1)} = w_2^{(1)T} x + b_2^{(1)}, \quad a_2^{(1)} = \sigma(z_2^{(1)})$$

$$z_3^{(1)} = w_3^{(1)T} x + b_3^{(1)}, \quad a_3^{(1)} = \sigma(z_3^{(1)})$$

$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}, \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

$$\begin{bmatrix} - & w_1^{(1)T} & - \\ - & w_2^{(1)T} & - \\ - & w_3^{(1)T} & - \\ - & w_4^{(1)T} & - \end{bmatrix}_{(4 \times 3)} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}_{(3 \times 1)} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}_{(4 \times 1)} = \begin{bmatrix} w_1^{(1)T} x + b_1^{(1)} \\ w_2^{(1)T} x + b_2^{(1)} \\ w_3^{(1)T} x + b_3^{(1)} \\ w_4^{(1)T} x + b_4^{(1)} \end{bmatrix} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix}_{(4 \times 1)}$$

\* After ~~above~~ ~~above~~ that we are going to discuss what happen ~~to~~ in output layer.



### Vectorizing Across Multiple Examples

$$\begin{aligned}
 x &\longrightarrow a^{(2)} = \bar{y} \\
 x^{(1)} &\longrightarrow a^{(2)(1)} = \bar{y}^{(1)} \\
 x^{(2)} &\longrightarrow a^{(2)(2)} = \bar{y}^{(2)} \\
 &\vdots \\
 x^{(m)} &\longrightarrow a^{(2)(m)} = \bar{y}^{(m)}
 \end{aligned}$$

$a^{(2)(i)}$   
 $\uparrow$  training example  $i$   
 Layer 2

for  $i=1$  to  $m$ :

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

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

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

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

$$Z^{(1)} = W^{(1)} X + b^{(1)}$$

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

$$Z^{(2)} = W^{(2)} A^{(1)} + b^{(2)}$$

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

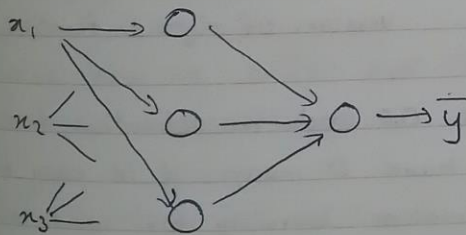
$$X = \begin{bmatrix} | & | & | & | \\ x^{[1]} & x^{[2]} & \dots & x^{[m]} \\ | & | & | & | \end{bmatrix} \quad (n, m)$$

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

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

hidden units  
Total examples

## Activation Function



Given  $x$ :

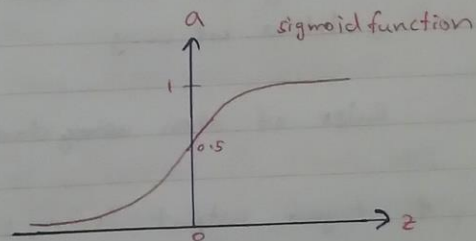
$$z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

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

sigmoid function  
(activation function)



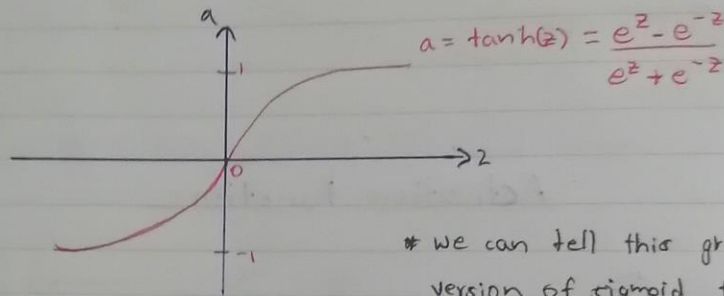


- \* But in general cases we can use different function called "g" of z. g could be a non linear function that may not be the sigmoid function.

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

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

- \* Activation function that almost always better than the sigmoid function. (tanh function or hyperbolic tangh function)



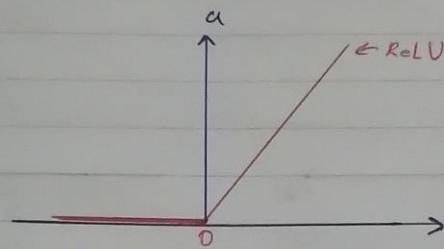
\* we can tell this graph is lifted version of sigmoid function.

- \* In sigmoid function's, mean is 0.5 and in activation function's mean is 0. If your data have 0 mean it is easy to centering your data and it makes learning for next layer a little bit easy.

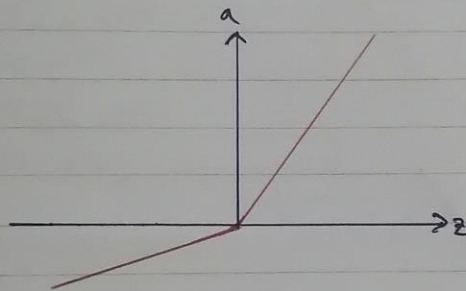
Rules of when ~~using~~ choosing activation function

- \* if your output is 0, 1 ~~or if~~ (if you using binary classification) then the sigmoid activation function is the most suitable choice for the output layer.
- \* If you are not sure what to use for hidden layer, Use the ReLU activation function

- \* But one disadvantage of ReLU is that the derivative is equal to 0.

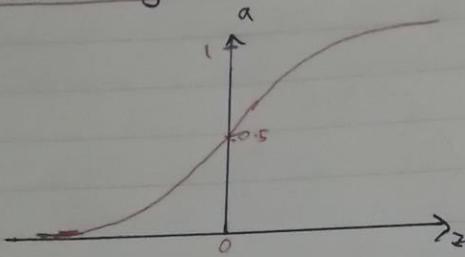


- \* When  $z$  is a negative we can use "Leaky ReLU".



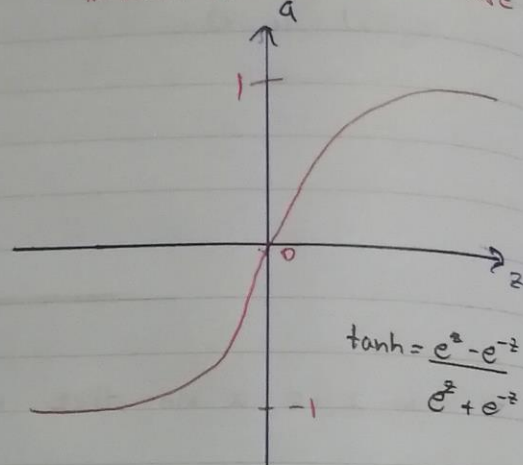
- \* Leaky ReLU is better ~~than~~ than ReLU. But we didn't use Leaky as much in practice.
- \* But advantage of ReLU and Leaky ReLU is that for a lot of the space of  $z$ , the derivative of the activation function, the slope of the activation function is very different from 0.
- \* So, in practice using the ReLU activation function, ~~with~~ your neural network will often learn faster than when using tanh or the sigmoid activation function.

## Summary

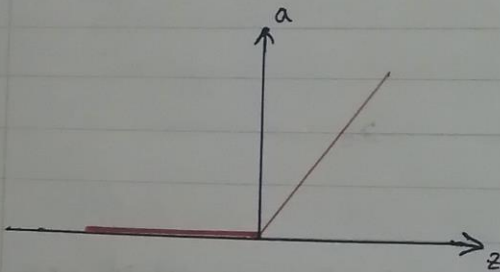


$$\text{Sigmoid} \rightarrow a = \frac{1}{1 + e^{-z}}$$

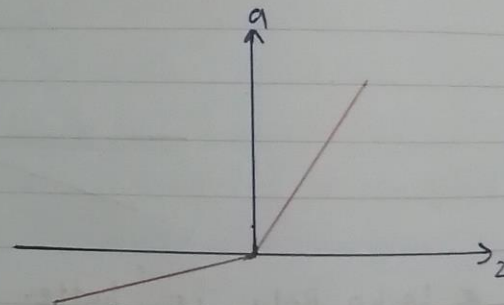
\* tanh is the best to use



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

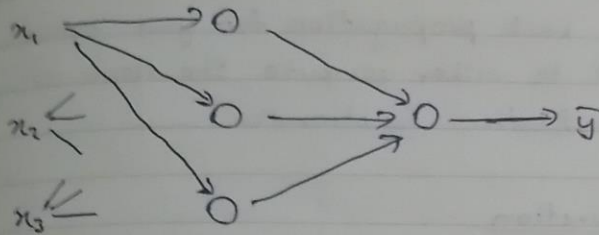


$$\text{ReLU} \rightarrow a = \max(0, z)$$



$$\text{Leaky ReLU} \rightarrow a = \max(0.01z, z)$$

Why do we need Non-Linear Activation Functions?



Given  $x$ :

Assume  $g(z) = z$

$$\begin{aligned} \text{for hidden layer} \quad \begin{cases} z^{[1]} = W^{[1]}x + b^{[1]} \\ a^{[1]} = g^{[1]}(z^{[1]}) \end{cases} & \longrightarrow \therefore a^{[1]} = z^{[1]} \\ \text{for output layer} \quad \begin{cases} z^{[2]} = W^{[2]}a^{[1]} + b^{[2]} \\ a^{[2]} = g^{[2]}(z^{[2]}) \end{cases} & \longrightarrow \therefore a^{[2]} = z^{[2]} \end{aligned}$$

$$\therefore a^{[1]} = z^{[1]} = W^{[1]}x + b^{[1]}$$

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

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

$$= (W^{[2]}W^{[1]})x + (W^{[2]}b^{[1]} + b^{[2]})$$

$$a^{[2]} = W'x + b' \quad \leftarrow \text{This means neural network outputting linear function of input.}$$

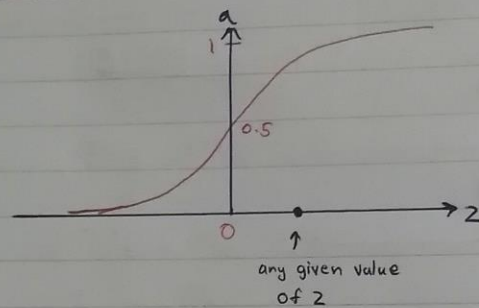
\* In above ~~prove~~ calculation we can see if we don't have ~~activation~~ activation function then no matter how many layers our neural network because all layers doing ~~linear~~ is just computing linear activation function. So you might not have any hidden layers.



## Derivatives of Activation Function

When you implement back propagation for your neural network, you need to either compute the slope or derivative of the activation functions.

### Sigmoid activation function



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

On above, given value of  $z$  will have some slope or derivative  
We can define that slope as

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

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

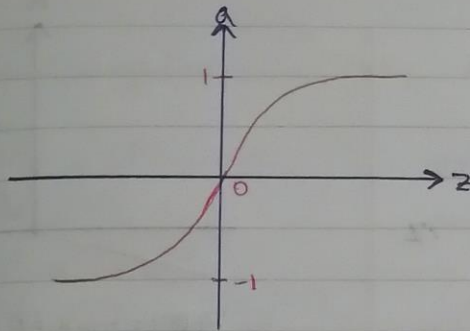
$$= g(z) (1 - g(z))$$

if  $z$  is very large  $\overset{z=10}{\text{then } g(z) \approx 1}$  and  $\frac{d}{dz} g(z) \approx 0$

if  $z$  is negative  $\overset{z=-10}{\text{then } g(z) \approx 0}$  and  $\frac{d}{dz} g(z) \approx 0$

if  $z$  is equal to 0 then  $g(z) = \frac{1}{2}$  and  $\frac{d}{dz} g(z) = \frac{1}{2} \left( 1 - \frac{1}{2} \right) = \frac{1}{4}$   
(0.5)

## Tanh activation function



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

\* we can define slope as

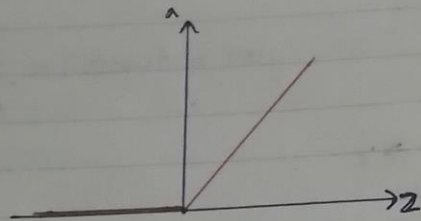
$$\begin{aligned} \frac{d}{dz} g(z) &= \text{slope of } g(z) \text{ at } z \\ &= 1 - (\tanh(z))^2 \end{aligned}$$

$$\begin{aligned} \text{if } z=10, \text{ then } \tanh(z) &\approx 1 \\ \frac{d}{dz} g(z) &\approx 0 \end{aligned}$$

$$\begin{aligned} \text{if } z=-10 \text{ then } \tanh(z) &\approx -1 \\ \frac{d}{dz} g(z) &\approx 0 \end{aligned}$$

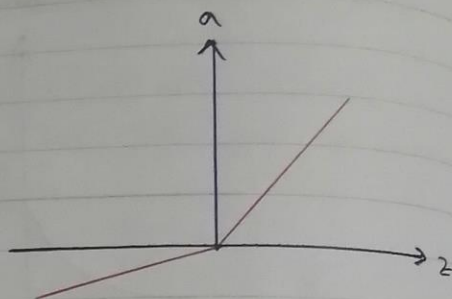
$$\begin{aligned} \text{if } z=0 \text{ then } \tanh(z) &= 0 \\ \frac{d}{dz} g(z) &= 1 \end{aligned}$$

## ReLU and Leaky ReLU



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

$$\frac{d}{dz} g(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$



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

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

## Gradient Descent for Neural Networks

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

$(n^{[1]}, n^{[2]})$     $(n^{[1]}, 1)$     $(n^{[2]}, n^{[1]})$     $(n^{[2]}, 1)$

$h_x = n^{[0]}, n^{[1]}, n^{[2]} = 1$

Cost function  $\rightarrow J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}) = \frac{1}{m} \sum_{i=1}^m L(\bar{y}, y)$

$\uparrow$   
 $a^{[2]}$

\* To train parameters, of our Algorithm, we have to perform gradient descent.

## Formulas for computing derivatives

### Forward Propagation

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

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

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

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

### Backward Propagation

$$dz^{[2]} = A^{[2]} - Y \quad Y = [y^{[1]}, y^{[2]}, \dots, y^{[m]}]$$

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

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

$$dz^{[1]} = \underbrace{W^{[2]T}}_{(n^{[1]}, m)} \underbrace{dz^{[2]}}_{(m)} \underset{\substack{\uparrow \\ \text{elementwise} \\ \text{Product}}}{*} \underbrace{g^{[1]'}(z^{[1]})}_{(n^{[1]}, m)}$$

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

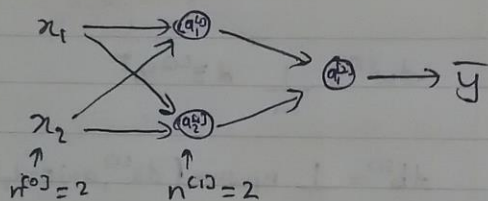
$$db^{[1]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=\overset{\text{True}}{n})$$



## Random Initialization

- \* When you change your neural network, it is important to initialize the weights randomly.
- \* For Logistic regression, it is okay to initialize the weights to 0. But for a neural network, initialize weights and parameters to 0 and then applied gradient descent doesn't work much.

What happens if you initialize weights to zero?



$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

- \* Initializing bias ( $b$ ) to 0 is actually okay but initializing  $w$  to 0 is a problem.
- \* If you initialized  $w$  to 0 then  $a_1^{[1]} = a_2^{[1]}$ . It means both of hidden nodes in ~~hidden~~ hidden layers are computing exactly same function.
- \* And then, when you compute backpropagation,  $\delta z_1^{[1]} = \delta z_2^{[1]}$ . It is also computing same function.