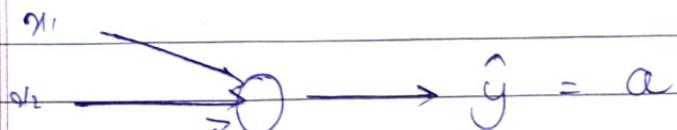
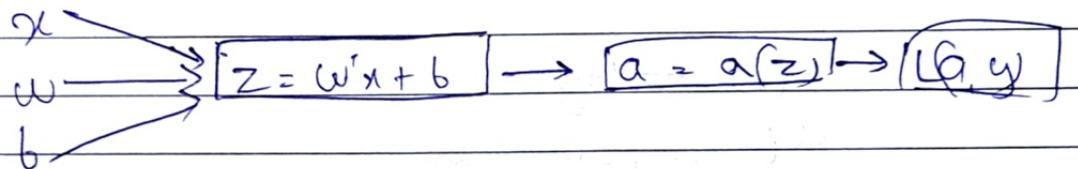


## Week 3

Vid 1: Overview + Vid 2

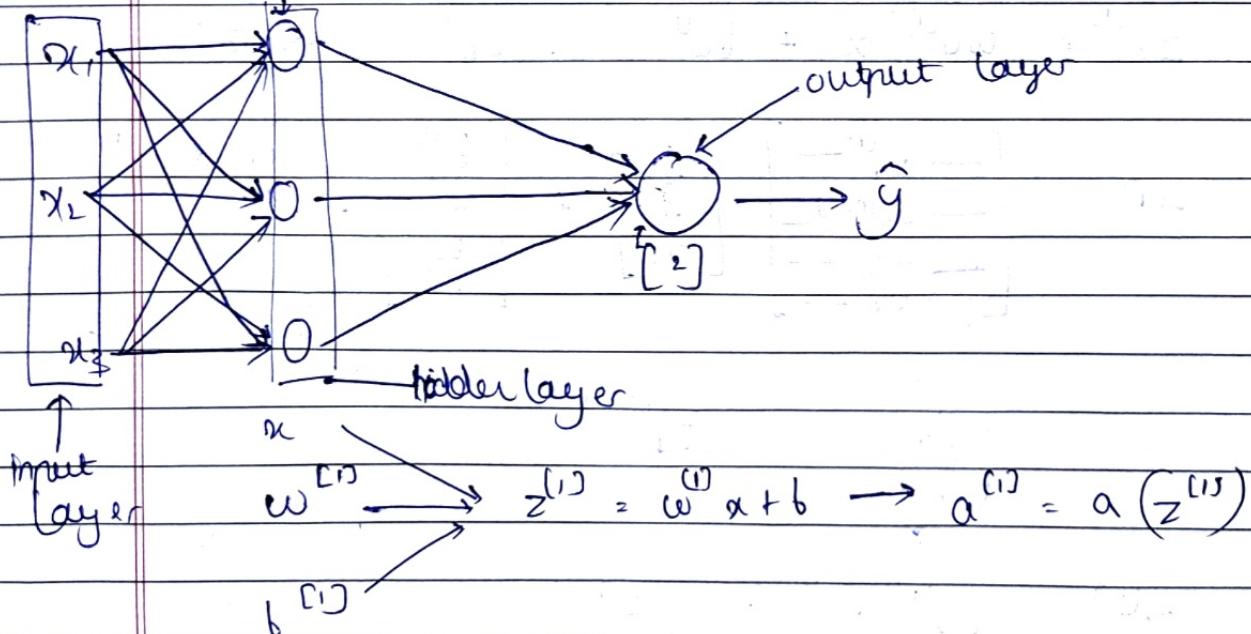


~~sigmoid function to implement sigmoid~~



neural networks: stacking of sigmoid units

$[1]$  — quantities associated with 1<sup>st</sup> layer



hidden layer - you can't see in training! set

activation - A  $\rightarrow$  values passed onto the subsequent layers

## Computing Output of Neural Network

$(j) \leftarrow$  layer

$a_i \leftarrow$  node in layer

$$z_1^{(1)} = w_1^{(1)\top} x + b_1^{(1)}$$

$$z_2^{(1)} = w_2^{(1)\top} x + b_2^{(1)}$$

$$z_3^{(1)} = w_3^{(1)\top} x + b_3^{(1)}$$

$$z_4^{(1)} = w_4^{(1)\top} x + b_4^{(1)}$$

$$a_1^{(1)} = a(z_1^{(1)})$$

$$a_2^{(1)} = a(z_2^{(1)})$$

$$a_3^{(1)} = a(z_3^{(1)})$$

$$a_4^{(1)} = a(z_4^{(1)})$$

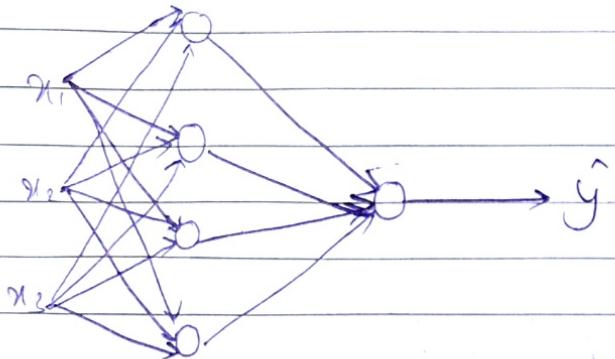
$$z^{(1)} = \begin{bmatrix} -w_1^{(1)\top} \\ -w_2^{(1)\top} \\ -w_3^{(1)\top} \\ -w_4^{(1)\top} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix} = \begin{bmatrix} w_1^{(1)\top} x + b_1^{(1)} \\ \vdots \\ \vdots \\ w_4^{(1)\top} x + b_4^{(1)} \end{bmatrix} = z^{(1)}$$

$$a^{(1)} = \begin{bmatrix} a_1^{(1)} \\ \vdots \\ a_4^{(1)} \end{bmatrix} = a(z^{(1)})$$

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

$$w^7 = w^{(2)}$$

$$a^{(2)} = a(z^{(2)})$$



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

$$\begin{aligned}
 z^{(1)} &= w^{(1)}x + b^{(1)} \\
 a^{(1)} &= \alpha(z^{(1)}) \\
 z^{(2)} &= w^{(2)}a^{(1)} + b^{(2)} \\
 a^{(2)} &= \alpha(z^{(2)})
 \end{aligned}$$

$$\begin{aligned}
 z^{(2)} &= w^{(2)}a^{(1)} + b^{(2)} \\
 a^{(2)} &= \alpha(z^{(2)})
 \end{aligned}
 \quad \text{gradient initialized}$$

for  $i = 1$  to  $m$

$$\begin{aligned}
 z^{(1)(i)} &= w^{(1)}x^{(i)} + b^{(1)} \\
 a^{(1)(i)} &= \alpha(z^{(1)(i)}) \\
 z^{(2)(i)} &= w^{(2)}a^{(1)(i)} + b^{(2)} \\
 a^{(2)(i)} &= \alpha(z^{(2)(i)})
 \end{aligned}$$

$$X = \begin{bmatrix} & & & \\ \uparrow & \uparrow & \uparrow & \uparrow \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \text{features} & & & \\ (n_x, m) & & & \end{bmatrix}$$

$$\begin{aligned} z^{(1)} &= W^{(1)} X + b^{(1)} \\ A^{(1)} &= \sigma(z^{(1)}) \\ z^{(2)} &= W^{(2)} A^{(1)} + b^{(2)} \\ A^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

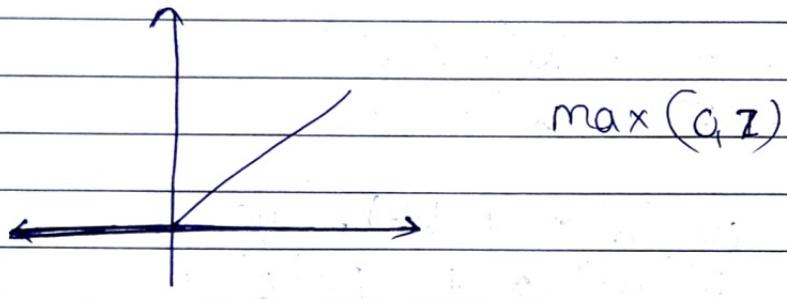
$$A^{(2)} = \left[ \begin{array}{c c c c} & & & \\ \uparrow & \uparrow & \uparrow & \uparrow \\ a^{(2)(1)} & a^{(2)(2)} & \dots & \\ \downarrow & \downarrow & \downarrow & \downarrow \\ \text{unit} & & & \end{array} \right]$$

$$W^{(1)} \left[ \begin{array}{c c c c} 1 & 1 & 1 & 1 \\ x^1 & x^2 & x^3 & \dots & x^m \\ \downarrow & \downarrow & \downarrow & & \downarrow \\ 1 & 1 & 1 & & 1 \end{array} \right]$$

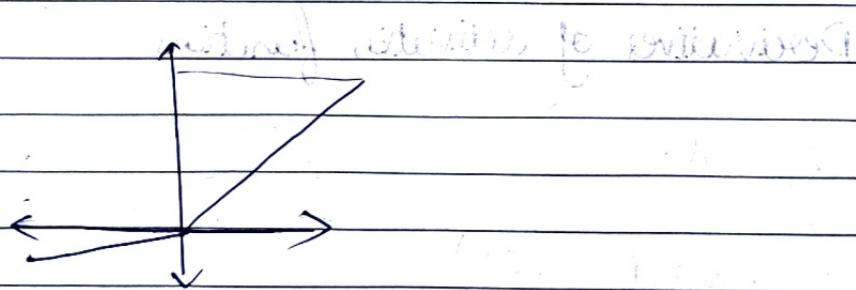
## Activation Functions

- \* So far we have worked with sigmoid function but it is not preferred.
- \* Its range is  $[0, 1]$
- \* Another activation function is tanh function whose range is  $[-1, 1]$
- \* Tanh function is better than sigmoid because the avg of its output is closer to 0, so it helps centering data better for next layer

- \* It cannot be used for binary classification where output is between 0 & 1 and tanh gives output between (-1, 1)
- \* The disadvantage of sigmoid or tanh function is that if input is too large / too small, the slope will near 0 which will cause gradient descent problem.
- \* ReLU - most preferred



- \* Leaky ReLU



- ② Instead of returning 0 for -ve values it returns a small negative number

In hidden layers using linear activation function  
is rare.

Linear activation function produces linear activation.  
It is used only in output layer.

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

$$a^{(1)} = z^{(1)}$$

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

$$a^{(2)} = z^{(2)}$$

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

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

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

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

## Derivatives of activation functions

Sigmoid:

$$\begin{aligned} g(z) &= g(z)(1 - g(z)) \\ &= a(1 - a) \end{aligned}$$

tanh:

$$\begin{aligned} g'(z) &= 1 - g(z)^2 \\ &= 1 - a^2 \end{aligned}$$

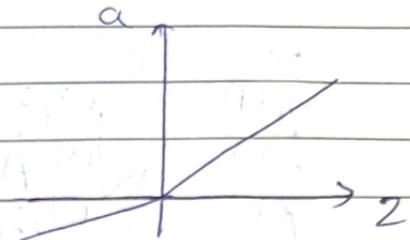
ReLU

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

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

at 0

Leaky ReLU



Leaky ReLU

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

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

Gradient Descent for Neural Networks.

Parameters:  $\omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)}$   
 $(n^{(0)}, n^{(1)})'$ ,  $(n^{(1)}, 1)', (n^{(2)}, n^{(1)})', (n^{(1)}, 1)'$

$n^{(0)} = n^{(0)}$  - input feature

$n^{(1)}$  - hidden units

$n^{(2)}$  - output unit

Cost function:  $J(\omega^{(1)}, b^{(1)}, \omega^{(2)}, b^{(2)})$

$$= \frac{1}{m} \sum_{i=1}^m L(\hat{y}_i, y_i)$$

$$d\omega^{(1)} = \frac{dJ}{d\omega^{(1)}}$$

$$db^{(1)} = \frac{dJ}{db^{(1)}}$$

$$\Delta \omega^{(1)} = \omega^{(1)} - \alpha d\omega^{(1)}$$

$$\Delta b^{(1)} = b^{(1)} - \alpha db^{(1)}$$

Page No.	
Date	

## Forward propagation | Back-Propagation Step

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

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

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

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

- because  
binary  
classifier

$$dZ^{(2)} = A^{(2)} - Y \text{ S.P}$$

$$dw^{(2)} = \frac{1}{m} dZ^{(2)} A^{(1)\top}$$

$$db^{(2)} = \frac{1}{m} \text{ np. sum}(dZ^{(2)}, ax=1, \text{ keepdim} = \text{true})$$

$$dZ^{(1)} = \underbrace{(W^{(2)\top} dZ^{(2)} * g^{(1)'} Z^{(1)}}_{(n^{(1)}, m)} \uparrow \underbrace{(n^{(1)}, m)}_{\text{element wise product}}$$

similarly forward ref forward training

## Random Initialisation

- \* In logistic regression, it wasn't important to initialise weights randomly.
- \* However, in NN - we initialise randomly.
- \* Initialising all weights to 0  $\rightarrow$  X - Symmetric hidden units  
Initialising bias ( $b$ ) to 0  $\rightarrow$  ✓

This is because:

- all hidden units will be identical & compute same function
- same update on iteration with gradient descent.

\*  $\therefore W$  is initialised randomly,

example:

$$w^{(1)} = \text{np.random.randn}(2, 2) * 0.1$$

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

:

$$\begin{aligned} w^{(2)} &= \text{np.random.randn}(1, 2) * 0.1 \\ b^{(2)} &= 0 \end{aligned}$$