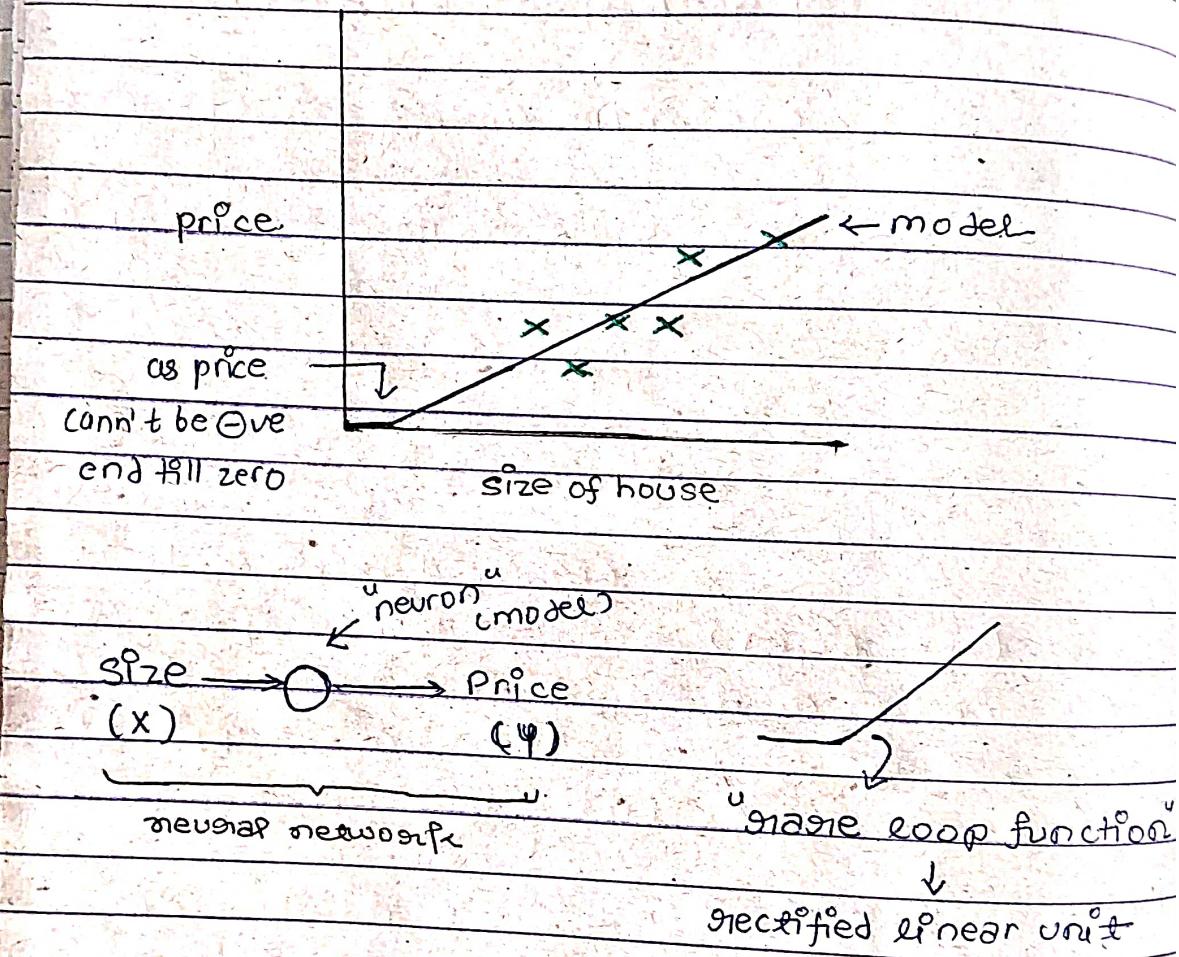


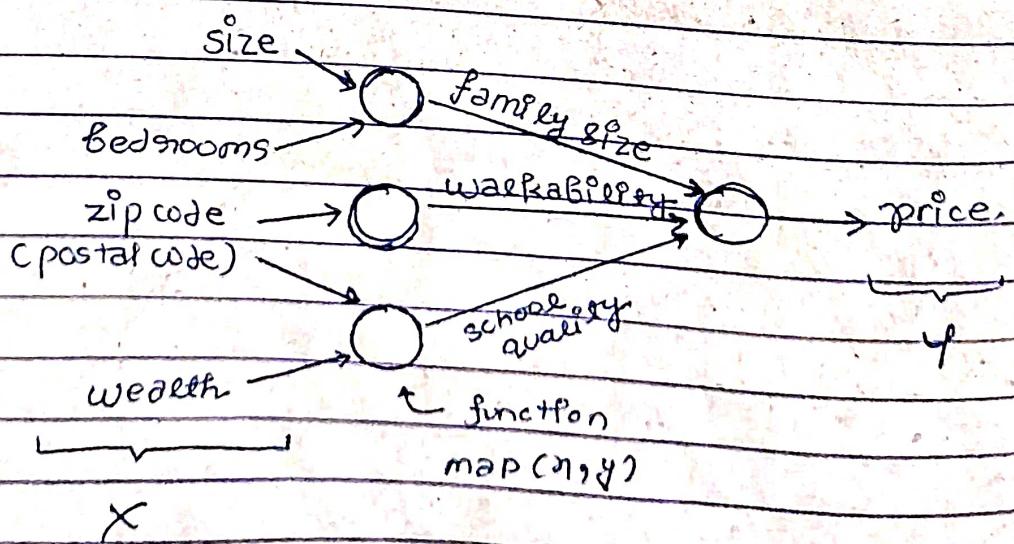
Artificial Neural Network

what is a neural network?

i.e. house-price prediction

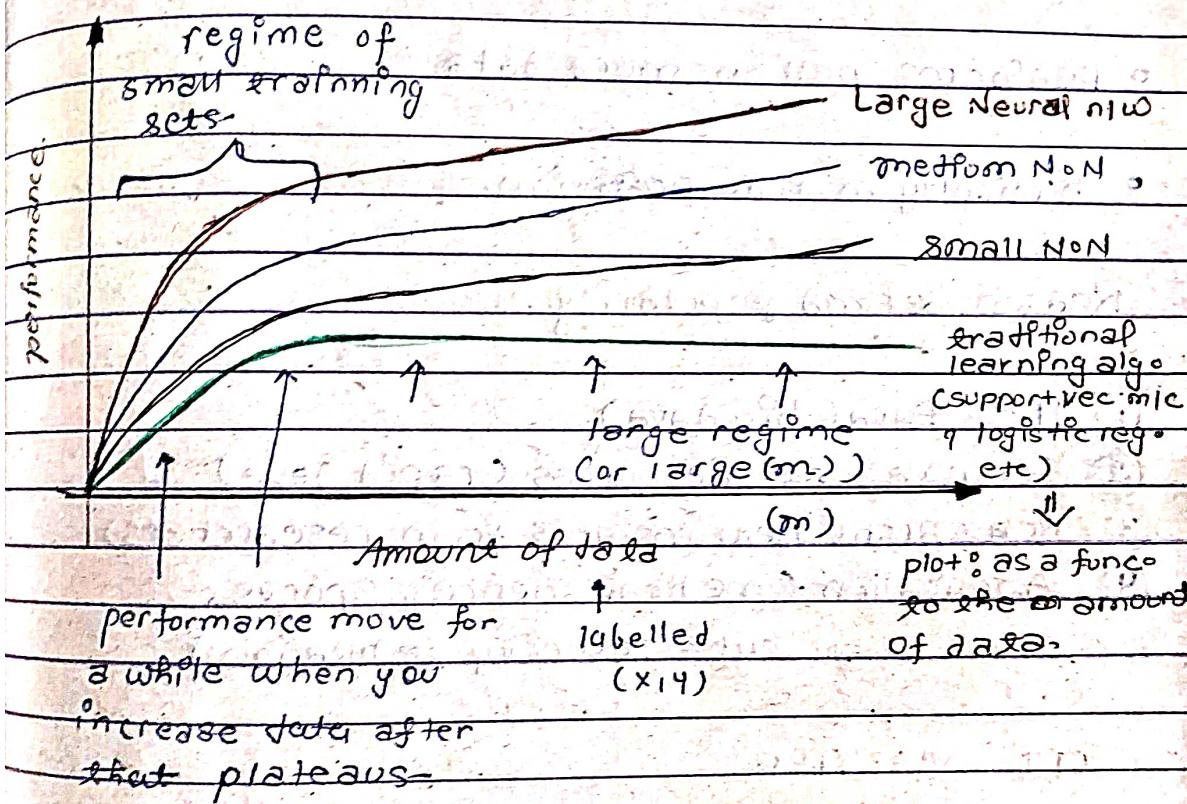


i.e. multiple-input features.

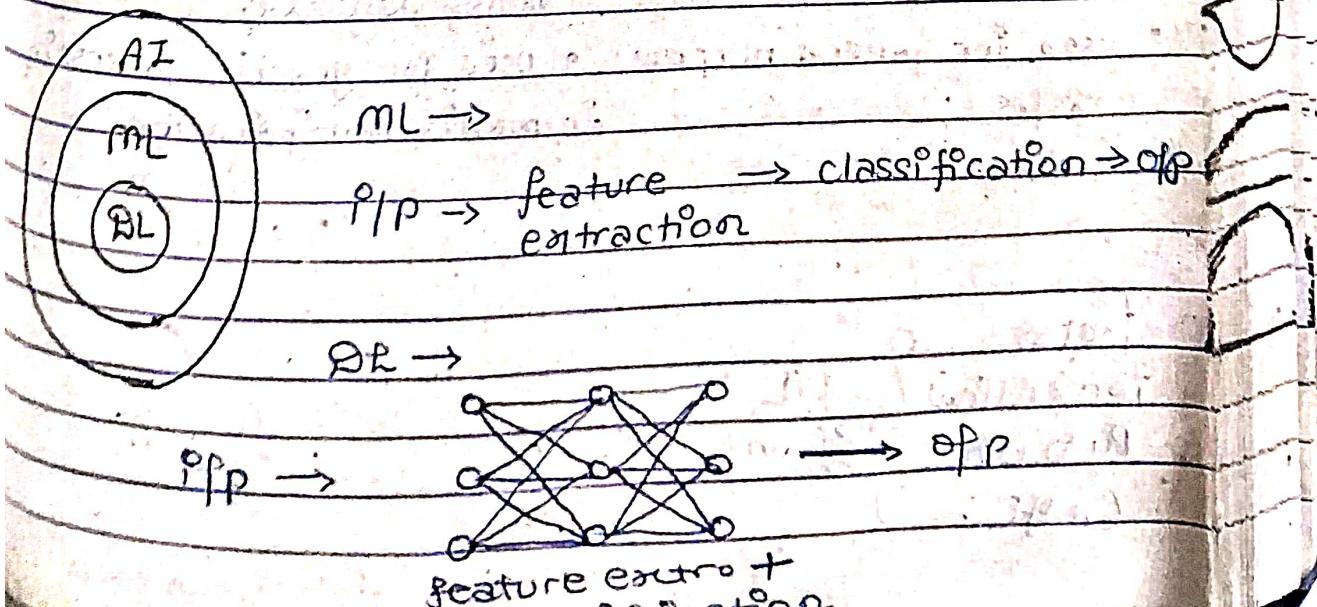


Neural networks are computational models that mimic the complex functions of the human brain.

The neural network consists of interconnected nodes or neurons that processes and learns from data, enabling tasks such as pattern recognition and decision making in machine learning.



• m: size of training set in training example



- It's a particular kind of ML that achieves great power & flexibility by ML to represent
- requires more computation power.
- Automatic feature extraction
- performs well on more data
- difficulty in interpretation.

Reason behind popularity

- ① Huge Data (big data)
- ② Computational Resources (rapid dev.)
- ③ Research (new insights from researchers)
- ④ Application (medical science, space, cybersecurity, education)

CPU vs GPU

CPU

- few complex cores
- focuses on doing one task as possible
- used for general purpose tasks

GPU

- hundreds of simpler cores
- possible of doing various tasks parallelly,
- used for graphic processing or matrix multiplication.

Frameworks

Tensorflow

Pytorch

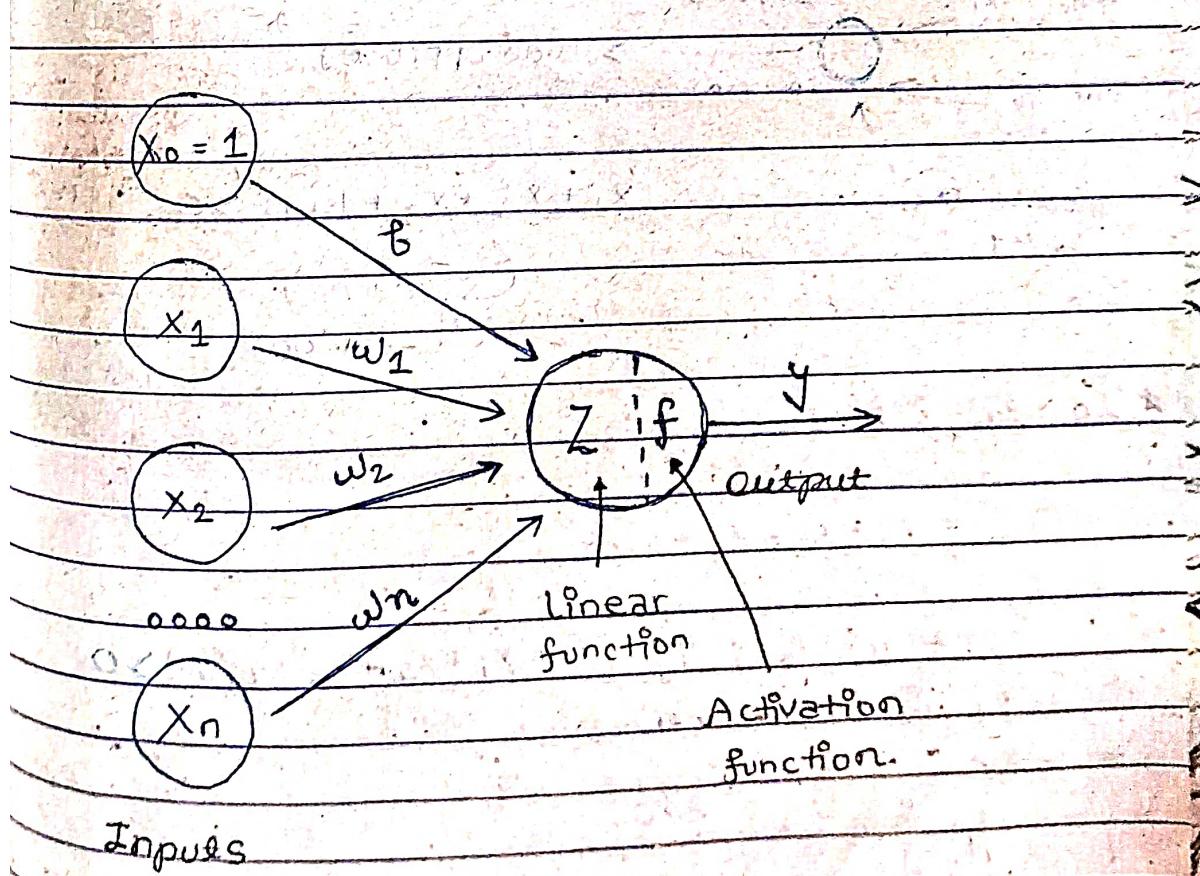
Caffe

DL

frameworks

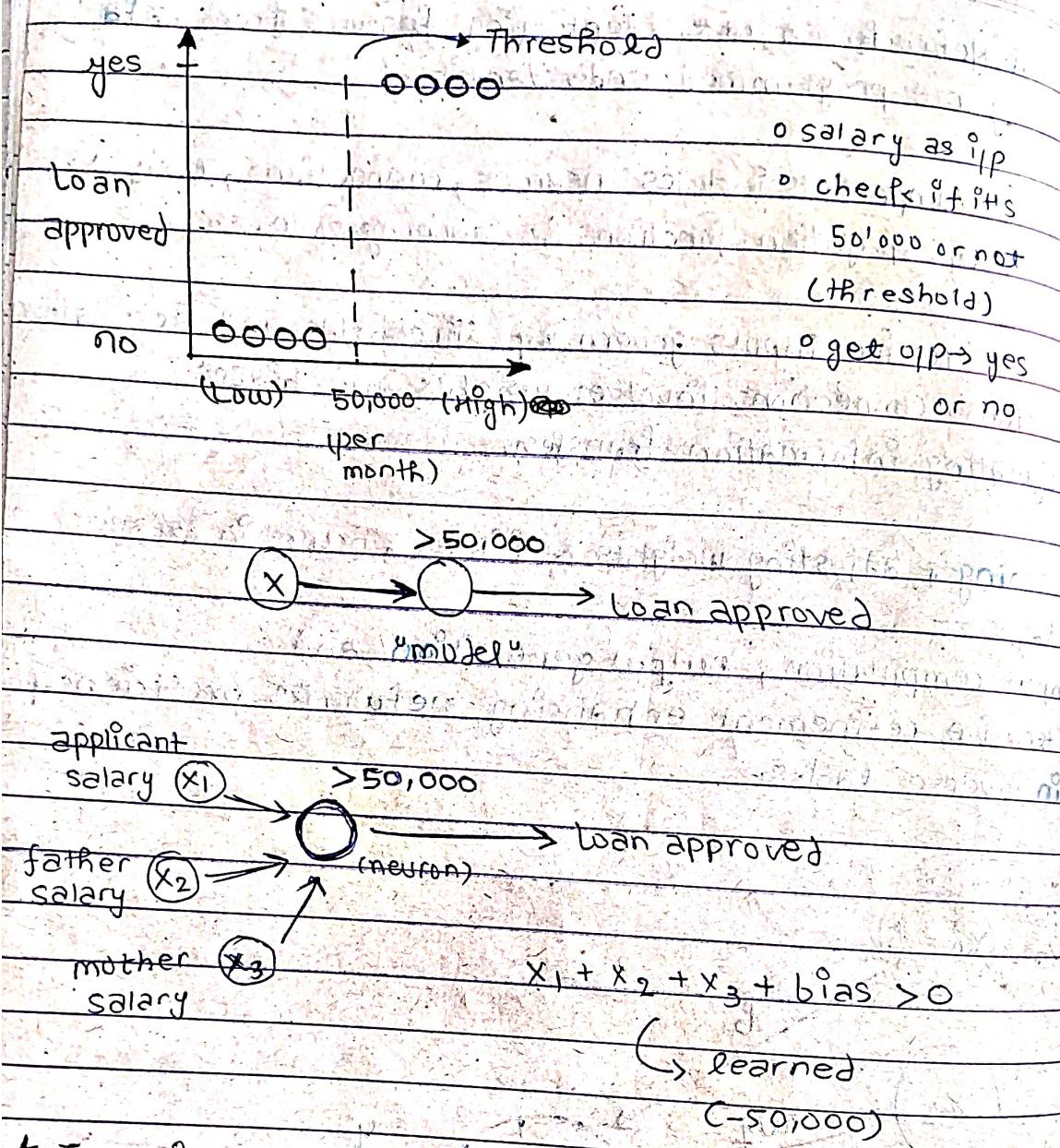
What are Neural Networks?

- Neural Networks extract identifying features from data, lacking pre-programmed understanding.
- Network components includes neurons, connections, weights, biases, propagation functions & learning rule.
- Neurons receive inputs govern by thresholds and activation functions. Connections involves weights and biases regulating information transfer.
- Learning, adjusting weights & biases, occurs in three stages — input computation, output generation and iterative refinement enhancing networks proficiency in diverse tasks.



★ Intuition behind perceptron

Draw
Page



★ Intuition Behind activation func.

$$z = x_1 + x_2 + x_3 + \text{bias}$$

O/P \rightarrow will be '1' if $(z = x_1 + x_2 + x_3 + \text{bias}) > 0$
it will be '0' otherwise

$$\text{O/P } \begin{cases} 1, & z > 0 \\ 0, & z \leq 0 \end{cases}$$

O/P \rightarrow step func(z) or O/P \rightarrow Step(z)

Derivatives

Date _____
Page _____

$$f(a) = 3a$$

$$a = 2 \quad f(a) = 6$$

$$a = 2.001 \quad f(a) = 6.003$$

$\frac{0.003}{0.001} = 3$ (height) / (width)
slope (derivative) of $f(a)$ at $a=2$ is 3

$$a = 5 \quad f(a) = 15$$

$$a = 5.001 \quad f(a) = 15.003$$

slope at ($a=5$) is $f'(a)$ also 3

$\frac{15.003 - 15}{5.001 - 5} = 3 = \frac{d}{da} f(a)$ if I nudge 'a' to the right a little bit

I expect $f(a)$ to go up by three times as much as nudge the value of 'a',

if $f'(a)$ of

$$f(a) = a^2$$

$$a = 2 \quad f(a) = 4$$

$$a = 2.001 \quad f(a) \approx [4.004] \\ (\approx 4.004001)$$

(+ slope = derivative) of $f(a)$ at $a=2$ is 4.

$$\frac{d}{da} f(a) = 4 \text{ when } a = 2$$

$$a = 5 \quad f(a) = 25$$

$$a = 5.001 \quad f(a) = 25.001$$

$$\frac{d}{da} f(a) = \frac{d}{da} a^2 = 2a$$

$$\frac{d}{da} f(a) = 10 \text{ when } a = 5$$

* for any given value of 'a'
you nudge up by (0.0001)

infinitesimal number you would expect $f(a)$ to go up by $"2a"$.

$$\textcircled{1} \quad f(a) = a^2 \quad \frac{d}{da} f(a) = 2a \quad a=2 \quad f(a)=4$$

(619)

$$a=1.99 \quad \frac{d}{da} f(a) = 2a \quad a=2.001 \quad f(a)=4.004$$

$$\textcircled{2} \quad f(a) = a^3 \quad \frac{d}{da} f(a) = 3a^2 \quad a=2 \quad f(a)=8$$

(619)

$$a=1.99 \quad \frac{d}{da} f(a) = 3a^2 \quad a=2.001 \quad f(a) \approx 8.012$$

(619)

$$3 \times 2^2 \rightarrow 12$$

when you nudge 'a' to the right little by tiny bit $f(a)$ should go up by 1.2

$$\textcircled{3} \quad f(a) = \log_e(a) \Rightarrow \frac{d}{da} f(a) = \frac{1}{a} \quad a=2, f(a)=0.6931$$

(619)

$$a=2.001, f(a)=0.693$$

(619)

or
en(a)

f

↓
 $\frac{d}{da} f(a) = 1/2$ go up by $1/2$.

0.0005

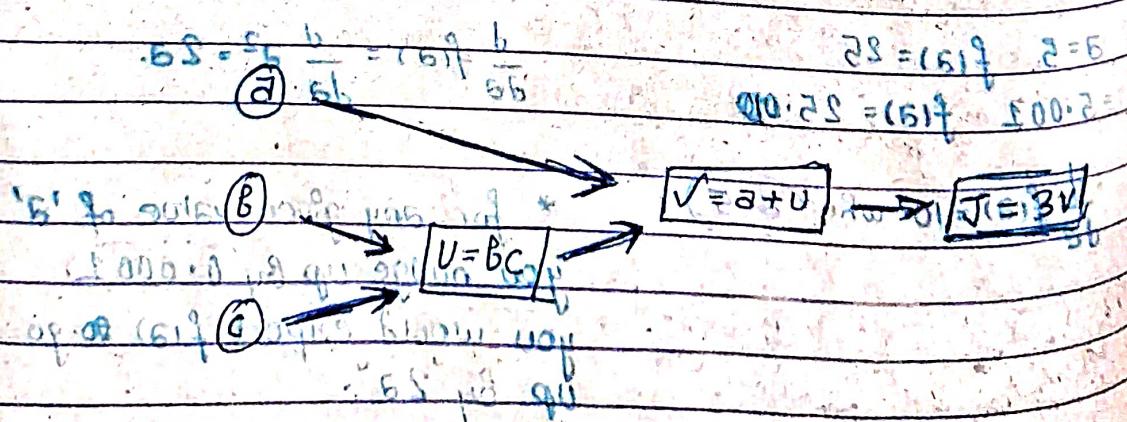
* Computation graph :
minimum value

let $J(a, b, c) = 3(a + b * c)$

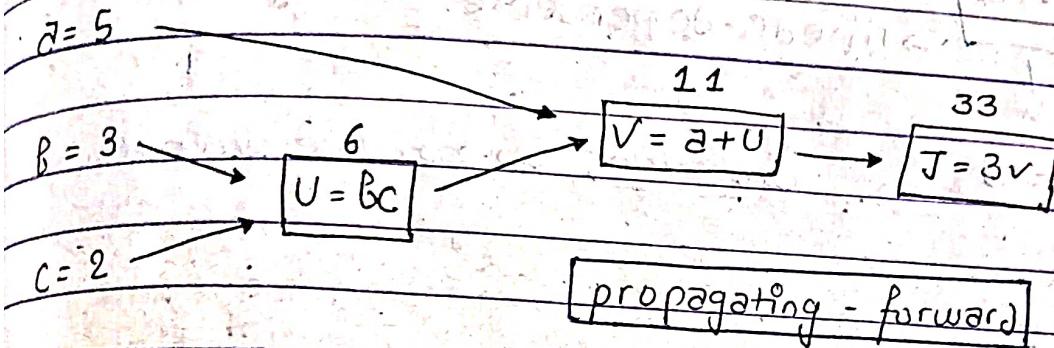
$$U = b \cdot c = 5 \quad \begin{matrix} U \\ \downarrow \end{matrix}$$

$$V = a + U = a + 5 \quad \begin{matrix} V \\ \downarrow \end{matrix}$$

$$J = 3V \quad \begin{matrix} J \\ \downarrow \end{matrix}$$



derivatives with computation graph



$\rightarrow \frac{dJ}{dv} = 0$ if we were to take this value (v) and change it little bit. How would the value of "J" change.

$$J = 3v$$

$$v = 11 \quad \text{pump up to } 11.001 \quad J \text{ bump up to } 33.003$$

Hence, $\frac{dJ}{dv} = 3$

$$\rightarrow \frac{dJ}{da} = 0 \quad \text{How does 'a' affect 'J'}$$

$$a = 5 \quad \text{bump to } 5.001$$

$$v = 11 \quad \text{bump to } 11.001$$

$$J = 33 \quad \text{bump to } 33.003$$

change 'a' \downarrow $v \downarrow$

changes $v \downarrow$

change 'J'

$$\text{as } \frac{dJ}{dv} = 3$$

$$\frac{\partial J}{\partial a} = \frac{\partial J}{\partial v} \cdot \frac{\partial v}{\partial a}$$

{chain-rule}

but 'a' is same as 'v'

$$a \leftarrow \frac{v}{b}, \frac{t_b}{v} \leftarrow \frac{t_b}{a}$$

$$\partial = 5 \text{ (to } 5.001)$$

$$\frac{dV}{d\alpha} = 1 \quad \text{as } (\sqrt{V} = 11.001) \rightarrow \text{linear-dependence.}$$

$$\frac{\partial J}{\partial V} = 3$$

outcome: to produce output

$$\begin{aligned} \partial = 5 & \\ \frac{\partial J}{\partial \alpha} \rightarrow \partial \alpha = 3 & \quad \text{---} \quad \sqrt{V} = \partial + U \quad \text{---} \quad J = 3V \\ B = 3 & \quad \text{---} \quad U = BC \quad \frac{\partial J}{\partial V} \rightarrow \partial V = 3 \\ C = 2 & \end{aligned}$$

let, we call
(for writing code)

[propagating-backward]

$$\begin{aligned} \rightarrow \frac{\partial J}{\partial U} &= 2 \quad U = 6 \xrightarrow{\text{bump to}} 6.001 \\ &\quad \sqrt{V} = 11.001 \xrightarrow{\text{---}} 11.001 \\ &\quad J = 33 \xrightarrow{\text{---}} 38.003 \end{aligned}$$

$$\begin{aligned} \frac{\partial J}{\partial V} \cdot \frac{\partial V}{\partial U} & \quad \text{(chain rule)} \\ + & \quad \text{(sum rule)} \\ 3 & \quad 1 \end{aligned}$$

$$\rightarrow \frac{\partial J}{\partial B} = 0 \quad \frac{\partial J}{\partial U} \cdot \frac{\partial U}{\partial B} \quad \text{(chain rule)}$$

$$B \rightarrow 3 \rightarrow 3.001$$

$$U = b \cdot c \rightarrow 6 \rightarrow 6.002 \quad (\text{as } c=2)$$

'U' has gone up twice as 'B' has gone up

$$\text{Pf} \quad \frac{\partial J}{\partial B} \Rightarrow \frac{\partial J}{\partial V} \cdot \frac{\partial V}{\partial B} \rightarrow 6$$

$$b=3 \rightarrow 3.001 \quad \text{as } c=2$$

$$1=b \cdot c = 6 \rightarrow 6.002$$

$$\begin{aligned} \hookrightarrow J &= 33.006 \rightarrow V = 11.002 \\ &\rightarrow J = 3V \end{aligned} \quad \left\{ \begin{array}{l} \partial S \frac{\partial J}{\partial b} = 6 \\ \end{array} \right.$$

outcome: which variable has impact to $\uparrow\downarrow$ efficiency

* Logistic Regression in Neural Networks

LOR Recap]

$$z = w^T \theta x + b \quad (h_{\theta}(x) = \theta^T x + \epsilon)$$

$$\hat{y} = \sigma(z) \quad (h_{\theta}(x) \rightarrow \text{predicted } \hat{y})$$

$$L(\hat{y}, y) = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

(1)

(2)

(3)

(4)

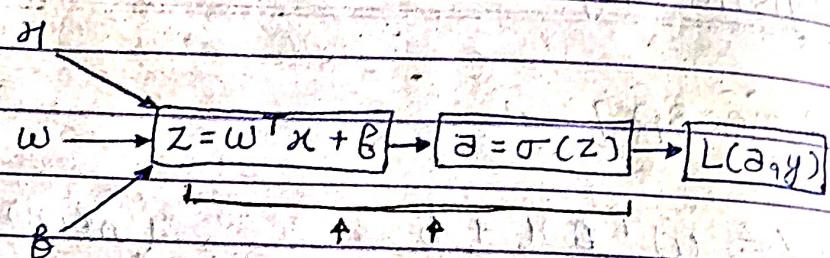
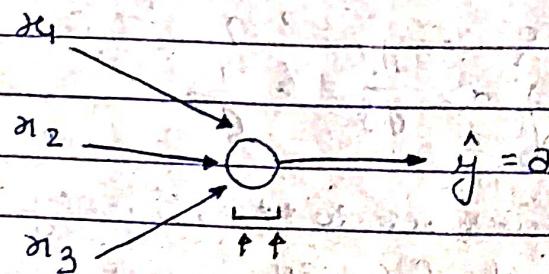
(5)

$$[1] z = w_1 x_1 + w_2 x_2 + b$$

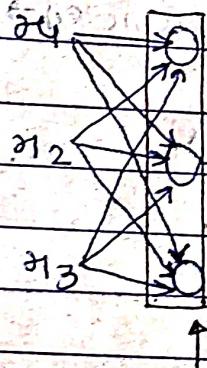
$$z = w_1 x_1 + w_2 x_2 + b \rightarrow \hat{y} = \sigma(z)$$

$$L(\hat{y}, y)$$

• Neural Network Representation?



$[1]$



superscript refer to belong to that layer

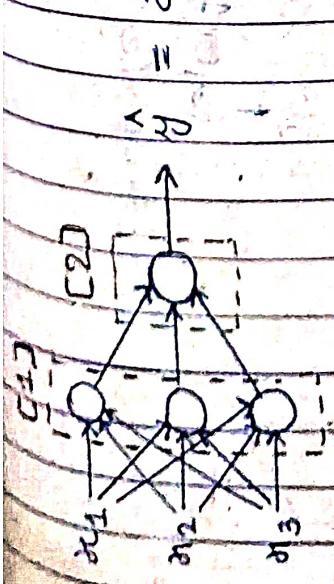
$[1]$

$$w^{[1]} \rightarrow z^{[1]} = w^{[1]} x_1 + b^{[1]}$$

$b^{[1]}$

$[2]$

$[3]$



$$z[c_1] = w[10] * x + b[10]$$

$$a[c_1] = \sigma(z[c_1])$$

$$z[c_2] = w[20] * x + b[20]$$

$$a[c_2] = \sigma(z[c_2])$$

$$w[c_2]$$

$$b[c_2]$$

$$x$$

$$z[c_3] = w[30] * x + b[30]$$

$$a[c_3] = \sigma(z[c_3])$$

$$z[c_4] = w[40] * x + b[40]$$

$$a[c_4] = \sigma(z[c_4])$$

$$w[c_4]$$

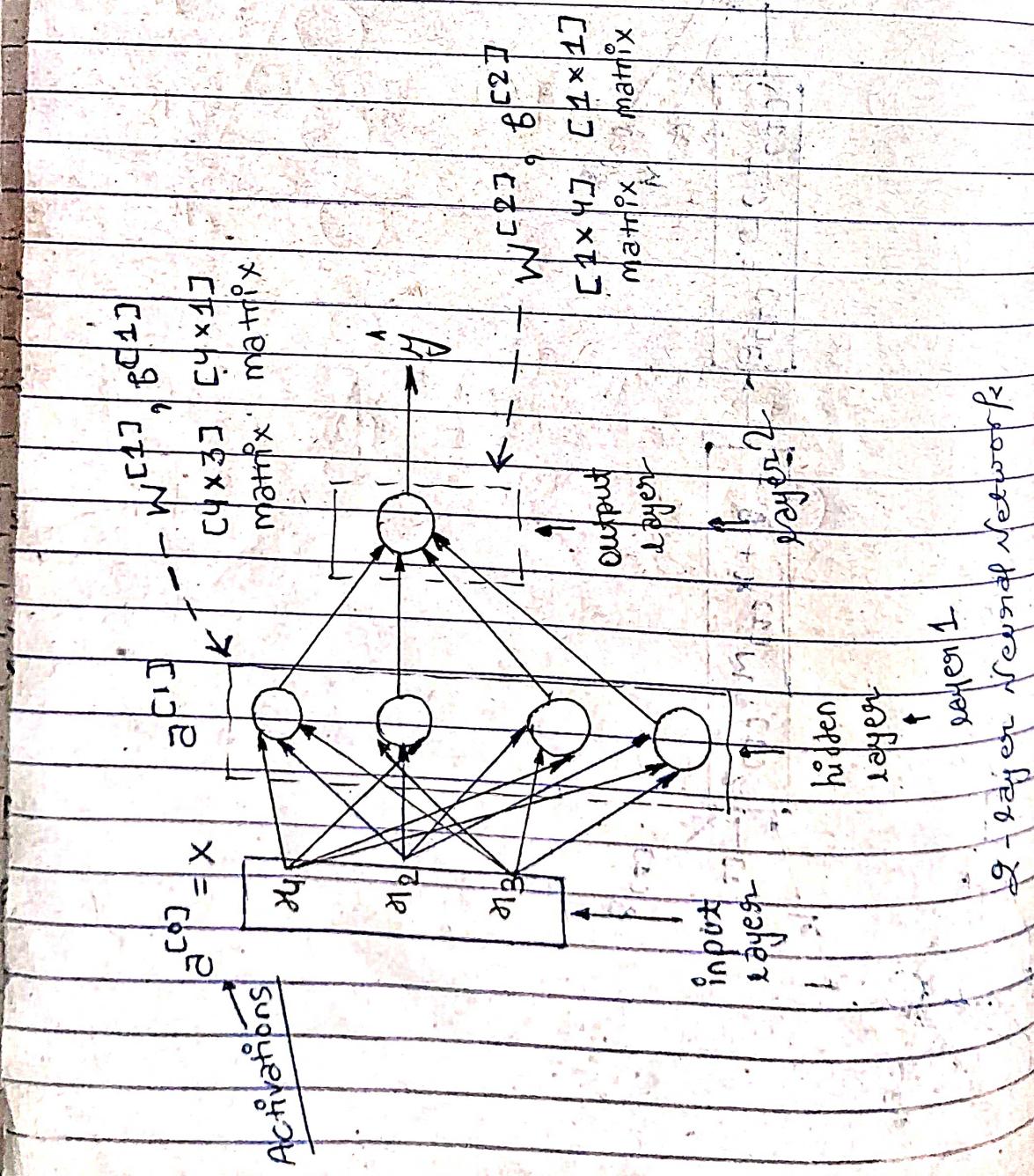
$$b[c_4]$$

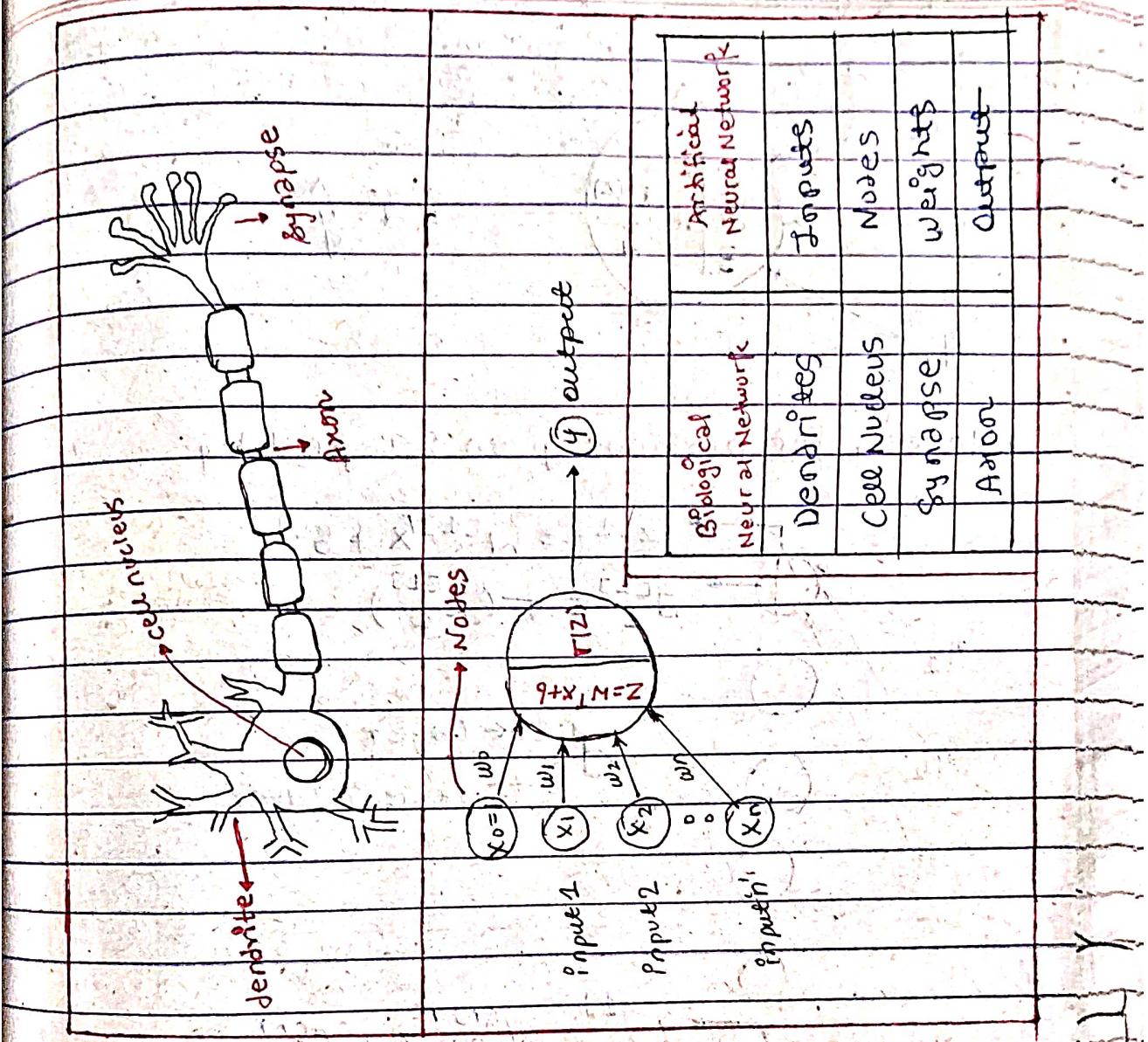
$$x$$

D do Page

o Hidden layer are responsible for extracting features and patterns from the input data which are then used to make predictions & classifications in the output layer.

- o \exists [Activations] →
 - the values that different players of the neural network are passing on to the next subsequent layer.
 - then come hidden layer generating activation function $\boxed{C1}$
 - then come hidden layer generating activation function $\boxed{C2}$
 - then output layer generate activation function

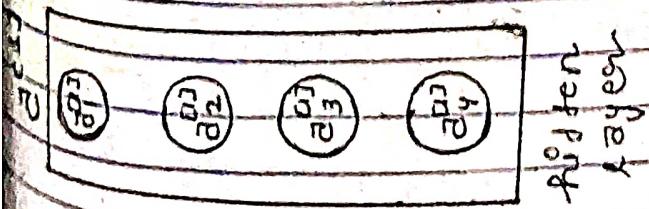




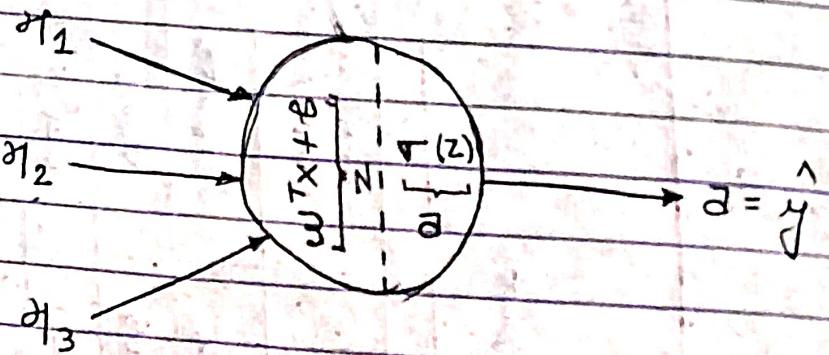
The nodes in this activation function $\sigma(z)$ generate values -

$$\sigma(z) = \begin{cases} 0 & z \leq 0 \\ 1 & z > 0 \end{cases}$$

4 dimensional vector



* Computing Neural Network's Output ;



Hence, for previous given 2-Layer NN.

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

$$\bar{o}_1^{[1]} = \sigma(z_1^{[1]})$$

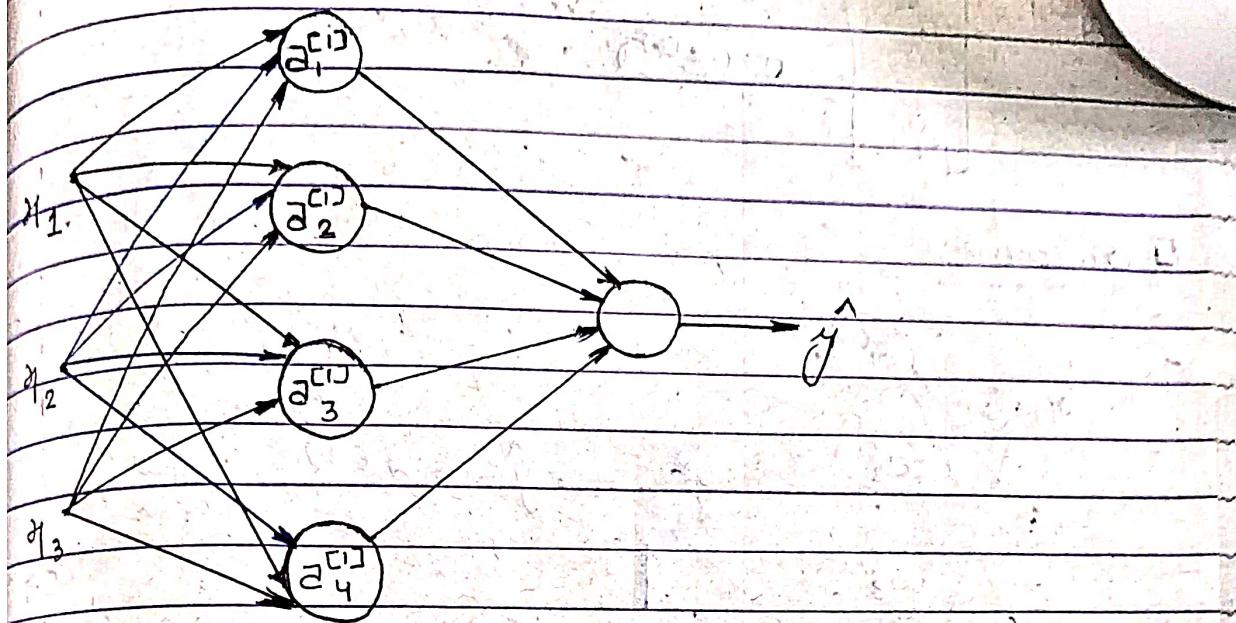
$\bar{o}_i^{[l]} \leftarrow$ Layer
 $\bar{o}_i^{[l]} \leftarrow$ node in layer

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

$$\bar{o}_2^{[1]} = \sigma(z_2^{[1]})$$

h_3

Similarly, perform with others...



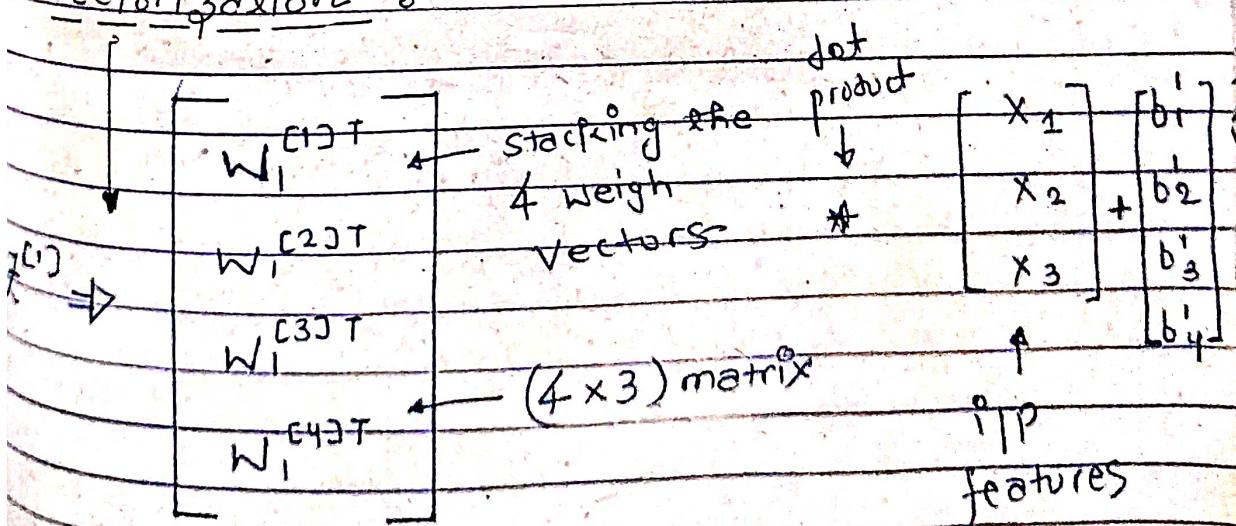
$$z_1^{(1)} = w_1^{(1)T} x + b_1^{(1)}, \quad a_1^{(1)} = \sigma(z_1^{(1)})$$

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

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

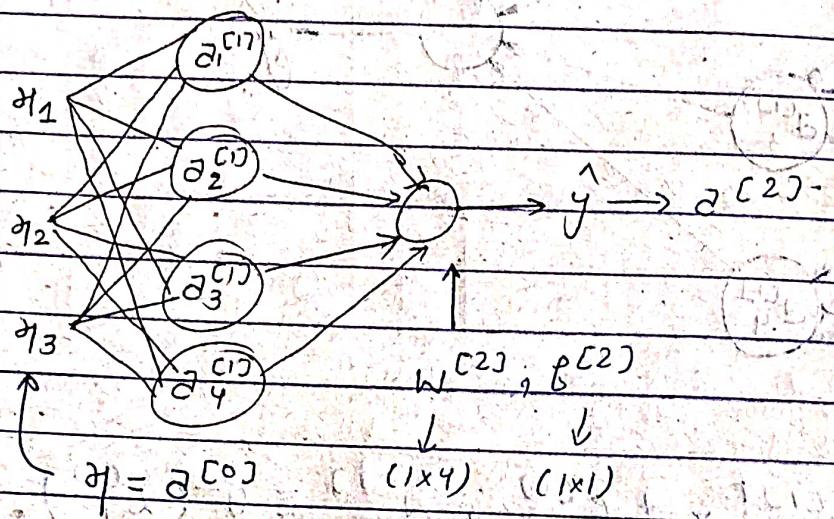
$$z_4^{(1)} = w_4^{(1)T} x + b_4^{(1)}, \quad a_4^{(1)} = \sigma(z_4^{(1)})$$

* Vectorization :



$$\tilde{a} = \begin{bmatrix} a^{[1]} \\ \vdots \\ a^{[4]} \end{bmatrix} = \sigma(z^{[1]})$$

5 outcome,



given input of

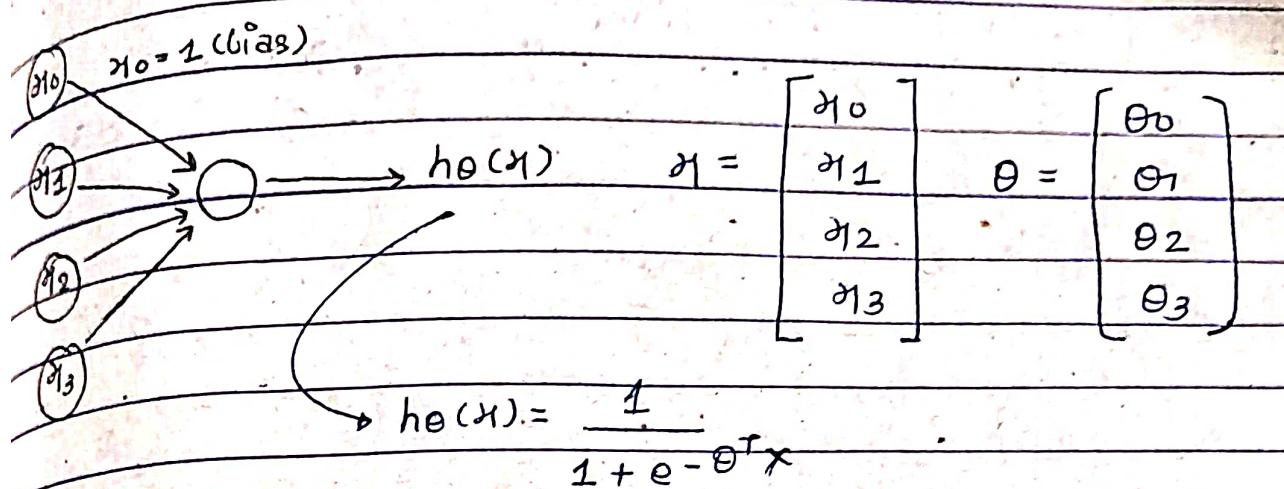
$$z^{[1]} = w^{[1]} a^{[0]} + b^{[1]} \\ (4 \times 1) \quad (4 \times 3) \quad (3 \times 1) \quad (4 \times 1)$$

$$d^{[1]} = \sigma(z^{[1]}) \\ (4 \times 1)$$

$$z^{[2]} = w^{[2]} d^{[1]} + b^{[2]} \\ (1 \times 1) \quad (1 \times 4) \quad (4 \times 1) \quad (1 \times 1)$$

$$d^{[2]} = \sigma(z^{[2]}) \\ (1 \times 1)$$

eurom Model o Logistic curve



Sigmoid (logistic) activation func: $g(z) = \frac{1}{1 + e^{-z}}$

⇒ Activation functions!

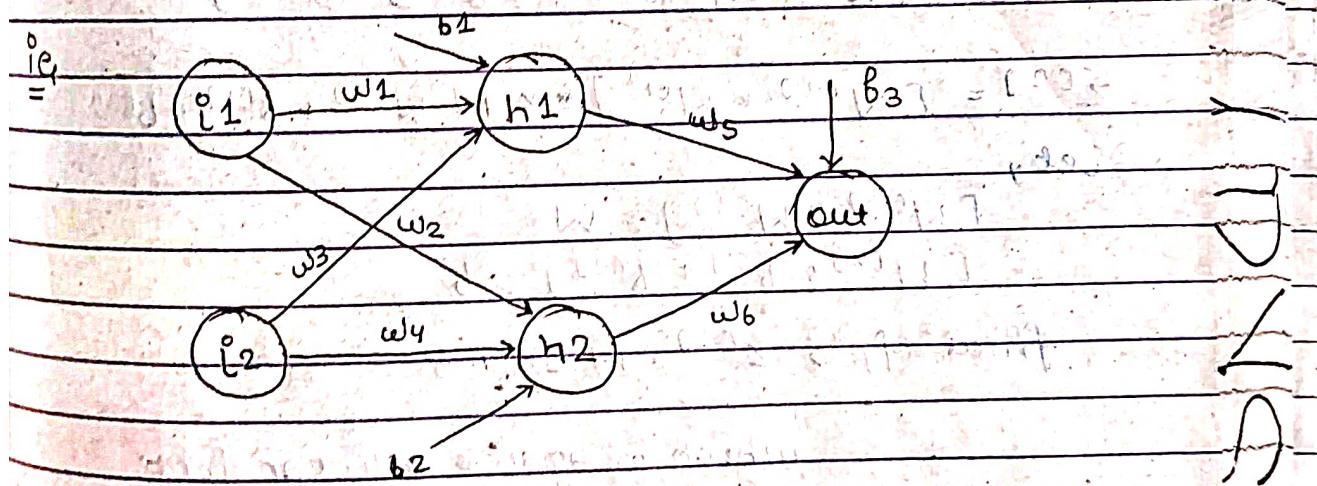
Why do we need activation functions?

We know that neural net has neurons that works in correspondence with weight, bias & their respective activation function.

In a neural net we would update the weight & biases of neurons on the basis of the error at the output. This process is known as Back-propagation

• Activation functions makes the back-propagation possible since the gradients are supplied along with the error to update the weight & biases

Q) Why do we need non-linear activation func?



hidden layer (layer 1)

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

• $z^{[1]}$ → vectorized off of layer 1

• $w^{[1]}$ → vectorized weights assigned to neurons of hidden layer i.e. $w_1^1, w_2^1, w_3^1, w_4^1$

→ w_1, w_2, w_3, w_4

- $x \rightarrow$ vectorized I/O features (i.e.) φ_1 and φ_2
- $b \rightarrow$ vectorized bias assigned to neurons in hidden layer i.e. B_1, B_2

$$\rightarrow b_1^1 \quad b_2^1$$

- $\varphi^{[1]} \rightarrow$ vectorized form of any linear func.

layer 2 (i.e. output) layer

$\varphi_{1/p}$ for layer 2 is o/p from layer 1

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

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

calculation at opp layer

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

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

Let,

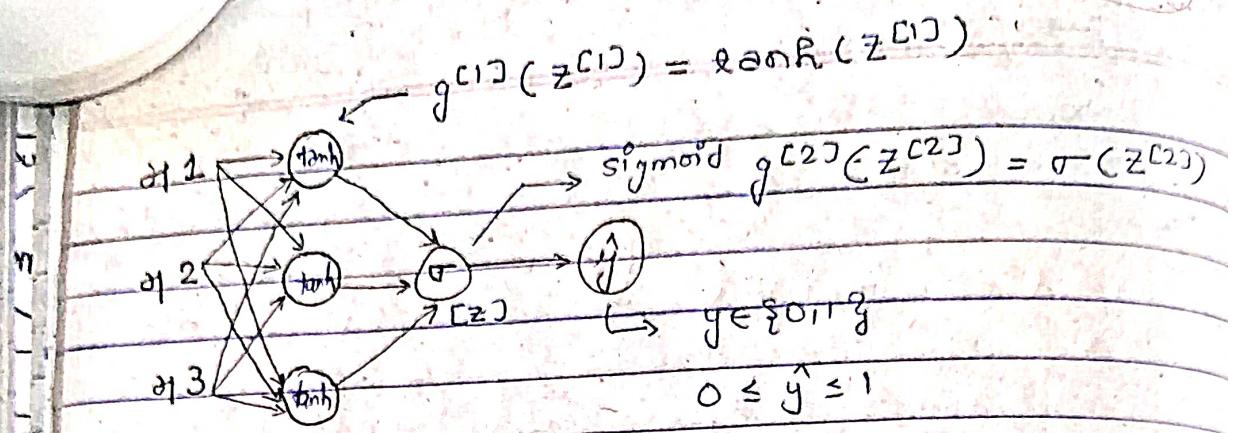
$$[W^{[2]} * W^{[1]}] = W$$

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

$$\text{final o/p} \rightarrow z^{[2]} = W * x + b$$

which is again a linear func
even after applying neural net

Hence, we can conclude that, doesn't matter how many hidden layer we attach in neural net, all layer will behave same way because, the composition of two linear func is linear itself



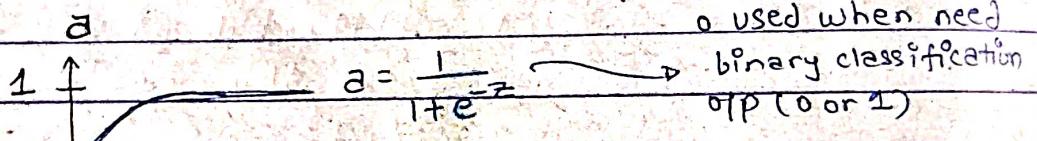
given α^0 :

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

$$\alpha^0 = \sigma(z^{[1]}) \rightarrow g^{[1]}(z^{[1]})$$

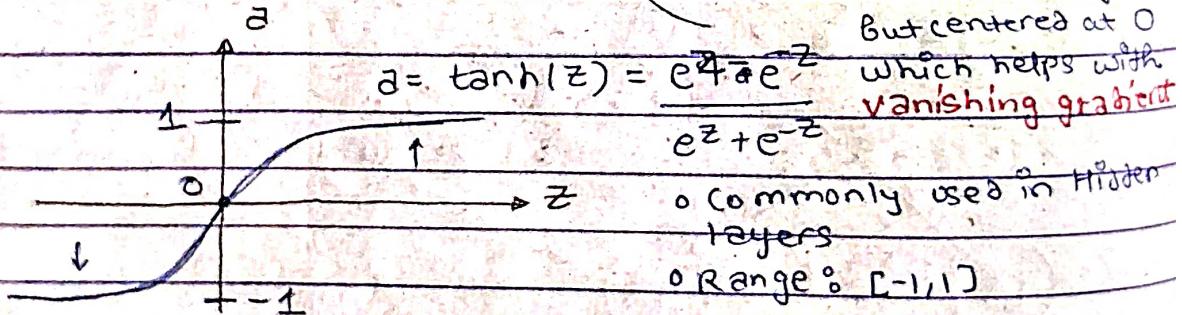
$$z^{[2]} = W^{[2]} \alpha^0 + b^{[2]}$$

$$\alpha^1 = \sigma(z^{[2]}) \rightarrow g^{[2]}(z^{[2]})$$

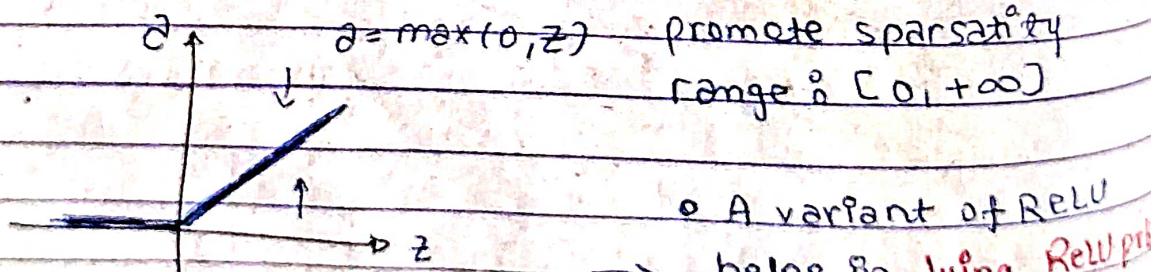


- Historically used in hidden layers not preferred due to vanishing gradient problem

- Similar to sigmoid but centered at 0 which helps with vanishing gradient



- fast computation & promote sparsity range in $[0, +\infty]$

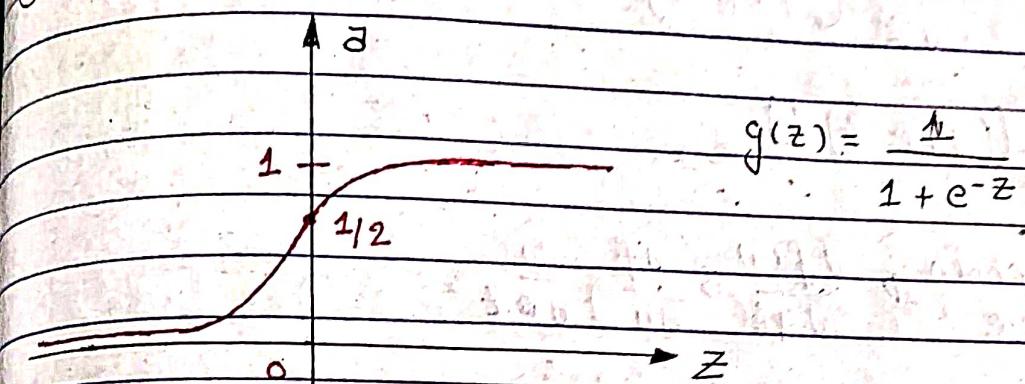


$a = \max(0.1z, z)$ allows a small non-zero gradient when the ReLU is zero.

- Help in training deeper networks
- Range: $[-\infty, +\infty]$

* Derivatives of Activation functions :

① Sigmoid Activation function



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

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

$$\Rightarrow g(z)(1-g(z)) \text{ or } z(1-z)$$

$$z=10 \rightarrow g(z) \approx 1 \quad g'(z) = z(1-z)$$

$$\frac{d}{dz} g(z) \approx 1(1-1) = 0.$$

$$z=-10 \rightarrow g(z) \approx 0$$

$$\frac{d}{dz} g(z) \approx 0(1-0) = 0$$

$$z=0 \rightarrow g(z) = \frac{1}{2}$$

$$\frac{d}{dz} g(z) \approx \frac{1}{2}(1-\frac{1}{2}) \Rightarrow \frac{1}{4}$$

derivation

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

$$\frac{d}{dz} ((1+e^{-z})^{-1})$$

Appley chain rule $\Rightarrow \frac{df(u)}{dz} \Rightarrow \frac{df}{du} \cdot \frac{du}{dz}$

$$u = (1+e^{-z}), f = u^{-1}$$

$$\Rightarrow \frac{d}{du} (u^{-1}) \frac{d}{dz} ((1+e^{-z}))$$

$\hookrightarrow -1/u^2$

$$\Rightarrow \frac{1}{(1+e^{-z})^2} \frac{d}{dz} (1+e^{-z})$$

$$\Rightarrow \frac{1}{(1+e^{-z})^2} \cdot e^{-z} \frac{d}{dz} (-z) \Rightarrow \frac{e^{-z}}{(1+e^{-z})^2}$$

generalize,

$$\frac{1+e^{-z}-1}{(1+e^{-z})^2} \Rightarrow \frac{1}{1+e^{-z}} - \frac{1}{(1+e^{-z})^2}$$

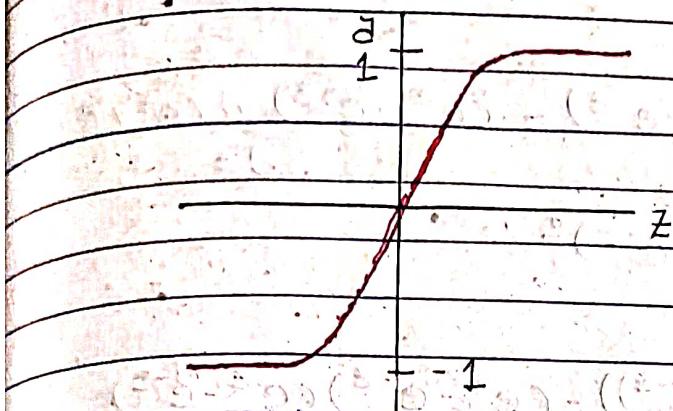
$$\Rightarrow \frac{1}{(1+e^{-z})} \cdot \left(1 - \frac{1}{1+e^{-z}} \right)$$

$$\Rightarrow g(z)(1-g(z))$$

or

$$(z \in 1-g(z))$$

② Tanh Activation function



$$g(z) = \tanh(z)$$

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

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

$$\Rightarrow 1 - (\tanh(z))^2$$

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

• $z = 10 \rightarrow \tanh(z) \approx 1$ • $z = -10 \rightarrow \tanh(z) \approx -1$

$g'(z) \approx 0$ $g'(z) \approx 0$

• $z = 0 \rightarrow \tanh(z) = 0$, $g'(z) = 1$

derivation,

The quotient rule: $\frac{d}{dx} \left(\frac{u}{v} \right) \Rightarrow \frac{v \left(\frac{du}{dx} \right) - u \left(\frac{dv}{dx} \right)}{v^2}$

By chain rule: $\frac{d}{dx} e^{ax} \Rightarrow e^{ax} \cdot \frac{d}{dx} (ax) \Rightarrow a e^{ax}$

As a particular case of this we have,

$$\frac{d}{dz} e^{-z} \Rightarrow -e^{-z}$$

$$g'(z) = \frac{d}{dz} (g(z))$$

$$(e^z + e^{-z}) \frac{d}{dz} (e^z - e^{-z}) - (e^z - e^{-z}) \frac{d}{dz} (e^z + e^{-z})$$

$$(e^z + e^{-z})^2$$

$$(e^z + e^{-z})(e^z - (-e^{-z})) - (e^z - e^{-z})(e^z + e^{-z})$$

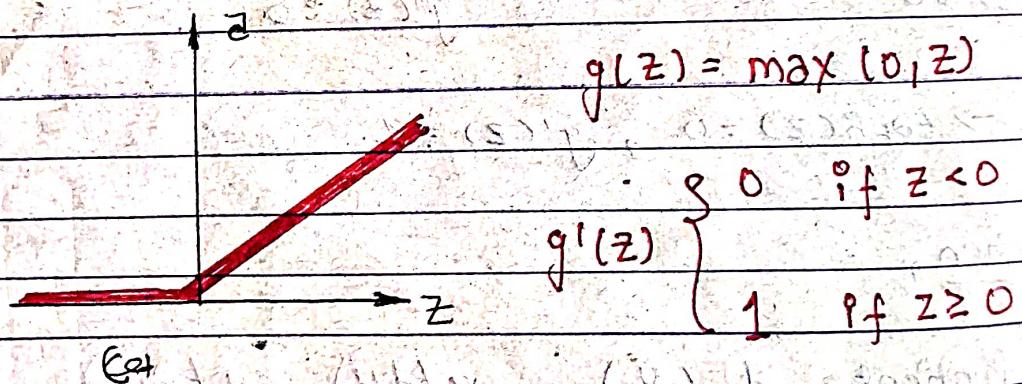
$$(e^z + e^{-z})^2$$

$$1 - (\tanh(z))^2$$

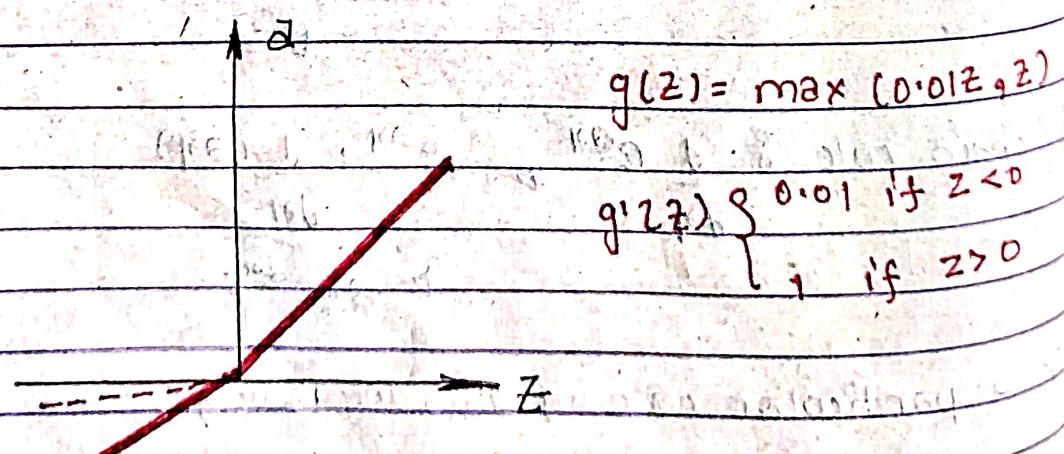
or

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

③ ReLU Activation function



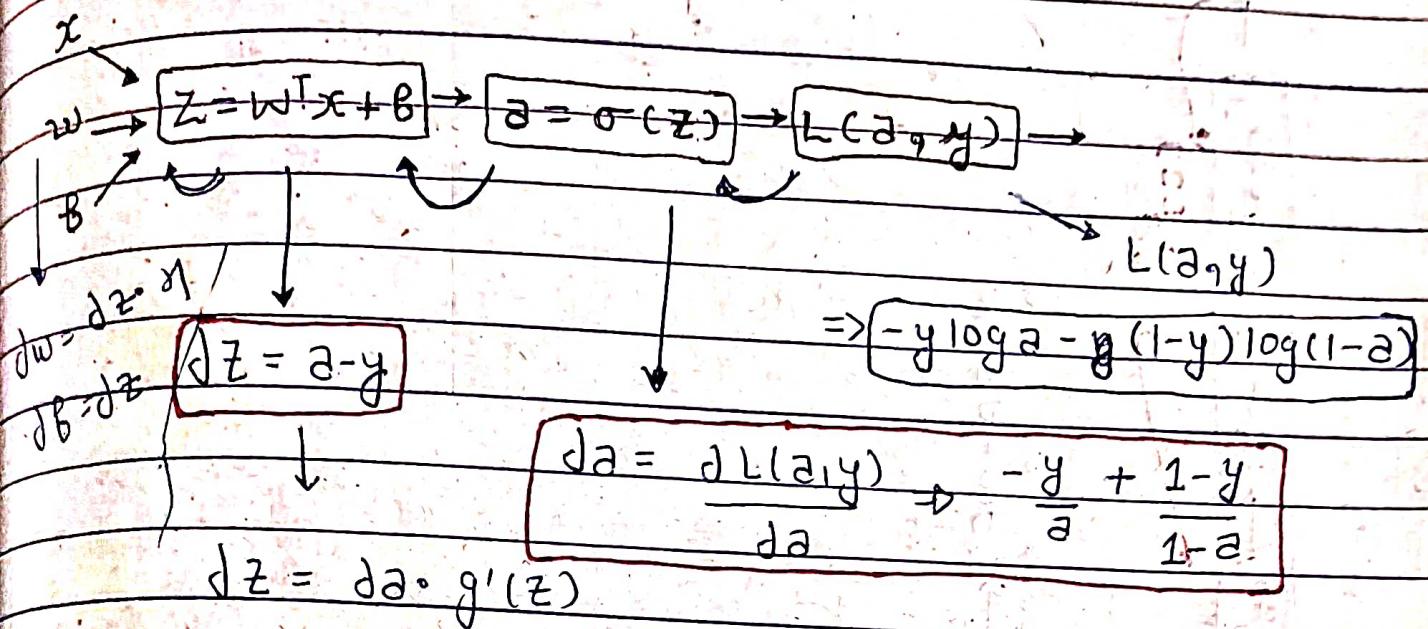
④ Leaky ReLU Activation function



Gradient Descent for Neural Networks

Date _____
Page _____

Computing gradients -



$$\{ g(z) = \sigma(z) \}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \left(\frac{\partial a}{\partial z} \right)$$

$$\downarrow \quad \downarrow \quad \downarrow$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{d}{dz} g(z) = g'(z)$$

$$\begin{aligned}
 & \boxed{\beta_{1,1}^2 = \beta_{1,2}^2} \\
 & \beta_{1,1}^2 = \beta_{1,2}^2 = N = M_{1,1}^2 = M_{1,2}^2 \\
 & \beta_{1,1}^2 = \beta_{1,2}^2 = M_{1,1}^2 = M_{1,2}^2 \\
 & \beta_{1,1}^2 = \beta_{1,2}^2 = M_{1,1}^2 = M_{1,2}^2 \\
 & \beta_{1,1}^2 = \beta_{1,2}^2 = M_{1,1}^2 = M_{1,2}^2
 \end{aligned}$$

$$\begin{array}{l}
 \boxed{\Delta C_{11} = C_{11}} \\
 \boxed{\Delta C_{12} = C_{12}} \\
 \boxed{\Delta C_{21} = C_{21}} \\
 \boxed{\Delta C_{22} = C_{22}} \\
 \Delta M_{11} = M_{11} \\
 \Delta M_{12} = M_{12} \\
 \Delta M_{21} = M_{21} \\
 \Delta M_{22} = M_{22} \\
 \Delta X_{11} = X_{11} \\
 \Delta X_{12} = X_{12} \\
 \Delta X_{21} = X_{21} \\
 \Delta X_{22} = X_{22} \\
 \Delta P_{11} = P_{11} \\
 \Delta P_{12} = P_{12} \\
 \Delta P_{21} = P_{21} \\
 \Delta P_{22} = P_{22} \\
 \Delta A = A \\
 \Delta Z = Z \\
 \Delta Y = Y \\
 \Delta \rho = \rho \\
 \Delta \theta = \theta
 \end{array}$$

Gradient descent

$$\delta z^{[2]} = A^{[2]} - y$$

$$\delta w^{[2]} = \frac{1}{m} \delta z^{[2]} A^{[1]T}$$

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

$$\delta z^{[1]} = W^{[2]T} \delta z^{[2]} * g^{[1]'}(z^{[1]})$$

$$\delta w^{[1]} = \frac{1}{m} \delta z^{[1]} X.T$$

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