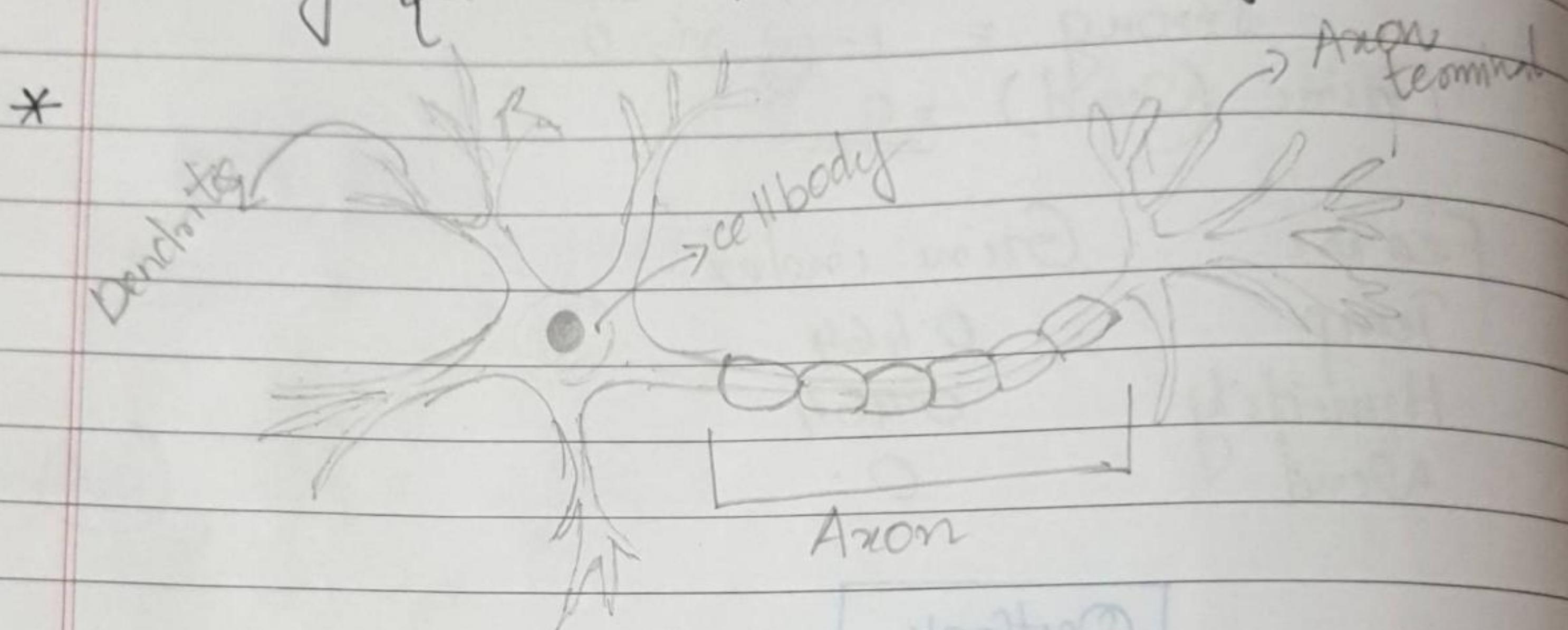
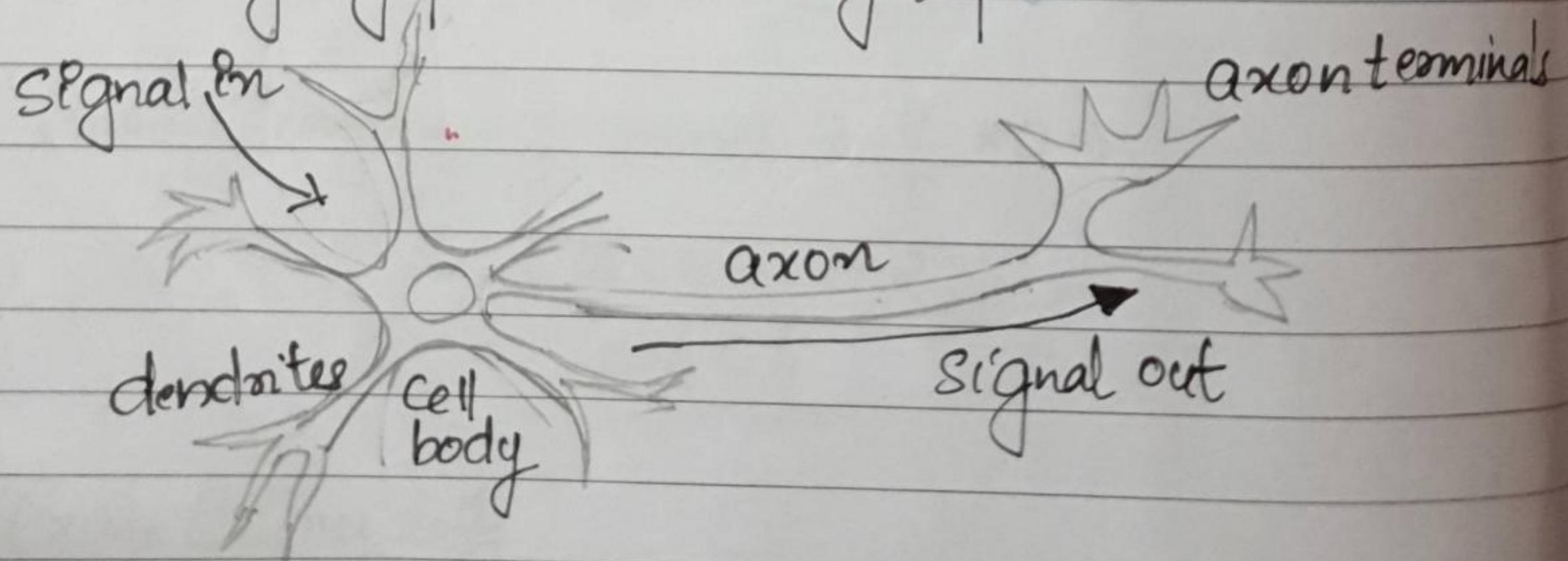


## Neural Networks:-

- \* Motivation behind neural network is human brain. Human brain contains billion of neurons which are connected to many other neurons to form a network so that if we sense anything, it recognizes and processes the output.

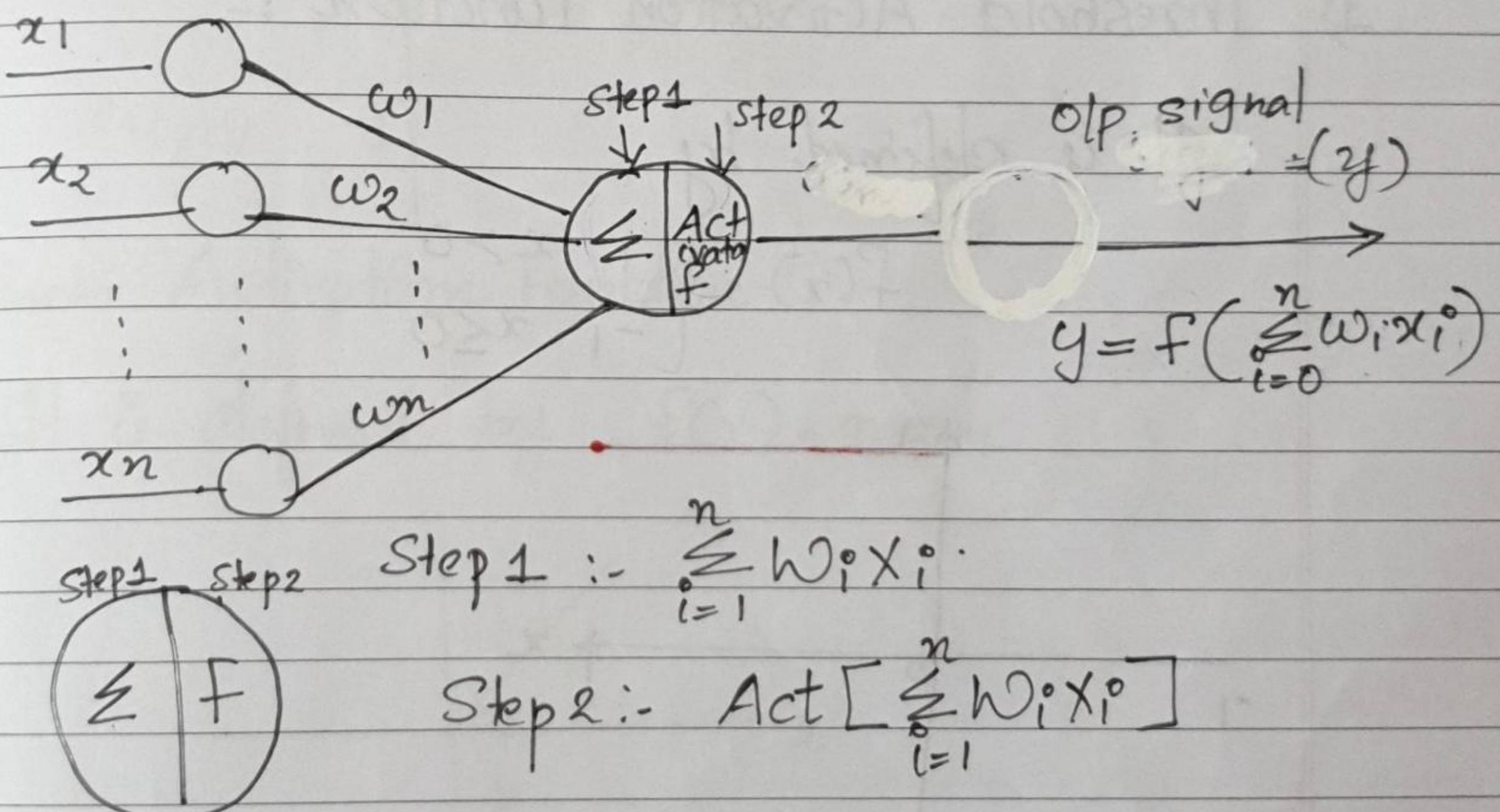


- \* Dendrites receive signals from other neurons.
- \* Cell body sums the incoming signals to generate input.
- \* When the sum reaches a threshold value, neuron fires and the signal travels down the axon to the other neurons through a tiny gap called Synapse.



## Artificial Neurons:-

- \* An artificial neuron is a mathematical function
- \* Artificial neurons are elementary units in an artificial neural network.
- \* Artificial neuron receives one or more inputs and sums them to produce an output.
- \* Each input is separately weighted, and the sum is passed through a function known as an activation function or transfer function.



$x_1, \dots, x_n$  = input signals

$w_1, w_2, \dots, w_n$  = weights associated with ip signals

$\sum$  = summation of input signals

$f$  = activation function which produces the output

$y$  = o/p signal

- \* Circles are the nodes of the neurons
- \*  $x_1, x_2, \dots, x_n$  are the input nodes
- \* Circle in the second layer outputs the value 'y' is called output node.

⑥

## Activation functions:-

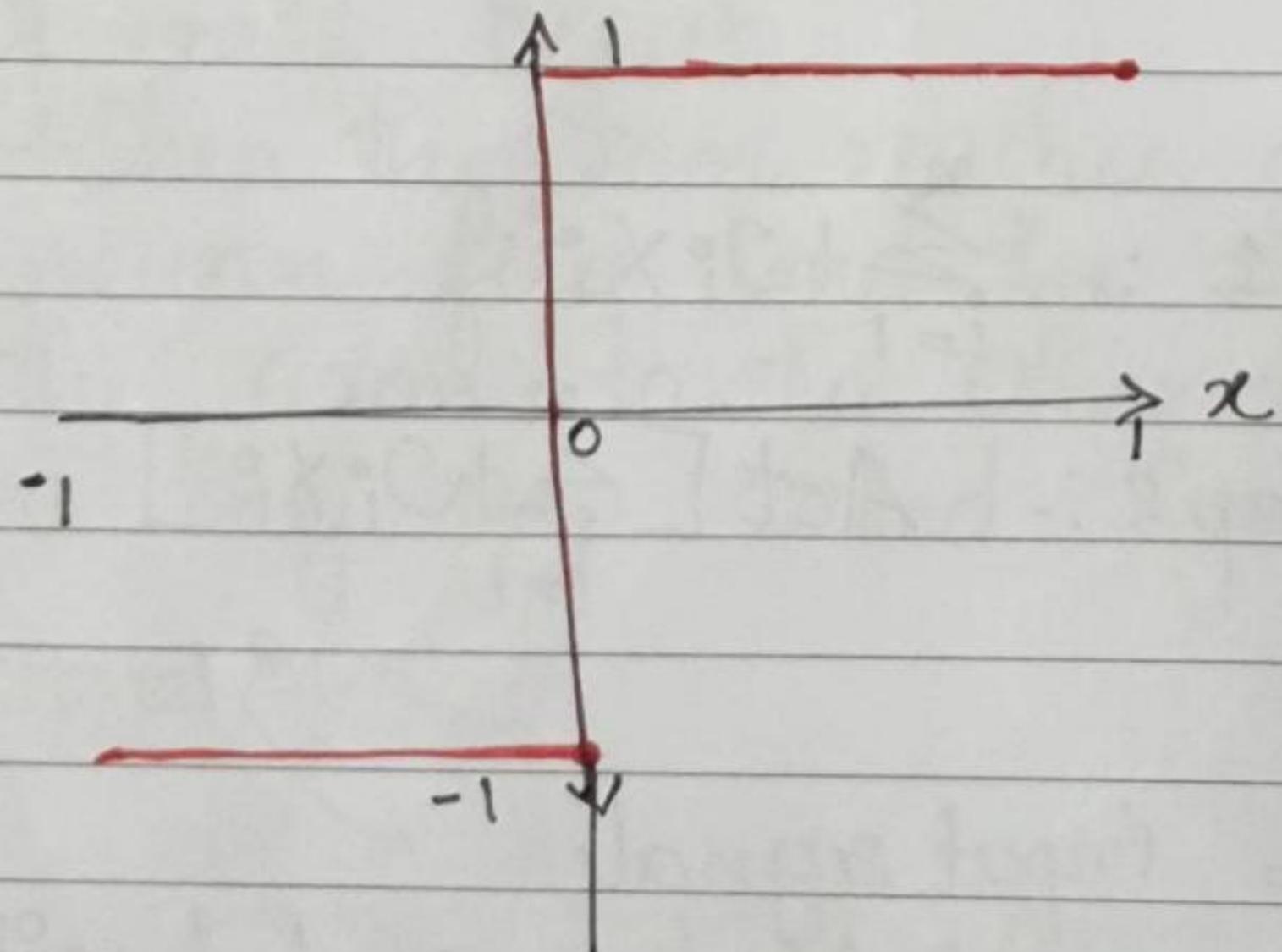
- \* The function which takes the incoming signals as input and produces the output signal is known as activation function

### i) Threshold Activation Function :-

It is defined by

$$f(x) = \begin{cases} 1 & x > 0 \\ -1 & x \leq 0 \end{cases}$$

Bipolar

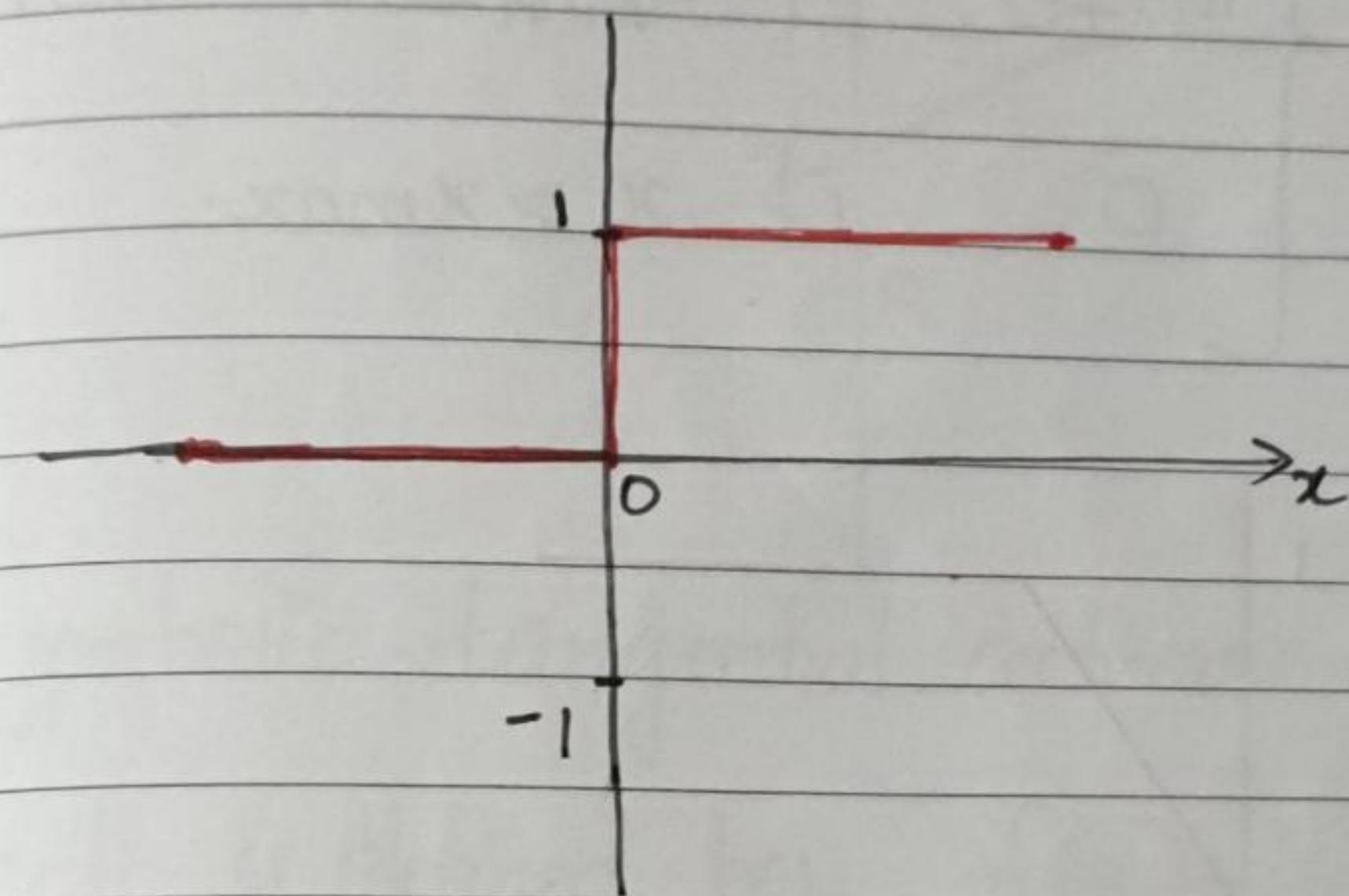


## 2) Unit Step Function:

By Sagar Binary

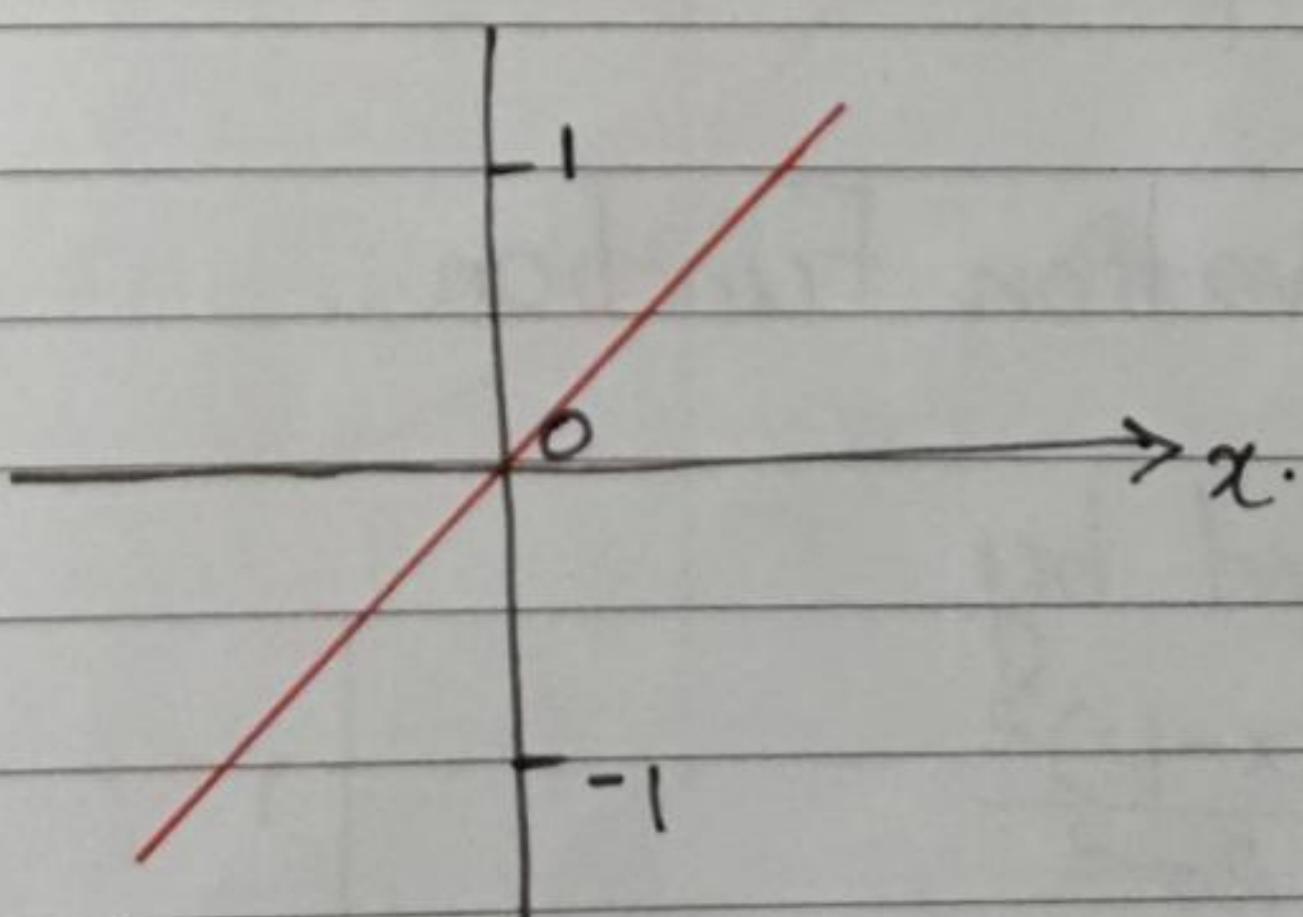
Sometimes the threshold activation function is also defined as a unit step function, <sup>in such cases</sup> then it is called unit step activation function.

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



## 3) Linear Activation Function:

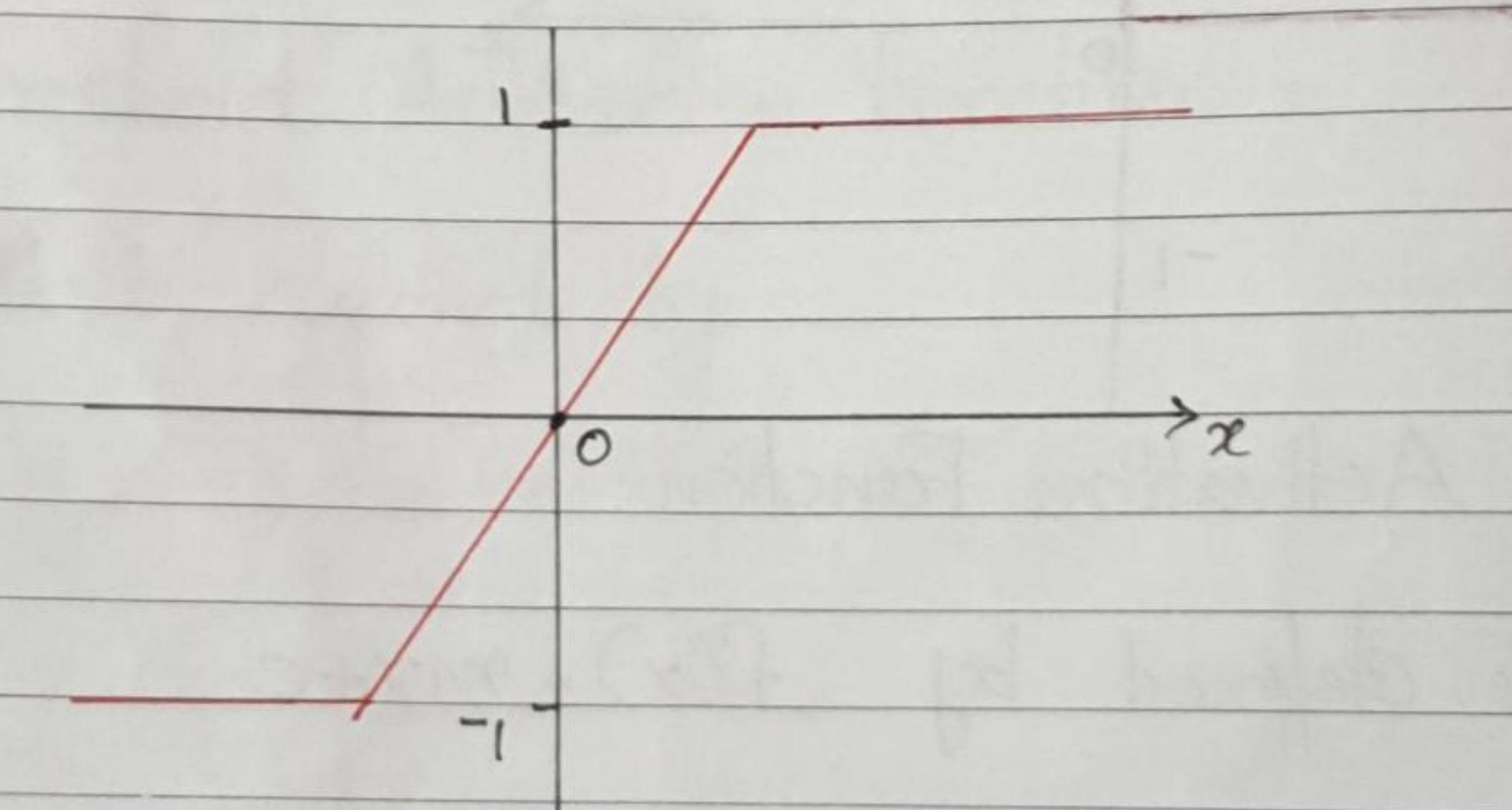
It is defined by  $f(x) = mx + c$



## 4) Piecewise / Saturated Linear Activation function

It is defined by

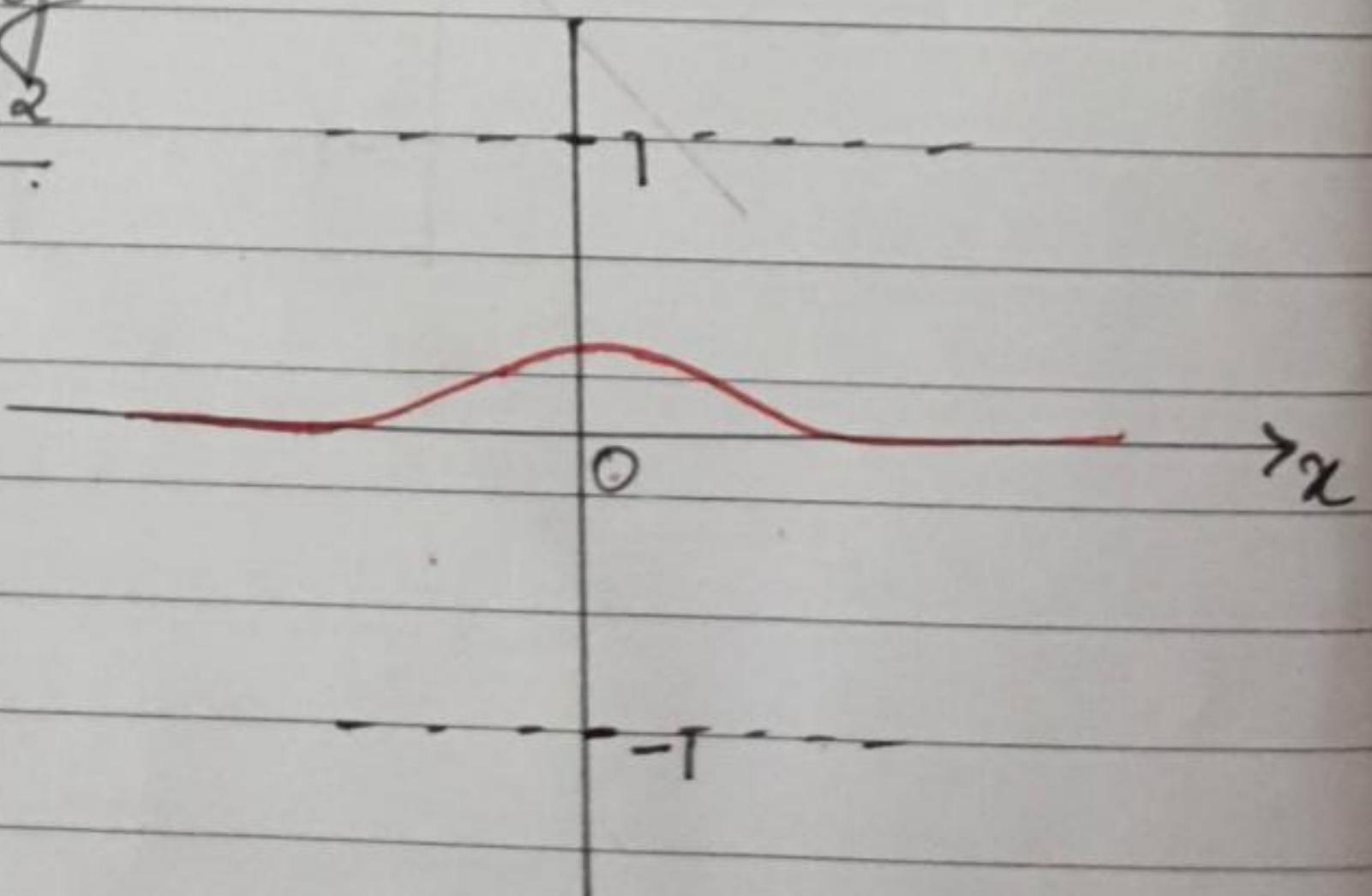
$$f(x) = \begin{cases} 0 & \text{if } x < x_{\min} \\ mx + c & \text{if } x_{\min} \leq x \leq x_{\max} \\ 0 & \text{if } x > x_{\max} \end{cases}$$



## 5) Gaussian Activation Function:

It is defined by

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



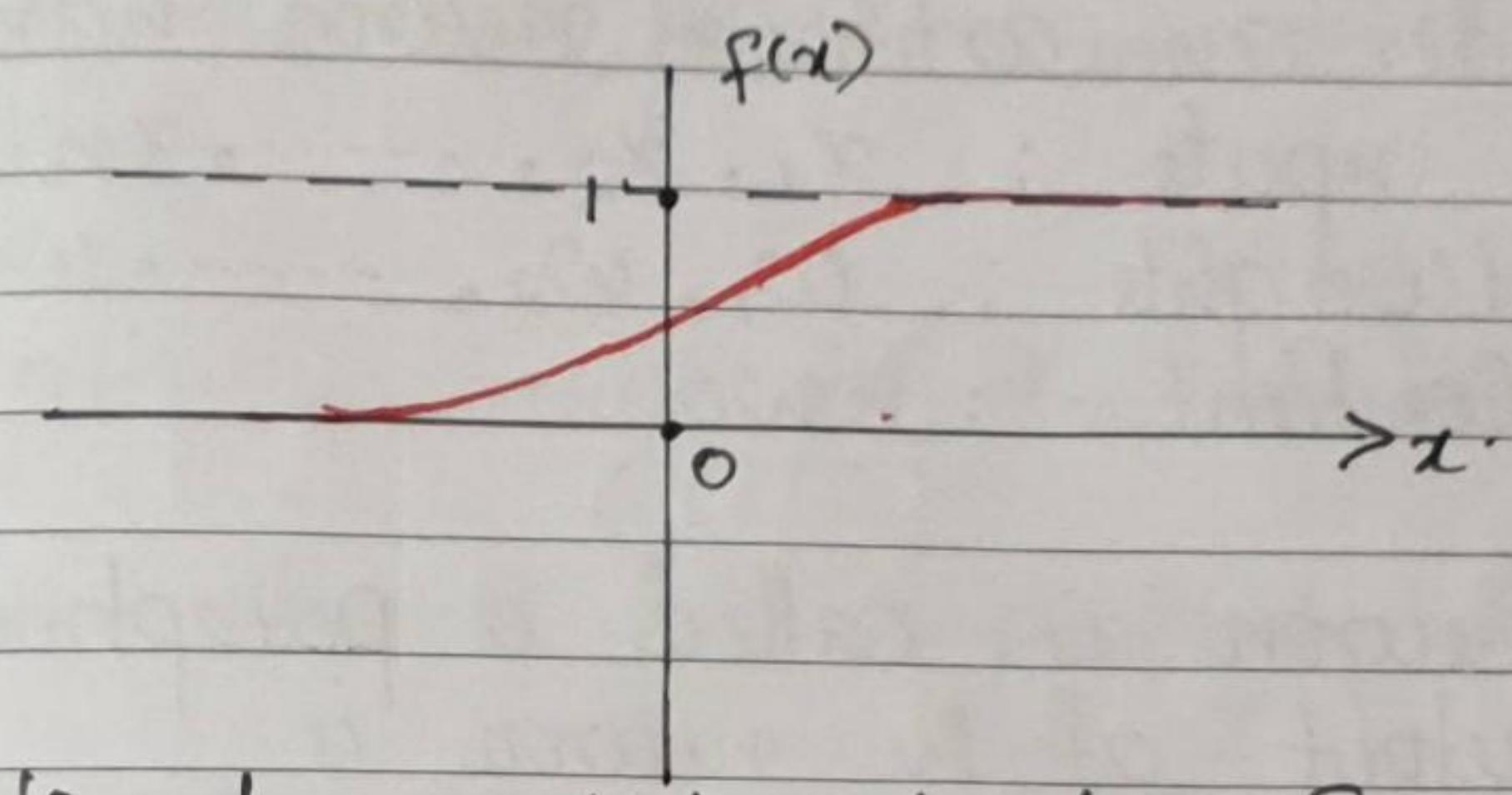
## 6) Sigmoid activation function [Logistic function]

- Most commonly used activation functions
- It is defined as

$$f(x) = \frac{1}{1+e^{-x}}$$

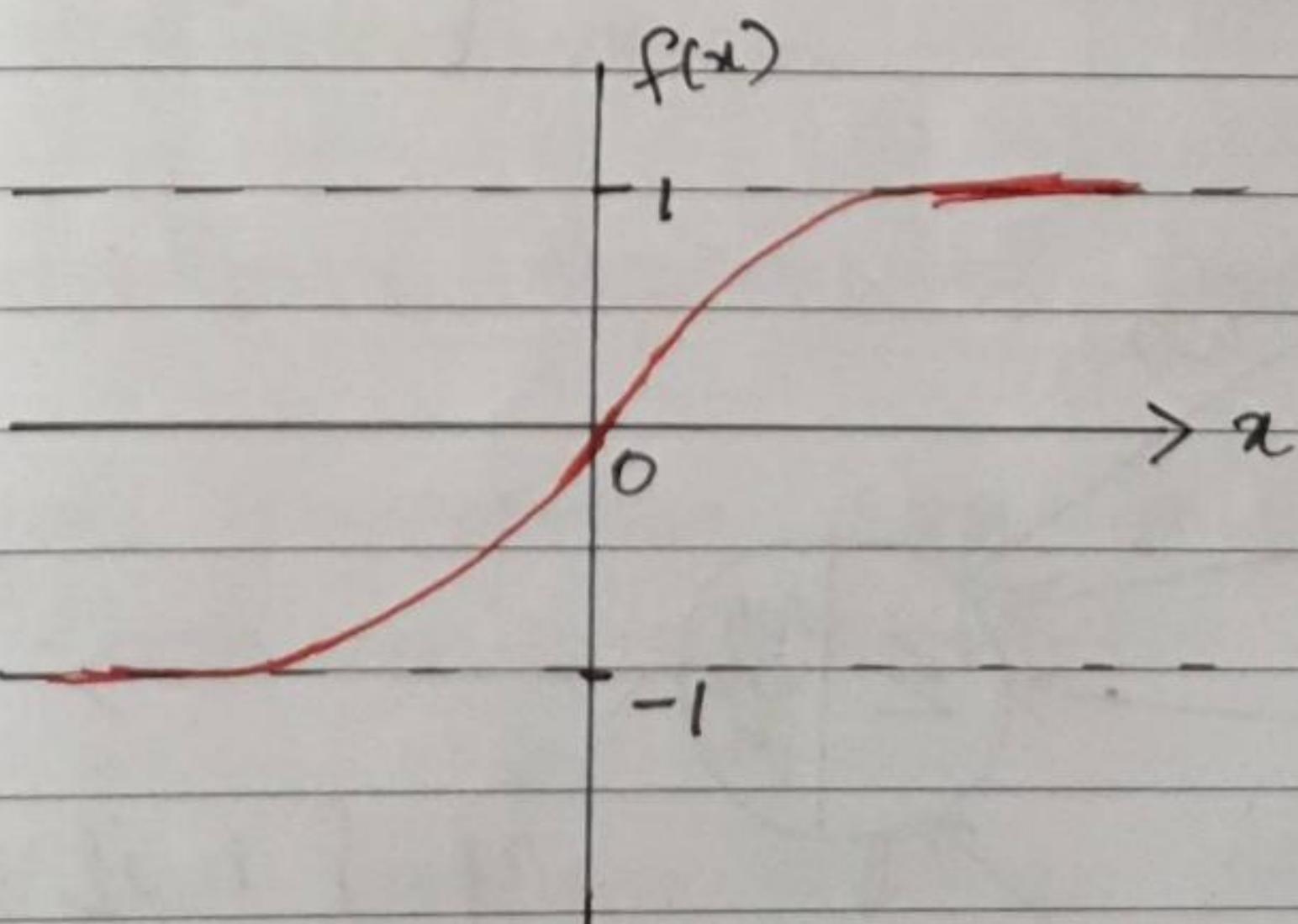
Bipolar:

$$\rightarrow \frac{2}{1+e^{-x}} - 1$$



## 7) Hyperbolic tangential activation function

- It is defined by  $f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$



(15b)

## Perceptron

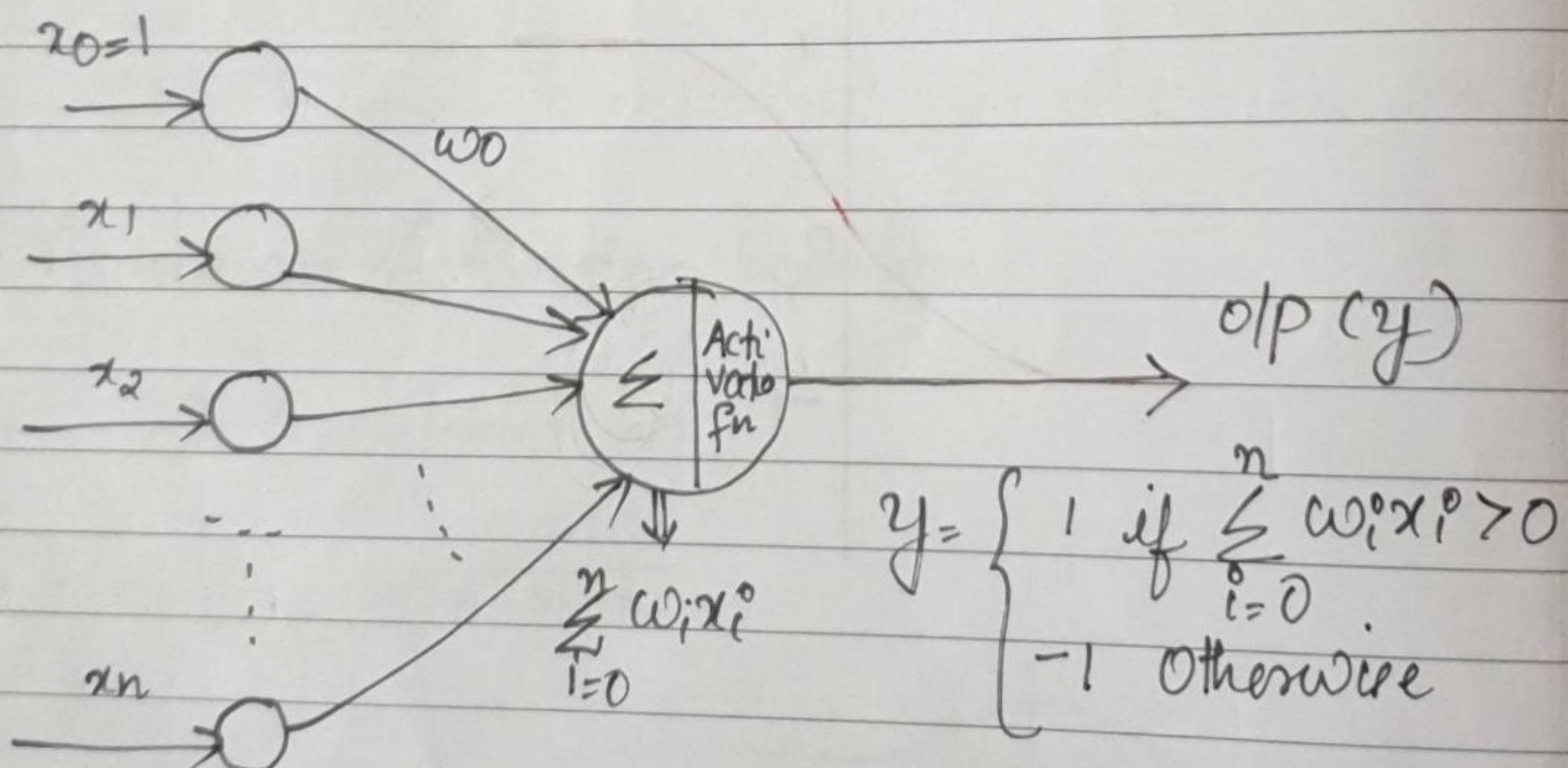
## Single-layer N.N.

- \* The perceptron is a special type of artificial neuron in which the activation function is the threshold function.
- \* Basic neural network building block
- \* Consider an artificial neuron having
  - inputs :  $x_1, x_2, \dots, x_n$
  - associated weights :  $w_1, w_2, \dots, w_n$
  - Constant :  $w_0$

The neuron is called a perceptron if the output of the neuron is

Output :-

$$O(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1x_1 + \dots + w_nx_n > 0 \\ -1 & \text{if } w_0 + w_1x_1 + \dots + w_nx_n \leq 0 \end{cases}$$



## Working of a perceptron:-

- (a) All inputs ' $x_i$ ' are multiplied by weights ' $w_i$ '.
- (b) Add all multiplied values and call it as weighted sum.
- (c) Apply the weighted sum to the activation function.
- (d) Output the result of the activation function.

weights :- Strength of a particular input

Basic Building blocks of Artificial neural network:-

- 1) Artificial Neuron
- 2) Weights
- 3) Bias → allows to shift the activation function by adding a constant.
- 4) Activation Function

## Representational power of perceptrons :-

A single perceptron can be used to represent many boolean functions.

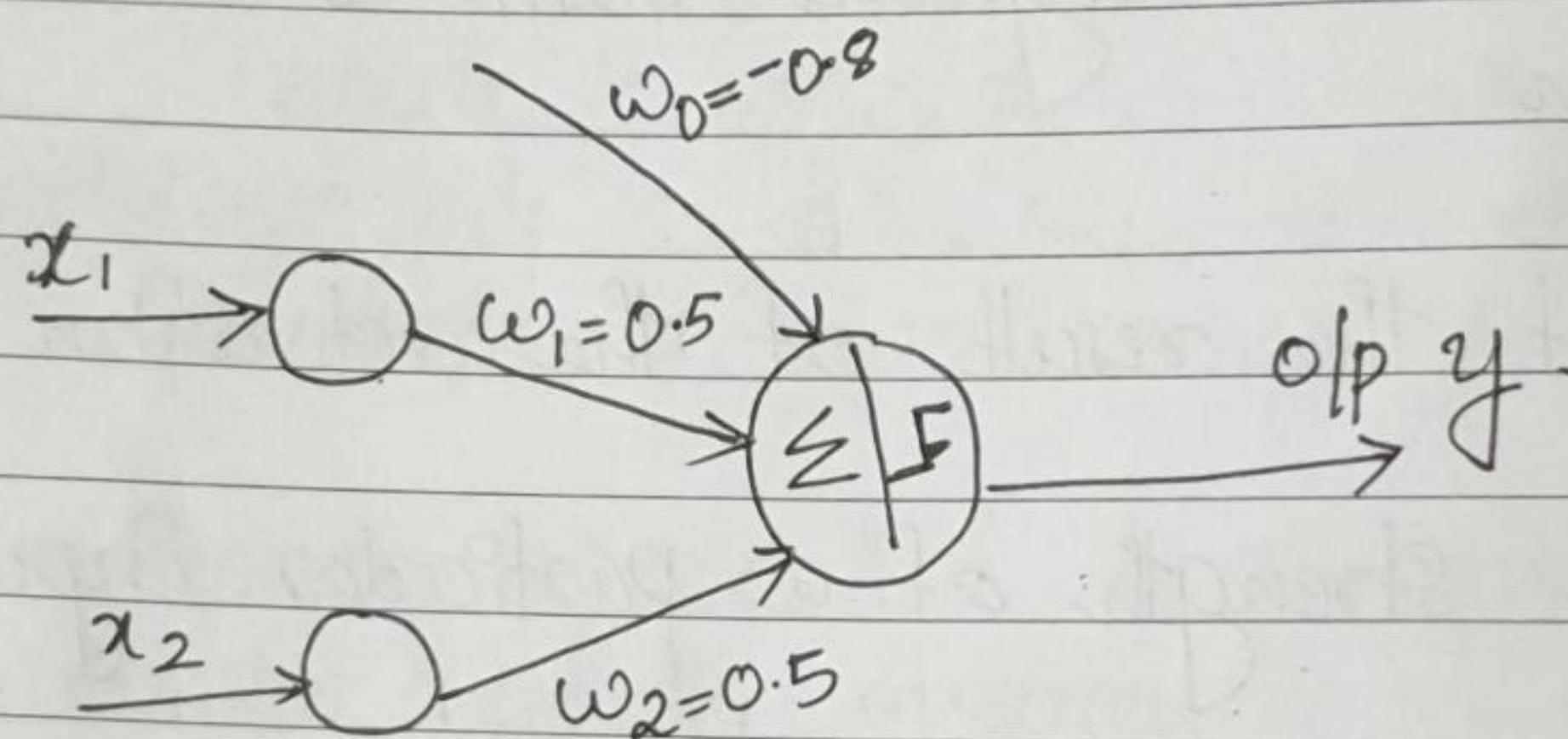
Eg:- If we assume boolean value of 1 (true) and -1 (false), then one way to represent a '2' input perceptron to implement AND function is to set weights.

$$w_0 = -0.8 ; w_1 = w_2 = 0.5$$

## AND Truth table

Perception value

A	B	T	
0	0	0	-1 (false)
1	0	0	-1 (false)
0	1	0	-1 (false)
1	1	1	1 (True)



$$(x_1, x_2) = (0, 0) \quad w_0 + 0.5 \times 0 + 0.5 \times 0 = -0.8 + 0.5 \times 0 + 0.5 \times 0 = -0.8 < 0 \Rightarrow -1$$

$$(x_1, x_2) = (1, 0) \quad w_0 + 0.5 \times x_1 + 0.5 \times x_2 = -0.8 + 0.5 \times 1 + 0.5 \times 0 = -0.8 + 0.5 = -0.3 < 0 \Rightarrow -1$$

$$(x_1, x_2) = (0, 1) \quad w_0 + 0.5 \times x_1 + 0.5 \times x_2 = -0.8 + 0.5 \times 0 + 0.5 \times 1 = -0.8 + 0.5 = -0.3 < 0 \Rightarrow -1$$

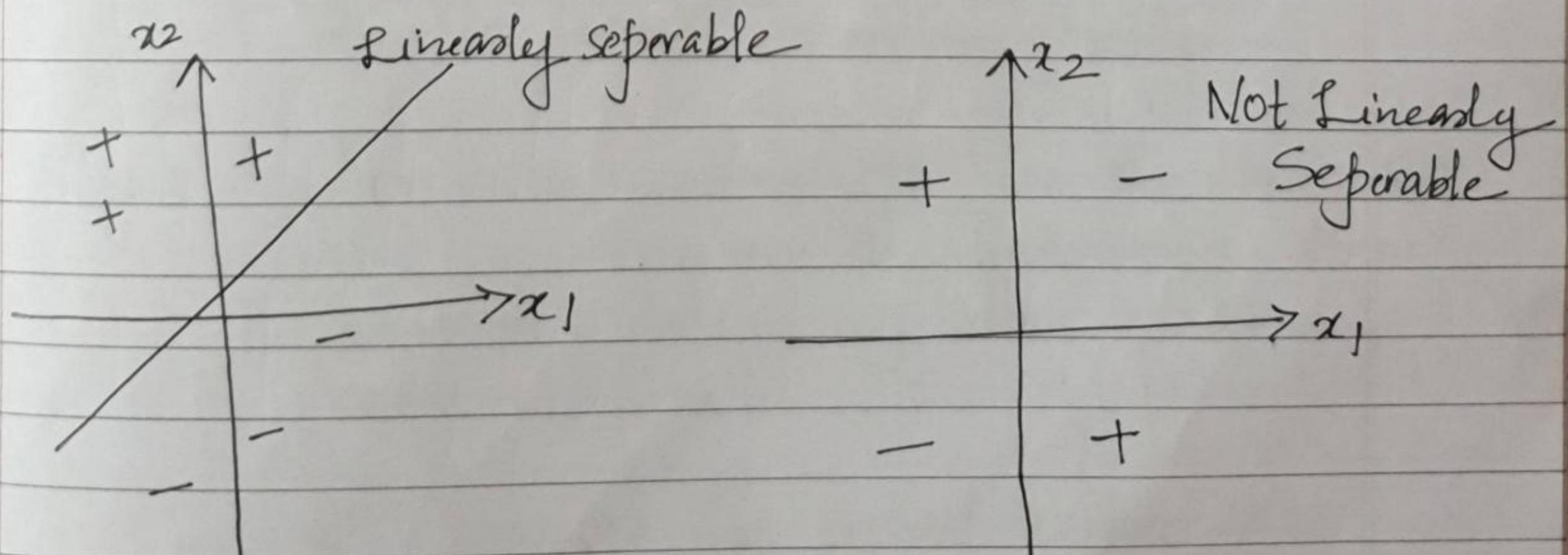
$$(x_1, x_2) = (1, 1) \quad w_0 + 0.5 \times x_1 + 0.5 \times x_2 = -0.8 + 0.5 + 0.5 = -0.8 + 1 = 0.2 > 0 \Rightarrow 1$$

Similarly OR function can be represented by modifying  $w_0 = -0.3$

- NAND function,  $w_0 = 0.8$
- NOR function,  $w_0 = 0.3$ .

Boolean function	$w_0$	$w_1$	$w_2$
$x_1 \text{ AND } x_2$	-0.8	0.5	0.5
$x_1 \text{ OR } x_2$	-0.3	0.5	0.5
$x_1 \text{ NAND } x_2$	0.8	-0.5	-0.5
$x_1 \text{ NOR } x_2$	0.3	-0.5	-0.5

- Perceptrons can represent all primitive boolean functions like AND, OR, NAND, and NOR.
- However XOR requires more than one single perceptron whose op is 1 if and only if  $x_1 \neq x_2$ . It represents the set of linearly non separable training examples.



# Perceptron Learning Algorithm

P  $\leftarrow$  inputs with label 1;

N  $\leftarrow$  P inputs with label 0;

Initialize  $w = [w_1, w_2, \dots, w_n]$  randomly

while ! Convergence do

    pick random  $x \in P \cup N$ ;

    if  $x \in P$  and  $\sum_{i=0}^n w_i * x_i < 0$  then

$w = w + x$ ;

    end

$w = n+1$  dimension

$x = n$  dimension

    if  $x \in N$  and  $\sum_{i=0}^n w_i * x_i \geq 0$  then

$w = w - x$ ;

    end

end

// the algorithm converges when all the  
inputs are classified correctly

# Perceptron Learning Algorithm

 $n$ 

: No. of input variables -

 $y$ 

: OLP from the perceptron

 $D$ 

: Training set of 'S' samples.

 $x = (x_{j0}, x_{j1}, \dots, x_{jn})$ :  $n$ -dimensional input $d_j^o$ : Desired olp value of the perceptron for the input  $x_j$  $x_{ji}$ : value of the  $i$ -th feature of the  $j$ -th training input $x_{j0}$ 

: 1

 $w_i$ : weight of the  $i$ -th input variable $w_{i(t)}$ : weight  $i$  at the  $t$ -th iteration.

Step 1 : Initialize the weights and the threshold.  
 (Weights may be initialized to '0' or a to a small random value.)

Step 2 : For each example  $j$  in the training set  $D$ , perform the following steps over the input  $x_j$  and desired output  $d_j^o$ .

a) Calculate the actual output

$$y_j(t) = f [w_{0(t)}x_{j0} + w_{1(t)}x_{j1} + w_{2(t)}x_{j2} + \dots + w_{n(t)}x_{jn}]$$

b) Update the weights :

$$w_i^o(t+1) = w_i^o(t) + (d_j - y_j(t))x_{ji}$$

for all features  $0 \leq i \leq n$ .

Step 3 : Step 2 is repeated until the iteration error

$$\frac{1}{S} \sum_{j=1}^S |d_j - y_j(t)|$$

is less than a user specified error threshold

$r$ ; or a predetermined no. of iterations have been completed.

$S$  = size of the sample set

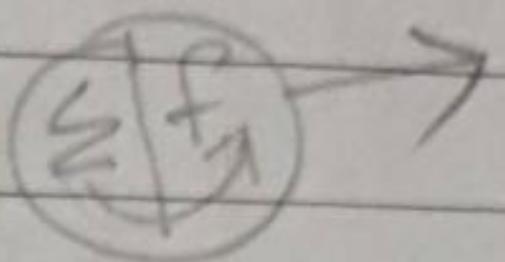
### Artificial Neural Networks:

- \* ANN is based on a collection of connected units called artificial neurons
- \* Each connection between artificial neurons can transmit a signal from one to another
- \* The artificial neuron that receives the signal can process it and then signal other artificial neurons connected to it.

## ★ Characteristics of ANN

### 1) The activation function

- \* This function defines how a neuron's combined input signals are transformed into a single output signal



### 2) The network topology

- \* describes the number of neurons in the model as well as the number of layers and the way they are connected
- \* It determines the complexity of the task

Different forms of topology can be differentiated by the following characteristics

#### a) The no: of layers :-

Input nodes:- Those nodes which receive unprocessed signals directly from the input data.

nodes are arranged in layers } (feed layer)

Output nodes:- Those nodes which generate the final predicted values.

Hidden nodes : a node that processes the signals from the input nodes (or other such nodes) prior to reaching output nodes

b) Whether information in the network is allowed to travel backward

Feedforward networks:-

Networks in which the input signal is fed continuously in one direction from connection to connection until it reaches the output layer called feedforward networks.

Feedback networks:-

Networks which allows signals to travel in both directions using loops are called feedback networks / Recurrent networks

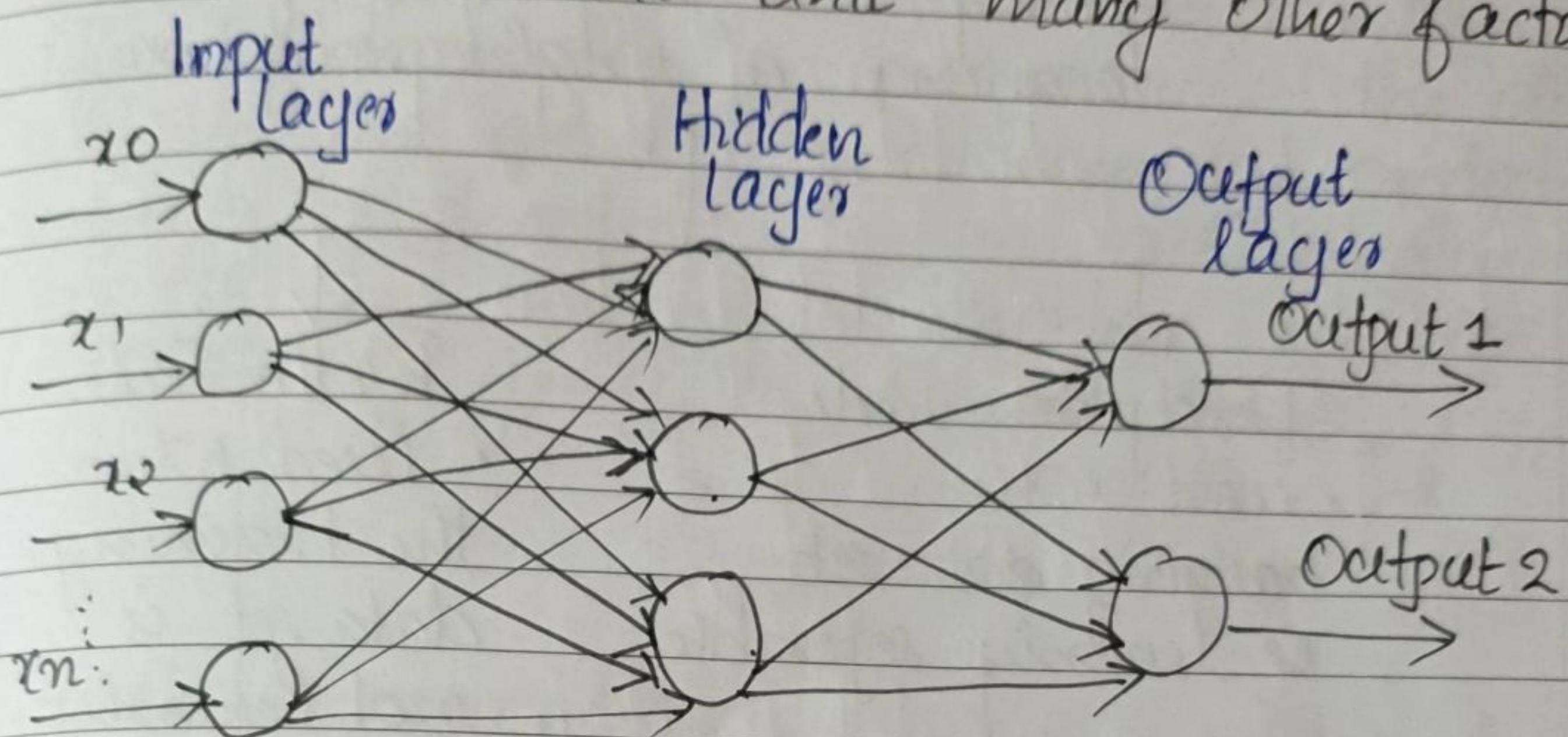
c) The number of nodes within each layer of the network

- Number of features in the input data  $\Rightarrow$  the number of input nodes

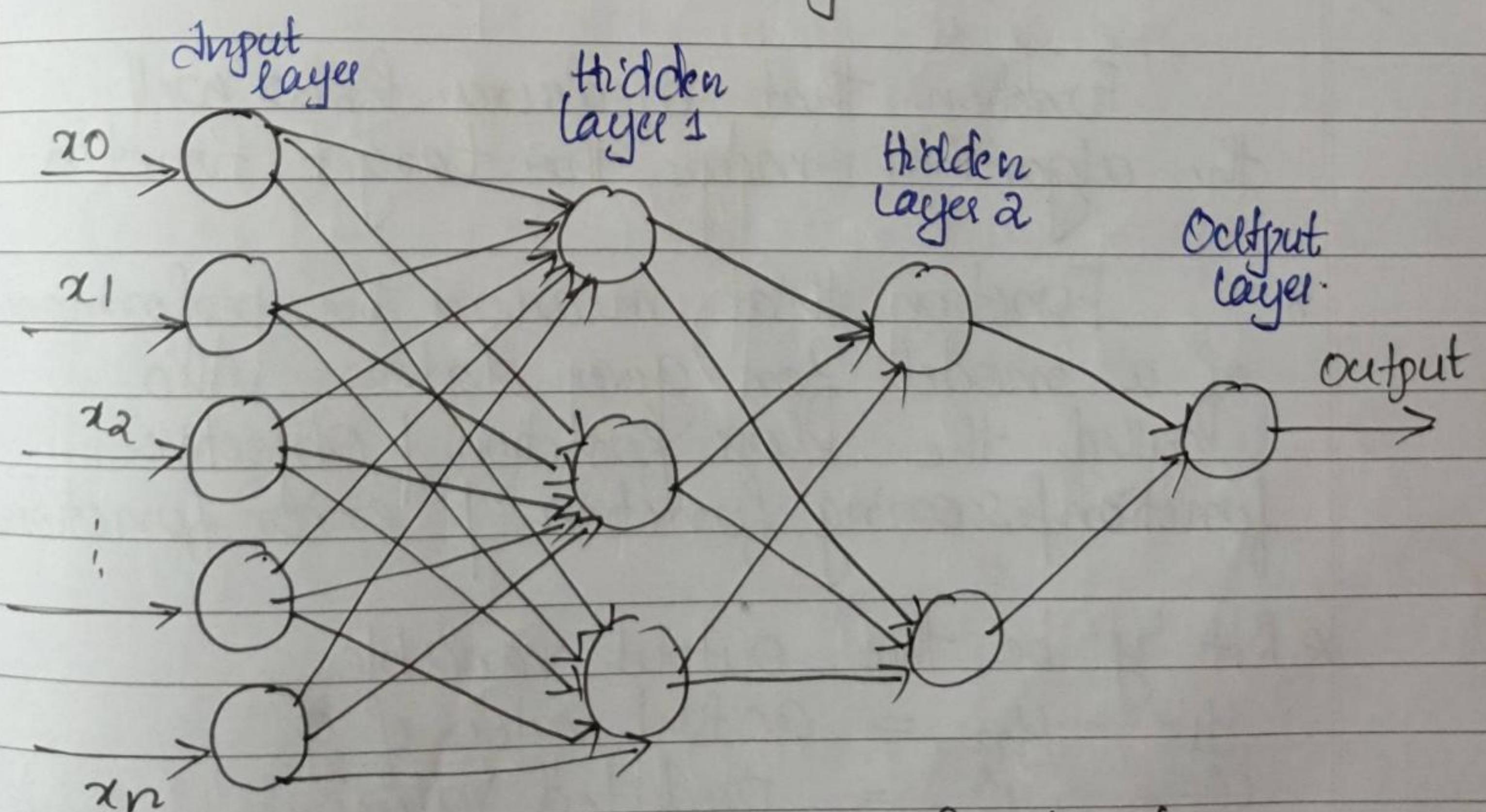
- Number of output nodes = no: of outcomes to be modeled

- Number of hidden nodes  
 $\hookrightarrow$  user to decide prior to training the model

→ The appropriate number depends on no: of input nodes, amount of training data, amount of noisy data, complexity of the learning task and many other factors.



Network with one hidden layer and two output nodes.



Network with two hidden layers.

### 3) The training algorithm

- algorithm specifies how connection weights are set in order
- Mainly two algorithms for learning a single perception

↓  
Perception Rule

- used when training data set is linearly separable

↓  
Delta Rule

- used when the training data set is not linearly separable

### 4) Cost Function:-

~~Function that measures how well the algorithm maps the target function~~

\* Function that measures the performance of a model for given data. Also called the Loss function | Objective function | scoring function | error function

\* Let 'y' be the output variable

$y_1, \dots, y_n$  = actual values of y

$\hat{y}_1, \dots, \hat{y}_n$  = predicted values

SSE                          MSE

• Sum of square of differences between the predicted and actual value of  $y$  denoted by SSE

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

• Mean of the sum of squares of the difference between the predicted and actual values of  $y$  denoted by MSE

$$MSE = \frac{1}{n} \cdot \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

## Back Propagation -

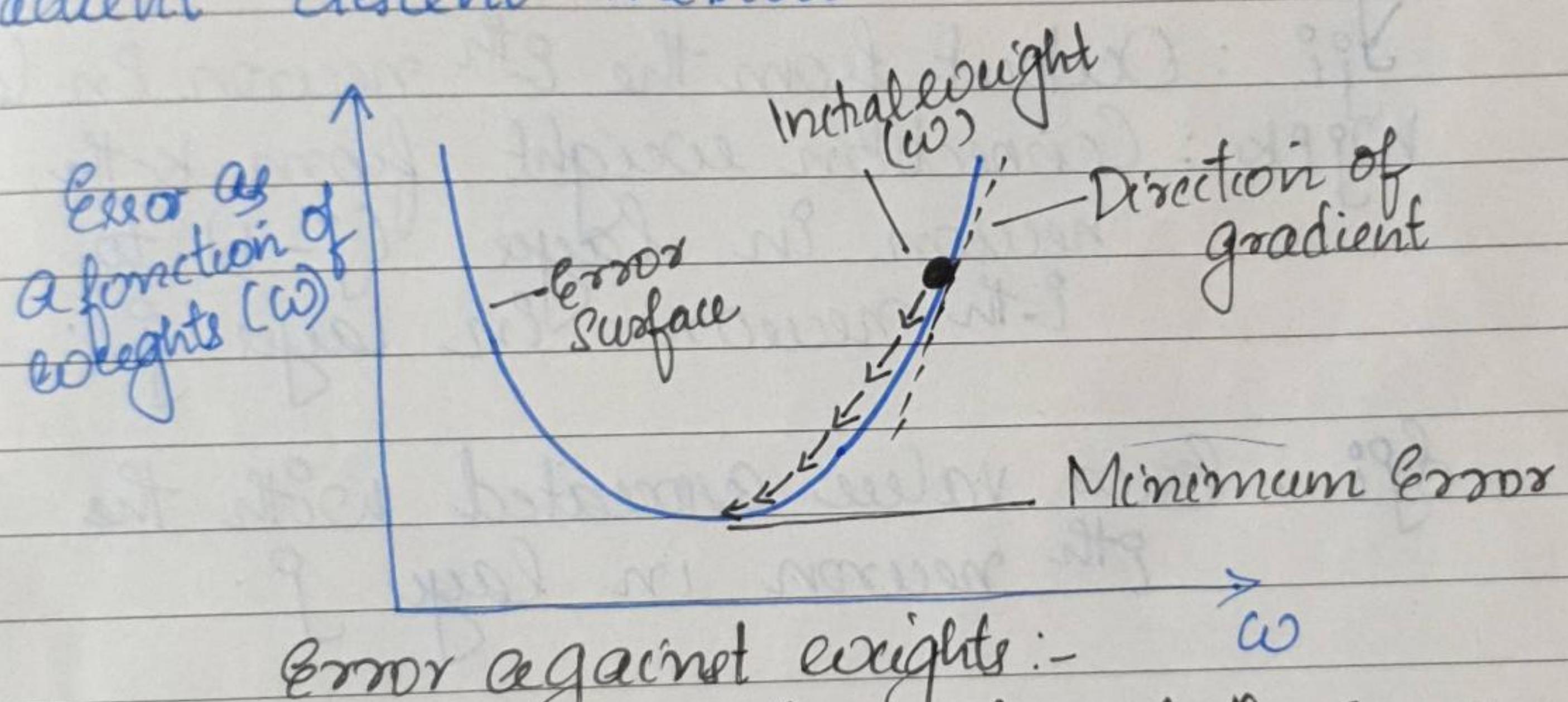
- \* Backpropagation is an algorithm for supervised learning of artificial neural networks using gradient descent.
  - \* Given an artificial neural network and an error function, the method calculates the gradient of the error function with respect to the neural network's weights.
- 1) Initially the weights are assigned random
  - 2) The algorithm iterates through many cycles of two processes called
    - (a) forward phase
    - (b) Backward phase
- (a) A forward phase:-
- \* Neurons are activated in sequence from the input layer to the output layers

- \* Apply neurons weights and activation function along the way. An output signal is produced from the final layer.

### (b) A backward phase:-

- \* Output signal compared with the true target value
- \* Difference in results ie ~~zero~~ error is propagated backwards in the network to modify the connection weights

3) The technique used to determine how much a weight can be changed is known as gradient descent method



Error against weights :-

- \* At any point on the surface of the curve, the gradient suggests how steeply the error will be reduced or increased for a change in the weight
- \* Algorithm helps to change the weights that result in the greatest reduction in error.

Algorithm:-

M :> No: of layers (excluding the input layer which is assigned the layer no: 0)

N<sub>j</sub> : No: of neurons in the j<sup>th</sup> layer.

X<sub>p</sub> = (X<sub>p1</sub>, X<sub>p2</sub>, ..., X<sub>pN<sub>0</sub></sub>) : p-th training sample  
T<sub>p</sub> = (T<sub>p1</sub>, T<sub>p2</sub>, ..., T<sub>pN<sub>M</sub></sub>) : Known output corresponding to p<sup>th</sup> training sample

O<sub>p</sub> = (O<sub>p1</sub>, O<sub>p2</sub>, ..., O<sub>pN<sub>M</sub></sub>) : Actual output to the p<sup>th</sup> training sample

y<sub>ji</sub> : Output from the i<sup>th</sup> neuron in layer j  
w<sub>jk</sub>: Connection weight from k-th neuron in layer (j-1) to i-th neuron in layer j.

e<sub>pi</sub> : Error value associated with the i<sup>th</sup> neuron in layer p.

Step 1 : Initialize connection weights into small random variables

Step 2: Present the p<sup>th</sup> sample input vector of pattern

$$X_p = (X_{p1}, X_{p2}, \dots, X_{pN_0})$$

and the corresponding op target

$$T_p = (T_{p1}, T_{p2}, \dots, T_{pN_M})$$

Step 3 : Pass the input values to the first layer i.e.  
For every input node 'i' in layer 0;  
 $y_{0i} = x_{pi}$

Step 4 : For every neuron 'i' in every layer  
 $j=1, 2, \dots, M$ . find the output

$$y_{ji} = f\left(\sum_{k=1}^{N_{j-1}} (y_{(j-1)k} w_{jk})\right)$$

## Algorithm:-

Step 1 : Initialize connection weights into small random variables.

Step 2 : Feed forward network

$$O_{ij} = \sum w_{ij} O_i + \theta$$

$\theta$  = Bias

• Apply Sigmoid Activation Function

$$O_j = \frac{1}{1 + e^{-\sum E_{O_1} + E_{O_2} - O_{ij}}}$$

Step 3 : Back Propagated network

→ Error calculation of output layer

$$E_O = O_j(1-O_j)(T_j - O_j)$$

$O_j$  = Output value

$T_j$  = Target value

→ Error Calculation of hidden layers

$$E_H = O_j(1-O_j) \sum_k E_k \cdot w_{pk}$$

Step 4 : Weights are updated by the equation

$$\Delta w_{ip} = \text{Learning rate} * E_{out} * O_{ij}$$

$$w_{ip} = w_{ip} + \Delta w_{ip}$$

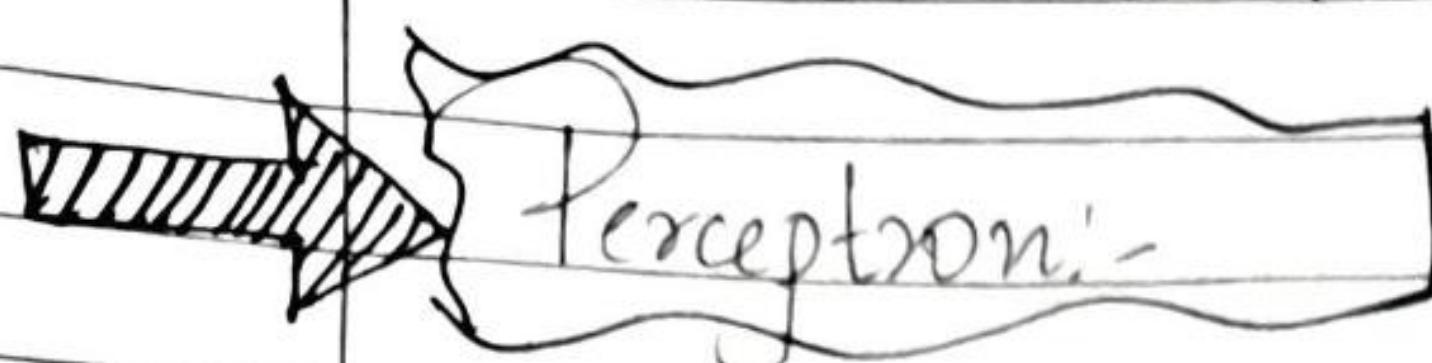
Biases are updated by

$$\Delta \theta_j = L * E$$

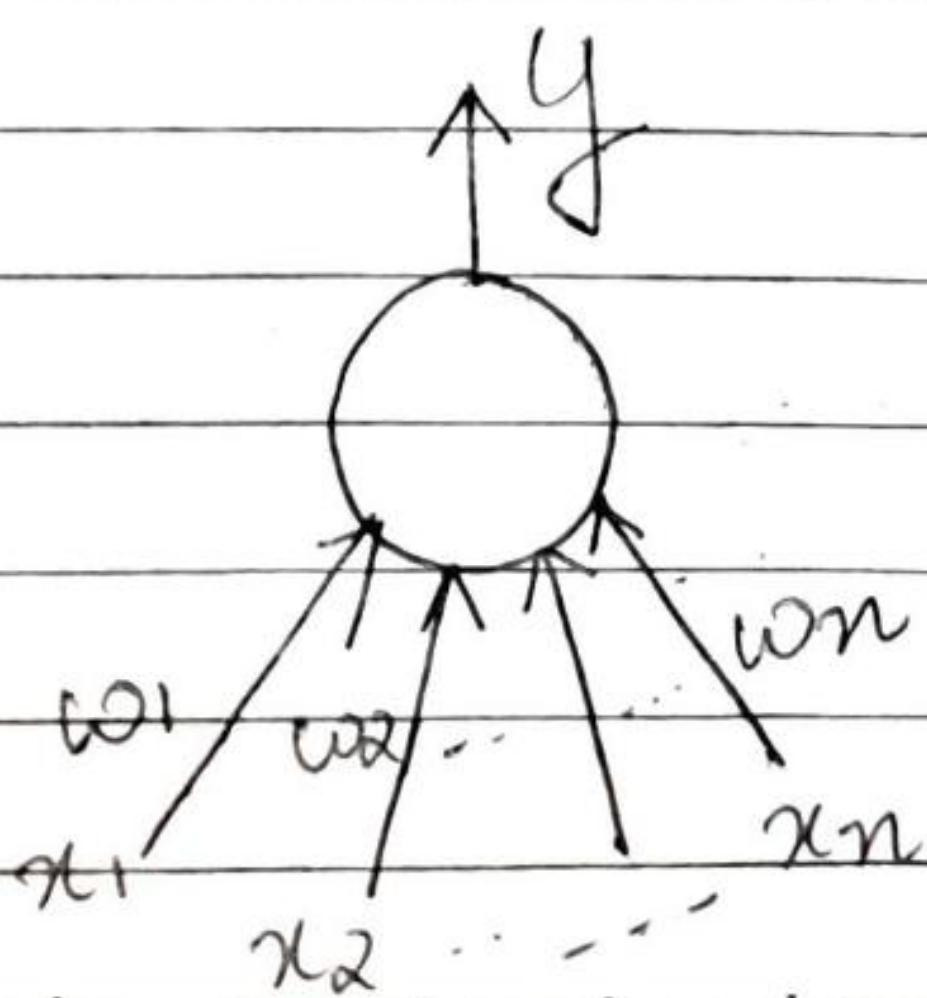
$$\theta_j = \theta_j + \Delta \theta_j$$

Step 3 and 4 are repeated until stopping  
Criteria meets.

## MODULE 1.



- \* McHilloch Pitts gives solution for boolean inputs, also hard code the threshold
- \* If it is non-boolean inputs, we need a way to learn our problem.
- \* All inputs are not equal means (weight) priority will be given for each input and based on that our decision lies on
- \* Perceptron will perform the decision making based on some learning algorithm and weights associated with each input
- \* Frank Rosenblatt, an American psychologist, proposed the classical perceptron model (1958)



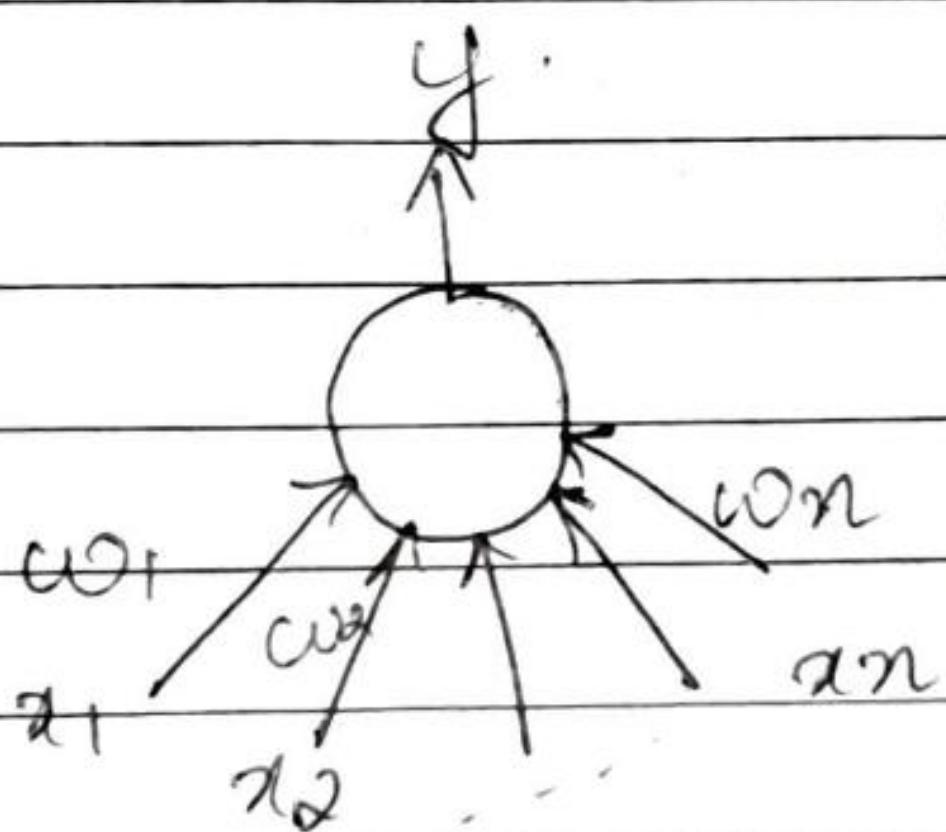
\* A more general computational model than McCulloch-Pitts Neurons.

\* Main Difference :

Introduction of numerical weights for inputs and a mechanism for learning these weights.

\* Inputs are no longer limited to Boolean values.

\* Refined and carefully analyzed by Minsky and Papert (1969) - their model is referred to as the Perceptron model.



$$\textcircled{1} \quad y = 1 \text{ if } \sum_{i=1}^n w_i^o * x_i^o \geq 0 \\ = 0 \text{ if } \sum_{i=1}^n w_i^o * x_i^o < 0$$

\textcircled{2} Rewriting the above,

$$y = 1 \text{ if } \sum_{i=1}^n w_i^o * x_i^o - \theta \geq 0$$

$$= 0 \text{ if } \sum_{i=1}^n w_i^o * x_i^o - \theta < 0$$

\textcircled{3}

A more accepted convention,

$$y = 1 \text{ if } \sum_{i=0}^n w_i^o * x_i^o \geq 0$$

$$= 0 \text{ if } \sum_{i=0}^n w_i^o * x_i^o < 0$$

$$\hookrightarrow \text{where } w_0^o = \theta \text{ and } x_0^o = 1$$

Eg:

Consider the task of predicting whether we would like a movie or not?

- \* Suppose, we base our decision based on 3 ips (binary)
- \* Based on our past viewing experience (data), we may give a high weight to Director Nolan as compared to other inputs.
- \* Specifically, even if the actor is not Matt Damon and the genre is not thriller we would still want to cross the threshold '0' by assigning a high weight to Director Nolan.
- \* '0' is called the bias as it represents the prior.
- \* A movie buff may have a very low threshold and may watch any movie irrespective of the genre, actors / director [0, 0]
- \* On the other hand, a selective viewer may only watch thrillers starring

Matt Damon and directed by Nolan  
[θ = 3]

\* The weights ( $w_1, w_2, \dots, w_n$ ) and the bias ( $w_0$ ) will depend on the data.

McCulloch-Pitts Neuron  
(assuming no inhibitory inputs)

$$y = 1 \text{ if } \sum_{i=0}^n w_i x_i \geq 0$$

$$= 0 \text{ if } \sum_{i=0}^n w_i x_i < 0.$$

Perceptron

$$y = 1 \text{ if } \sum_{i=0}^n w_i x_i \geq 0.$$

$$= 0 \text{ if } \sum_{i=0}^n w_i x_i < 0.$$

Continuee perceptron book an. fns.)

# Algorithm: Perceptron Learning Algorithm:-

P  $\leftarrow$  inputs with label 1

N  $\leftarrow$  inputs with label 0

Initialize  $w = [w_0, w_1, w_2, \dots, w_n]$  randomly;

while !Convergence do (until not making any more errors on the training data)

Pick a random  $x \in P \cup N$ ;

if  $x \in P$  and  $\sum_{i=0}^n w_i * x_i < 0$  then

$w = w + x$ ;

end

if  $x \in N$  and  $\sum_{i=0}^n w_i * x_i > 0$  then

$w = w - x$

end

end

↳ Consider 2 vectors ' $w$ ' and ' $x$ '

$$w = [w_0, w_1, \dots, w_n]$$

$$x = [x_0, x_1, \dots, x_n]$$

$$w \cdot x = w^T x = \sum_{i=0}^n w_i * x_i$$

↳ Rewriting the perceptron rule as

$$y = 1, \text{ if } w^T x \geq 0$$

$$= 0, \text{ if } w^T x < 0$$

\* Finding the line  $w^T x = 0$ , which divides the input space into two halves.

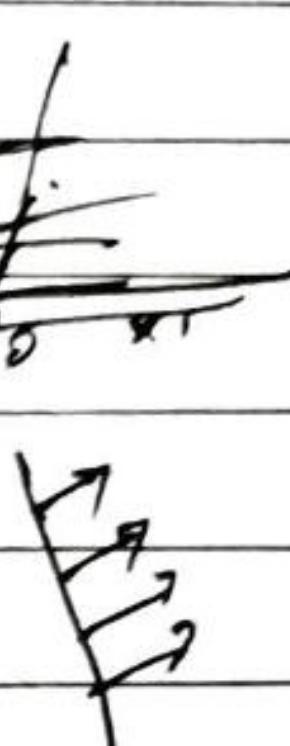
i) Every point ' $x$ ' on this line satisfies the equation  $w^T x = 0$

$$\text{eg: } x_1 + x_2 = 0$$

$$\text{points } (1, -1), (2, -2)$$

$w^T x$  is dot product ( $w \cdot x = 0$ ),

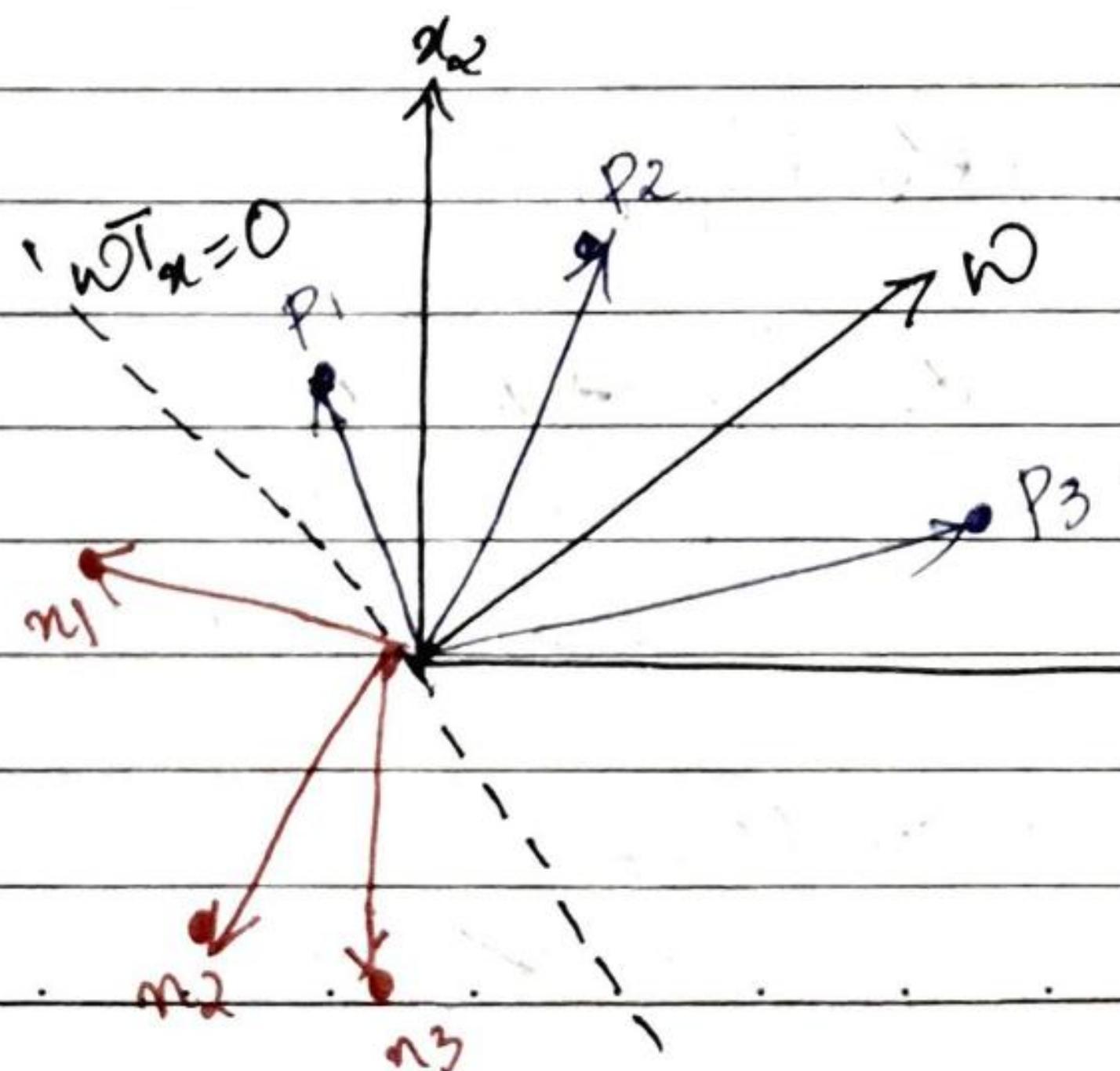
then they are orthogonal means every points lie on the line & perpendicular to the weight



[Dot product of 2 vectors is equal to the product of the magnitude of the 2 vectors and the cosecant of the angle b/w the two vectors]

i.e. The angle b/w point on line and the 'w' is  $90^\circ$ . So

$$\cos \alpha = \frac{w^T x}{\|w\| \|x\|} = 0$$



[angle b/w  $w \cdot p_1, p_2, p_3$  and  $w$  is less than  $90^\circ$ ]

$\cos \alpha < 0$  (less than  $90^\circ$ ) to be on the +ve side.

$\alpha_1$  [angle b/w  $w \cdot n_1, n_2, n_3$  and  $w$  are greater than  $90^\circ$ ]

$\cos \alpha > 0$  ( $> 90^\circ$ ) to be on the -ve side.

$P \leftarrow$  inputs with label 1;  
 $N \leftarrow$  inputs with label 0;  
 Initialize  $w$  randomly;

$$\begin{aligned} 95 > 90 \\ \cos 95^\circ = -0.022 \end{aligned}$$

$$\cos 90^\circ = 0.13$$

$$\cos 95^\circ < \cos 90^\circ$$

$$\text{as } 95^\circ > 90^\circ$$

while ! Convergence do

Pick Random  $x \in P \cup N$ ;

if  $x \in P$  and  $w \cdot x < 0$  then

$$w = w + x;$$

end

if  $x \in N$  and  $w \cdot x \geq 0$  then

$$w = w - x;$$

end

|| the algorithm converges when all the inputs are classified correctly.

\* For  $x \in P$  if  $w \cdot x < 0$ , then if mean the angle ( $\alpha$ ) b/w  $w$  and  $x$  is greater than  $90^\circ$  [bcz we want  $\alpha$  to be less than  $90^\circ$ ]

\* Compute  $w_{\text{new}} = w + x$ .

$$\cos(\alpha_{\text{new}}) \propto w^T x$$

$$\propto (w + x)^T x$$

$$\propto w^T x + x^T x$$

$$\propto \cos \alpha + x^T x$$

$$\cos(\alpha_{\text{new}}) > \cos \alpha$$

$$\text{i.e. } (\alpha_{\text{new}}) \text{ i.e. } < 90^\circ$$

\* For  $x \in N$ , if  $w \cdot x \geq 0$  then it means that the angle ( $\alpha$ ) b/w the  $x$  and the current  $w$  is less than  $90^\circ$ . (But we want  $\alpha'$  to be greater than  $90^\circ$ )

$$\begin{aligned}\cos(\alpha_{\text{new}}) &\propto w_{\text{new}}^T x \\ &\propto (w - x)^T x \\ &\propto w^T x - x^T x + \text{value} \\ &\propto \cos\alpha - x^T x\end{aligned}$$

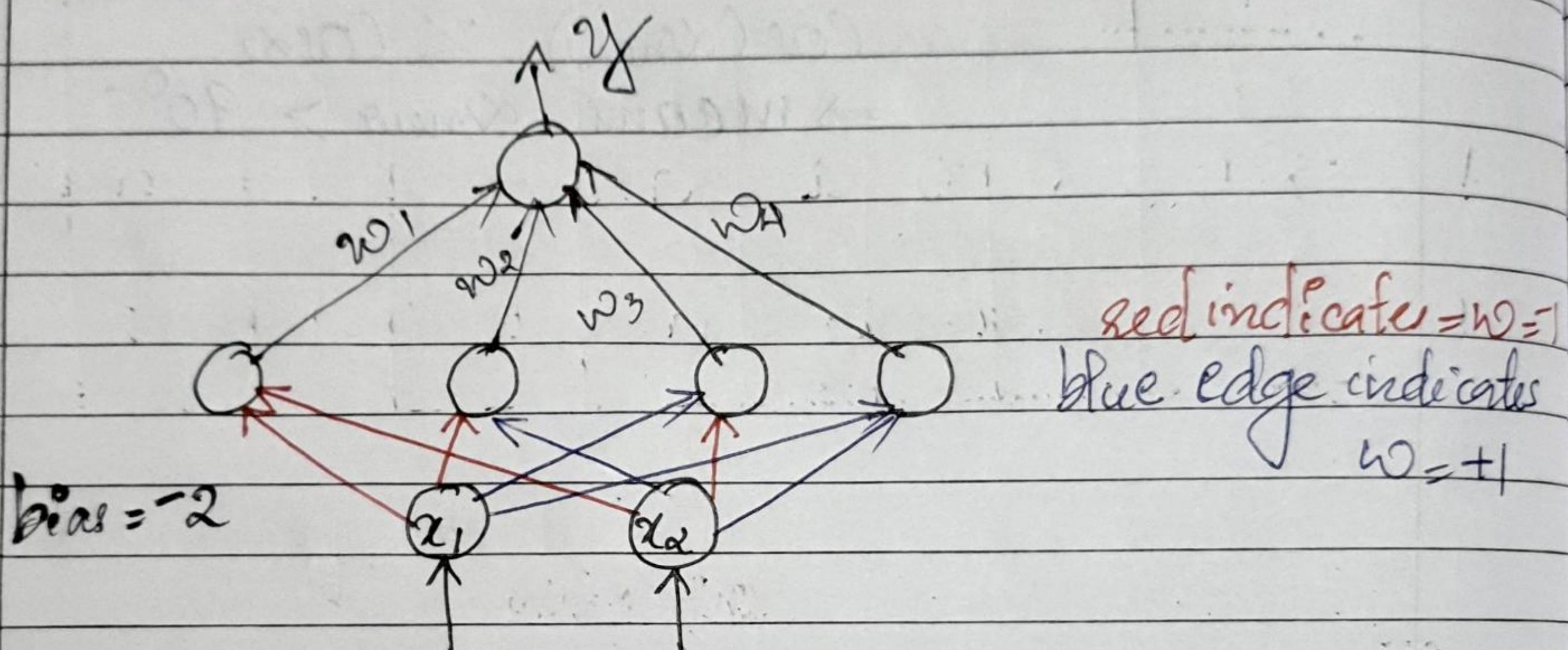
$$\begin{aligned}\cos(\alpha_{\text{new}}) &< \cos\alpha \\ \Rightarrow \text{means } \alpha_{\text{new}} &> 90^\circ.\end{aligned}$$

# Network of Perceptrons :-

OR

## Multilayer Perceptrons:-

\*



- \* Consider 2 inputs and 4 Perceptrons
- \* Each input is connected to all the 4 perceptrons with specific weights
- \* The bias ( $w_0$ ) of each perceptron is -2  
(ie each perceptron will fire only if the weighted sum of its inputs is  $\geq -2$ )
- \* Each of these perceptrons is connected to an output perceptron by weights (which need to be learned)

## Terminology :-

Input layer  
hidden layer  
Op layer

- \* Ops of 4 perceptrons in the hidden layer are denoted by  $h_1, h_2, h_3, h_4$
- \* The red and blue edges are called layer 1 weights
- \*  $w_1, w_2, w_3, w_4$  are called layer 2 weights
- \* This net can be used to implement any boolean functions (linearly separable or not)
- \* Perceptron 1 will fire when the input is  $(-1, -1)$
- \* Perceptron 2 "  $\rightarrow (-1, 1)$
- \* Perceptron 3 "  $\rightarrow (1, -1)$
- \* Perceptron 4 "  $\rightarrow (1, 1)$

eg XOR

Let  $w_0$  be the bias o/p of the neuron  
(i.e., it will fire if  
 $\sum_{i=1}^4 w_i h_i \geq w_0$ )

$x_1$	$x_2$	XOR	$h_1$	$h_2$	$h_3$	$h_4$	$\sum_{i=1}^4 w_i h_i \geq w_0$
0	0	0	1	0	0	0	$w_1$
0	1	1	0	1	0	0	$w_2$
1	0	1	0	0	1	0	$w_3$
1	1	0	0	0	0	1	$w_4$

\* To implement XOR:

$$w_1 < w_0$$

$$w_2 \geq w_0$$

$$w_3 \geq w_0$$

$$w_4 < w_0.$$

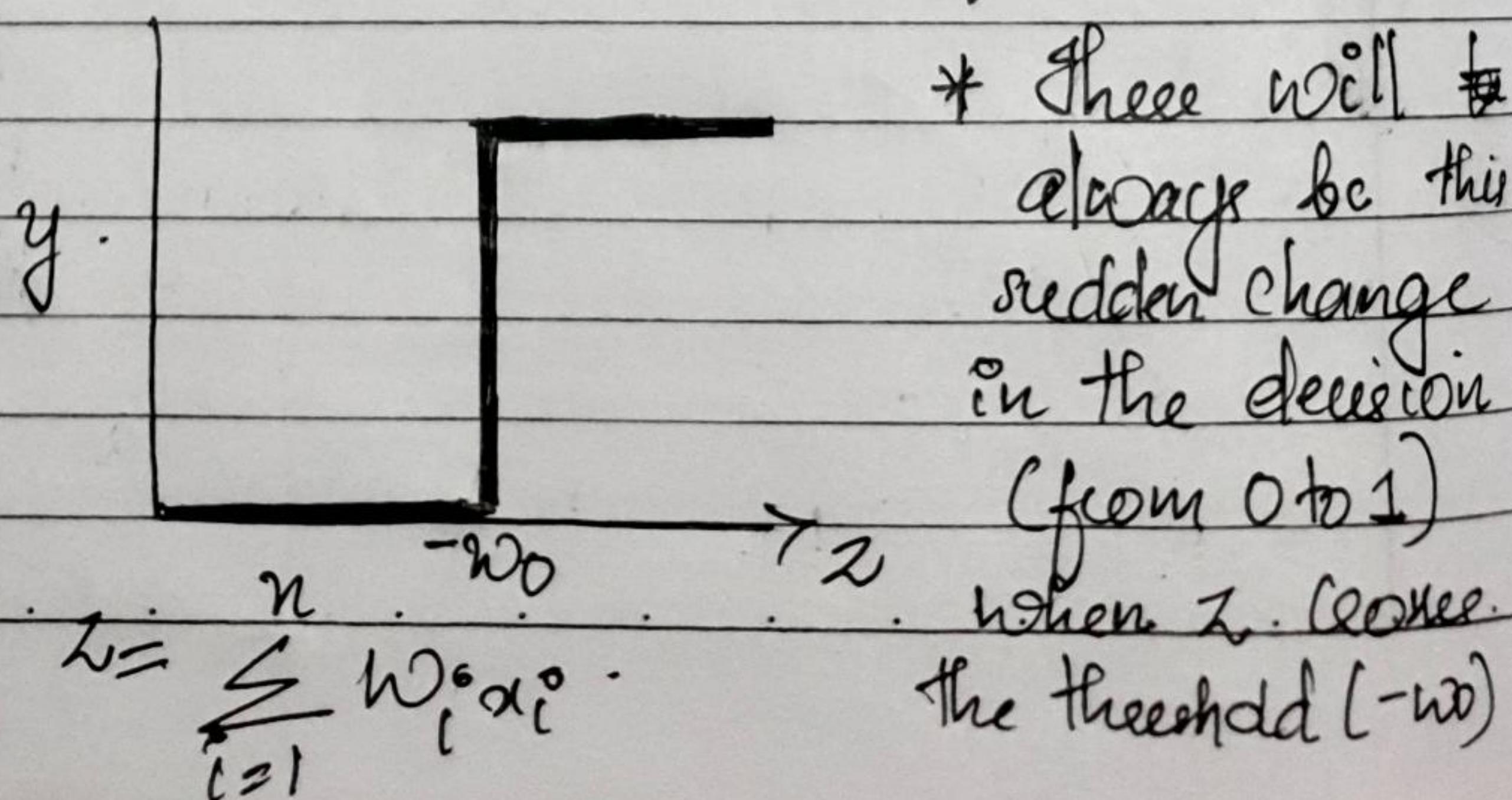
\* Each ' $w_i$ ' is now responsible for one of the 4 possible inputs, and can be adjusted to get the desired o/p for that i/p.

- \* If 3<sup>ops</sup>:  $b_{in} = 3$  and hidden layer contains  $2^3$  neurons i.e (8)
- \* If n<sup>ops</sup> hL will  $2^n$  neurons
- \* Any boolean fn of 'n' inputs can be represented exactly by a nn of perceptions containing 1 hidden layer with  $2^n$  perceptions and one op layer containing 1 perception
- \* As 'n' increases the no: of perceptions in the hidden layers obviously increases exponentially.

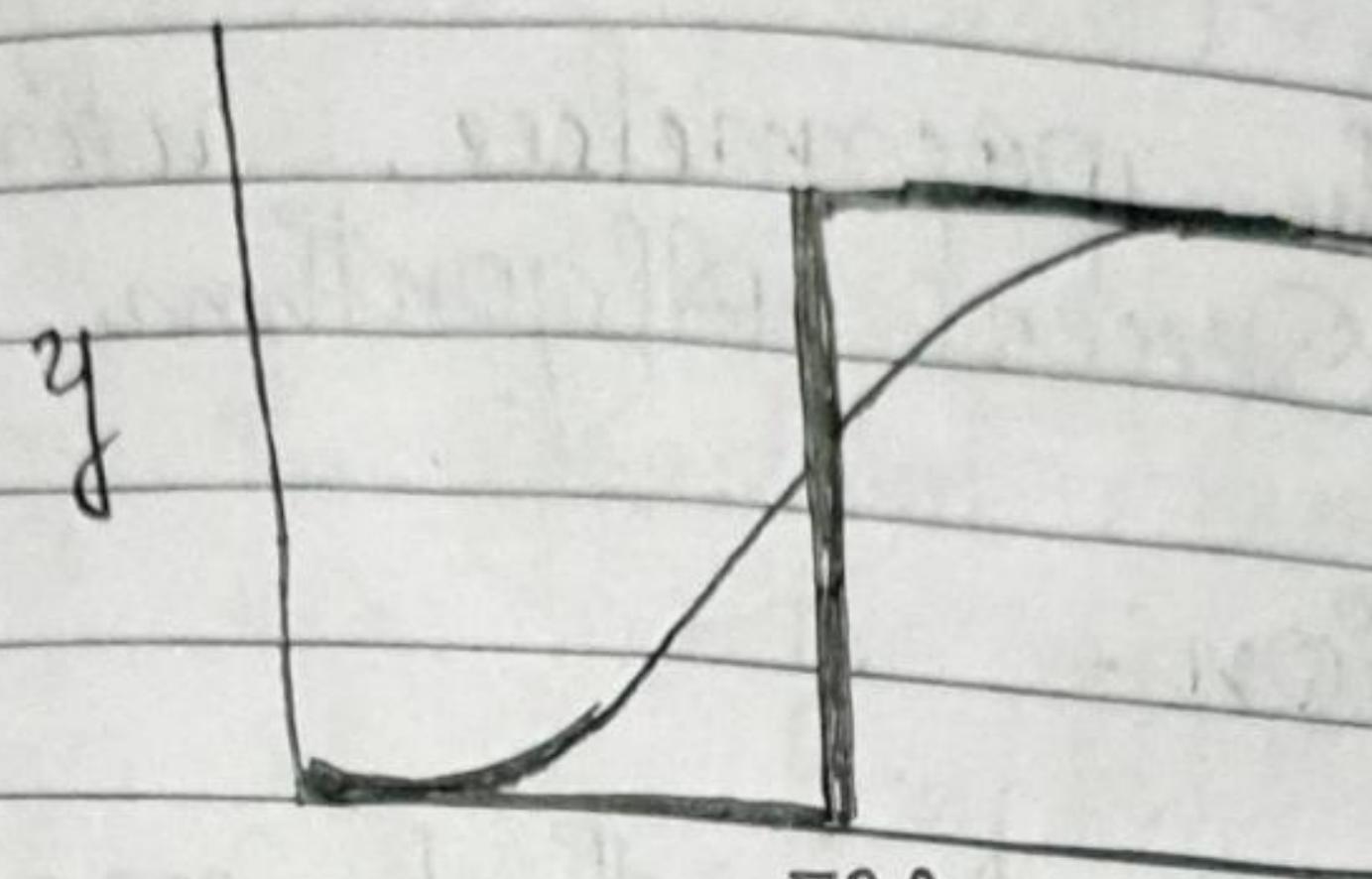
## Sigmoid Neuron:-

- \* A sigmoid neuron is similar to the perceptron neuron, for every input ' $x_i$ ' it has weight ' $w_i$ ' associated with the input.
- \* The weights indicate the importance of the input in the decision making process.
- \* The output from the sigmoid is not 0 or 1 like the perceptron model instead it is a real value b/w 0-1 which can be interpreted as a probability.
- \* The most commonly used sigmoid fn is the logistic function, which has the characteristics of an 'S' shaped curve.
- \* It is an abstract function of the form  $y = f(x)$  where  $x \in \mathbb{R}^n$  and  $y \in \mathbb{R}$ .

↳ Perceptron function itself which behaves like a step function



\* For most real world application we would expect a smoother decision fn which gradually changes from 0 to 1



One form of the sigmoid fn called the logistic fn.

$$y = \frac{1}{1 + e^{-(w_0 + \sum_{i=1}^n w_i x_i)}}$$

when  $w^T x = x_1$ ,

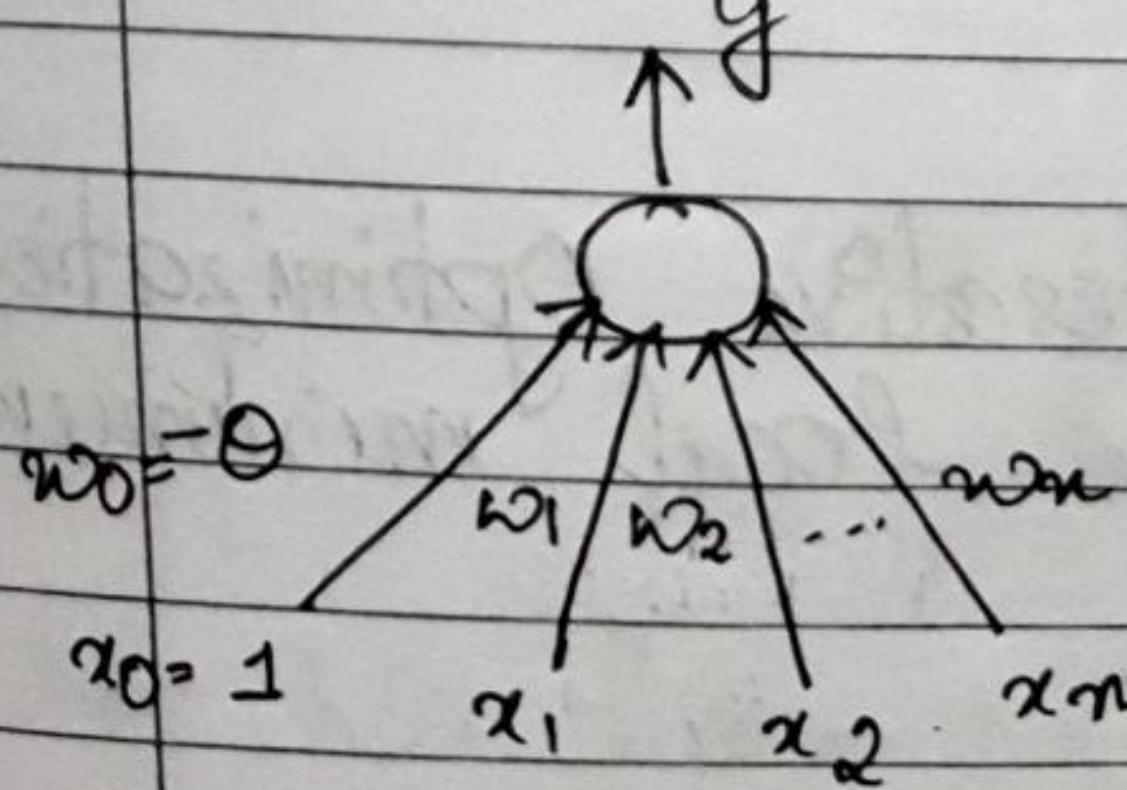
$$y = 1$$

when  $w^T x = -x_1$

$$y = 0$$

## Perception

\* Not smooth, not continuous and not differentiable

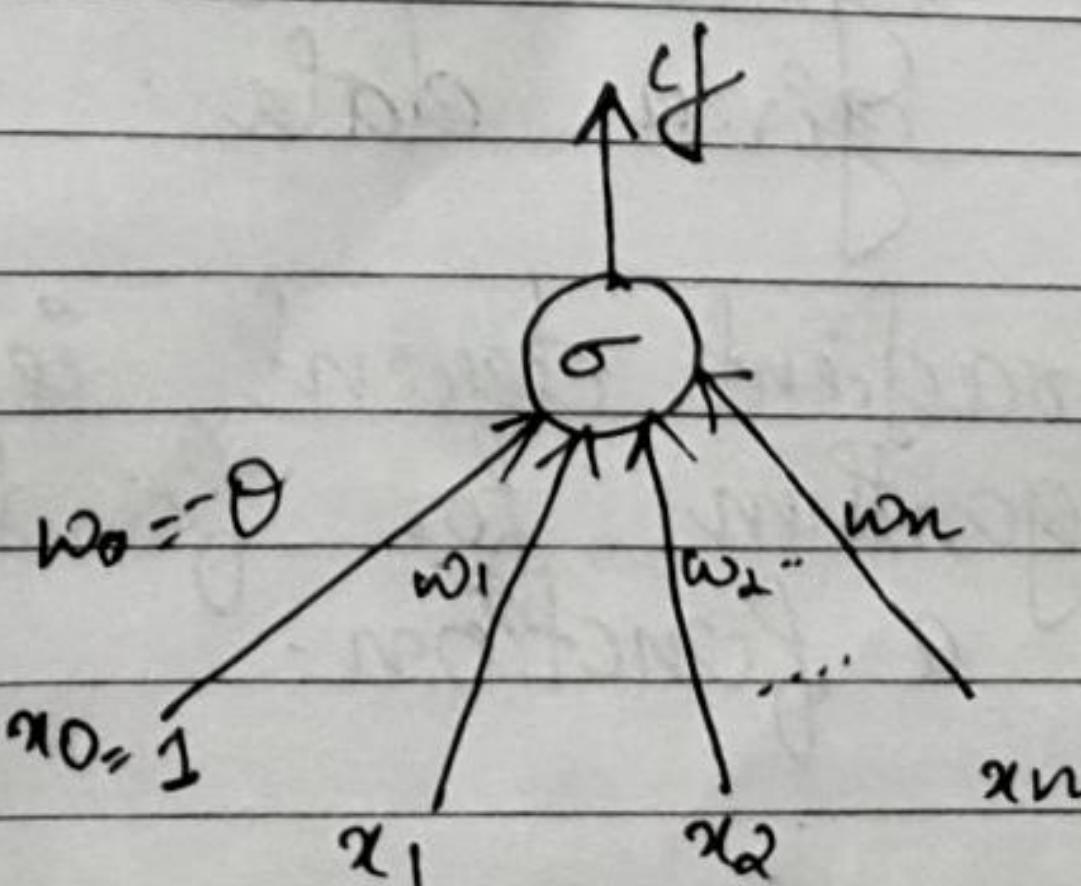


$$y = 1 \text{ if } \sum_{i=0}^n w_i * x_i \geq 0$$

$$= 0 \text{ if } \sum_{i=0}^n w_i * x_i < 0$$

## Sigmoid (Logistic) Neuron

\* Smooth, Continuous and differentiable



$$y = \frac{1}{1 + e^{-\sum_{i=0}^n w_i x_i}}$$

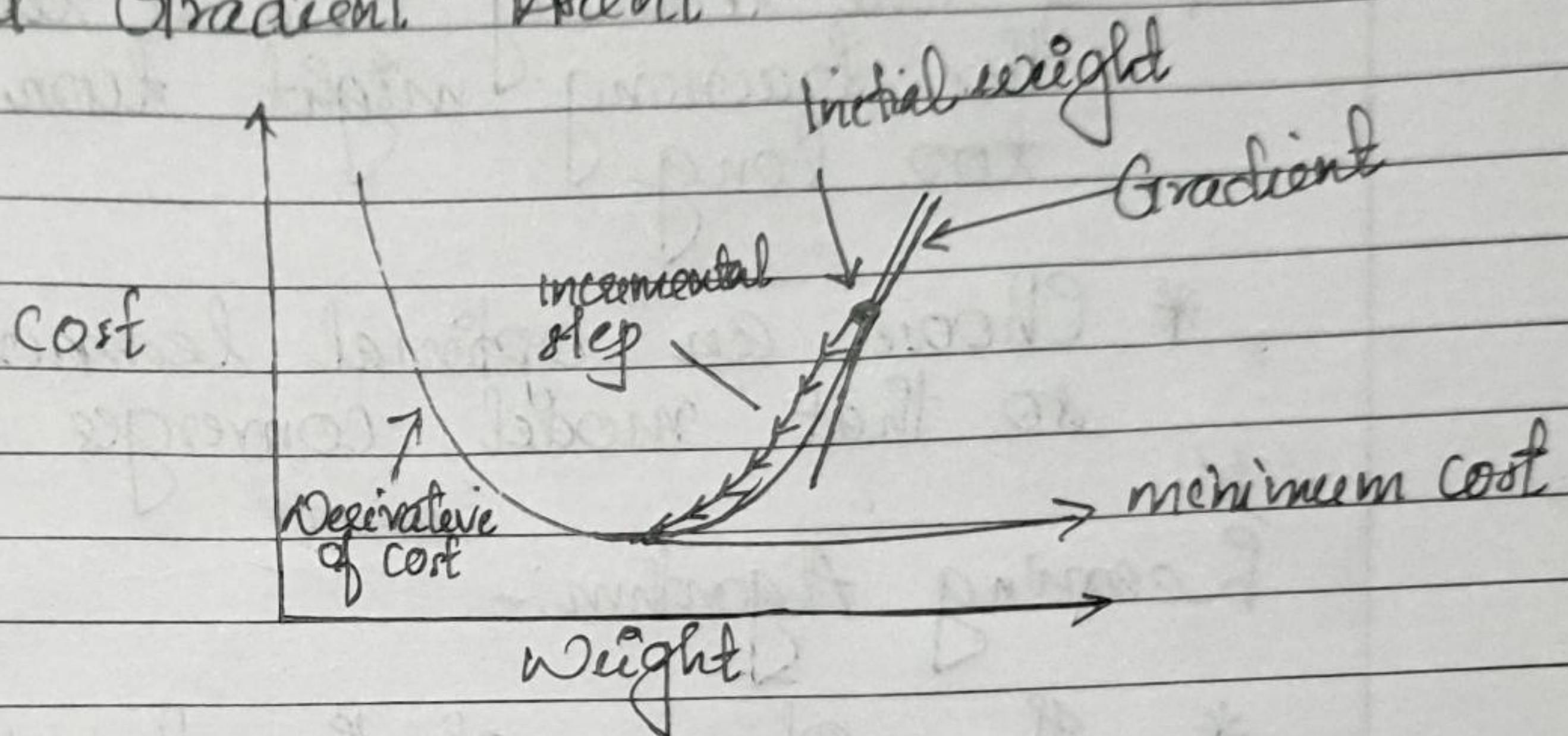
## Gradient Descent Algorithm

- \* We had an algorithm for learning the weights of a perceptron, we also need a way of learning the weights of a sigmoid neuron.
- \* To learn the parameters, using Gradient Descent Algorithm.

### Cost Function:

- It is a function that measures the performance of a model for every given data.
  - Cost function quantifies the error b/w predicted values and expected values and presents it in the form of a single real number.
  - After making a hypothesis with initial parameters, we calculate the cost fn.
  - Goal is to reduce the cost fn; for that modify the parameters by using the Gradient Descent Algorithm over the given data.
- \* Gradient descent is an iterative Optimization Algorithm for finding the local minimum of a function.

- \* To find the local minimum of a fn using gradient descent, takes steps proportional to the negative of the gradient (move away from the gradient) of the fn at the current point
  - If take steps proportional to the tve of the gradient (moving towards the gradient), we will approach a local maximum of the function. This procedure is called Gradient Ascent.



- \* The goal of the gradient descent algorithm is to minimize the given function
  - (1) Compute the gradient (slope), the first order derivative of the fn at that point
  - (2) Make a step in the direction opposite to the gradient

- \* We have now the direction to move in, next is to decide the size of the step

Learning rate ( $\eta$ ) - a tuning parameter  
in the optimization process  
- It decides the length of the steps.

- \* If the learning rate is too high,  
it might overshoot the minimum and  
keep bouncing without reaching the  
minimum
- \* If the learning rate is too small,  
the training might turn out to be  
too long
- \* Choose an optimal learning rate,  
so that model converges to the minimum.

### Learning Algorithm:-

\* The objective of the learning algorithm  
is to determine the best possible values  
for the parameters ( $w$  and  $b$ ), such  
that the overall loss of the model is  
minimized as much as possible

Initialize  $w, b$   
Iterate over data :

repeated  
until converge

Compute  $y$   
Compute  $L(w, b)$

$$w_{t+1} = w_t - \eta \Delta w_t$$
$$b_{t+1} = b_t - \eta \Delta b_t$$

, where  $\Delta w_t = \frac{\partial L(w, b)}{\partial w}$

at  $w=w_t, b=b_t$

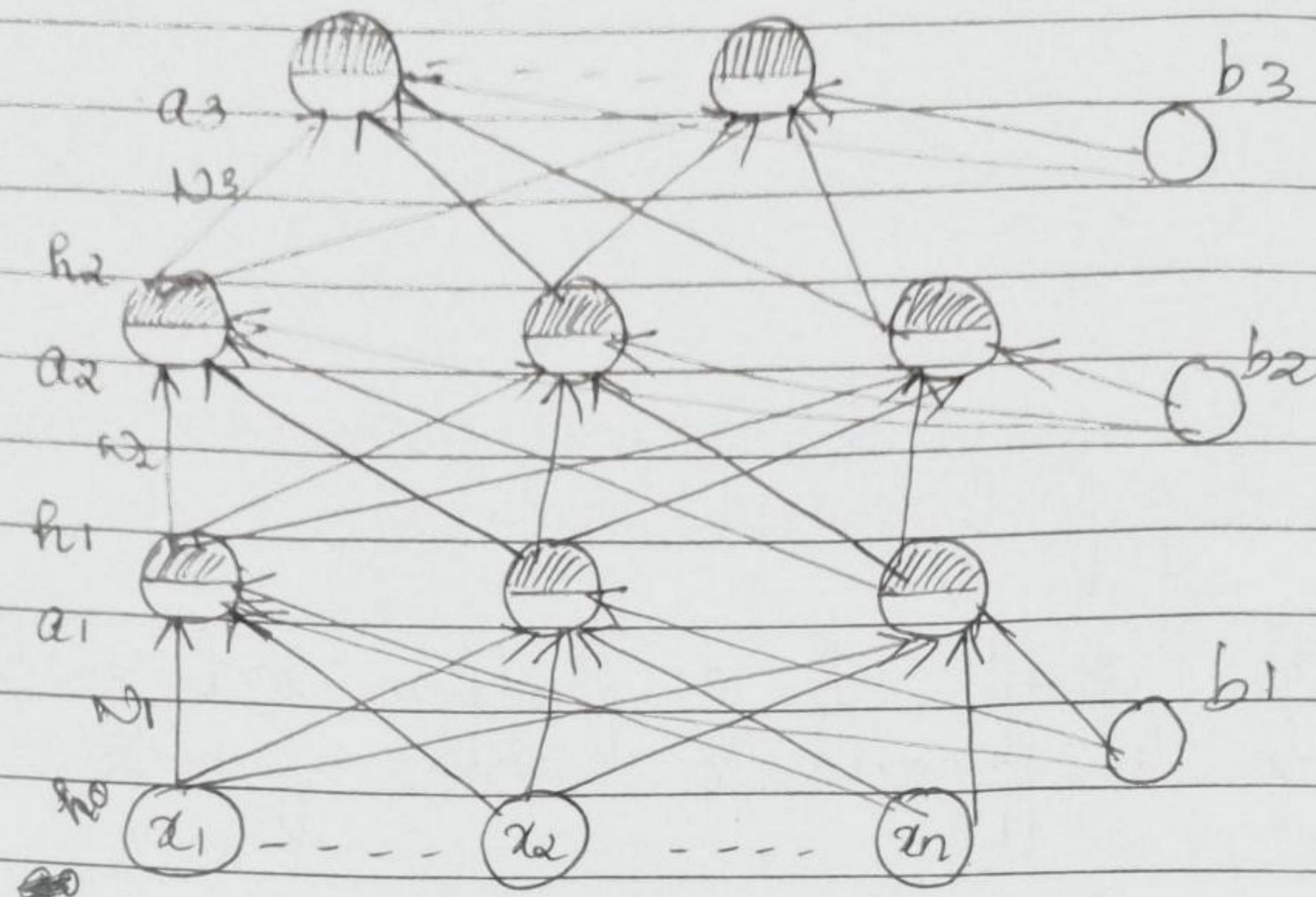
$\Delta b_t = \frac{\partial L(w, b)}{\partial b}$  at

$w=w_t, b=b_t$

## Feed Forward Neural Network:-

- \* Feed forward neural networks are also called deep feedforward networks.
- \* The goal of a feedforward network is to approximate some function  $f^*$ .  
eg:- for a classifier,  $y = f^*(x)$  maps an input ' $x$ ' to a category ' $y$ '.
- \* A feedforward net defines a mapping  $y = f(x, \theta)$   
 $\uparrow$   
 $\rightarrow \theta$  is the parameter to learn that results in the best function approximation.
- \* Such models are called feedforward bcz information flows through the function being evaluated from ' $x$ ', through the intermediate computations used to define ' $f$ ' and finally to the op ' $y$ '.

$$h_L = \hat{y} = f(x)$$



- \* The input to the n/w is a  $n$ -dimensional vector
- \* The n/w contains ' $L-1$ ' hidden layers ( $2$  in the above case) having ' $n$ ' neurons each
- \* Finally, there is one output layer containing ' $k$ ' neurons (corresponding to ' $k$ ' classes)
- \* Each neuron in the hidden layers and

and output layer can be split into two parts:

- Pre activation and
- Activation

\* Each neuron in the hidden / input layer can be called the ' $i^{th}$ ' layer and op layer can be called the ' $L^{th}$ ' layer

\*  $W_i \in \mathbb{R}^{n \times n}$  and  $b_i \in \mathbb{R}^n$  are the weight and bias b/w layers  $i-1$  and  $i$  ( $0 \leq i \leq L$ )

\*  $W_L \in \mathbb{R}^{n \times k}$  and  $b_L \in \mathbb{R}^k$  are the weight and bias b/w the last hidden layer and the op layer

\* The pre-activation at layer ' $i$ ' is given by

$$a_i(x) = b_i + W_i h_{i-1}(x)$$

e.g.:  $a_1 = b_1 + W_1 h_0$

$$\begin{bmatrix} a_{11} \\ a_{12} \\ a_{13} \end{bmatrix} = \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} + \begin{bmatrix} w_{111} & w_{112} & w_{113} \\ w_{121} & w_{122} & w_{123} \\ w_{131} & w_{132} & w_{133} \end{bmatrix} \begin{bmatrix} h_{01} = x_1 \\ h_{02} = x_2 \\ h_{03} = x_3 \end{bmatrix}$$

$$= \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} + \begin{bmatrix} w_{111}x_1 + w_{112}x_2 + w_{113}x_3 \\ w_{121}x_1 + w_{122}x_2 + w_{123}x_3 \\ w_{131}x_1 + w_{132}x_2 + w_{133}x_3 \end{bmatrix} = \begin{bmatrix} \sum w_{11i}x_i + b_{11} \\ \sum w_{12i}x_i + b_{12} \\ \sum w_{13i}x_i + b_{13} \end{bmatrix}$$

\* The activation at layer ' $i$ ' is given by  $h_i^o(x) = g(a_i^o(x))$

where 'g' is called the activation function (eg:- logistic, tanh, linear etc)

\* The o/p layer ' $i$ ' is given by

$$f(x) \cdot h_i^o(x) = O(a_i^o(x))$$

where 'O' is the o/p activation fn  
(eg:- softmax, linear etc)

Data:  $\{x_i^o, y_i^o\}_{i=1}^N$

Model:

$$y_i^o = f(x_i^o) = O(w_3^3 g(w_2^3 g(w_1^3 g(w_0^3 x + b_1) + b_2) + b_3))$$

Parameters:  $\theta = w_0, \dots, w_N, b_1, b_2, \dots, b_N$

Algorithm: Gradient Descent with Backpropagation

Objective loss / error fn: say,  $J(\theta)$

$$\min \frac{1}{N} \sum_{i=1}^N (y_i^o - y_i)^2$$

In general  $\min L(\theta)$ , where  $L(\theta)$  is some fn of parameters

# Output Functions and Loss Functions

	Output	
Op Activation	Real values	Probabilities
Loss Function	Squared Error	Cross Entropy

Computing gradients with respect to the  
Output Units :-

$$\frac{\partial L(O)}{\partial w_{111}} = \frac{\partial L(O)}{\partial y} * \frac{\partial y}{\partial a_3} \frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}$$

Talk to the output layer      Talk to the previous hidden layer

Talk to the weight directly       $\frac{\partial a_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_{111}}$

[ $y$  is the predicted op]

> Talk now to the weights

\* A computational graph where the nodes or variables correspond to operations and edges correspond to operations can feed their value into other operations and feed their value into other operations.

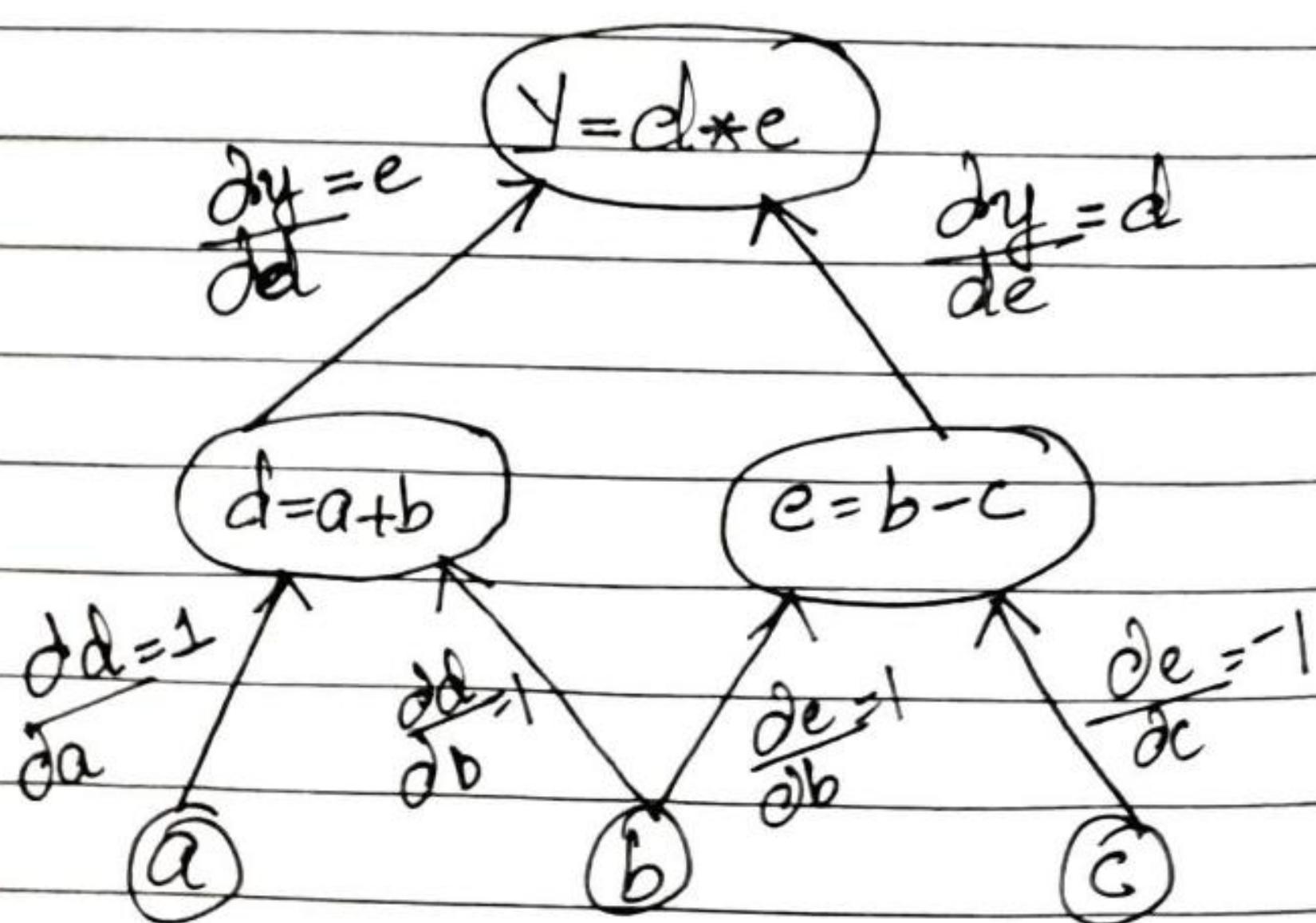
## Computational Graphs :-

\* Computational graph is a directed graph which is used for expressing and evaluating the mathematical expression.

\* Computation of the neural netw are organized in terms of a forward pass or forward propagation step in which we compute the O/P of the neural netw followed by a backward pass or backward propagation step, which we use to compute gradients / derivatives.

\* To understand derivatives in a computation graph, the key is to understand how a change in one variable brings change on the variable that depends on it.

eg:-



\* We can easily apply chain rule to evaluate

partial derivatives of final op variable w.r.t.  
input variable  $a, b, c$ .

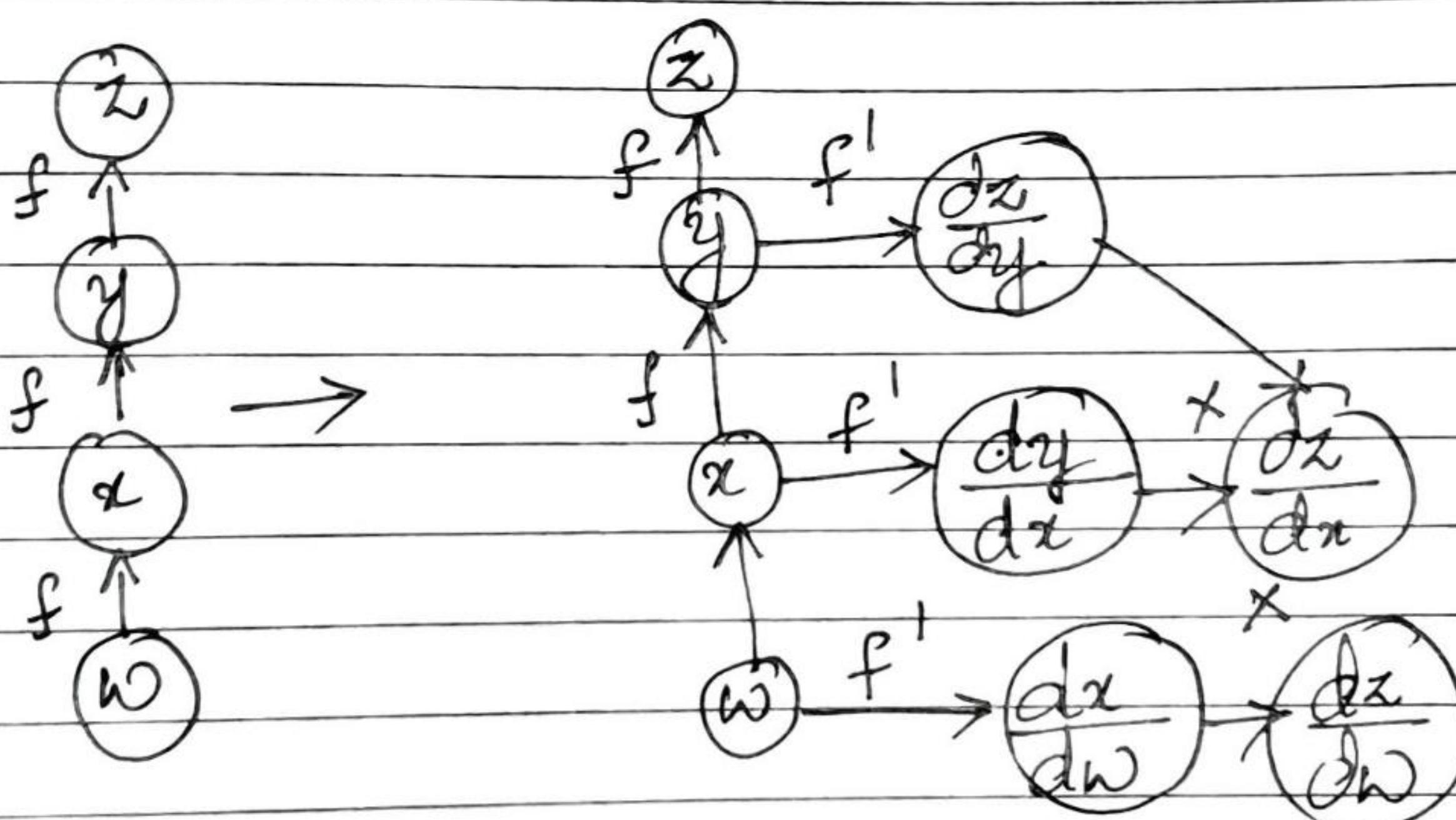
$$\frac{\partial y}{\partial a} = \frac{\partial y}{\partial d} \times \frac{\partial d}{\partial a} = ex_1 = e$$

$$\frac{\partial y}{\partial b} = \frac{\partial y}{\partial d} \times \frac{\partial d}{\partial b} = ex_1 - e$$

$$\frac{\partial y}{\partial c} = \frac{\partial y}{\partial e} \times \frac{\partial e}{\partial c} = dx - 1 = -d$$

\* These computational graphs make it easier  
to get the derivatives using backpropagation

eg:



~~$z = f(f(y))$~~

$$z = f(f(f(w)))$$