

Artificial Neural Network

ANN is based on a unit called
a Perceptron.

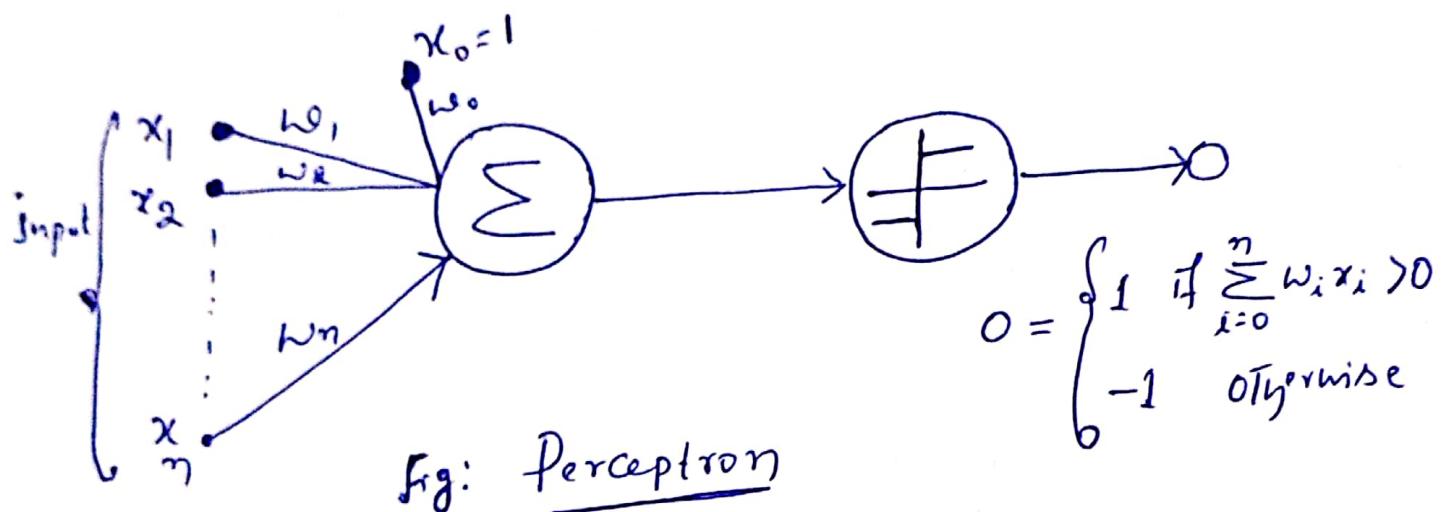


Fig: Perceptron

→ A perceptron takes a vector

of real-valued inputs

&

Calculates a linear Combination
of these inputs

Then outputs a 1 if the result is
greater than
some
threshold

$$O(x_1, x_2, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise} \end{cases}$$

Inhere

Each w_i is Weight T_i at

Determines the Contribution
of input x_i to the
perceptron output.

~~Weight~~

x_0 = Additional Constant Input

$$\underbrace{\downarrow}_{\text{bias}} = 1$$

$\underbrace{\text{bias}}$ Decides the position
of the Decision Surface
in the n -Dimensional
feature Space.

$$o(x_1, \dots, x_n) = \begin{cases} 1 & \sum_{i=0}^n w_i x_i > 0 \\ -1 & \text{otherwise.} \end{cases}$$

In Vector form

$$o(\vec{x}) = \begin{cases} 1 & \vec{w} \cdot \vec{x} > 0 \\ -1 & \text{otherwise} \end{cases}$$

Note: Learning a perceptron involves

Choosing values for the weight

$$w_0 = -\frac{w_1}{w_2}$$

Note:- If Data is linearly
Separable & Binary classification,
(2 class problem)

Then we can implement

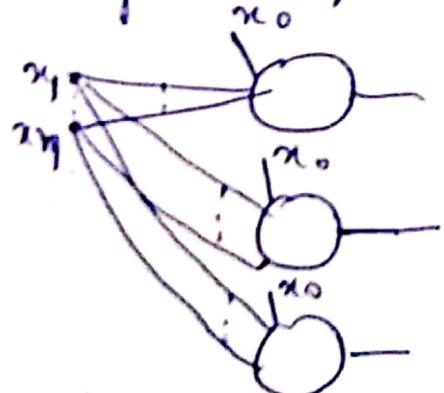
using single perceptron

&

If Data is Linearly Separable
& Multi-Class problem

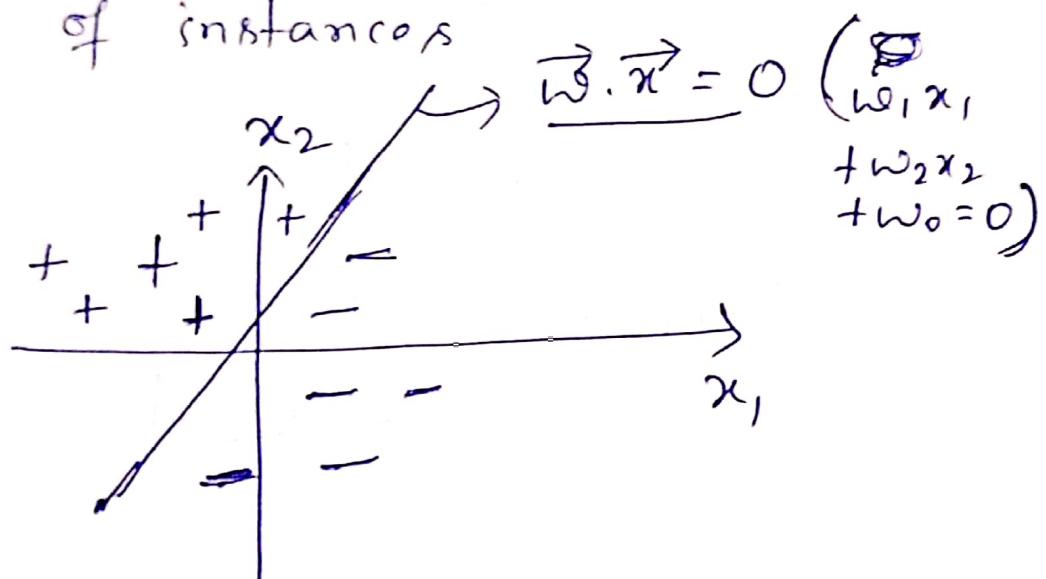
then we need a more
than one perceptron in

a single layer



Representation power of Perceptron

Perceptron represents a hyperplane Decision Surface in the n -Dimensional Space of instances



Perceptron outputs a 1 for the instances lying on one side of hyperplane.

$$\text{i.e } \vec{w} \cdot \vec{x} > 0$$

$$(w_1x_1 + w_2x_2 + w_0) > 0$$

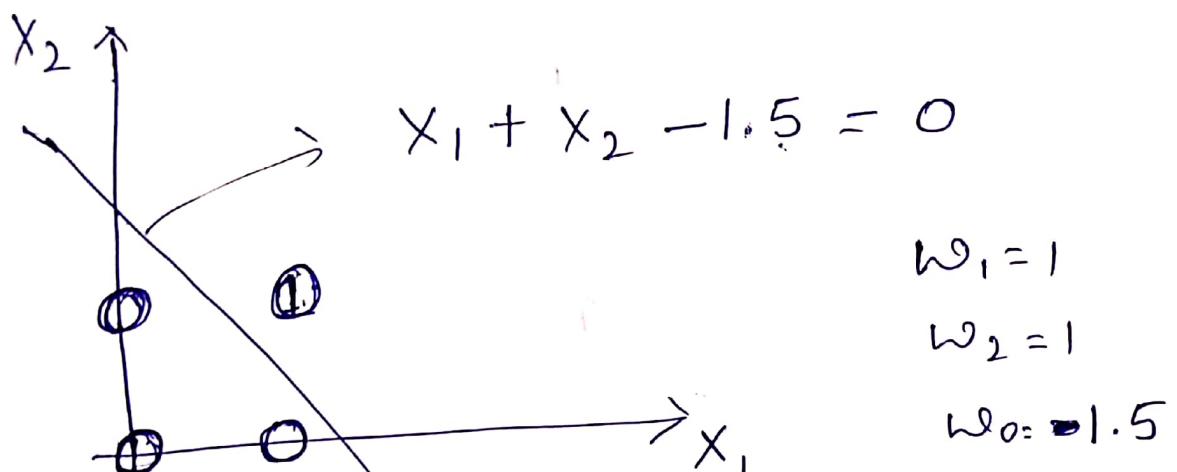
Perceptron outputs a -1 for the instances lying on other side of hyperplane

$$\text{i.e } \vec{w} \cdot \vec{x} < 0$$

$$(w_1x_1 + w_2x_2 + w_0) < 0$$

AND Function

x_1	x_2	y
0	0	0
0	1	0
1	0	0
1	1	1



$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

Feature Vector (\vec{x})

$$\begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

Weight Vector

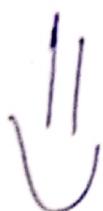
$$\begin{bmatrix} -1.5 \\ 1 \\ 1 \end{bmatrix}$$

Feature matrix $X^T =$

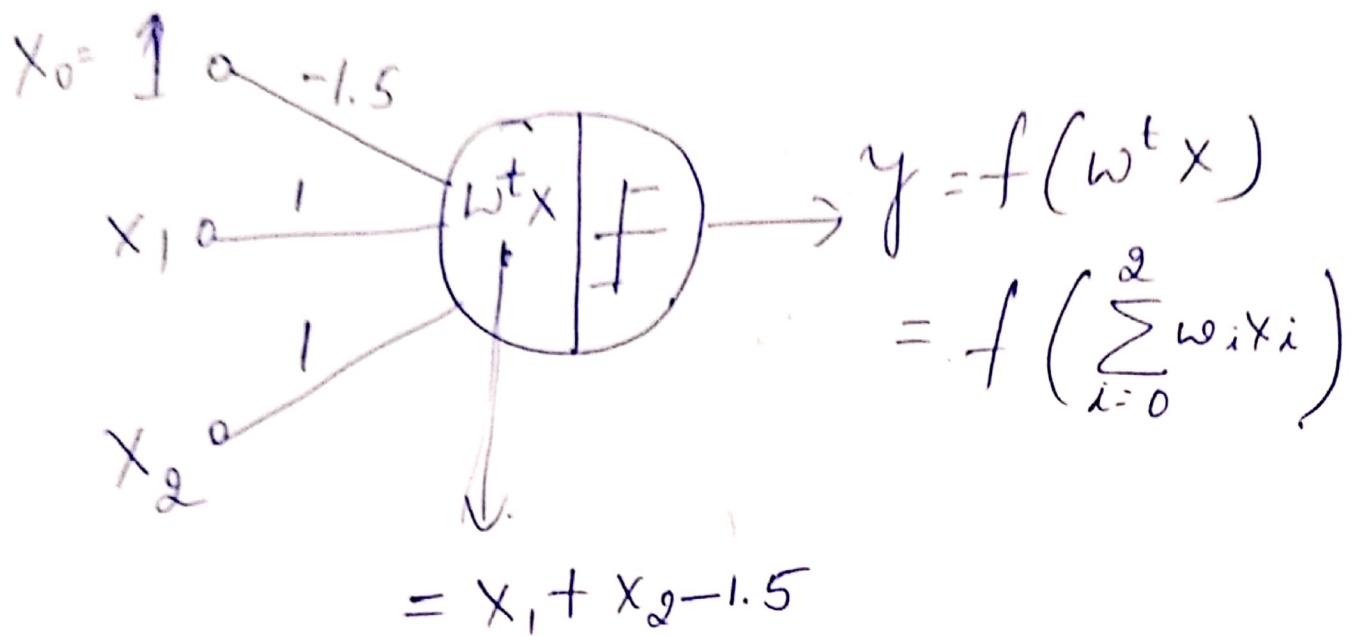
$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix}$$

$$x^T w = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} -2 & 1 & 5 \\ 1 & 1 & 1 \end{bmatrix}.$$

$$= \begin{bmatrix} -1.5 \\ -0.5 \\ -0.5 \\ 0.5 \end{bmatrix}$$

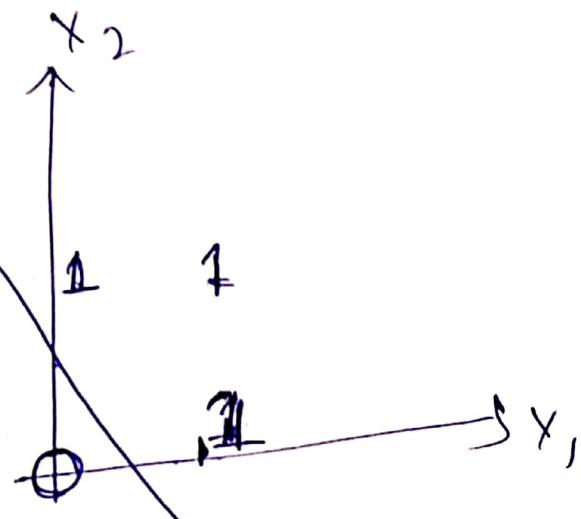


$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$



OR Function

x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	1

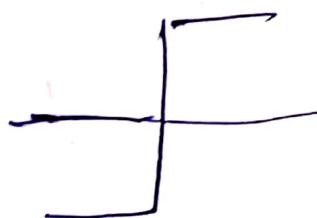


$$x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \quad w = \begin{bmatrix} -0.5 \\ 1 \\ 1 \end{bmatrix} \quad \rightarrow x_1 + x_2 - 0.5 = 0$$

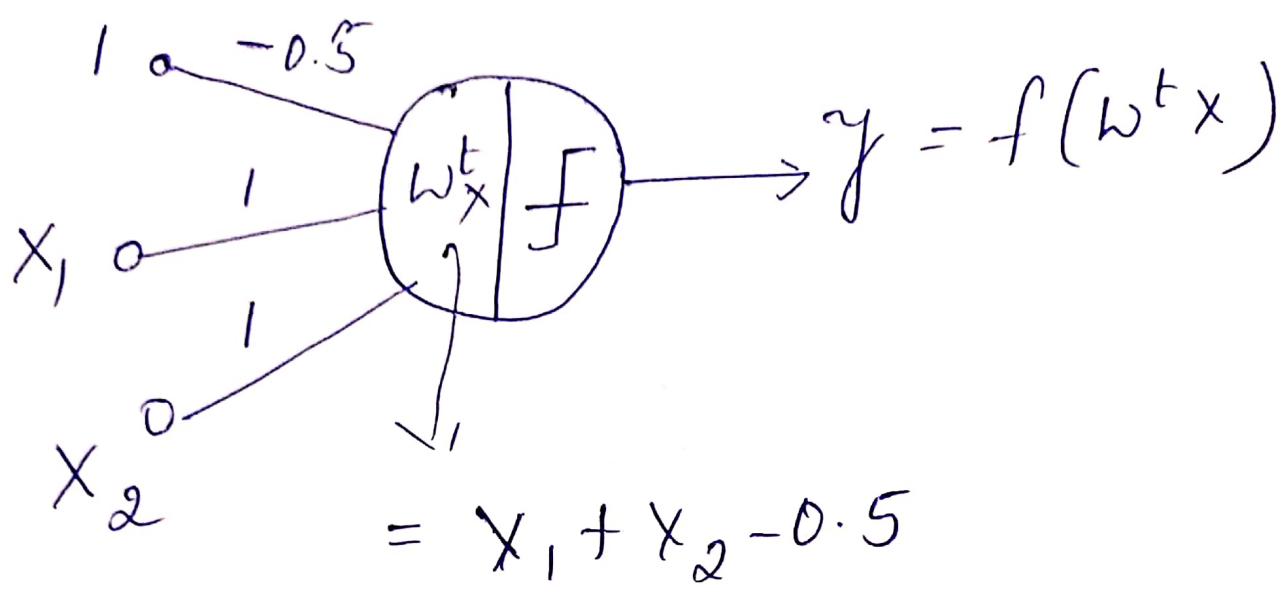
$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$X^T \omega = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \cdot \begin{bmatrix} -0.5 \\ 1 \\ 1 \end{bmatrix}$$

$$= \begin{bmatrix} -0.5 \\ 0.5 \\ 0.5 \\ 1.5 \end{bmatrix}$$



$$\begin{bmatrix} 0 \\ 1 \\ 1 \end{bmatrix}$$

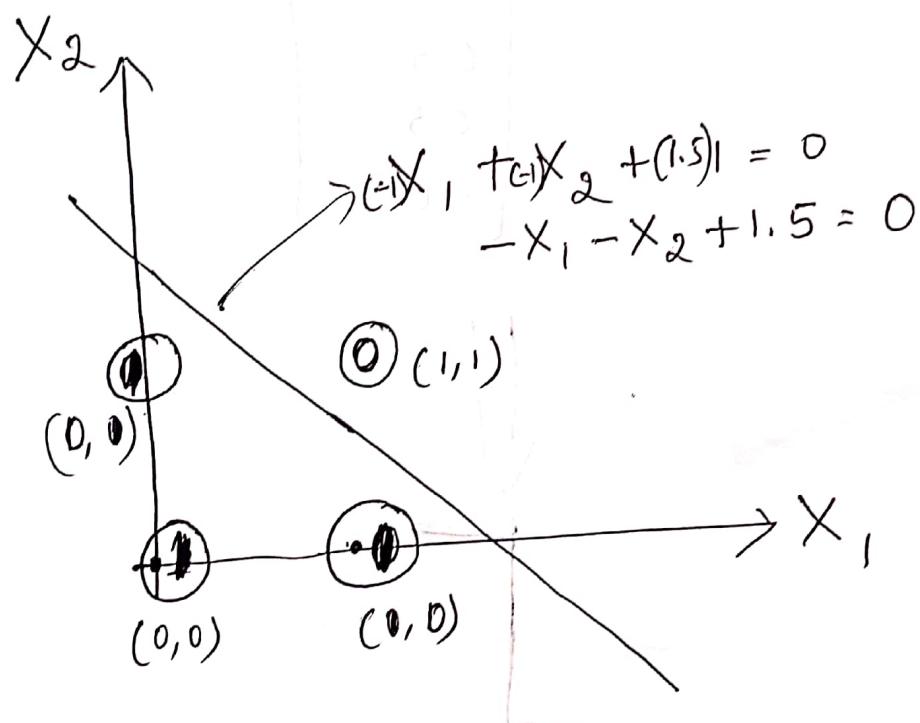


NAND Function

X_1	X_2	Y
0	0	1
0	1	1
1	0	1
1	1	0

where $y = \overline{X_1 \cdot X_2}$

i.e., Complement
of AND operation



Since it is complement of AND operation,

Weight Vector = $\vec{w} = \begin{bmatrix} -1 \\ -1 \\ 1.5 \end{bmatrix}$

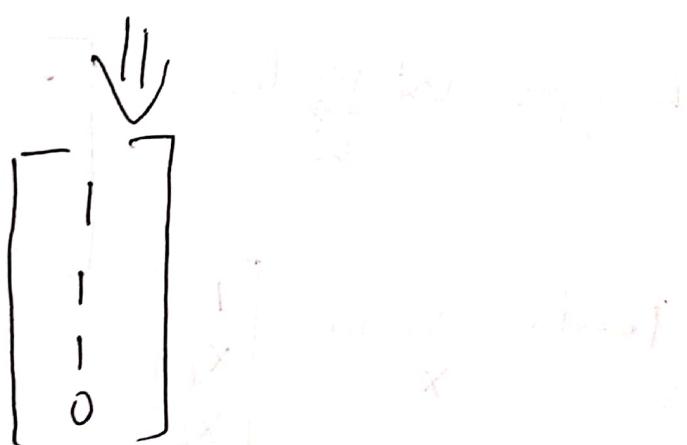
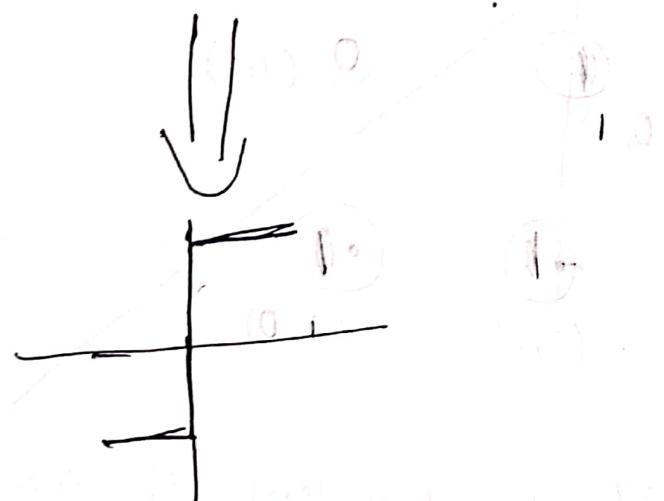
Feature Vector = $\vec{x} = \begin{bmatrix} 1 \\ X_1 \\ X_2 \end{bmatrix}$

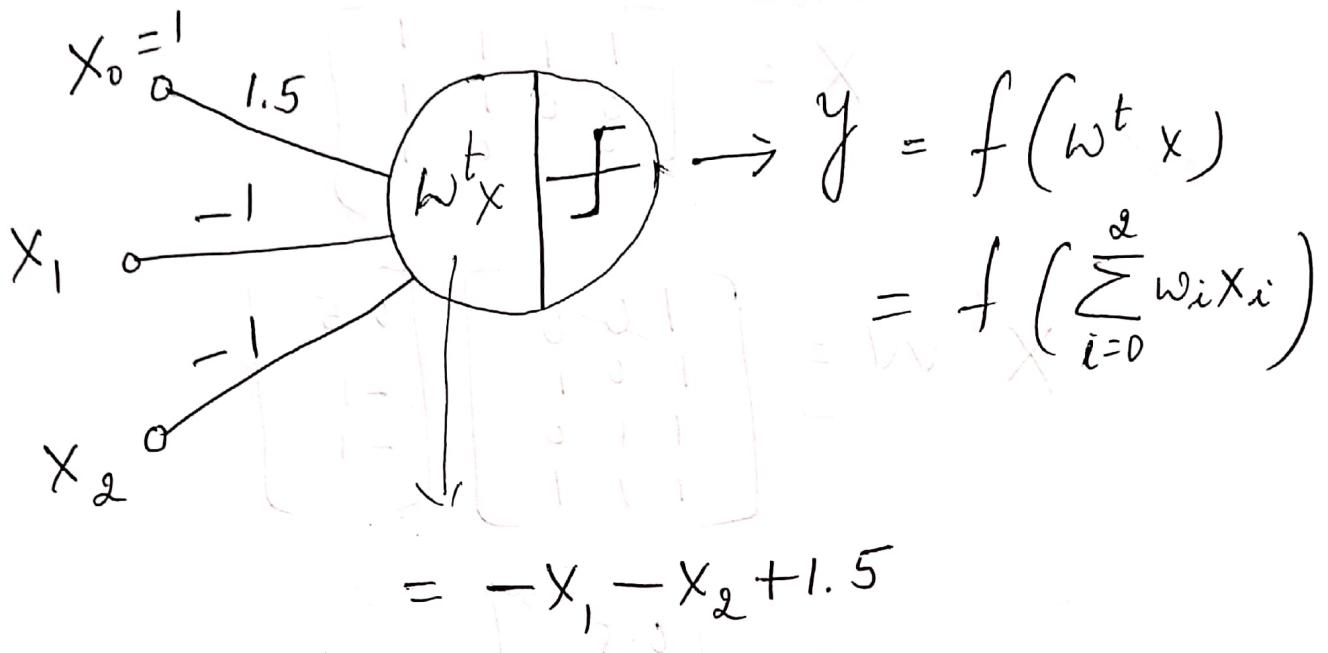
Feature Matrix

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$X^t W = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1.5 \\ -1 \\ -1 \end{bmatrix}$$

$$= \begin{bmatrix} 1.5 \\ 0.5 \\ 0.5 \\ -0.5 \end{bmatrix}$$



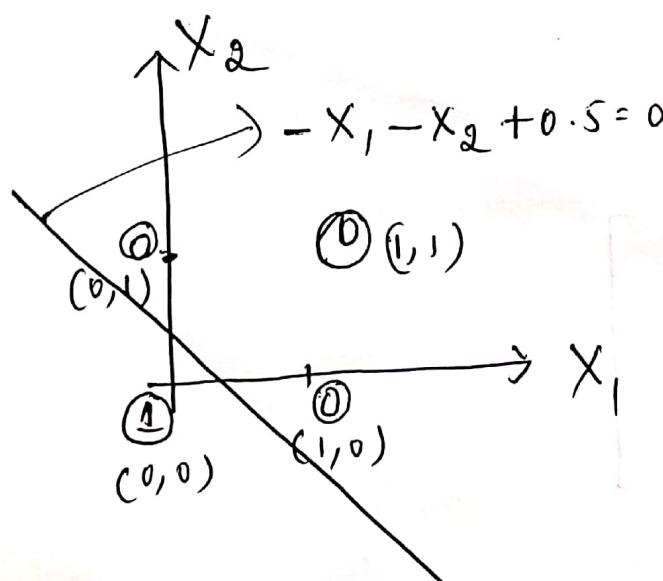


NOR operation

\$x_1\$	\$x_2\$	\$y\$
0	0	1
0	1	0
1	0	0
1	1	0

$$y = \overline{x_1 + x_2}$$

Complement of
OR operation



Since it is
complement of
OR operations

$$\vec{w} = \begin{bmatrix} 0.5 \\ -1 \\ -1 \end{bmatrix}$$

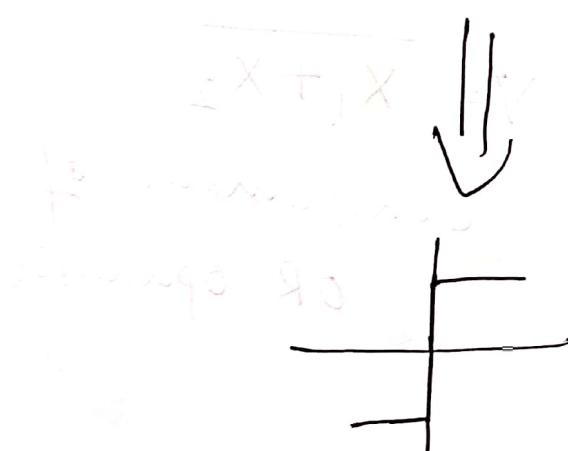
$$\text{Feature Vector} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix}$$

Feature Matrix

$$X = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 \end{bmatrix}$$

$$X^T w = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 0.5 \\ -1 \\ -1 \end{bmatrix}$$

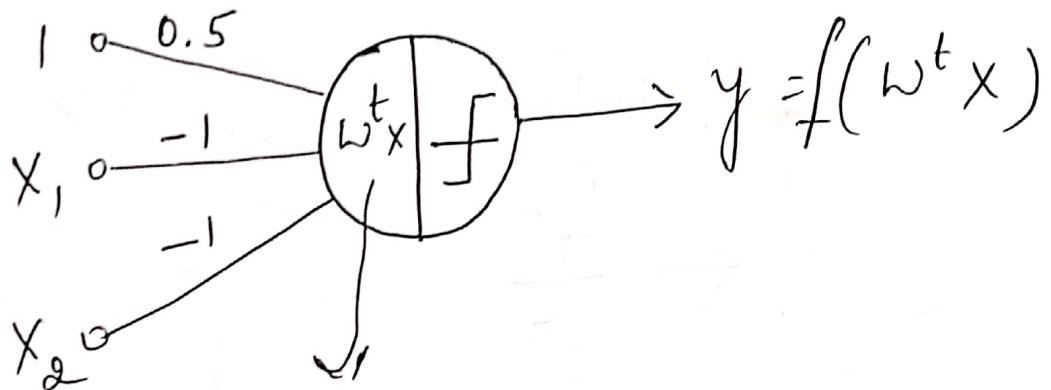
$$= \begin{bmatrix} 0.5 \\ -0.5 \\ -0.5 \\ -1.5 \end{bmatrix}$$



$$X^T \cdot X^T \cdot X$$

$$\begin{bmatrix} 4 & 2 & 2 & 2 \\ 2 & 6 & 2 & 2 \\ 2 & 2 & 6 & 2 \\ 2 & 2 & 2 & 4 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

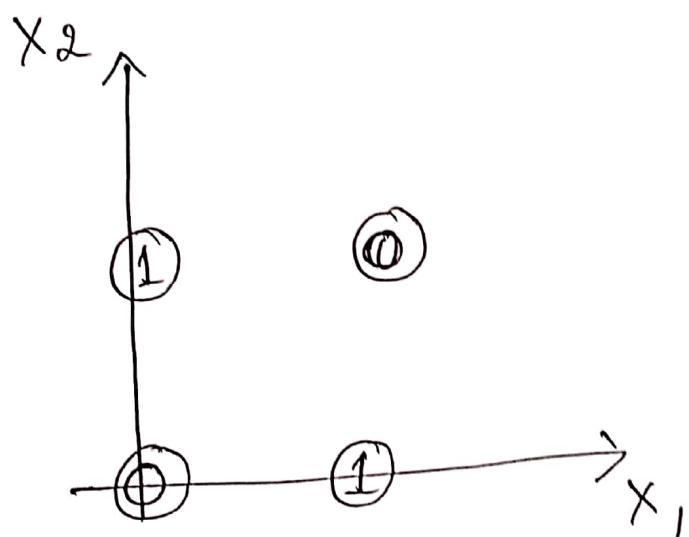


$$= -x_1 - x_2 + 0.5$$

XOR

$$X_1 \oplus X_2$$

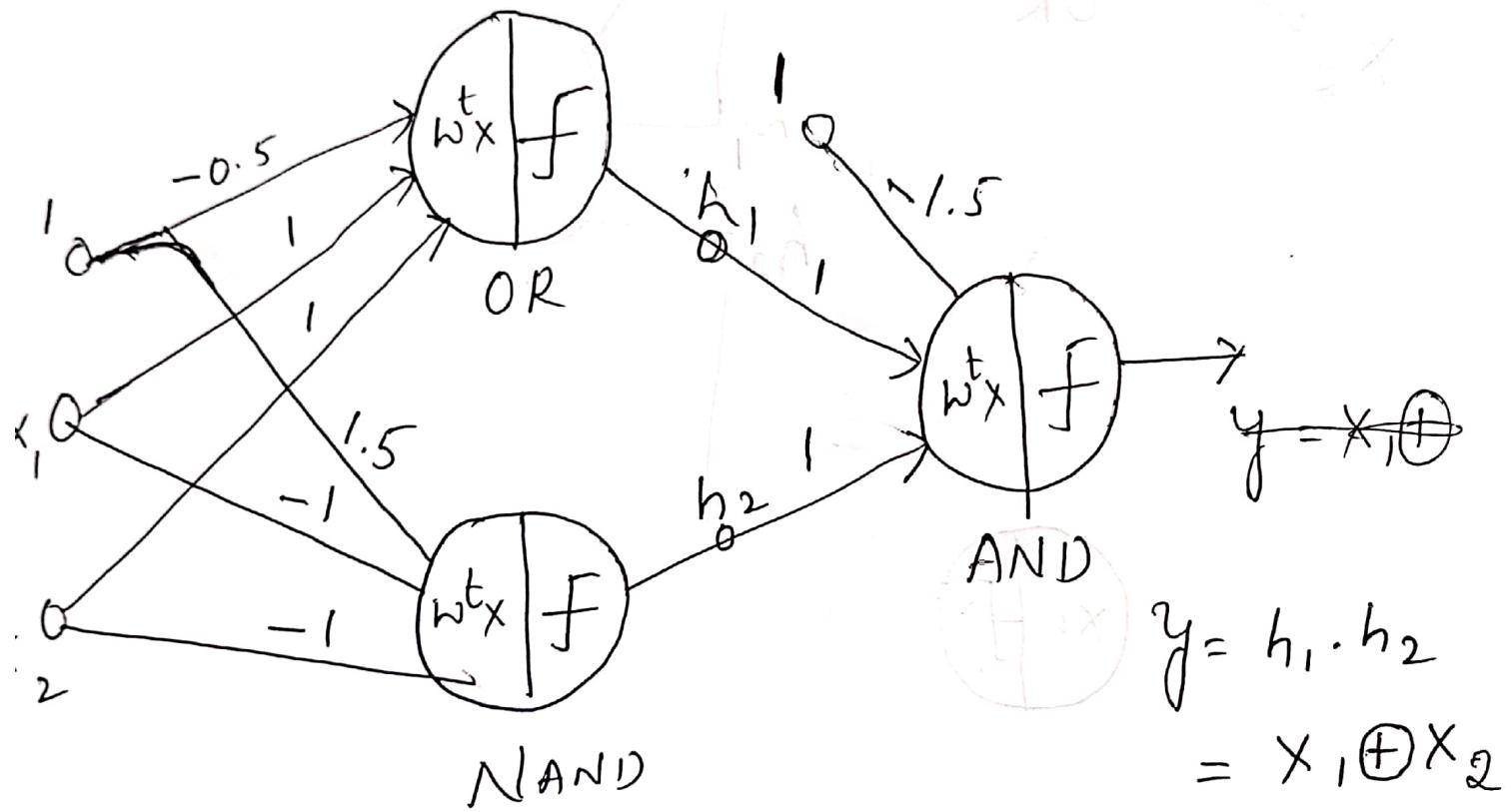
X_1	X_2	Y
0	0	0
0	1	1
1	0	1
1	1	0



$$X_1 \oplus X_2 = (X_1 + X_2) \cdot (\overline{X}_1 + \overline{X}_2)$$

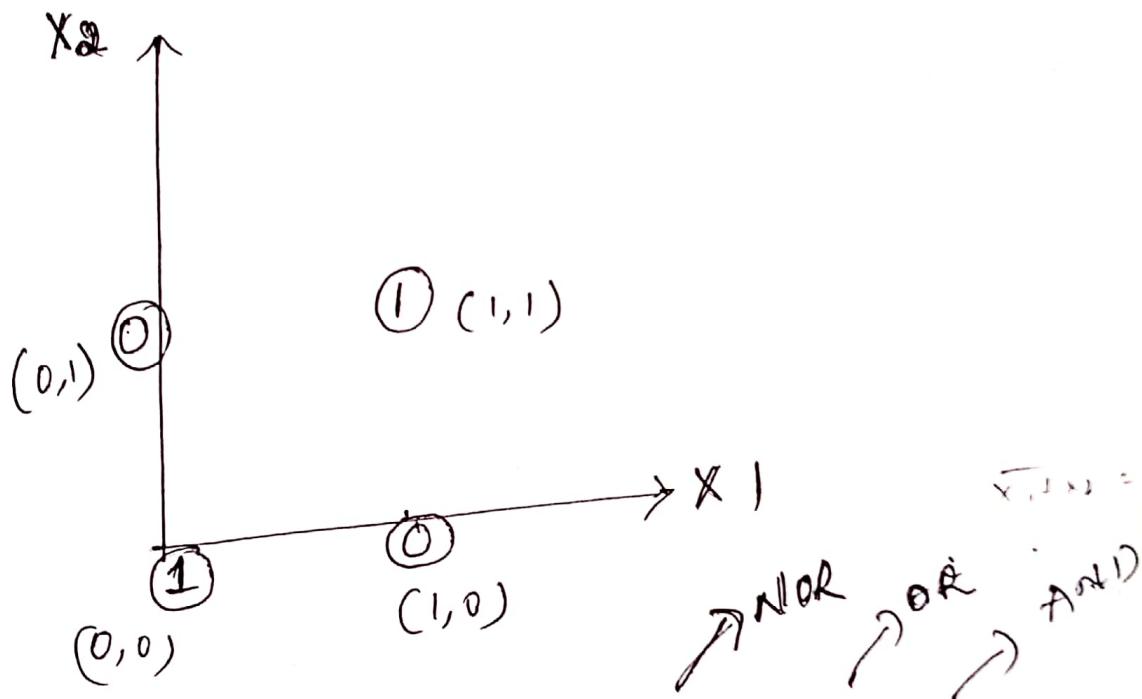
$\curvearrowright X_1 \cdot X_2$

X_1	X_2	$h_1 = (X_1 + X_2)$	$h_2 = \overline{X}_1 + \overline{X}_2$	$h_1 \cdot h_2 = \cancel{X_1} \cdot \cancel{X_2}$ $= X_1 \oplus X_2$
0	0	0	1	0
0	1	1	0	1
1	0	1	0	1
1	1	1	0	0



XNOR : $(\overline{x_1 \oplus x_2})$

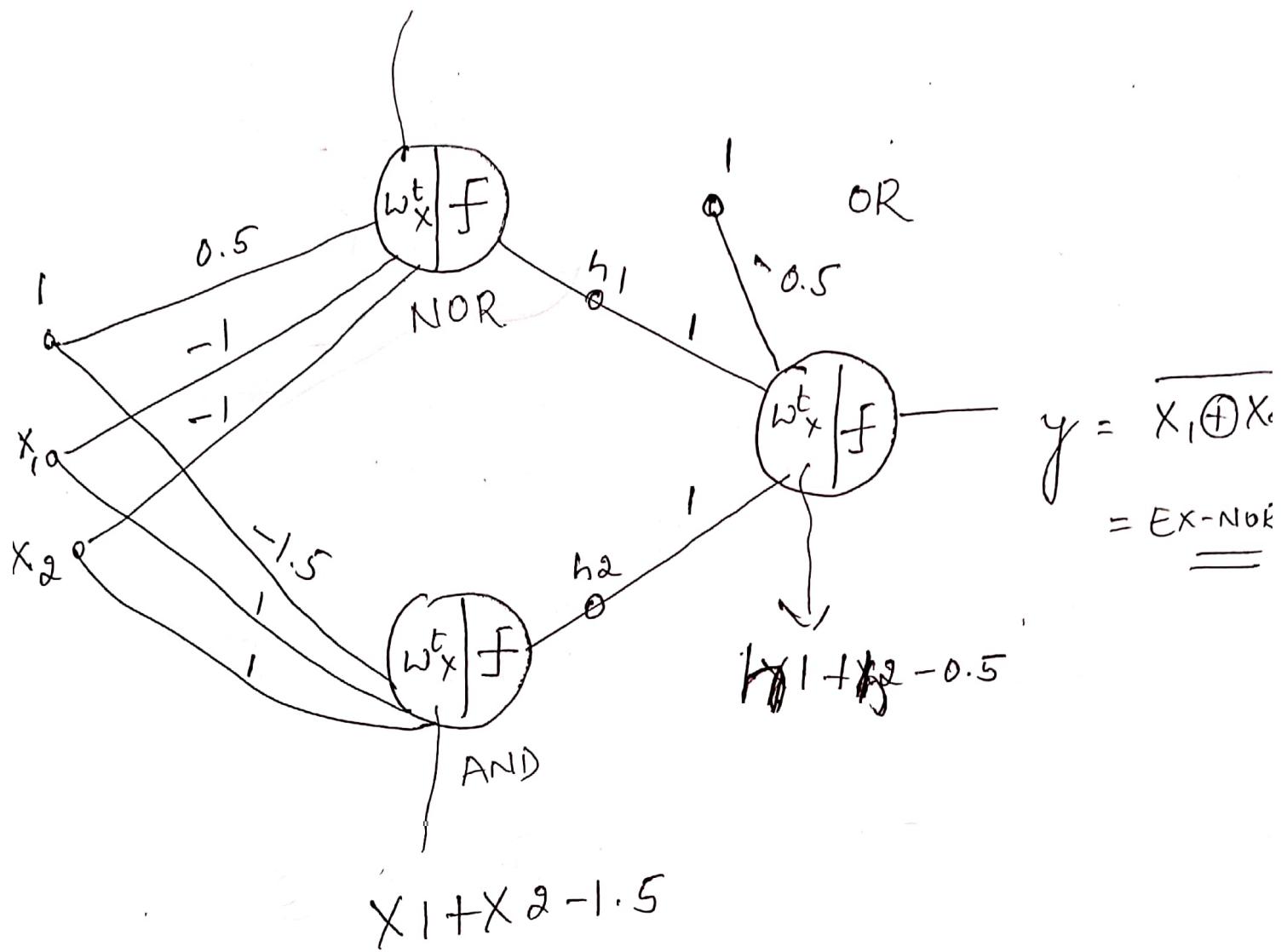
x_1	x_2	y
0	0	1
0	1	0
1	0	0
1	1	1



$$\overline{x_1 \oplus x_2} = (\overline{x_1+x_2}) + (x_1 \cdot x_2)$$

x_1	x_2	$h_1 = \overline{x_1+x_2}$	$h_2 = x_1 \cdot x_2$	$h_1 + h_2 = \overline{x_1 \oplus x_2}$
0	0	1	0	1
0	1	0	0	0
1	0	0	0	0
1	1	0	1	1

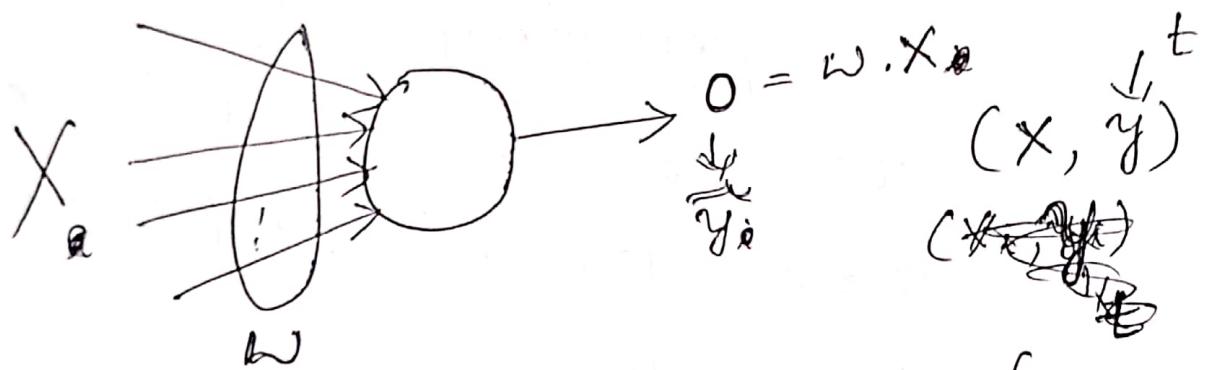
$$-x_1 - x_2 + 0.5$$



Perceptron Training Rule

Ques

How to Learn weights for
a single perceptron?



Procedure

- Begin with random weights for weight vector
- Then apply iteratively apply the perceptron to each training example
- Modify the perceptron weights whenever perceptron misclassifies an example.
- This process is repeated, iterating through the training examples as many ~~as~~ times as needed until

Perception classifies all training examples correctly.

 Weights are modified at each step according to the Perceptron training rule,
which revises the weights w_i associated with input x_i according to rule.

$$w_i \leftarrow w_i + \Delta w_i$$

Where $\Delta w_i = \eta(t - o)x_i$

Where t = target output for the current training example

o = output generated by the perceptron

η = learning rate
(Convergence rate)
↳ positive constant

Note:- The role of the learning rate (η) is to moderate the degree to which weights are changed at each step.

It is usually set to small value (say e.g.: 0.1).

~~Ques.~~ ~~Ans.~~

Case 1

Note:- If training example is correctly classified

Then $(t - o)$ is zero

& Δw_i is zero.

i.e. no weights are updated.

Case 2: If training example is misclassified.

$$t = 1, o = -1 \quad \checkmark$$

Then weights must be altered to increase the value of $\vec{w} \cdot \vec{x}$.

Example:

if $x_i = 0.8$, $\eta = 0.1$,

$t = 1$ & $o = -1$

Then $\Delta w_i = \eta(t - o)x_i$

$$\begin{aligned} &= 0.1(1 - (-1))0.8 \\ &= \underline{\underline{0.16}} \end{aligned}$$

i.e., if $x_i > 0$, Then increasing
 w_i will bring the perception
closer to correctly classify

This Example. \rightarrow point

if $t = -1$, and $o = 1$, Then
Weights associated with positive
 x_i will be Decreased.

Note:- The above learning procedure can be proven to converge within a finite number of applications of perceptron training rule to a weight vector that correctly classifies all training examples provided the training examples are linearly separable

* Gradient Descent and Delta Rule

(weight update rule)

→ Limitation of Perception training rule

→ If the training examples are linearly separable, perception training rule can fail to converge.

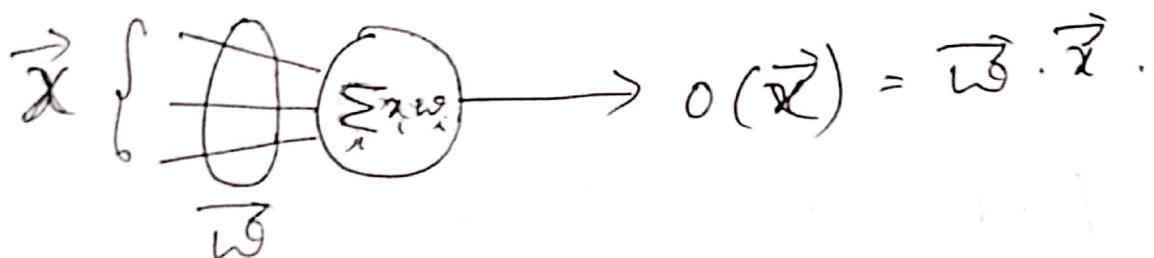
Delta rule is designed to overcome this difficulty.

→ If the training examples are not linearly separable, the Delta rule converges toward a best-fit approximation to the target concept.

→ Delta rule uses Gradient descent

Approach to find learning weight
vectors that best fit the training
examples.

In order to understand Delta rule,
consider the task of training an
unthresholded perceptron,
i.e., linear unit



From Training Error (Loss)

Can be measured in many
ways

④ Squared Error Loss

④ Cross Entropy Loss

Consider Squared Error Loss
relative to training examples
& is given by

$$E(\vec{w}) \equiv \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2 \quad \boxed{1}$$

where $D \rightarrow$ Set of Training Examples
 $t_d \rightarrow$ target off for training example d.

$o_d \rightarrow$ output of linear unit
for training example d.

Now, we have to

Minimize E Let weight vector \vec{w}

Solution is * Gradient descent search

Gradient descent search.

determines a weight vector that

minimizes E by

* Starting with arbitrary initial weight vector

* Modifying weight vector in small steps.

* At each step weight vector is altered in the direction that produces steepest descent along the error surface.

④ This process continues until the global & minimum error is reached.

∴ Gradient of E with respect to \vec{w} , $\nabla E(\vec{w})$

$$\left(\frac{\partial E}{\partial \vec{w}} \right)$$

is given by

$$\boxed{\nabla E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]}$$

Since the Gradient specifies the direction of steepest increase of E ,

Training rule for Gradient descent is

$$\boxed{\vec{w} \leftarrow \vec{w} + \Delta \vec{w}} \rightarrow ③$$

$$\text{where } \Delta \vec{w} = -\eta \nabla E(\vec{w})$$

→ ④

η is a positive constant called
The learning rate,
It determines the Step Size
in the Gradient Descent
Search.

The -ve sign is used to move
the weight vector in the
direction that decreases E.

Eq(4) can also be written as

$$\boxed{\nabla w_i = -\eta \frac{\partial E}{\partial w_i}} \rightarrow ⑤$$

Thus

$$\begin{aligned} \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \left[\sum_{d \in D} (t_d - o_d)^2 \right] \\ &= \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - o_d)^2 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{1}{2} \sum_{d \in D} \ell(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d) \\
 &= \sum_{d \in D} (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot \vec{x}_d) \\
 &= \sum_{d \in D} (t_d - o_d) (-x_{id})
 \end{aligned}$$

~~where x_{id} denotes the single input component x_i for the training example d .~~

where x_{id} denotes the single input component x_i for the training example d .

Then weight update rule is

$$w_i \leftarrow w_i + \Delta w_i \quad (6)$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$= -\eta (o_d - t)$$

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) x_{id}$$

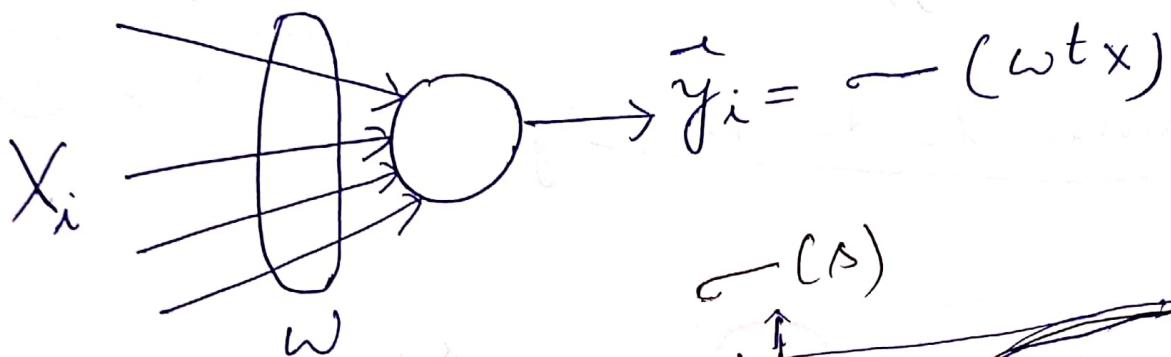
Substitute $\hat{L}(\theta)$ in $L(\theta)$

Weight update rule for gradient
descent in G_{new}

$$w_i \leftarrow w_i + \eta \sum_{d \in D} (t_d - Q_d) x_{id}$$

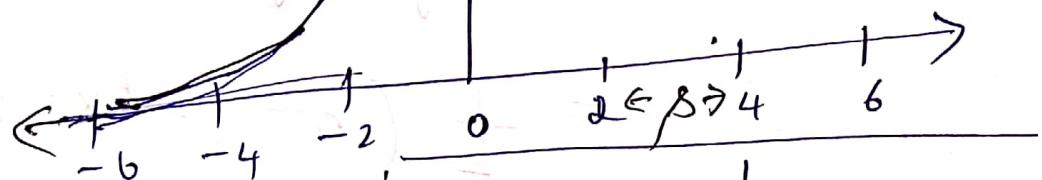
~~for Stochastic optimization~~
Batch optimization

Single Layer Network - Single output
having non-linearity



Sigmoid
function

$$\sigma(\beta)$$



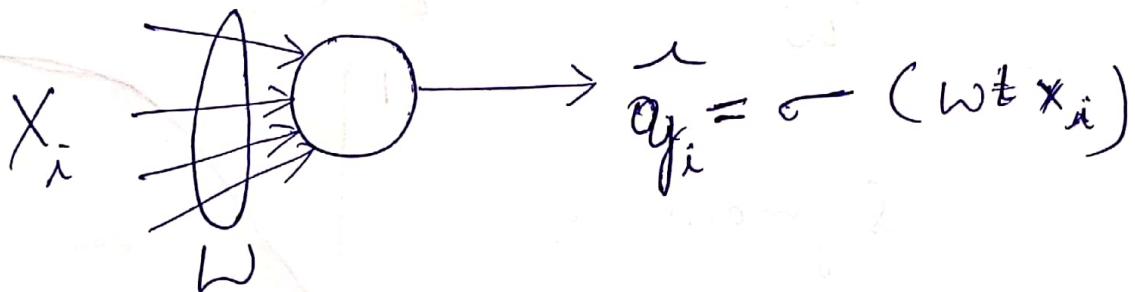
$$\sigma(\beta) = \frac{1}{1 + e^{-\beta}}$$

- (*) as β approaches ∞
 $\sigma(\beta)$ approaches 1
- (**) as β approaches $-\infty$
 $\sigma(\beta)$ approaches 0

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$

$$\frac{\partial \sigma(s)}{\partial s} = \sigma(s)(1 - \sigma(s))$$

Now to do training



$$y_i = \sigma(w^T x_i)$$

~~del~~ ~~del~~ ~~del~~

$$E = \frac{1}{2} (\hat{y}_i - y_i)^2$$

$$= \frac{1}{2} (\sigma(w^T x_i) - y_i)^2$$

$$\nabla_w E = \frac{\partial E}{\partial w} = \frac{\partial}{\partial w} \frac{1}{2} (\sigma(w^T x_i) - y_i)^2$$

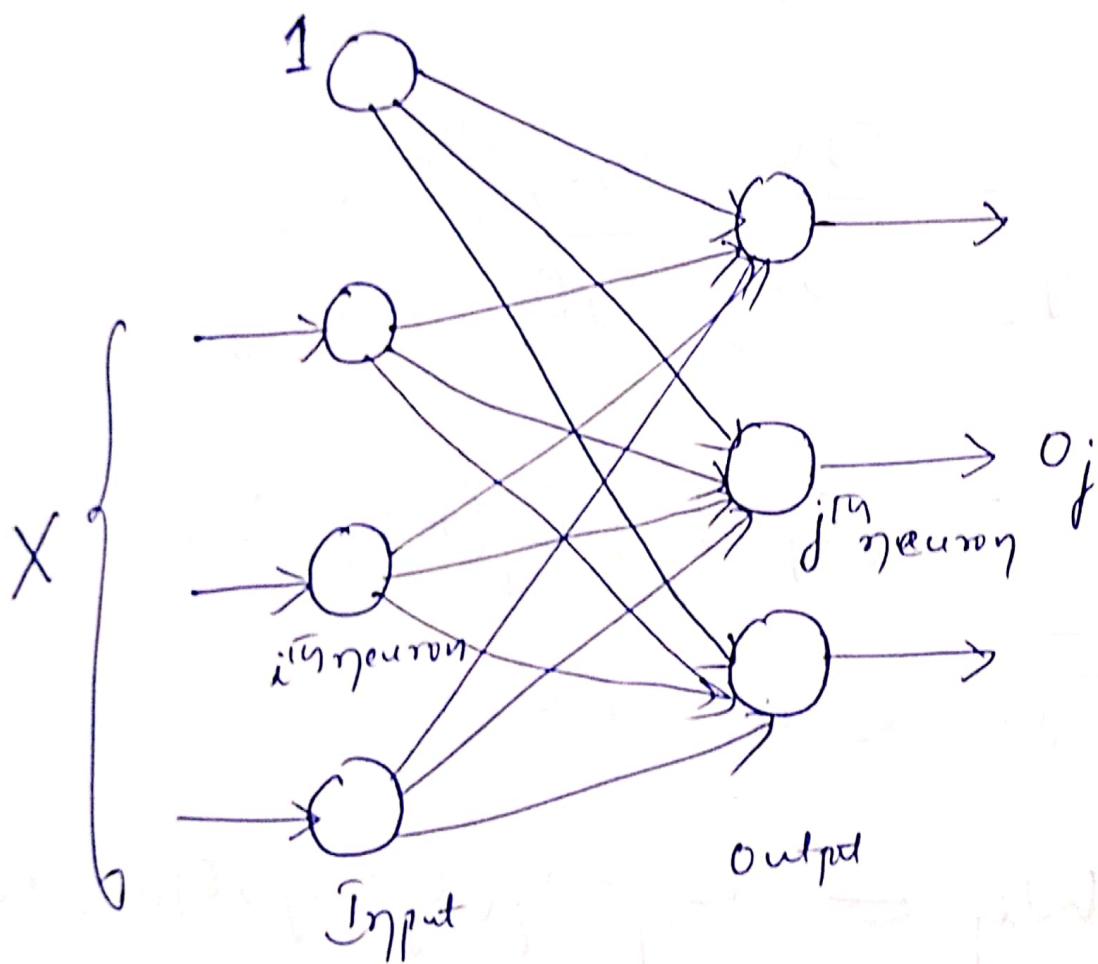
$$= \frac{\partial}{\partial w} \left[\frac{1}{2} \left(\frac{1}{1 + e^{-w^T x_i}} - y_i \right)^2 \right]$$

$$\nabla_w E = \hat{y}_i (1 - \hat{y}_i) (\hat{y}_i - y_i) x_i$$

Weight update rule \Rightarrow

$$W \leftarrow W - \eta \hat{y}_i (1 - \hat{y}_i) (\hat{y} - y_i) X_i$$

Back propagation Learning: Single layer
Multiple output



$$o_j = \frac{1}{1 + e^{-\theta_j}}$$

$$\theta_j = \sum_{i=1}^D w_{ij} x_i$$

$$E = \frac{1}{2} \sum_{j=1}^M (o_j - t_j)^2$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial o_j} \cdot \frac{\partial o_j}{\partial o_j} \cdot \frac{\partial o_j}{\partial w_{ij}}$$

$$= (o_j - t_j) o_j (1 - o_j) x_i$$

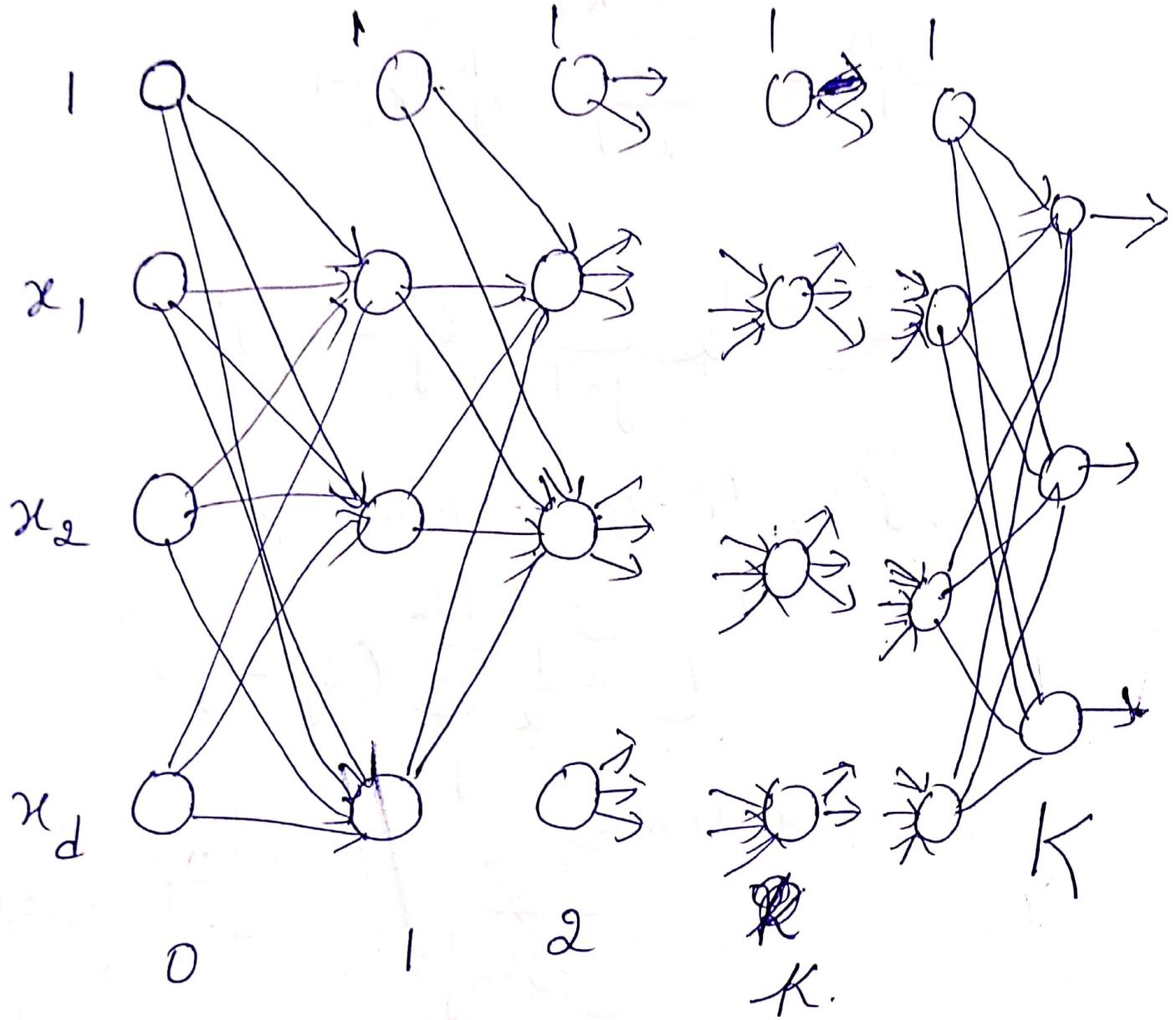
(Weight update rule \Rightarrow)

$$w_{ij} \leftarrow w_{ij} - \gamma (o_j - t_j) o_j (1 - o_j) x_i$$

$$W = M \times D$$

Every row corresponds to a weight vector of a particular class

Multi Layer perceptron



$M_K \rightarrow$ No. of Nodes in

the K^{th} layer.

$$O_j^K = \frac{1}{1 + e^{-\theta_j^K}}$$

$$\theta_j^K = \sum_{i=1}^{M_{K-1}} w_{ij}^K x_i^{K-1}$$

$$E = \frac{1}{2} \sum_{j=1}^{M_K} (O_j^K - t_j)^2$$

Find w_{ij}^K to minimize E .

$$E = \frac{1}{2} \sum_{j=1}^{M_K} (O_j^K - t_j)^2$$

Gradient Descent

$$\frac{\partial E}{\partial w_{ij}^K}$$

$$\frac{\partial O_j^K}{\partial w_{ij}^K}$$

$$\frac{\partial E}{\partial w_{ij}^K} = \frac{\partial E}{\partial O_j^K} \frac{\partial O_j^K}{\partial \theta_j^K} \frac{\partial \theta_j^K}{\partial w_{ij}^K}$$

$$\left. \begin{aligned} O_j^K &= \frac{1}{1 + e^{-\theta_j^K}} \\ \theta_j^K &= \sum_{i=1}^{M_{K-1}} w_{ij}^K O_i^{K-1} \end{aligned} \right\}$$

$$= (O_j^K - t_j) O_j^K (1 - O_j^K) O_i^{K-1}$$

Let ;

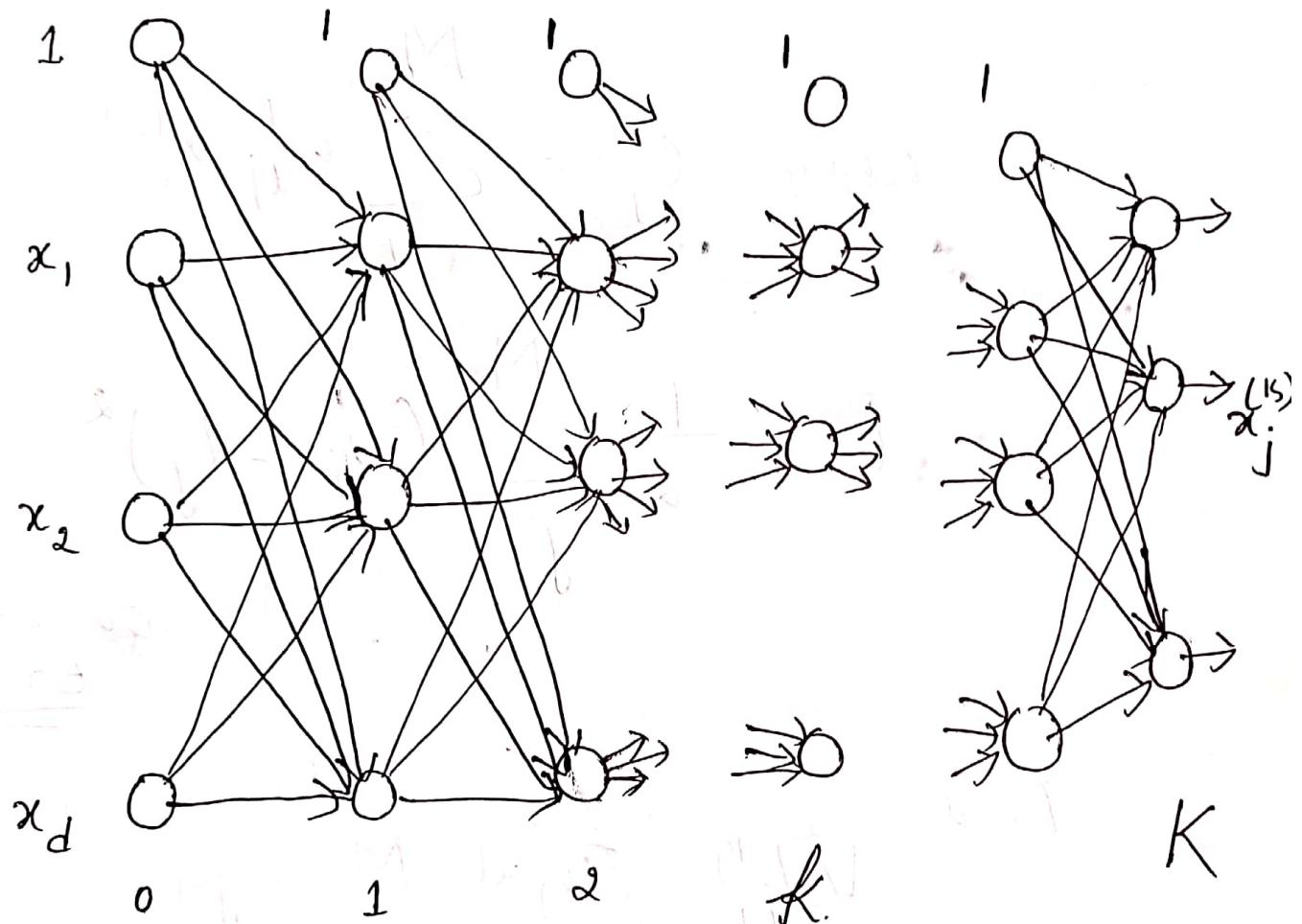
$$\delta_j^{ls} = o_j^{ls}(1 - o_j^{ls})(o_j^{ls} - t_j)$$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}^{ls}} = \delta_j^k o_i^{k-1}$$

Then Weight updation rule
for O/P layer.

$$w_{ij}^k \leftarrow w_{ij}^{ls} - \eta \delta_j^k o_i^{k-1}$$

Multilayer Perceptron: Backpropagation Learning



$M_K \rightarrow$ No. of nodes in
The K^{th} layer

L15 T

$$O_j^{15} = \frac{1}{1 + e^{-\theta_j^{15}}}$$

where $O_j^{15} = \sum_{i=1}^{M_{15-1}} w_{ij}^{15} x_i^{15-1}$

$$E = \frac{1}{2} \sum_{j=1}^{M_{15}} (O_j^{15} - t_j)^2$$

Squared Error

Find

w_{ij}^{15} that minimizes

$$E = \frac{1}{2} \sum_{j=1}^{M_{15}} (O_j^{15} - t_j)^2$$

Squared Error

Gradient Descent

$$\frac{\partial E}{\partial w_{ij}^k}$$

Using chain rule

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{\partial E}{\partial o_j^k} \frac{\partial o_j^k}{\partial \theta_j^k} \frac{\partial \theta_j^k}{\partial w_{ij}^k}$$

$$= (o_j^k - t_j)$$

$$= (o_j^k - t_j) o_j^k (1 - o_j^k) o_i^{k-1}$$

$$o_j^k = \frac{1}{1 + e^{-o_j^k}}$$

$$o_j^k = \sum_{i=1}^{M_{k-1}} w_{ij}^k o_i^{k-1}$$

Let $s_j^k = o_j^k (1 - o_j^k) (o_j^k - t_j)$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}^k} = s_j^k o_i^{k-1}$$

Weight updating rule for output layer.

$$w_{ij}^k \leftarrow w_{ij}^k - \eta s_j^k o_i^{k-1}$$

For hidden layers

$$E = \frac{1}{2} \sum_{j=1}^{M_{k-1}} (o_j^{k-1} - t_j)^2$$

Let p^{th} neuron in $(k-2)^{\text{th}}$ layer

$$\downarrow w_{pi}^{k-1}$$

i^{th} neuron in $(k-1)^{\text{th}}$ layer

$$\downarrow w_{ij}^k$$

j^{th} neuron in k^{th} layer

$$\frac{\partial E}{\partial w_{pi}^{k-1}} = \frac{\partial E}{\partial o_i^{k-1}} \frac{\partial o_i^{k-1}}{\partial w_{pi}^{k-1}}$$

$$= \frac{\partial E}{\partial o_i^{k-1}} \cdot \frac{\partial o_i^{k-1}}{\partial o_i^{k-1}} \frac{\partial o_i^{k-1}}{\partial w_{pi}^{k-1}}$$

$$= \frac{\partial E}{\partial o_i^{k-1}} \cdot o_i^{k-1} (1 - o_i^{k-1}) \cdot o_p^{k-2}$$

$$o_i^{k-1} = \frac{1}{1 + e^{-o_i^{k-1}}}$$

$$o_i^{k-1} = \sum_{p=1}^{M_{k-2}} w_{pi}^{k-1} o_p^{k-2}$$

$$\frac{\partial E}{\partial o_i^{k-1}} = \frac{\partial E}{\partial o_j^k} \cdot \frac{\partial o_j^k}{\partial o_i^k} \cdot \frac{\partial o_i^k}{\partial o_i^{k-1}}$$

$$= \sum_{j=1}^{M_{15}} (o_j^k - t_j) \cdot o_j^k (1 - o_j^k)$$

$$E = \frac{1}{2} \sum_{j=1}^{M_{15}} (o_j^k - t_j)^2$$

$$= \sum_{j=1}^{M_{15}} (o_j^k - t_j) o_j^k (1 - o_j^k) w_{ij}^k$$

$$o_j^k = \frac{1}{1 + e^{-o_j^k}}$$

$$= \sum_{j=1}^{M_{15}} o_j^k w_{ij}^k$$

$$o_j^k = \sum_{i=1}^{M_{15-1}} w_{ij}^k o_i^{k-1}$$

$$s_j^k = o_j^k (1 - o_j^k)$$

$$(o_j^k - t_j)$$

$$\frac{\partial E}{\partial w_{pi}^{k-1}} = \frac{\partial E}{\partial x_i^{k-1}} \cdot o_i^{k-1} (1 - o_i^{k-1}) \cdot o_p^{k-2}$$

$$\frac{\partial E}{\partial w_{pi}^{k-1}} = O_i^{k-1} (1 - O_i^{k-1}) \cdot O_p^{k-2} \sum_{j=1}^{M_k} \partial_j^k w_{ij}^k$$

Putting $\delta_i^{k-1} = O_i^{k-1} (1 - O_i^{k-1}) \sum_{j=1}^{M_k} \partial_j^k w_{ij}^k$

Then

Weight updating rule for

Last but output layer

$$w_{pi}^{k-1} \leftarrow w_{pi}^{k-1} - \eta \delta_i^{k-1} O_p^{k-2}$$

In General

For any hidden layer weight w_{ij}^k

$$\text{putting } \delta_i^k = O_i^k (1 - O_i^k) \sum_{j=1}^{M_{k+1}} \partial_j^{k+1} w_{ij}^{k+1}$$

Weight updating rule

$$w_{ij}^k \leftarrow w_{ij}^k - \eta \delta_i^k O_i^{k-1}$$