

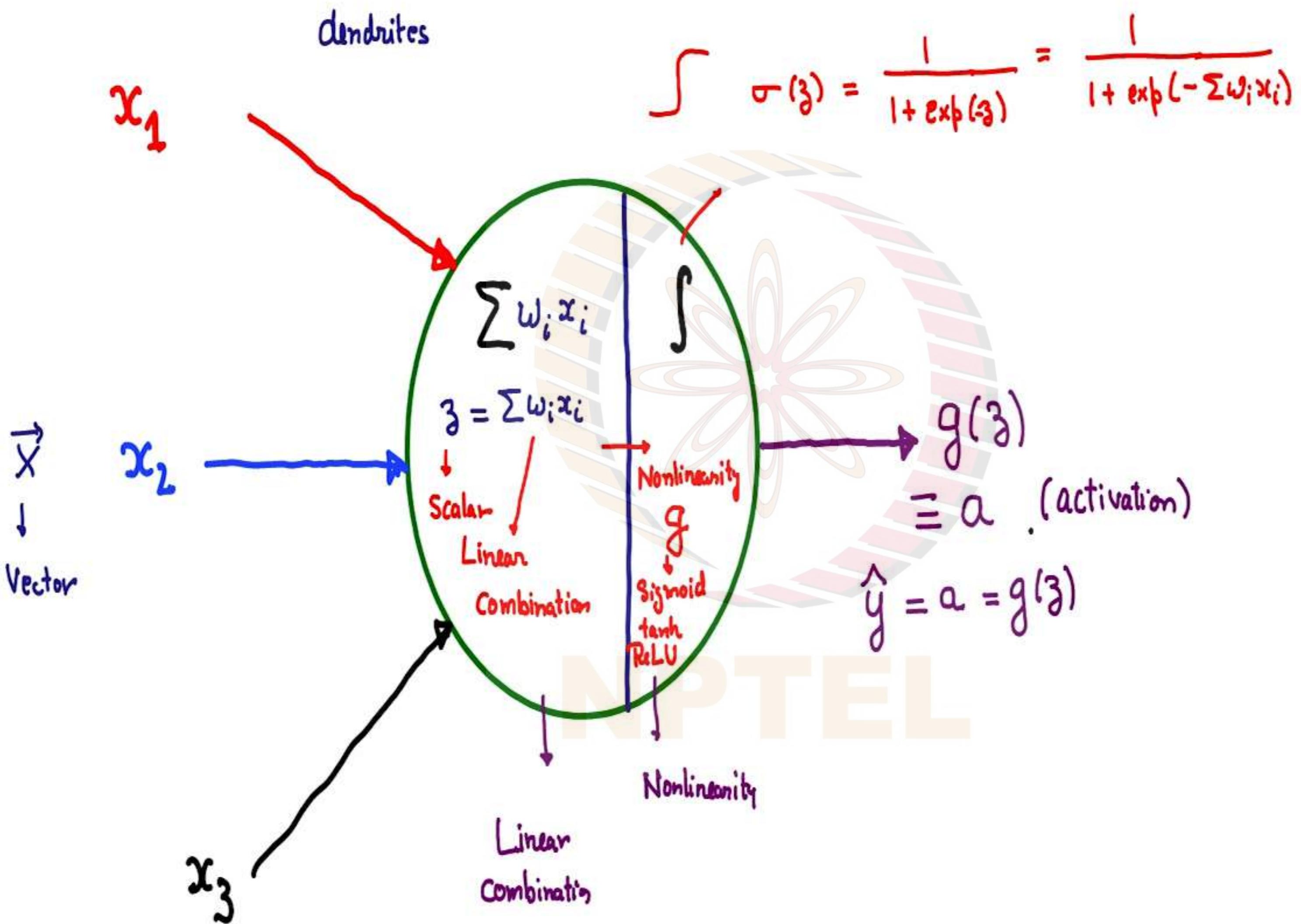
## 1. Structure of an artificial neuron

a. Linear combination

b. Nonlinear activation

Also used in  
logistic regression

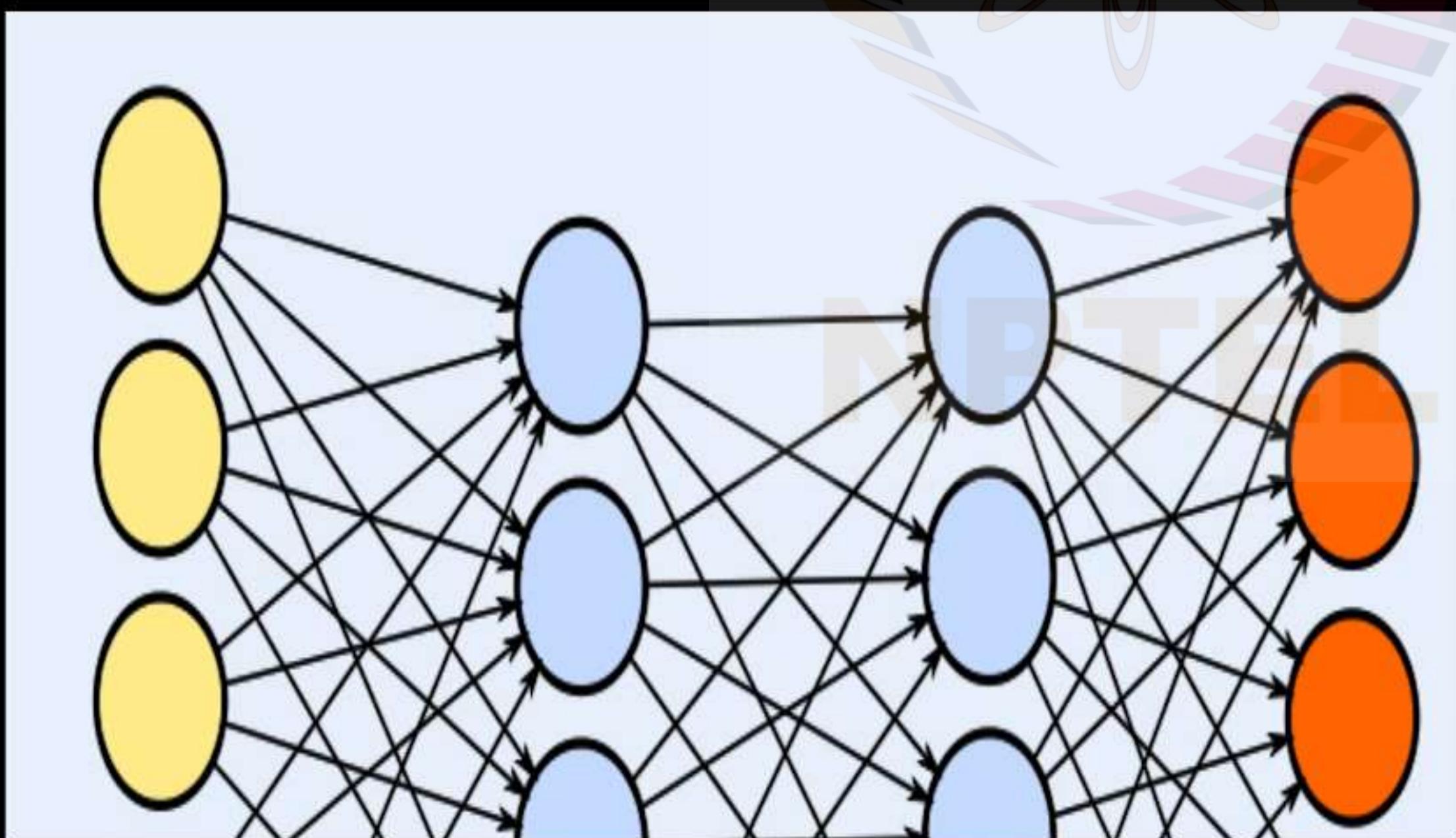
NPTEL

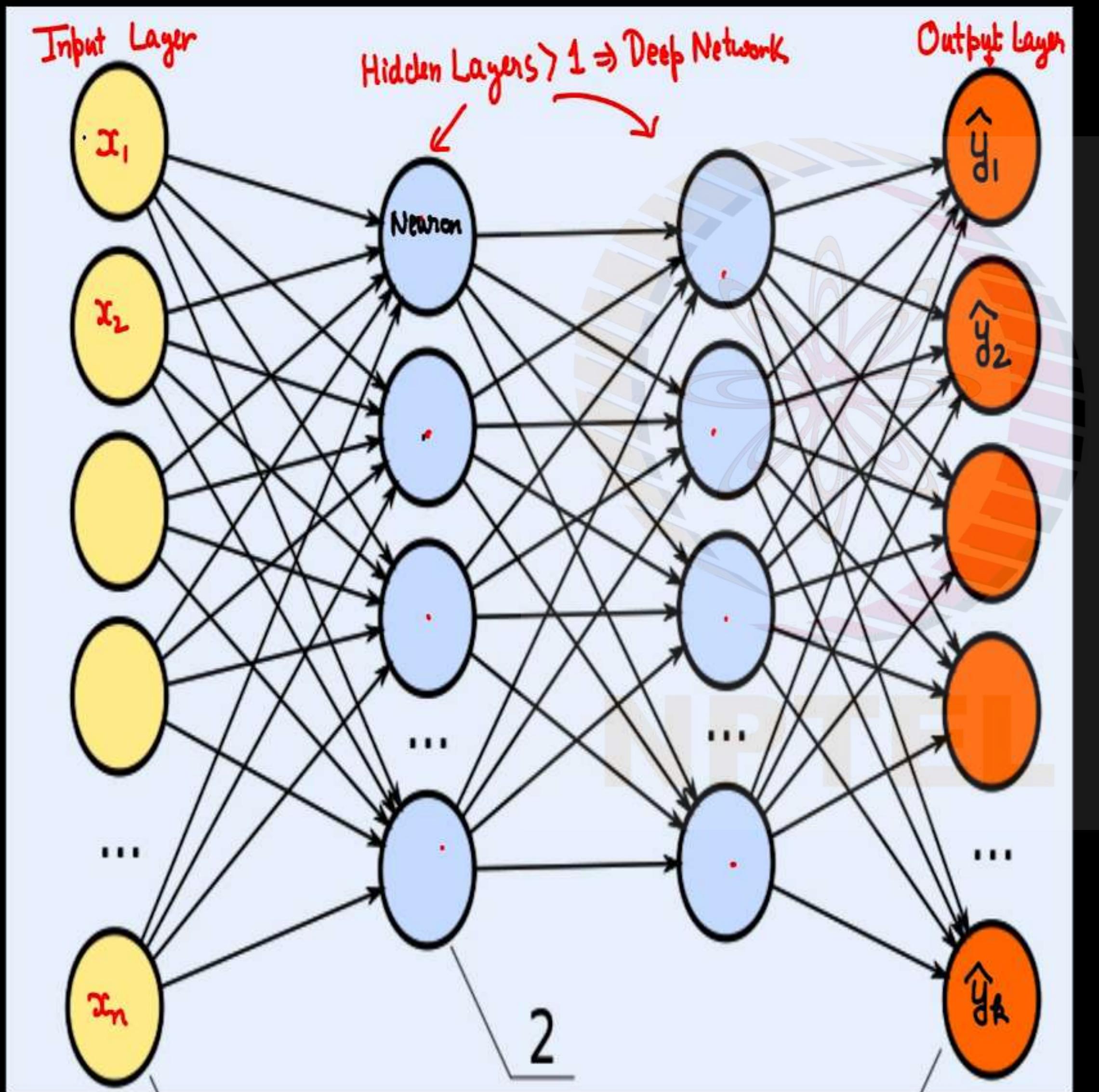


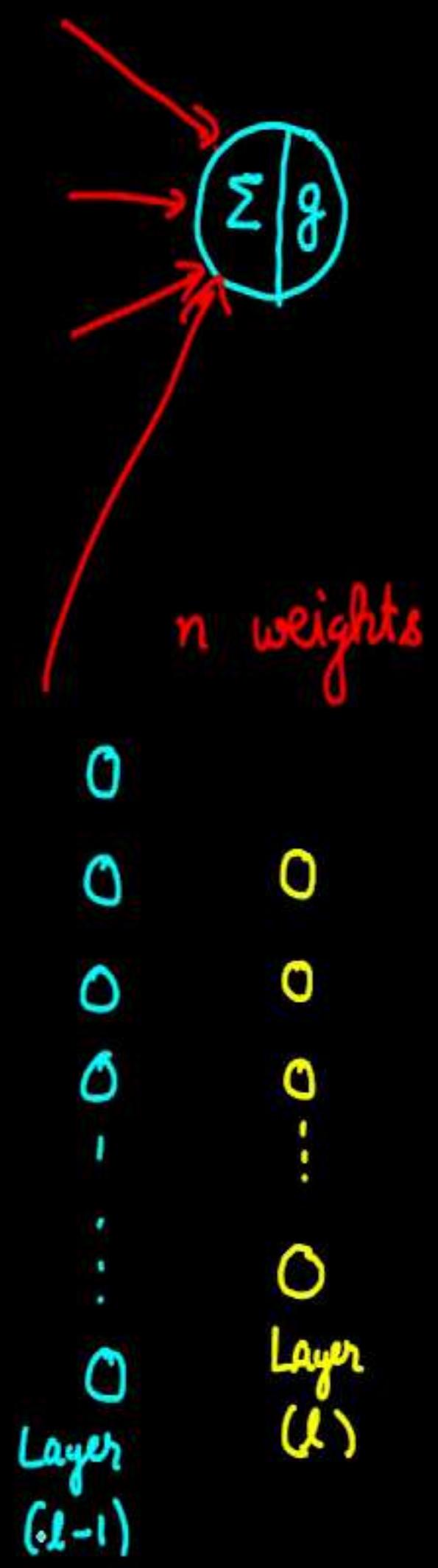
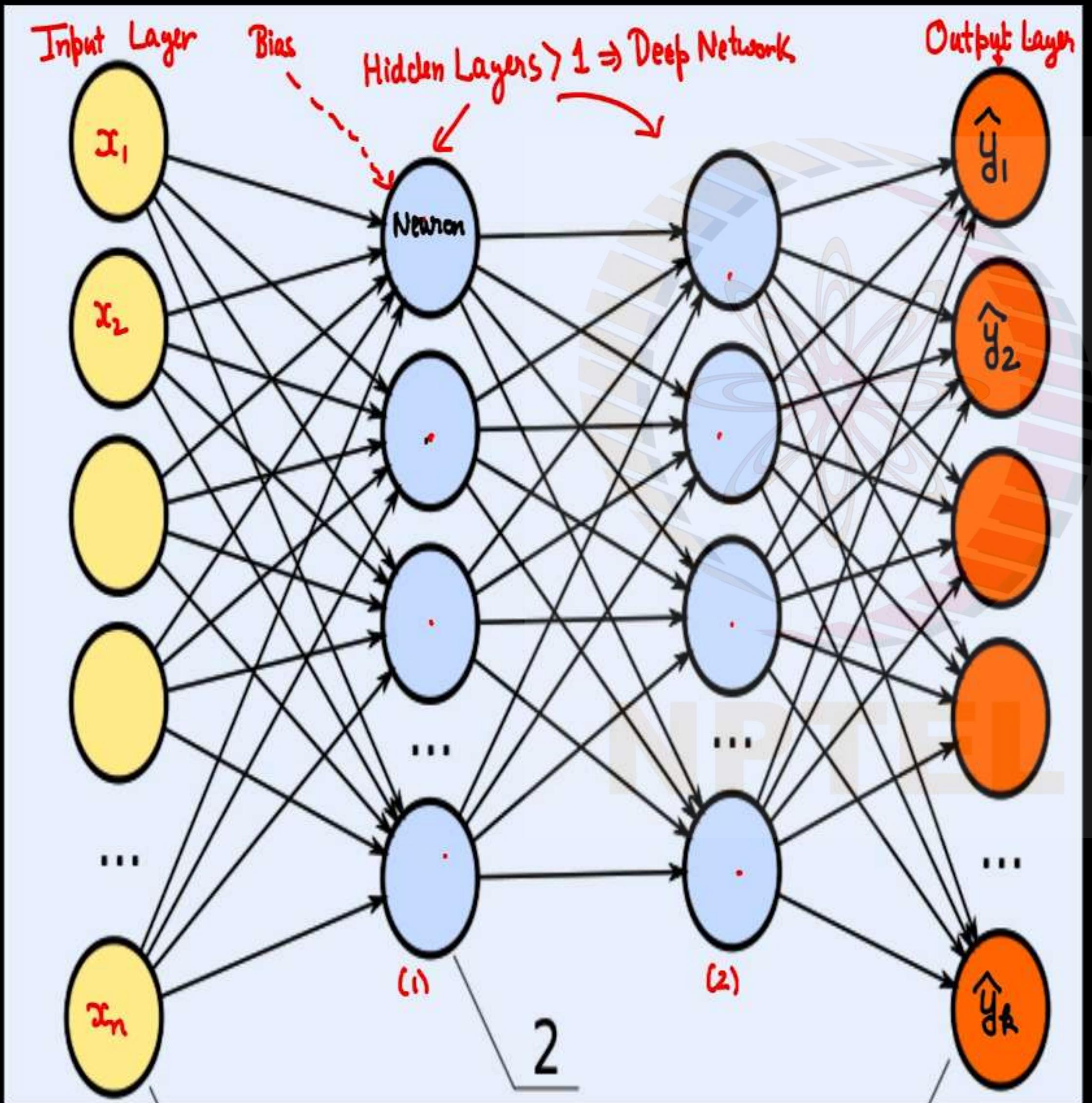
# Feedforward Network

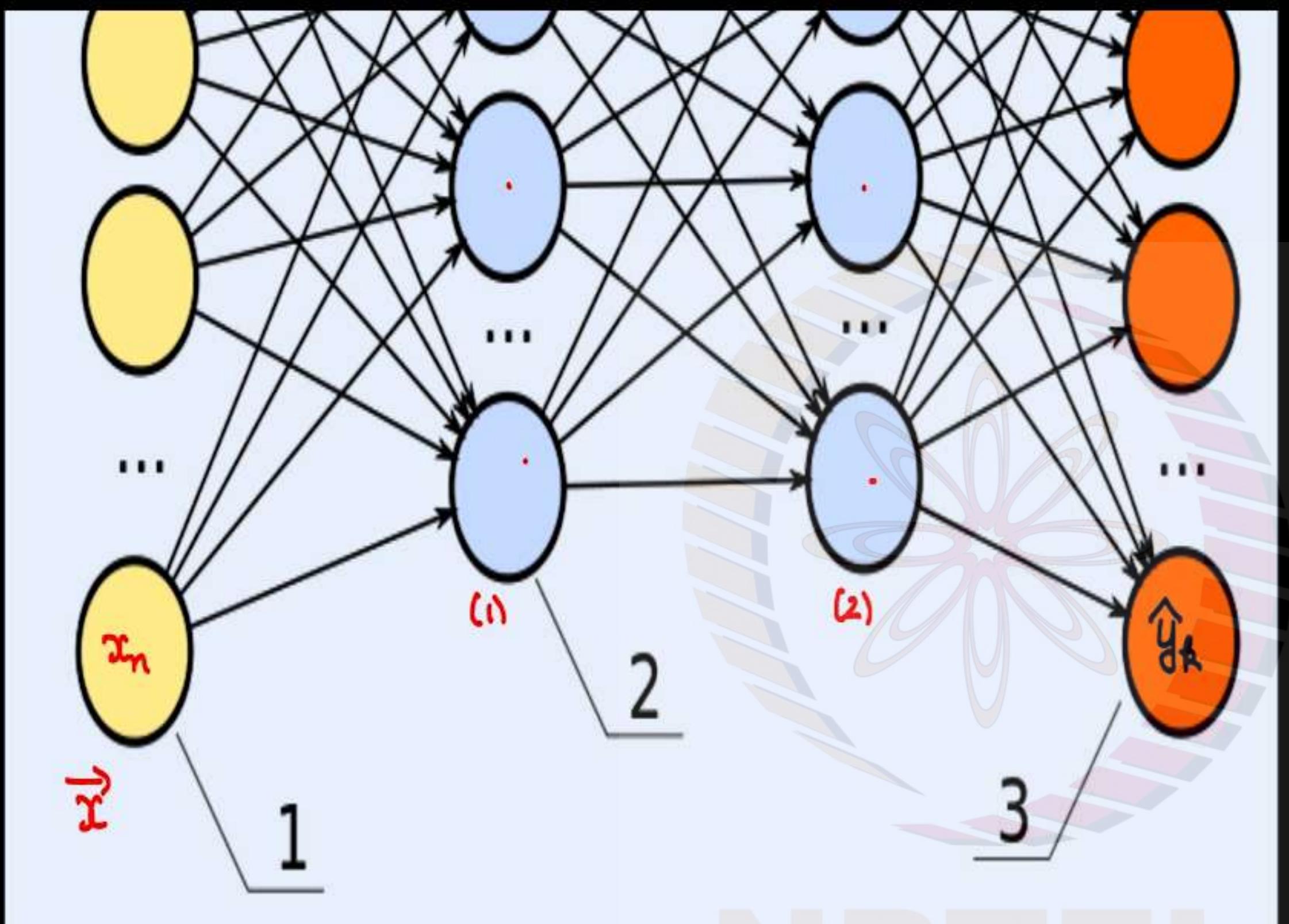
Neural Network

Note Title





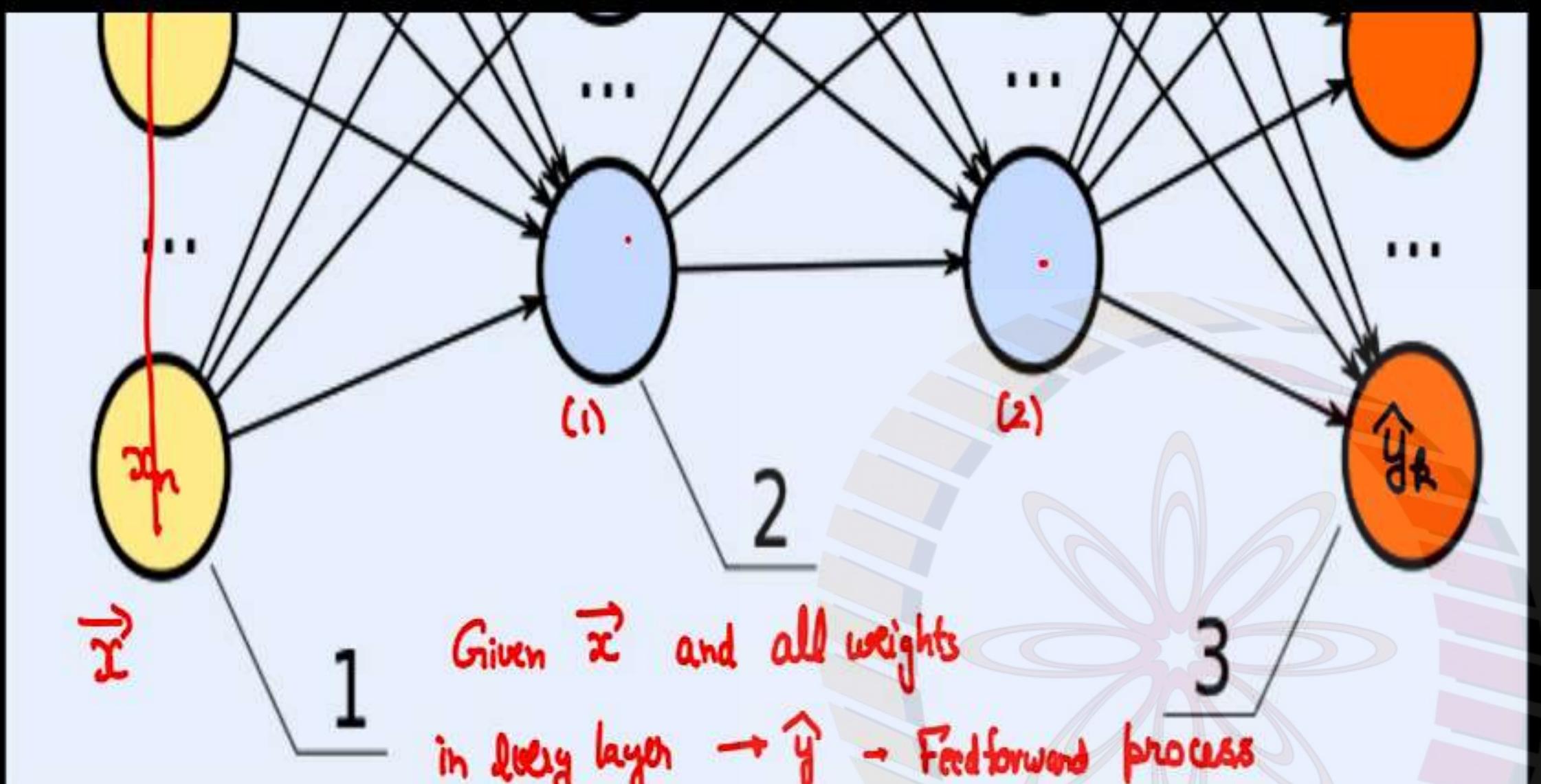




$i$	$n$ weights
0	0
0	0
0	0
0	0
1	:
⋮	⋮
0	Layer $(l)$
0	Layer $(l-1)$
$N$	$M$
	neurons

⇒ Total weights required

assuming every neuron is connected to every other neuron =  $N \times M + \text{Bias}$

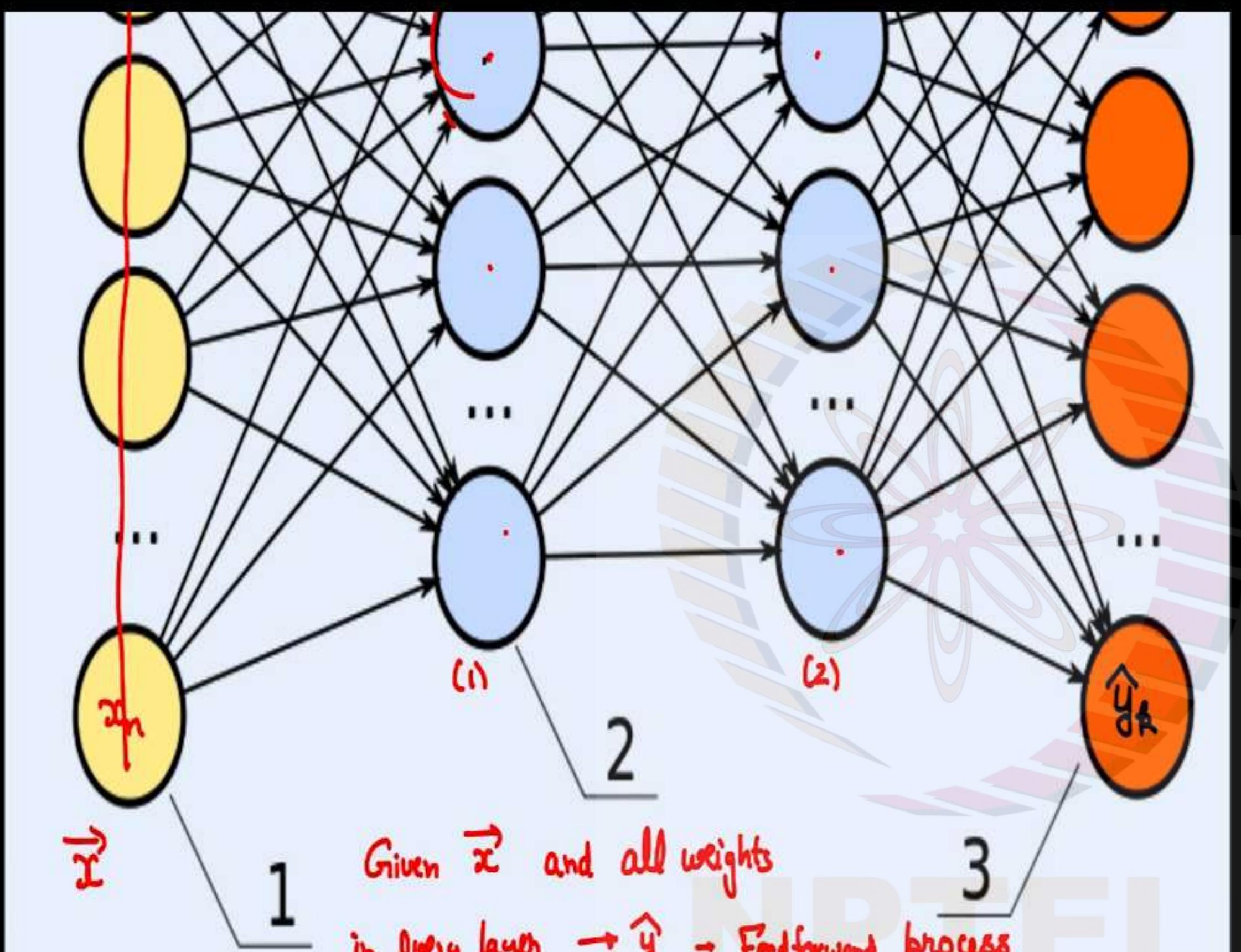


0	0-
0	0-
1	:
⋮	⋮
0-1	Layer (l)
0	Layer (l-1)
N	M
neurons	neurons

⇒ Total weights required

Fully Connected network { assuming every neuron is connected to every other neuron =  $N \times M + \text{Bias} \rightarrow M$  =  $N \times N + M$

NPTEL



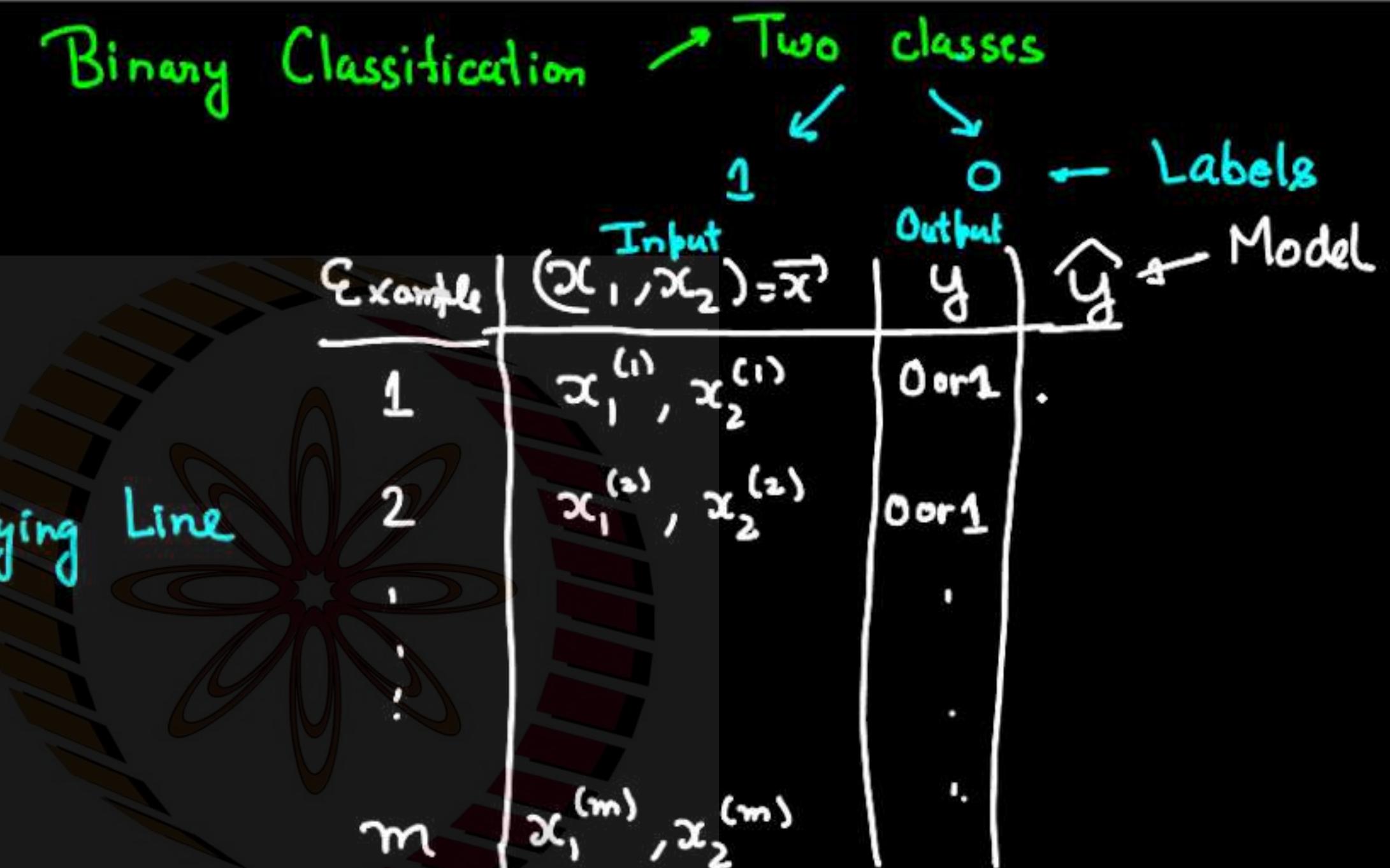
/ *n weights*  
 0 0 0- 0-  
 0 0 0- 0-  
 0 0 0- 0-  
 0 0 0- 0-  
 : : : :  
 0- 1 Layer (l)  
 0 Layer (l-1)  
 N M  
 neurons neurons

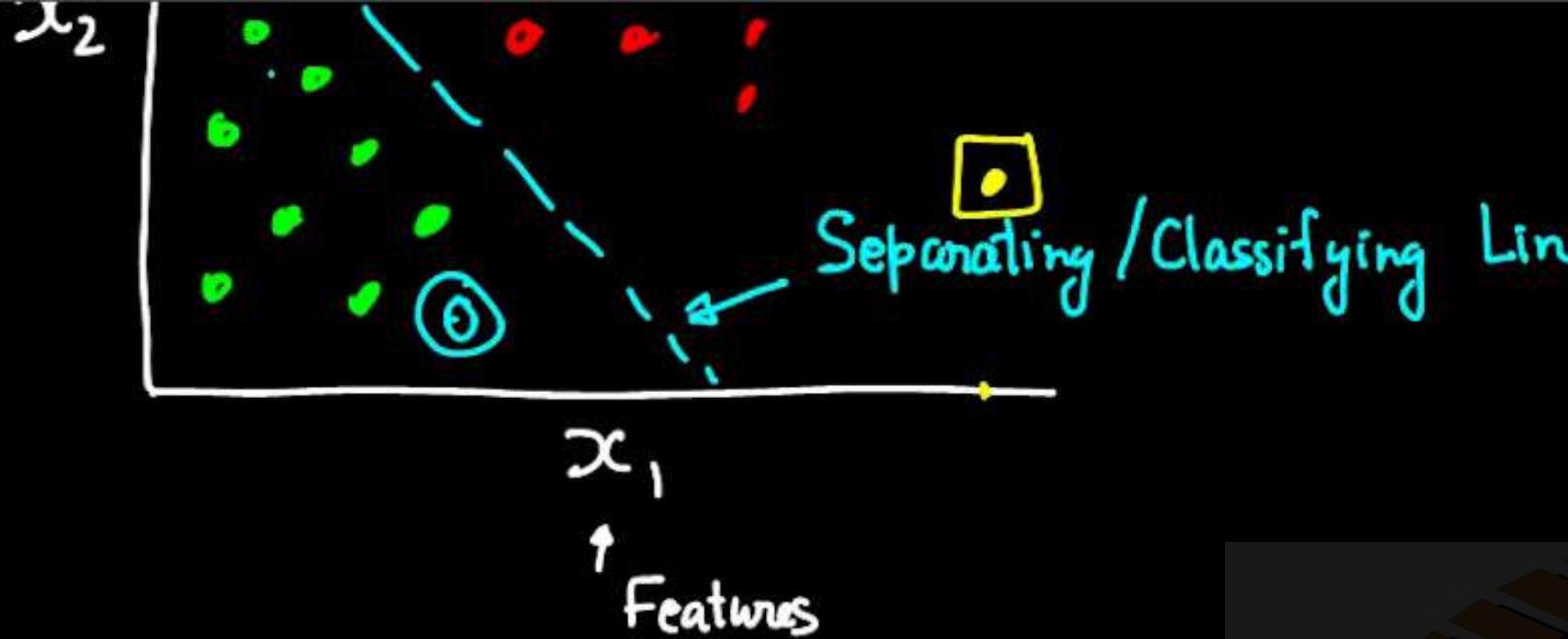
⇒ Total weights required

Fully  
Connected networks  
(FC)

{ assuming every neuron is connected to fully  
 other neuron =  $N \times M + \text{Bias} \rightarrow M$   
 =  $N \times N + M$

# Note Title





1	$x_1^{(1)}, x_2^{(1)}$	0 or 1
2	$x_1^{(2)}, x_2^{(2)}$	0 or 1
.	.	.
$m$	$x_1^{(m)}, x_2^{(m)}$	.

If we try simple linear regression

$\hat{y} = w_0 + w_1 x_1 + w_2 x_2$

Will not lie between 0 and 1

# NPTEL

$$\hat{y} = \omega_0 + \omega_1 x_1 + \omega_2 x_2$$

↑

Will not lie between 0 and 1

Simple idea : "Squish" all data  $\rightarrow [0, 1]$

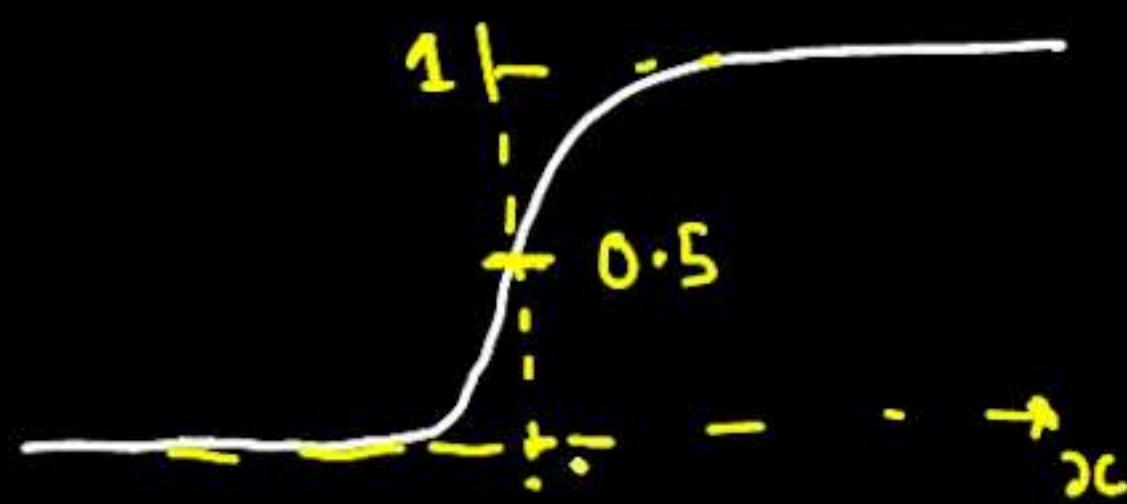
$\sigma \leftarrow$  Sigmoid function

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

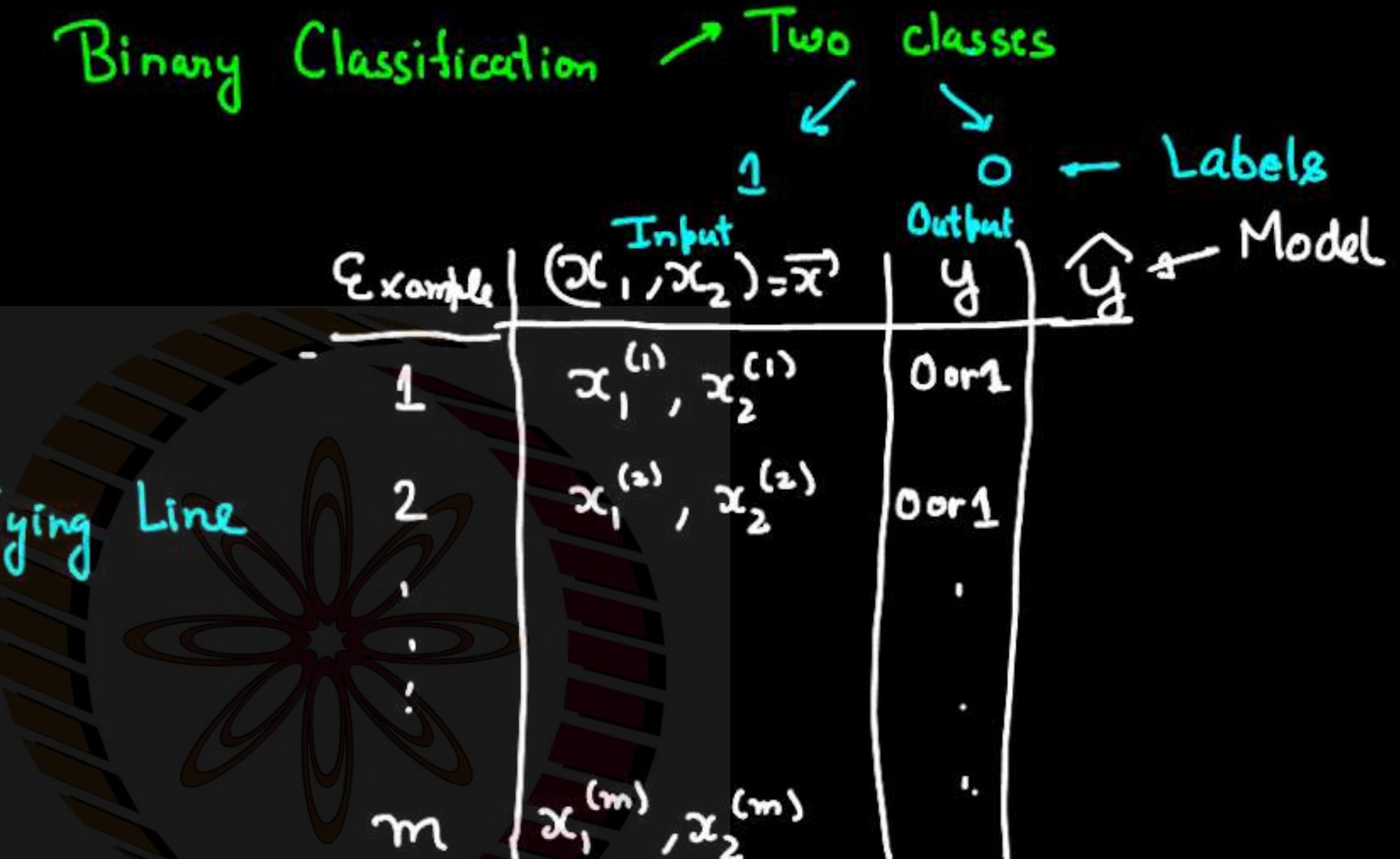
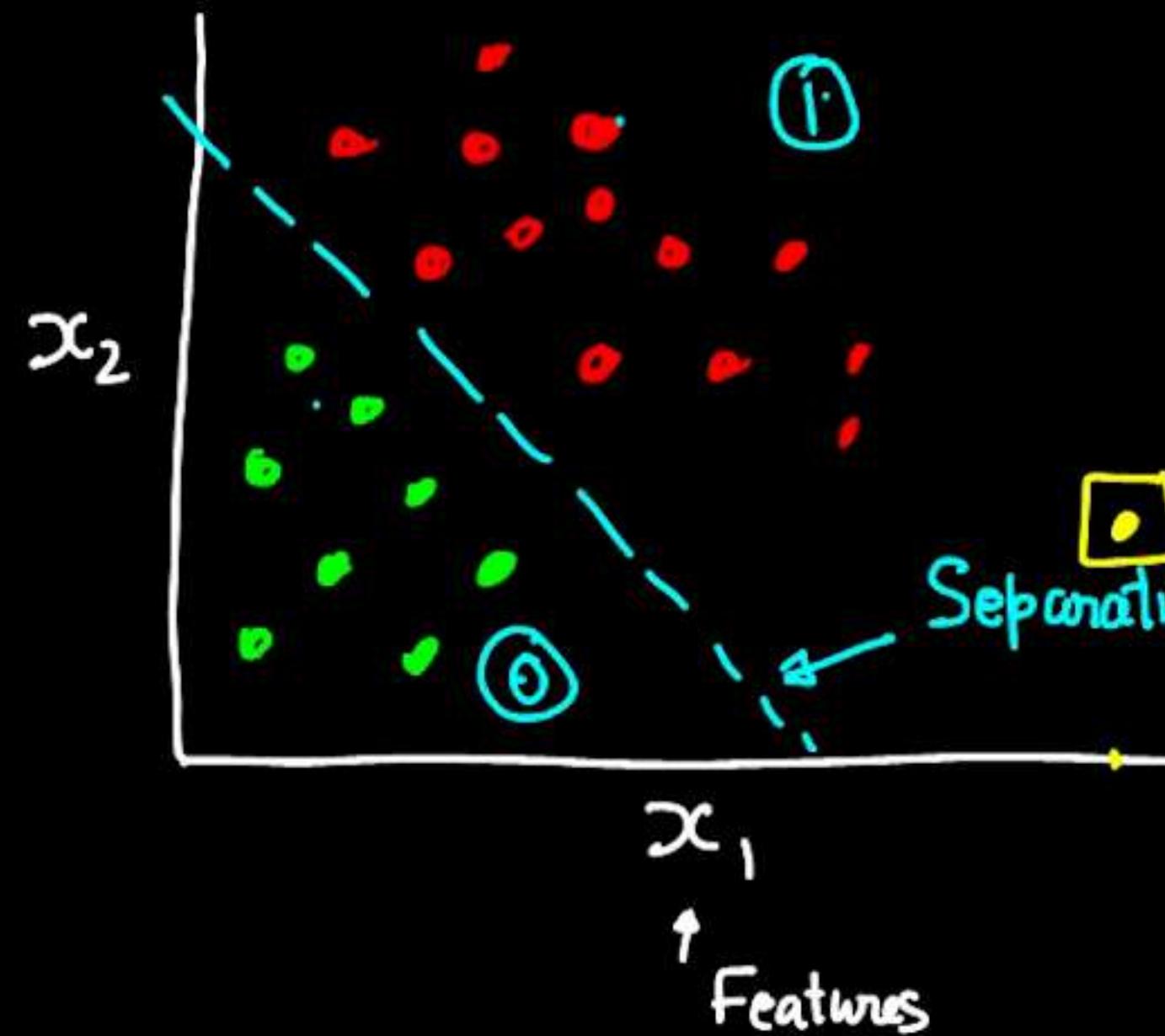
As  $z \rightarrow \infty$ ,  $\sigma(z) \rightarrow 1$

As  $z \rightarrow -\infty$ ,  $\sigma(z) \rightarrow 0$

At  $z = 0$ ,  $\sigma(z) = 0.5$



# Note Title



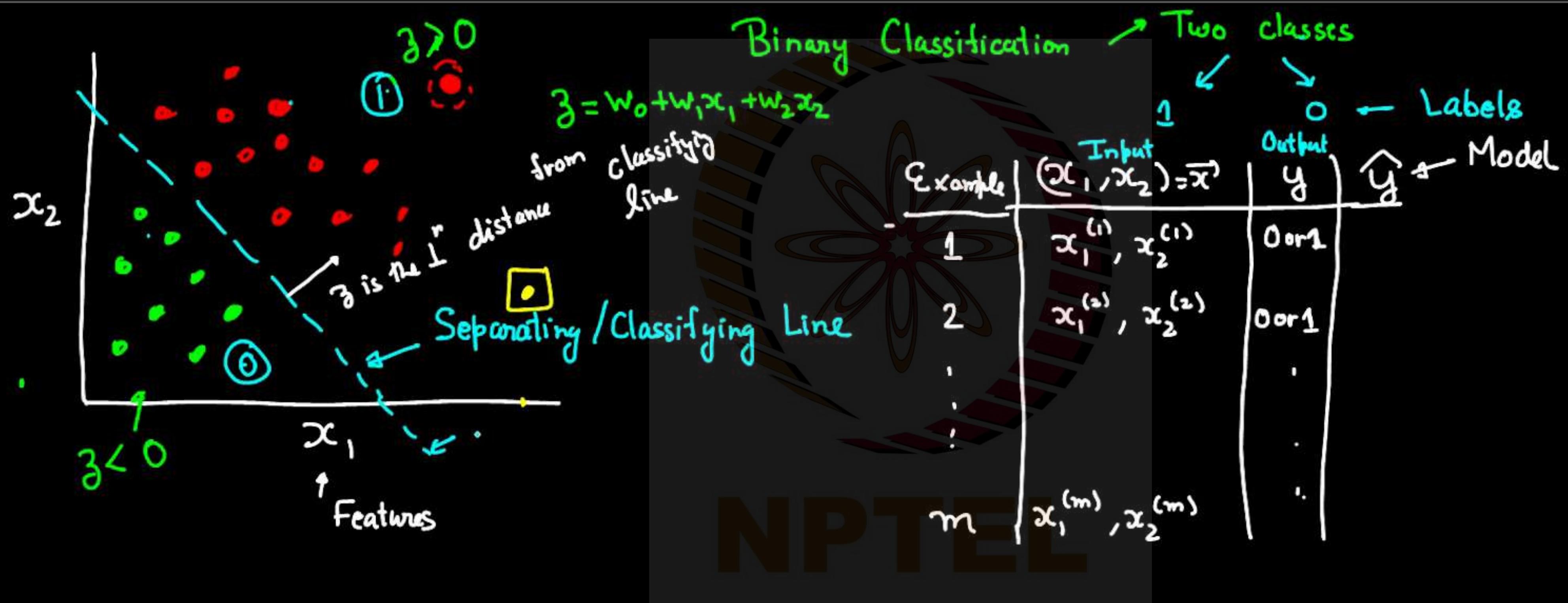
If we try simple linear regression

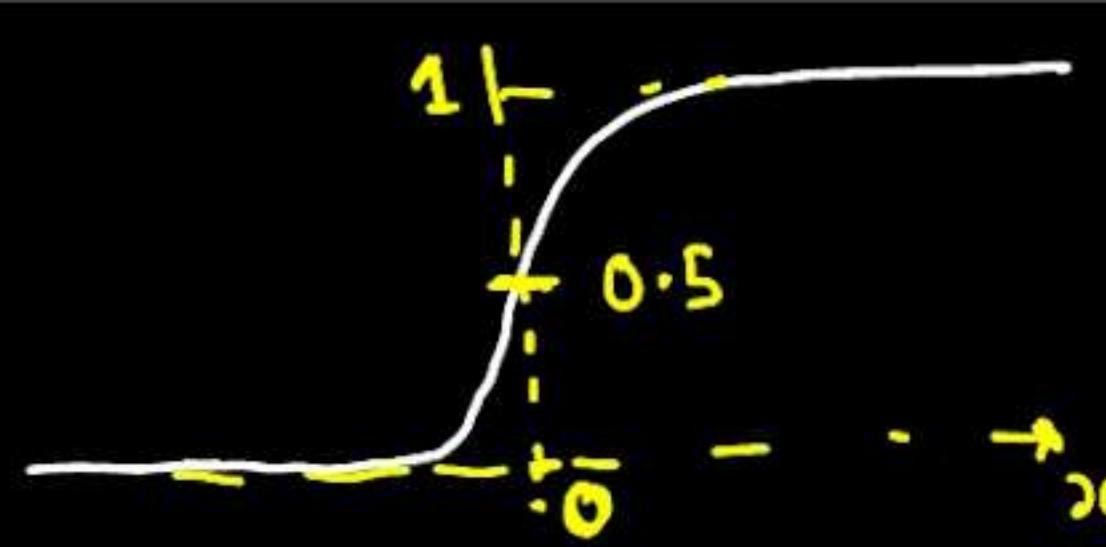
$$\hookrightarrow \hat{y} = w_0 + w_1 x_1 + w_2 x_2$$

# Lecture 10

## (Classification)

Note Title



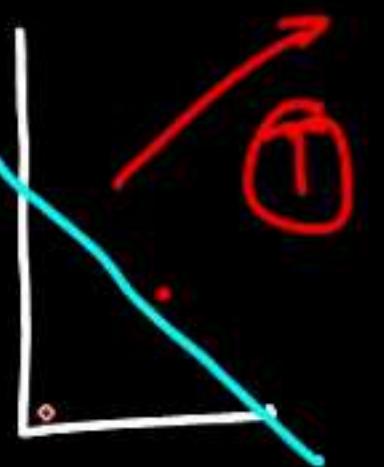


Monotonic function

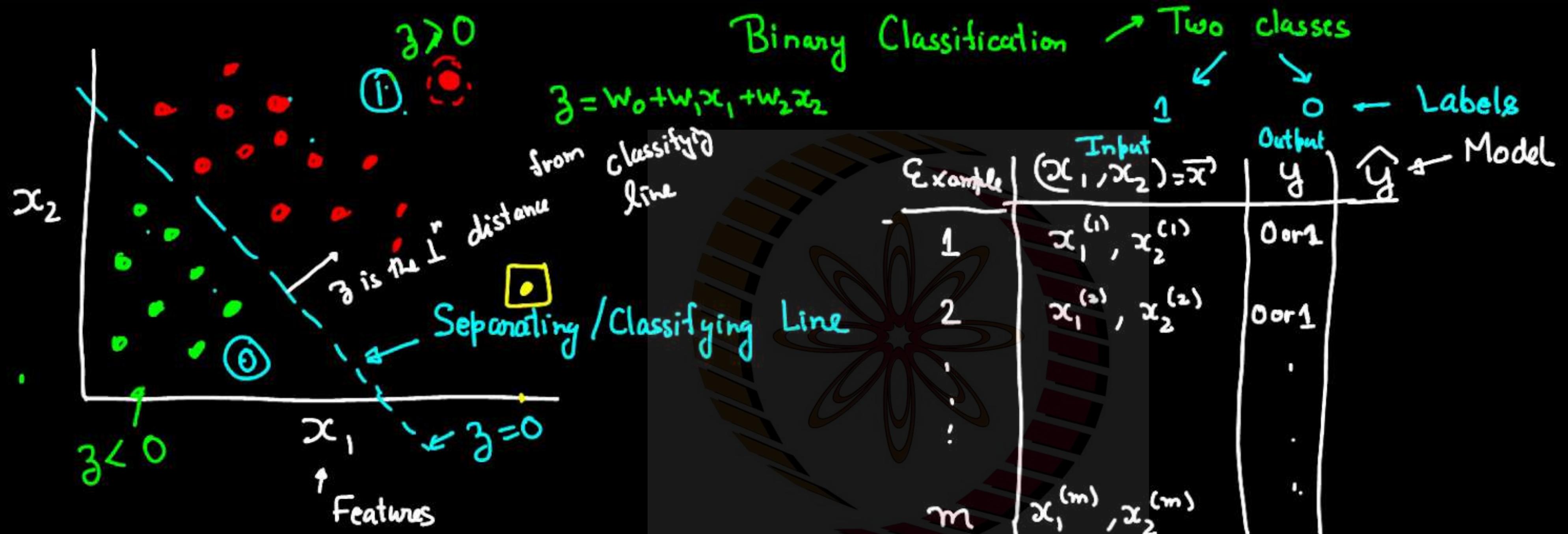
$$\hat{y} = \sigma(\underbrace{\omega_0 + \omega_1 x_1 + \omega_2 x_2}_{\beta}) \Rightarrow \hat{y} \in [0, 1]$$

Interpret  $\hat{y} = P(y=1|x)$

$$\hat{y} = \sigma(\beta) = P(y=1|x)$$



## Note Title



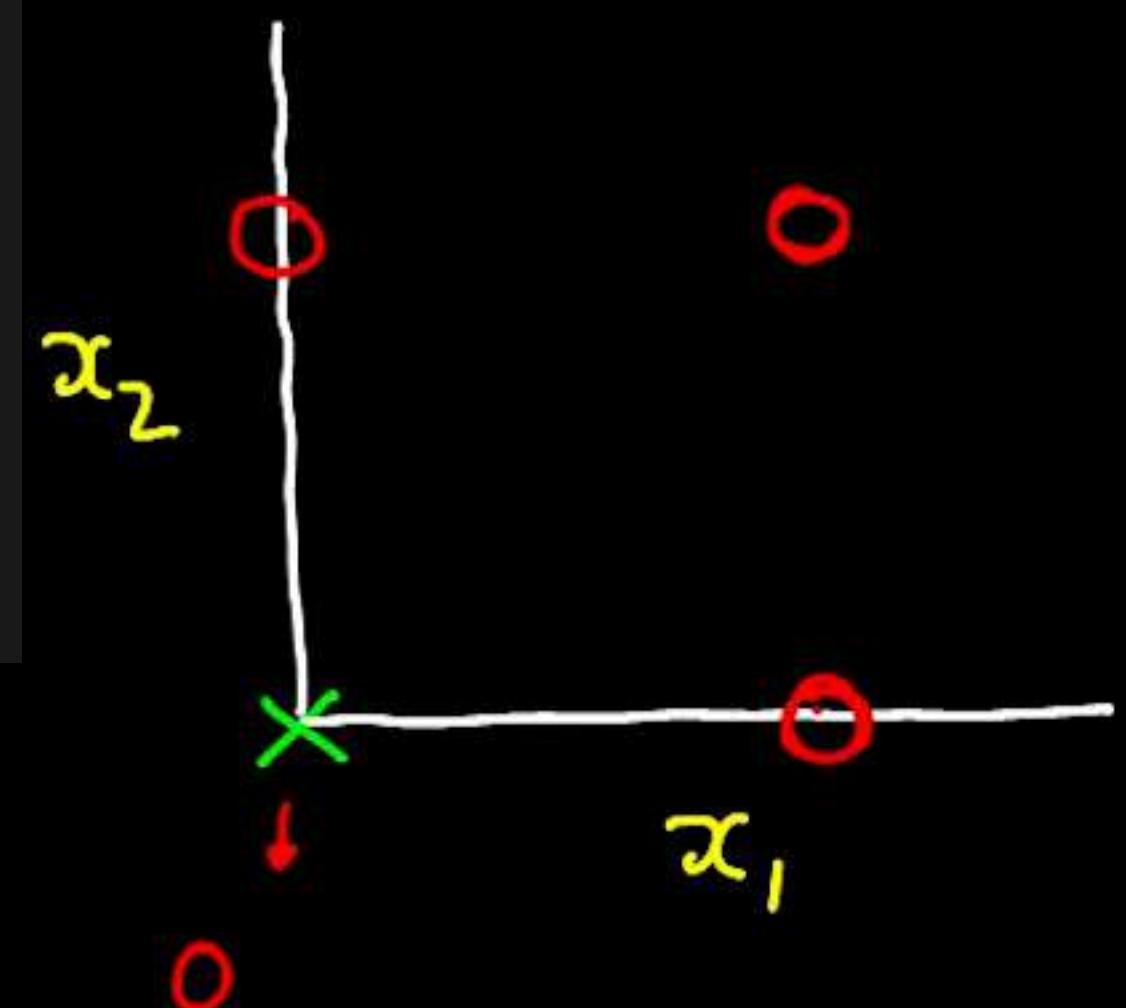
If we try simple linear regression

## Note Title

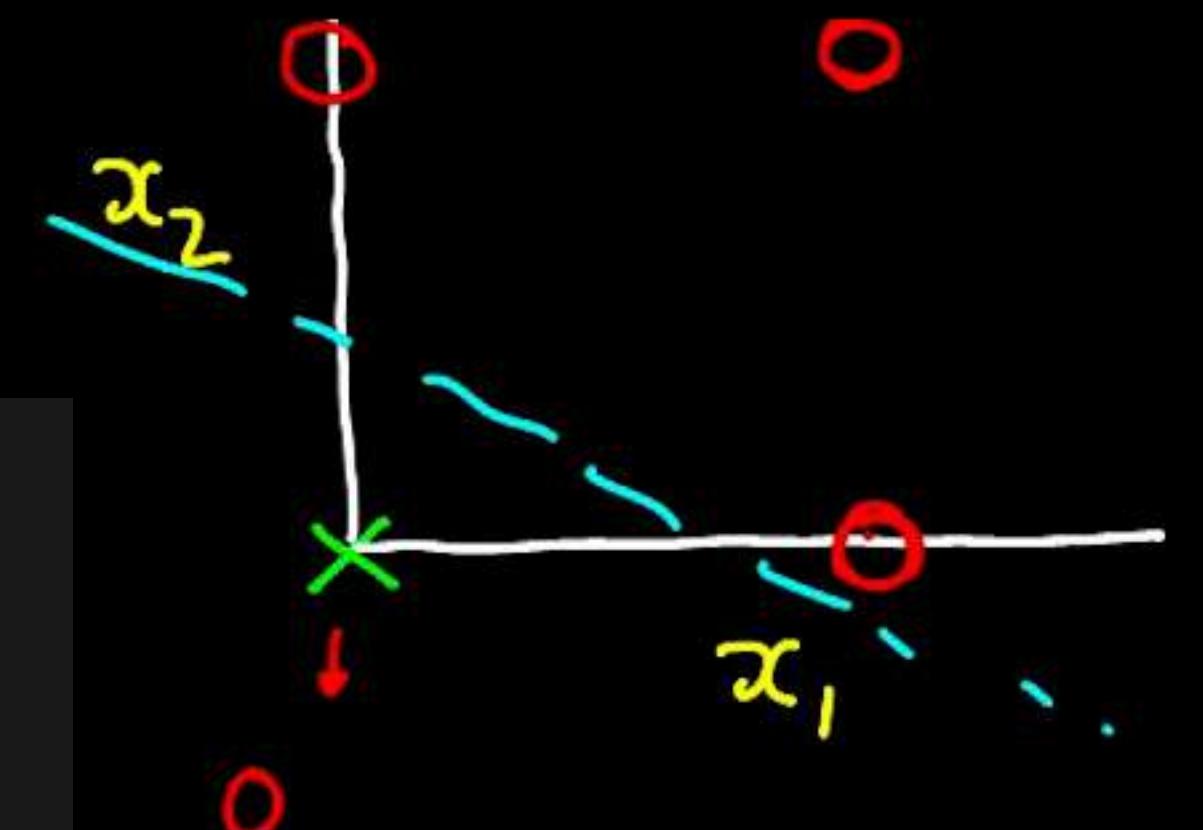
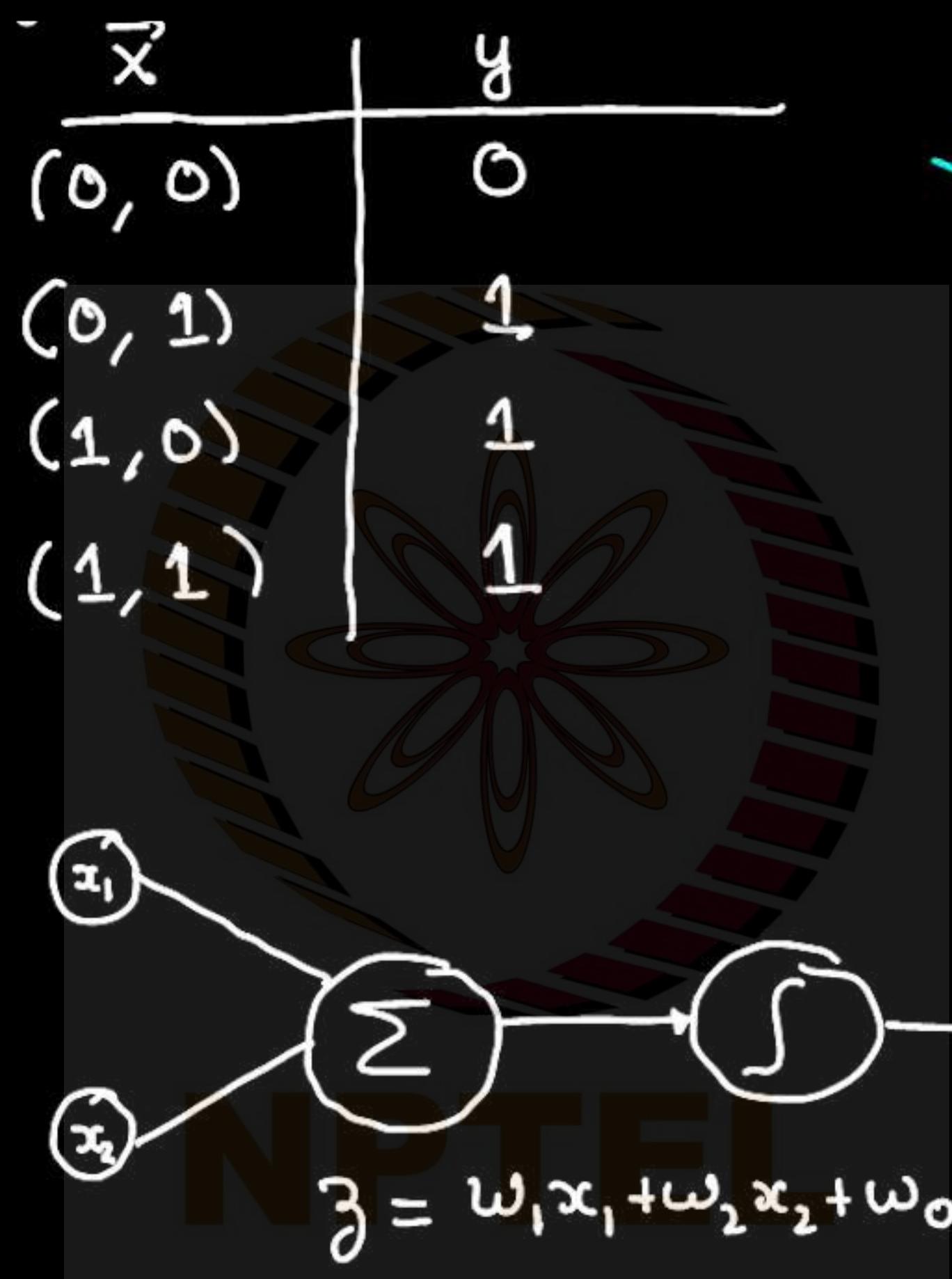
OR gate

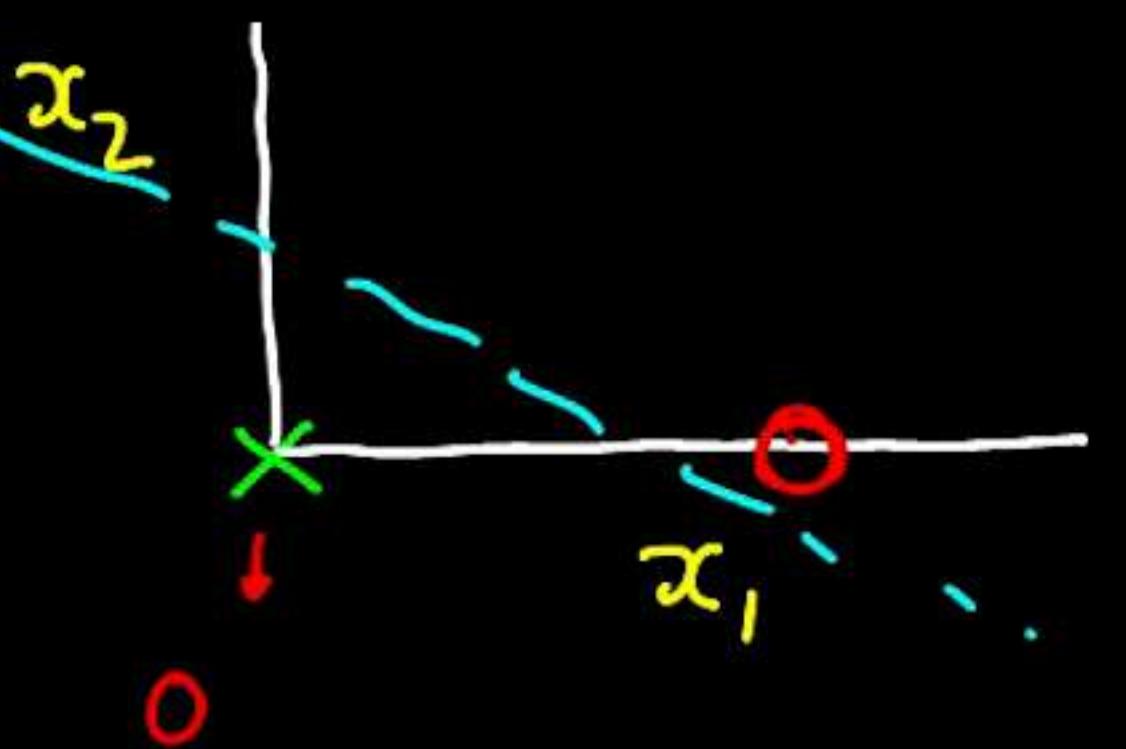
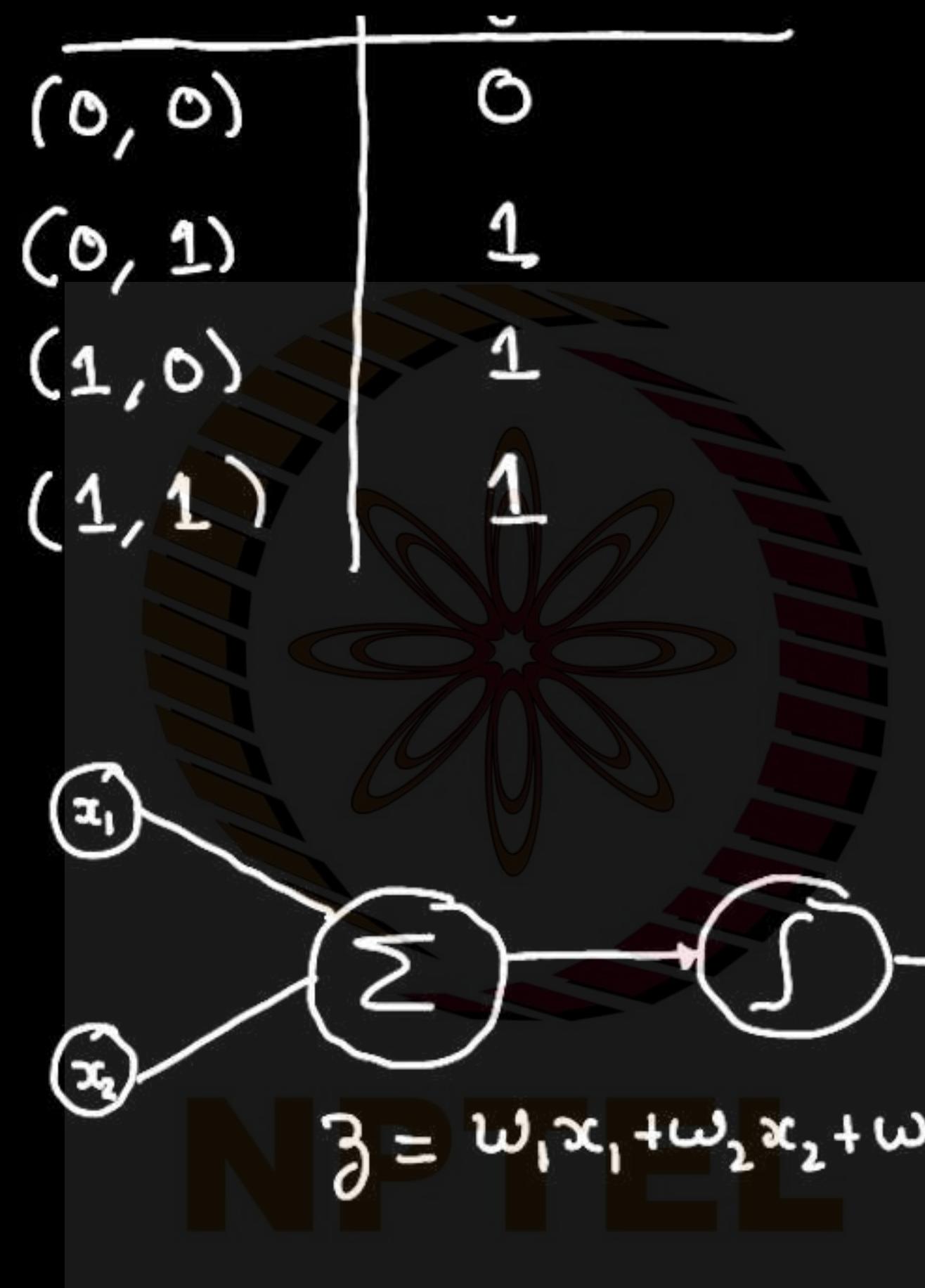
$OR(x_1, x_2)$

$\vec{x}$	y
(0, 0)	0
(0, 1)	1
(1, 0)	1
(1, 1)	1



$$\begin{cases} 0 \rightarrow 1 \\ x \rightarrow 0 \end{cases}$$



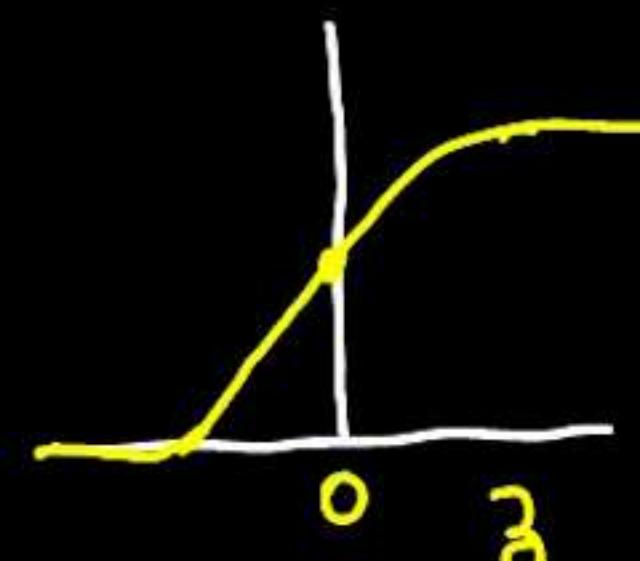


$$\begin{aligned} & 1 \text{ if } \hat{y} \geq 0.5 \\ & 0 \text{ if } \hat{y} < 0.5 \end{aligned}$$

What weights  $w_0, w_1, w_2$  ( $\vec{w}$ ) will classify same as OR gate

Note Title

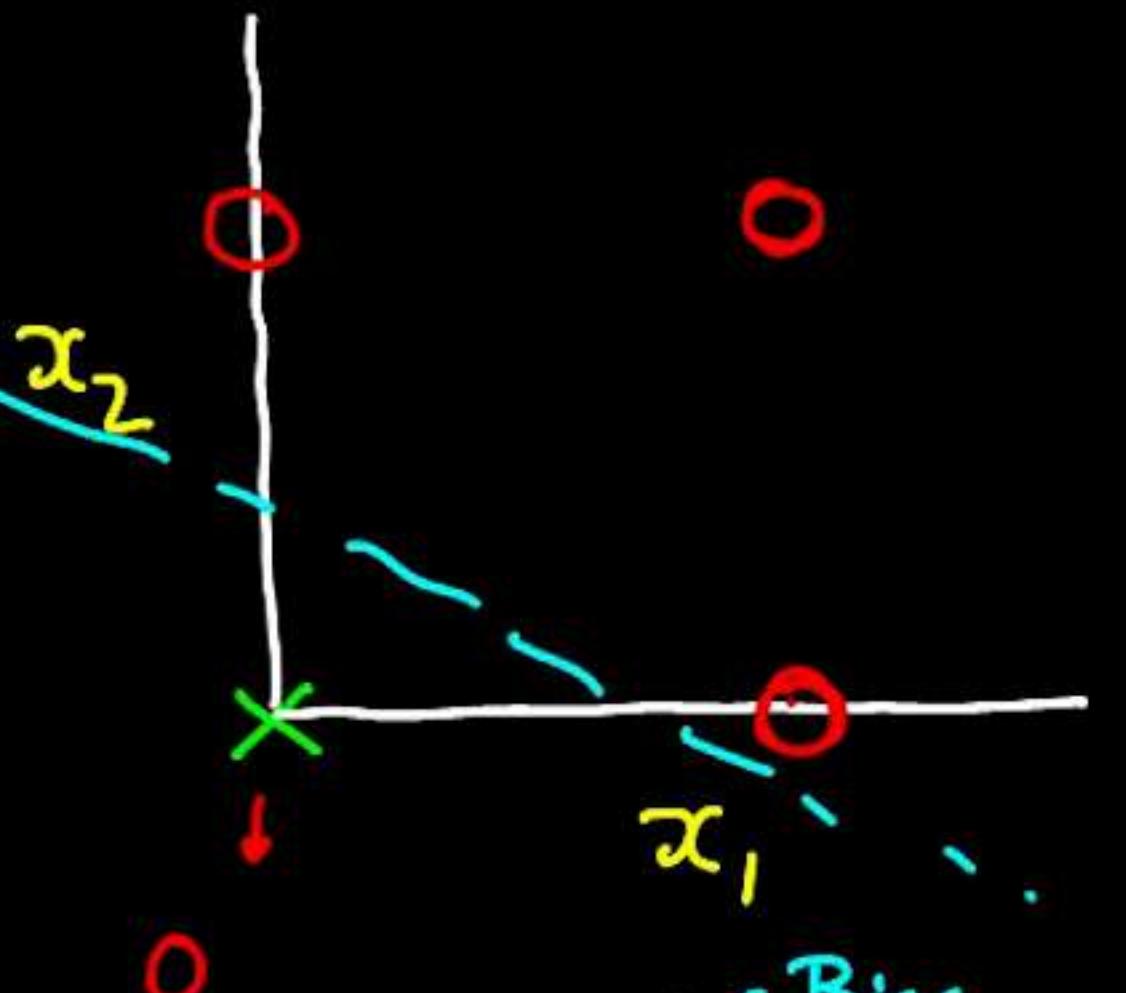
$$\boxed{\begin{aligned} w_0 &= -1 \\ w_1 &= 2 \\ w_2 &= 2 \end{aligned}}$$



OR gate

$OR(x_1, x_2)$

$\vec{x}$	$y$	$\beta$
$(0, 0)$	0	$-ve$
$(0, 1)$	1	$+ve$
$(1, 0)$	1	$+ve$
$(1, 1)$	1	$+ve$



$$\begin{cases} 0 \rightarrow 1 \\ x \rightarrow 0 \end{cases}$$

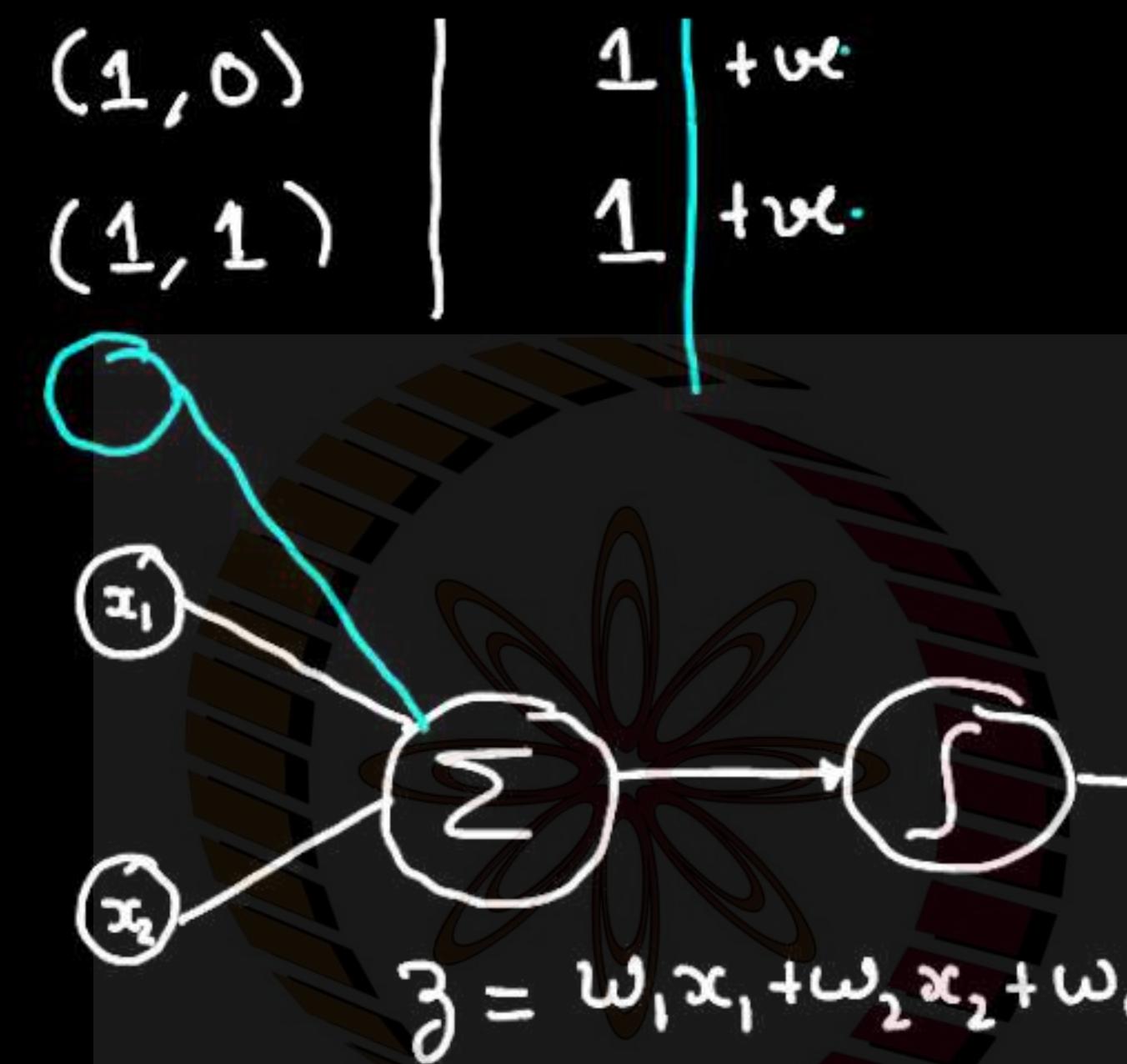
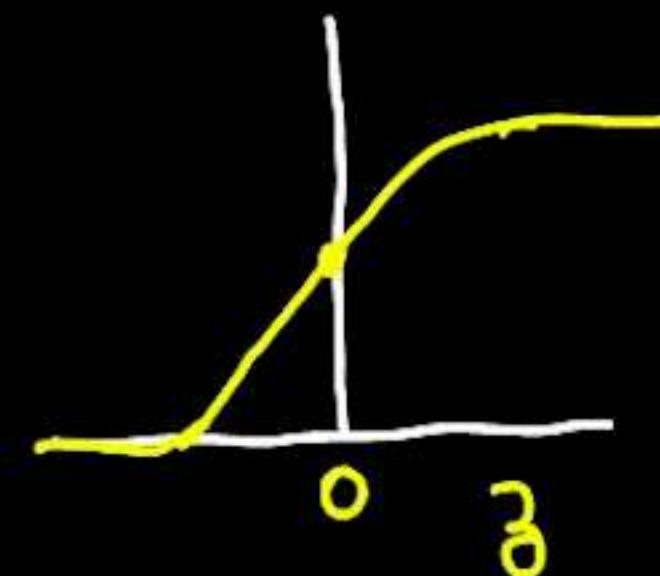
$$\begin{matrix} w_0 + w_2 > 0 \\ \downarrow \quad \downarrow \\ -1 \quad 2 \end{matrix}$$

$$\beta = w_0 + w_1 x_1 + w_2 x_2$$

1 if  $\hat{y} \geq 0.5$

0 if  $\hat{y} < 0.5$

$$\begin{cases} w_2 = -1 \\ w_1 = 2 \end{cases}$$



Inputs:  $x_1$ ,  $x_2$

Output:  $\hat{y}$

$\beta = w_0 + w_1 x_1 + w_2 x_2$

$\hat{y} = \sigma(\beta)$

Decision rule:

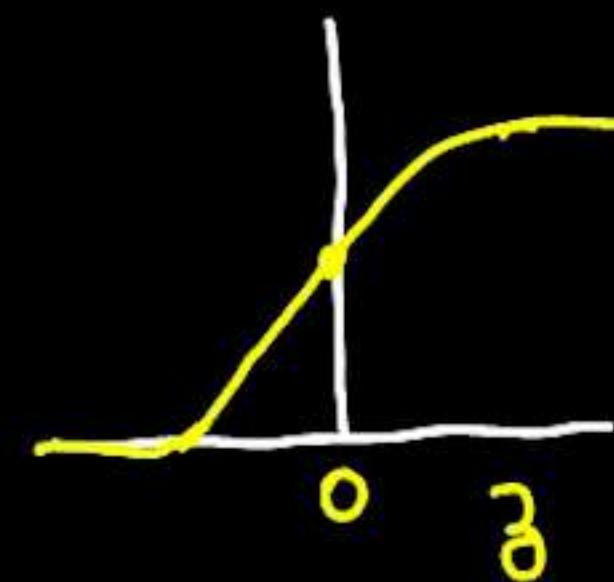
- $1$  if  $\hat{y} \geq 0.5$
- $0$  if  $\hat{y} < 0.5$

What weights  $w_0, w_1, w_2$  ( $\vec{w}$ ) will classify same as OR gate

- NOTE : (1)  $w_0$  is essential  $\rightarrow$  Need for bias unit
- (2)  $w_0, w_1, w_2$  are not unique

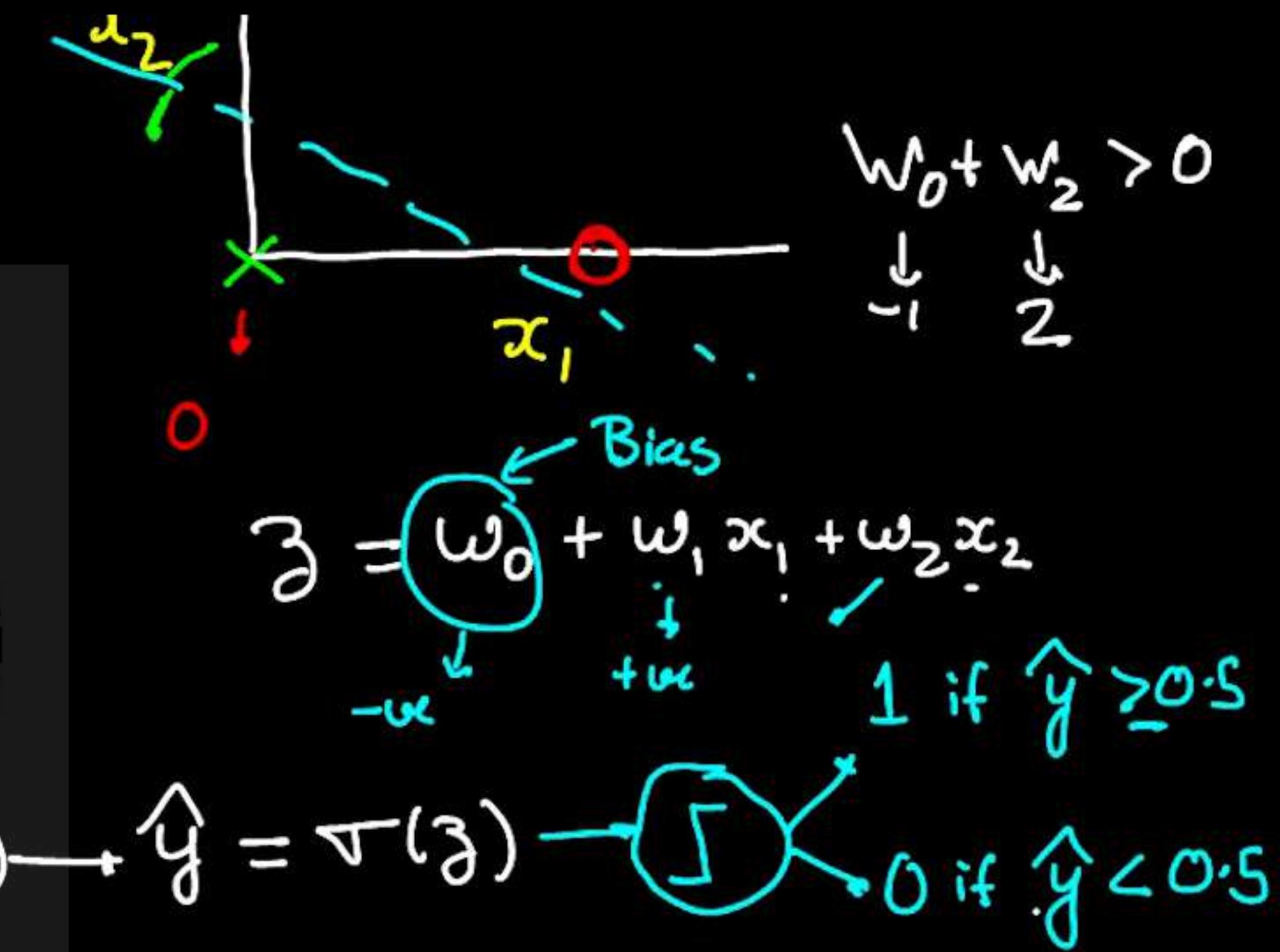
↳ The weights of logistic regression

$$\begin{cases} w_0 = -1 \\ w_1 = 2 \\ w_2 = 2 \end{cases}$$



$(0, 0)$	$0$	$-ve$
$(0, 1)$	$1$	$+ve$
$(1, 0)$	$1$	$+ve$
$(1, 1)$	$1$	$+ve$

$\beta = w_1 x_1 + w_2 x_2 + w_0$



What weights  $w_0, w_1, w_2$  ( $\vec{w}$ ) will classify same as OR gate

NOTE : (1)  $w_0$  is essential  $\rightarrow$  Need for bias unit

— 0 3

(x<sub>2</sub>)

$$\vec{z} = w_1 x_1 + w_2 x_2 + w_0$$

What weights  $w_0, w_1, w_2$  ( $\vec{w}$ ) will classify same as OR gate

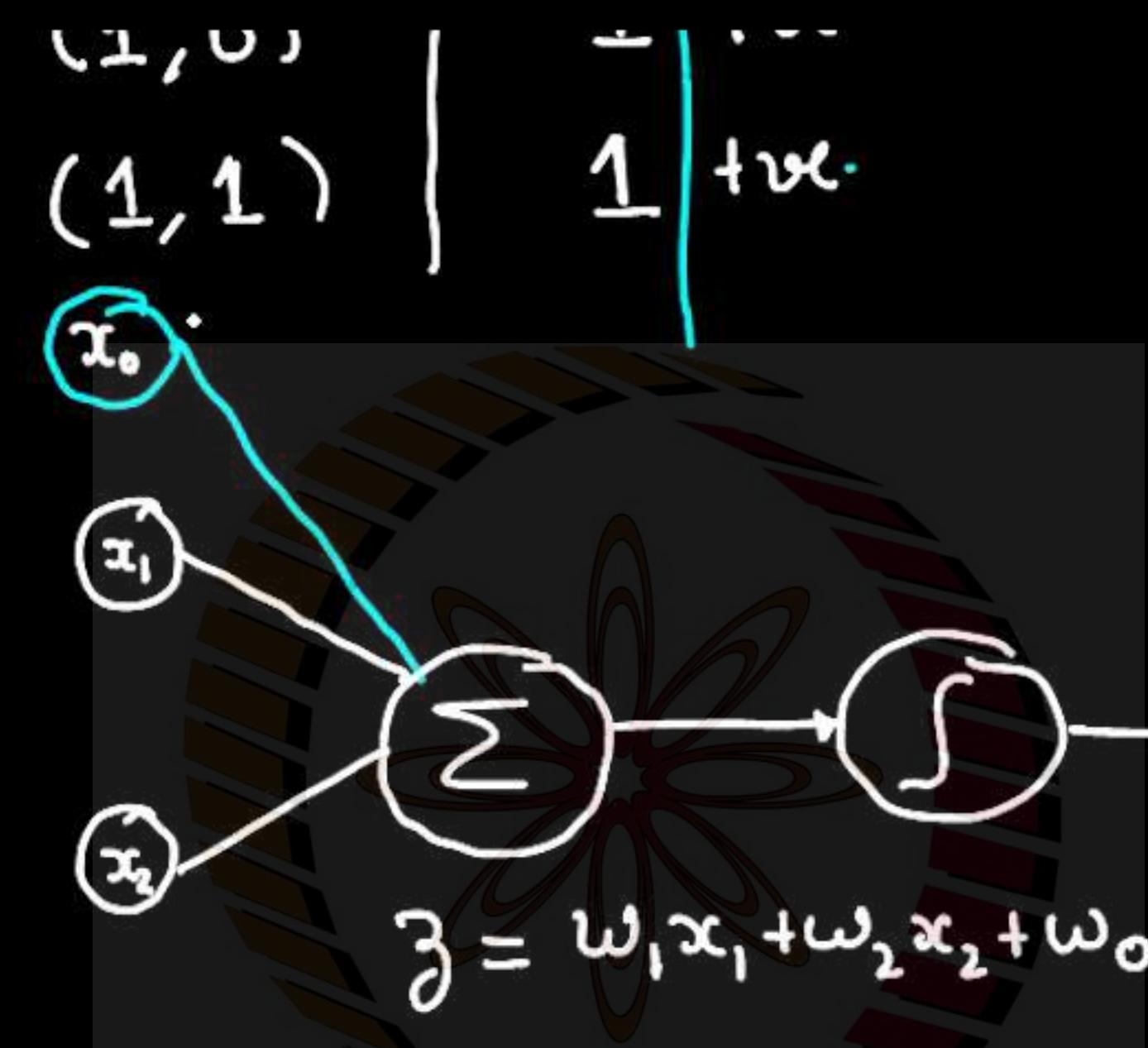
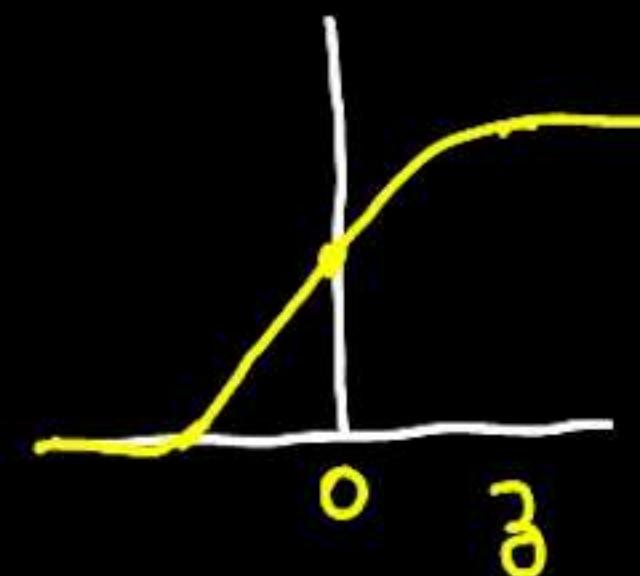
NOTE : (1)  $w_0$  is essential  $\rightarrow$  Need for bias unit

(2)  $w_0, w_1, w_2$  are not unique

↳ The results of logistic regression  
need not be unique

$$w_0 + w_1 x_1 + w_2 x_2 = \vec{z} .$$

$$\boxed{w_1 = 2}$$



$\hat{\beta} = w_0 + w_1 x_1 + w_2 x_2$

$\hat{y} = \sigma(\hat{\beta})$

$\begin{cases} 1 & \text{if } \hat{y} \geq 0.5 \\ 0 & \text{if } \hat{y} < 0.5 \end{cases}$

What weights  $w_0, w_1, w_2$  ( $\vec{w}$ ) will classify same as OR gate

NOTE : (1)  $w_0$  is essential  $\rightarrow$  Need for bias unit

(2)  $w_0, w_1, w_2$  are not unique

↳ The results of logistic regression

## Note Title

NOR, AND, NAND gates

$$\text{NOR} = \text{NOT(OR)}$$

$x_1, x_2$	$y$
0, 0	1
0, 1	0
1, 0	0
1, 1	0

$$\underbrace{w_0 = -1, w_1 = 2, w_2 = 2}_{\text{OR gate}}$$

## Note Title

$$\text{NOR} = \text{NOT(OR)}$$

AND

$x_1, x_2$	$y$
0, 0	0
0, 1	0
1, 0	0
1, 1	1

$x_1, x_2$	$y$
0, 0	1
0, 1	0
1, 0	0
1, 1	0

$$w_0 = -1, w_1 = 2, w_2 = 2$$

OR gate

$$w_0 = 1, w_1 = -2, w_2 = -2$$

NOR gate

$$w_0 = -3, w_1 = 2, w_2 = 2$$

AND gate

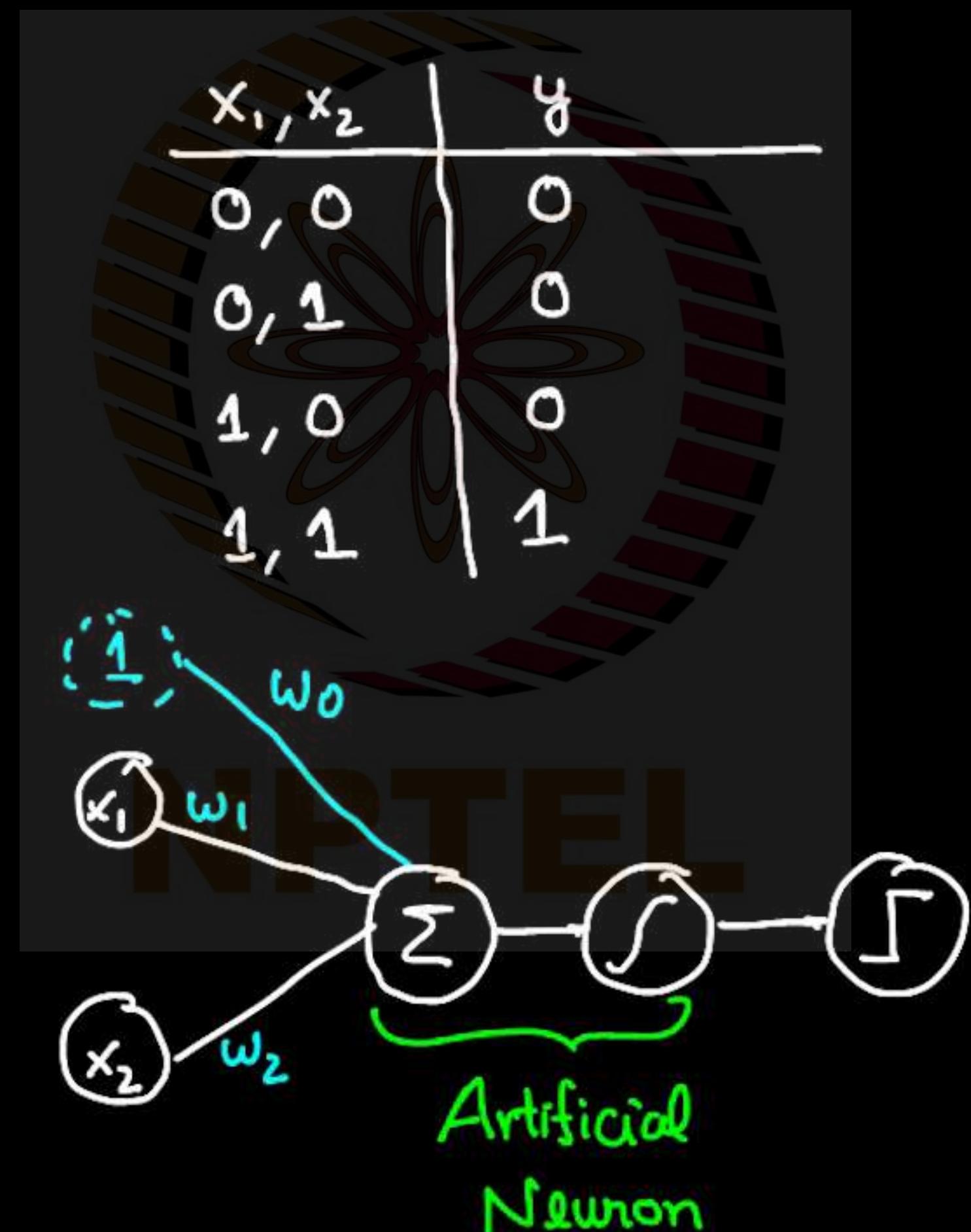
$$w_0 = 3, w_1 = -2, w_2 = -2$$

NAND gate

1, 0		0
1, 1		0

$w_0 = 1, w_1 = -4, w_2 = -4$   
 NOR gate

AND

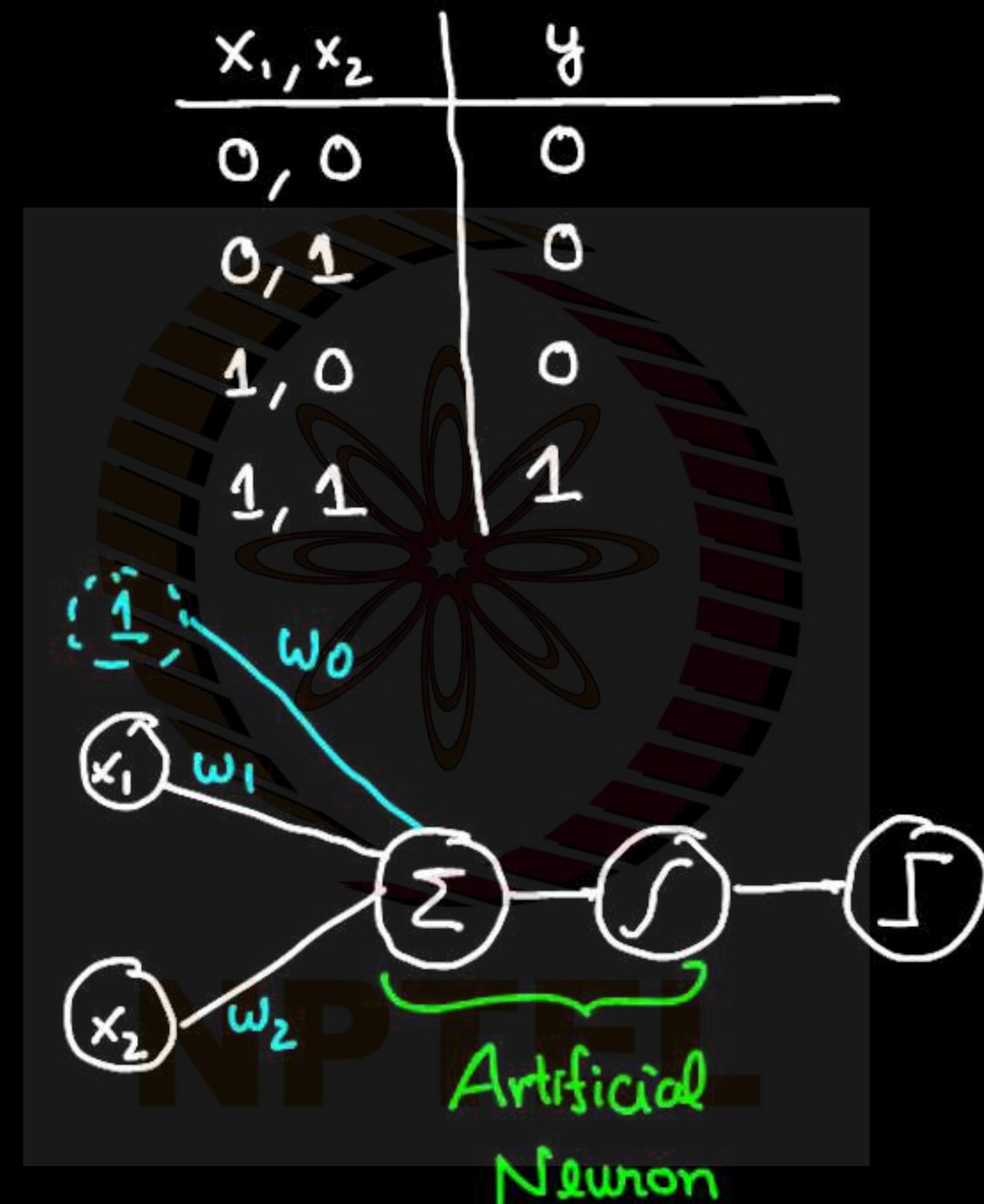


$w_0 = -3, w_1 = 2, w_2 = 2$   
 AND gate

$w_0 = 3, w_1 = -2, w_2 = -2$   
 NAND gate

Can all logic circuits/gates  
 be represented by this  
 simple network architecture?

AND



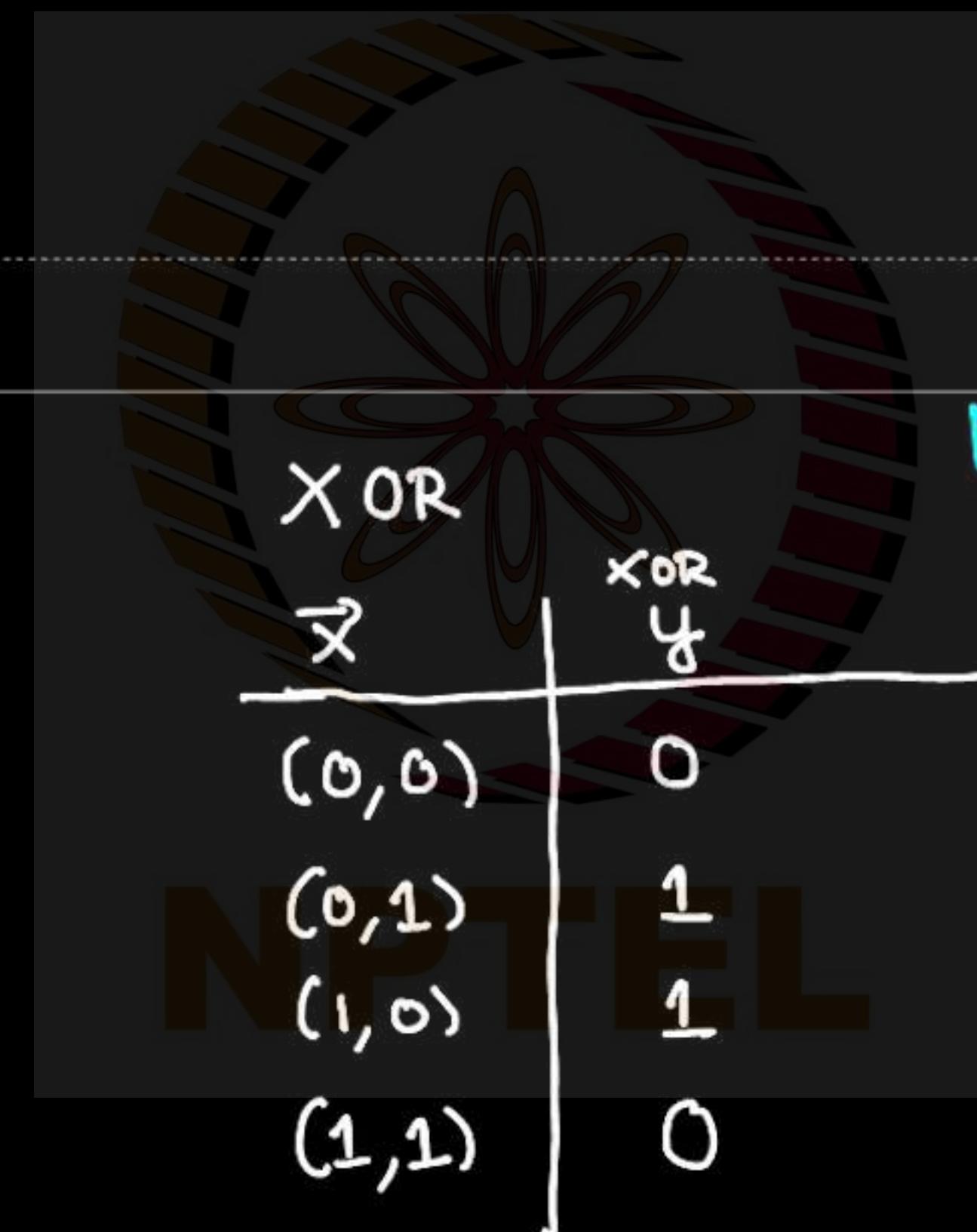
$$\underbrace{w_0 = -3, w_1 = 2, w_2 = 2}_{\text{AND gate}}$$

$$\underbrace{w_0 = 3, w_1 = -2, w_2 = -2}_{\text{NAND gate}}$$

Can all logic circuits/gates  
be represented by this  
Simple network architecture?

## Note Title

### XOR GATE

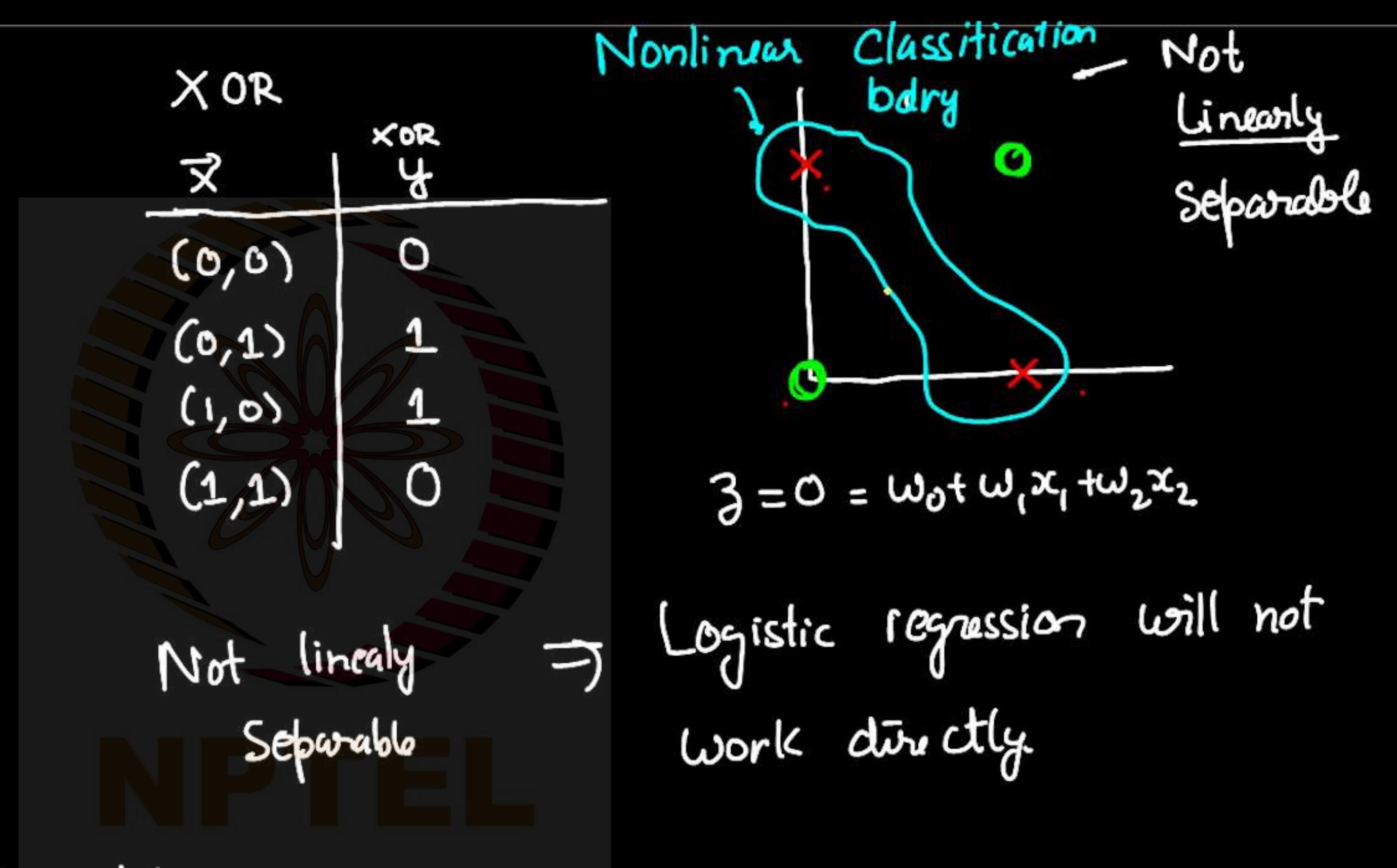


Nonlinear Classification  
bdry



$$\mathcal{Z} = 0 = \omega_0 + \omega_1 x_1 + \omega_2 x_2$$

Not  
Linearly  
Separable

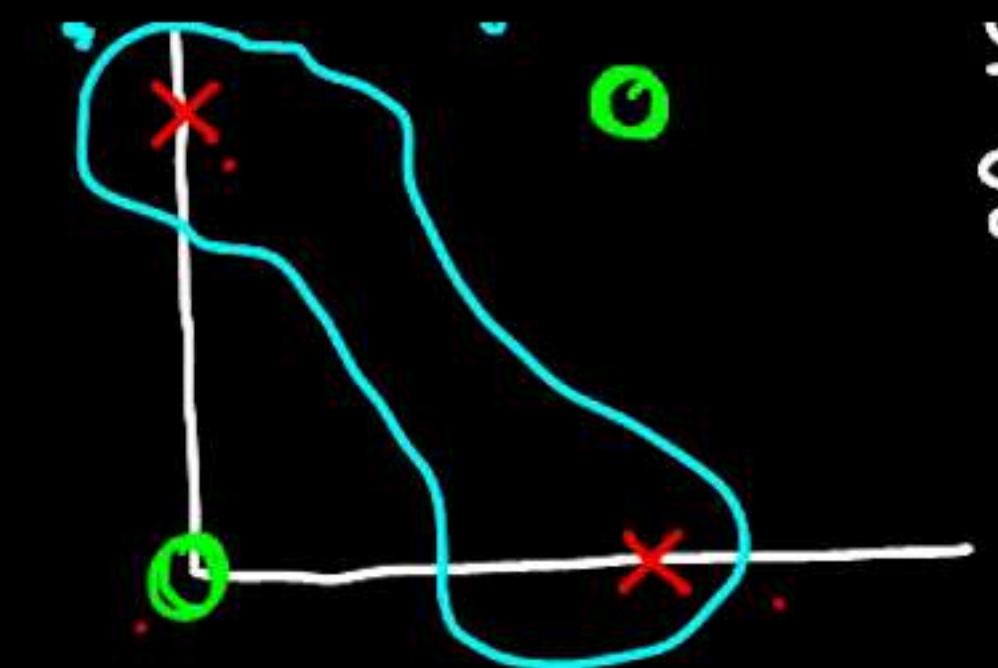


Two possible solutions

- ① Nonlinear features → "Kernel" trick

$\vec{x}$	xOR
(0, 0)	0
(0, 1)	1
(1, 0)	1
(1, 1)	0

Not linearly Separable



Nonlinear  
Separable

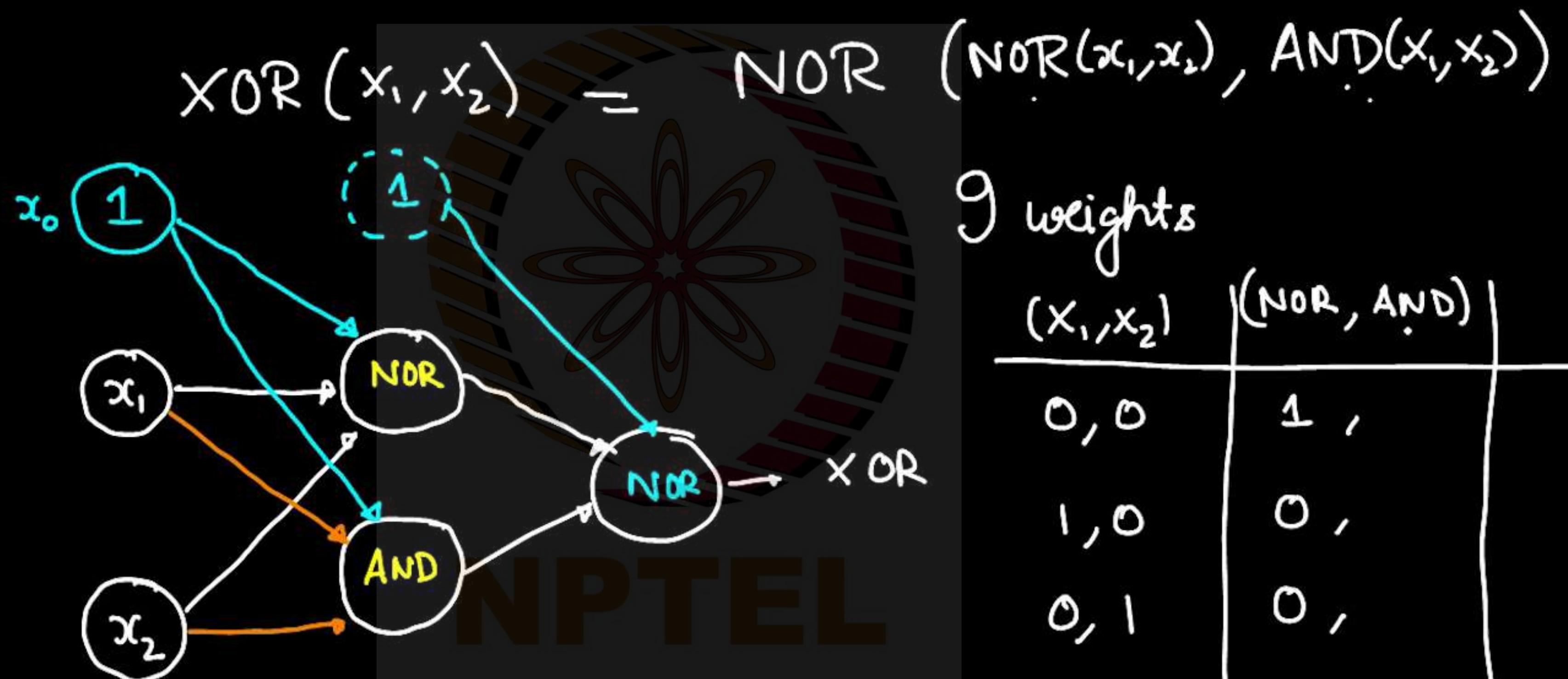
$$z = 0 = w_0 + w_1 x_1 + w_2 x_2$$

Logistic regression will not work directly.

Two possible solutions

- ① Nonlinear features  $\rightarrow$  "Kernel" trick
- ② Add extra layers  $\rightarrow$  "Deep" networks.

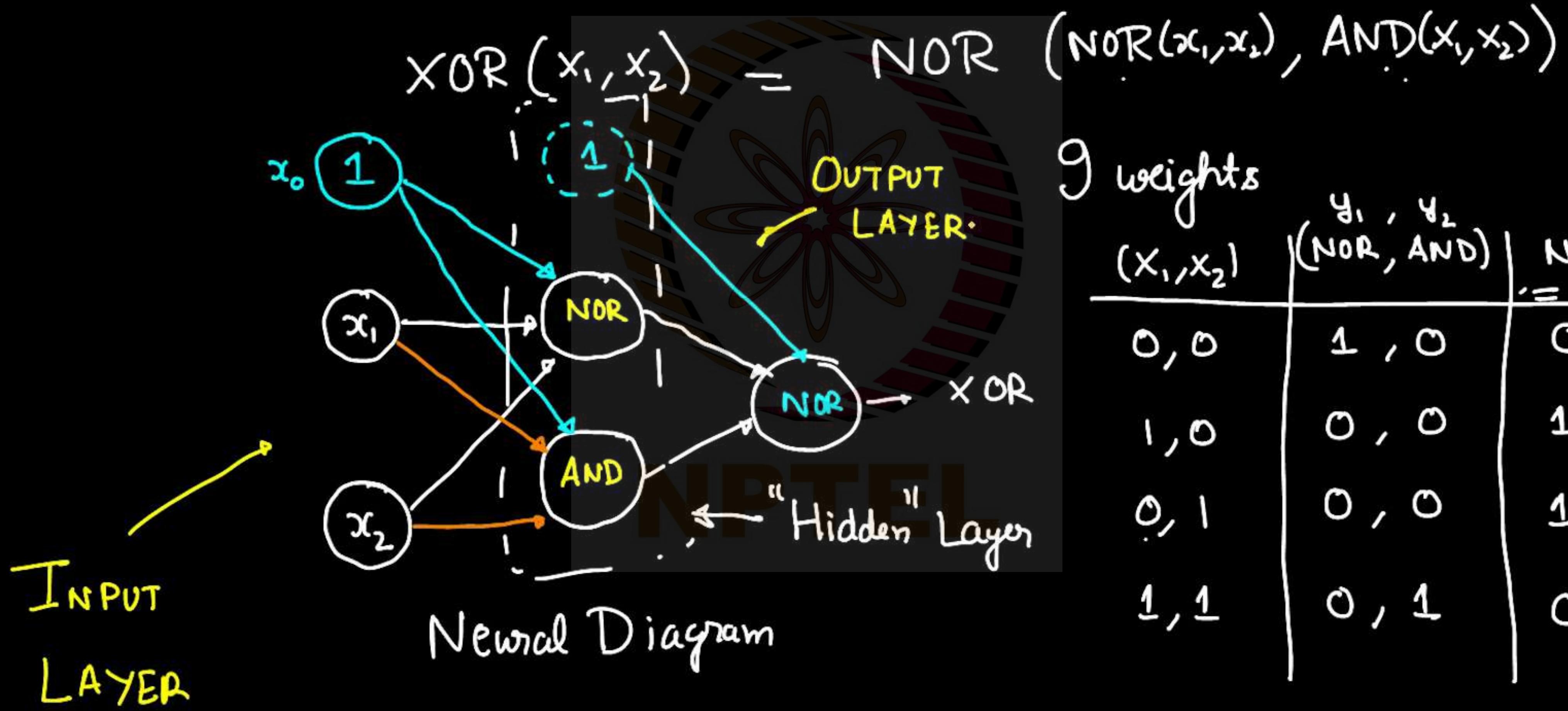
② Add extra layers  $\rightarrow$  "Deep" networks.



9 weights

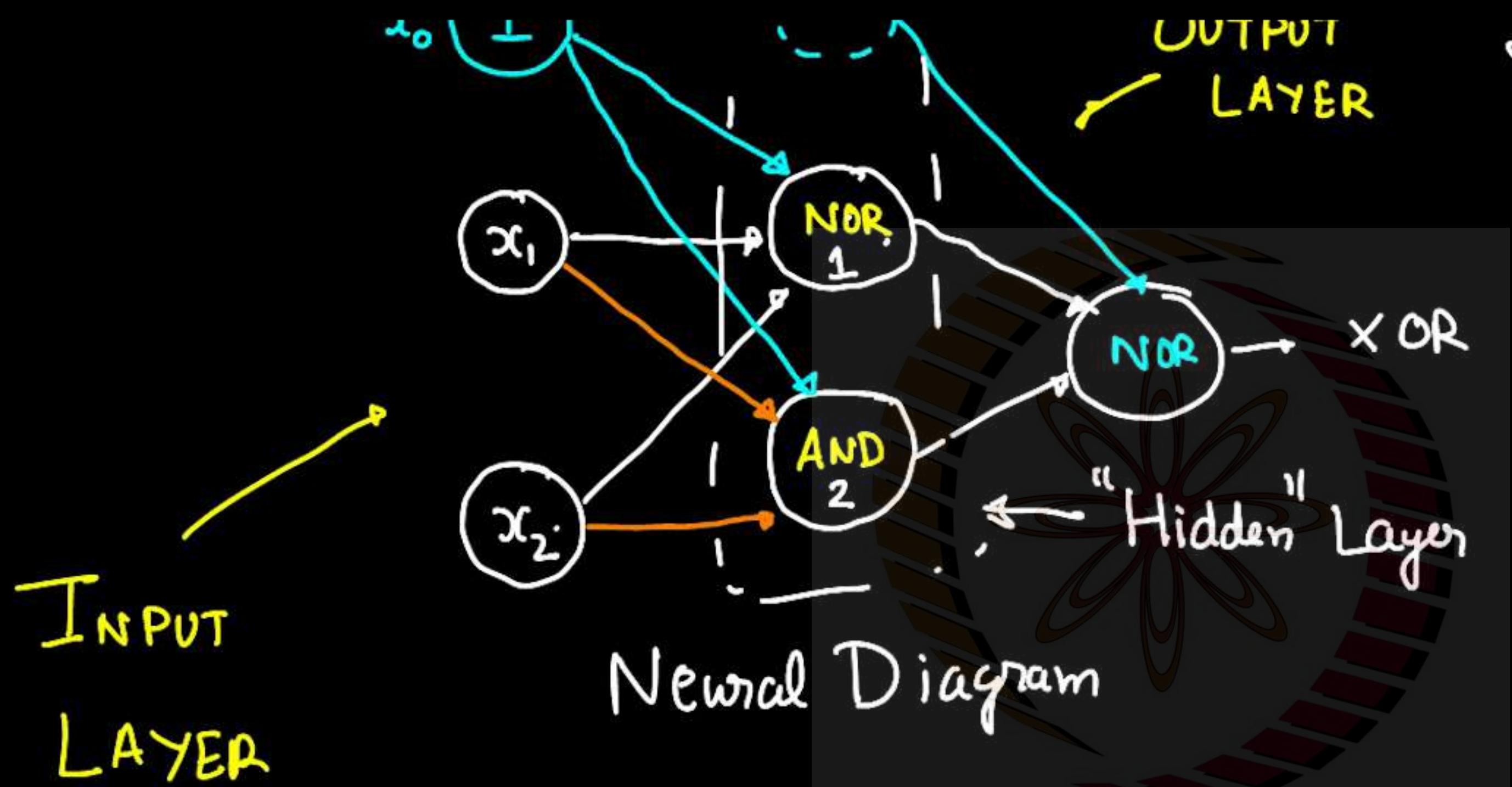
$(x_1, x_2)$	$(\text{NOR}, \text{AND})$
0, 0	1,
1, 0	0,
0, 1	0,
1, 1	0,

② Add extra layers → "Deep" networks.



9 weights

$(x_1, x_2)$	$y_1, y_2$ $(\text{NOR}, \text{AND})$	$\text{NOR}(y_1, y_2)$ $= \text{XOR}$
0, 0	1, 0	0
1, 0	0, 0	1
0, 1	0, 0	1
1, 1	0, 1	0



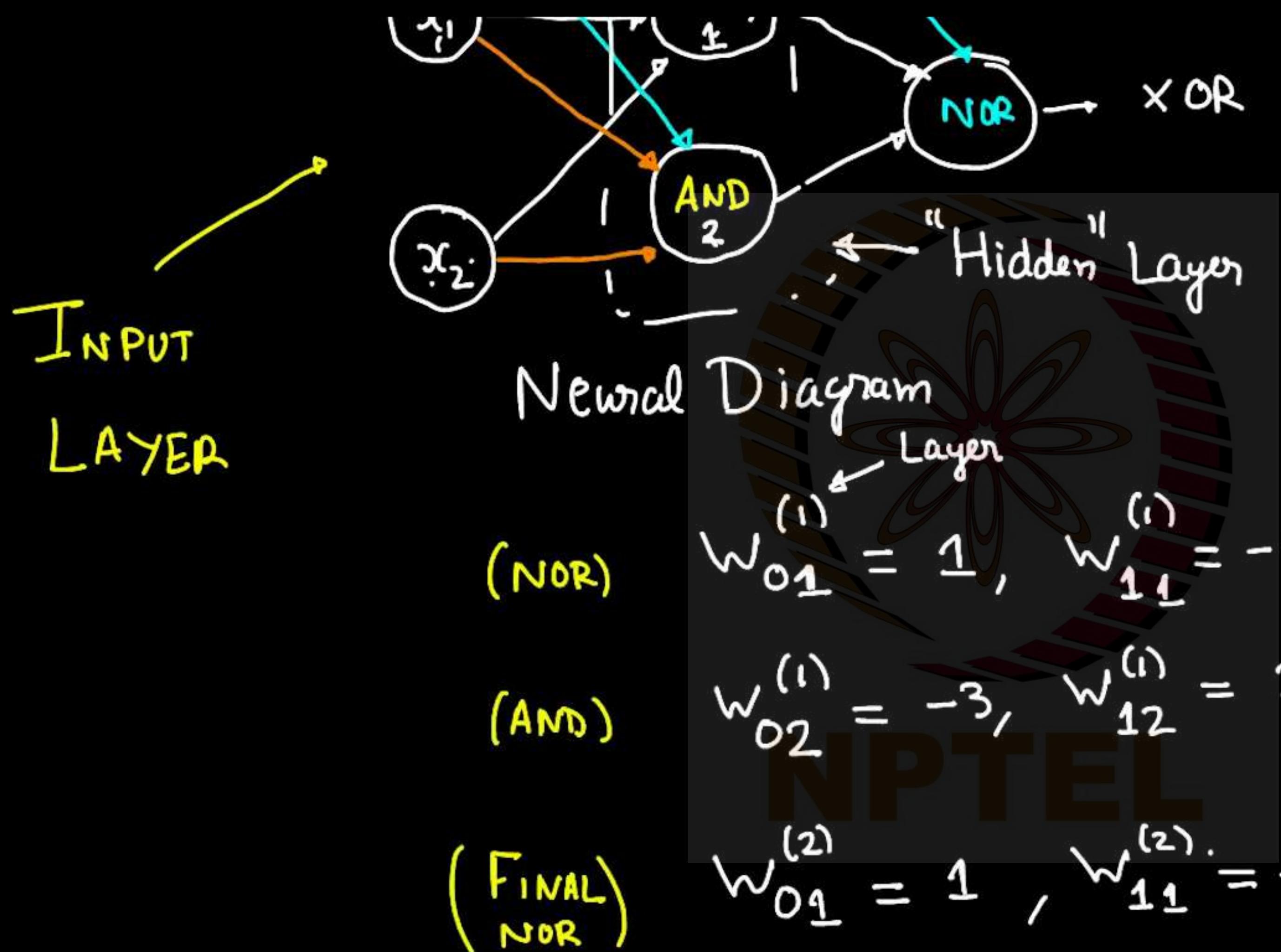
(NOR)

$$w_{01} = 1, w_{11} = -2, w_{21} = -2$$

(AND)

$$w_{02} = -3, w_{12} = 2, w_{22} = 2$$

Weights		$y_1, y_2$ (NOR, AND)	$\text{NOR}(y_1, y_2)$ $= \text{XOR}$
$(x_1, x_2)$	$0, 0$	$1, 0$	0
	$1, 0$	$0, 0$	1
	$0, 1$	$0, 0$	1
	$1, 1$	$0, 1$	0



0, 0	1, 0	0
1, 0	0, 0	1
0, 1	0, 0	1
1, 1	0, 1	0

Layer (1) weights  
6 weights in total

$$\left. \begin{array}{l} w_{01}^{(1)} = 1, \quad w_{11}^{(1)} = -2, \quad w_{21}^{(1)} = -2 \\ w_{02}^{(1)} = -3, \quad w_{12}^{(1)} = 2, \quad w_{22}^{(1)} = 2 \end{array} \right\}$$

$$w_{01}^{(2)} = 1, \quad w_{11}^{(2)} = -2, \quad w_{21}^{(2)} = -2$$

(NOR)  $w_{01}^{(1)} = 1, w_{11}^{(1)} = -2, w_{21}^{(1)} = -2 \quad \left. \begin{array}{l} \text{Layer (1) weights} \\ 6 \text{ weights in total} \end{array} \right\}$

(AND)  $w_{02}^{(1)} = -3, w_{12}^{(1)} = 2, w_{22}^{(1)} = 2$

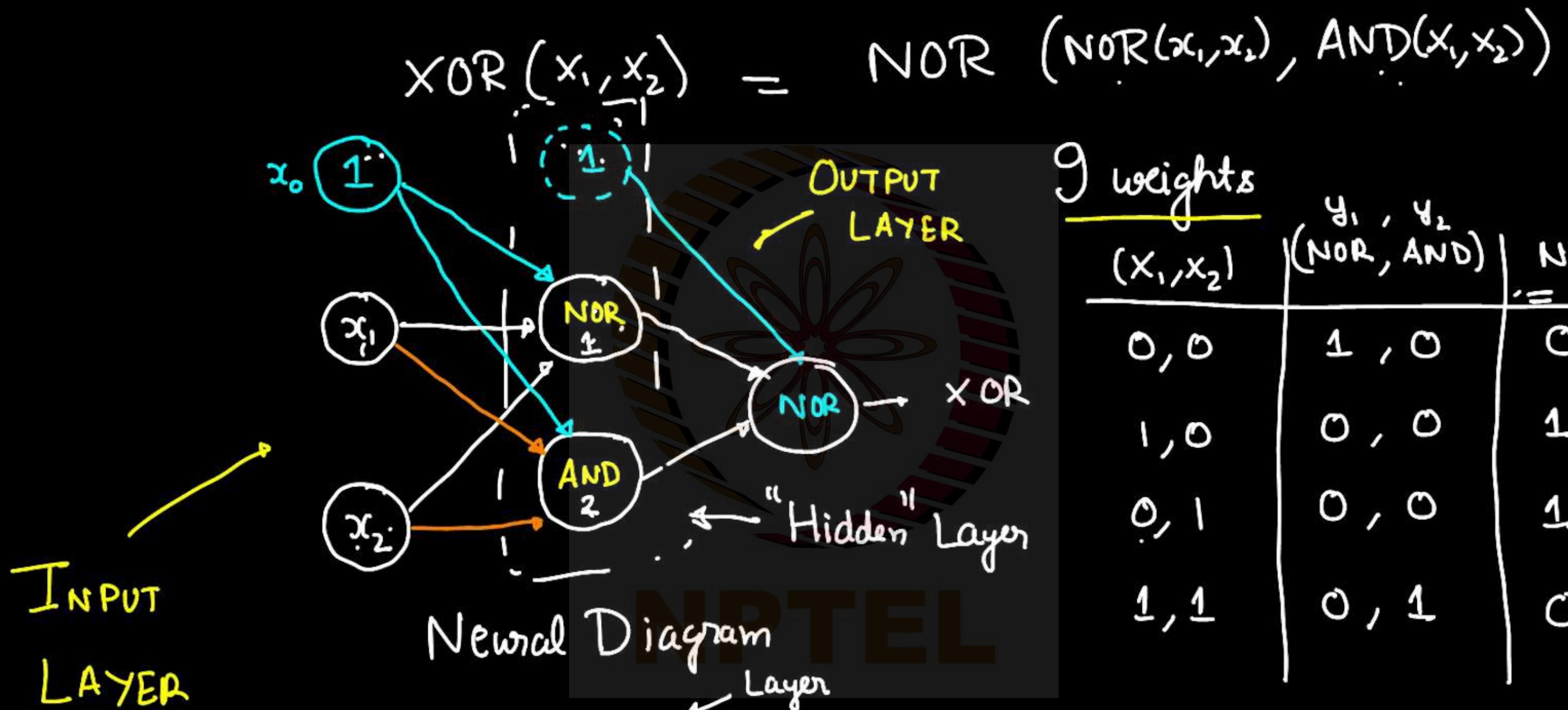
(FINAL NOR)  $w_{01}^{(2)} = 1, w_{11}^{(2)} = -2, w_{21}^{(2)} = -2$

Universal Approximation

Thm : Any function can be expressed

as a neural network with one hidden

layer to desired accuracy



9 weights

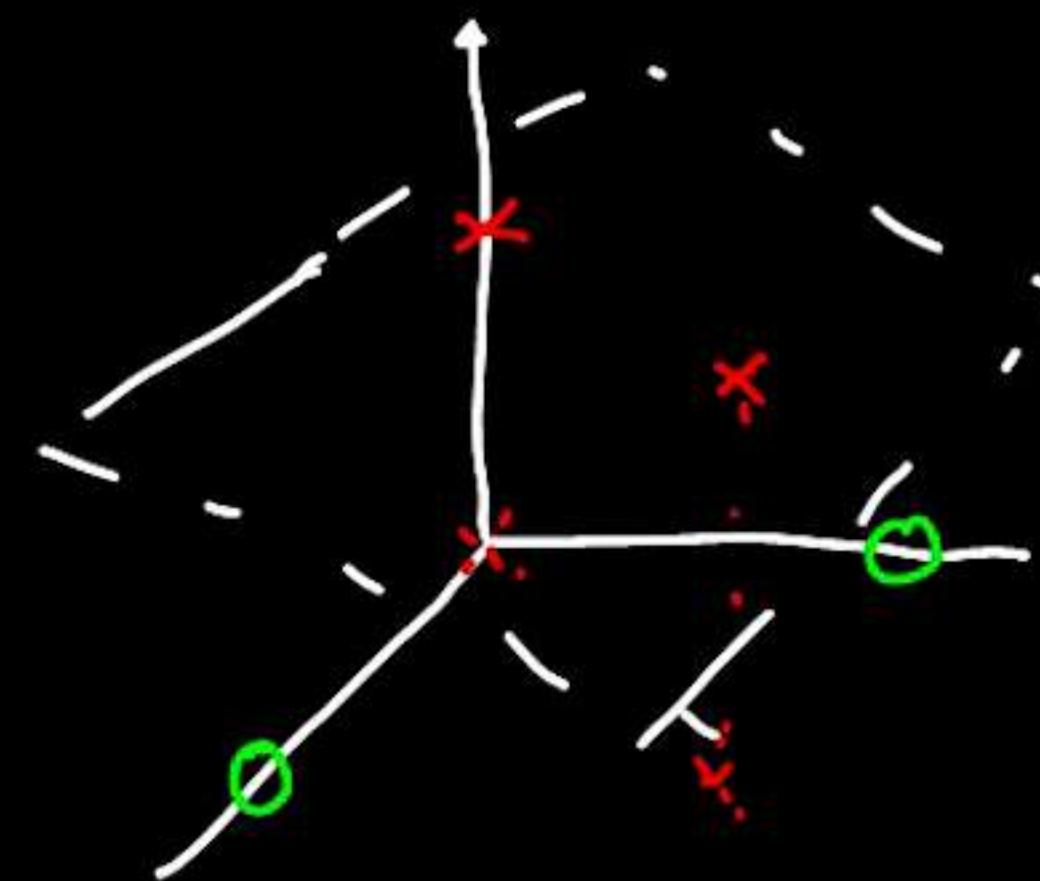
$(x_1, x_2)$	$y_1, y_2$ $(\text{NOR}, \text{AND})$	$\text{NOR}(y_1, y_2)$ $= \text{XOR}$
0, 0	1, 0	0
1, 0	0, 0	1
0, 1	0, 0	1
1, 1	0, 1	0

Layer (1) weights  
6 weights in total

$$(NOR) \quad w_{01}^{(1)} = 1, \quad w_{11}^{(1)} = -2, \quad w_{21}^{(1)} = -2 \quad \left. \right\} \text{Layer (1) weights}$$

$$(\text{AND}) \quad w_{01}^{(1)} = -3, \quad w_{11}^{(1)} = 2, \quad w_{22}^{(1)} = 2 \quad \left. \right\} 6 \text{ weights in total}$$

# Note Title



Linearly separable in a higher dimension.

Two possible solutions



Not linearly Separable  $\Rightarrow$

Nonlinear classification bdry  
Not linearly separable

$$\mathbf{z} = \mathbf{0} = w_0 + w_1 x_1 + w_2 x_2$$

Logistic regression will not work directly.

Note Title

## Classification

Multinomial Classification,

$$K > 2$$

No of classes

$0, 1, \dots, 9$   
 $K = 10$  class  
.

Multinomial Classification,

$$K > 2$$

↓

$0, 1, \dots, 9$

No of classes  $K = 10$  class

What we need

horses, cats, dogs

① How do we represent the output class (label)?

↳ One hot vector ↗

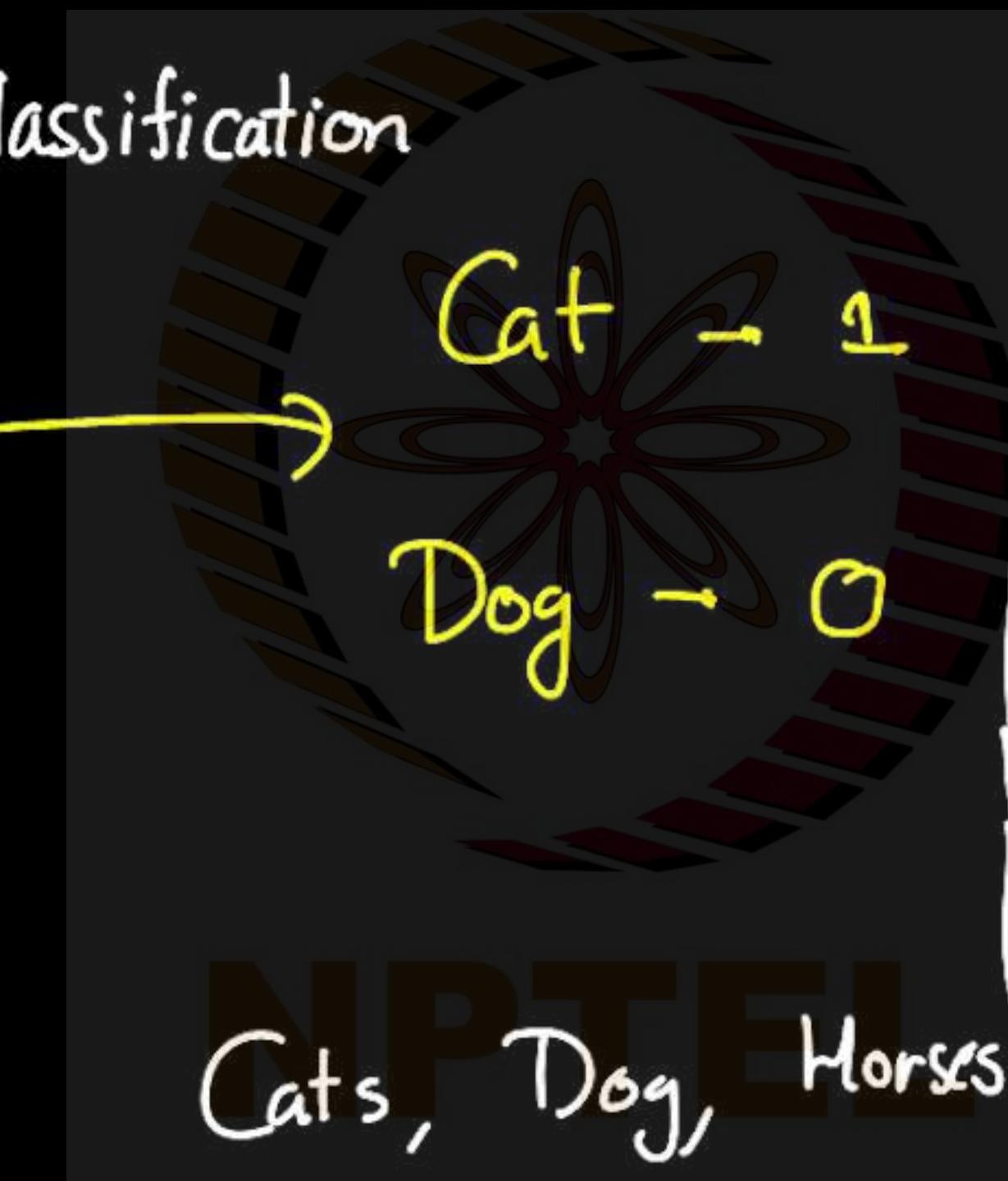
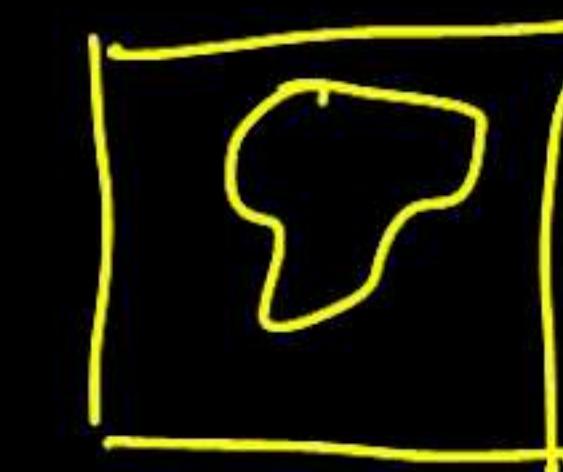
② What happens in the final layer for nonlinearity?

↳ Softmax ↗

③ Loss function

# Note Title

## Multinomial Classification



Vector

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

Cat

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

Dog

Vector -

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

Horse

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

1  
+  
1

vog

u

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_k \end{bmatrix} \rightarrow \text{prob that output belongs to class 1} \rightarrow \in [0, 1]$$

Vector -

$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$	,	$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$	,	$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$
Dog	Horse	Cat		

One 'hot'

$$\hat{y} = \begin{bmatrix} 0.9 \\ 0.01 \\ 0.09 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_k \end{bmatrix}$$

that output  
belongs to class 1  
 $\rightarrow \in [0, 1]$

One hot

$$\hat{y} = \begin{bmatrix} 0.9 \\ 0.01 \\ 0.09 \end{bmatrix}$$

$$\hat{y}_k = P(\text{belongs to class } k)$$
$$\hat{y}_k = P(y_k = 1 \mid \text{Image/Input } x)$$



## Differentiating the Sigmoid

$$\begin{aligned}\sigma(z) &= \frac{1}{1 + e^{-z}} \\ \frac{d\sigma}{dz} &= \sigma'(z) = -\frac{1}{(1 + e^{-z})^2} \cdot \frac{d e^{-z}}{dz} \\ &= \frac{1}{(1 + e^{-z})^2} \cdot e^{-z} \\ &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}\end{aligned}$$



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

$$\begin{aligned}\frac{d\sigma}{dz} &= \sigma'(z) = \frac{-1}{(1 + e^{-z})^2} \frac{d e^{-z}}{dz} \\ &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}} \\ &= \frac{1}{1 + e^{-z}} \cdot \frac{e^{-z}}{1 + e^{-z}}\end{aligned}$$

$$\sigma'(z) = \sigma(z) [1 - \sigma(z)]$$



File Edit View Insert Actions Tools Help



$$\hat{y} = \sigma(w \cdot x) \rightarrow [1 \ x_1 \ \dots \ x_n]$$

 $J = ?$ 

What about least squares? → For classification?

$$J_{LS} = \frac{1}{2} \sum_{i=0}^n (y_i - \hat{y}_i)^2$$

↑ So  $\hat{y} = 1$

→ Not a good cost function

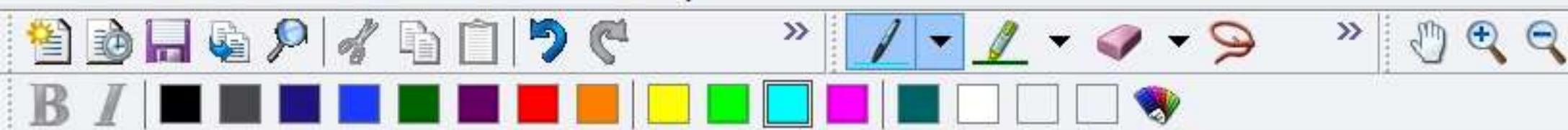
Cost incurred for misclassification

is low

$x$	$y$	$\hat{y}$
0, 1	[0, 1]	



File Edit View Insert Actions Tools Help



NPTEL

$$J = \frac{1}{2} [0 \cdot 0 + 1 \cdot 1]$$

Soy  $\hat{y} = 1$

→ Not a good cost function

Cost incurred for misclassification  
is low

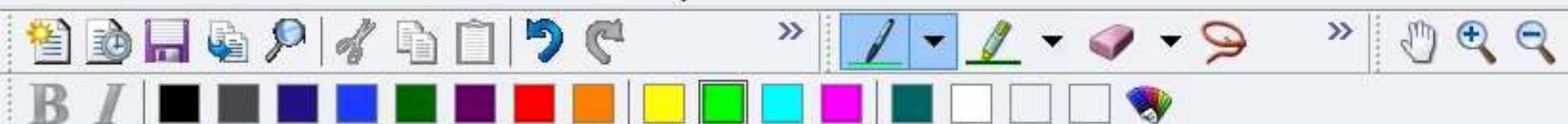
We use Binary Cross Entropy Cost Function

$$J = - \{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \}$$

Desirable properties : (1)  $J=0$  if  $y=\hat{y}$

(2)  $J$  should be very high for misclassification

Consistency → (3)  $J \geq 0$



We use

Binary Cross Entropy Cost Function

$y$  is either 0 or 1

$$J = - \{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \}$$

Desirable properties

Consistency  $\rightarrow$  (3)

: (1)  $J=0$  if  $y=\hat{y}$

(2)  $J$  should be very high for misclassification

$$J \geq 0 \checkmark$$

$y=0, \hat{y}=0 \leftarrow$  close to 0.



$$\text{J} = - \left\{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right\}$$

Desirable properties

Consistency

: (1)  $\text{J}=0$  if  $y=\hat{y}$  ✓

(2)  $\text{J}$  should be very high for misclassification

→ (3)  $\text{J} > 0$  ✓

$y=0$ ,  $\hat{y}=0$  ← close to 0.  $\Rightarrow \text{J} \approx 0$

$y=1$ ,  $\hat{y} \approx 1$   $\Rightarrow \text{J} \approx 0$

File Edit View Insert Actions Tools Help



We use Binary Cross Entropy Cost Function  $\approx 1$

$$J = - \left\{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right\}$$

$y$  is either 0 or 1

Desirable properties:

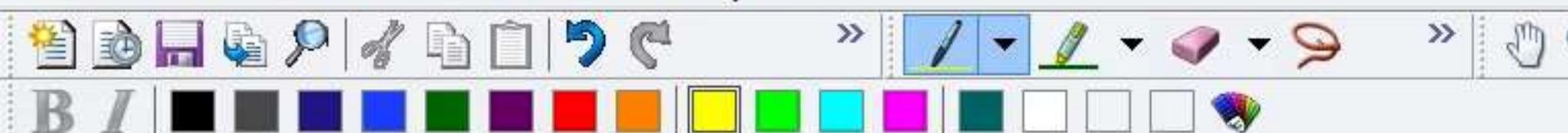
- (1)  $J=0$
- (2)  $J$  should be very high for misclassification
- (3)  $J \geq 0$  ✓

Consistency → (3)

$y = 0, \hat{y} = 0 \leftarrow$  close to 0.  $\Rightarrow J \approx 0$

$y = 1, \hat{y} \approx 1 \Rightarrow J \approx 0$

IS low



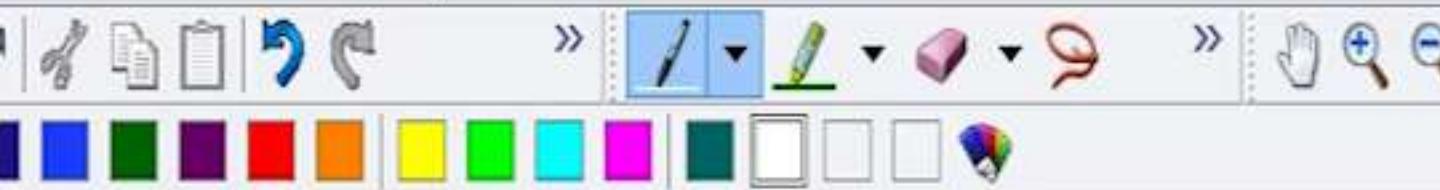
Consistency  $\rightarrow (3)$   $J \geq 0 \checkmark$

$y = 0$ ,  $\hat{y} = 0 \leftarrow$  close to 0.  $\Rightarrow J \approx 0$

$y = 1$ ,  $\hat{y} \approx 1 \Rightarrow J \approx 0$

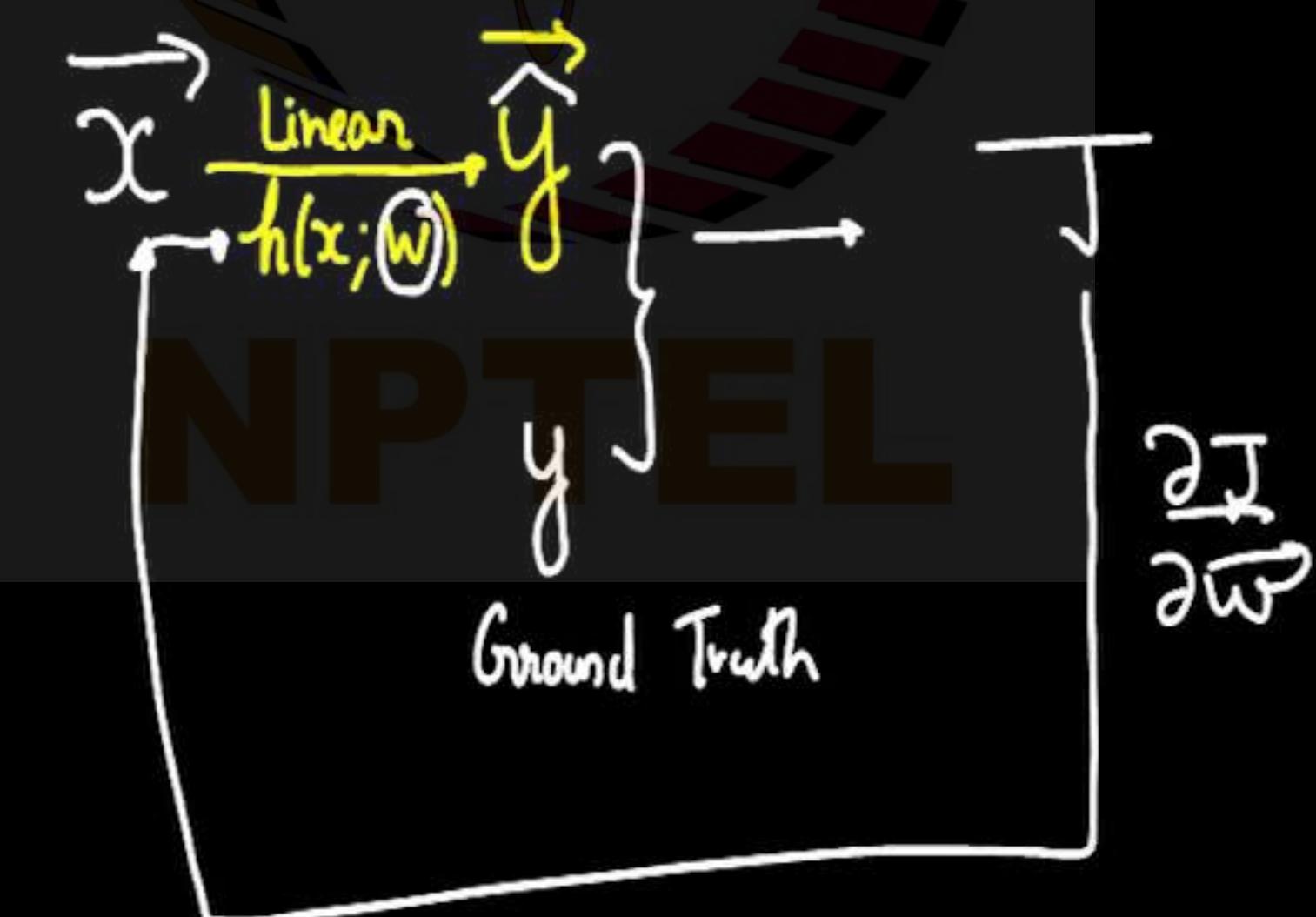
$y = 0$ ,  $\hat{y} \approx 1 \Rightarrow J \rightarrow \infty$

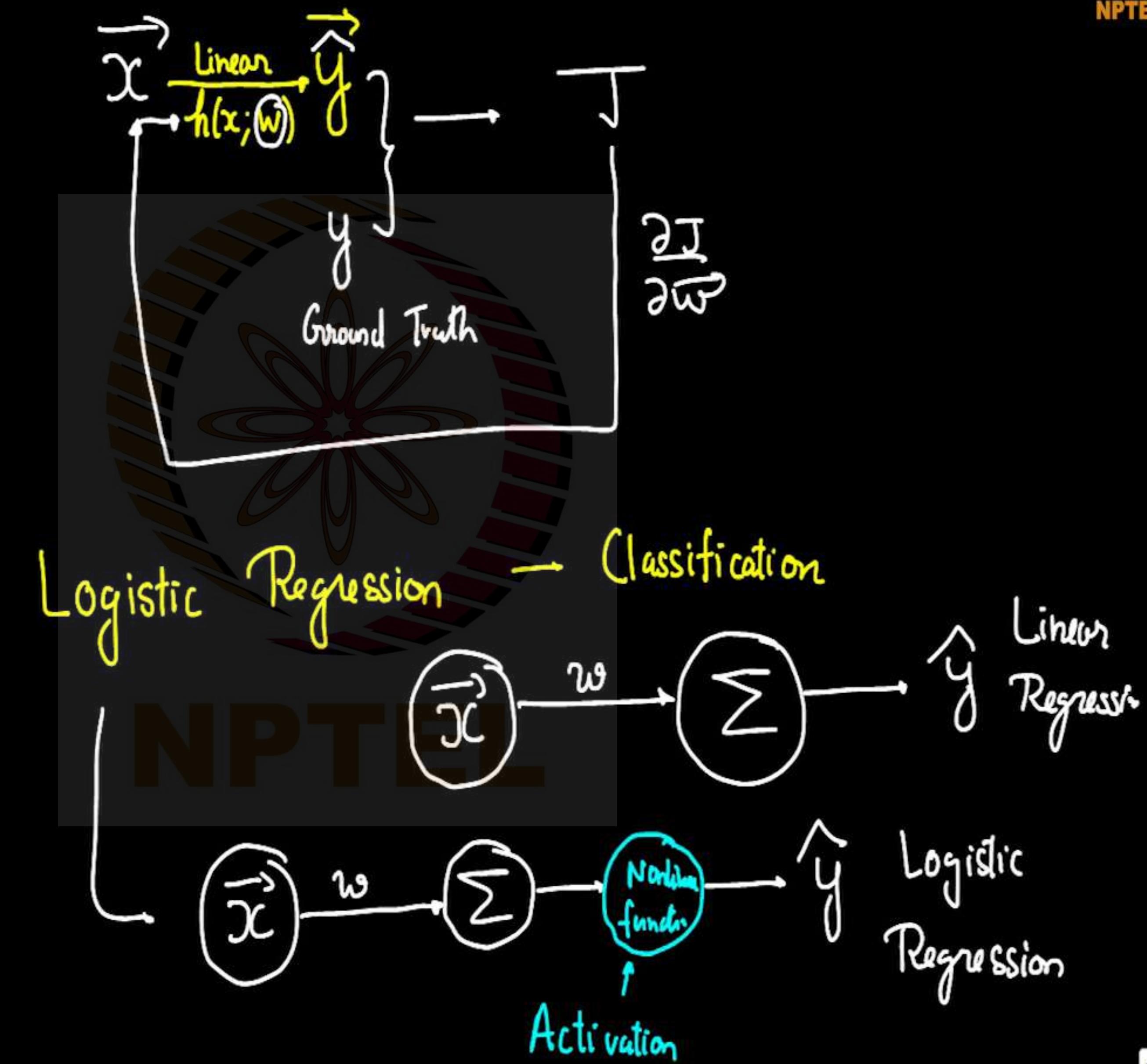
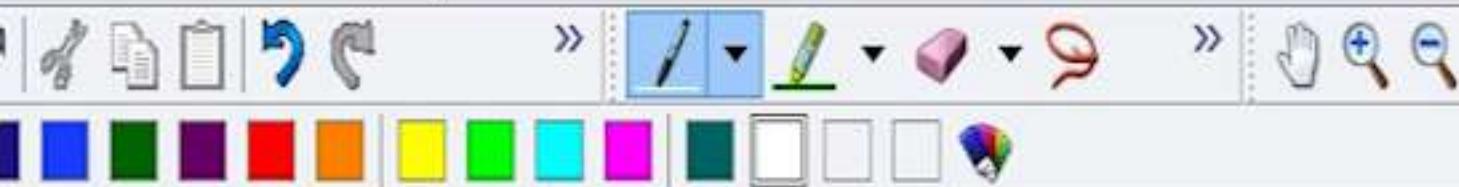
Check  $y = 1$ ,  $\hat{y} \approx 0 \Rightarrow J \rightarrow \infty$



Note Title

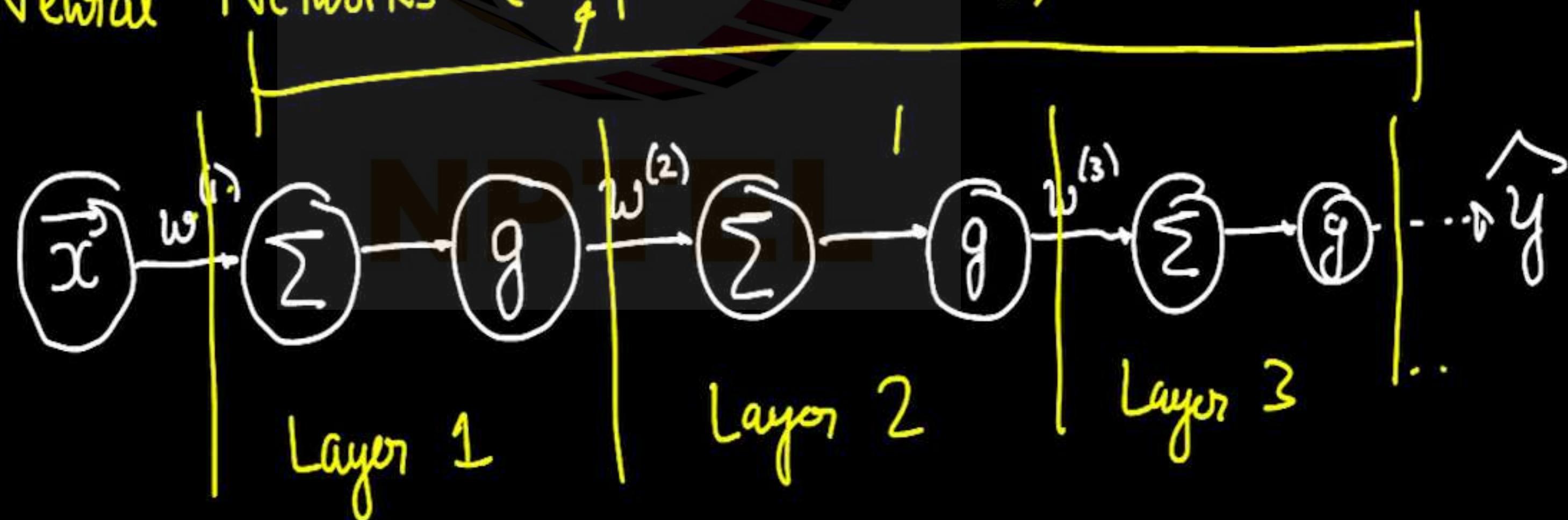
## Introduction to

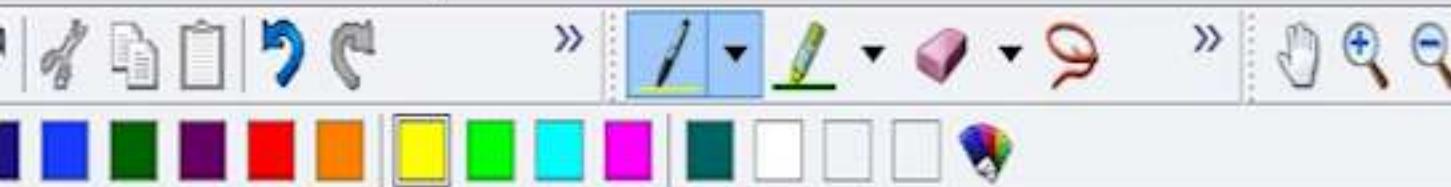






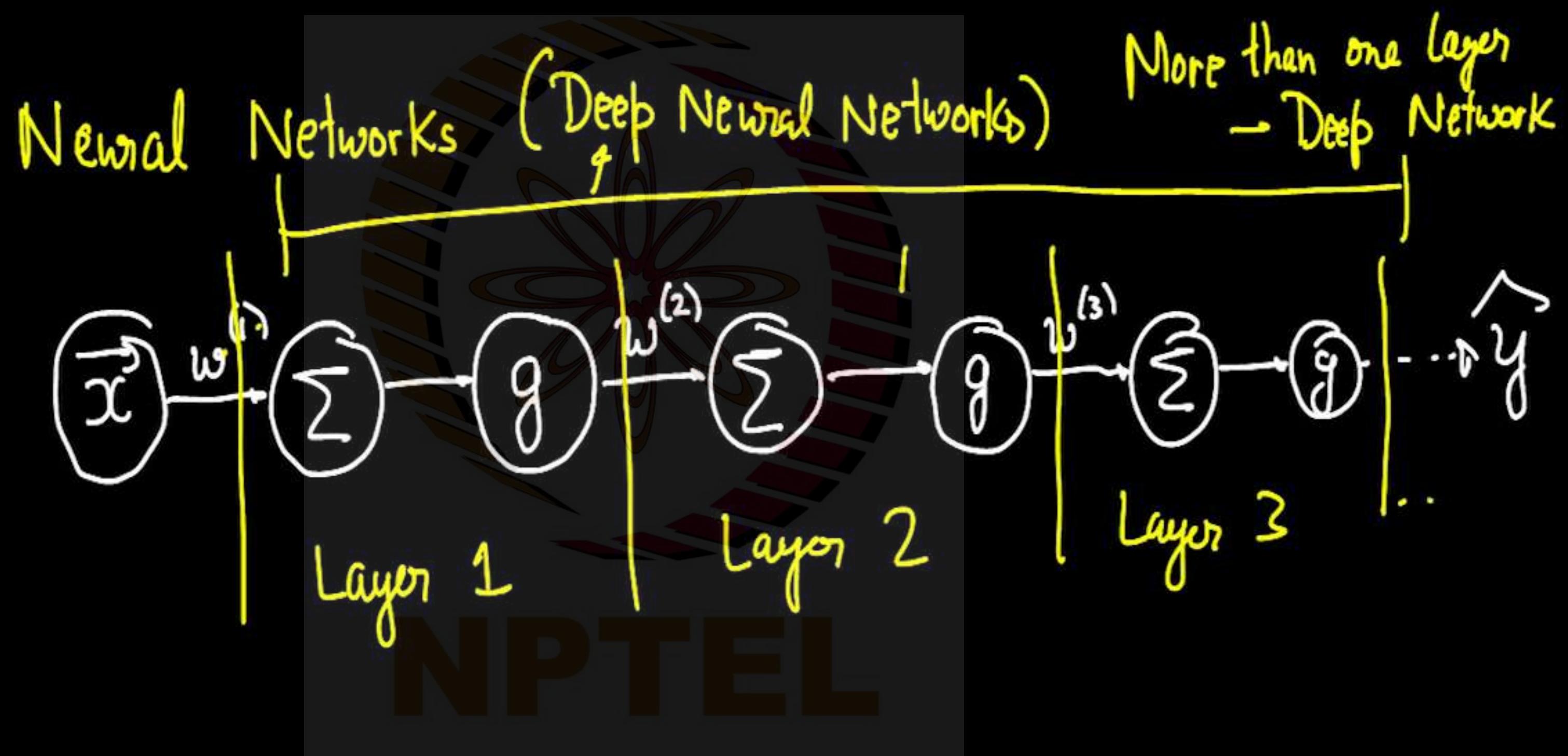
## Neural Networks (Deep Neural Networks)

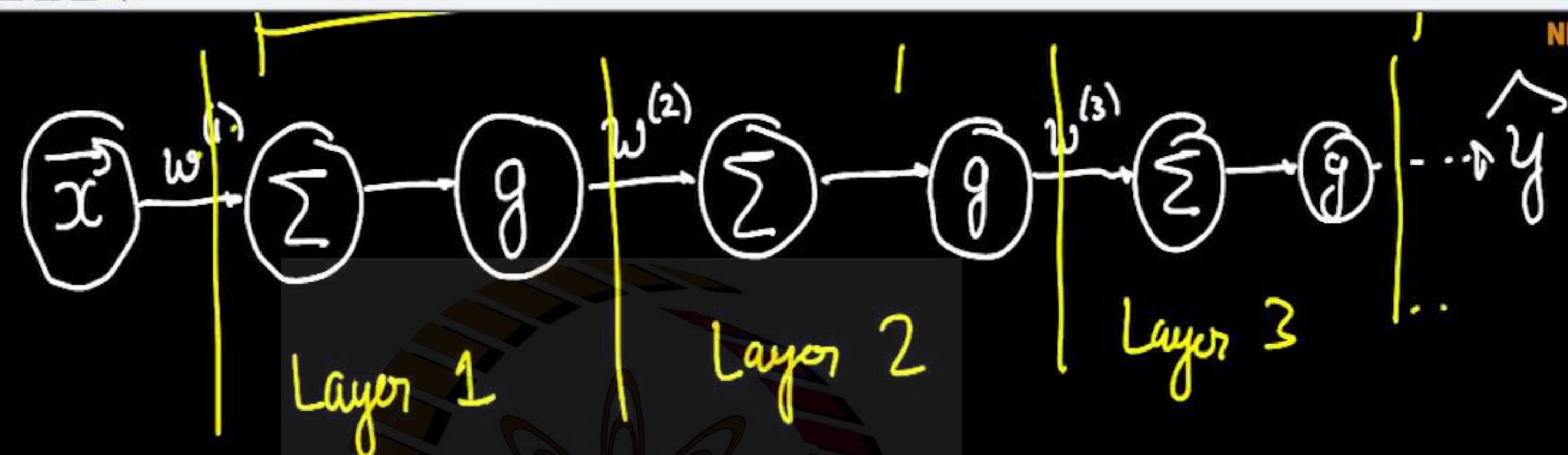
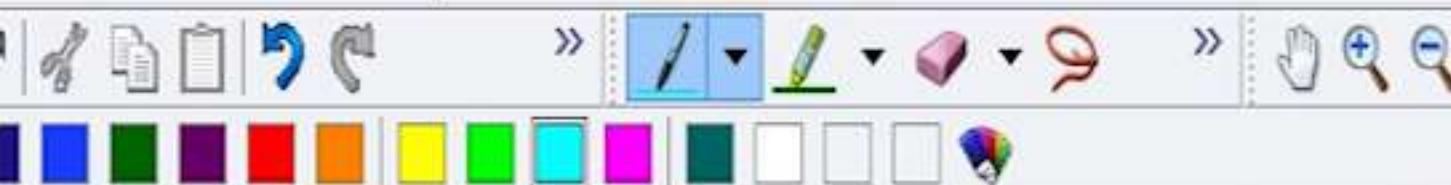




Tatyasaheb

NPTEL





Pay attention to  $\Gamma$  How do we assign numbers?

$\hat{y}$  a) How do we characterize the output  $y, \hat{y}$

$x \xrightarrow[w]{h(x)} \hat{y}$  b) What is the forward model? [Which nonlinear function?  $g$ ]

$\hat{y}, y \rightarrow J$  c) What is the loss function  $J$ ?



Layer 1 | Layer L | Layer -

Pay attention to

How do we assign numbers?

$\hat{y}$

a) How do we characterize the output  $y, \hat{y}$

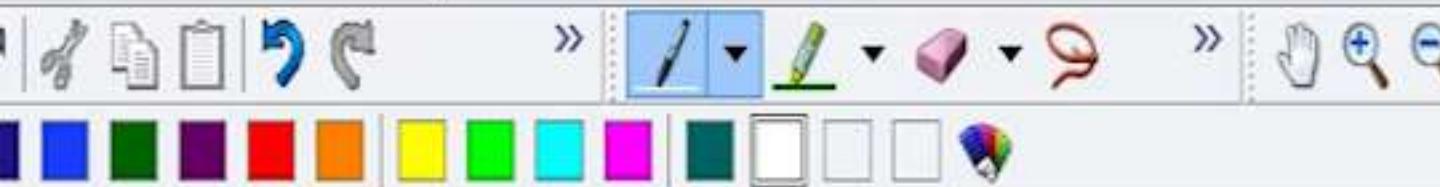
$x \xrightarrow{h(x)} \hat{y}$

b) What is the forward model? [Which nonlinear function?  $g$ ]

$\hat{y}, y \rightarrow J$  c) What is the loss function  $J$ ?

d) How do we calculate  $\frac{\partial J}{\partial w}$ ? Gradient

e) How do we use  $\frac{\partial J}{\partial w}$  to find better  $w$



'Pay attention TO'

$\hat{y}$

a) How do we characterize the output  $y, \hat{y}$  (Representation)

$x \xrightarrow[w]{h(x)} \hat{y}$

b) What is the forward model? [Which nonlinear function? g]

$\hat{y}, y \rightarrow J$

c) What is the loss function  $J$ ?

$\frac{\partial J}{\partial w}$

d) How do we calculate  $\frac{\partial J}{\partial w}$ ? Gradient

e) How do we use  $\frac{\partial J}{\partial w}$  to find better  $w$  } Optimization

$$w \leftarrow w - \alpha \frac{\partial J}{\partial w}$$



# Logistic Regression

Note Title

$K > 2$

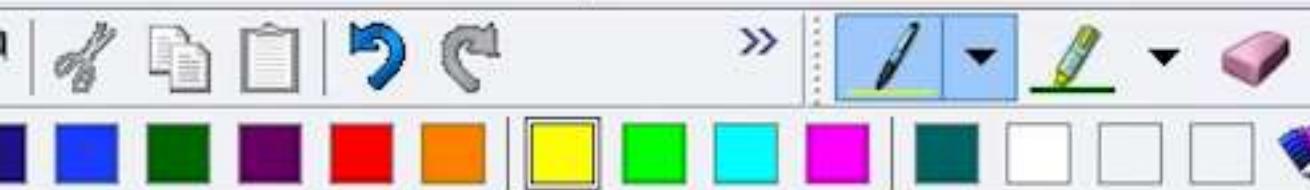
Represent

$K = 3$

$$\hat{y} = \begin{bmatrix} 0.75 \\ 0.1 \\ 0.15 \end{bmatrix}$$

One hot vector

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



Represent  $\hat{y}$  → One hot vector

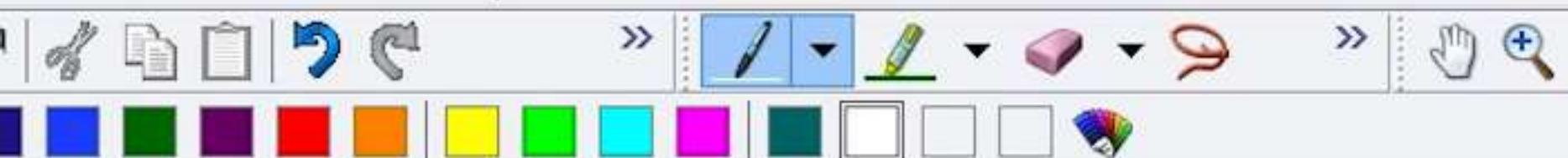
K = 3

$$\hat{y} = \begin{bmatrix} 0.75 \\ 0.1 \\ 0.15 \end{bmatrix}$$

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

What is the nonlinearity that will achieve classification?

NPTEL



$$k = 3$$

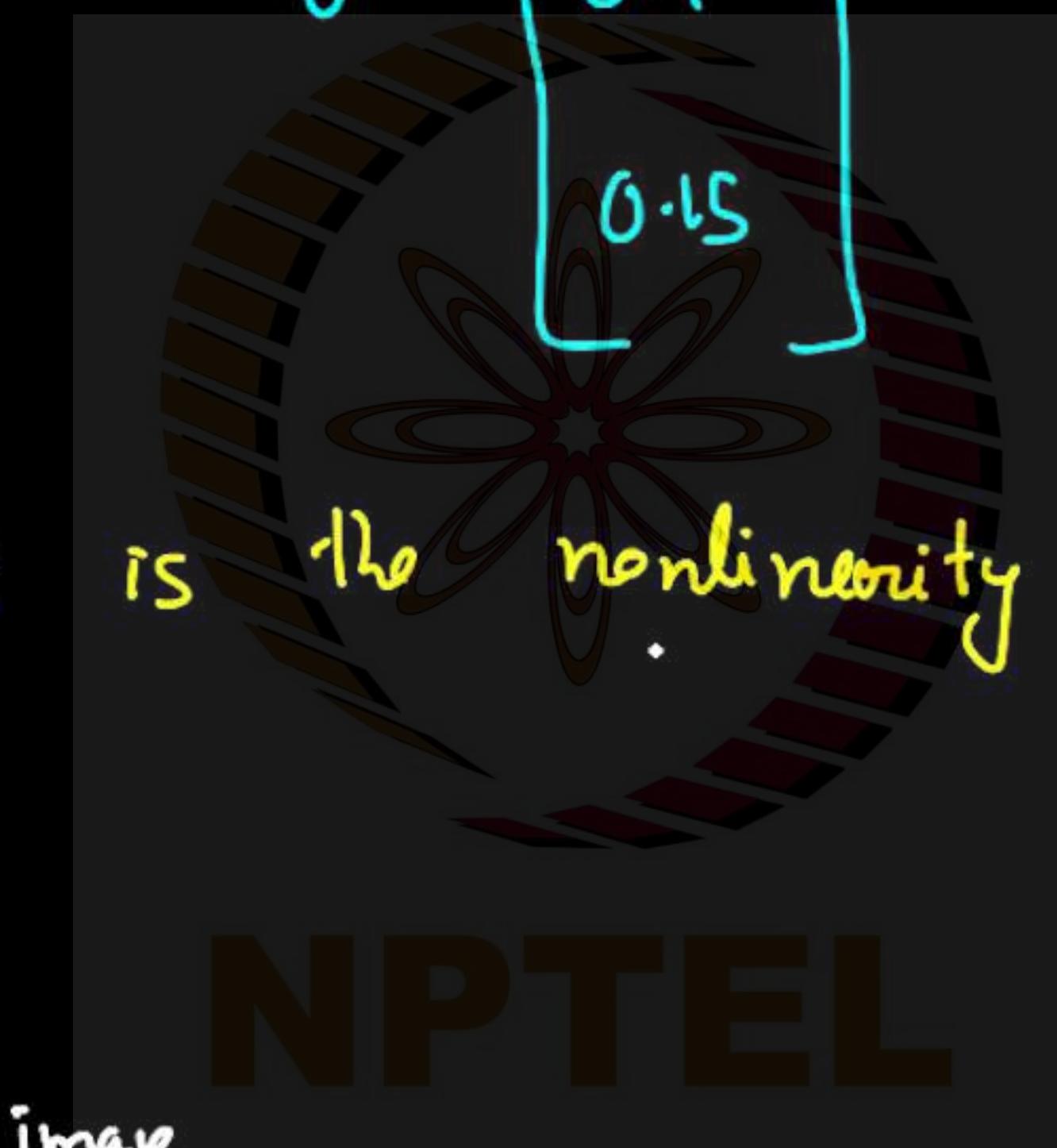
$$\hat{y} = \begin{bmatrix} 0.75 \\ 0.1 \\ 0.15 \end{bmatrix}$$

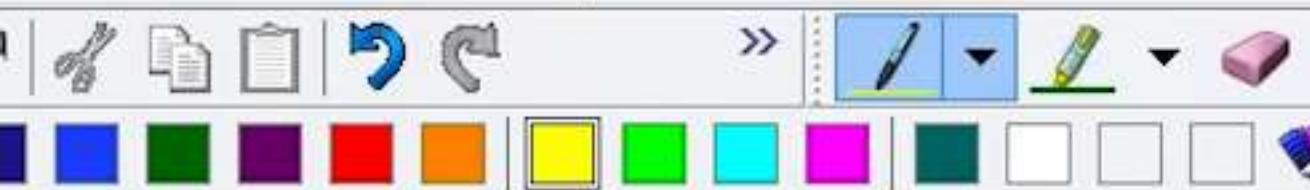
$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

What is the nonlinearity that will achieve classification?

$\vec{x}$  is an image

60x60 grayscale image





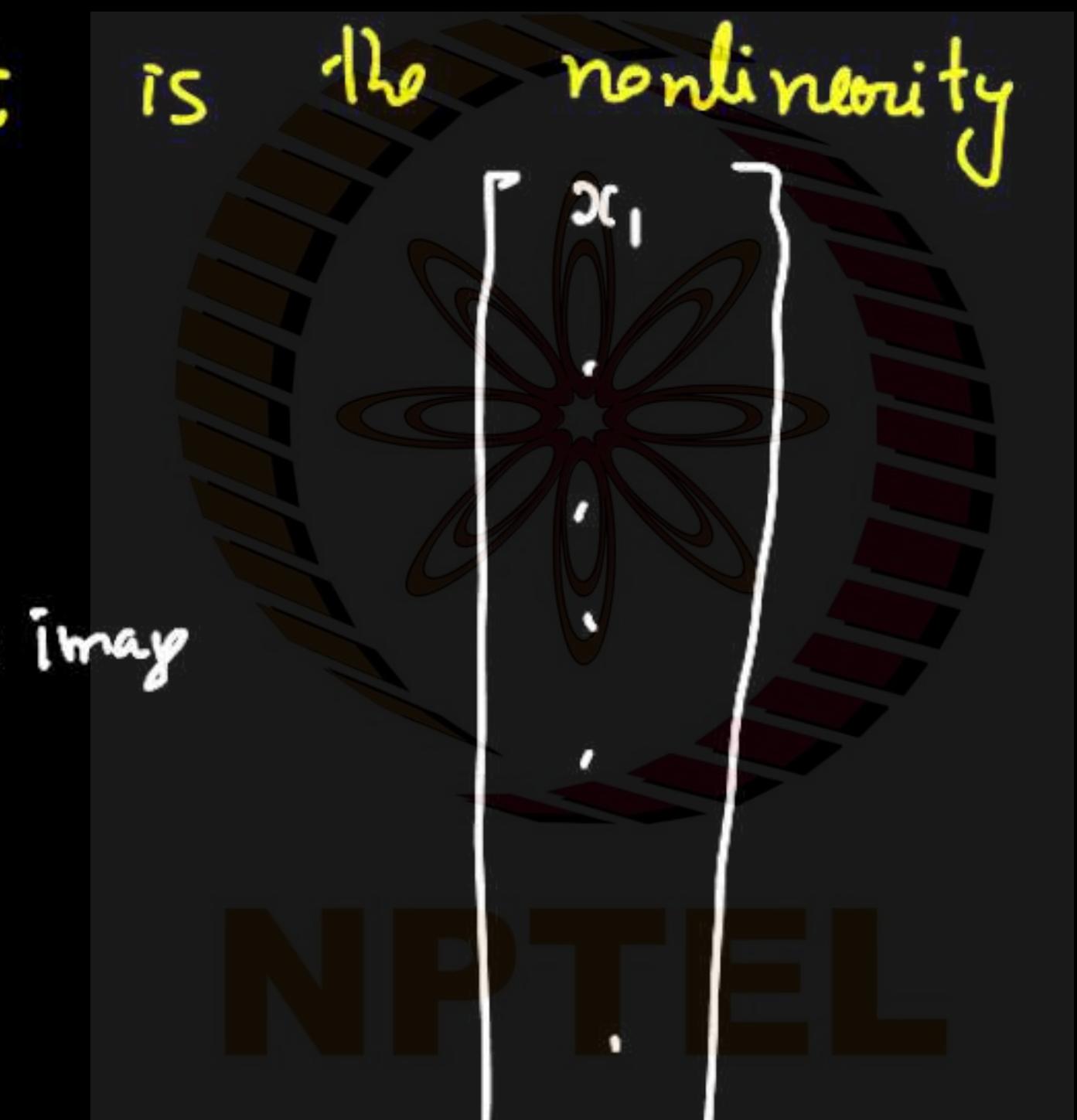
Cat, dog, horse

$\vec{x}$  is an image

60x60 grayscale image

$$\vec{x} \rightarrow \Sigma$$

What is the nonlinearity that will achieve classification?

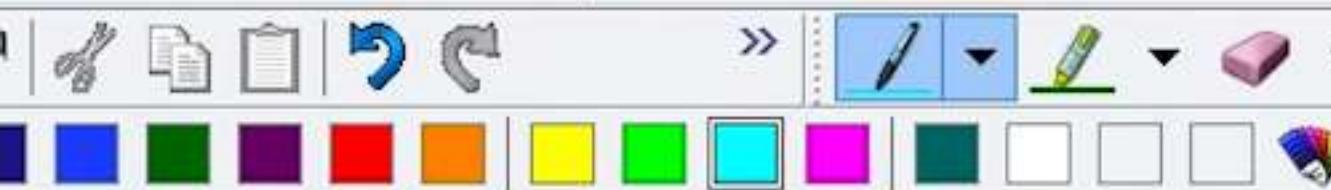


$x_1$

$x_n$

⋮

$\hat{y}$  - 0  
0  
0

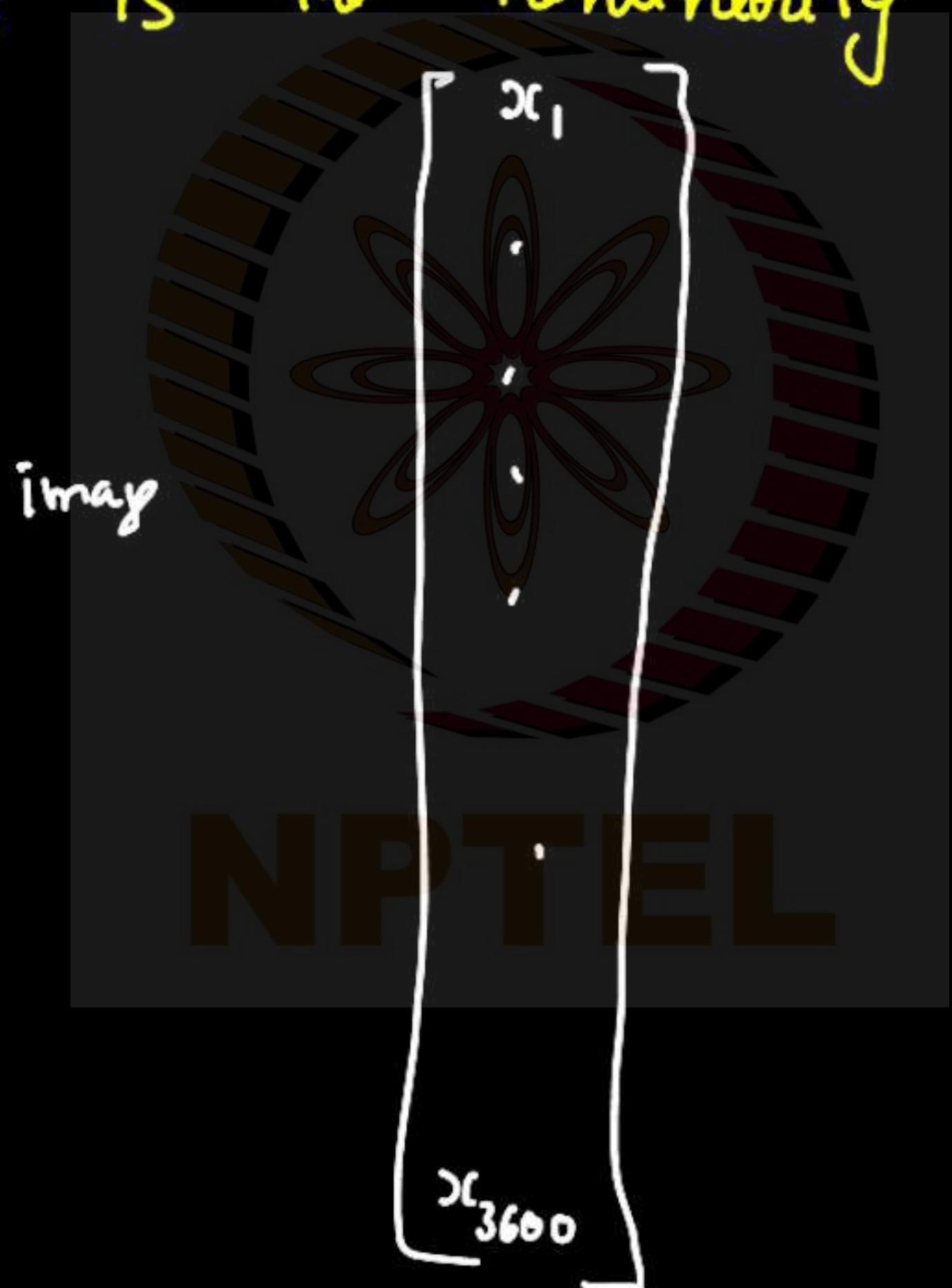


Cat, dog, horse

$\vec{x}$  is an image

60x60 grayscale image

$\vec{x} \rightarrow \Sigma$



What is the nonlinearity that will achieve classification?

$x_1$

$x_n$

$\hat{y}$  - 0  
0  
0

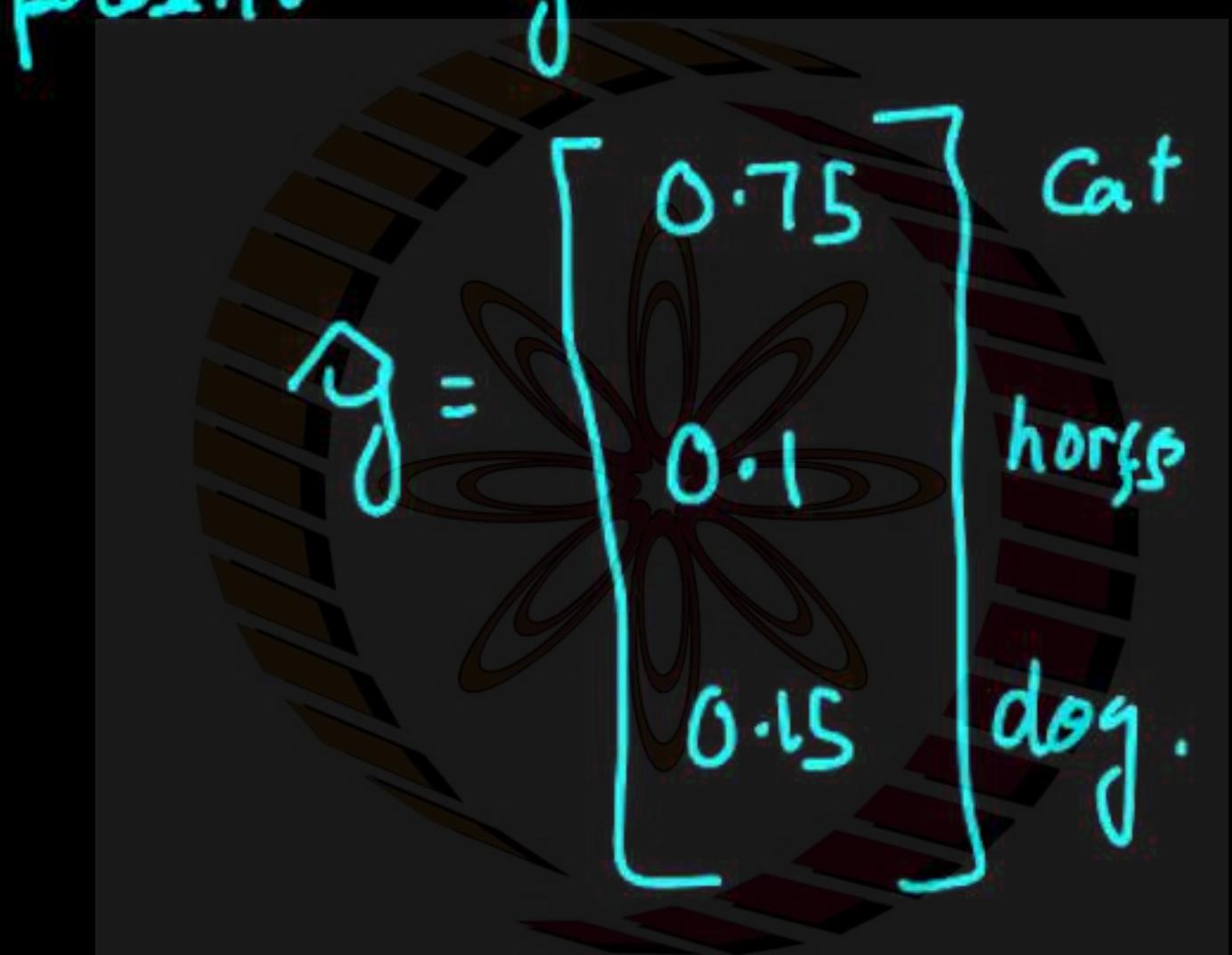
0 - prob



$K > 2$

Represent  $\hat{y} \rightarrow$  One hot vector

$K = 3$



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

What is the nonlinearity that will achieve classification?

Cat, dog, horse

$\vec{x}$  is an image

60x60 grayscale image

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}$$

( $x_i$ )

$\hat{y}$  - O - prob  
O



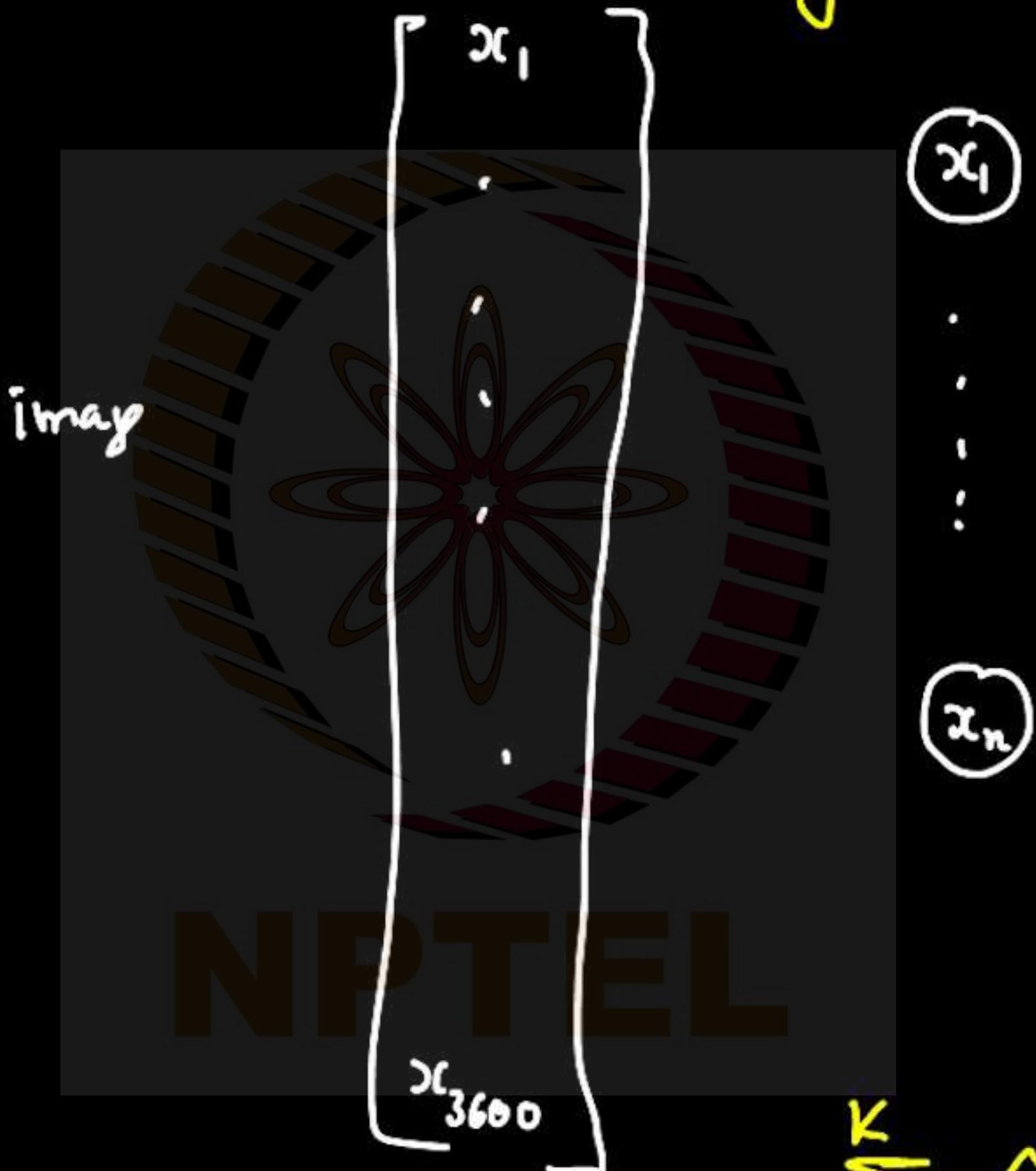
What is the nonlinearity that will achieve classification?

Cat, dog, horse

$\vec{x}$  is an image

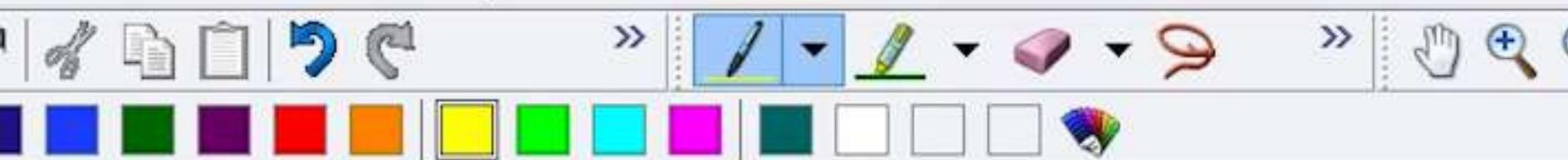
60x60 grayscale image

$\vec{x} \rightarrow \Sigma$



$$\sum_{k=1}^K \hat{y}_k = 1 - \text{Required}$$

$\hat{y}_j$  - 0 - prob  
0  
0

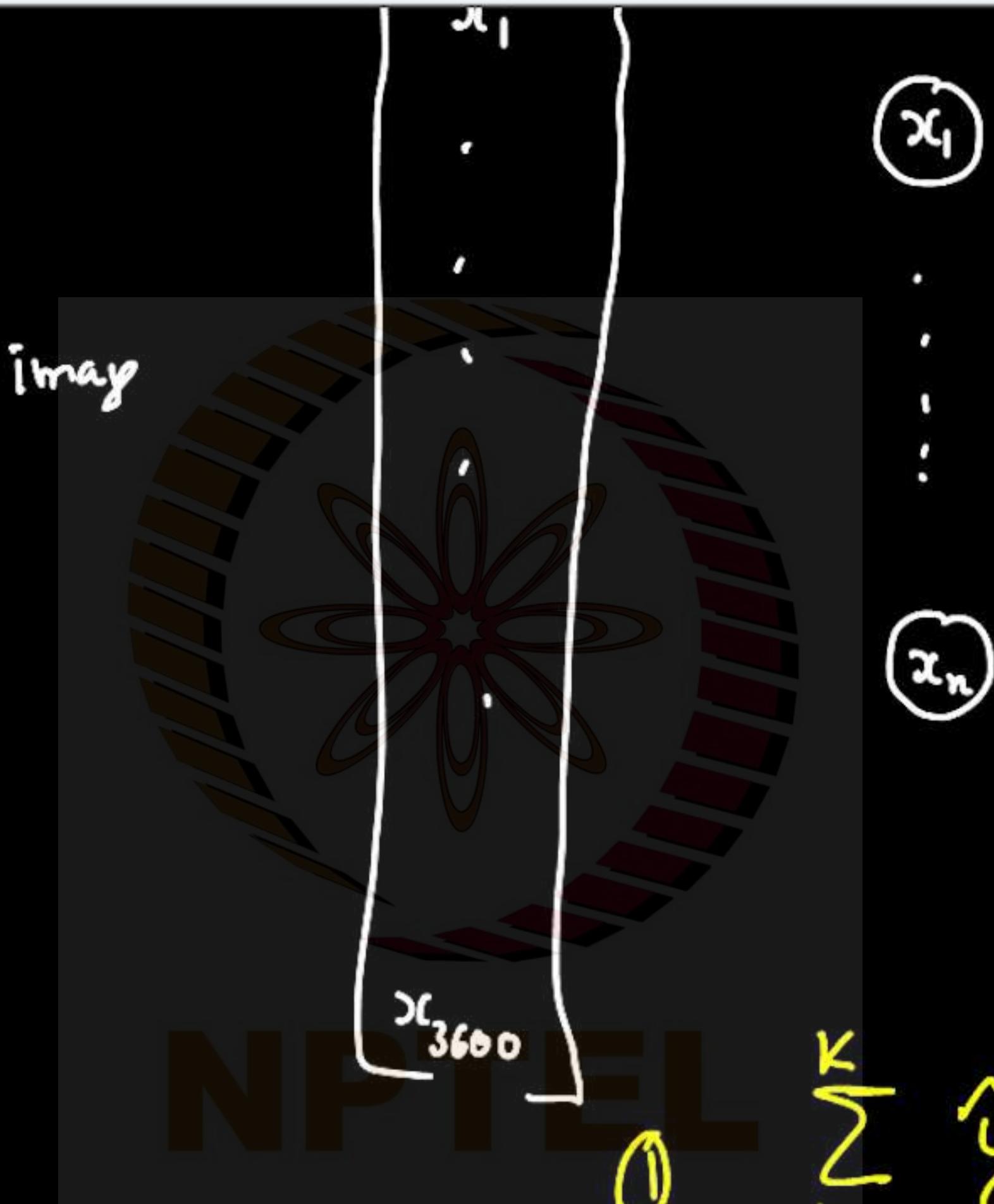


Cat, dog, horse

$\vec{x}$  is an image

60x60 grayscale image

$\vec{x} \rightarrow \Sigma$



$$\textcircled{1} \quad \sum_{k=1}^K \hat{y}_k = 1 - \text{Required}$$

$$\textcircled{2} \quad \text{All } \hat{y}_k \in [0, 1]$$



U.S. J.W.J.

L U J

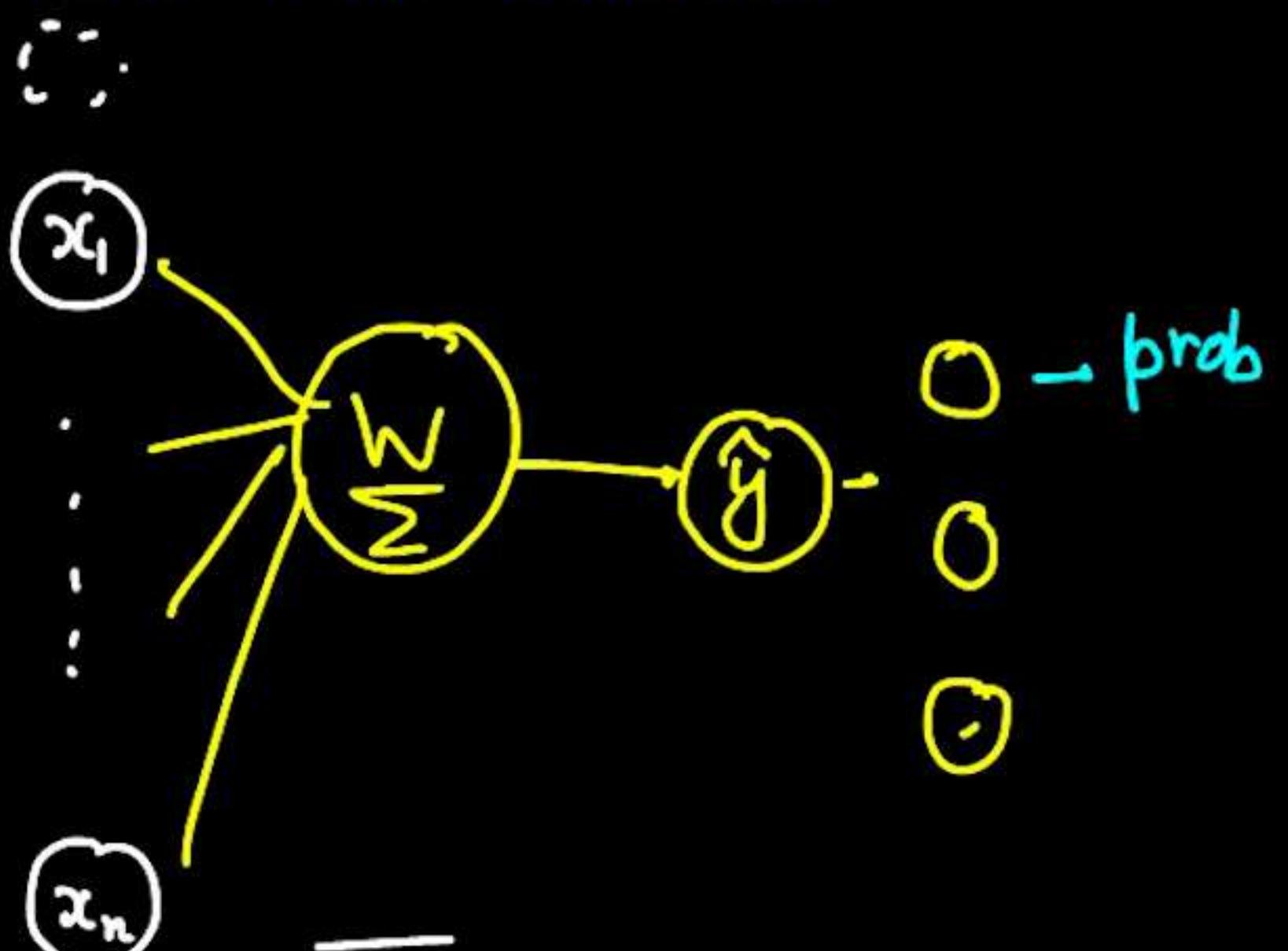
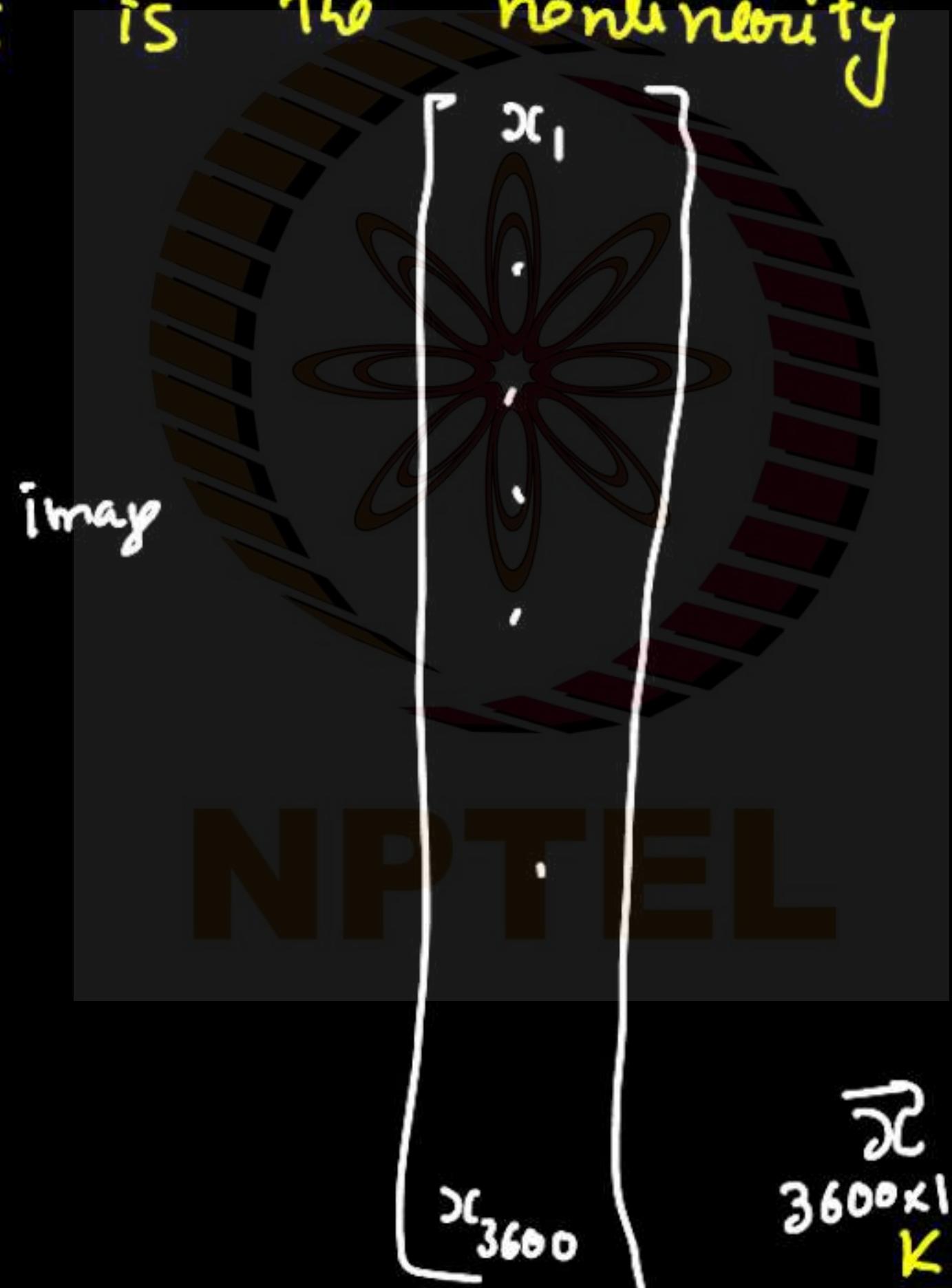
What is the nonlinearity that will achieve classification?

Cat, dog, horse

$\vec{x}$  is an image

60x60 grayscale image

$\vec{x} \rightarrow \Sigma$

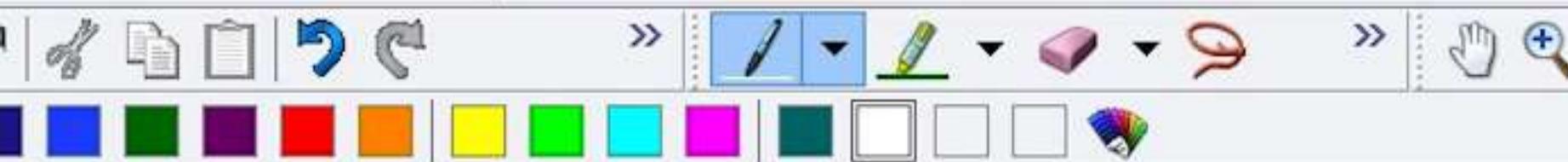


$\bar{W}$  matrix

$\vec{x} \xrightarrow{W} \vec{y}$   
 $3600 \times 1 \quad 3 \times 3600 \quad 3 \times 1$

$W_{3600 \times 3600 \times 1} \rightarrow \vec{y}$   
 $3 \times 1$

$$\sum \hat{y}_k = 1 - \text{Required}$$


 $\vec{x} \rightarrow \Sigma$ 

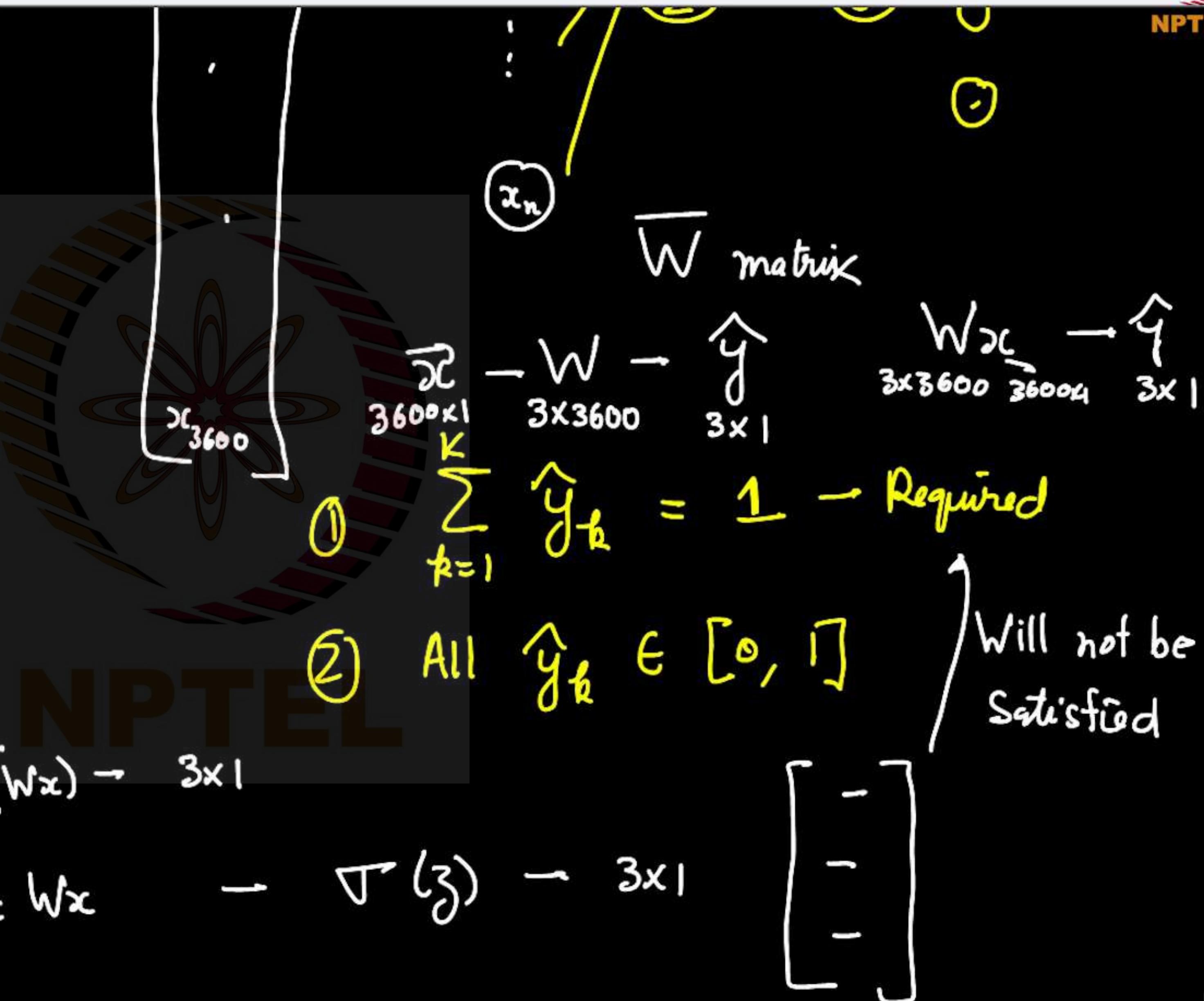
Logistic Regression

$\tau - [0, 1]$

$$\tau(wx) \rightarrow 3 \times 1$$

$$z = wx \rightarrow \tau(z) \rightarrow 3 \times 1$$

1





## Logistic Regression

$\sigma : [0, 1]$

$$\begin{bmatrix} \mathbf{x}_{3600} \end{bmatrix} \xrightarrow[3600 \times 1]{\text{Input}} \mathbf{W} \xrightarrow[3 \times 3600]{\text{Matrix}} \begin{bmatrix} \hat{y} \end{bmatrix} \xrightarrow[3 \times 1]{\text{Output}}$$

$\sum_{k=1}^K \hat{y}_k = 1 - \text{Required}$

② All  $\hat{y}_k \in [0, 1]$

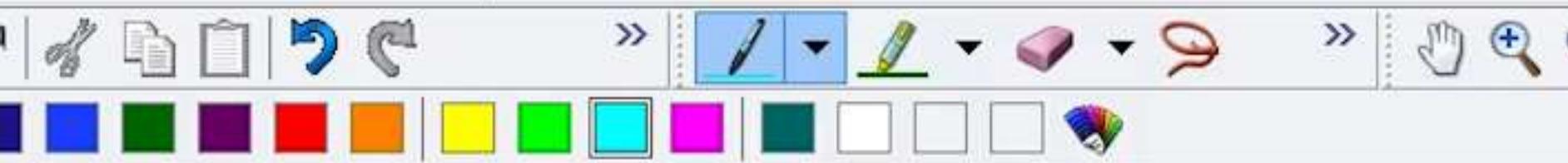
Will not be satisfied

$$\begin{aligned} \sigma(\mathbf{w}\mathbf{x}) &\rightarrow 3 \times 1 \\ \mathbf{z} = \mathbf{w}\mathbf{x} &\rightarrow \sigma(\mathbf{z}) \rightarrow 3 \times 1 \end{aligned}$$

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

Softmax function.

$$\text{Softmax}(\mathbf{z}_i) = \frac{\exp(\mathbf{z}_i)}{\sum_{j=1}^K \exp(\mathbf{z}_j)}$$



$$\sum_{j=1} \exp(\beta_j)$$

$$x \xrightarrow[W]{Wx} z$$

$x$   
 $3600 \times 1$

$z$   
 $3 \times 1$

$$\hat{y} = \text{Softmax} \left( \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \right) = \begin{bmatrix} \text{Softmax}(\beta_1) \\ \text{Softmax}(\beta_2) \\ \text{Softmax}(\beta_3) \end{bmatrix} = \begin{bmatrix} e^{\beta_1} \\ e^{\beta_1} \\ e^{\beta_3} \end{bmatrix} \frac{1}{\sum}$$



$$\sum_{j=1} \exp(\beta_j)$$

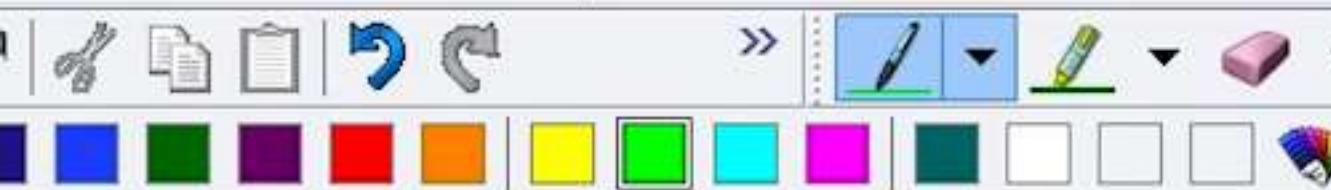
$$x \xrightarrow[W]{Wx} z$$

$3600 \times 1$        $3 \times 1$

$$\hat{y} = \text{Softmax} \left( \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \right) = \begin{bmatrix} \text{Softmax}(\beta_1) \\ \text{Softmax}(\beta_2) \\ \text{Softmax}(\beta_3) \end{bmatrix} = \begin{bmatrix} e^{\beta_1} \\ e^{\beta_2} \\ e^{\beta_3} \end{bmatrix} \frac{1}{D}$$

$$D = e^{\beta_1} + e^{\beta_2} + e^{\beta_3}$$

$$\sum_{k=1}^3 \hat{y}_k = 1$$



$y = \boxed{f^{-1}(z)}$

Binary Logistic

Regression

$k = 2$

Sigmoid

Multinomial

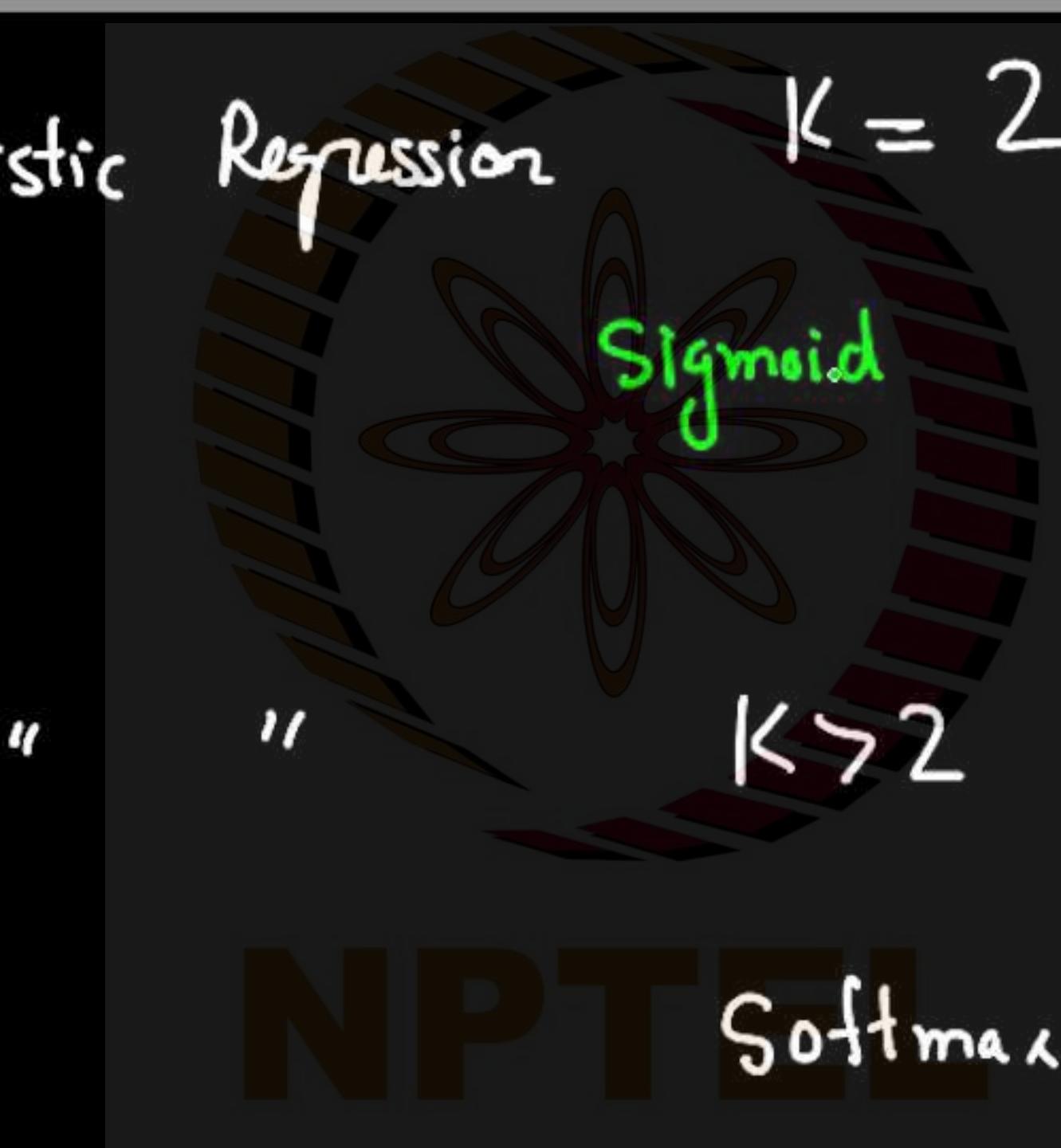
"

$k > 2$

$\hat{y}$  (One-hot vector)

$\hat{y}$  (0 or 1)

$$J = - [y \ln \hat{y} + (1-y) \ln (1-\hat{y})]$$





## Binary Logistic Regression

$K = 2$

$\hat{y}$  (0 or 1)

$$J = -[y \ln \hat{y} + (1-y) \ln(1-\hat{y})]$$

Multinomial "

Sigmoid

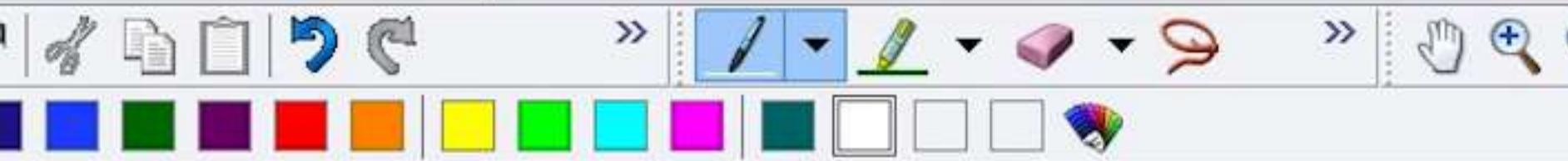
"  $K > 2$

Softmax

$\hat{y}$  (One-hot vector)

Loss function for  $K > 2$ ?

$$J = - \sum_{k=1}^K y_k \ln \hat{y}_k$$


 $K = 2$ 

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

$$\hat{y}_2 = (1 - \hat{y}_1)$$

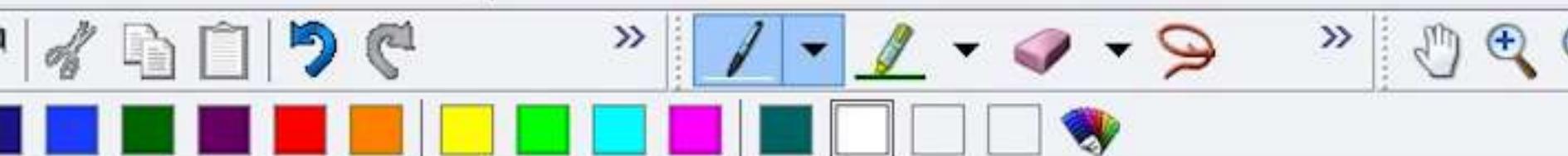
$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

$$y_2 = (1 - y_1)$$

$$J = - \sum_{k=1}^2 y_k \ln \hat{y}_{k\alpha} = - [y_1 \ln \hat{y}_1 + y_2 \ln \hat{y}_2]$$

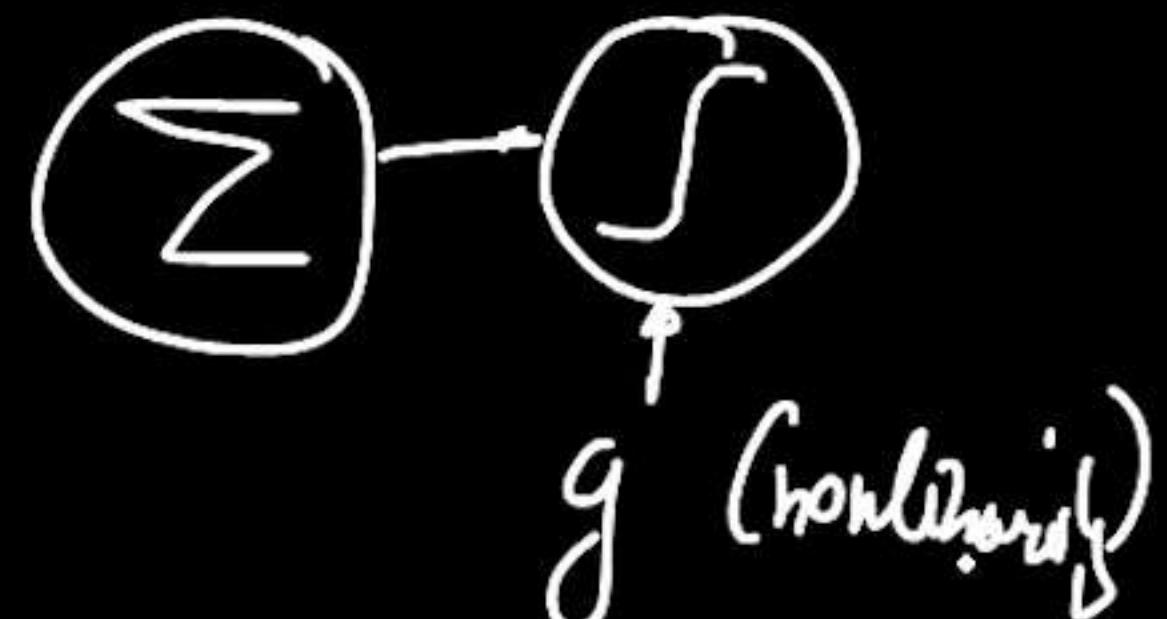
$$= - [y_1 \ln \hat{y}_1 + (1 - y_1) \ln (1 - \hat{y}_1)]$$

↑  
Binary Cross-Entropy



Multinomial

"

 $K > 2 \quad \hat{y}$  (One-hot vector)

Loss function for

 $K > 2 ?$ 

Softmax

General

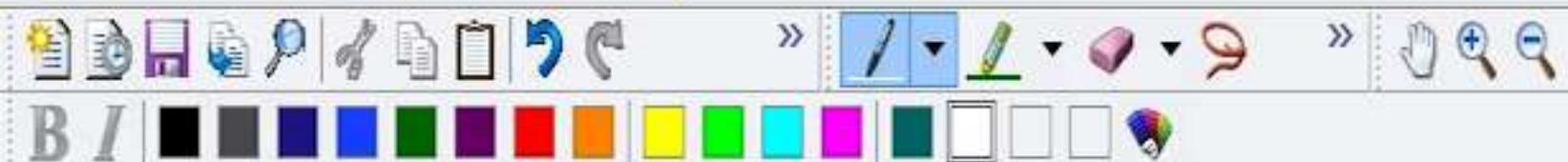
Cross - Entropy  
loss function.  
for  $K > 2$ .

$$J = - \sum_{k=1}^K y_k \ln \hat{y}_k$$

 $K = 2$ 

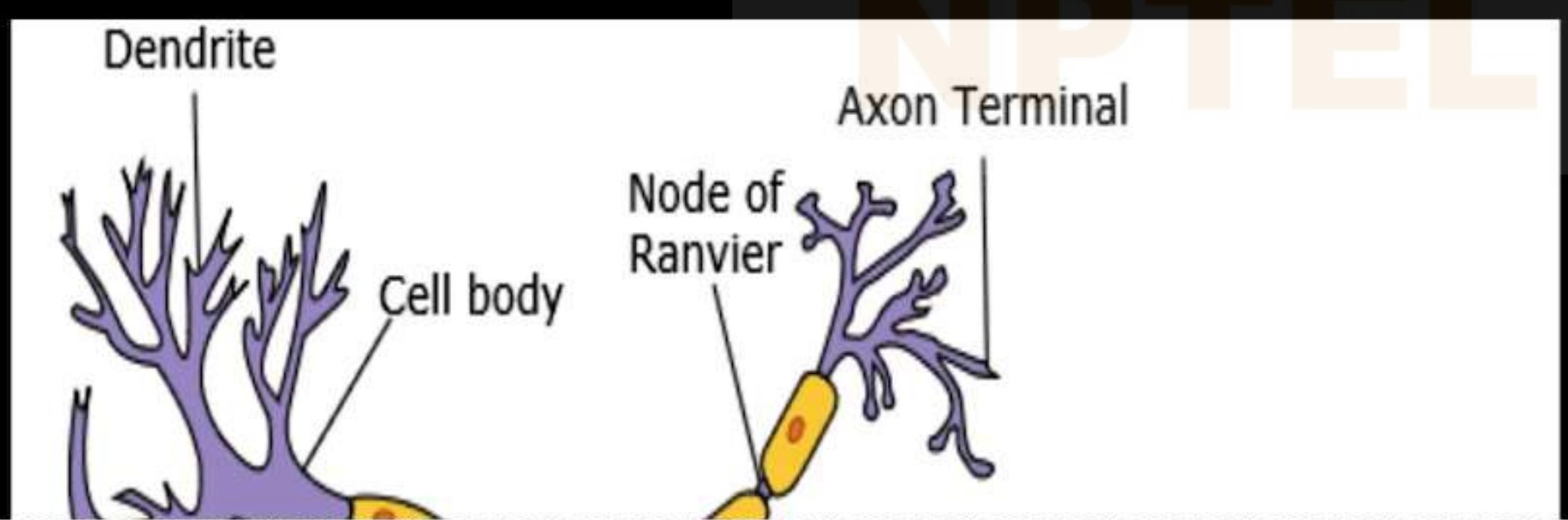
$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \end{bmatrix}$$

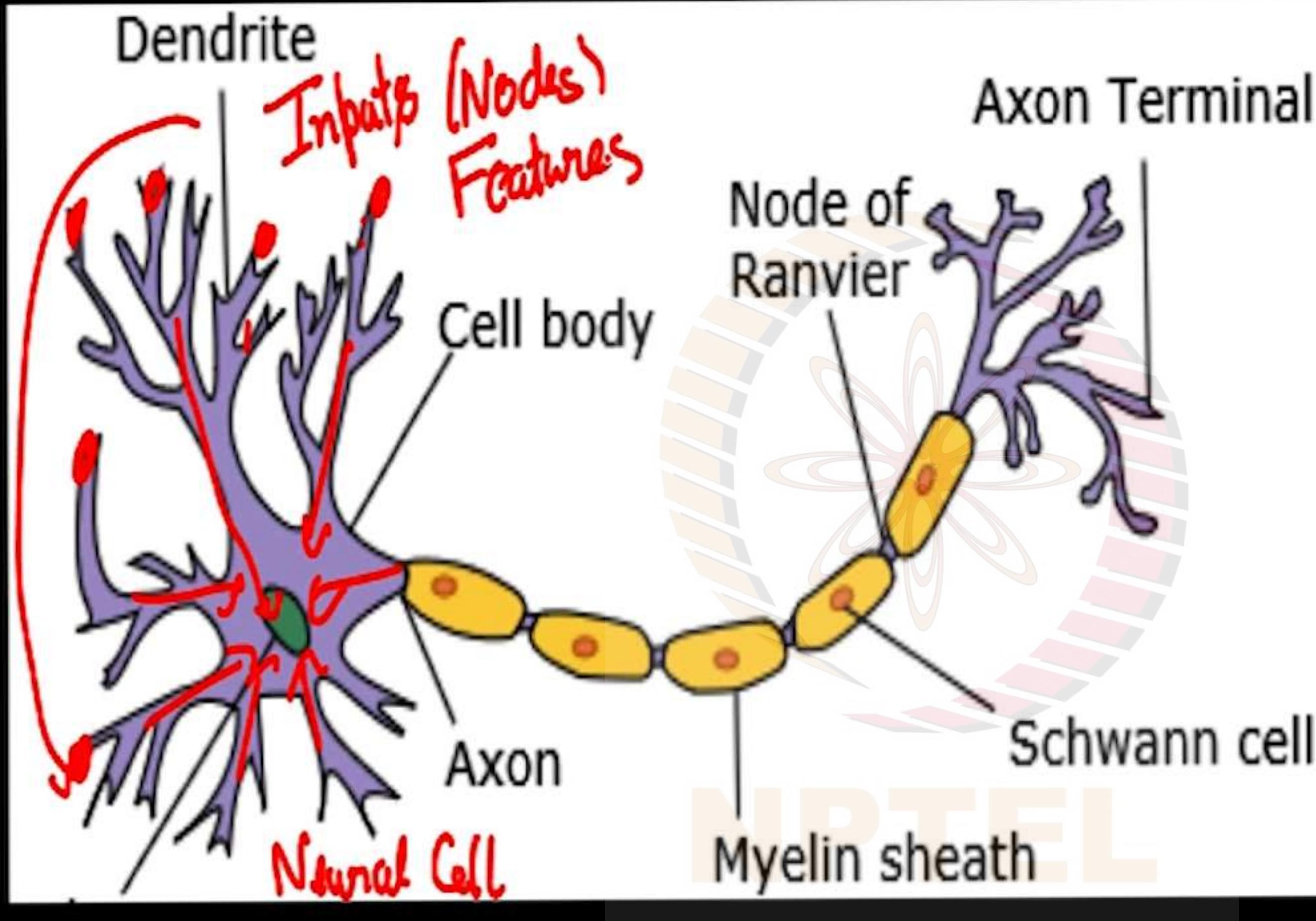
$$y = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$

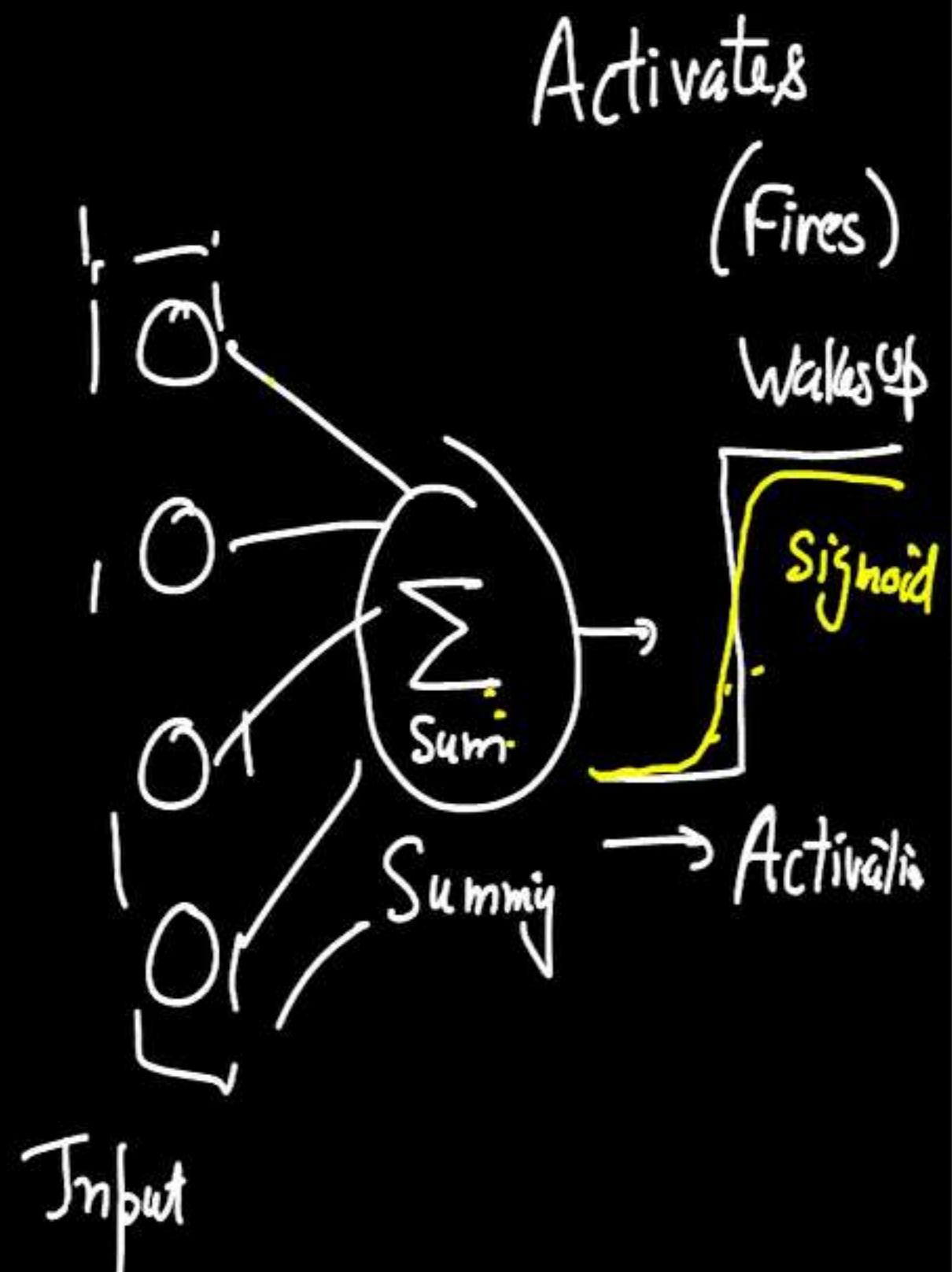
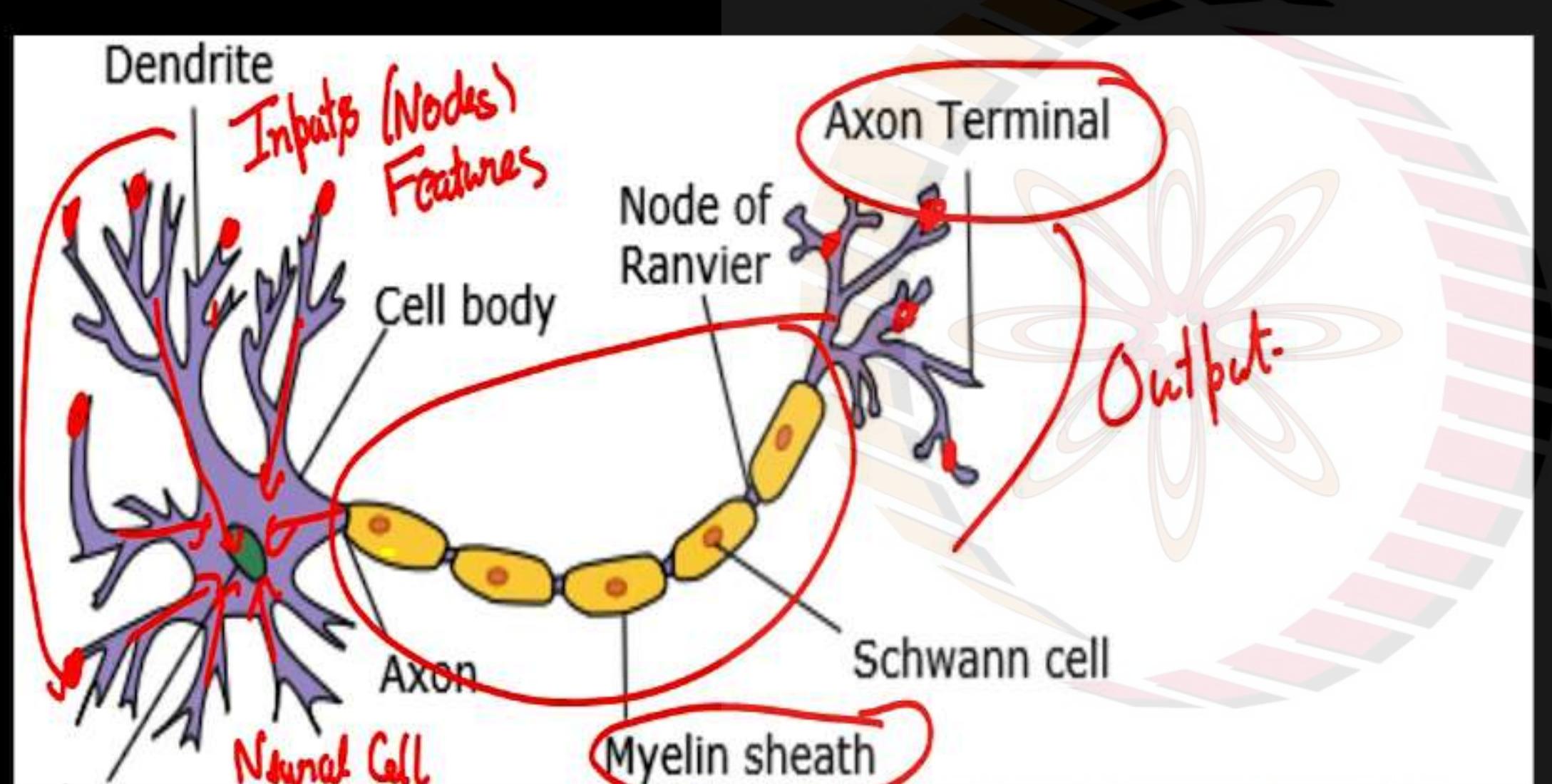


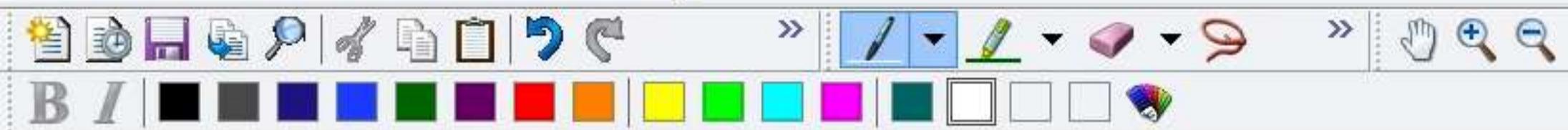
Note Title

# Biological Neuron







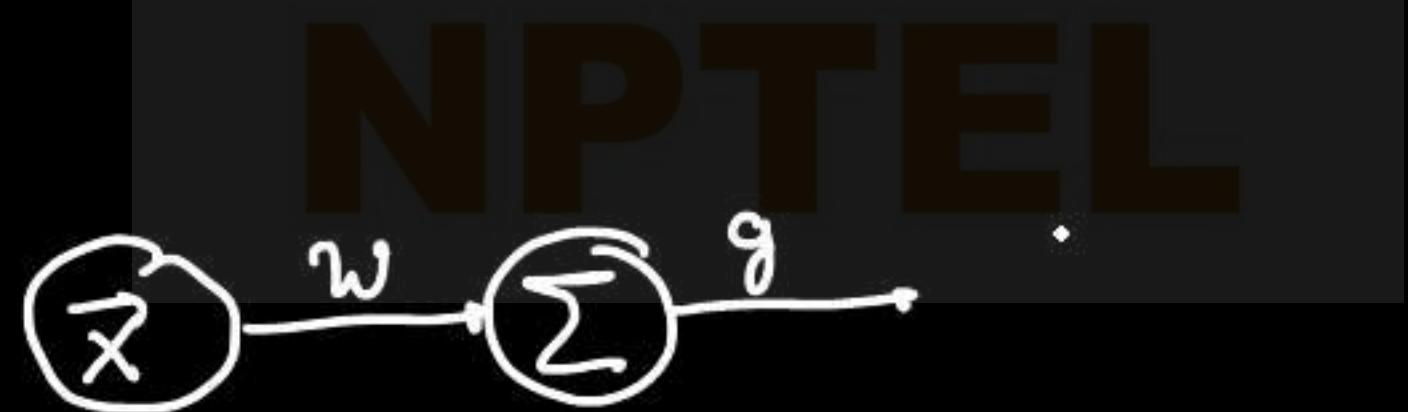
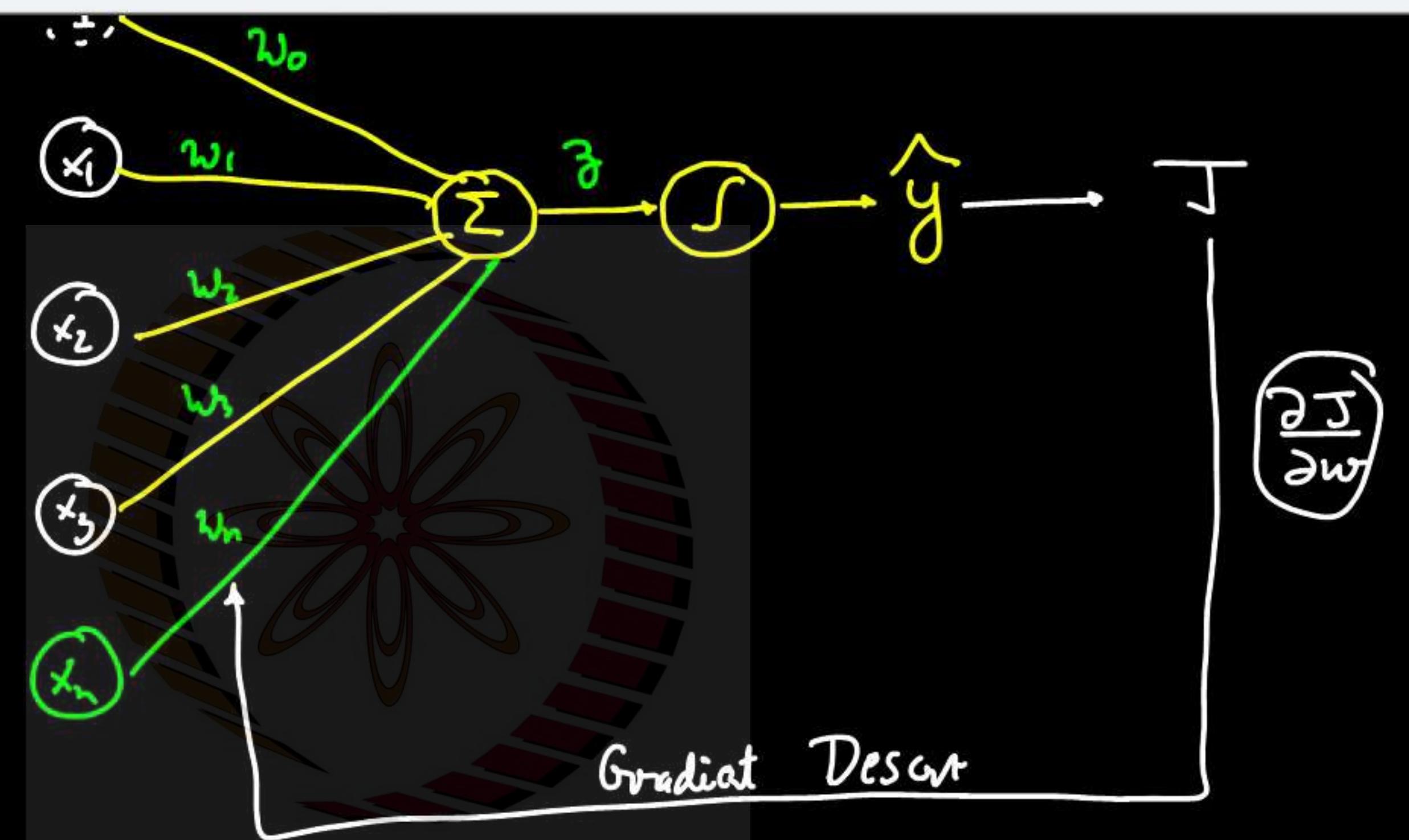
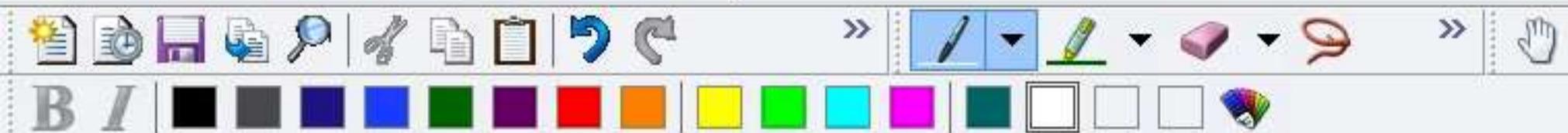


Note Title

Gradient of Logistic Regression

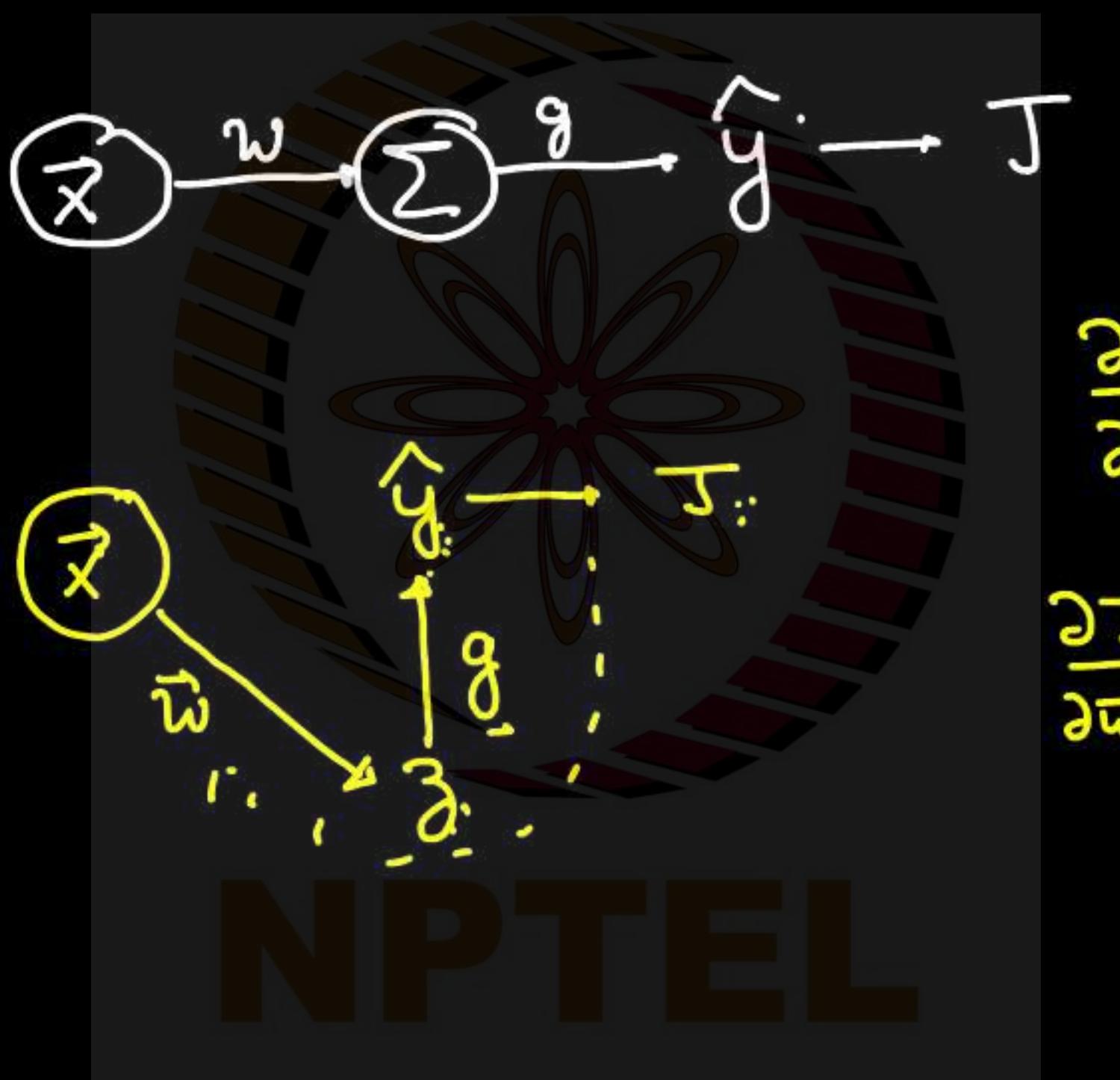
OR, AND, etc → Guessing for weights

NPTEL



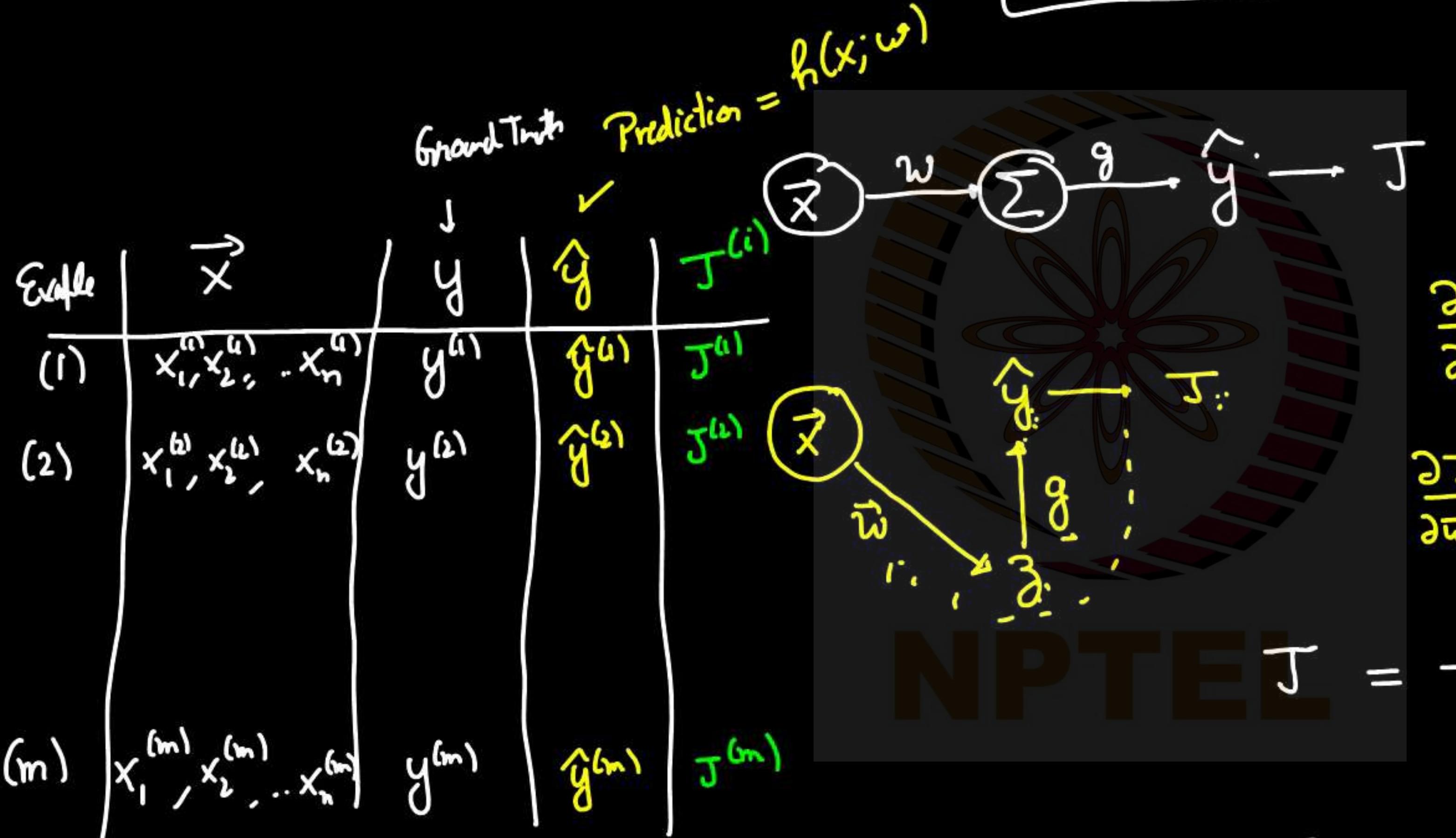


## Gradiat Descent



$$\frac{\partial J}{\partial \vec{w}} = ?$$

$$\frac{\partial J}{\partial \vec{w}} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial g} \frac{\partial g}{\partial \vec{w}} \rightarrow \text{Chain Rule}$$



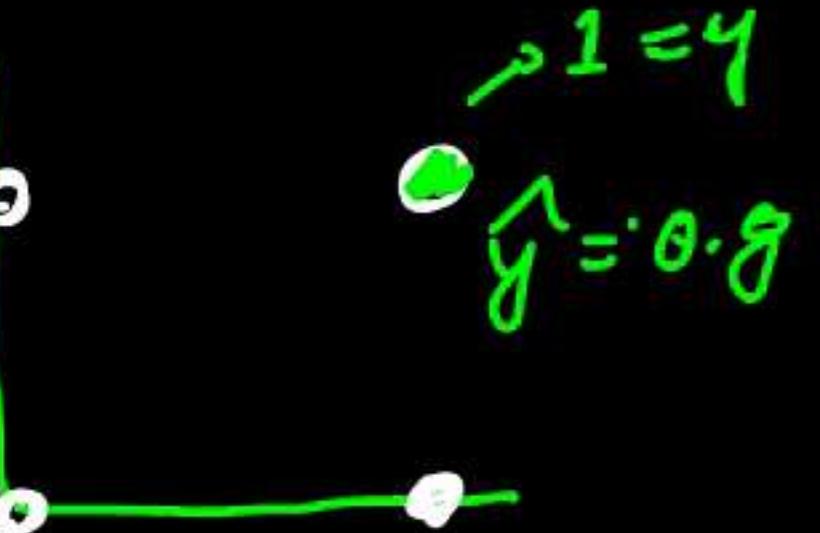
Gradient Descent

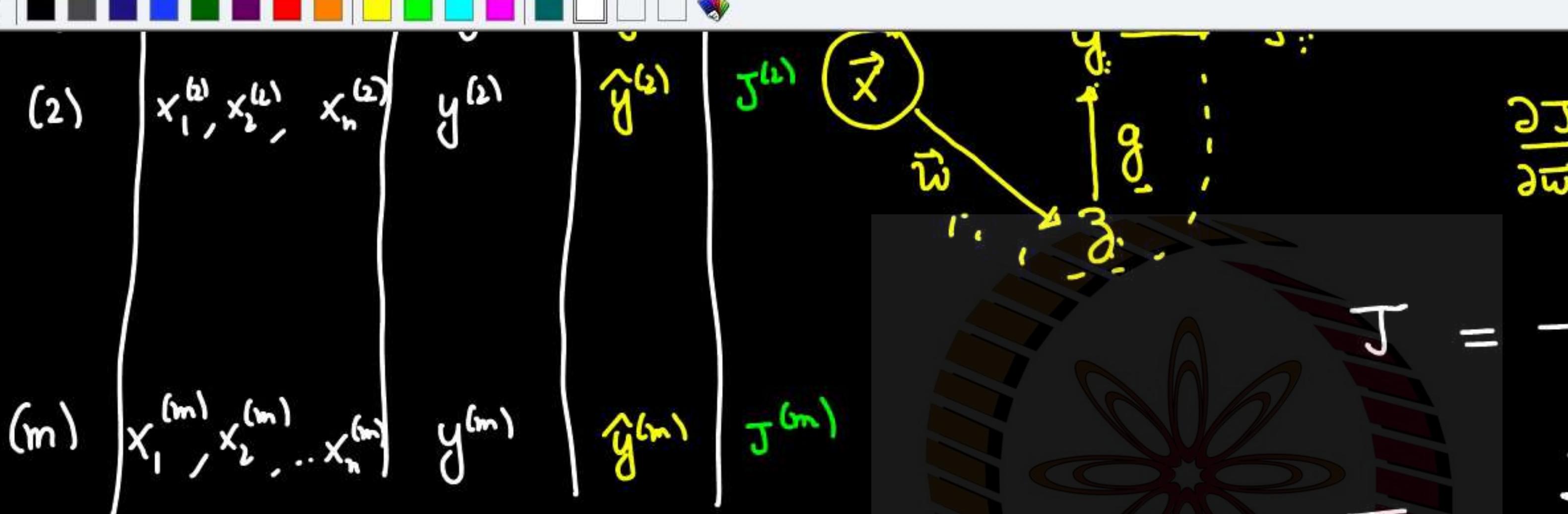
$$\frac{\partial J}{\partial w} = ?$$

$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w} \rightarrow \text{Chain Rule}$$

$$J = - \sum_{i=1}^m y^{(i)} \ln \hat{y}^{(i)} + (1-y^{(i)}) \ln (1-\hat{y}^{(i)})$$

$$J = - \sum_{i=1}^m J^{(i)}$$





$$\frac{\partial J}{\partial w} = \frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w} \rightarrow \text{Chain Rule}$$

$$J = - \sum_{i=1}^m y_i \ln \hat{y}^{(i)} + (1-y_i) \ln (1-\hat{y}^{(i)})$$

$$J = \sum_{i=1}^m J^{(i)}$$

$$\sum_{i=1}^m \frac{\partial J^{(i)}}{\partial w}$$

$$\frac{\partial J^{(i)}}{\partial w} = \frac{\partial J^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z} \frac{\partial z}{\partial w}$$



File Edit View Insert Actions Tools Help



NPTEL

$$\frac{\partial J}{\partial \vec{w}} = \sum_{i=1}^m \frac{\partial J^{(i)}}{\partial \vec{w}}$$

$$\frac{\partial J^{(i)}}{\partial \vec{w}} = \frac{\partial J^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z} \frac{\partial z}{\partial \vec{w}}$$

$$\hat{y} = \sigma(z)$$

$$J = - \left\{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right\}$$

$$\frac{\partial J}{\partial \hat{y}} = - \left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{1-\hat{y}} \right\}$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$= \hat{y}(1-\hat{y})$$



$$J = -\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\}$$

$$\frac{\partial J}{\partial \hat{y}} = -\left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\}$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z))$$

$$= \hat{y}(1-\hat{y})$$

$$\underbrace{\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z}}_{\text{- Error}} = -\left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\} \hat{y}(1-\hat{y})$$

$$\frac{\partial J}{\partial z} = -\underbrace{\{y - \hat{y}\}}_{\text{- Error}} \quad (\text{Check!})$$



$$\frac{\partial J}{\partial \vec{w}} = \sum_{i=1}^m \frac{\partial J^{(i)}}{\partial \vec{w}}$$

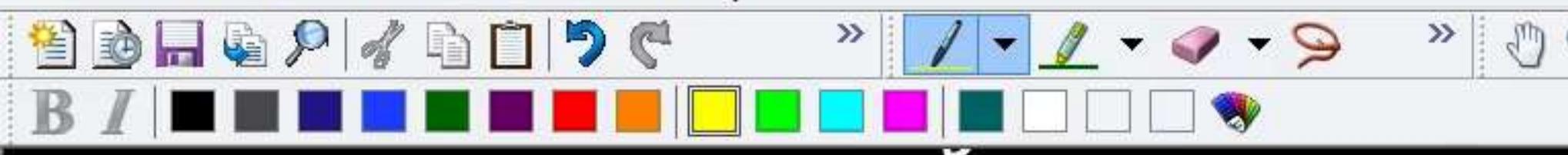
$$\frac{\partial J^{(i)}}{\partial \vec{w}} = \frac{\partial J^{(i)}}{\partial \hat{y}^{(i)}} \frac{\partial \hat{y}^{(i)}}{\partial z} \underbrace{\frac{\partial z}{\partial \vec{w}}}_{\text{Includes } w_0}$$

$$J = - \left\{ y \ln \hat{y} + (1-y) \ln (1-\hat{y}) \right\}$$

$$\frac{\partial J}{\partial \hat{y}} = - \left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\}$$

$$\begin{aligned} \frac{\partial \hat{y}}{\partial z} &= \frac{\partial \sigma(z)}{\partial z} = \sigma(z)(1-\sigma(z)) \\ &= \hat{y}(1-\hat{y}) \end{aligned}$$

$$\frac{\partial J}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} = - \left\{ \frac{y}{\hat{y}} - \frac{(1-y)}{(1-\hat{y})} \right\} \hat{y}(1-\hat{y})$$



0 0

$$\underbrace{\frac{\partial J}{\partial \hat{y}}}_{\text{Error}} \frac{\partial \hat{y}}{\partial \beta} = - \left\{ \frac{y - \hat{y}}{\hat{y}} \right\} \hat{y}^{(1-\hat{y})}$$

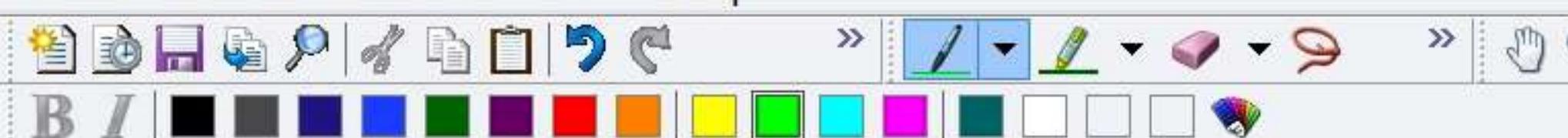
$$\frac{\partial J}{\partial \beta} = - \left\{ y - \hat{y} \right\} \quad (\text{Check !})$$

- Error

$\frac{\partial \beta}{\partial \vec{w}}$  where  $\beta = \vec{x}^T \vec{w}$ .

$$\Rightarrow \boxed{\frac{\partial \beta}{\partial \vec{w}} = \vec{x}}$$

NPTEL



$$\frac{\partial J}{\partial \beta_3} = - \{ y - \hat{y} \} \quad (\text{Check!})$$

- Error

$$\Rightarrow \boxed{\frac{\partial J}{\partial \vec{w}} = \vec{x}}$$

$$\Rightarrow \frac{\partial J^{(i)}}{\partial \vec{w}} = \frac{\partial J}{\partial \beta_3} \cdot \frac{\partial \beta_3}{\partial \vec{w}} = - \{ y^{(i)} - \hat{y}^{(i)} \} \vec{x}$$

= - Error · Input

Exactly the  
Same as Linear Reg  
with Least Squares

$$\frac{\partial J}{\partial \vec{w}} = - \sum_{i=1}^m \{ y^{(i)} - \hat{y}^{(i)} \} \vec{x}.$$

Logistic Regression  
with  
Binary Cross-Entropy



Exactly the  
Same as Linear Reg  
with Least Squares

$$\frac{\partial J}{\partial \vec{w}} = -\sum_{i=1}^m \{y^{(i)} - \hat{y}^{(i)}\} \vec{x}.$$

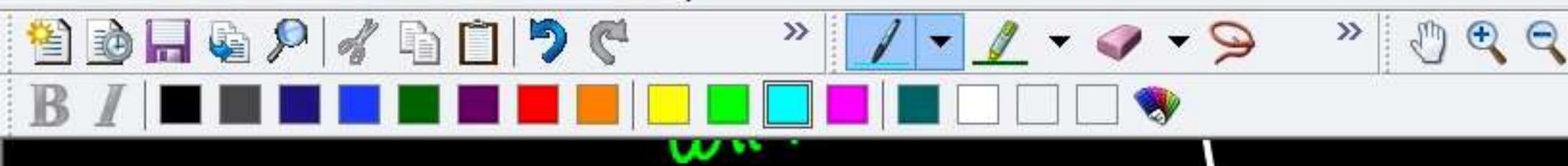
= - Error · Input

Logistic Regression  
with  
Binary Cross-Entropy

For Logistic Regression, the loss function  $J$  is not  
Convex

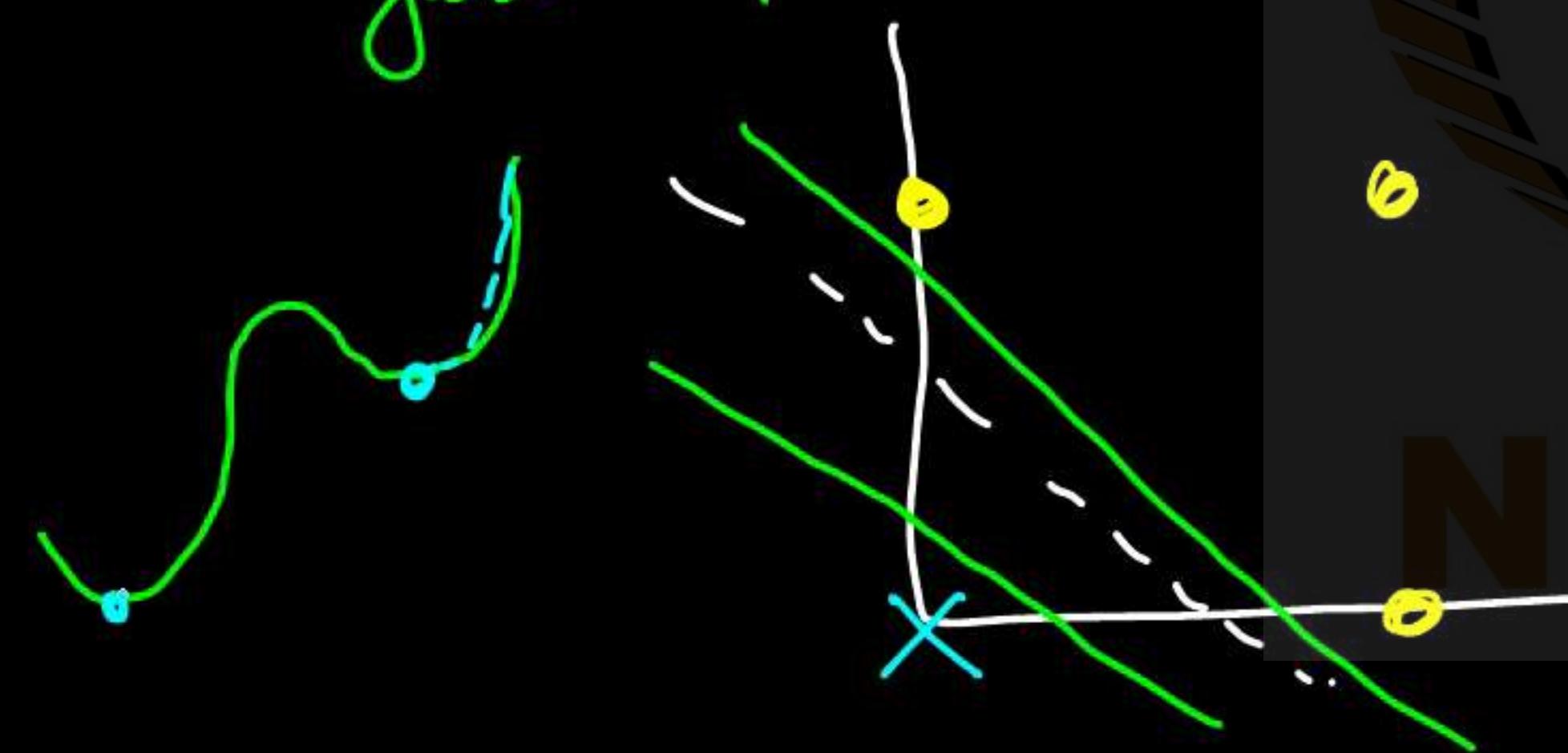
i.e. the minimum is, in general, not unique

File Edit View Insert Actions Tools Help



For Logistic Regression, the loss function  $J$  is not convex

Does not necessarily give global optimum

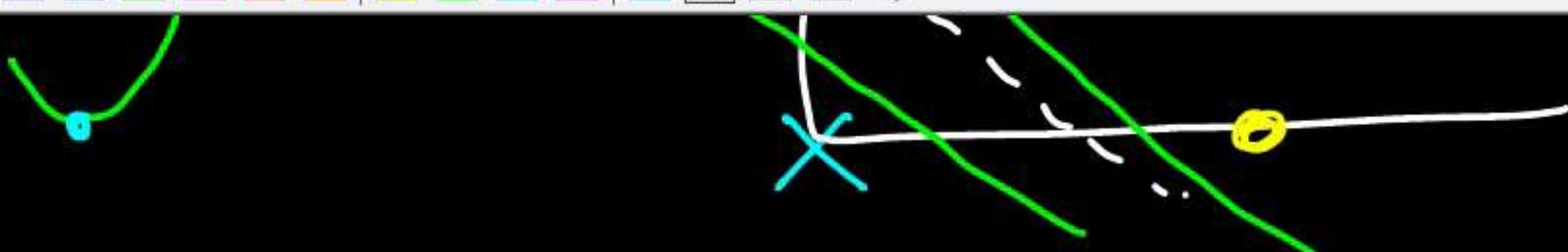


i.e. the minimum

is, in general, not unique

Many possible classifying lines

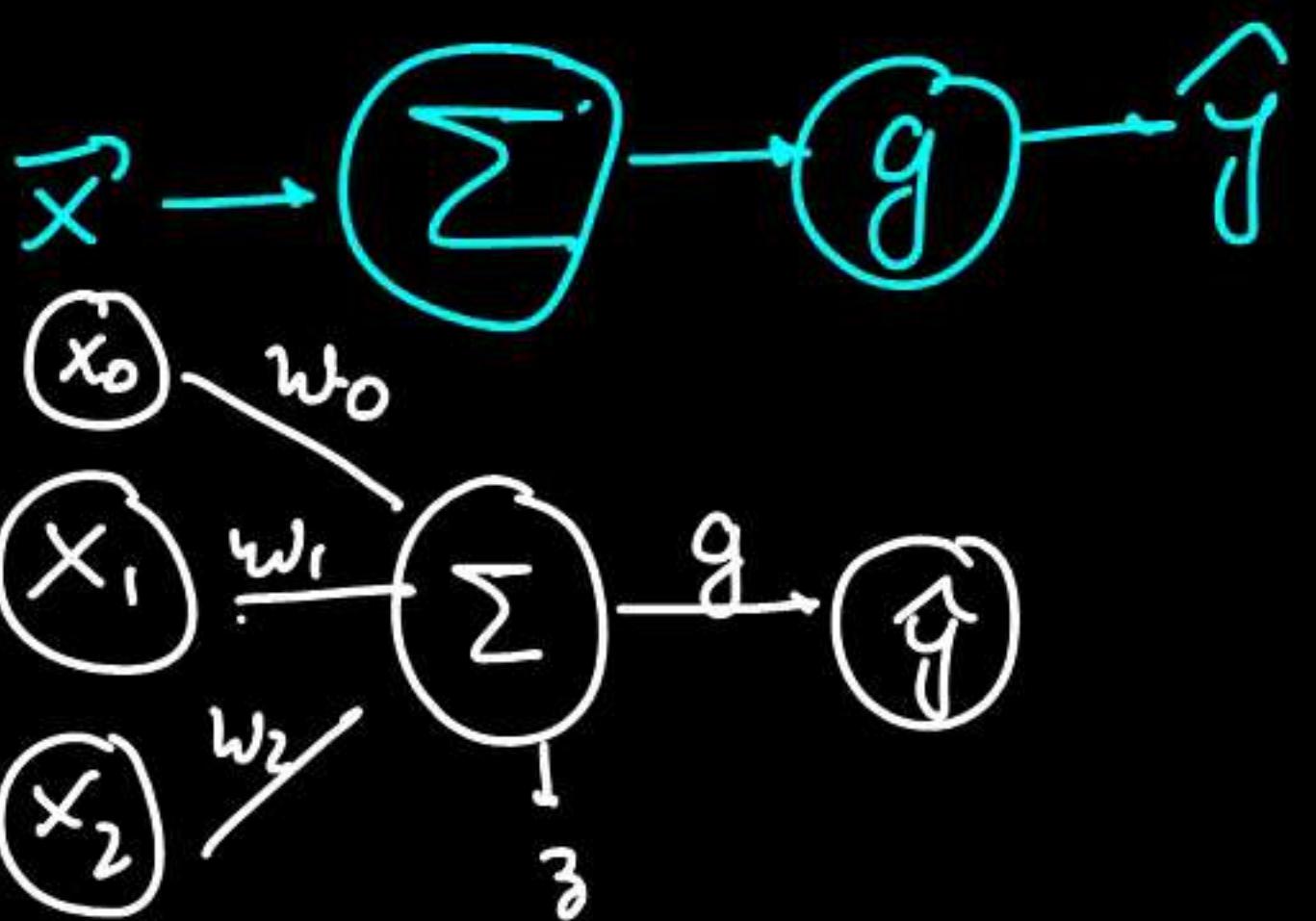
→ Logistic Regression does not give unique solutions



→ Logistic Regression does not give unique solutions

	$\vec{x}$	$y$
1	(0, 0)	0
2	(0, 1)	1
3	(1, 0)	1
4	(1, 1)	1

$$x_1^2, x_2^2$$



$$\begin{aligned} z &= \vec{w} \cdot \vec{x} \\ y &= g(z) \end{aligned} \quad \left. \begin{array}{l} \text{vectorized exprn.} \\ \text{ } \end{array} \right\}$$



HOME PLOTS APPS LIVE EDITOR INSERT VIEW

New Save Find Files Go To ↑ Text Aa Heading 1 B I U M Code % X Run Section Run and Advance Section Break Run to End Run Step Stop

FILE NAVIGATE TEXT CODE SECTION RUN

MATLAB Drive

CURRENT FOLDER LogisticModel.m ORgate.m GeneralizedLogisticRegression mlx +

Name

- NPTEL
- Published (my site)
- GeneralizedLinearRegression.ml
- LinearRegressionPolyFit.mlx
- LinearRegressionSimple.mlx

WORKSPACE

NAME	VALUE	SIZE	CLASS
------	-------	------	-------

**Generalized Function for Logistic Regression**

This function takes in the data and fits it using linear regression for some given learning rate and stopping criterion

**Inputs**

$x$  : The input vector/matrix. We assume that  $x$  is a vector of size  $n$  i.e  
 $x = [x_1, x_2, \dots, x_n]$

$y$ : The output vector (Ground truth)

$\alpha$  : The learning rate  $\alpha$

$\epsilon$  : The stopping criterion  $\epsilon$

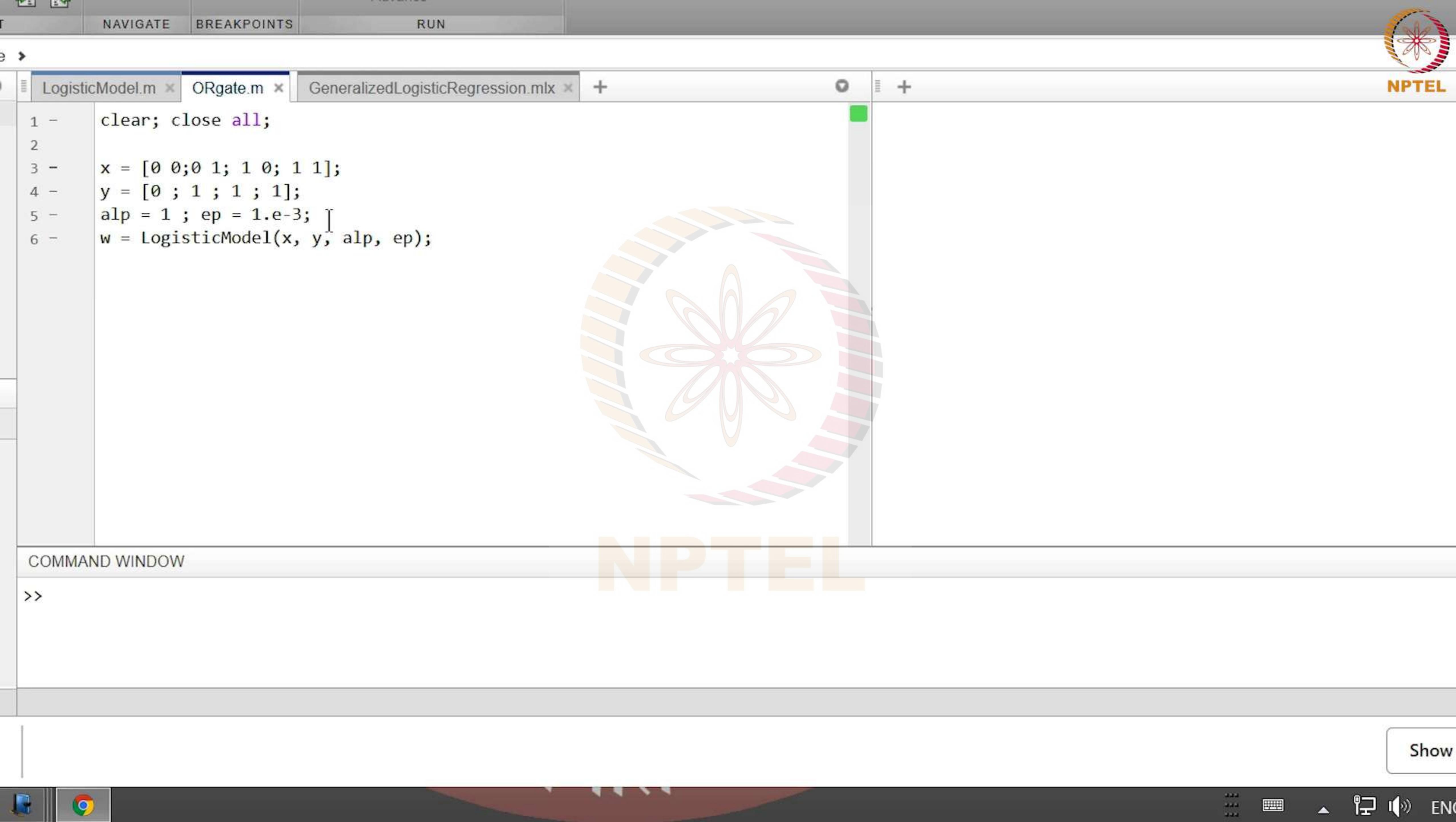
**Output**

$w$  : The weight vector  $w = [w_0, w_1, w_2, \dots, w_n]$

COMMAND WINDOW

>>

Note1.jnt Show all



Show

LogisticModel.m x ORgate.m x

GeneralizedLogisticRegression mlx x

```
1 function w = GeneralizedLogisticRegression(x, y, alp, ep)
```

Determine size of incoming data

m : Number of examples

n : Number of features

```
2 [m,n] = size(x);
```

Make initial guess for the w vector  $w = [w_0, w_1, w_2, \dots, w_n]^T$

```
[
```

```
3 w = rand(n+1,1);
```

## Iterating for w using gradient descent

The hypothesis function is now

COMMAND WINDOW

&gt;&gt;

Show

LogisticModel.m x ORgate.m x

GeneralizedLogisticRegression mlx x

The hypothesis function is now

$$\hat{y} = \sigma(w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n) = w \bullet x = \sigma(w^T x)$$

Define the sigmoid first

```
4 sigmoid = @(z) 1./(1+exp(-z));  
5  
6 for i = 1:m  
7 z = [1 x(i,:)]*w; %[1 x] is augmenting x with x0 = 1 to achieve  
8 yh(i,1) = sigmoid(z);  
9 end
```

**Exercise :** You can also try out an alternate, vectorized form of the above calculation without the for loop

```
10 J(1) = -1*sum(y.*log(yh)+(1.-y).*log(1.-yh)); %Initial value of J -  
err is a variable which denotes the appropriate error variable which we wish to
```

COMMAND WINDOW

&gt;&gt;

Show

LogisticModel.m x ORgate.m x

GeneralizedLogisticRegression mlx x

The gradient has  $(n+1)$  components. That is  $\nabla_w J$  is a  $(n+1) \times 1$  vector

$$\frac{\partial J}{\partial w_0} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})$$

$$\frac{\partial J}{\partial w_1} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_1^{(i)}$$

...

$$\frac{\partial J}{\partial w_r} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_r^{(i)}$$

...

$$\frac{\partial J}{\partial w_n} = \frac{1}{m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)}) x_n^{(i)}$$

16

```
nred_err = vh - v:
```

COMMAND WINDOW

&gt;&gt;

Show

LogisticModel.m x ORgate.m x

GeneralizedLogisticRegression mlx x

```
29
30     err = J(1iter)-J(1iter-1));
31     %err = norm(alp*DJ);
```

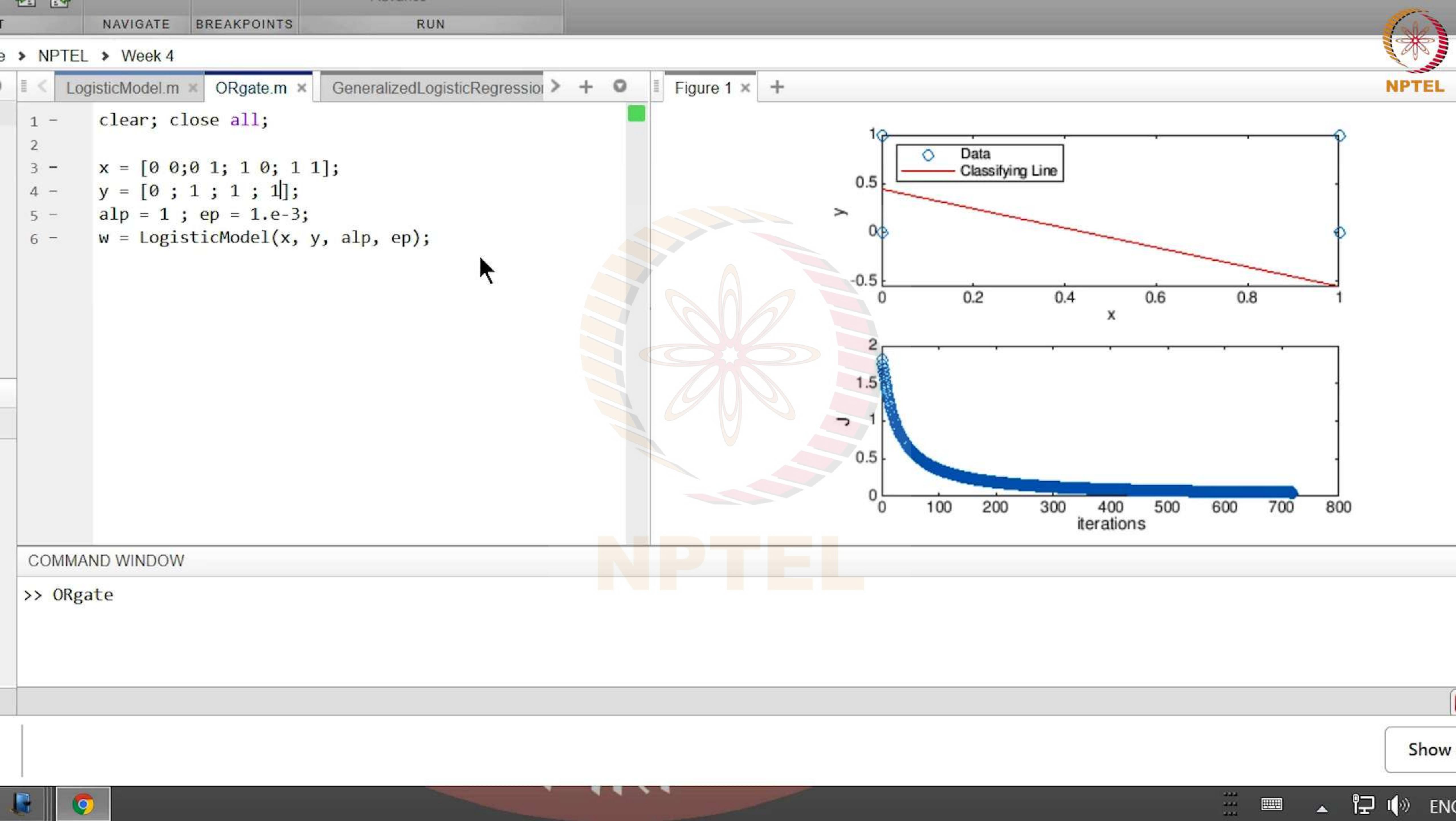
## Create Plots

```
32 subplot(211)
33 plot(x(:,1),y, 'o',x(:,1),yh, 'rx');
34 xlabel('x')
35 ylabel('y')
36 legend({'Data','Logistic'},'Location', 'NorthWest')
37 drawnow;
38 subplot(212)
39 plot([1:iter],J([1:iter]),'-o');
40 xlabel('iterations');
41 ylabel('J')
42 drawnow;
```

COMMAND WINDOW

&gt;&gt;

Show

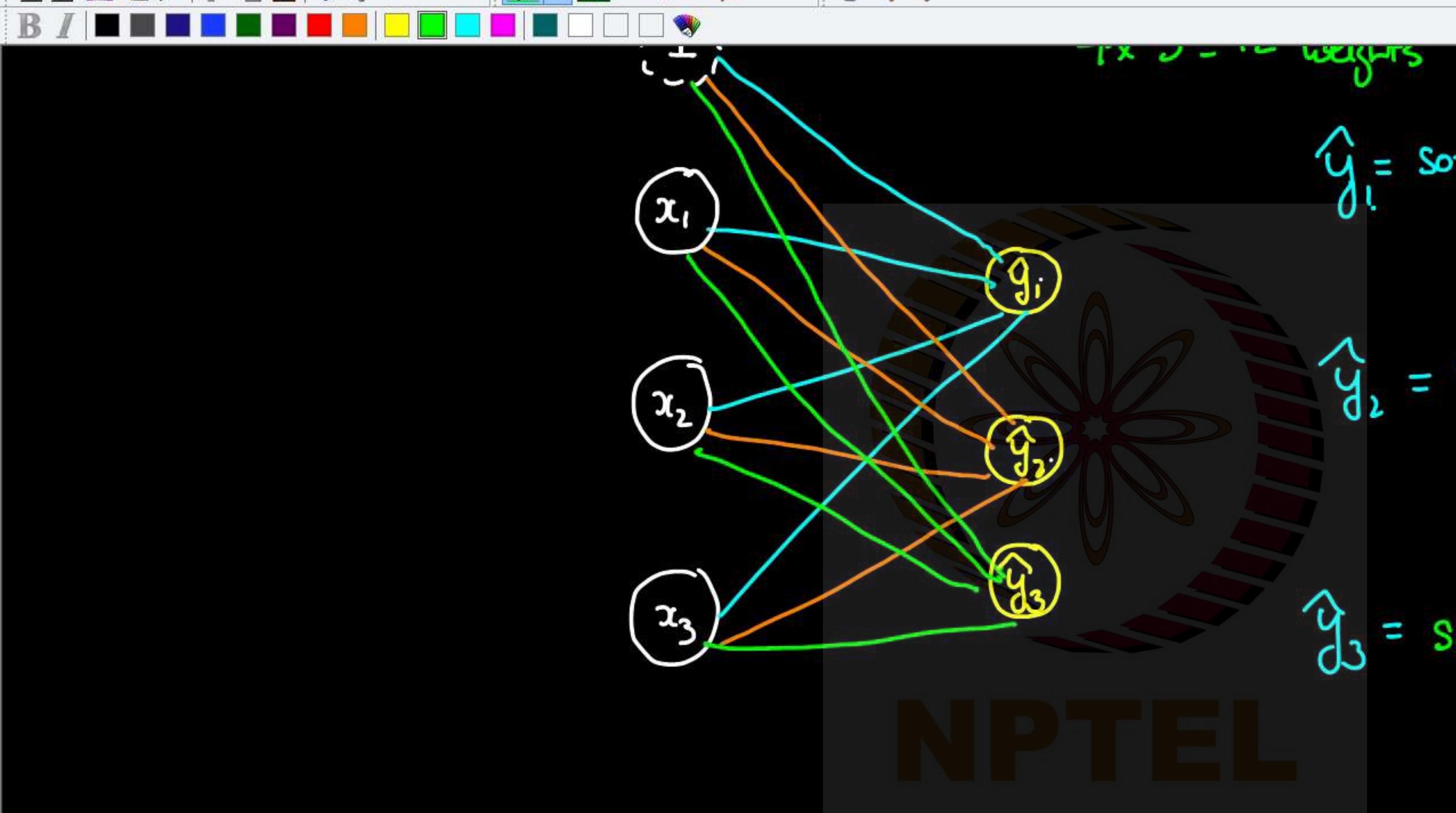




File Edit View Insert Actions Tools Help



NPTEL



$$\hat{y}_1 = \text{softmax} (\omega_{01} + \omega_{11} x_1 + \omega_{21} x_2 + \omega_{31} x_3)$$

$$\hat{y}_2 = \text{softmax} (\omega_{02} + \omega_{12} x_1 + \omega_{22} x_2 + \omega_{32} x_3)$$

$$\hat{y}_3 = \text{softmax} (\omega_{03} + \omega_{13} x_1 + \omega_{23} x_2 + \omega_{33} x_3)$$



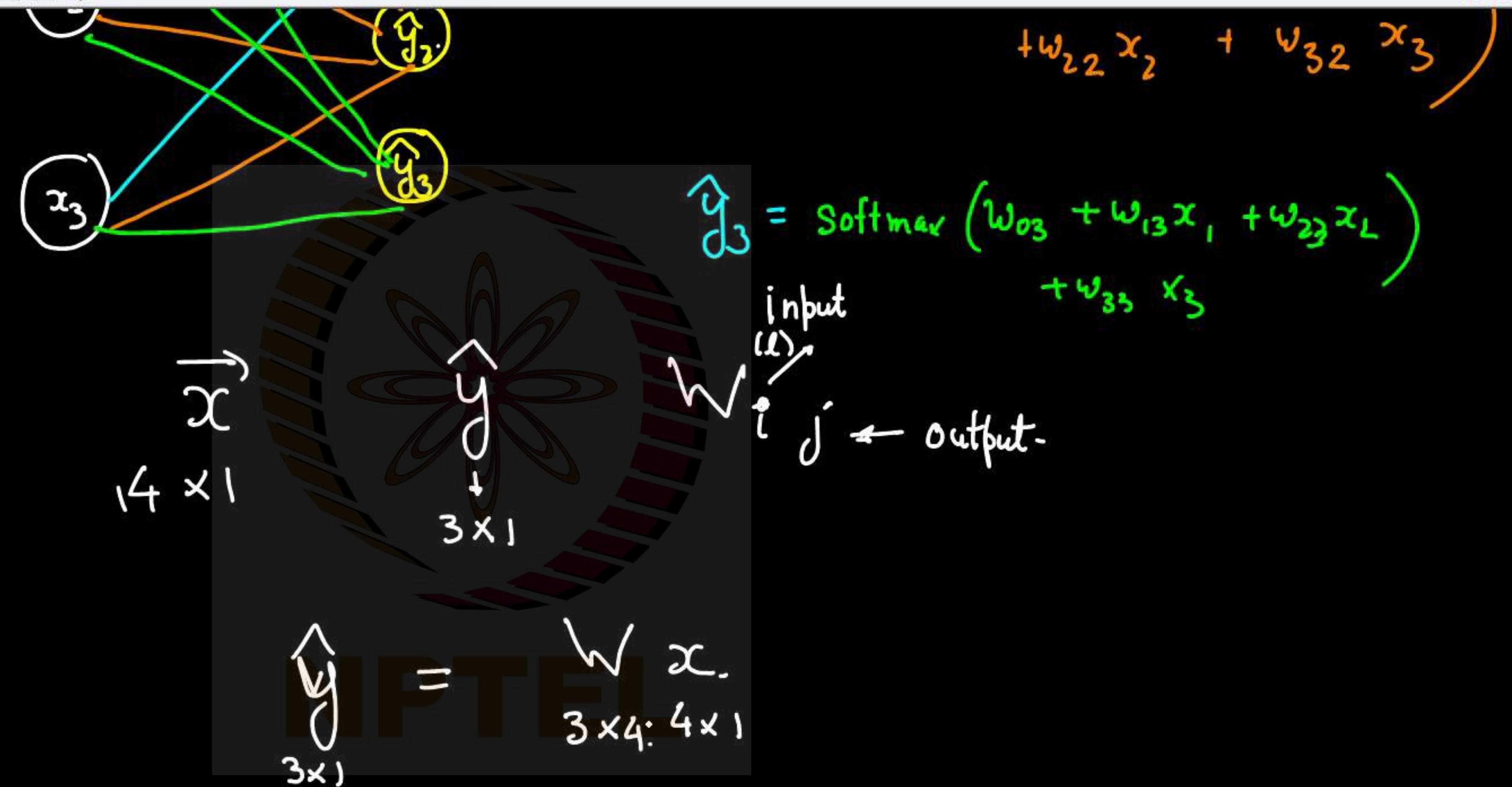
File Edit View Insert Actions Tools Help



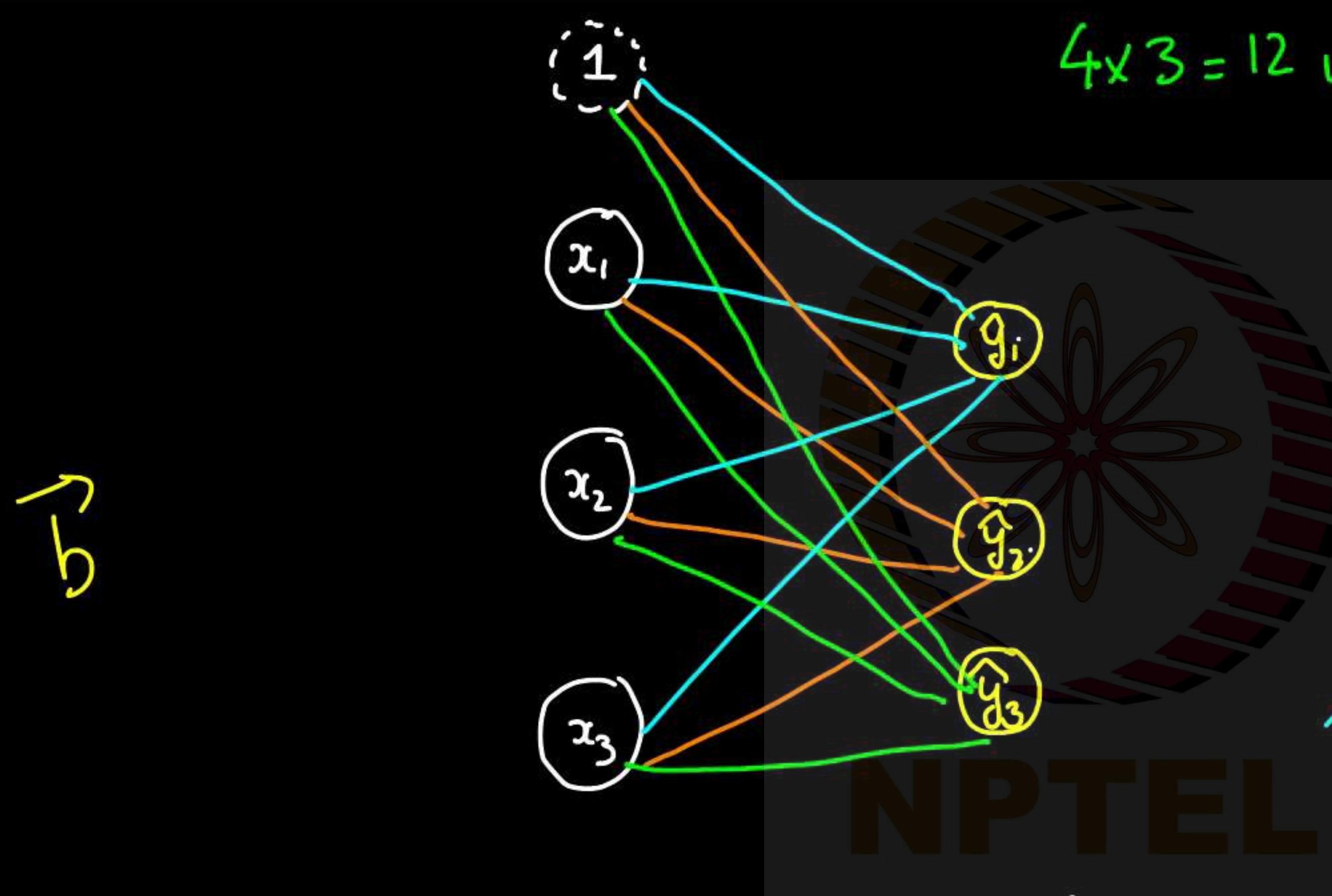
B I



NPTEL



File Edit View Insert Actions Tools Help



$$\hat{y}_1 = \text{softmax} (\omega_{01} + \omega_{11} x_1 + \omega_{21} x_2 + \omega_{31} x_3)$$

$$\hat{y}_2 = \text{softmax} (\omega_{02} + \omega_{12} x_1 + \omega_{22} x_2 + \omega_{32} x_3)$$

$$\hat{y}_3 = \text{softmax} (\omega_{03} + \omega_{13} x_1 + \omega_{23} x_2 + \omega_{33} x_3)$$

input

$w_i^{(l)}$

$j \leftarrow \text{output}$

$\hat{y}$

$3 \times 1$

1

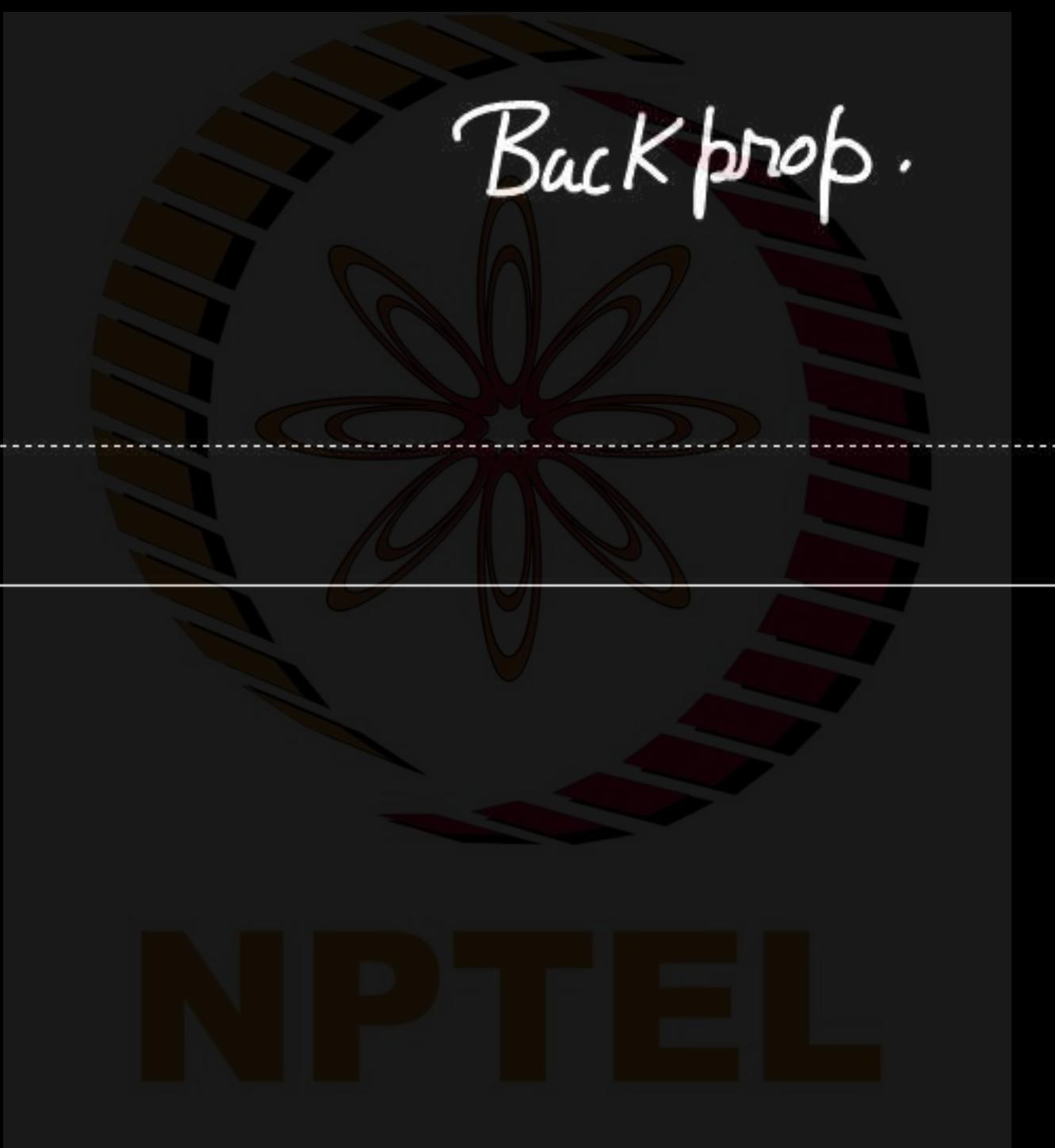


File Edit View Insert Actions Tools Help



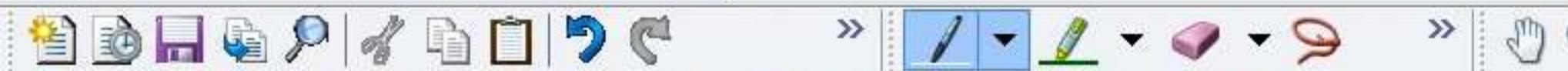
NPTEL

Note Title

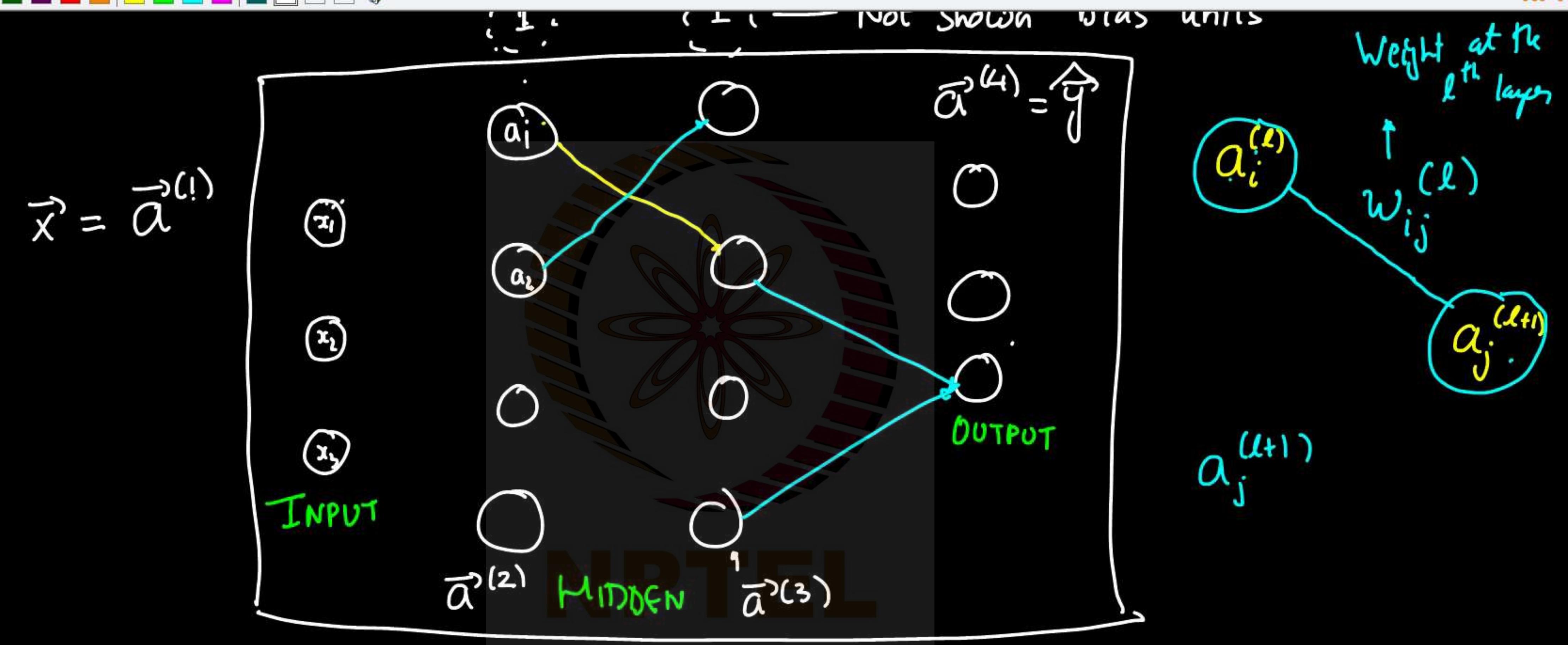




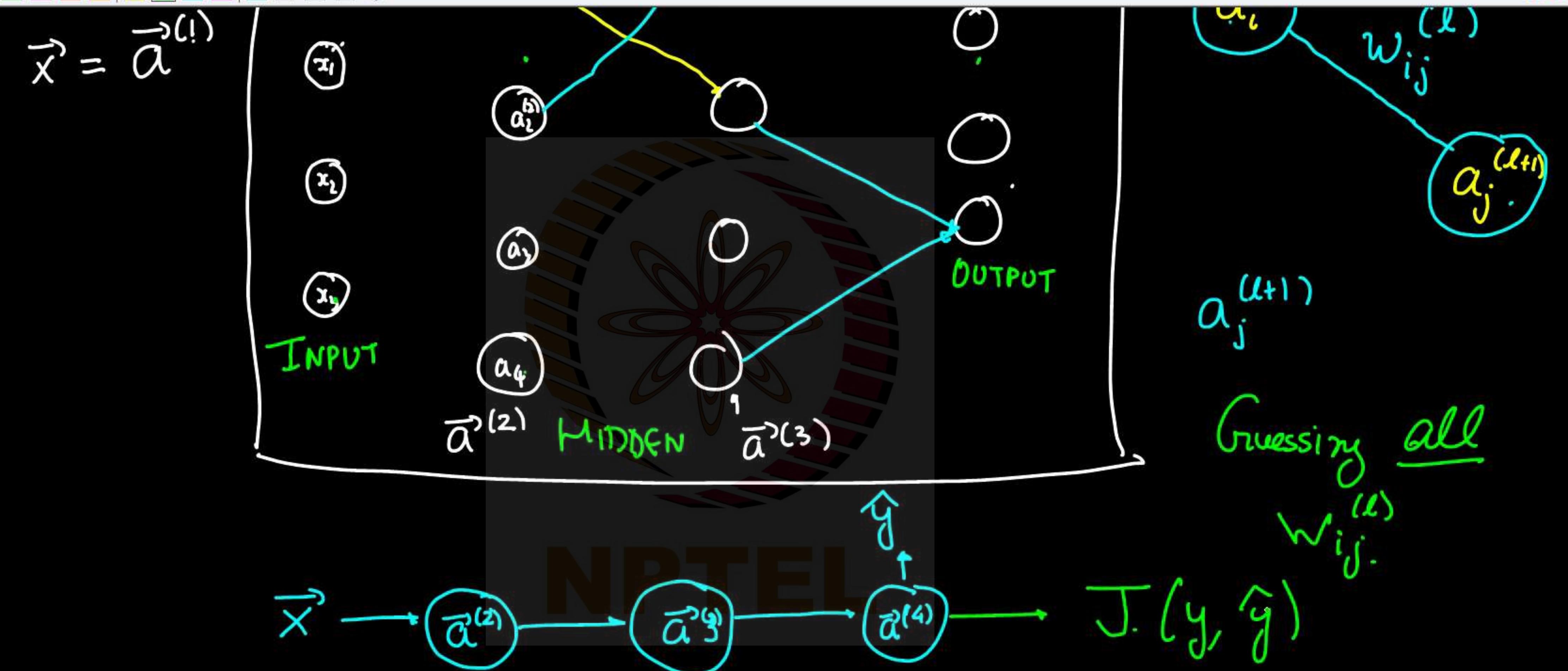
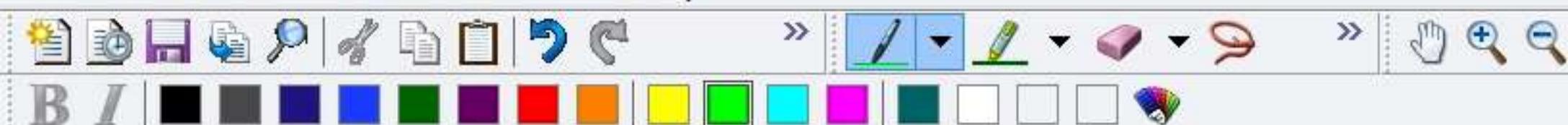
File Edit View Insert Actions Tools Help



NPTEL



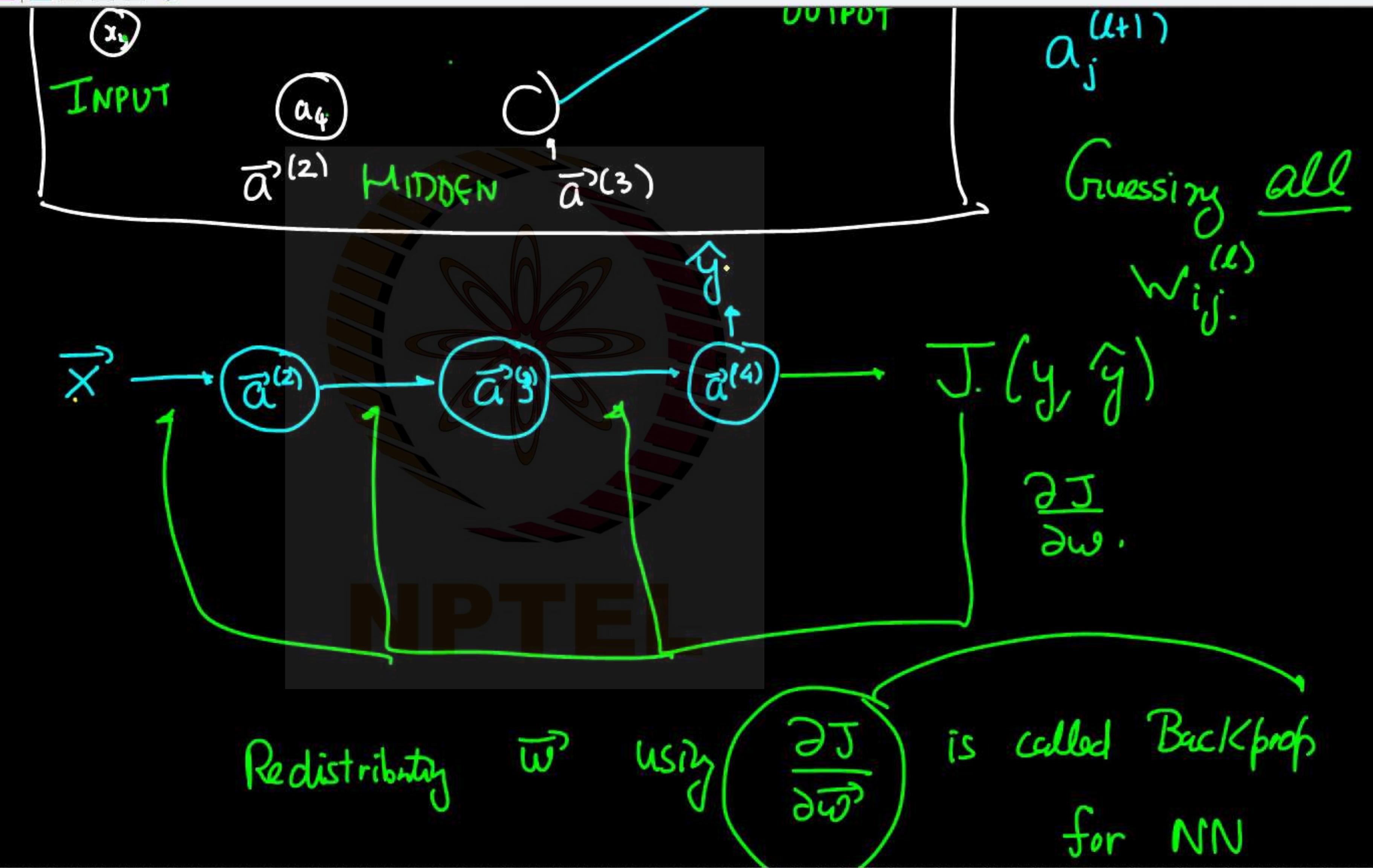
File Edit View Insert Actions Tools Help



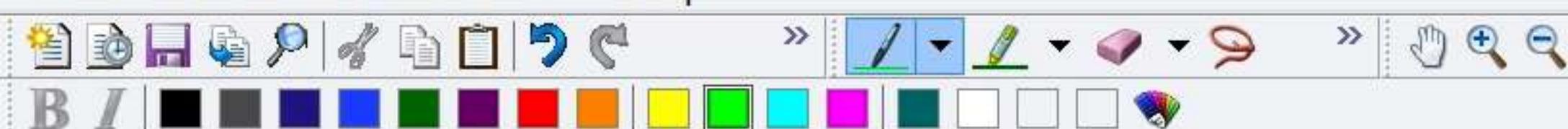
File Edit View Insert Actions Tools Help



B I



File Edit View Insert Actions Tools Help



Redistributing  $\vec{w}$  using  $\left( \frac{\partial J}{\partial \vec{w}} \right)$  is called Backprop for NN

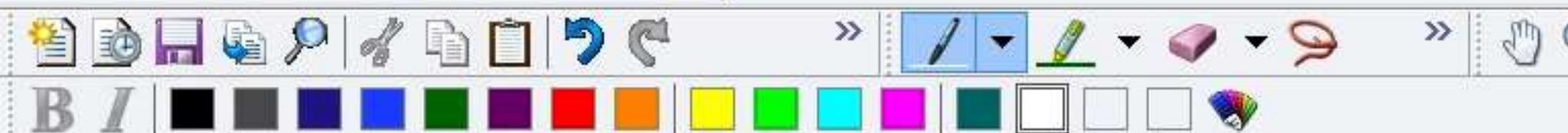
Simple Method

Coding

- ① Guess  $\vec{w}$   $\xrightarrow[\text{Perturbation}]{\text{Forward}}$   $\hat{y}$ . Calculate  $J(\vec{w})$
- ② Guess  $\vec{w} + \Delta \vec{w}$   $\xrightarrow{\text{Forward}}$   $\hat{y} + \delta \hat{y}$ . "  $J(\vec{w} + \Delta \vec{w})$

$$\frac{\partial J}{\partial w_1} \approx \frac{J(\vec{w} + \Delta \vec{w}) - J(\vec{w})}{\Delta w_1} : \text{Finite Difference Method}$$

Experiments



Simple Method

Coding

- ① Guess  $\vec{w}$   $\xrightarrow{\text{Forward}}$   $\hat{y}$ . Calculate  $J(\vec{w})$
- ② Guess  $\vec{w} + \Delta\vec{w}$   $\xrightarrow{\text{Forward}}$   $\hat{y} + \delta\hat{y}$ . "  $J(\vec{w} + \Delta\vec{w})$

$$\frac{\partial J}{\partial w_1} \approx \frac{J(\vec{w} + \Delta\vec{w}) - J(\vec{w})}{\Delta w_1} : \text{Finite Difference Method}$$

Very Expensive

Backprop (Minton)

Automatic Differentiation

Chain Rule



$$\frac{\partial J}{\partial w_1} \approx \frac{J(\vec{w} + \Delta \vec{w}) - J(\vec{w})}{\Delta w_1} : \text{Finite Difference Method}$$



Automatic Differentiation

Chain Rule.

Tensorflow, MATLAB have back prop routines



Very expensive

Back prop

( Minton )

Automatic Differentiation

Chain Rule.

Tensorflow, MATLAB have back prop routines

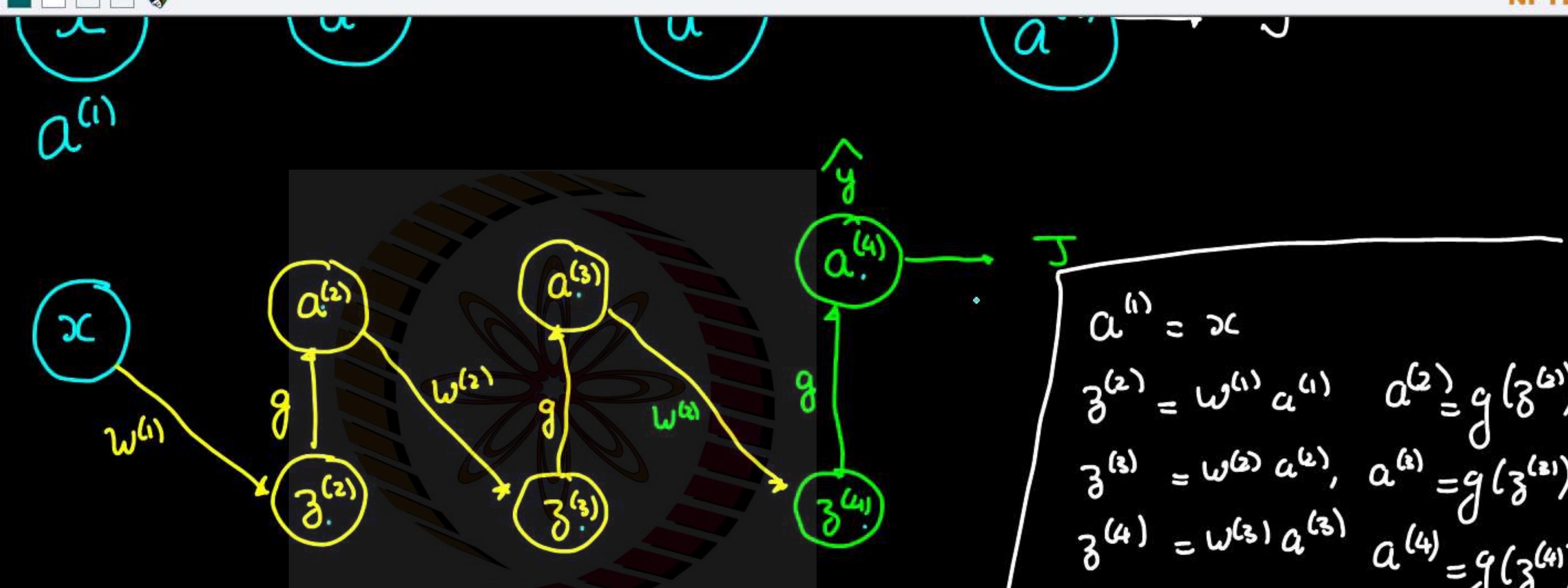
Toy derivation ( Simplified derivation )

NPTEL

File Edit View Insert Actions Tools Help



Ignoring the bias term



$$\frac{\partial J}{\partial w^{(1)}} = ?$$

$$\frac{\partial J}{\partial w^{(2)}} = ? \quad \frac{\partial J}{\partial w^{(3)}} = ?$$

$$\hat{y} = a^{(4)}$$

$$\begin{aligned}
 a^{(1)} &= x \\
 z^{(2)} &= w^{(1)} a^{(1)} \quad a^{(2)} = g(z^{(2)}) \\
 z^{(3)} &= w^{(2)} a^{(2)}, \quad a^{(3)} = g(z^{(3)}) \\
 z^{(4)} &= w^{(3)} a^{(3)} \quad a^{(4)} = g(z^{(4)})
 \end{aligned}$$

File Edit View Insert Actions Tools Help

NPTEL

$\delta^{(2)} = y - \hat{y}$

$\delta^{(3)} = w^{(3)} \delta^{(2)}$

$\delta^{(4)} = w^{(4)} \delta^{(3)}$

$\hat{y} = a^{(4)}$

$J = -\{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\}$

$g = \text{Sigmoid}$

$\frac{\partial J}{\partial w^{(1)}} = ?$

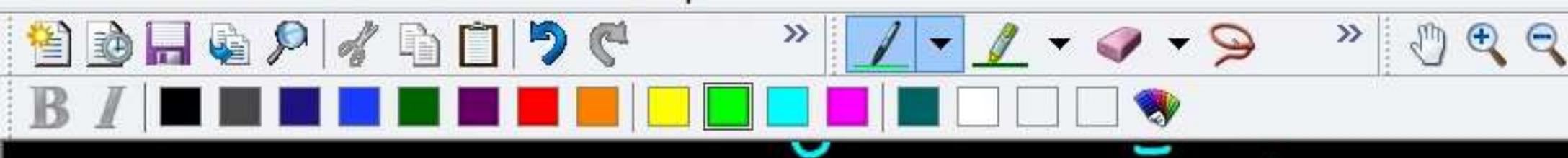
$\frac{\partial J}{\partial w^{(2)}} = ?$

$\frac{\partial J}{\partial w^{(3)}} = \frac{\partial J}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial \delta^{(4)}} \frac{\partial \delta^{(4)}}{\partial w^{(3)}}$

$\frac{\partial J}{\partial \delta^{(4)}} = -\{y - \hat{y}\}$



File Edit View Insert Actions Tools Help



$$\frac{\partial J}{\partial w^{(1)}}$$

$$\frac{\partial J}{\partial w^{(2)}}$$

$$\frac{\partial J}{\partial w^{(3)}}$$

$$J = - \{y \ln \hat{y} + (1-y) \ln (1-\hat{y})\}$$

$\hat{y}$  = Sigmoid

$$\frac{\partial J}{\partial w^{(3)}} = \frac{\partial J}{\partial a^{(4)}} \underbrace{\frac{\partial a^{(4)}}{\partial z^{(4)}}}_{\frac{\partial J}{\partial z^{(4)}}} \frac{\partial z^{(4)}}{\partial w^{(3)}}$$

$$\frac{\partial J}{\partial z^{(4)}} = -\{y - \hat{y}\}.$$

$$a^{(3)}$$

Denotes Error in activation

$$\frac{\partial J}{\partial w^{(3)}} = -\{y - a^{(4)}\} a^{(3)}$$

Error in output activation.

$$\frac{\partial J}{\partial z^{(4)}} = \delta^{(4)} \cdot \frac{\partial J}{\partial z^{(4)}}$$

$$\frac{\partial J}{\partial z^{(4)}} = \delta^{(4)} \cdot \frac{\partial J}{\partial z^{(4)}}$$



B I

 $\partial w^{(3)}$ 

$$\frac{\partial J}{\partial w^{(2)}} = \frac{\partial J}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial z^{(4)}} \frac{\partial z^{(4)}}{\partial a^{(3)}} \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$$\frac{\partial J}{\partial z^{(3)}} \equiv \delta^{(3)}$$

(\*) (\*)

$$\frac{\partial z^{(3)}}{\partial w^{(2)}}$$

 $a^{(2)}$ 

$$\boxed{\frac{\partial J}{\partial w^{(l)}} = \delta^{(l+1)} a^{(l)}}$$

$$z^{(3)} = w^{(2)} a^{(2)}$$

We Know  $\delta^{(1)}$

Don't "  $\delta^{(3)}, \delta^{(2)}$

We Know  $\delta^{(1)}$

We need to find  $\delta^{(3)}$

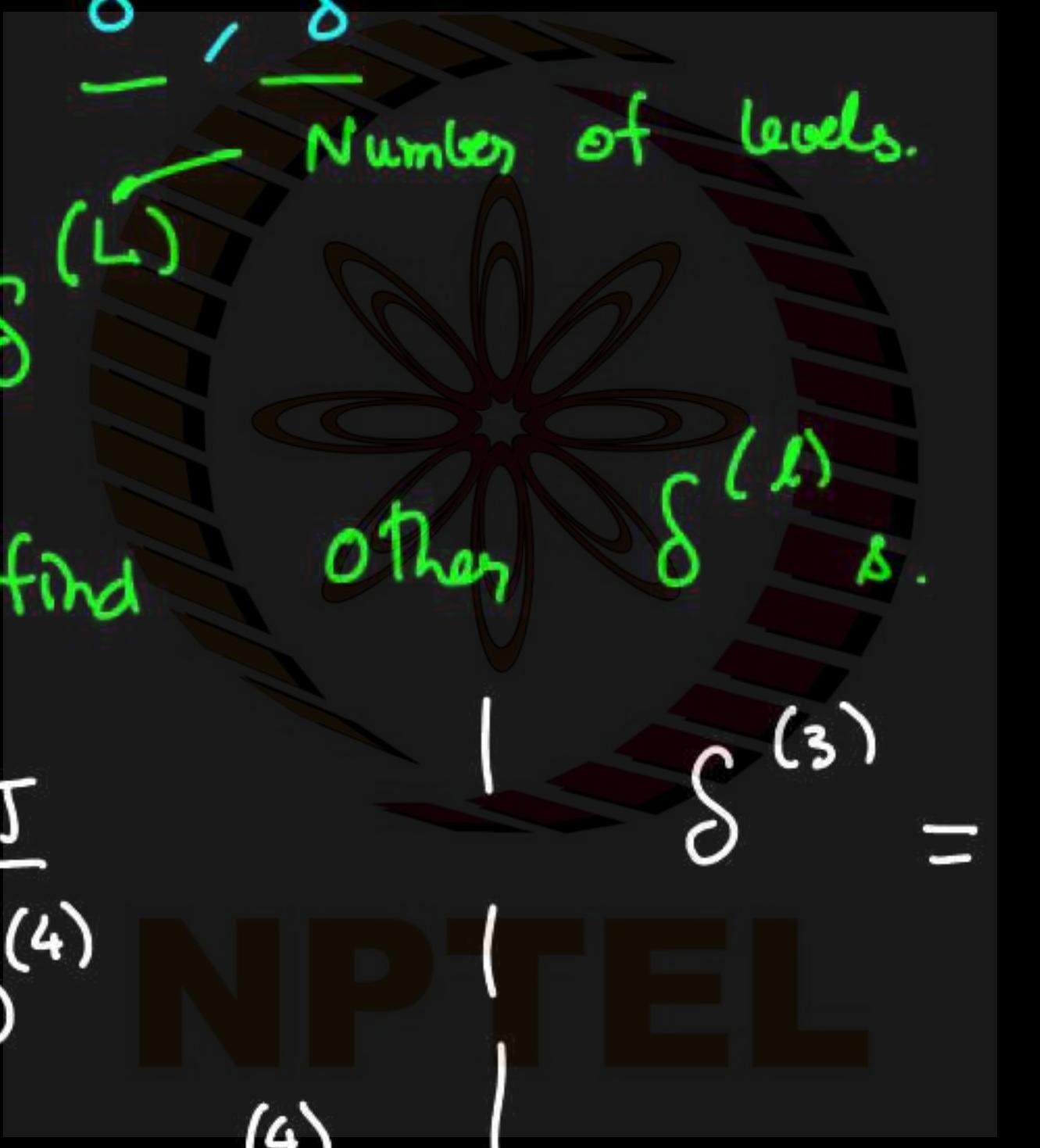
$\delta^{(4)} = \frac{\partial J}{\partial z^{(4)}}$

$= \frac{\partial J}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial z^{(4)}}$

$\delta^{(3)} = \frac{\partial J}{\partial z^{(3)}}$

$= \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}}$

$\overline{\partial w^{(1)}} - - -$





We need to find  $\sigma_{ij}$  when  $a_i$ .

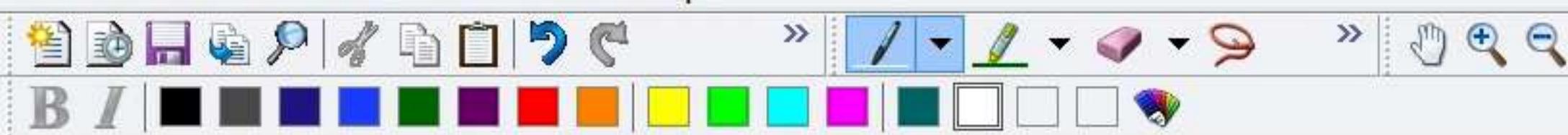
$$\sigma^{(4)} = \frac{\partial J}{\partial z^{(4)}}$$

$$= \frac{\partial J}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial z^{(4)}}$$

$$\sigma^{(3)} = \frac{\partial J}{\partial z^{(3)}}$$

$$= \underbrace{\frac{\partial J}{\partial a^{(4)}} \frac{\partial a^{(4)}}{\partial z^{(4)}}}_{\sigma^{(4)}} \left[ \frac{\partial z^{(4)}}{\partial a^{(3)}} \right] \frac{\partial a^{(3)}}{\partial z^{(3)}}$$

$w^{(3)}$   
↑



$$\delta^{(3)} = \delta^{(4)} w^{(3)} g'(z^{(3)})$$

Exact to  
Machine precision  
↑

$$\textcircled{a} \stackrel{\textcircled{g}}{\rightarrow} \delta^{(4)} \\ - (y - \hat{y})$$

$\parallel$   $\delta^{(l)} = \delta^{(l+1)} w^{(l)} g'(z^{(l)})$

Backprop

$$\left[ \frac{\partial J}{\partial w^{(l)}} = \delta^{(l+1)} a^{(l)} \right]$$

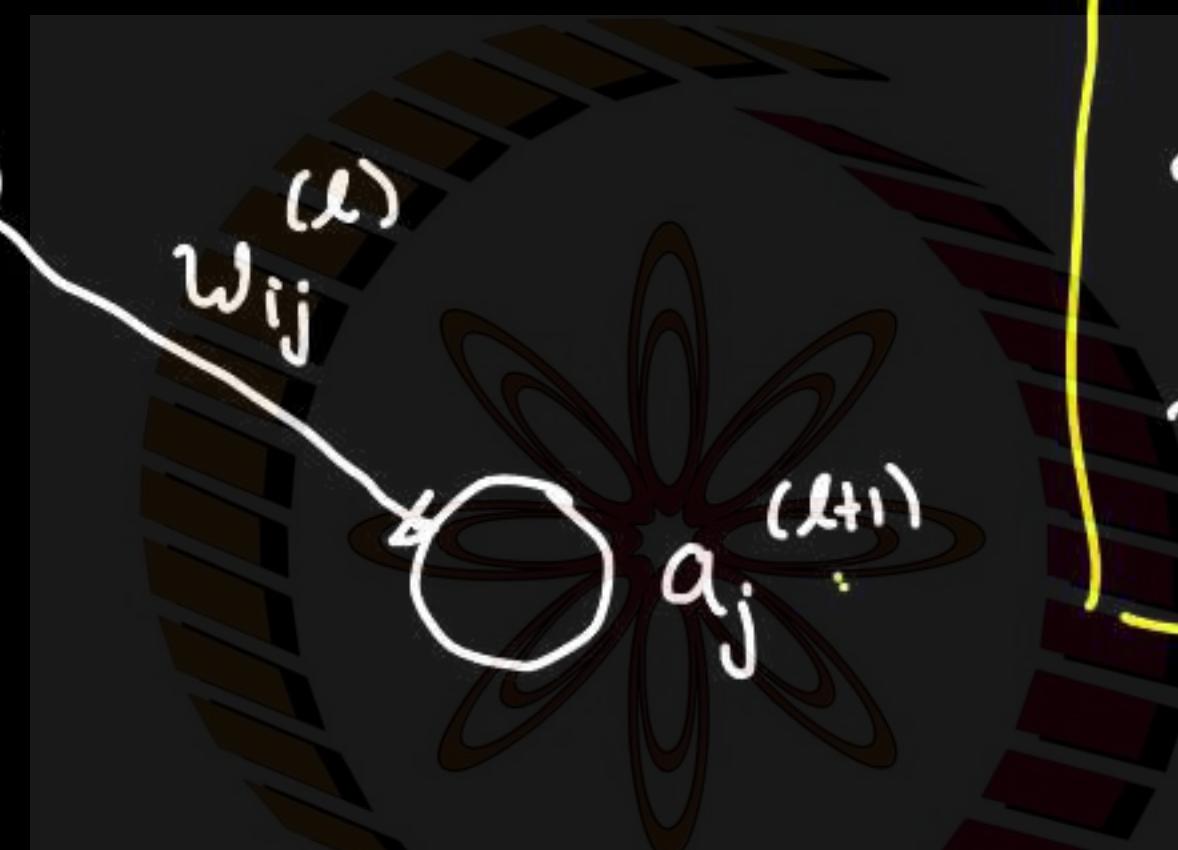
Chain Rule

All gradients can be computed.



General, vectorized expressions.

$$a_i^{(l)}$$

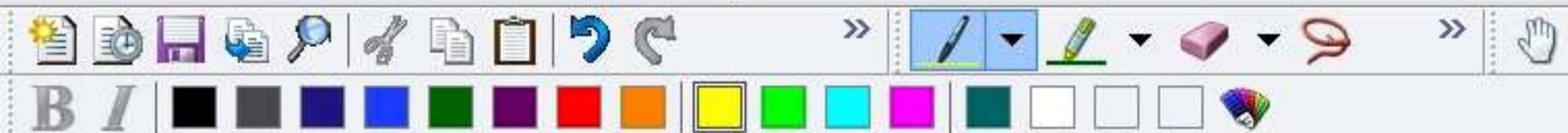


$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \delta_j^{(l+1)} a_i^{(l)}$$

$$\frac{\partial J}{\partial w_{ij}^{(l)}} = \delta_j^{(l+1)} a_i^{(l)}$$

$$\delta^{(l+1)} = \underbrace{\left[ W^{(l)} \delta^{(l)} \right]}_{\substack{\text{Matrix} \\ \text{Vector}}} \odot \underbrace{\left[ g'(\vec{z}^{(l)}) \right]}_{\text{Vector}}$$

Hadamard Product (Element wise Multiplication)



Graph

File Edit View Insert Actions Tools Help

B I

NPTEL

# Week 5

## Note Title

Linearly Separable Case  
Logistic Regression → Two classes  
Multinomial " " K classes

$\hat{y} = \nabla (\vec{w} \cdot \vec{x}) \leftarrow j = \text{Binary class}$

$\sum_{i=1}^m -\{y_i - \hat{y}_i\}x_i$

Softmax  $(\vec{w} \cdot \vec{x})$

1 / 1

File Edit View Insert Actions Tools Help



Multinomial " " ↗ K classes



$$\sum_{i=1}^m -\{y_i \hat{y}_i\} x_i$$

$$\hat{y} = \nabla (\vec{w} \cdot \vec{x}) \leftarrow J = \text{Binary cross}$$

$$\text{Softmax } (\vec{w} \cdot \vec{x})$$

XOR → Extra Layers.

$$(x) \rightarrow (f)$$

Universal Approx Thm.

Neural Networks → More Than One hidden layer

'Deep' Learning > 1 " "

 $i=1$ SOFTmax ( $w \cdot z$ )

Training Deep Network

Satwate / No Training.

XOR → Extra Layers.

$\times \rightarrow \oplus$

Universal Approx Thm.

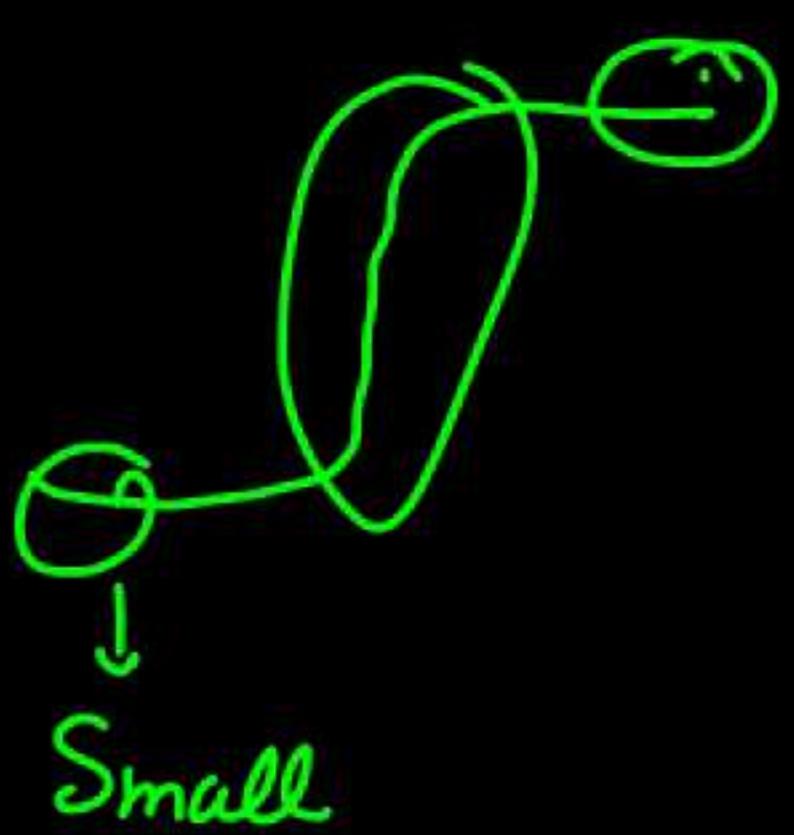
Neural Networks → More Than One hidden layer  
 'Deep' Learning  $> 1$  "

$$\delta^{(l+1)} = \delta^{(l)} w^{(l)}$$

$$\delta^{(1)} \approx 0.1 \delta^{(2)}$$

$$\delta^{(2)} \approx 0.1 \delta^{(3)}$$

$$\delta^{(l+1)} = \delta^{(l)} w^{(l)} g'(z^{(l)})$$



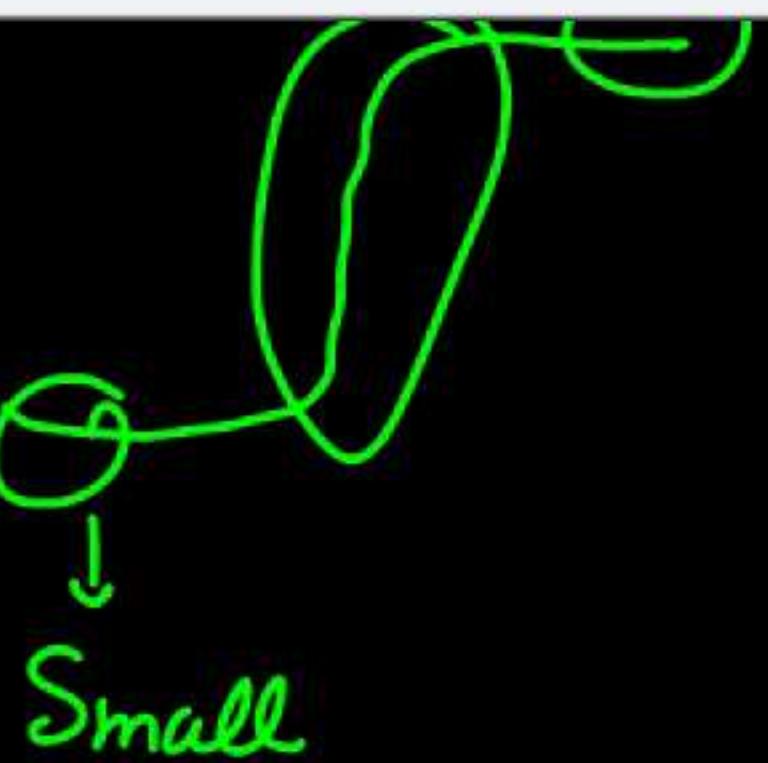
Satwate /

Back prob

$$\delta^{(l+1)} \approx 0.1 \delta^{(4)}$$

$$\delta^{(l)} \approx 0.1 \delta^{(3)}$$

$$\delta^{(l+1)} = \delta^{(l)} w^{(l)} g'(z^{(l)})$$



Did not cover

- ① How do we initialize  $\vec{w}$ ?
- ② How do we determine # of layers # of neurons
- ③ Nonlinearity?

Sigmoid

tanh

Hyperparameters

Linear

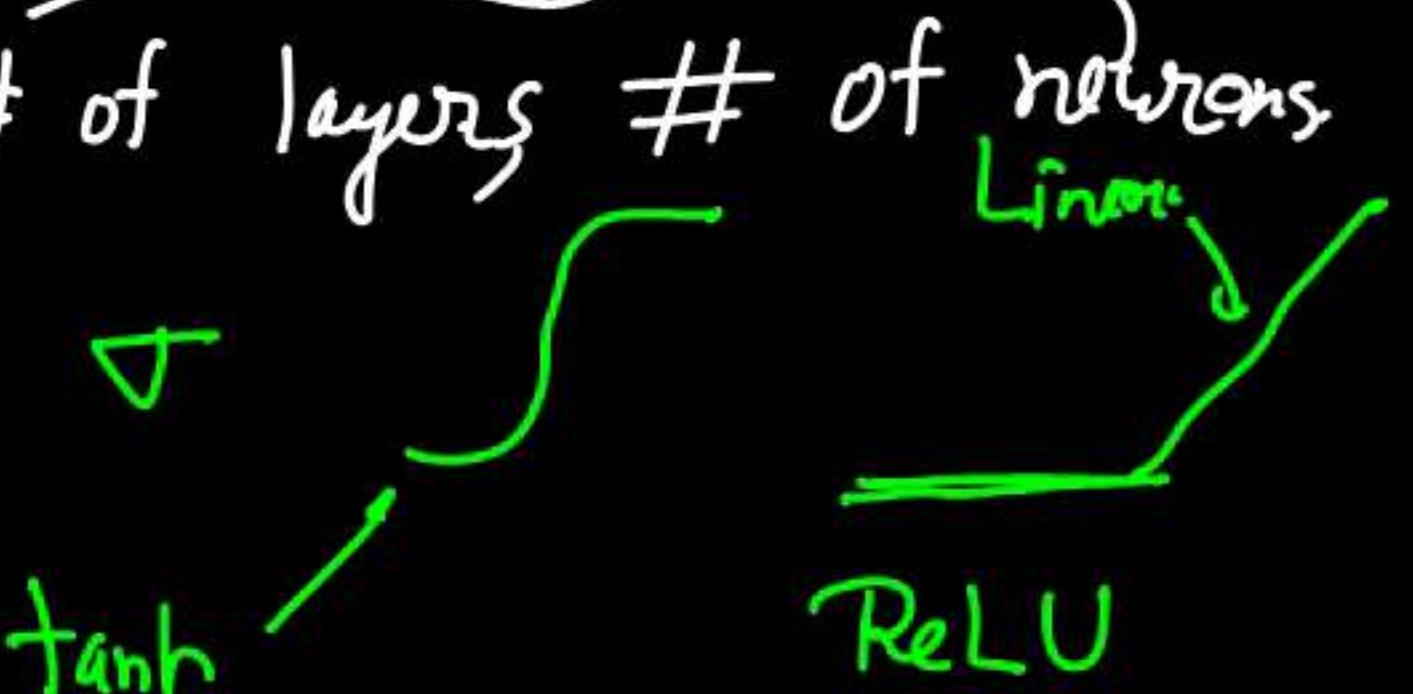
ReLU



① How do we initialize  $\vec{w}$ ?

- ② How do we determine  
Nonlinearity?
- ③ Sigmoid

Hyperparameters



Convolutional Neural Networks (CNNs)

NPTEL

ANNs for vision problems

$\begin{matrix} 0 & 0 \\ 0 & 0 \\ \vdots & \vdots \end{matrix}$

$3600 \quad 3600 \quad (3600)^2$  weights