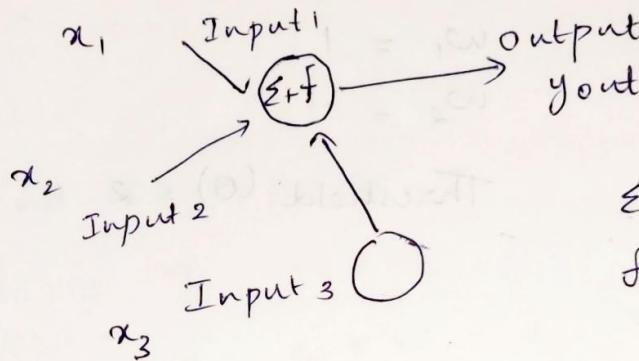


7/9

Deep Learning

## Basics of Neural network

- MP Neuron model - First foundation
- Single Layer perceptron neural network
- MP Neuron - synapsis, dendrites, soma  
n number of inputs



$\Sigma$  - Summation  
 $f$  - Activation  
 - ReLU  
 - sigmoid  
 - Softmax

## MIP neuron model

AND gate

- accepts only binary input  
and binary output

OR gate

## And gate

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

$$y = x_1 \cdot w_1 + x_2 \cdot w_2$$

$$f(x) = x_1 \cdot w_1 + x_2 \cdot w_2$$

$$f(x) = \sum_{i=1}^n x_i \cdot w_i$$

Summation

Activation function ( $f$ )

- Binary
- Sigmoid
- Softmax
- ReLU
- unit step
- tanh

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

0 - Threshold values

Construct the m/p neuron model for AND gate

$x_1$	$x_2$	$y$	$w_1 = 1$
Input		Output	$w_2 = 1$
0	0	0	
0	1	0	Threshold (0) = 2
1	0	0	
1	1	1	

First observation / First Input instance  
 $x_1 = 0 \quad x_2 = 0 \quad w_1 = 1 \quad w_2 = 1 \Rightarrow$  Parameters

Summation

$$\begin{aligned} f(x) &= \sum_{i=1}^n x_i w_i \\ &= x_1 w_1 + x_2 w_2 \\ &= 0(1) + 0(1) = 0 \end{aligned}$$

Activation function

$$\begin{aligned} g(x) &= \begin{cases} 1, & f(x) \geq 0 \\ 0, & f(x) < 0 \end{cases} \\ &= \begin{cases} 1, & \geq 2 \\ 0, & < 2 \end{cases} \end{aligned}$$

$$g(f(x)) = y_{\text{out}} = 0$$

$x_1$	$x_2$	$y$ actual	$y_{out}$ predicted
0	0	0	0
0	1	0	0
1	0	0	0

Second Observation

$$y = x_i * w_i$$

$$\begin{aligned} y &= x_1 * w_1 + x_2 * w_2 \\ &= 0(1) + 1(1) \end{aligned}$$

$$f(x) = 1$$

$$\text{Activation } y_{out} = 0$$

Third observation

$$x_1 = 1 \quad x_2 = 0 \quad w_1 = w_2 = 1 \quad \theta = 0$$

$$\begin{aligned} \sum_{i=1}^n x_i w_i \\ &= x_1 w_1 + x_2 w_2 \\ &= (1)(1) + (0)(0) \end{aligned}$$

$$\therefore f(x) = 1$$

$$\text{Activation } y_{out} = 0$$

Fourth observation

$$x_1 = 1 ; \quad x_2 = 1 ; \quad w_1 = w_2 = 1$$

$$f(x) = x_1 w_1 + x_2 w_2 \\ = (1 \times 1) + (1 \times 1)$$

$$f(x) = 2$$

$$y_{out} = 1$$

Construct a MP model for OR gate  $w_1, w_2 = 2$   
threshold = 1

$x_1$	$x_2$	$y$	$w$	No Learning, No error
0	0	0		
0	1	1		
1	0	1		
1	1	1		

I:  $f(x) = x_1 w_1 + x_2 w_2 \\ = 0(2) + 0(2) = 0$

Activation fn:  $\begin{cases} 1, & \geq 1 \\ 0, & < 1 \end{cases}$

$$y_{out} = 0$$

II  $f(x) = x_1 w_1 + x_2 w_2 \\ = 0(2) + 1(2) = 2$

Activation fn:  $\begin{cases} 1, & \geq 1 \\ 0, & < 1 \end{cases}$

$$y_{out} = 1$$

$$\text{III } f(x) = 1(2) + 0(2) = 2$$

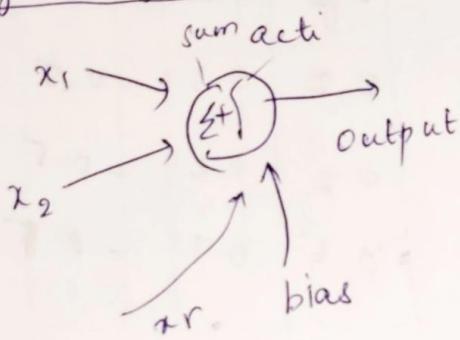
$x_1$	$x_2$	Actual	Predicted
0	0	0	0
0	1	1	1
1	0	1	1
1	1	1	1

$$y_{out} = 1$$

$$\text{IV } f(x) = 1(2) + 1(2) = 4$$

$$y_{out} = 1$$

### Single Layer perceptron (Delta rule)



- Only one computational Layer

error = Actual - Predicted

Learning rate ( $\alpha$ )  $\rightarrow$  0 to 1

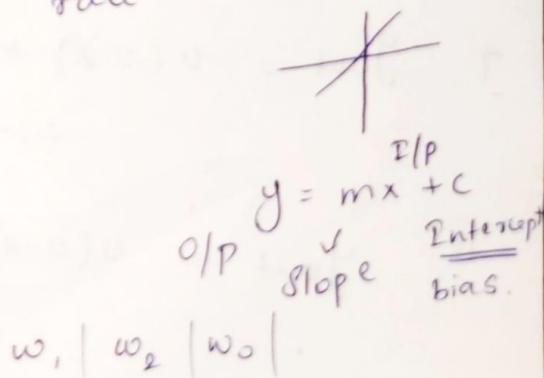
0.4, 0.5, 0.6 - stable range learning rates.

$$f(x) = y_{out} = x_i \cdot w_i$$

becomes:  $x_i \times \frac{\text{learning rate}}{\text{rate}} \times \frac{\text{error}}{\text{rate}} \times w_i$

AND gate

$x_1$	$x_2$	$x_0$	$y$	$y_{out}$	error
actual					-predicted



Construct single layer perceptron for AND gate

$$\alpha = 0.4 \quad w_1 = 0.2 \quad w_2 = 0.3 \quad w_0 = -0.5 \quad x_0 = 1 \quad \text{bias input}$$

(Learning rate)

Activation function (unit step)

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

AND gate

$$w_{\text{new}} = w_{\text{old}} + \frac{\text{learning rate} \times \text{error} \times x_i}{\text{rate}}$$

Weight updation

$x_1$	$x_2$	$x_0$	Actual y	Predicted $y_{\text{out}}$	Error $a - P$	$w_1$	$w_2$	$w_0$
0	0	1	0	0	0	0.2	0.3	-0.5
0	1	1	0	0	0	0.2	0.3	-0.5
1	0	1	0	0	0	0.2	0.3	-0.5
1	1	1	1	0	1	0.6	0.7	-0.5

$$f(x) = x_i w_i + (\text{bias}) (\text{input weight})$$

I.  $y_{\text{out}} = 0(0.2) + 0(0.3) + 1(-0.5)$   
= -0.5  $\Rightarrow$  0

II.  $y_{\text{out}} = 0(0.2) + 1(0.3) + 1(-0.5)$   
= 0.8 - 0.5 = -0.2  $\Rightarrow$  0

III.  $y_{\text{out}} = 1(0.2) + 0(0.3) + 1(-0.5)$   
= -0.3  $\Rightarrow$  0

IV.  $y_{\text{out}} = 1(0.2) + 0(0.3) + 1(-0.5)$   
= 0  $\Rightarrow$  0

$$w_1 = 0.2 + 0.4(1)(1) = 0.6$$

$$w_2 = 0.3 + 0.4(1)(1) = 0.7$$

$$w_0 = -0.5 + 0.4(1)(1)$$

$$= -0.1$$

OR gate	$x_1$	$x_2$	$x_0$	Actual y	Predicted y <sub>out</sub>	Error a-p	$w_1$	$w_2$	$w_0$
	0	0	1	0	0	0	0.2	0.3	-0.5
	0	1	1	1	0	1	0.2	0.7	-0.1
	1	0	1	1	1	0	0.2	0.7	-0.1
	1	1	1	1	1	0	0.2	0.7	-0.1

I  $y_{out} = 0(0.2) + 0(0.3) + 1(-0.5)$   
 $= -0.5 \Rightarrow 0$

II  $y_{out} = 0(0.2) + 1(0.3) + 1(-0.5)$   
 $= -0.2 \Rightarrow 0$

$$w_1 = 0.2 + 0.4(1)(0)$$

$$= 0.2$$

$$w_2 = 0.3 + 0.4(1)(1)$$

$$= 0.7$$

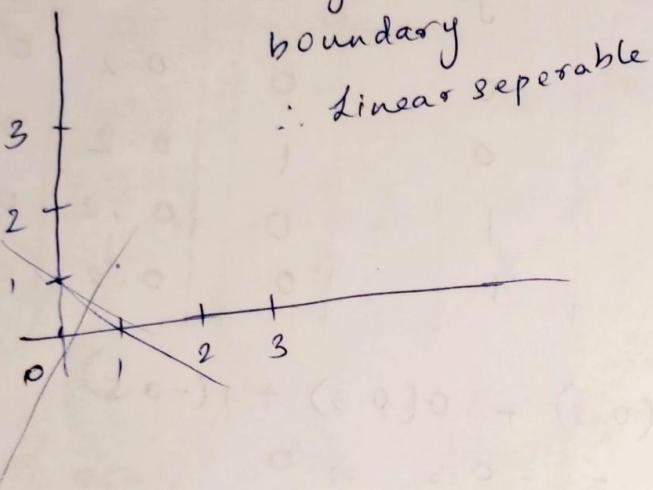
$$w_0 = -0.5 + 0.4(1)(1)$$

$$= -0.1$$

III  $y_{out} = 1(0.2) + 0(0.3) + 1(-0.1)$   
 $= 0.3 \Rightarrow 1$

IV  $y_{out} = 1(0.2) + 1(0.7) + 1(-0.1) = 0.8 \Rightarrow 1$

Linear separable		AND	OR	XOR	$\rightarrow$ Not
0	0	0	0	0	linear separable
0	1	0	1	1	
1	0	0	1	1	
1	1	1	1	0	



10.09.22

- Neural Network
- MC neural network
  - Hebbian Learning
  - Single layer Perceptron

### Hebbian Learning

AND gate :

$x_1$	$x_2$	O/P
0	0	0
0	1	0
1	0	0
1	1	1

SLP

$$w_{\text{new}} = \frac{w_{\text{old}}}{w_{\text{old}}} + \Delta w \quad \begin{matrix} \text{x input} \\ \text{learning} \\ \text{rate} \end{matrix}$$

$$w_{\text{new}} = w_{\text{old}} + \Delta w_i$$

$$\Delta w_i = x_i \times y_i$$

Actual  
O/P

We use bipolar input and bipolar o/p instead  
of binary i/p and o/p

[Else  $x_i y_i = 0$   
if  $x_i = 0$ ]

$x_1$	$x_2$	O/P
-1	-1	-1
-1	1	-1
1	-1	-1
1	1	1

Activation function  $\begin{cases} 1, & \geq 2 \\ 0, & < 2 \end{cases}$

$$w_1 = w_2 = 0$$

$x_1$	$x_2$	Bias ( $x_b$ )	Actual O/P	$\Delta w_1$	$\Delta w_2$	$\Delta w_b$	$w_1, w_2, w_b$
-1	-1	1	-1	1	1	-1	1 1 -1
-1	1	1	-1	1	-1	-1	2 0 -2
1	-1	1	-1	-1	1	-1	+1 1 3
1	1	1	1	1	1	1	2 2 -2

2, 1, -2

Q.

Construct a Hebbian Learning Rule for OR gate.

OR gate:

$x_1$	$x_2$	O/P
-1	-1	0
-1	1	1
1	-1	1
1	1	1

Activation function  $\begin{cases} 1, & \geq 2 \\ 0, & < 2 \end{cases}$

$$w_1 = w_2 = 0$$

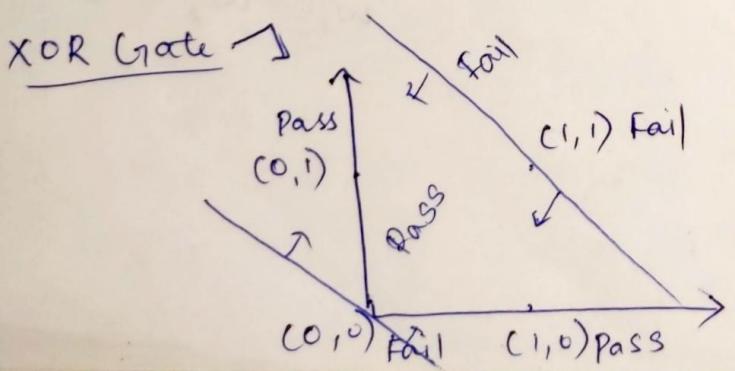
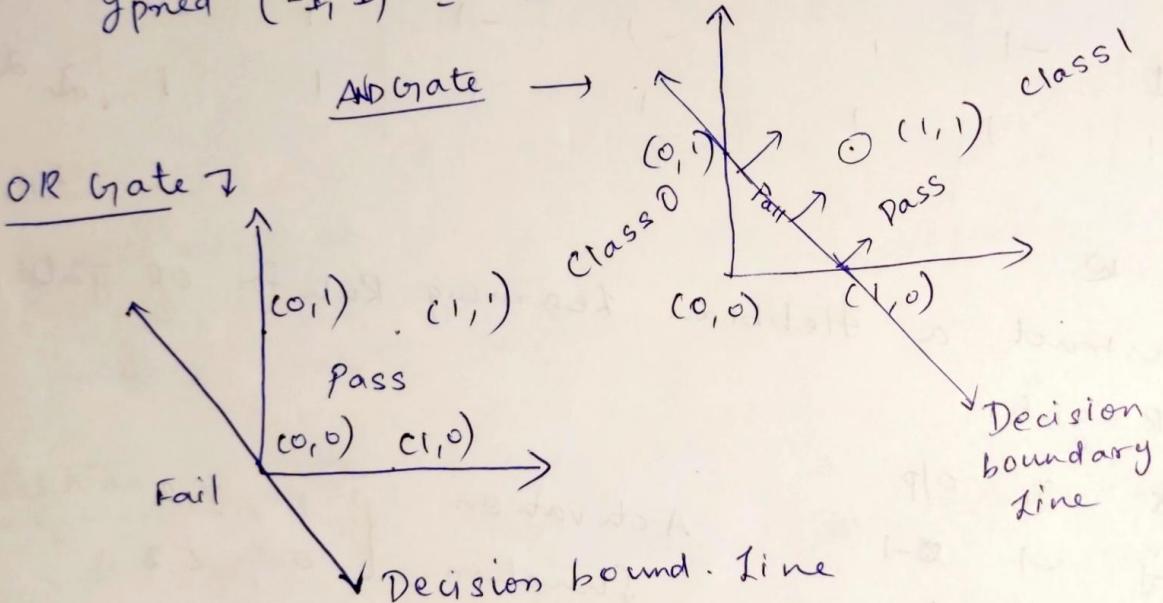
$x_1$	$x_2$	$x_b$	Actual O/P	$\Delta w_1$	$\Delta w_2$	$\Delta w_b$	$w_1$	$w_2$	$w_b$
-1	-1	1	-1	1	1	-1	1	1	-1
-1	1	1	1	-1	1	1	0	2	0
1	-1	1	1	1	-1	1	1	1	1
1	1	1	1	1	1	1	2	2	2

$$y_{pred}(-1, -1) = -2 - 2 + 2 = -2 \Rightarrow -1$$

$$y_{pred}(-1, 1) = -2 + 2 + 2 = 2 \Rightarrow 1$$

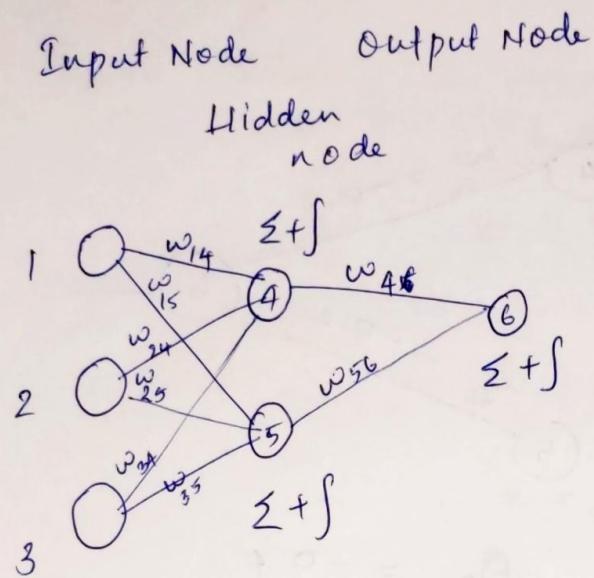
$$y_{pred}(1, -1) = 2 - 2 + 2 = 2 \Rightarrow 1$$

$$y_{pred}(1, 1) = 2 + 2 + 2 = 6 \Rightarrow 1$$



- 2 Decision boundary lines
- Not able to construct a single decision boundary line to separate 2 classes
- Non linearly separable

## Multi layer perceptron



$$T_j = \sum x_i \cdot w_{ij}$$

Summation

$$I_4 = x_1 \cdot w_{14} + x_2 \cdot w_{24} + x_3 \cdot w_{34} + \theta_4$$

- computations on hidden & o/p
- 2 layer feed forward neural network
- Summation ( $I_4, I_5, I_6$ )  $\Sigma +$
- Activation ( $O_4, O_5, O_6$ )
- Weight values ( $w_{14}, \dots, w_{56}$ )
- Error value ( $E_4, E_5, E_6$ )
- Bias value ( $\theta_4, \theta_5, \theta_6$ )  $\Sigma +$

Back propagation

Activation:

$$O_4 = \frac{1}{1 + e^{-I_4}}$$

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

Error value

Target = 1

$$\text{O/p Node } E_{O_6} = O_6 (1 - O_6) (T_j - O_6)$$

$T_j$  - Target / Actual output.

Hidden Node

$$E_{O_5} = O_5 (1 - O_5) E_{O_6} \cdot w_{56}$$

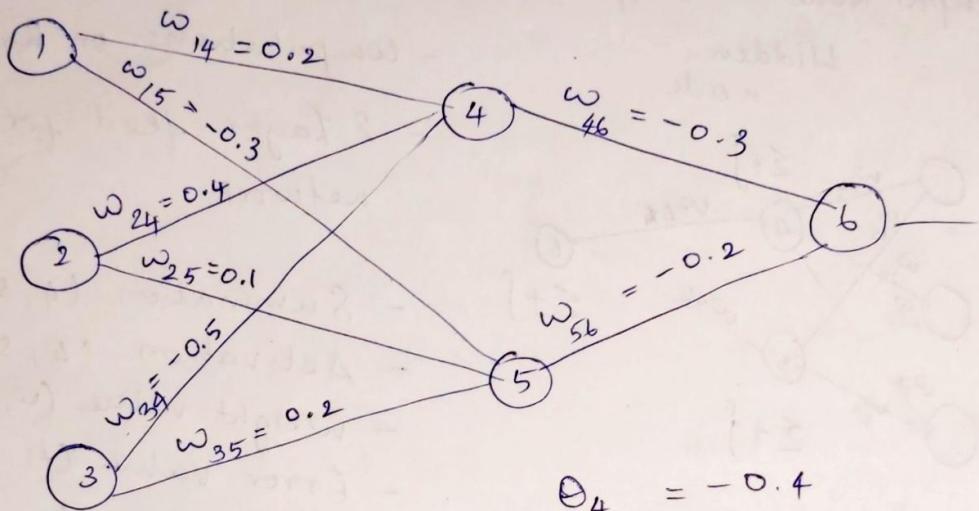
$$E_{O_4} = O_4 (1 - O_4) E_{O_6} \cdot w_{46}$$

Weight value

New weight

$$\begin{aligned} &= \text{Old weight} + \Delta w_{ij} \\ \Delta w_{ij} &= \text{learning rate} \times \text{error}_j \times x_i \end{aligned}$$

Construct a multi layer Perceptron



$$\begin{aligned} x_1 &= 1 \\ x_2 &= 0 \\ x_3 &= 1 \end{aligned}$$

$$\begin{aligned} \theta_4 &= -0.4 \\ \theta_5 &= 0.2 \\ \theta_6 &= 0.1 \end{aligned}$$

$$\begin{aligned} I_4 &= 1(0.2) + 0(0.4) + 1(-0.5) + -0.4 \\ &= 0.2 - 0.5 - 0.4 = -0.7 \end{aligned}$$

$$\begin{aligned} I_5 &= 1(-0.3) + 0(0.1) + 1(0.2) + 0.2 \\ &= 0.1 \end{aligned}$$

$$I_6 = \theta_4 = \frac{1}{1+e^{-I_4}} = \frac{1}{1+e^{+0.7}} = \frac{0.454}{0.332}$$

$$\text{as } \theta_5 = \frac{1}{1+e^{-I_5}} = \frac{1}{1+e^{-0.1}} = 0.525$$

$$\begin{aligned} I_6 &= 0.332(-0.3) + 0.525(-0.2) + 0.1 \\ &= -0.1046 \end{aligned}$$

$$Err_6 = O_6(1-O_6)(T-O_6)$$

$$O_6 = \frac{1}{1+e^{-I_6}} = \frac{1}{1+e^{0.1046}} = 0.473$$

$$\begin{aligned} Err_6 &= 0.473(1-0.473)(1-0.473) \\ &= 0.131 \end{aligned}$$

$$\begin{aligned} Err_5 &= O_5(1-O_5)err_6 \times w_{56} \\ &= 0.525(1-0.525)(0.131)(-0.2) \\ &= -0.00653 \end{aligned}$$

$$\begin{aligned} Err_4 &= 0.332(1-0.332)(0.131)(-0.3) \\ &= -0.0087 \end{aligned}$$

$$\text{New } w_{ij} = \text{Old } w_{ij} + \eta \times err_j \times x_i$$

$$w_{46} = -0.3 + 0.9 \times 0.131 \times 0.332 = -0.26$$

$$w_{56} = -0.2 + 0.9 \times 0.131 \times 0.525 = -0.138$$

$$w_{35} = 0.2 + 0.9 \times -0.0065 \times 1 = 0.194$$

$$w_{34} = -0.5 + 0.9 \times -0.0087 \times 1 = -0.507$$

$$w_{25} = 0.1 + 0.9 \times -0.0065 \times 0 = 0.1$$

$$w_{24} = 0.4 + 0.9 \times -0.0087 \times 0 = 0.4$$

$$w_{15} = -0.3 + 0.9 \times -0.0065 \times 1 = -0.305$$

$$w_{14} = 0.2 + 0.9 \times -0.0087 \times 1 = 0.19$$

$$O_5 = 0.2 + 0.9 \times -0.0065 = 0.194$$

$$O_6 = 0.1 + 0.9 \times 0.131 = 0.2179$$

$$O_4 = -0.4 + 0.9 \times -0.0087 = -0.407$$

$$I_4 = 1(0.193) + 0(0.4) + 1(-0.507) - 0.407 \\ = -0.721$$

$$O_4 = \frac{1}{1 + e^{-0.721}} = 0.327$$

$$I_5 = 1(-0.305) + 0(0.1) + 1(0.195) + 0.195 \\ = 0.085$$

$$O_5 = \frac{1}{1 + e^{-0.085}} = 0.521$$

$$I_6 = 0.327(-0.261) + 0.521(-0.138) + 0.218 \\ = 0.061$$

$$O_6 = \frac{1}{1 + e^{0.061}} = 0.515$$

$$\text{err}_6 = 0.515(1-0.515)(1-0.515) = 0.121$$

$$\text{err}_5 = 0.521(1-0.521)(0.121)(-0.138) = -0.004$$

$$\text{err}_4 = 0.327(1-0.327)(0.121)(-0.261) = -0.006$$

$$w_{46} = -0.261 + (0.9)(0.121)(0.327) = -0.225$$

$$w_{56} = -0.138 + (0.9)(0.121)(0.521) = -0.081$$

$$w_{14} = 0.193 + (0.9)(-0.006) \times 1 = 0.188$$

$$w_{15} = -0.305 + (0.9)(-0.004)(1) = -0.309$$

$$w_{24} = 0.4 + (0.9)(-0.006)(0) = 0.4$$

$$w_{25} = 0.1 + (0.9)(-0.004)(0) = 0.1$$

$$w_{34} = -0.507 + (0.9)(-0.006)(1) = -0.512$$

$$w_{35} = 0.195 + (0.9)(-0.004)(1) = 0.191$$

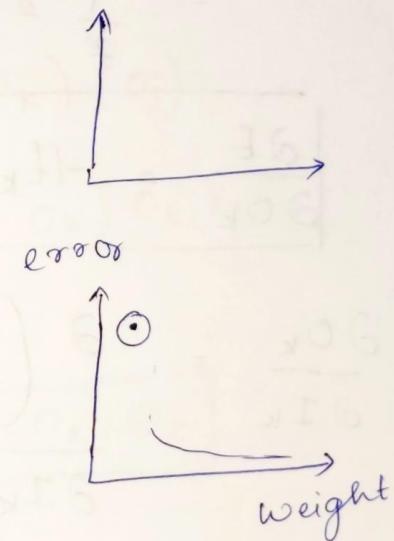
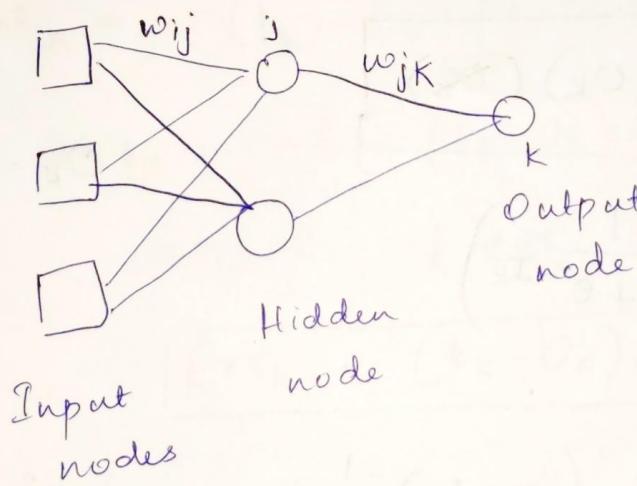
$$\theta_4 = -0.407 + (0.9)(-0.006) = -0.412$$

$$\theta_5 = 0.193 + (0.9)(-0.004) = 0.191$$

$$\theta_6 = 0.218 + (0.9)(0.121) = 0.327$$

12/09/22 Backpropagation with gradient Descent

Weight updation Technique



$$\Delta w \propto -\frac{\partial E}{\partial w}$$

$$\Delta w = -L \left( \frac{\partial E}{\partial w} \right)$$

$$w_1 = 0.2$$

$$w_2 = 0.8$$

$$w_3 = 0.5$$

$$\text{Error} = t - \text{actual}$$

objective : reduce  
the error

Change the weight

w.r.t E and weight

L - Learning rate

For error  $\epsilon_{jk} \Rightarrow$  check  $w_{jk}$

$$\Delta w_{jk} = \frac{\partial E}{\partial o_k} \cdot \frac{\partial o_k}{\partial I_k} \cdot \frac{\partial I_k}{\partial w_{jk}}$$

Chain rule (Summation, and activation)

$$\begin{aligned} \frac{\partial E}{\partial o_k} &= \frac{\partial \left( \frac{1}{2} (t_k - o_k)^2 \right)}{\partial o_k} & E &= \text{sum-squared error} \\ &= \frac{1}{2} 2(t_k - o_k) (-1) (\cancel{*}) & (t_k - o_k)^2 & \\ &\boxed{\frac{\partial E}{\partial o_k} = -(t_k - o_k) (\cancel{*})} & = \text{half sum squared error} & \\ &&& = \frac{1}{2} (t_k - o_k)^2 \end{aligned}$$

$$o_k = \frac{1}{1 + e^{-I_k}}$$

$$\begin{aligned} \frac{\partial o_k}{\partial I_k} &= \frac{\partial \left( \frac{1}{1 + e^{-I_k}} \right)}{\partial I_k} \\ &= \frac{\partial (1 + e^{-I_k})^{-1}}{\partial I_k} = (-1)(1 + e^{-I_k})^{-2} (e^{-I_k})(-1) \\ &= \frac{e^{-I_k}}{(1 + e^{-I_k})^2} \\ &= \frac{e^{-I_k}}{1 + e^{-I_k}} \cdot \frac{1}{1 + e^{-I_k}} = \frac{e^{-I_k}}{1 + e^{-I_k}} \cdot o_k \end{aligned}$$

$$= \left( 1 - \frac{1}{1 + e^{-I_k}} \right) \cdot O_k$$

$$\boxed{\frac{\partial O_k}{\partial I_k} = (1 - O_k) O_k}$$

$$\frac{\partial I_k}{\partial w_{jk}} = \frac{\partial (O_j \cdot w_{jk})}{\partial w_{jk}}$$

$$\boxed{\frac{\partial I_k}{\partial w_{jk}} = O_j}$$

$$\Delta w_{jk} = -(t_k - O_k) (1 - O_k) O_k O_j$$

$$\Delta w_{ijk} = -l (- (t_k - O_k) (1 - O_k) O_k) \cdot O_j$$

$$\Delta w_{ijk} = l e^{w_k} \cdot O_j$$

$$\boxed{Error_k = (t_k - O_k) (1 - O_k) O_k}$$

$$\Delta w \propto - \left( \frac{\partial E}{\partial w} \right)$$

$$\Delta w_{ij} = -l \left( \frac{\partial E}{\partial w_{ij}} \right)$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial O_j} \times \frac{\partial O_j}{\partial I_j} \times \frac{\partial I_j}{\partial w_{ij}}$$

$$\frac{\partial E}{\partial O_j} = \frac{\partial E}{\partial O_k} \cdot \frac{\partial O_k}{\partial I_k} \cdot \frac{\partial I_k}{\partial O_j}$$

$$\frac{\partial E}{\partial O_k} = \frac{\partial \left( \frac{1}{2} (t_k - O_k)^2 \right)}{\partial O_k} = (t_k - O_k)$$

$$\frac{\partial O_k}{\partial I_k} = (1 - O_k) O_k$$

$$\frac{\partial I_k}{\partial O_j} = \frac{\partial (\omega_{jk} \times O_j)}{\partial O_j}$$

$$\boxed{\frac{\partial I_k}{\partial O_j} = \omega_{jk}}$$

$$\begin{aligned} \frac{\partial O_j}{\partial I_j} &= \frac{\partial (\omega_{ij} \frac{1}{1 + e^{-I_j}})}{\partial I_j} \\ &= \frac{\partial (1 + e^{-I_j})^{-1}}{\partial I_j} \\ &= (-1) (1 + e^{-I_j})^{-2} (e^{-I_j})(-1) \\ &= \frac{e^{-I_j}}{(1 + e^{-I_j})^2} \end{aligned}$$

$$= \frac{e^{-I_j}}{1 + e^{-I_j}} \times \frac{1}{1 + e^{-I_j}}$$

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

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

$$= (1 - o_j)(o_j)$$

$$\frac{\partial I_j}{\partial w_{ij}} = \frac{\partial (o_i \times w_{ij})}{\partial w_{ij}}$$

$$\frac{\partial I_j}{\partial w_{ij}} = o_i$$

$$= -l \left( -((t_k - o_k)(1 - o_k)(o_k) \times (w_{jk}) \times (1 - o_j)(o_j)) \right) \frac{(o_i)}{(o_i)}$$

$$= l \left( \underline{(t_k - o_k)(1 - o_k)(o_k)} \times w_{jk} \times (1 - o_j)(o_j)(o_i) \right)$$

$$= l \text{err}_k w_{jk} \times (1 - o_j)(o_j)(o_i)$$

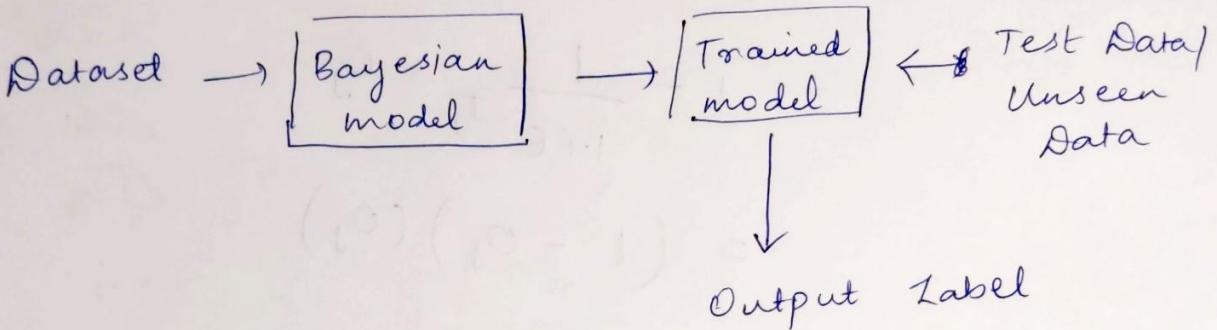
$$\text{Err}_j = (w_{jk})(1 - o_j)(o_j)$$

14/09/22

## Maximum likelihood Estimation

### Bayesian statistics

- Statistics model to predict the event
- Supervised learning technique / classification algorithm



### Bayesian statistics

#### ↳ Bayes Theorem

- Hypothesis - Bayesian
- Event - ML Statistics
- Class Label - ML

Attributes - Evidence

- Prior knowledge

- Likelihood

CGPA / Joboffer

### Naive Bayes theorem

- All attributes are independent to each other

$P(H/Evidence)$

Probability of hypothesis  
given Evidence

Class membership probability value

Posterior probability Conditional probability

$P(H/x) =$

Hypothesis  
 $x = \text{Evidence}$

$$\frac{P(x/H) \cdot P(H)}{P(x)}$$

→ Prior knowledge  
→ Evidence  
→ Likelihood

$$P(H/x_1, x_2, x_3 \dots x_n) = P(x_1/H) \times P(x_2/H) \times \dots P(x_n/H) \times P(H)$$

Posterior probability

$$= \text{Likelihood} \times \text{Prior probability}$$

$P(H) \rightarrow P(c)$   
(Class label)

CGPA	Interactivity	Practical knowledge	Communication skill	Job offer
1. $\geq 9$	Yes	Excellent	Good	Yes
2. $\geq 8$	No	Good	Moderate	Yes
3. $\geq 9$	No	Average	Poor	No
4. $\leq 8$	No	Average	Good	No
5. $\geq 8$	Yes	Good	Moderate	Yes
6. $\geq 9$	Yes	Good	Moderate	Yes
7. $\leq 8$	Yes	Good	Poor	No
8. $\leq 8, \geq 9$	Yes No	Excellent	Poor	Yes <del>No</del> *
9. $\geq 8$	Yes	Excellent	Good	Yes *
10. $\geq 8$	Yes	Good Avg	Good	Yes
$\geq 9$	Yes	Avg	Good	?

$$P(\text{Job offer} = \text{Yes}) = \frac{7}{10}$$

$$P(\text{Job offer} = \text{No}) = \frac{3}{10}$$

CGPA 4

Job offer = Yes

Job offer = No

 $\geq 9$ 

3

1

 $\geq 8$ 

4

-

 $< 8$ 

0

2

CGPA

Job offer = Yes

Job offer = No

 $\geq 9$ 

3/7

Y<sub>3</sub> $\geq 8$ 

4/7

2/3 -

 $< 8$ 

0

2/3

Interactivity

Job offer = Yes

Job offer = No

Yes

5

1

No

2

2

Likelihood

Yes

5/7

Y<sub>3</sub>

No

2/7

2/3

Practical

Yes

No

Excellent

2/7

0/3

Good

4/7

1/3

Avg

1/7

2/3

Comm.

Yes

No

Good

4/7

Y<sub>3</sub>

Moderate

3/7

0/3

Poor

0/7

2/3

$$P(\text{Yes} \mid \text{testdata}) = \frac{3}{7} \times \frac{5}{7} \times \frac{2}{7} \times \frac{4}{7} \times \frac{7}{10}$$

$$= 0.0174$$

$$P(\text{No} \mid \text{testdata}) = 0.0074$$

Transaction ID	Age	Income	Student	Credit rating	Buys comp
1	Youth	High	No	Fair	No
2	Youth	H	N	E	No
3	Middle	H	W/o	F	Y
4	Snr	M	No	F	Y
5	Snr	L	Yes	E	No
6	Snr	L	Yes	E	Y.
7	Middle	L	Yes	F	No
8	Youth	M	No	F	Yes
9	Youth	L	Yes	F	Yes
10	Snr	M	Yes	F	Yes
11	Youth	M	Yes	E	Yes
12	Middle	M	No	F	Yes
13	Middle	H	Yes	Exc	No
14	Snr	M	No		

Test data ( $x$ )  $x = \text{age} = \text{youth}, \text{income} = \text{medium},$   
 $\text{student} = \text{yes}, \text{credit rating} = \text{fair}$

$$P(\text{Buys computer} = \text{Yes}) = \frac{9}{14}$$

$$P(\text{Buys computer} = \text{No}) = \frac{5}{14}$$

<u>Age</u>	Yes	Yes No
		Yes No
Youth	2/9	3/5
Middle	4/9	0/5
Sn	3/9	2/5

<u>Income</u>	Yes	No
		No
High	2/9	2/5
Medium	4/9	2/5
Low	3/9	1/5

<u>One Student</u>	Yes	No
		No
Yes	6/9	1/5
No	3/9	4/5

<u>Credit rating</u>	Yes	No
		No
Excellent	3/9	3/5
Fair	6/9	2/5

Given a data tuple having the systems values  
 system 26-30 46K 50K dept, age, salary  
 resp.

	Dept	Status	age	salary	Count
1.	Sales	Senior	31..35	46-50K	30
2.	Sales	Junior	26-30	26-30K	40
3.	Sales	Jnr	31..35	31-35K	40
4.	Systems	Jnr	21-25	46-50K	20
5.	Systems	Snr	31-35	66-70K	5
6.	Systems	Jnr	26-30	46-50K	3
7.	Systems	Snr	41-45	66-70K	3
8.	Marketing	Snr	36-40	46-50K	10
9.	Marketing	Jnr	31-35	41-45K	4
10.	Secretary	Snr	46-50	36-40K	4
11.	Secretary	Jnr	26-30	26-30K	6

→ Capacity, underfitting, overfitting

↙  
 Training error

Test error / generalisation error

Generalisation

Bias = Training error

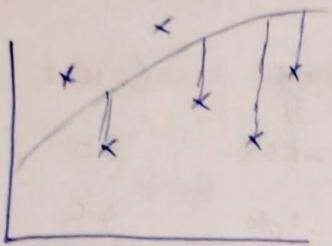
Variance = Testing error

Underfitting : - High bias  
 - High variance

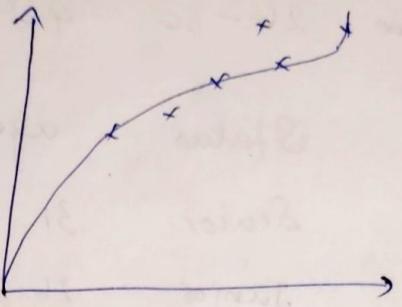
Overttting:

- Low bias
- High variance
- Generalised model
- Low bias
- Low variance

## Underfitting

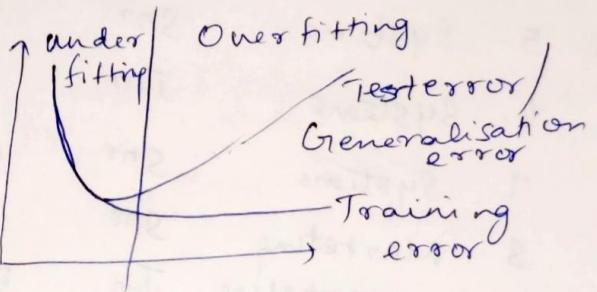
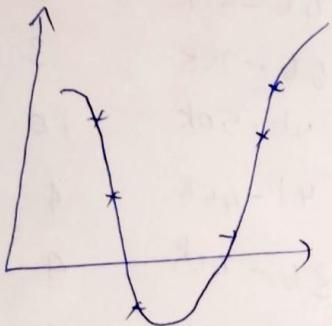


## Generalised Model / capacity



## Overfitting

(Polynomial degree)



$$\text{Solution : } P(\text{Snr}) = \frac{52}{165}$$

$$P(\text{Jnr}) = \frac{113}{165}$$

Department :

	Snr	Jnr
Sales	30/52	80/113
System	8/52	23/113
Marketing	10/52	4/113
Secretary	4/52	6/113

Age :

	Snr	Jnr
31 - 35	35/52	44/113
26 - 30	0	49/113
21 - 25	0	20/113
41 - 45	3/52	0
36 - 40	10/52	0
46 - 50	4/52	0

## Salary:

	Snr	Jnr
46-50	40/52	23/113
26-30	0	46/113
31-35	0	40/113
66-70	8/52	0
41-45	0	4/113
36-40	4/52	0

$$P\left(\frac{\text{Snr}}{\text{Systems}}, 26-30, 46\text{k} - 50\text{k}\right)$$

$$= \frac{8}{52} \times 0 \times \frac{40}{52} \times \frac{52}{165} = 0$$

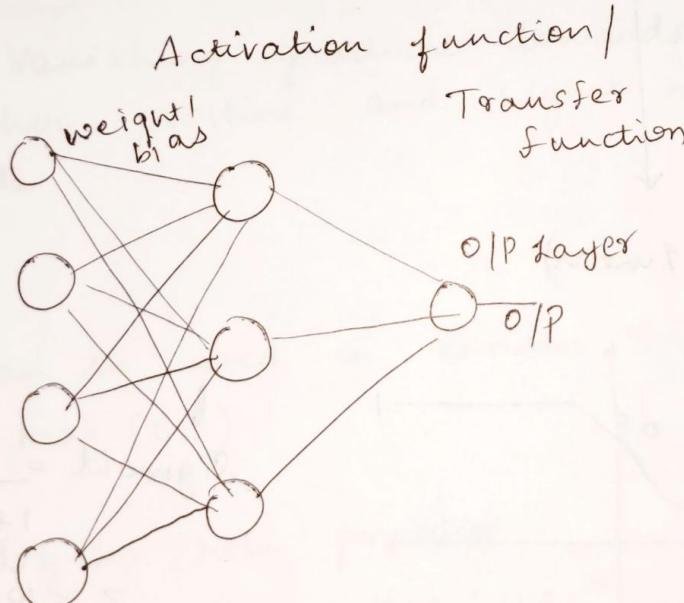
$$P\left(\frac{\text{Jnr}}{\text{Systems}}, 26-30, 46\text{k} - 50\text{k}\right)$$

$$= \frac{23}{113} \times \frac{49}{113} \times \frac{23}{113} \times \frac{113}{165}$$

$$= 0.012$$

∴ Prediction: Jnr.

24109  
Input layer

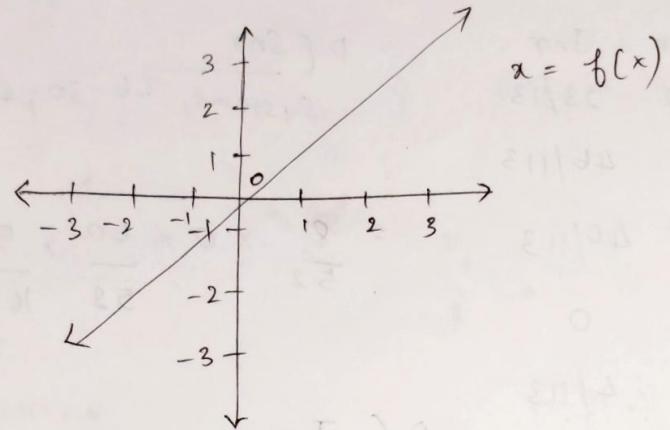


- (1) Sigmoid
- (2) Tanh
- (3) Relu
- (4) Threshold
- (5) Binary step

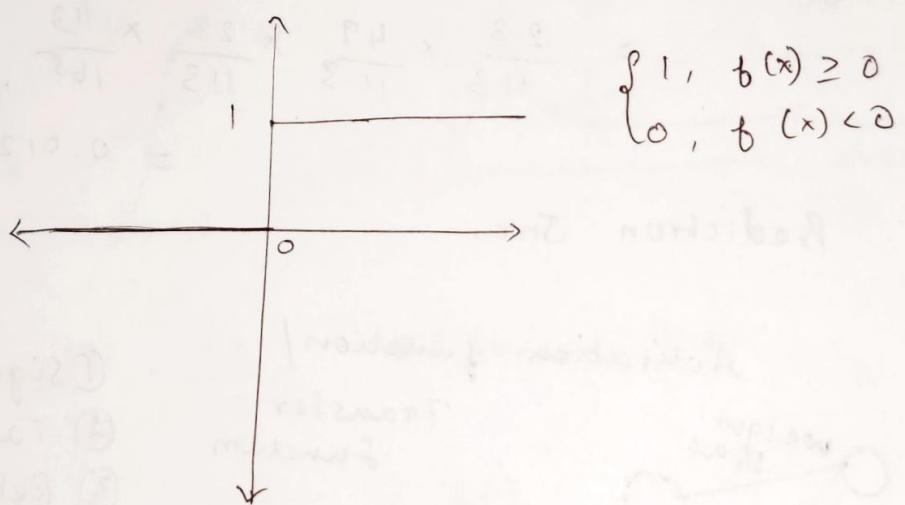
/ Squashing Function

- ↳ Linear
- ↳ Linear / Identity
- ↳ Binary step
- ↳ Non Linear

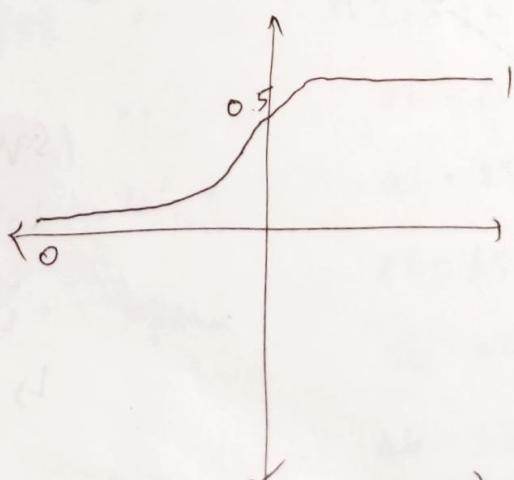
## ① Linear function



## ② Binary Step activation fn.



## ③ Sigmoid (Non-Linear)



$$\text{Sigmoid} = \frac{1}{1 + e^{-z}}$$

$$z = w \times \text{input} + \text{bias}$$

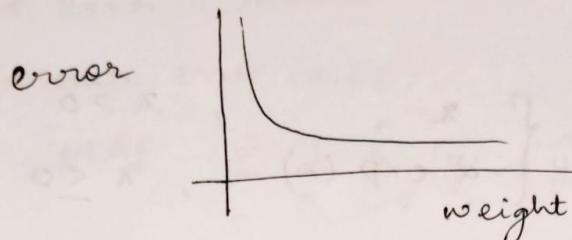
→ between (0 to 1)

→ can be used in o/p layer

→ In binary classification prblm

### Drawback:

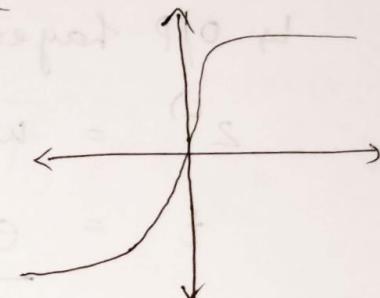
- Vanishing gradient / Saturation problem
- No weight updation for higher values, rate of change becomes 0.



- Mean is around 0.5  $\Rightarrow$  not zero-centric

### Tanh $\rightarrow$ Tan hyperbolic function

- Non linear activation function
- Range can be -1 to 1
- Mean is zero-centric



### Drawback:

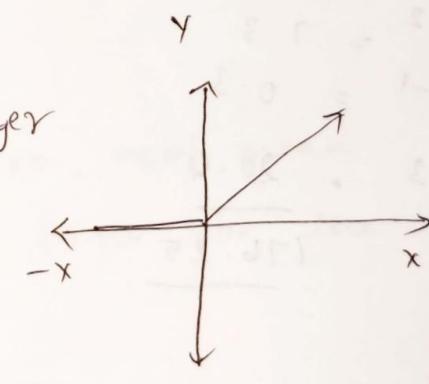
- Vanishing gradient towards higher positive and higher negative side

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

### Relu

- Can be used in hidden layer

$$\max(0, x)$$



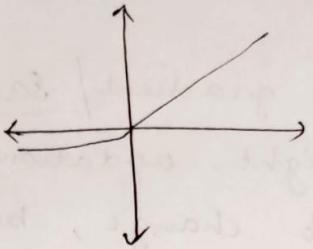
- Dying ReLU problem

(Towards negative side, weight is not updated)

- Overcome - (1) Leaky ReLU (3) Exponential ReLU  
(2) Parametric ReLU

## Leaky Relu

$$\max(0.01x, x)$$



## Parametric Relu

$$\max(\alpha x, x)$$

$\alpha$  - Hyper Parameter

$$\text{Exponential Relu} : \begin{cases} x & , x > 0 \\ \alpha \cdot \exp(x) - 1 & , x \leq 0 \end{cases}$$

## Multi class classification problem (4 classes)

↳ Softmax Activation function

↳ O/P layer

$$z^{(l)} = w \cdot x + b$$

$$\begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}$$

$$t = \frac{e^{z^{(l)}}}{e^5 + e^2 + e^{-1} + e^3} \xrightarrow{\text{normalize}}$$

$$e^5 = 148.4$$

$$\textcircled{i} \text{ First class : } \frac{148.4}{176.25}$$

$$e^2 = 7.3$$

$$\textcircled{ii} : \frac{7.3}{176.25}$$

$$e^{-1} = 0.3$$

$$\textcircled{iii} : \frac{0.3}{176.25}$$

$$e^3 = \frac{20.9}{176.25}$$

$$\textcircled{iv} : \frac{20}{176.25}$$

$$\begin{bmatrix} 0.84 \\ 0.064 \\ 0.002 \\ 0.113 \end{bmatrix}$$

$\therefore \textcircled{i}$  is class label

## Softmax Activation Function

j - number of classes

$$t = \sum_{j=1}^4 \frac{e^j}{e_j}$$

(1) Linear / Identity fn  
HW application

26/09/22 Error / Residuals

Standard error = Actual - Predicted (Observed)

Mean Absolute error (MAE)

$$MAE = \frac{1}{n} \sum_{i=0}^n |y_i - \hat{y}_i|$$

$y_i$  = Actual  
 $\hat{y}_i$  = Predicted

Mean Square Error (MSE)

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

Root Mean squared error (RMSE)

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

Actual Predicted

8	4
2	16
1	1

Find out the mean squared error, root mean squared error, mean absolute error for the given example

Items	Actual Sales		
1	80	5	120
2	90		
3	100		
4	110		

$$t = \frac{4}{\sum_{j=1}^4 \frac{e^j}{e^j}}$$

(1) Linear / Identity fn  
HW application

26/09/22 Error / Residuals

Standard error = Actual - Predicted (Observed)

Mean Absolute error (MAE)

$$MAE = \frac{1}{n} \sum_{i=0}^n |y_i - \hat{y}_i| \quad y_i = \text{Actual}$$

$\hat{y}_i = \text{Predicted}$

Mean Square Error (MSE)

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

Root Mean squared error (RMSE)

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

Actual Predicted

8	4
2	16
1	1

Find out the mean squared error, root mean squared error, mean absolute error for the given example

Items	Actual Sales	5	120
1	80		
2	90		
3	100		
4	110		

The value predicted by the model is 75

Find the MAE, RMSE, MSE

$$MAE = \frac{1}{5} [(80-75) + (90-75) + (100-75) + (110-75) + (120-75)]$$

$$= \frac{1}{5} (5 + 15 + 25 + 35 + 45)$$

$$= \frac{125}{5} = 25$$

$$MSE = \frac{1}{5} [5^2 + 15^2 + 25^2 + 35^2 + 45^2] = 825$$

$$RMSE = 28.7$$

Actual	Predicted
--------	-----------

$$P_1 \quad 80 \quad 75$$

$$P_2 \quad 75 \quad 85$$

$$MAE = \frac{1}{2} [(80-75) + (85-75)]$$

$$= \frac{1}{2} [5 + 10] = 7.5$$

$$MSE = \frac{1}{2} [25 + 100] = 62.5$$

$$RMSE = 7.90$$

→ Find the sales for week 6 using linear regression

Weeks (x) Sales

$$1 \quad 1000$$

$$2 \quad 3000$$

$$3 \quad 4000$$

$$4 \quad 8000$$

$$5 \quad ?$$

$$6 \quad ?$$

$$m = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$$c = \frac{\sum y - m \sum x}{n}$$

$$a_1 = \frac{(\bar{xy}) - \bar{x}\bar{y}}{(\bar{x^2}) - (\bar{x})^2}$$

$$a_0 = \bar{y} - a_1 \bar{x}$$

x	y	xy	$x^2$
1	1000	1000	1
2	3000	6000	4
3	4000	12K	9
4	8000	32K	16
—	—	—	—
10	16000	51000	—

$$\frac{(12750 - (2.5)(4000)}{7.5 - 6.25}$$

$$m = \frac{12750}{1.25} = 10,192$$

$$c = \frac{4000 - 10,192(2.5)}{-25,476}$$

$$y = 10192x - 25476$$

$$6 = 11700 \rightarrow 16100 = 35676$$

②	$x_1$	$x_2$	$y$
1	4	1	
2	5	6	
3	8	8	
4	2	12	

$$\beta = (x^T x)^{-1} x^T y$$

$$x^T x = \begin{bmatrix} 4 & 10 & 19 \\ 10 & 30 & 46 \\ 19 & 46 & 109 \end{bmatrix}$$

When we use Eigen vectors, value

Differentiation

Partial derivative

why we use activation fn.

Linear  $\rightarrow$  Non Linear  $\rightarrow$

Activation

SVD Decomposition

- P.

Vectors / Tensors

Basic <sup>temp</sup> operations

29/09

Loss function

$\rightarrow$  Classification

$\rightarrow$  Regression

$\hookrightarrow$  Mean square error

$\hookrightarrow$  Mean absolute error

$\hookrightarrow$  Huber loss

O/P : Binary class

Multi class

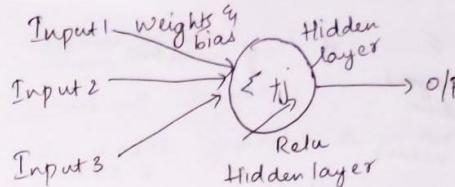
Binary cross

Entropy loss

Log loss

(Two class)

Categorical cross Entropy Loss



ReLU variants of  
relu

$\hookrightarrow$  Leaky

$\hookrightarrow$  PReLU

$\hookrightarrow$  Exponential ReLU

Error value /  
cost function /  
loss function

Binary cross entropy loss / log loss

$$\text{Binary cross entropy} = -\frac{1}{n} \sum_{i=1}^n y_i \log p + (1-y) \log (1-p)$$

$y$  represents  
actual / target  
value

$y$ : Actual / target  
value

$p$ : Predicted  
probability / estimated  
value

Sigmoid ( $f(x)$ )

$\hookrightarrow$  0 to 1

$0.5 \rightarrow 1$

below  $\rightarrow 0$

S.No	$x_1$	$x_2$	$x_3$	$\Sigma$ Weighted sum (Prob. value)	$y$	$g$	Loss
1.	0.1	0.5	0.2	0.26	1	0	0.585
2.	0.2	0.3	0.1	0.20	0	0	0.09
3.	0.7	0.4	0.2	0.48	1	0	0.31
4.	0.1	0.4	0.3	0.30	0	0	0.15
							$\frac{\Sigma}{4} 0.285$

Instances:

$$-y \log p + (1-y) \log (1-p)$$

$$-1 \log (0.26) + (1-1) \log (1-0.26) \\ = 0.585$$

Instance 2:

$$-0 \log p + (1-0) \log (1-0.20) \\ = -1 \log (0.80) = 0.096$$

Instance 3:

$$-1 \log (0.48) \\ = 0.318$$

Instance 4:

$$0.154$$

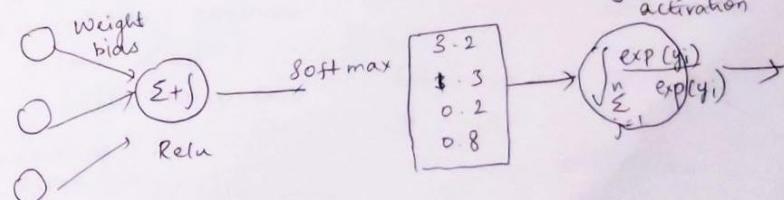
∴ Binary cross entropy  
loss = 0.285

### Optimization

- Gradient descent
- Stochastic gradient descent
- Adadelta
- Ada grad
- RMS prop

Reg No	Actual value	Estimated probability	LOSS
1	1	0.94	0.026
2	1	0.90	0.045
3	1	0.78	0.107
4	0	0.586	0.356
5	0	0.51	0.309
6	1	0.47	0.327
7	1	0.32	0.494
8	0	0.10	0.045
9			
			$\frac{\Sigma}{n} 0.213$

### Categorical cross Entropy



Estimated O/P      Actual O/P

0.775	1
0.116	0
0.039	0
0.070	0

### Categorical cross entropy

$$\text{Loss} = - \sum_{i=1}^n [T_i \log(P_i)]$$

↑                          ↑  
 Actual target      Estimated  
 Probability value

$$= - [1 \times \log(0.775) + 0 \times \log(0.116) + 0 \times \log(0.039) + 0 \times \log(0.070)] = 0.116$$

Iteration 2

Estimated O/P      Actual O/P

0.938	1
0.028	0
0.013	0
0.023	0

Loss  
= 0.027

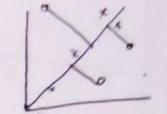
Mean square error  $\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$

Mean absolute error  $= \frac{1}{n} \sum_{i=1}^n |y - \hat{y}|$

Huber loss

$$\cdot = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & |y - \hat{y}| \leq \delta \\ |y - \hat{y}| - \frac{1}{2} \delta^2 & |y - \hat{y}| > \delta \end{cases}$$

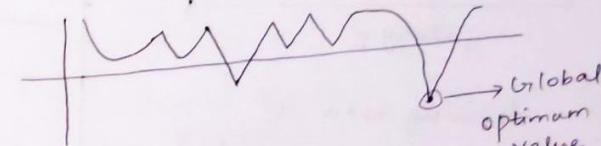
Drawback of mean<sup>2</sup>  
→ Very sensitive to outliers



(high residuals)<sup>2</sup>  
↳ high

Advantages

- Applying global optimum value



Huber loss

$\delta$  - Hyperparameter value  
To overcome Drawbacks of MSE

Regularisation

Solve  
Overfitting problem

- Low bias
- High variance

Prerequisite - Bias and Variance

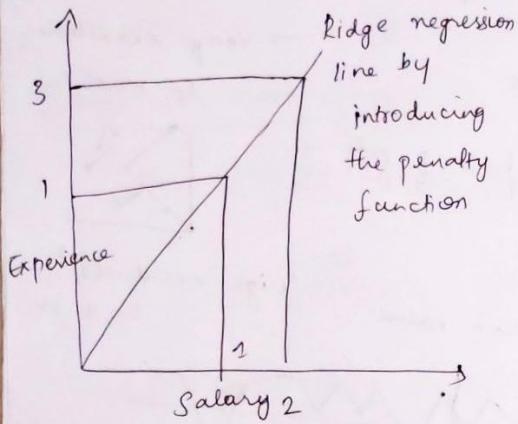
Underfitting

- High bias
- High variance

Generalised model

- Low bias
- Low variance

## Ridge regression and Lasso regression



Training error = 0

Testing error = high

Generalise by reducing the slope

$$\text{Ridge estimator} = \frac{1}{n} \sum \text{Sum of residuals} + \lambda \times \text{slope}^2$$

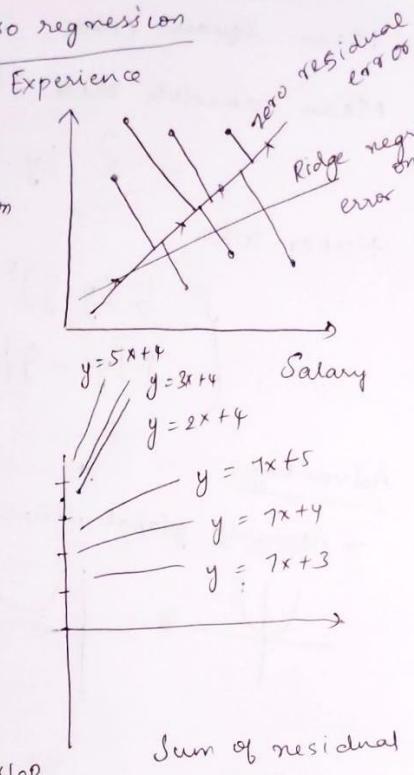
$$= \sum_{i=1}^n (y - \hat{y})^2 + \lambda \times \text{slope}^2$$

$$\text{Example } y = 1.3x + 0.4$$

$$m = 1.3$$

$$y \text{ intercept } c = 0.4$$

$$= 0 + 1 + (1.3)^2 = 2.69$$



Slope = 0.8

Lasso regression - Sum of residual +  $\lambda + (\text{slope})$

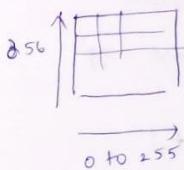
L2

## Convolution Neural Network (CNN)

Kernels / Filters / Feature detectors are used to determine the features.

Image

256 x 256



Filter / Feature detector / Kernel

## Building block of CNN

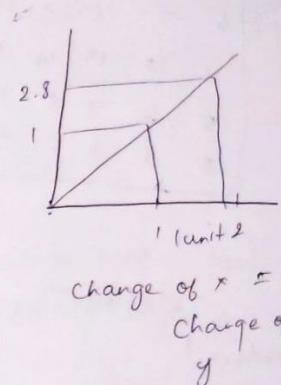
→ Convolution

- Kernel / Filter / Feature detector, Padding, Stride

→ Pooling

- Max pooling, Average pooling

→ Flatten ANN



$5 \times 5$  Image convolution

1	2	3	4	5
6	9	8	9	10
7	12	13	14	15
5	6	7	8	9
2	3	4	5	6

kernel/filter

1	2	3
4	5	6
7	8	9

x

411	456	501
372	417	462

=

stride = 1

Move filter by

Padding - Border preserving

n - size of image

f - filter size

p - padding

s - stride value

Feature maps /

Activation maps /  
convolved feature

Size of feature map

$$= \left\lceil \frac{n - f + 2p}{s} \right\rceil + 1$$

$$\text{If padding} = 0, n = 5, f = 3, s = 1$$

$$FM = \frac{5 - 3 + 2(0)}{1} + 1 = 3$$

One  $5 \times 5 \rightarrow$  overfitting

$4 \times 4 \rightarrow$  Heavy computation

Max pooling, Average pooling, min pooling - Reduce dimensions

↓  
Most common

After  
Pooling

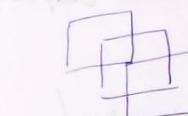
11	22
33	44

→

11
22
33
44

ANN - Classification

Image  
 $256 \times 256 \times 3$



x filter  
( $3 \times 3$ )

Feature map  
 $\frac{256 \times 256 \times 3}{3}$   
 $\frac{256}{2} \times 127 \times 3$

Convolution

$127 \times 127 \times 32$

x filter  
( $3 \times 3$ )  
 $\times 32$

Convolution  
 $62 \times 62 \times 64$

$80 \times 30 \times 64$

51,600

-1 2 -1  
+1 3 +1  
-1 4 -1

Horizontal  
edge

-1 2 -1  
+1 3 +1  
-1 2 -1

Vertical  
edge

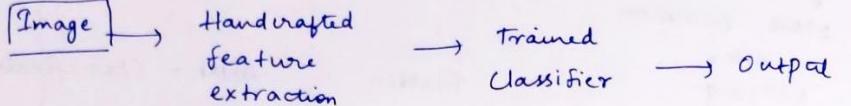
$256 \times 256 \times 3$

$\frac{256 \times 3}{1} = 254$

+ 1

$254 - 3$

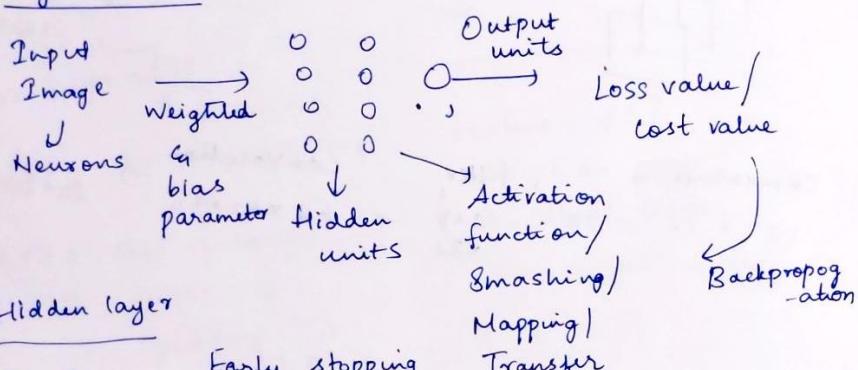
## Machine Learning



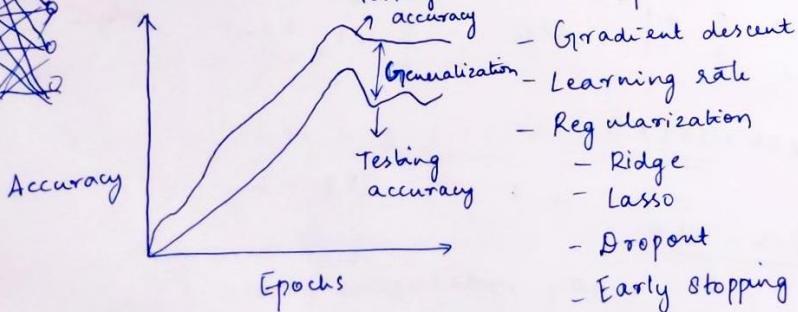
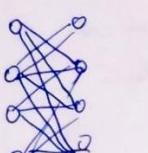
### CNN



### Key concepts



### Hidden layer



## Hyperparameters

- No. of hidden layers
- No. of neurons
- Learning rate
- Activation functions
- Optimisers
- Dropout values
- Weight initialization

← ANN

### CNN

- No. of convolution layers
- Stride
- Padding
- = No. of Kernels
- Pooling layer → Average / Global Maximum
- Batch size
- Epochs

Filter - Receptive field

Convolution Layers [1 - 10]

Pooling operation [Max, Min, Globals]

Hidden layers [1 - 25]

Hidden units [32 - 512]

Activation functions [ReLU, Leaky ReLU, FReLU, Exponential ReLU, tanh]

Optimizers = [Gradient descent, Stochastic Gradient Descent, SGD with momentum,

Adagrad]

RMS prop

Adam]

Learning rate: [0.4 to 0.8]

Epochs: [200 to 500]

Dropout: [0.4 to 0.7)

Weight / Bias

problems  
on CNN and  
ANN

construct  
CNN

Grid search

= [RandomSearch /  
Genetic Algorithm /  
Particle Swarm Algorithm]