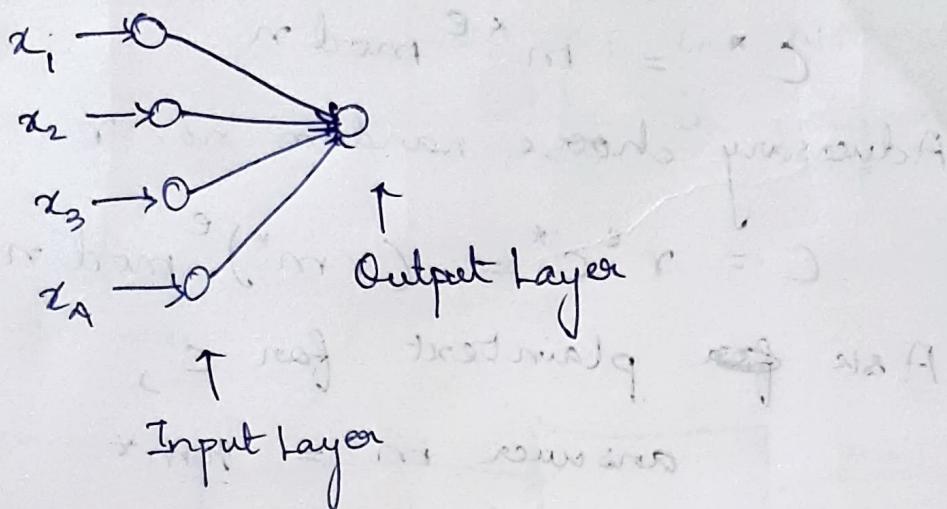


# Deep Learning

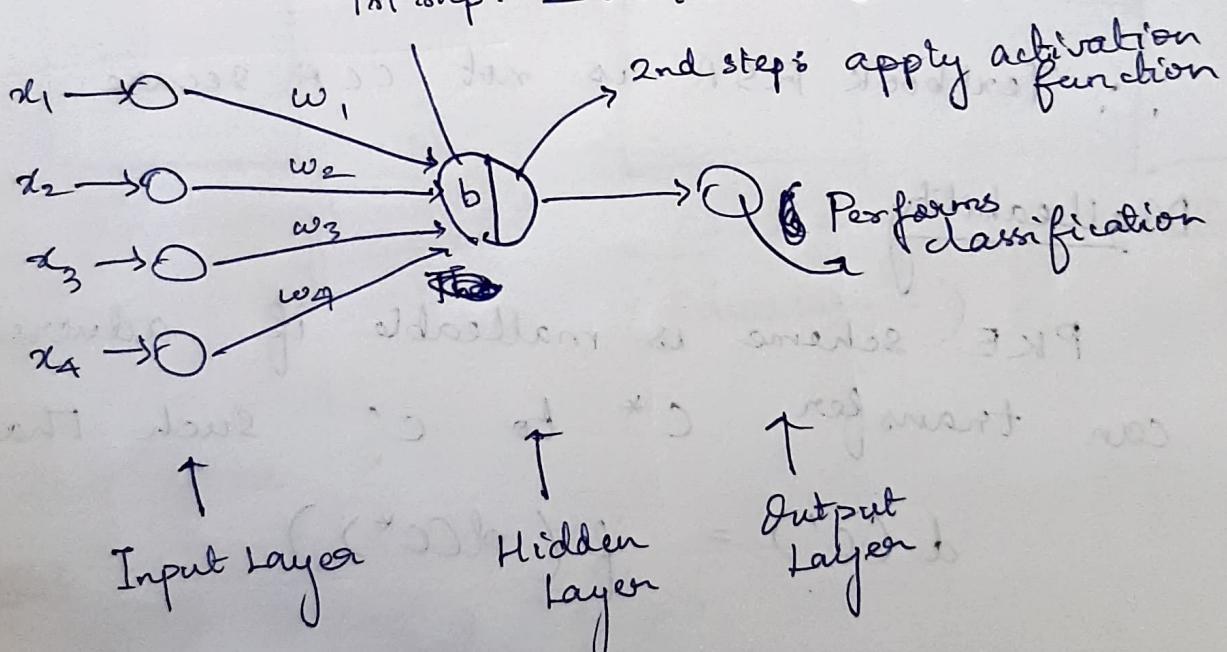
## Perceptron

### 1. Single Layer perceptron



## Artificial Neural Network

$$\text{1st step: } \sum w_i x_i + b \text{ (bias)}$$



## Forward propagation

Input  $\rightarrow$  Input layer  $\rightarrow$  Hidden layer  $\rightarrow$  Output layer  $\rightarrow$  Output

# Deep Learning

## ⑩ Single layer perceptron

1. No hidden layer
2. No activation function
3. Only applicable when data is linearly separable

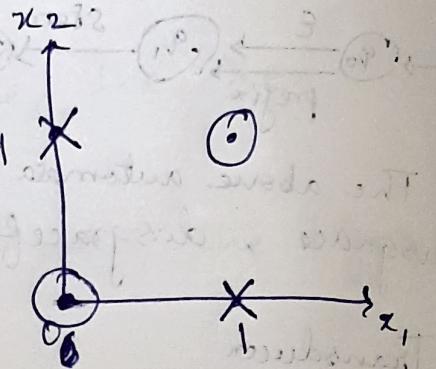
## Multi Layer perception

Hidden layer is present  
Activation function

3. Can be applied for non-linear data

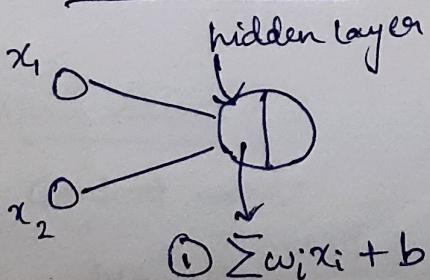
XOR

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



The points cannot be separated by a linear line. This data is not linearly separable.

## Activation Function



## ② Activation function

Types of activation functions:

- (i) Linear
- (ii) Non-Linear

## Linear function

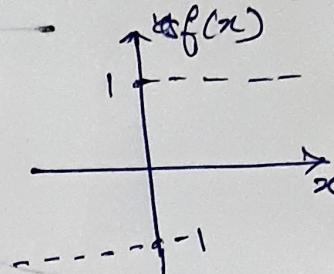
### (i) Step function

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

Domain =  $(-\infty, \infty)$  Range =  $\{0, 1\}$

### (ii) Signum function

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



Domain =  $(-\infty, \infty)$

Range =  $\{-1, 1\}$

### (iii) Linear function

$$f(z) = z$$

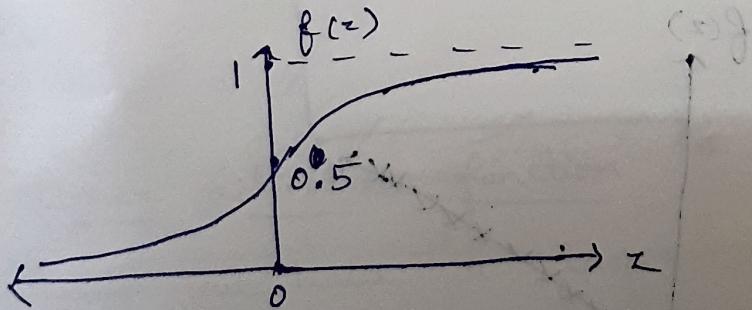
Domain =  $(-\infty, \infty)$

Range =  $(-\infty, \infty)$

## Non-linear function

### (i) Sigmoid function

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



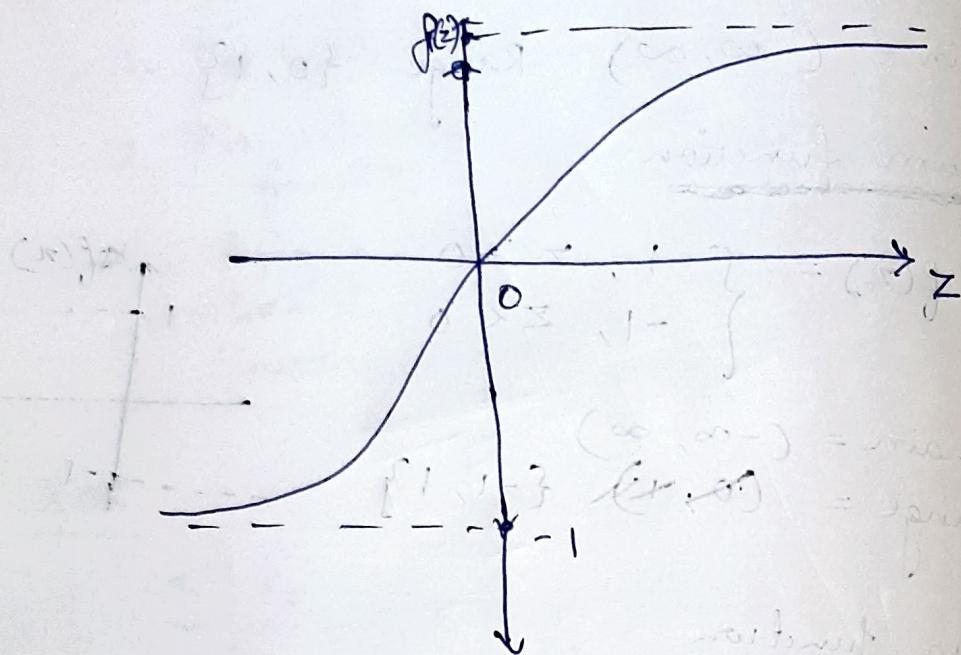
Domain =  $(-\infty, \infty)$

Range =  $(0, 1)$

Since this function doesn't touch 0, it is called non-zero central activation function.

(ii) tanh activation function / Hyperbolic tangent function

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

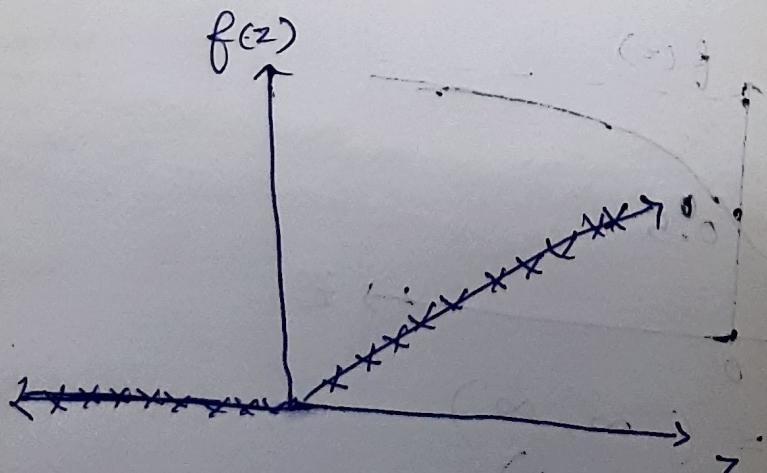


$$\text{Range} = (-1, 1)$$

Since the graph crosses  $0$ , this is called zero central activation function.

(iii) ReLU (Rectified Linear Unit)

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

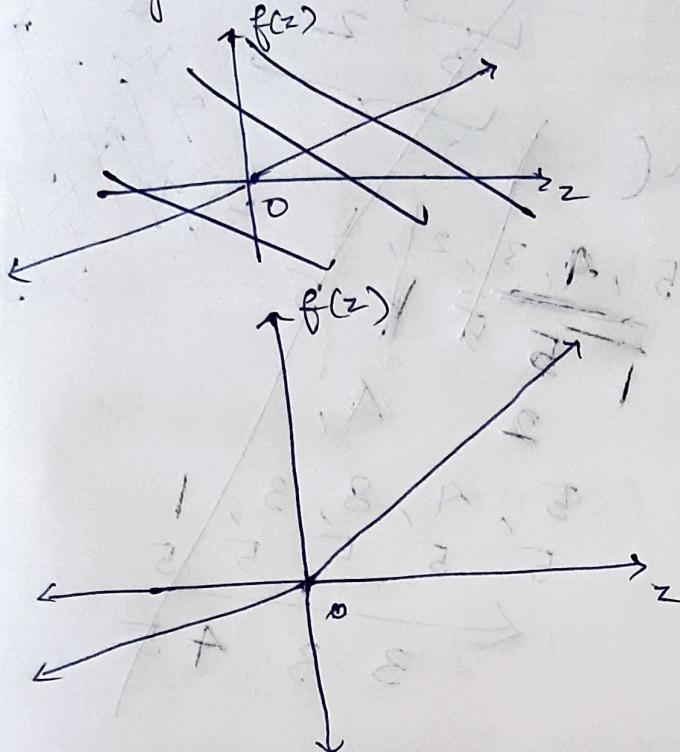


$$\text{Range} = (0, \infty)$$

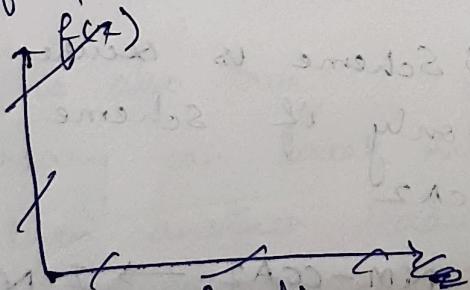
#### (iv) Leaky ReLU

$$f(z) = \begin{cases} z, & z \geq 0 \\ az, & z < 0 \text{ and } a \text{ is any constant} \end{cases}$$

$$\text{Range} = (-\infty, \infty)$$



Note:  
For Multi-class classification, we use soft-max activation function



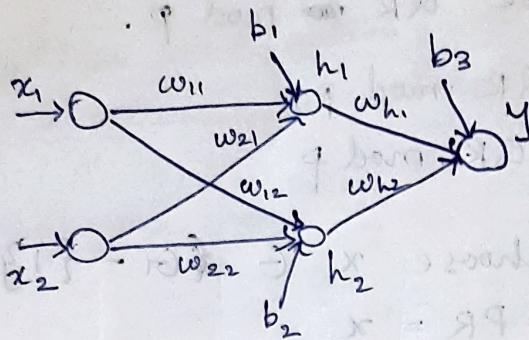
#### Softmax activation function

Suppose we have \$n\$ classes

$$f(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

$$f(z_n) = \frac{e^{z_n}}{e^{z_1} + e^{z_2} + \dots + e^{z_n}}$$

# Deep Learning



$$(h_1)_{in} = w_{11} \cdot x_1 + w_{21} \cdot x_2 + b_1$$

$$(h_1)_{out} = \text{Activation}((h_1)_{in})$$

$$(h_2)_{in} = w_{12} \cdot x_1 + w_{22} \cdot x_2 + b_2$$

$$(h_2)_{out} = \text{Activation}((h_2)_{in})$$

$$y_{in} = w_{h1} * (h_1)_{out} + w_{h2} * (h_2)_{out}$$

$$y_{out} = \text{Activation}(y_{in})$$

## Problem

$x_1$	$x_2$	$y_{out}$
0	0	0
0	1	1
1	0	1
1	1	0

$$w_{h1} = w_{21} = 2$$

$$w_{12} = w_{22} = -2$$

$$[b_1, b_2] = [-1, 3]$$

$$w_{h1} = w_{h2} = 2$$

$$(0, 0)$$

$$b_3 = -3$$

$$(h_1)_{in} = -1$$

$$(h_1)_{out} = \frac{1}{1 + e^{-(h_1)_{in}}} = \frac{1}{1 + e^{-(-1)}} = 0.26 \approx 0$$

$$(h_2)_{in} = 3$$

$$(h_2)_{out} = \frac{1}{1 + e^{-3}} = 0.95 \approx 1$$

$$y_{in} = -1, y_{out} = \text{sigmoid}(-1) = 0.26 \approx 0$$

For  $(0, 0)$ ,  $y_{out} = 0$

$(0, 1)$ :

$$y_{(h_1)\text{in}} = 1, (h_1)_{\text{out}} = 0.75 \approx 1$$

$$(h_2)_{\text{in}} = 1, (h_2)_{\text{out}} = 0.75 \approx 1$$

$$y_{in} = 1, y_{out} = 0.75 \approx 1$$

For  $(0, 1)$ :  $y_{out} = 1$

$(1, 0)$ :

$$(h_1)_{\text{in}} = 1, (h_1)_{\text{out}}$$

Same as  $(0, 1)$

$(1, 1)$ :

$$(h_1)_{\text{in}} = -1, (h_1)_{\text{out}} = 0.26 \approx 0$$

$$(h_2)_{\text{in}} = 3, (h_2)_{\text{out}} = 0.95 \approx 1$$

$$y_{in} = -1, y_{out} = 0.26 \approx 0$$

For  $(1, 1)$ :  $y_{out} = 0$

$x_1 \quad x_2$

0

0

0

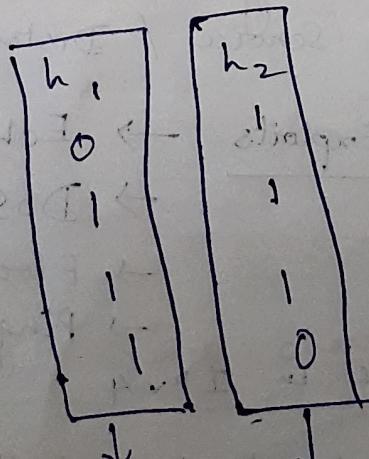
1

1

0

1

1



Similar to  
OR gate

Similar to  
NAND gate

Use ReLU for this problem:

$$w_{11} = w_{12} = 1, w_{21} = w_{22} = -1, b_1 = 0, b_2 = 0$$

~~b<sub>1</sub>, b<sub>2</sub>~~ = 0, -1, (0, 0) not

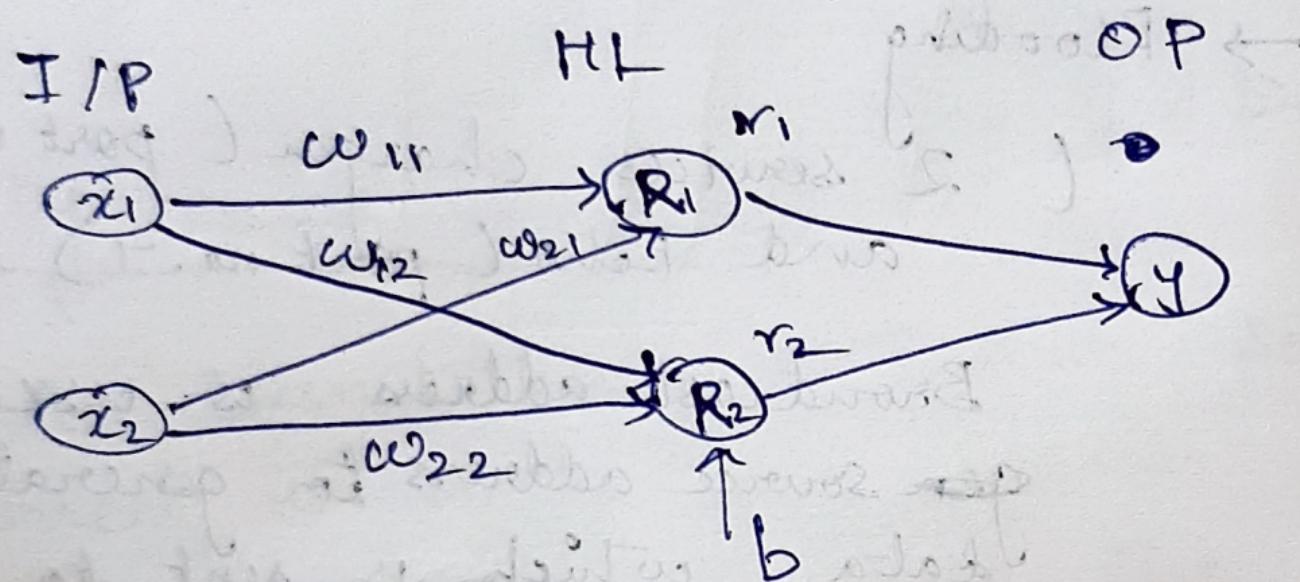
$$w_{h1}, w_{h2} = 1, -2$$

$$b_3 = 0 \quad (1, 1), 1 = 0(1, 1)$$

$x_1$	$x_2$	$(h_1)_{in}$	$(h_1)_{out}$	$(h_2)_{in}$	$(h_2)_{out}$	$y_{in}$	$y_{out}$
0	0	0	0	0	0	0	0
0	1	1	1	0	0	1	1
1	0	1	1	0	0	1	1
1	1	2	2	1	1	0	0

# Deep learning

## Forward Neural Network (Multi-Layer Perception)



$y = f(x, \theta)$ , where

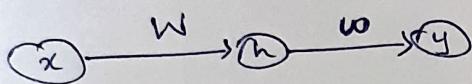
$\theta$  - parameters updated during learning

training

$$\begin{bmatrix} w_{11} & w_{21} \\ w_{12} & w_{22} \end{bmatrix}_{(2 \times 2)} \times \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_{(2 \times 1)} + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}_{(2 \times 1)} = \begin{bmatrix} w_{11}x_1 + w_{21}x_2 + b_1 \\ w_{12}x_1 + w_{22}x_2 + b_2 \end{bmatrix}$$

Weight matrix - (curr-layer  $\times$  prev-layer)

Representing weights as matrices, we get:



$$w = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$$

$$w = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0$$

$$x = \begin{bmatrix} x_1 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}, y = \begin{bmatrix} x_2 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}^T$$

$$xw + c = \begin{bmatrix} 0 & 0 \\ 1 & 1 \\ 1 & 1 \\ 2 & 2 \end{bmatrix}$$

$$\text{Add } c = ((xw)^T + c)^T = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}$$

Apply ReLU,

we get,

$$H =$$

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

Input space =  $\{(0,0); (0,1), (1,0), (1,1)\}$

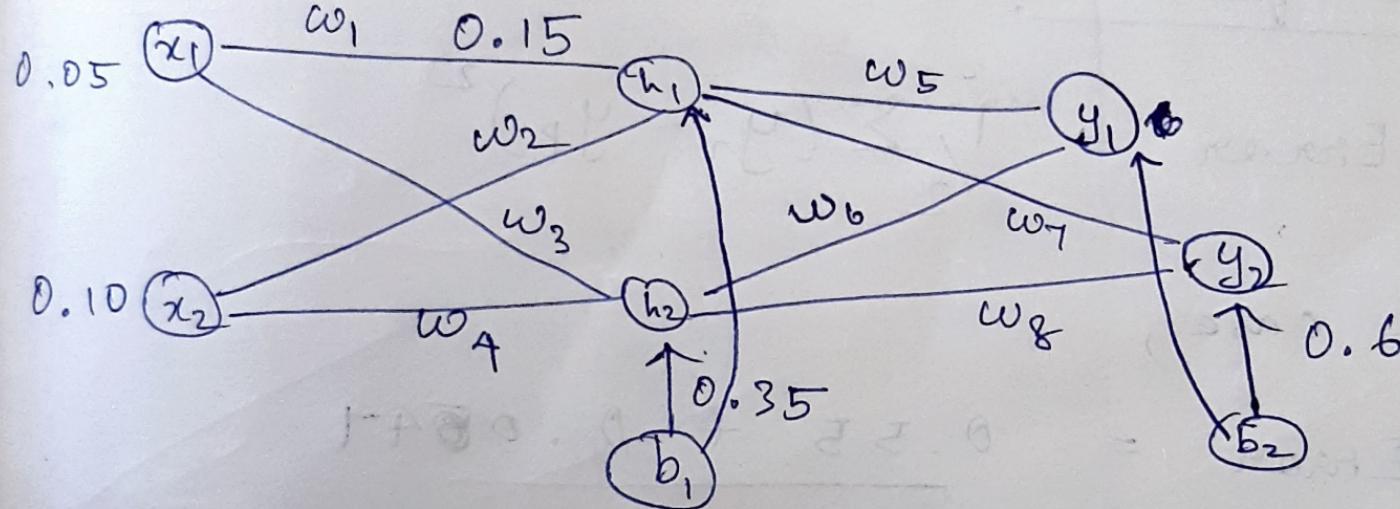
Hidden  
Layer space =  $\{(0,0); (1,0), (2,1)\}$

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

Note

We are getting target output the same as actual output because, the weights and biases are proper.

# Deep Learning



$$w_1 = 0.15$$

$$w_2 = 0.20$$

$$w_3 = 0.25$$

$$w_4 = 0.30$$

$$w_5 = 0.4$$

$$w_6 = 0.45$$

$$w_7 = 0.5$$

$$w_8 = 0.55$$

$$\begin{aligned}
 & \left[ \begin{matrix} 0.05 & 0.10 \\ \end{matrix} \right] \left[ \begin{matrix} 0.15 & -0.25 \\ 0.20 & 0.30 \\ \end{matrix} \right] + \left[ \begin{matrix} 0.35 & 0.35 \\ \end{matrix} \right] \\
 (\mathbf{x}) &= \left[ \begin{matrix} 0.0275 & 0.0425 \\ \end{matrix} \right] + \\
 & \quad \left[ \begin{matrix} 0.35 & 0.35 \\ \end{matrix} \right] \\
 &= \left[ \begin{matrix} 0.3775 & 0.3925 \\ h_1 & h_2 \end{matrix} \right]
 \end{aligned}$$

$$\text{Sigmoid} \left( \left[ \begin{matrix} 0.3775 & 0.3925 \\ \end{matrix} \right] \right) = \left[ \begin{matrix} 0.5933 & 0.5969 \\ r_1 & r_2 \end{matrix} \right]$$

$$\begin{aligned}
 \left[ \begin{matrix} 0.5933 & 0.5969 \\ \end{matrix} \right] \left[ \begin{matrix} 0.4 & 0.5 \\ 0.45 & 0.55 \\ \end{matrix} \right] &= \left[ \begin{matrix} 0.505 & 0.624 \\ \end{matrix} \right] \\
 + \left[ \begin{matrix} 0.60 & 0.60 \\ \end{matrix} \right] &= \left[ \begin{matrix} 1.105 & 1.224 \\ \end{matrix} \right]
 \end{aligned}$$

$$\text{Sigmoid} \left( \left[ \begin{matrix} 1.105 & 1.224 \\ \end{matrix} \right] \right) = \left[ \begin{matrix} 0.7514 & 0.7729 \\ y_1 & y_2 \end{matrix} \right]$$

Actual values	Predicted values
$y_1$	0.01
$y_2$	0.99

### Mean Square Error

$$\text{Error}^2 = \frac{1}{2} \sum (y_t - y_p)^2$$

In our case,

$$\begin{aligned}
 \text{Error} &= \frac{0.55 + 0.0041}{2} \\
 &= 0.2984
 \end{aligned}$$

## Error functions

For regression,

$$\text{Mean Square Error} = \frac{1}{2m} \sum (T_i - P_i)^2$$

For binary class classification problem

$$\text{Cross Entropy function} = -[y_i \log a_i + (1-y_i) \log(1-a_i)]$$

## Deep Learning

The objective in training a model is not only to reduce training error but also to minimize generalization error. This is done because models with large generalization error ~~will~~ tend to overfit.

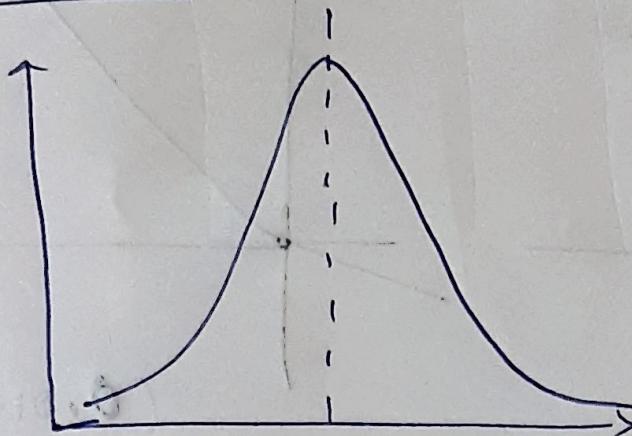
$$\therefore \text{Total cost} = \text{Training cost function} + \text{Regularization}$$

Common activation functions  
in Output layer

- Linear
- Sigmoid
- Softmax

Activation function based on different distributions (Output Layer)

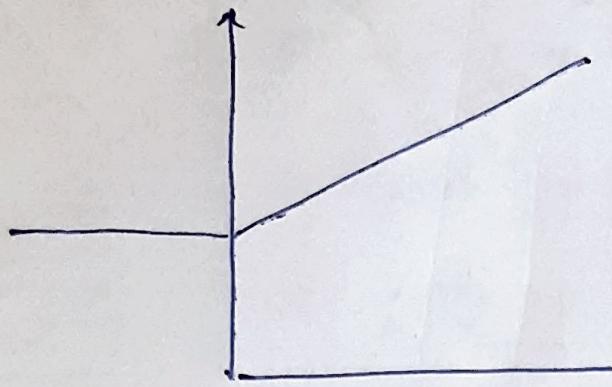
Normal distribution



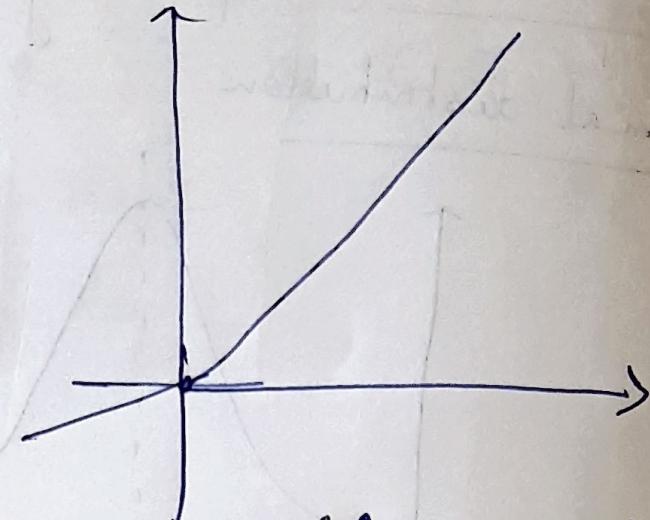
Linear activation function is preferred

<u>Activation function for output layer</u>	<u>Problem to be used for</u>
Linear	Output is continuous value (Regression problem)
Sigmoid	Binary class problem classification
Softmax	Multi class classification

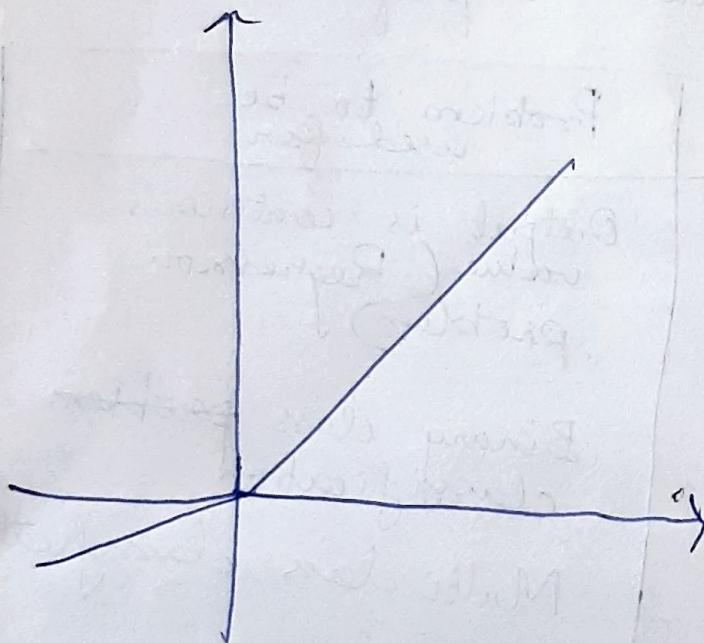
## Variants of ReLU



Absolute value  
rectification



( $0.01z, z$ )  
Leaky ReLU



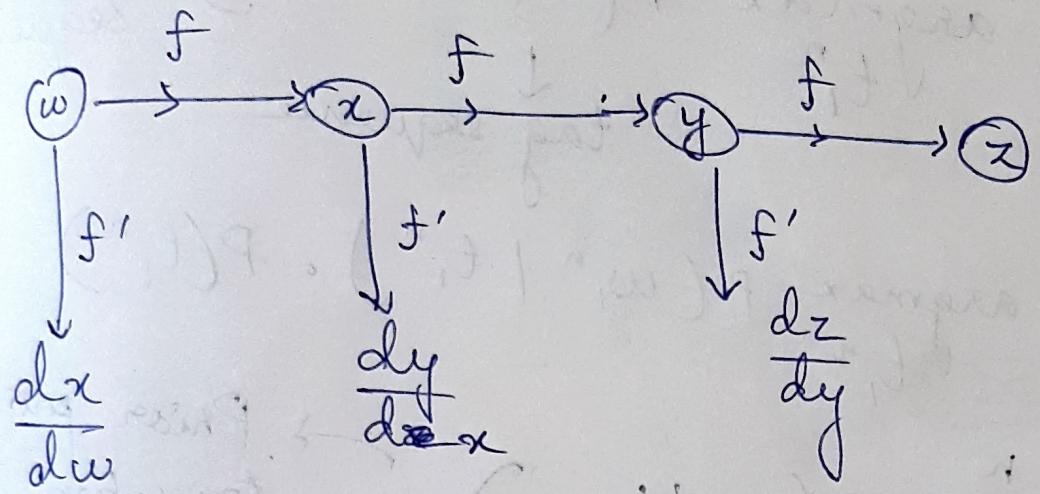
( $az, z$ )

Parametric ReLU

Maxout

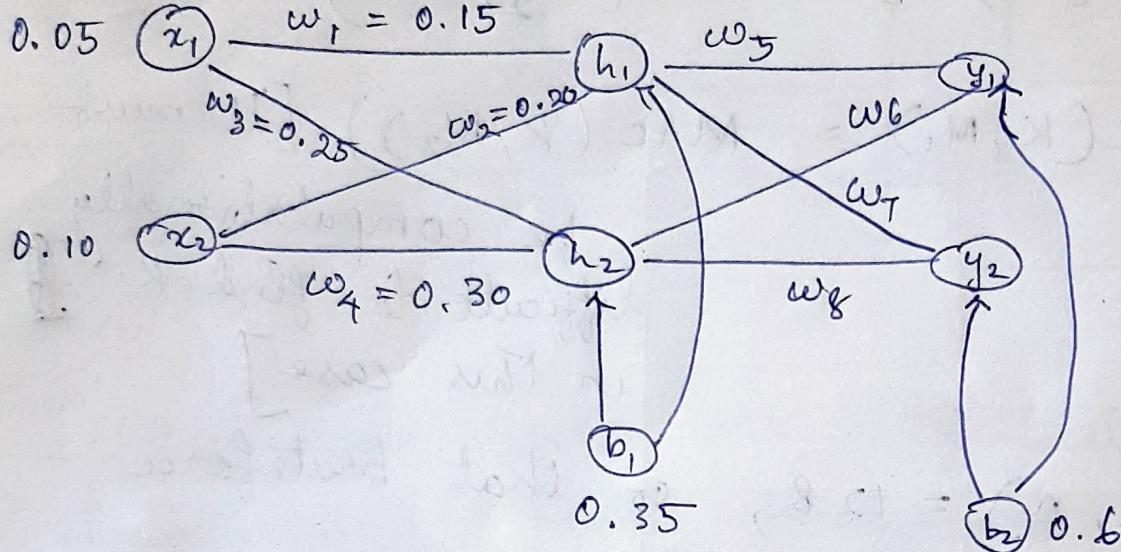
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

# Deep Learning



$$\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$$

$$\frac{dy}{dw} = \frac{dy}{dx} \cdot \frac{dx}{dw}$$



$$\begin{aligned}
 w_5 &= 0.4 \\
 w_6 &= 0.45 \\
 w_7 &= 0.5 \\
 w_8 &= 0.55
 \end{aligned}$$

Target values

$y_1$	$y_2$
0.01	0.99

Error = 0.2984 (not acceptable)

## Deep learning

Update weights :

$$w = w - \alpha \nabla, \quad \nabla = \frac{\partial \text{Error total}}{\partial w}$$

$\alpha$  - learning rate

$$E_{\text{total}} = \frac{1}{2} (T_1 - \text{out } y_1)^2 + \frac{1}{2} (T_2 - \text{out } y_2)^2$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = (T_1 - \text{out } y_1) \frac{\partial \text{out } y_1}{\partial w_5} + (T_2 - \text{out } y_2) \frac{\partial \text{out } y_2}{\partial w_5}$$

$$\frac{d E_{\text{total}}}{d w_5} = \frac{d E_{\text{total}}}{d \text{out } y_1} \times \frac{d \text{out } y_1}{d y_1} \times \frac{d y_1}{d w_5}$$

$$\frac{d E_{\text{total}}}{d \text{out } y_1} = \cancel{\text{out } y_1 - T_1} (\text{out } y_1 - T_1) = 0.7414$$

$$\frac{d \text{out } y_1}{d y_1} = \frac{-1 (e^{y_1}) (-1)}{(1 + e^{y_1})^2} = \frac{e^{y_1}}{(1 + e^{y_1})^2}$$

$$= \frac{\cancel{0.3312}}{(1 + \cancel{0.3312})^2}$$

$$\frac{d y_1}{d w_5} = \text{out } H_1 = 0.5933$$

0.1869

Reference:

$$T_1 = 0.01, \cancel{\text{out } y_1} = 0.7514$$

$$T_2 = 0.99 \quad \text{out } y_2 = 0.7729$$

$$\alpha = 0.5 \quad y_1 = 1.105$$

$$\frac{d E_{\text{total}}}{d w_5} = 0.0822$$

$$w_5 = w_5 - \alpha \frac{d E_{\text{total}}}{d w_5}$$

$$= 0.4 - (0.5)(0.0822)$$

$$= 0.3589$$

### Regularization

- to avoid overfitting, we use regularization parameter

Note:

Overfitting is a situation where model performs well on training but poorly on testing.

We have <sup>many</sup> ~~2~~ methods of regularization

- (i) L2  $\Rightarrow |slope|^2$  is Regularization parameter
- (ii) L1  $\Rightarrow |slope|$  " "
- (iii) Data Augmentation
- (iv) Early stopping criterion
- (v) Parameter sharing
- (vi) Dropout
- (vii) Boosting & Bagging

# Deep Learning

## Convolutional Network

Convolution operation is performed by applying a filter to the input.

E.g I/P  $\rightarrow$  1, 2, 3, 4, 5

Filter  $\rightarrow$  0 1 0

Multiply an element in a window with its respective element in the filter.

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} \quad 4 \quad 5$$

$\downarrow$

$$1^*0 + 2^*1 + 3^*0 = 2$$

$$1 \quad \begin{bmatrix} 2 & 3 & 4 \end{bmatrix} \quad 5$$

$\downarrow$

$$3$$

$$\text{Output} = 2, 3, 4$$

But we are missing the edge elements 1 and 5, so we add padding to the input.

New input is : 0 1 2 3 4 5 0

New output after padding : 1 2 3 4 5

No. of padding to add depends on the dimension of the filter.

If filter is of size  $s \times t$ ,

No. of rows to be added =  $s-1$

No. of cols to be added =  $t-1$

Note:

$\begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$  - filter can be used to detect lines in an image.

Case 1:

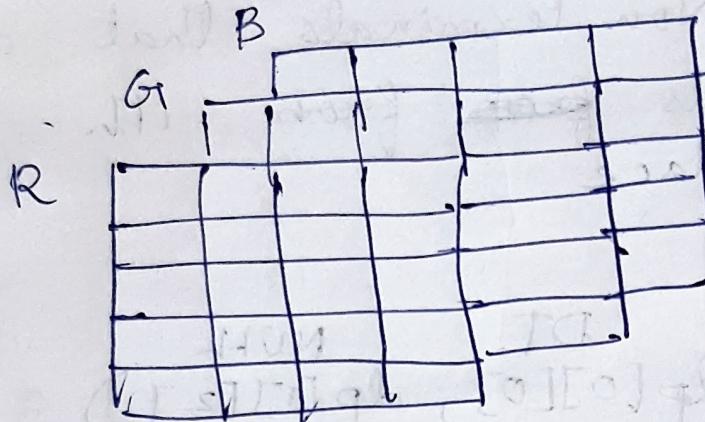
Single filter & Single image (only 1 output)

Case 2:

Multiple filters & Single image  
(as many outputs as filters)

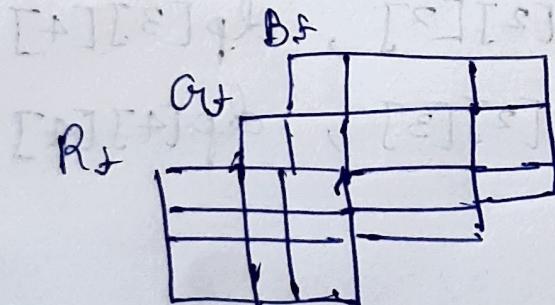
## Deep Learning

Case 3 (Multi-Channel <sup>input</sup> Image)



$m \times n \times 3$

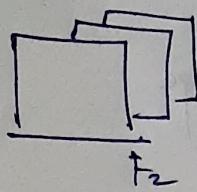
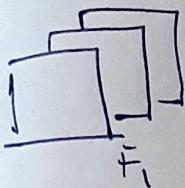
~~Filter with 3  $\delta$~~   
We use a 3 channel filter



R<sub>f</sub>, G<sub>f</sub>, B<sub>f</sub> will be applied to R, G, B respectively.

## case 4 (Multiple filter with Multiple channel)

multiple filters are applied on a single channel.



$F_1, F_2$  are separate filters, each of them being Multi-Channel filters. Output will be same as number of filters.

Note:

Output obtained by applying filters ~~are~~ called are nothing but features.

### Types of convolution

- \* stride
- \* Transpose
- \* Dilation

#### Stride ( $x, y$ )

- controls how the filter moves
- parameter  $x$  shows how much the filter moves in the horizontal direction
- parameter  $y$  in vertical direction

#### Transpose

- ~~use~~ to get output of greater size than image
- adds padding around the image

#### Dilation

- adds ~~one~~ zeroes in between pixels
- used for upscaling an image

In CNN, at each Neuron, the following operation is performed:

$$\text{Activation} (\sum F_i * I_i + b_i)$$

Reute

$F_i$  - filter  
 $I_i$  - image

Pooling

- There are 2 types of pooling: ~~max~~ and Max and Average pooling

Representation of Convolution ~~Altered~~ operation

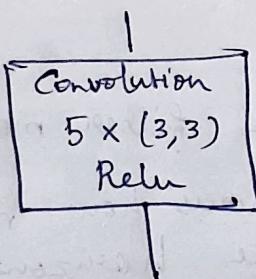
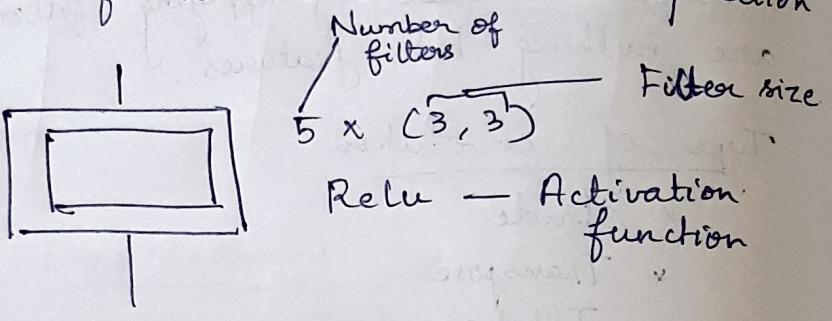


Image (600x600)

Conv  
 $8 \times (5, 5)$   
 $S = (2, 2)$   
 $P = (2, 2)$

→

Output  
( $300 \times 300$ )  
layer

Convolution  
 $40 \times (3, 3)$   
 $S = (3, 3)$   
 $P = (1, 1)$

4 output layers  
( $100 \times 100$ )