

Hebbian Learning Rule :- It is one of the first and easiest learning rules in ANN. It is a single layer network, i.e. has one input layer and one output layer. The input layer can have many ~~input~~ units say n . The output layer ~~can have many~~ has one unit. Hebbian rule works by updating the weights between neurons in the neural network for each training sample. It is unsupervised learning rule.

Hebbian Learning Algorithm,

1. Set all weights to zero; $w_i = 0$ for $i=1$ to n and bias to zero.
2. For each input vector, s (input vector) : t (target output pair), repeat steps 3-5.
3. Set activations for input units with the input vector $x_i = s_i$ for $i=1$ to n .
4. Set the corresponding output value to the output neuron, i.e. $y = t$.
5. Update the weight and bias by applying Hebb rule for all $i=1$ to n :

$$w_i(\text{new}) = w_i(\text{old}) + x_i y$$

$$b(\text{new}) = b(\text{old}) + y$$

Implementing AND Gate				Target
x	x_1	x_2	b	y
x_1	-1	-1	1	y_1
x_2	-1	1	1	y_2
x_3	1	-1	1	y_3
x_4	1	1	1	y_4

there are 4 training samples, so there will be 4 iterations.
Also activation function used here is Bipolar Sigmoidal
function so the range is [-1, 1];

Step 1: Set the weight and bias to zero, $w = [0 \ 0 \ 0]^T$,
and $b = 0$

Step 2: Set input vector $x_i = s_i$ for $i = 1 \dots 4$

$$x_1 = [-1 -1 1]^T$$

$$x_2 = [-1 1 1]^T$$

$$x_3 = [1 -1 1]^T$$

$$x_4 = [1 1 1]^T$$

Step 3: output value is set to $y = t$.

Step 4: modifying weights using Hebbian rule.

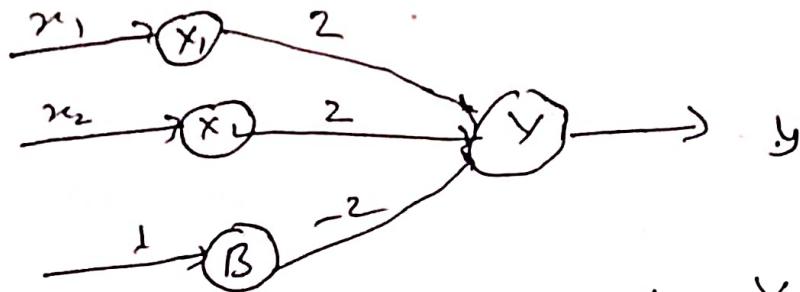
$$\begin{aligned} w(\text{new}) &= w(\text{old}) + \alpha_t y_i \\ &= [0 \ 0 \ 0]^T + [-1 -1 1]^T \cdot [-1] \\ &= [1 \ 1 -1]^T \end{aligned}$$

Second Iteration:

$$\begin{aligned} w(\text{new}) &= [1 \ 1 -1]^T + [1 -1 1]^T \cdot [1] \\ &= [2 \ 0 -2]^T \end{aligned}$$

$$\begin{aligned} \text{Third Iteration } w(\text{new}) &= [2 \ 0 -2]^T + [1 1 1]^T \cdot [1] \\ &= [1 \ 1 -3]^T \end{aligned}$$

Fourth iteration:- $\omega(\text{new}) = [1 \ 1 \ -3]^T + [1 \ 1 \ 1]^T (1)$
 $= [2 \ 2 \ -2]^T$



for $x_1 = -1, x_2 = -2, b = 1, y = (-1)2 + (-2)2 + 1(-2) = -6$

for $x_1 = -1, x_2 = 1, b = 1, y = (-1)2 + 1 \cdot 2 + 1(-2) = -2$

for $x_1 = 1, x_2 = -1, b = 1, y = 1 \cdot 2 + (-1)2 + 1(-2) = -2$

for $x_1 = 1, x_2 = 1, b = 1, y = 1 \cdot 2 + 1 \cdot 2 + 1(-2) = 2$

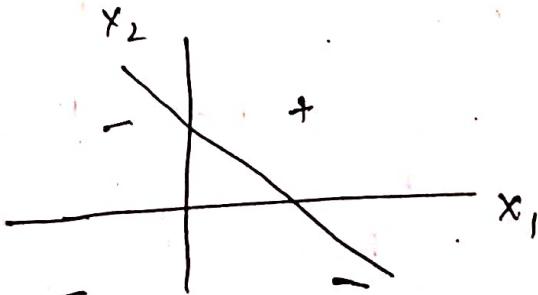
Decision boundary:

$$2x_1 + 2x_2 - 2b = y$$

Replacing y with 0, $2x_1 + 2x_2 - 2b = 0$

since bias; $b = 1, \text{ so } 2x_1 + 2x_2 - 2 \cdot 1 = 0$

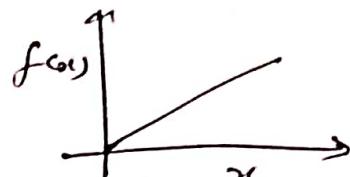
$$x_2 = 1 - x_1$$



Activation Functions - the activation function is used to calculate the output response of a neuron. The sum of signal is applied to with an activation to obtain the response. For neurons in same layer, same activation functions are used. There may be linear as well as nonlinear activation functions. The nonlinear activation functions are used in a multilayered.

(i) Identity Function

$$f(x) = x; \text{ for all } x.$$

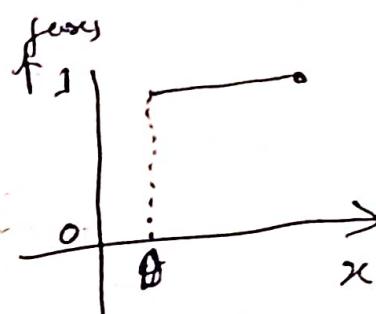


(ii) Binary Steps Function

BAB10329

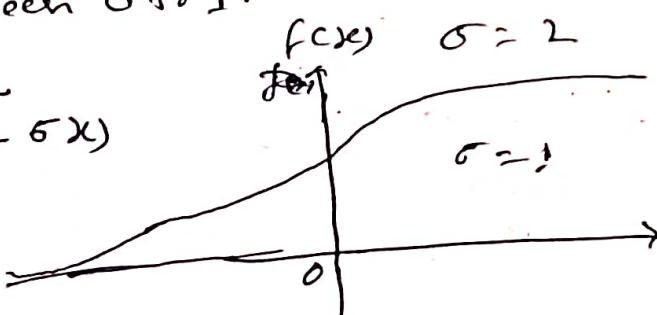
$$f(x) = \begin{cases} 1 & \text{if } f(x) \geq 0 \\ 0 & \text{if } f(x) < 0 \end{cases}$$

where θ is some threshold value.



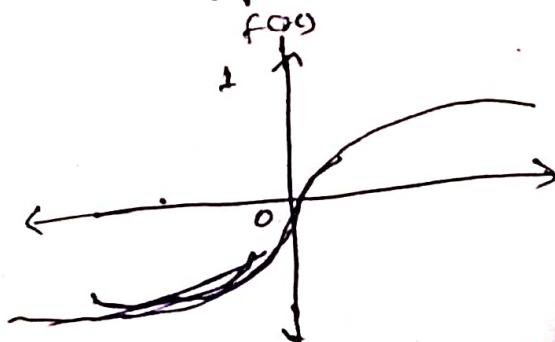
(iii) (A) Binary Sigmoidal Function - this is also called logistic function. It ranges between 0 to 1.

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



(B) Bipolar Sigmoidal Function - the desired range here is between +1 and -1. This function is hyperbolic tangent function.

$$\begin{aligned} b(x) &= 2f(x) - 1 = \\ &= \frac{1 - \exp(-\sigma x)}{1 + \exp(-\sigma x)} \end{aligned}$$



Perception Learning Rule:

For the perception learning rule, the learning signal is the difference between the desired and actual neuron's response. This type of learning is supervised.

Training Algorithm

Step 1: Initialize weights and bias (initially it can be zero). Set learning rate α (0 to 1).

Step 2: While stopping condition is false do steps 3-7.

Step 3: for each training pair (s, t) do step 4-6.

Step 4: Set activations of input units:

$$x_i = s_i \text{ for } i=1 \text{ to } n.$$

Step 5: Compute the output unit response

$$y_{in} = b + \sum_i x_i w_i$$

The activation function used is

$$y = f(y_{in}) = \begin{cases} + & \text{if } y_{in} > 0 \\ 0 & \text{if } -\theta \leq y_{in} \leq \theta \\ -1 & \text{if } y_{in} \leq -\theta \end{cases}$$

Step 6: The weights and bias are updated if the target is not equal to the output response.

if $t \neq y$ and one value of x_i is not zero

$$w_{i(\text{new})} = w_{i(\text{old})} + \alpha t x_i$$

$$b(\text{new}) = b(\text{old}) + \alpha t$$

$$\text{else } w_{i(\text{new})} = w_{i(\text{old})}$$

$$b(\text{new}) = b(\text{old})$$

Step 7: Test for stopping condition.

For AND Gate Function :-

x_1	x_2	b	$+/-$
-1	-1	1	-1
-1	1	1	-1
1	-1	1	-1
1	1	1	1

Step 1: Initialize weights $w_1 = w_2 = 0$ and $b = 0$, $t = 1$, $\theta = 0$

Step 2: for input pair $(-1, -1)$ do step 3-5

Step 3: Set activation of units

$$x_i = (-1, -1)$$

$$\text{Step 4: } y_{in} = b + \sum x_i w_i = 0 + -1 \times 0 + -1 \times 0 = 0$$

$$y = f(y_{in}) = 0$$

Step 5: Since $t \neq y$ so

$$\begin{aligned} w_{i(\text{new})} &= w_i(\text{old}) + \alpha t x_i \\ &= 0 + 1 \times (-1) [-1 \quad -1] \\ &= [1 \quad 1] \end{aligned}$$

$$b_{\text{new}} = b_{\text{old}} + \alpha t = 0 + 1 \cdot 1 = 1$$

The new weights and bias are $[1 \quad 1 \quad 1]$

$$x_2 = (-1, 1), y_{in} = 1 + -1 \times 1 + 1 \times 1 = 1, y = f(y_{in}) = 1$$

$$\text{so } y_{in} = 1, \text{ since } t = 1$$

$$\begin{aligned} w_{\text{new}} &= [1 \quad 1] + 1 \cdot (-1) [-1 \quad 1] \\ &= [1 \quad 1] + [1 \quad -1] = [2 \quad 0] \end{aligned}$$

$$b_{\text{new}} = 1 + 1 \cdot (-1) = 0,$$

$$x_3 = (1, -1), y_{in} = 0 + [1 \quad -1] \cdot [2 \quad 0] = 0 + 2 \times 2 + -1 \times 0 = 4$$

$$= 0 + 2 \times 2 + -1 \times 0 = 4$$

$$y = f(y_{in}) = f(4) = 1, b_{\text{new}} = 0 + 1 \cdot (-1) = -1$$

$$x_4 = (1, 1), y_{in} = -1 + t$$

Since $y = 1 \neq t$,

$$\omega_{new} = [2 \ 0]$$

~~if~~ $y \neq t$, since $y = 1, t = -1$,

$$\omega_{new} = [2 \ 0] + 1 \cdot (-1) [1 \ 1]$$

$$= [1 \ 1] \oplus$$

$$b_{new} = -1 + 0 = -1$$

* for $x_4 = (1, 1), y_{in} = -1 + 1 \times 1 + 1 \times 1 = 1$

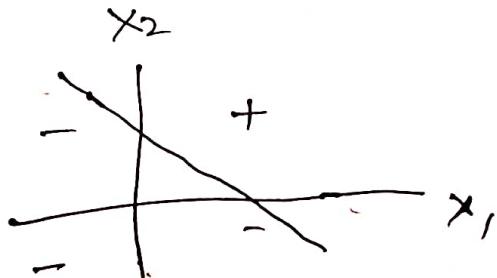
$$y = f(y_{in}) = f(1) = 1$$

Now $y = 1 = t$, so updated weight = [1 1]

and bias = -1

$$\Rightarrow -1 + x_1 + x_2 = 0$$

$$x_2 = 1 - x_1$$



Associative Memory Network :-

Associative neural nets are single layer nets in which the weights are determined to store an asset of pattern associations.

Hebb Rule for Pattern Association :-

Step 1: Initialize all weight ($i=1, \dots, n, j=1, 2, \dots, m$)
 $w_{ij} = 0$

Step 2: for each training input-target output vector, $\delta: t, d$

Step 3-5

Step 3: Set activation for input units to present training input ($i=1, 2, \dots, n$)

$$x_i = s_i$$

Step 4: Set activations for output units to current target output ($j=1, \dots, m$)

$$y_j = t_j$$

Step 5: Adjust the weights

$$w_{ij}(\text{new}) = w_{ij}(\text{old}) + \alpha_i y_j$$

Delta Rule for Pattern Association :-

$$\Delta w_i = \alpha (t - y_i) x_i$$

Delta Rule for Several Output Units :-

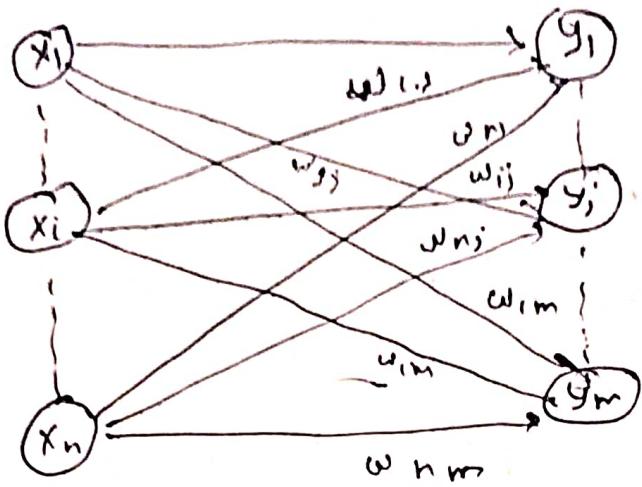
$$\Delta w_{ij} = \alpha (t_j - y_{inj}) x_i$$

Extended Delta Rule

$$\Delta w_{ij} = \alpha (t_i - y_j) x_i f'(y_{in})$$

Hetero Associative Memory Neural Network :-

Associative memory neural networks are networks in which the weights are determined in such a way that the net can store a set of p pattern associations. Hetero associative networks are static networks. Non linear or delay operation can be done using hetero associative networks. The weights may be found using the Hebb rule or delta rule.



Hetero Associative Neural Net.

Application Algorithm:

- Step 1: Weights are initialized using Hebb or delta rule.
- Step 2: for each input vector do steps 3 to 5
Set the activations for input layer units equal to
- Step 3: the current vector x_i
- Step 4: Compute the net input to the output units
- $$Y_{inj} = \sum_{i=1}^n x_i w_{ij}$$
- Step 5: Determine the activation of the output units
- $$Y_j = \begin{cases} 1, & \text{if } Y_{inj} > 0 \\ 0, & \text{if } Y_{inj} = 0 \\ -1, & \text{if } Y_{inj} \leq 0 \end{cases}$$

if the target

$$\textcircled{6} \quad Y_{in} = \sum x_i w_{ii}$$

$$= [1 \ 1 -1 -1] \begin{bmatrix} 1 & 1 & -1 & -1 \\ 1 & 1 & -1 & -1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & 1 & 1 \end{bmatrix}$$

$$Y_{in} = [4 \ 4 -4 -4]$$

$$y_2 = f(Y_{in}) = f[4 \ 4 -4 -4] \\ = [1 \ 1 -1 -1]$$

The response vector is same as the input vector, we can say that the input vector vector is recognized as a known vector.

Bidirectional Associative Memory :

In this network net iterates by sending a signal back and forth between the two layers until each neuron remains constant for several steps. Three form of BAM are

- (1) Binary
- (2) Bipolar
- (3) Continuous

Activation function for binary associative memory is

$$y_{ij} = \begin{cases} 1 & \text{if } Y_{inj} \geq 0 \\ \omega_j & \text{if } Y_{inj} = 0 \\ 0 & \text{if } Y_{inj} < 0 \end{cases}$$

$$x_i = \begin{cases} 1 & \text{if } X_{inj} > 0 \\ x_i & \text{if } X_{inj} = 0 \\ 0 & \text{if } X_{inj} < 0 \end{cases}$$

For bipolar vector, the activation function for Y-layer is

$$y_j = \begin{cases} 1 & \text{if } y_{inj} \geq v_j \\ y_j & \text{if } y_{inj} = v_j \\ 0 & \text{if } y_{inj} < v_j \end{cases}$$

$$x_i = \begin{cases} 1 & \text{if } x_{ini} \geq v_i \\ x_i & \text{if } x_{ini} = v_i \\ 0 & \text{if } x_{ini} < v_i \end{cases}$$

For Continuous BAM.

Y-layer $f(y_{inj}) = \frac{1}{1 + \exp(-y_{inj})}$
if bias is calculated $y_{ini} = b_i + \sum x_i w_{ij}$

X-layer $f(x_{ini}) = \frac{1}{1 + \exp(-x_{ini})}$

if bias is calculated

$$x_{ini} = b_i + \sum y_j w_{ij}$$

The memory capacity of BAM is $\min(n, m)$
where,

n is the number of units in X-layers

m is the number of units in Y-layers.

Training Algorithm for Medialine

Step 1: Initialize weights and bias, set learning rate.

$v_1 = v_2 = 0.5$ and $b_3 = 0.5$, other weights may be small random values.

Step 2: When stopping condition is false do steps 3-9.

Step 3: For each bipolar training pair s_i, t_i , do steps 4-9.

Step 4: Set activations of input units:

$$x_i = s_i \text{ for } i=1 \text{ to } n$$

Step 5: Calculate the net input of hidden adaline units

$$Z_{in1} = b_1 + x_1 w_{11} + x_2 w_{21}$$

$$Z_{in2} = b_2 + x_1 w_{12} + x_2 w_{22}$$

Step 6: Find the output of hidden adaline units using activation function for Z_1, Z_2 and y is given by

$$Z_1 = f(Z_{in1})$$

$$Z_2 = f(Z_{in2})$$

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

Step 7: Calculate net input to output

$$Y_{in} = b_3 + Z_1 v_1 + Z_2 v_2$$

Apply activation to get the output of ret.

$$Y = f(Y_{in})$$

Step 8: Find the error and do weight updation

if $t = 4$, no weight updation

if $t \neq 4$, then,

if $t = 1$, then update weight on z_i ; unit whose net input is closest to 0.

$w_{ii}(\text{new}) = w_{ii}(\text{old}) + \alpha (1 - z_{inj})^{2\ell_i}$

$b_i(\text{new}) = b_i(\text{old}) + \alpha (1 - z_{in})$

if $t = -1$, then update weights on all unit z_k which have positive net input

$w_{ik}(\text{new}) = w_{ik}(\text{old}) + \alpha (-1 - z_{ink}) x_i$

$b_k(\text{new}) = b_k(\text{old}) + \alpha (1 - z_{ink})$

Step 9: Test for the stopping condition.

Delta Learning Rule (Widrow-Hoff Rule or Least mean square rule)

" The mean square error for a particular training pattern is

$$E = \sum_j (t_j - y_{inj})^2$$

The gradient of E is a vector consisting of the partial derivative of E with respect to each of the weight.

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_j (t_j - y_{inj})^2$$

$$= 2(t_j - y_{inj}) \frac{\partial y_{inj}}{\partial w_{ij}}$$

$$= -2(t_j - y_{inj})x_i$$

Hence error will be rapidly reduced depending upon the given learning by adjusting the weight according to the delta rule.

$$\Delta w_{ij} = \alpha (t_j - y_{inj}) x_i$$

Training Algorithm :-

Step 1: Initialize weights (not zero but small random values are used). Set learning rate α .

Step 2: While stopping condition is false do step 3-7.

Step 3: For each bipolar training pair (s, t) perform 4-6.

Step 4: Set activation of input unit $x_i = s_i$ for $i = 1$ to n .

Step 5: Compute $y_{in} = b + \sum x_i w_i$

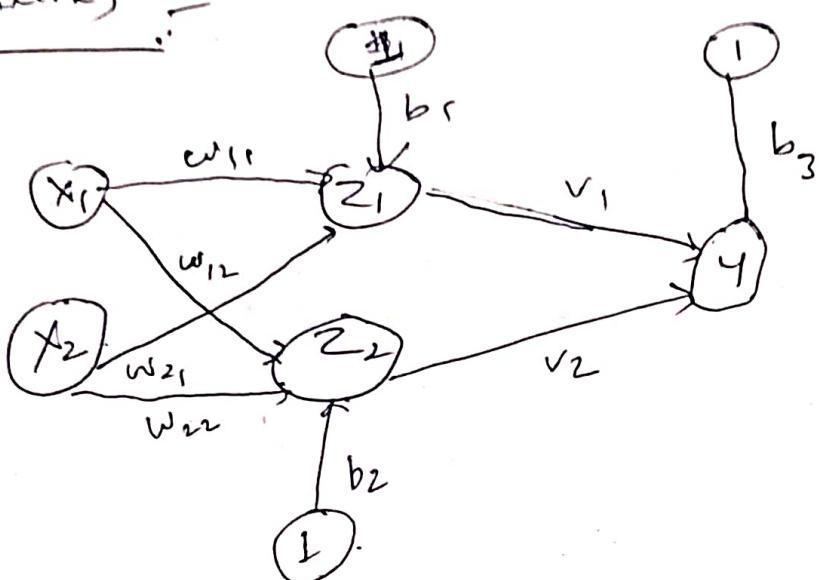
Step 6: Update bias and weight, $i = 1$ to n

$$w_{(new)} = w_{(old)} + \alpha (t - y_{in}) x_i$$

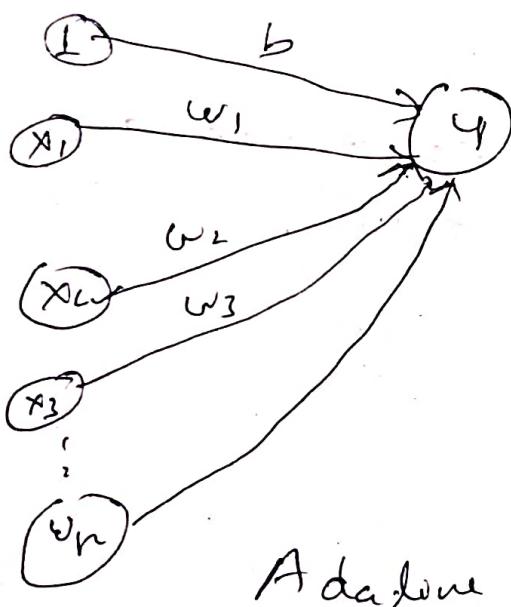
$$b_{(new)} = b_{(old)} + \alpha (t - y_{in})$$

Step 7: Test for stopping condition.

Adaptive Linear Neuron (Adaline) and Multilayered Adaline (Madalone) :-



Madalone



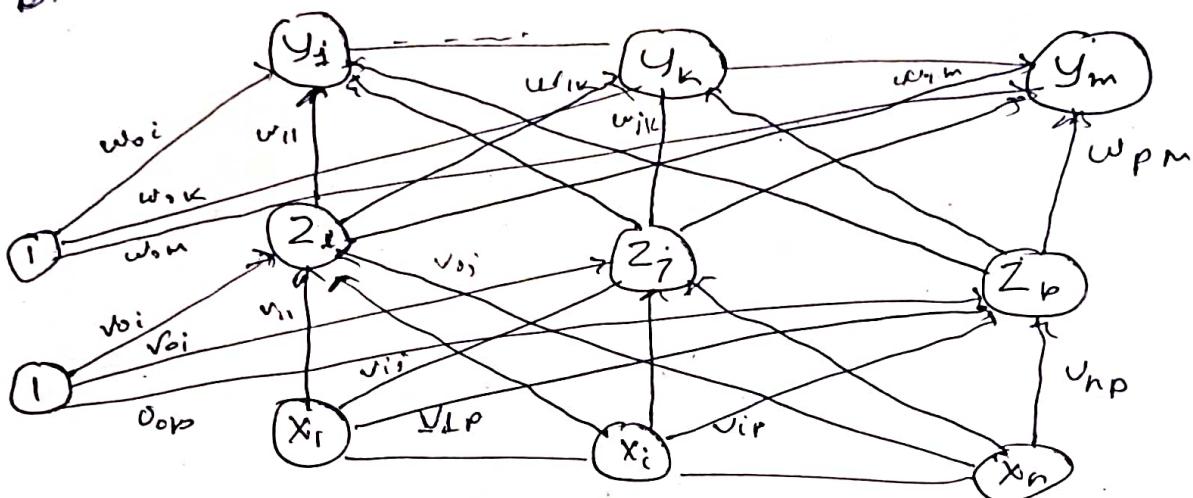
Adaline

Adaline and madalone network uses the Least mean square error (Lms error) rule. It uses error delta rule to update the weight.

Back Propagation Networks (BP Network)

It is a multi-layer forward network using extended gradient-descent based delta-learning rule, commonly known as back propagation (of errors) rule. The network is trained by supervised learning method.

Q.



Derivation

Consider an arbitrary activation function $f(x)$. The derivative of activation function is

Let

$$y_{ik} = \sum_i z_i w_{ik}$$

$$z_{inj} = \sum_i v_{ij} x_i$$

$$y_{ik} = f(z_{inj})$$

The error to be minimized is $E = 0.5 \sum_k (t_k - y_k)^2$

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} \left[0.5 \sum_k (t_k - y_k)^2 \right]$$

$$= -(t_k - y_k) \frac{\partial}{\partial w_{jk}} f(z_{inj})$$

$$= -(t_k - y_k) f'(z_{inj}) z_i$$

Let us define $\delta_k = -(t_k - y_k) f'(y_{ink})$

Weights on connections to the hidden unit z_j .

$$\frac{\partial E}{\partial v_{ij}} = - \sum_k (t_k - y_k) \frac{\partial}{\partial v_{ij}} y_k$$

$$= - \sum_k (t_k - y_k) f(y_{ink}) \frac{\partial}{\partial v_{ij}} y_{ink}$$

$$= - \sum_k \delta_k \frac{\partial}{\partial v_{ij}} y_{ink}$$

$$Now = - \sum_k \delta_k \frac{\partial}{\partial v_{ij}} (\sum_j z_j \cdot w_{jk})$$

$$= - \cancel{\sum_k \delta_k \frac{\partial}{\partial v_{ij}}} = - \sum_k \delta_k w_{jk} \frac{\partial}{\partial v_{ij}} z_j$$

$$= - \sum_k \delta_k w_{jk} \frac{\partial}{\partial v_{ij}} f(z_{ink})$$

$$= - \sum_k \delta_k w_{jk} f'(z_{ink}) (x_i)$$

$$\delta_j = - \sum_k \delta_k w_{jk} f'(z_{ink})$$

The weight update for output unit is given by

$$\Delta w_{jk} = -\alpha \frac{\partial E}{\partial w_{jk}}$$

$$= -\alpha \delta_k z_j$$

The weight update for hidden unit is given by

$$\Delta v_{ij} = -\alpha \frac{\partial E}{\partial v_{ij}}$$

$$= -\alpha \delta_j x_i$$

This is generalized delta rule used for training purpose.

Back Propagation Algorithm

- (I) Initialize weights to small random values
- (II) while stopping condition is false do step 3-10
- (III) for each training pair do steps 4-9
- (IV) Each input units receives the input signal x_i and transmits this signal to all hidden units.
- (V) Each hidden unit (z_j , $j=1, \dots, P$) sums its weighted input signals.

$$Z_{inj} = V_0 + \sum_{i=1}^n x_i V_{ij}$$

Applying activation function

$$z_j = f(Z_{inj})$$

and sends this signal to all units in the layer above i.e. output units.

- (VI) Each output unit (y_k , $k=1, \dots, m$) sums its weighted input signals

$$y_k = w_{0k} + \sum_{j=1}^P z_j w_{jk}$$

and applies its activation function to calculate the output signals.

$$y_k = f(y_k)$$

- (VII) Each output units (y_k , $k=1, \dots, m$) receives a target pattern corresponding to an input pattern. error information term is calculated as

$$\delta_k = (t_k - y_k) f'(y_k)$$

Step 8: Each hidden unit (z_j , $j=1, \dots, n$) sums its delta inputs from units in the layer above.

$$\delta_{inj} = \sum_{k=1}^m \delta_k w_{jk}$$

The error information term is calculated as

$$\delta_j = \delta_{inj} f'(z_{inj})$$

Step 9: Each output unit (y_k , $k=1, \dots, m$) updates its bias and weight ($j=0, \dots, p$)

The weight correction term is given by

$$\Delta w_{jk} = d\delta_k z_j$$

and bias correction term is given by

$$\Delta w_{0k} = d\delta_k$$

Therefore $w_{jk}(\text{new}) = w_{jk}(\text{old}) + \Delta w_{jk}$

$$w_{0k}(\text{new}) = w_{0k}(\text{old}) + \Delta w_{0k}$$

Each hidden units (z_j , $j=1, \dots, p$) updates its bias and weights ($i=0, \dots, n$)

The weight correction term

$$\Delta v_{ij} = d\delta_j x_i$$

The bias correction term

$$\Delta v_{0j} = d\delta_j$$

Therefore, $v_{ij}(\text{new}) = v_{ij}(\text{old}) + \Delta v_{ij}$,

$$v_{0j}(\text{new}) = \cancel{v_{0j}(\text{old})} \cdot v_{0j}(\text{old}) + \Delta v_{0j}$$

Step 10: Test the stopping condition.

Auto Associative Memory Network

In auto associative net the training input and target output vectors should be identical or same. The auto associative net training is often called storing the vectors, which may be binary or bipolar in nature.

Training Algorithm

① Initialize all weights, $i=1, \dots, n$, $j=1, \dots, h$

$$w_{ij} = 0_i$$

② for each vector to be stored follow steps 3-4.

③ Set activation for each input unit $i=1, \dots, n$

$$x_{i1} = s_i$$

④ Set activation for each output unit $j=1, 2, \dots, h$

$$y_{j1} = s_j$$

⑤ Adjust the weight for $i=1, \dots, n$ and $j=1, \dots, h$

$$w_{ij}(\text{new}) = w_{ij}(\text{old}) + x_{ij} y_j$$

Example: Use the Hebb rule to store the vector $(1 \ 1 \ -1 \ -1)$

in an auto associative neural net

⑥ Find weight matrix matrix

⑦ Test the input vector $x = (1 \ 1 \ -1 \ -1)$

Sol. ⑥ $w = \begin{pmatrix} 1 & 1 & -1 & -1 \\ 1 & 1 & -1 & -1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & 1 & 1 \end{pmatrix} [1 \ 1 \ -1 \ -1] = \begin{pmatrix} 1 & 1 & -1 & -1 \\ 1 & 1 & -1 & -1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & 1 & 1 \end{pmatrix}$

Associative memory :- Pattern association is the process of forming association between related patterns. Associative memory net can be seen as a simplified model of a human brain which can associate similar patterns. Associative neural nets are high single layer nets in which the weights are ~~updated~~ determined to store an array of pattern associations.

Example :- A hetero associative net is trained by Hebb outer product rule for input row vectors $s = (x_1, x_2, x_3, x_4)$

$$t_3 = (1 \ 0) \text{ to output vectors } t = (t_1, t_2) \cdot \text{Find the weight}$$

$$t_4 = (1 \ 0) \text{ to output vectors } t = (t_1, t_2) \cdot \text{Find the weight}$$

$$S_3 = (0 \ 0 \ 1 \ 1) \quad S_4 = (0 \ 1 \ 0 \ 0) \quad S_1 = (1 \ 1 \ 0 \ 0) \quad S_2 = (1 \ 1 \ 1 \ 0)$$

Sol. 1. Initialize weight to 0

$$w_1 = S_1^T t_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \end{bmatrix} [1 \ 0] = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$w_2 = S_2^T t_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} [0 \ 1] = \begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$w_3 = S_3^T t_3 = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} [1 \ 0] = \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}$$

$$w_4 = S_4^T t_4 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} [1 \ 0] = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$w = w_1 + w_2 + w_3 + w_4 = \begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 1 & 1 \\ 1 & 0 \end{bmatrix}$$

Logistic Regression, Logistic regression is an example of supervised learning. It is used to calculate or predict the probability of a binary (yes/no) event occurring. In linear regression, the outcome is continuous and can be any possible value. However in the case of logistic regression, the predicted outcome is discrete and restricted to a limited number of values.

Hence, linear regression is an example of a regression model and logistic regression is an example of a classification model.

Logistic regressions are three types

- ① Binary logistic regression ② multinomial logistic regression.
- ③ Ordinal logistic regression.

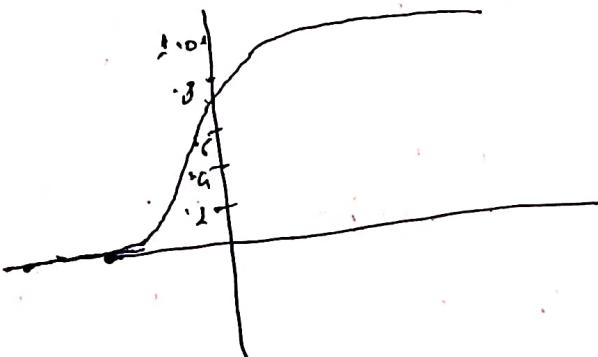
The logistic function is a simple S-shaped curve used to convert data into a value between 0 and 1.

Model, output = 0, 1

$$\text{Hypothesis} \Rightarrow Z = wX + B$$

$$h_\theta(x) = \text{sigmoid}(Z)$$

$$\text{sigmoid}(t) = \frac{1}{1 + e^{-t}}$$



Sigmoid Activation Function

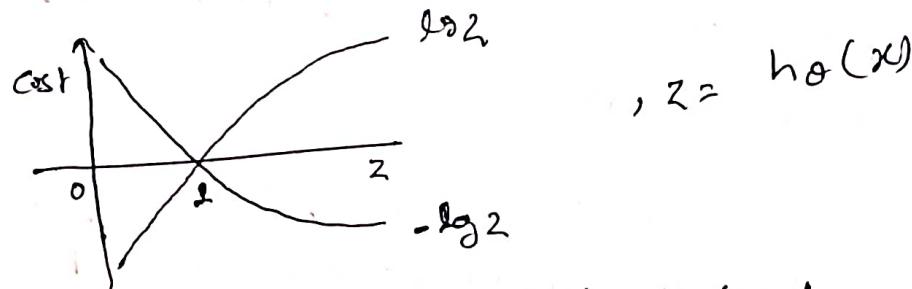
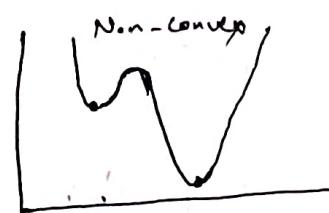
Decision Boundary → To predict which class a data belongs, a threshold can be set. Based upon this threshold, the obtained estimated probability is classified into classes. Decision boundary can be linear or non-linear.

Say, if predicted_value ≥ 0.5 , then classify email as spam else as not spam.

Cost Function \rightarrow

$$\text{Cost}(h_\theta(x), y(\text{actual})) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

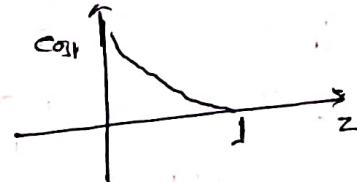
Linear regression uses mean squared error as its cost function. If this is used for logistic regression, then it will be a non-convex function of parameter θ . Gradient descent will converge into global minimum only if the function is convex.



$$\text{Cost}(h_\theta(x), y) = \begin{cases} -\log(h_\theta(x)) & \text{if } y=1 \\ -\log(1-h_\theta(x)) & \text{if } y=0 \end{cases}$$

if $y=1$

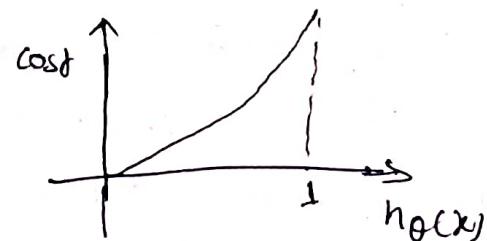
$$\text{cost}(h_\theta(x), y) = -\log(h_\theta(x))$$



if $\text{cost}=0 \Rightarrow y=1 \Rightarrow h_\theta(x)=1$

$\text{cost}=\infty$ for $h_\theta(x)=0$

if $h_\theta(x)=0$, it is similar to predicting $P(y=1|x; \theta)=0$



if $\text{cost}=0 \Rightarrow h_\theta(x)=0 \Rightarrow y=0$

$\text{cost}=\infty \Rightarrow h_\theta(x)=1$

if $h_\theta(x)=1$, it is similar to predicting $P(y=0|x; \theta)=0$

Why this cost function:-

Let us consider

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

\hat{y} is the probability that $Y=1$, given x

$$1 - \hat{y} = P(Y=0|x)$$

$$P(Y|X) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\Rightarrow \log P(Y|X) = \log(\hat{y}^y + (1-\hat{y})^{(1-y)})$$

$$\Rightarrow \log P(Y|X) = \log(\hat{y}^y \cdot (1-\hat{y})^{(1-y)})$$

$$\Rightarrow y \log \hat{y} + (1-y) \log(1-\hat{y})$$

$$\Rightarrow -L(\hat{y}, y)$$

$$\Rightarrow \log P(Y|X) = -L(\hat{y}, y)$$

This negative function is because when we train, we need to maximize the probability by minimizing loss function. Decreasing the cost will increase the maximum likelihood assuming that samples are drawn from an identically independent distribution.

Now

$$Z = w_1x_1 + w_2x_2 + b \rightarrow \begin{cases} \hat{y} = a = \sigma(z) \\ \hat{a} = \hat{y} \end{cases} \rightarrow L(\hat{y}, y)$$

$$w_1 \Rightarrow \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1}$$

$$\frac{\partial L}{\partial a} = \frac{\partial}{\partial a} (-y \log a - (1-y) \log(1-a))$$

$$\frac{\partial L}{\partial a} = \frac{-y}{a} - (-1) \frac{(1-y)}{(1-a)}$$

$$\frac{\partial a}{\partial z} = a(1-a)$$

$$\frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial L}{\partial w_1} = \left(\left[\frac{-y}{a} + \frac{(1-y)}{1-a} \right] a(1-a) \right) \cdot x_1 \\ = (a-y)x_1$$

update for w_1 ,

$$\frac{\partial L}{\partial w_1} = (a-y) \cdot x_1$$

$$\Rightarrow w_1 = w_1 - \alpha \frac{\partial L}{\partial w_1}$$

Similarly, for all parameters

$$w_i = w_i - \alpha \frac{\partial L}{\partial w_i} \quad i=1, 2, \dots, m$$

m = no. of parameters

$$b = b - \alpha \frac{\partial L}{\partial b}$$

$$\text{where, } \frac{\partial L}{\partial b} = (a-y)$$

Convolutional Neural Network

- ① Used in image recognition

-1	-1	1	-1	-1
-1	1	-1	0	-1
-1	1	1	-1	-1
-1	-1	1	0	
-1	-1	1		

Disadvantages of ANN using for Image classification

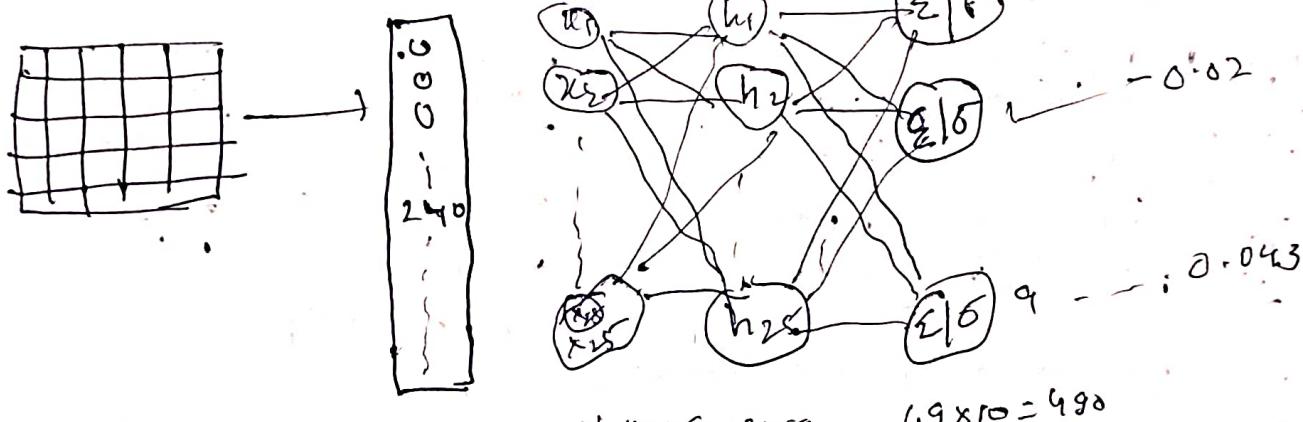
- ① Too much computation
- ② Treats local pixels same as pixels far apart
- ③ Sensitive to location of an object in an image.

Suppose we want to recognize on a 5×5 grid

- ④ CNN can recognize variety but ANN cannot.

e.g. 9 9 9

matrix is flattened to one dimensional vector



- ⑤ CNN is invariant of position.

Suppose image size $\approx 1920 \times 1080 \times 3$ (Cow image)

First layer neurons $= 1920 \times 1080 \times 3 = 6$ million

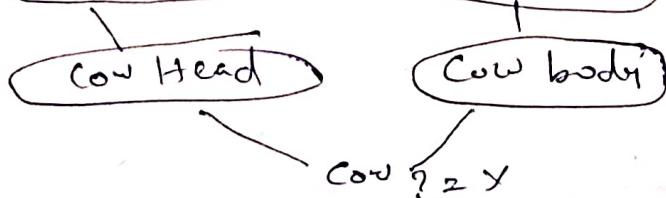
Hidden layer neurons \approx Let say $= 4$ million.

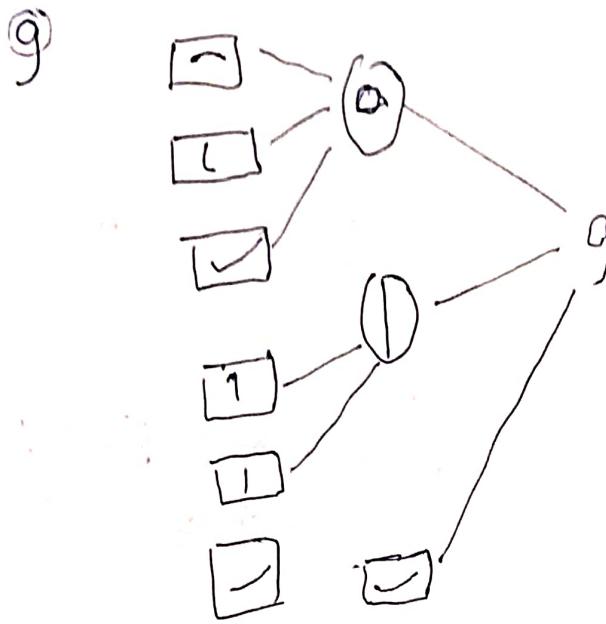
Weights between Input layer and Hidden layer $= 6 \times 4 = 24$ million.

CNN is basically used for feature extraction.

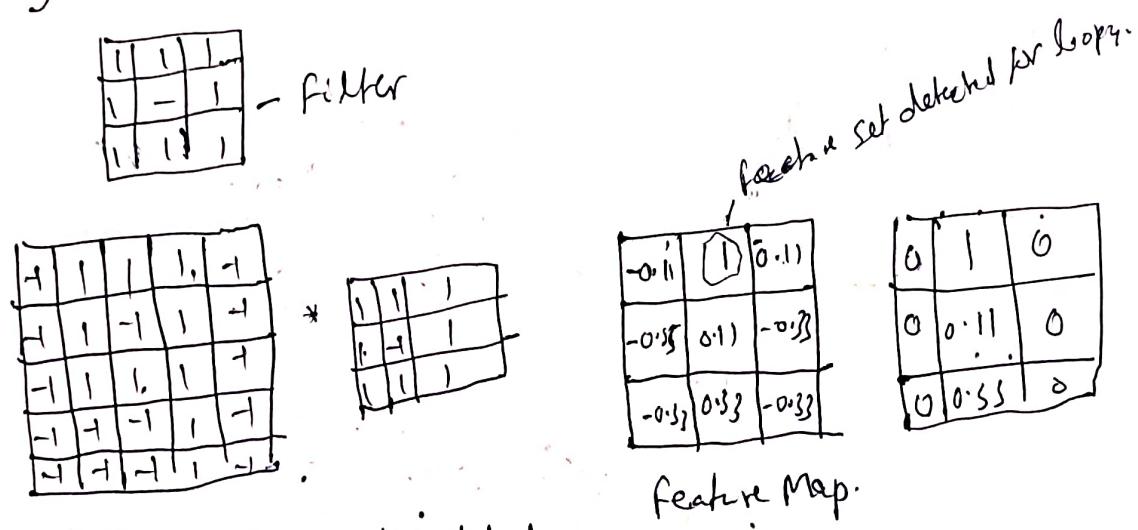
For recognition of cow image features are:

- ⑥ (eyes, ears, nose), (head, legs)



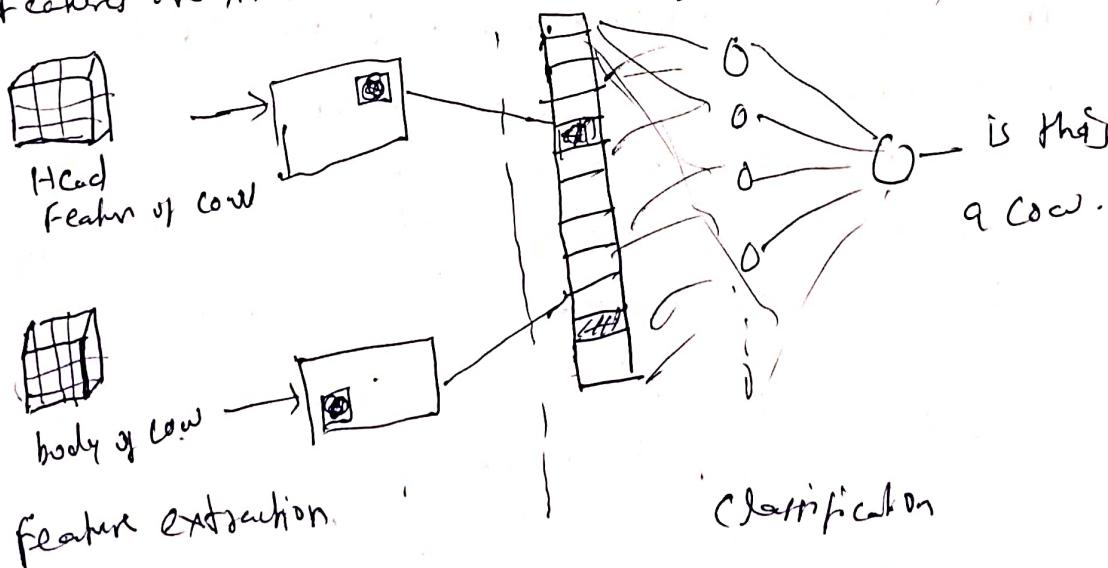


g) we apply filter for loopy image



- ① Filters are the feature detectors
- ② filters are Location invariant.

Features are stacked and creates 2D volume.



Dimensions are minimized using Pooling. After convolution we perform pooling to reduce the dimension.

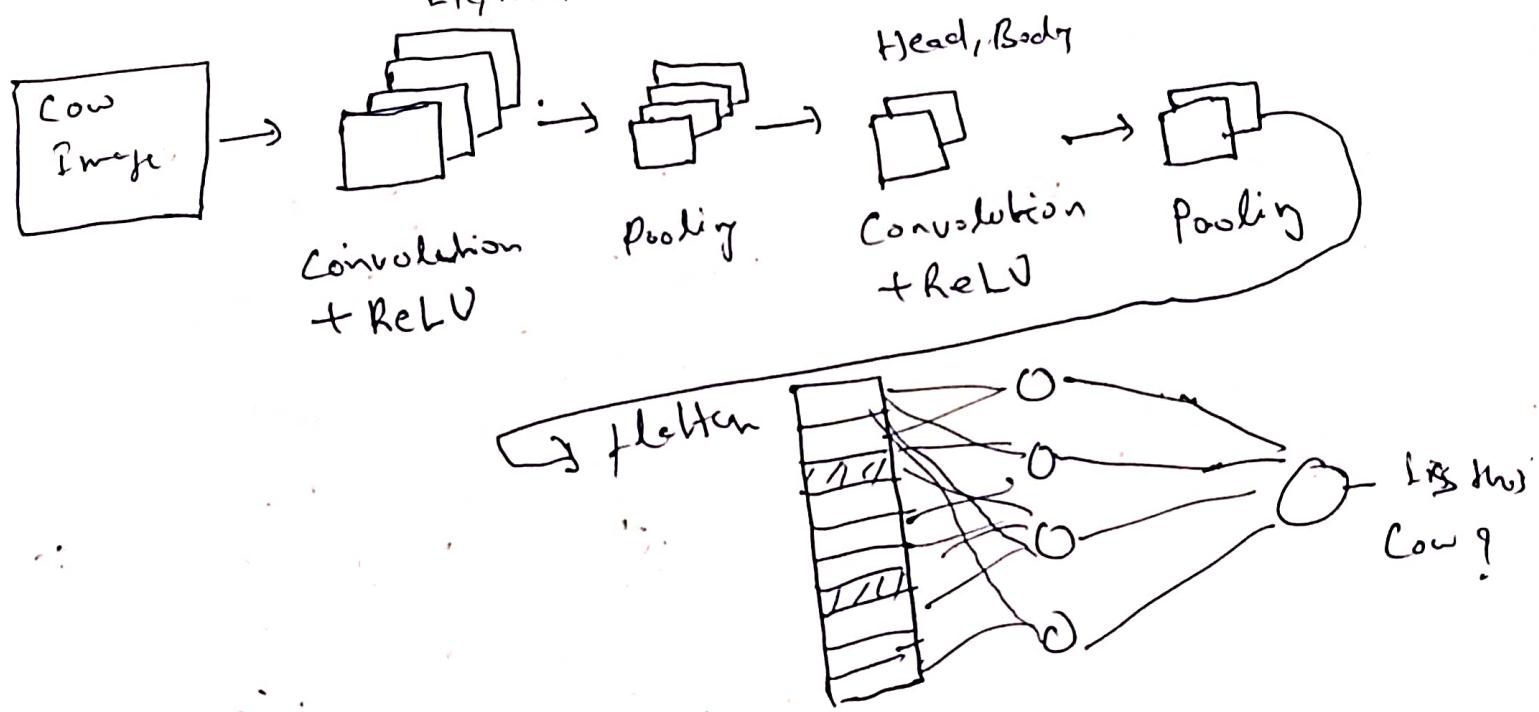
- (i) Max Pooling
- (ii) Average Pooling.

Benefits of Poolings are

- (i) Reduces dimensions & computations
- (ii) Reduces overfitting as there are less parameters
- (iii) Model is tolerant towards variations, distortions.

ReLU → used for non linear activation functions.

Eye, nose, ears etc



Benefits of Convolution

- Connections sparsity reduces overfitting
- Convolution + Pooling gives location invariant feature detection
- Parameter sharing

Benefits of ReLU

- Introduces non linearity
- Speeds up training, faster to compute

Benefits of Pooling

- ① Reduces dimension, overfitting
- ② makes the model tolerant towards small distortion and variations.

CNN can not take care of rotation and scale.

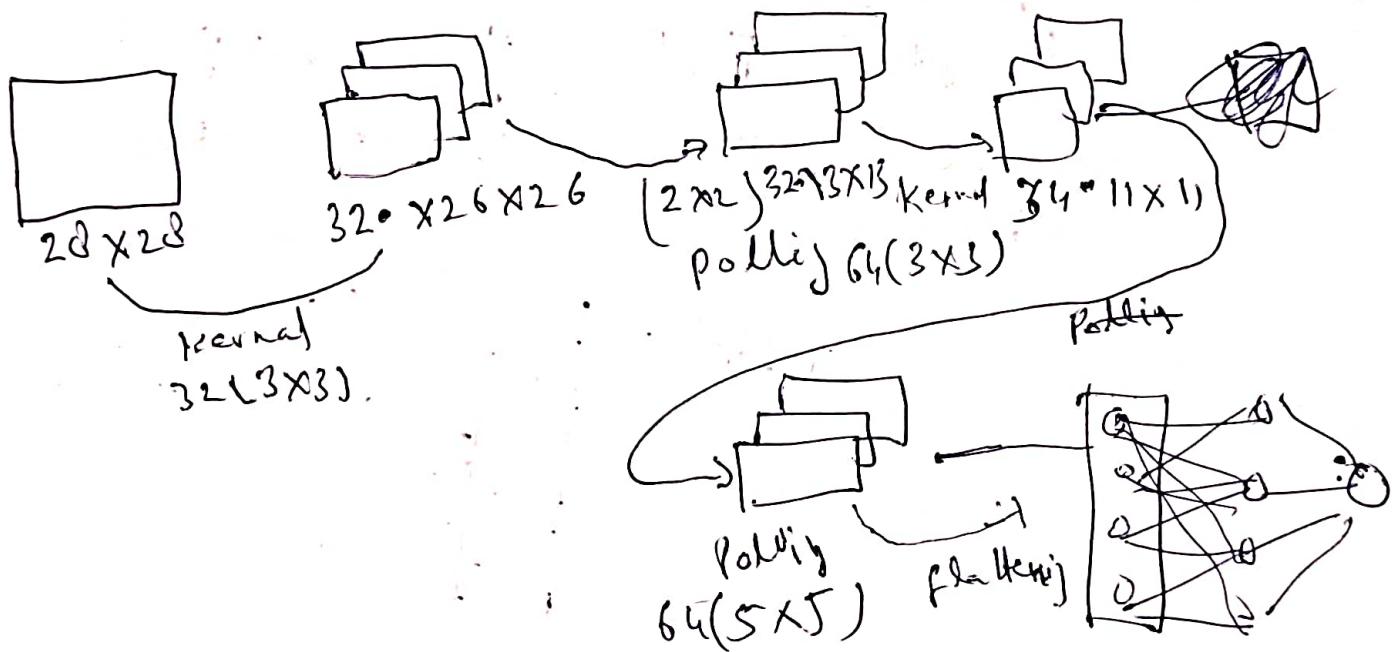
CNN has the following parts

① Convolution + ReLU

② Padding

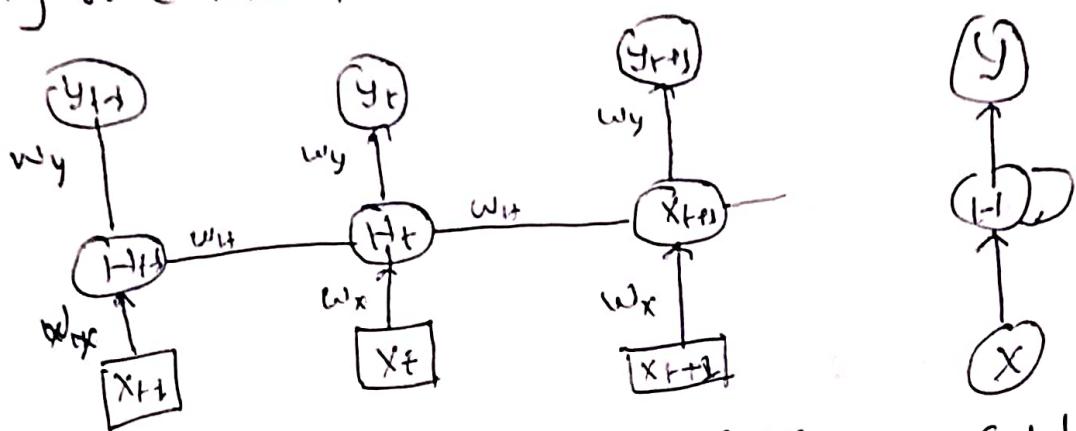
③ Pooling

④ Stride



LSTM

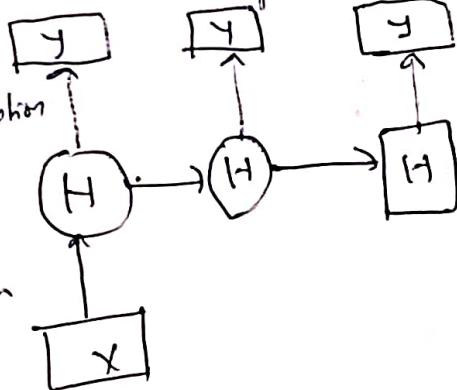
RNN → used for time series analysis which has some memory to connect previous state to next state.



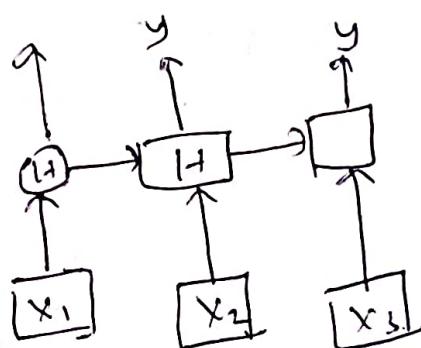
Unfold Representation of RNN

Fold RNN

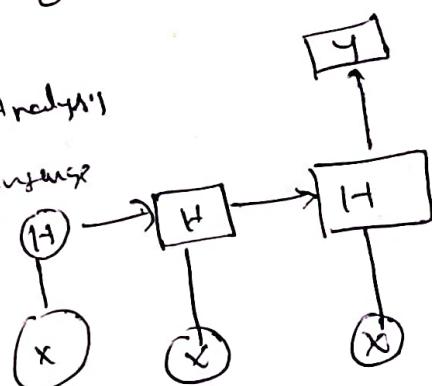
Single RNN
one row



one to many

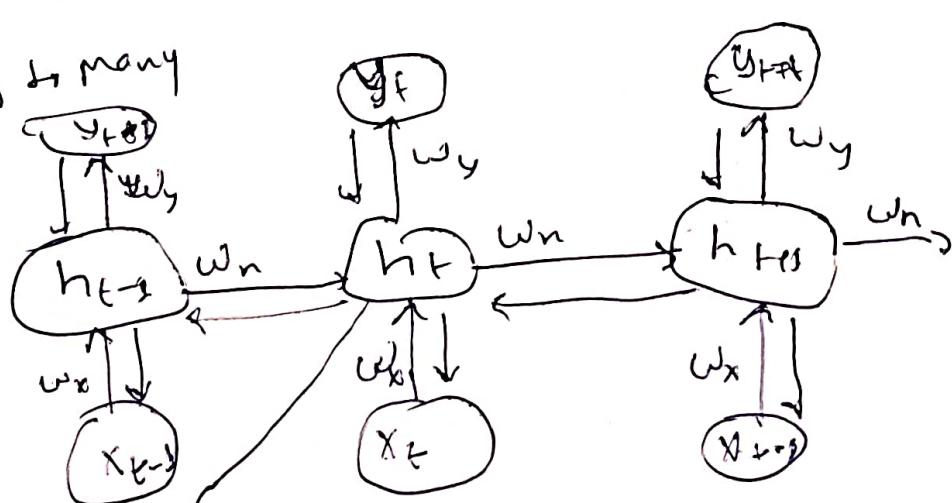


many to many



- use
 (1) Sequential Analysis
 (2) Sentence Length Recognition
 (3) Language Translation
 (4) Encoder and Decoder

many to many



$$z_n^{(t)} = [x_t w_x + h_{t-1} w_n] + b_n$$

$$n_t = \alpha(2^{(t)})$$

$$z_y^{(t)} = [h_t * w_y] + b_y$$

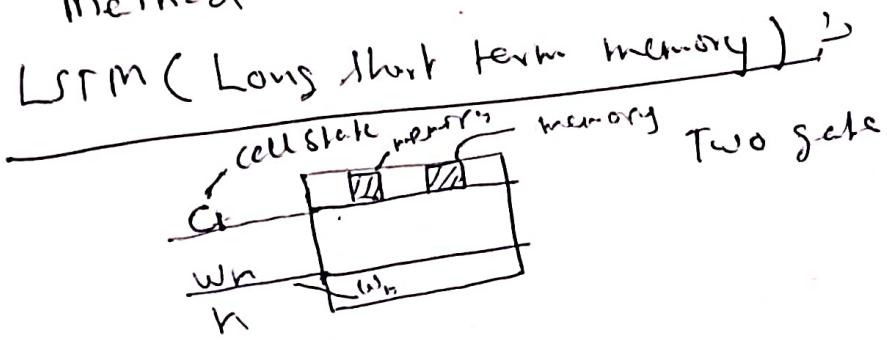
$$y_t = \alpha(z_y^{(t)})$$

$$L = \sum_i L_i$$

$$\Delta w = h \frac{\partial L}{\partial w}$$

$$w = w + \Delta w$$

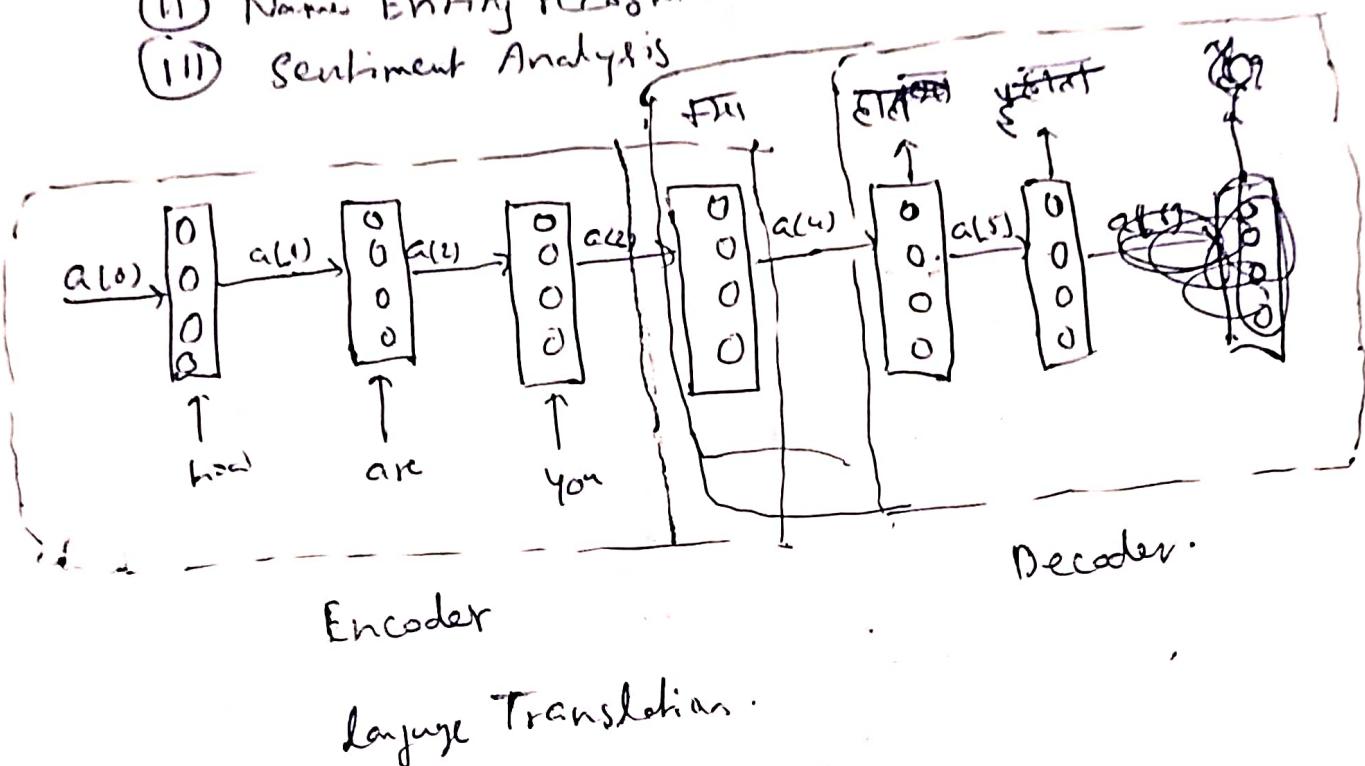
RNN suffers vanishing gradient problem. To shoot out the problem of vanishing gradient we use LSTM method.



one, Memory is to remove unnecessary memory, information
another memory is to add important information or memory.

Recurrent Neural Network :-

- USES
- (I) Translation from one language to other Chuncking of neurons become a problem
 - (II) Name Entity recognition
 - (III) Sentiment Analysis



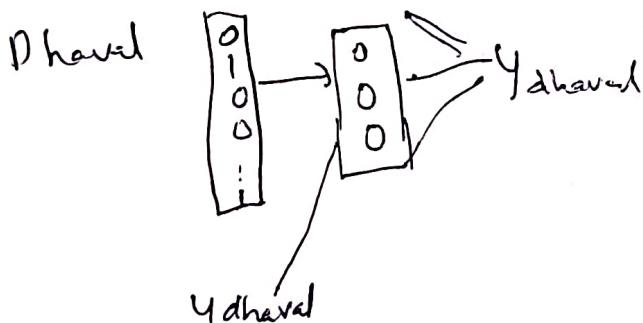
Problem W.M ANN :-

- (I) Variable size of input / output Neurons
- (II) Too much computation

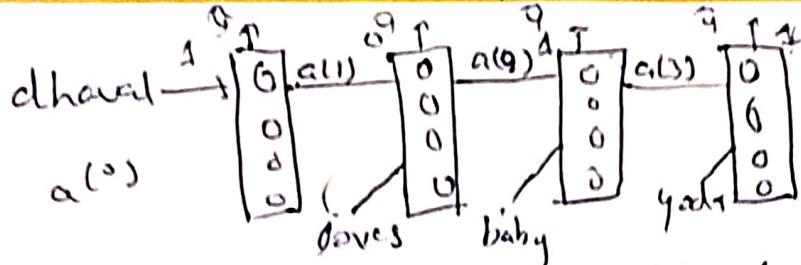
Named Entity Recognition

for example "Dhaval loves baby Yoda"

Entity	1	0	1	1
--------	---	---	---	---



It needs memory

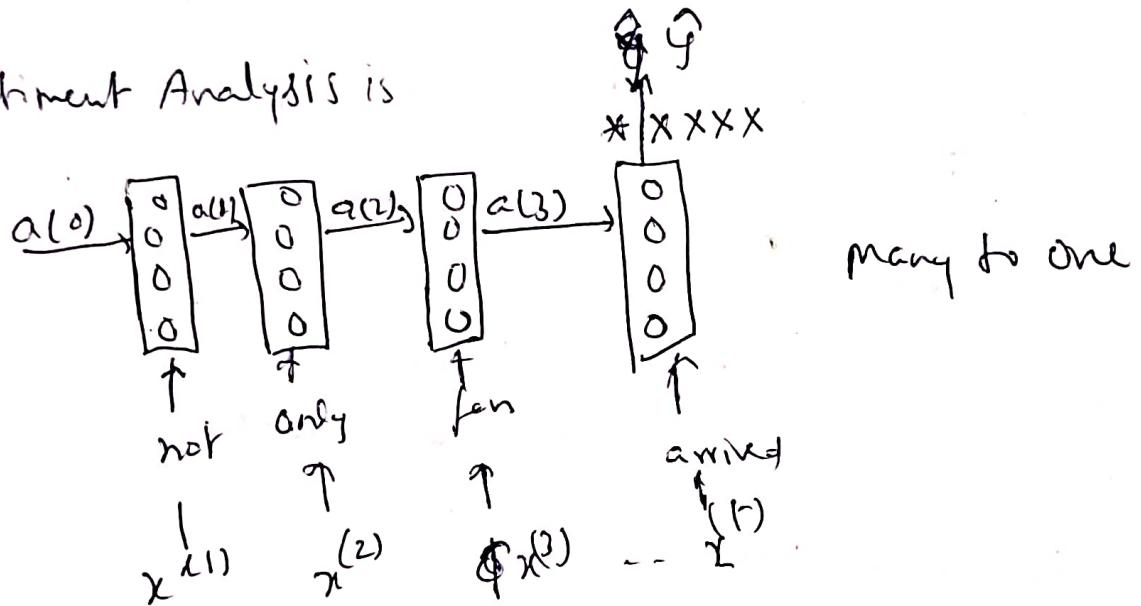


there is only one hidden layers. This is time branch.

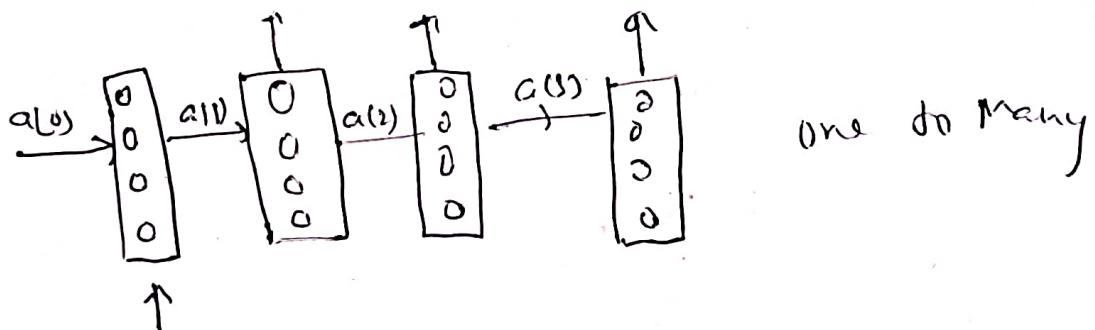


$$\text{Total loss} = \text{Loss 1} + \text{Loss 2} + \text{Loss 3} + \text{Loss 4}$$

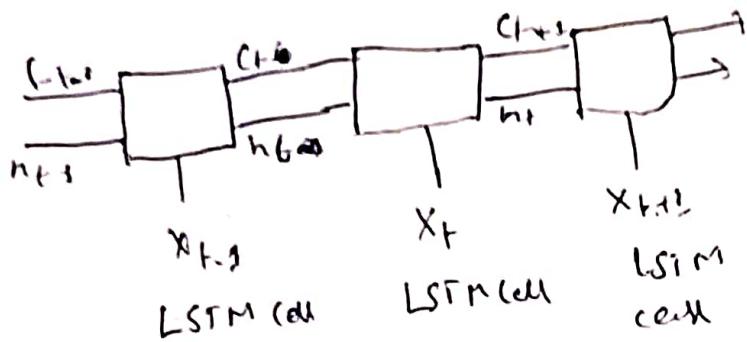
Sentiment Analysis is



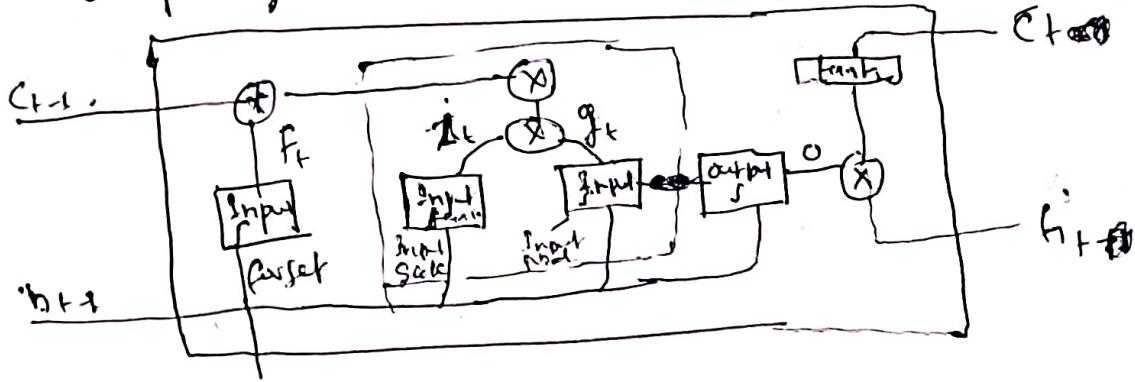
Writing Poem



LSTM \Rightarrow LSTM is used to solve the problem of vanishing gradient.



- ① A Simple RNN cell
- ② cell state Long term memory
- ③ forget gate
- ④ Input gate
- ⑤ Output gate



$$f_t = \text{sigmoid}(w_f h_{t-1} + w_f x_t + b_f)$$

$$C_t^f = C_{t-1} * f_t$$

$$i_t = \text{sigmoid}(w_i h_{t-1} + w_i x_t + b_i)$$

$$g_t = \text{tanh}(w_g h_{t-1} + w_g x_t + b_g)$$

tanh activation function

$$= \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

~~$$C_t^i = i_t * g_t$$~~

$$C_t = C_t^f + C_t^i$$

$$O_t = \sigma((W_{oh} \cdot h_{t-1}) + W_{ox} \cdot x_t) + b$$

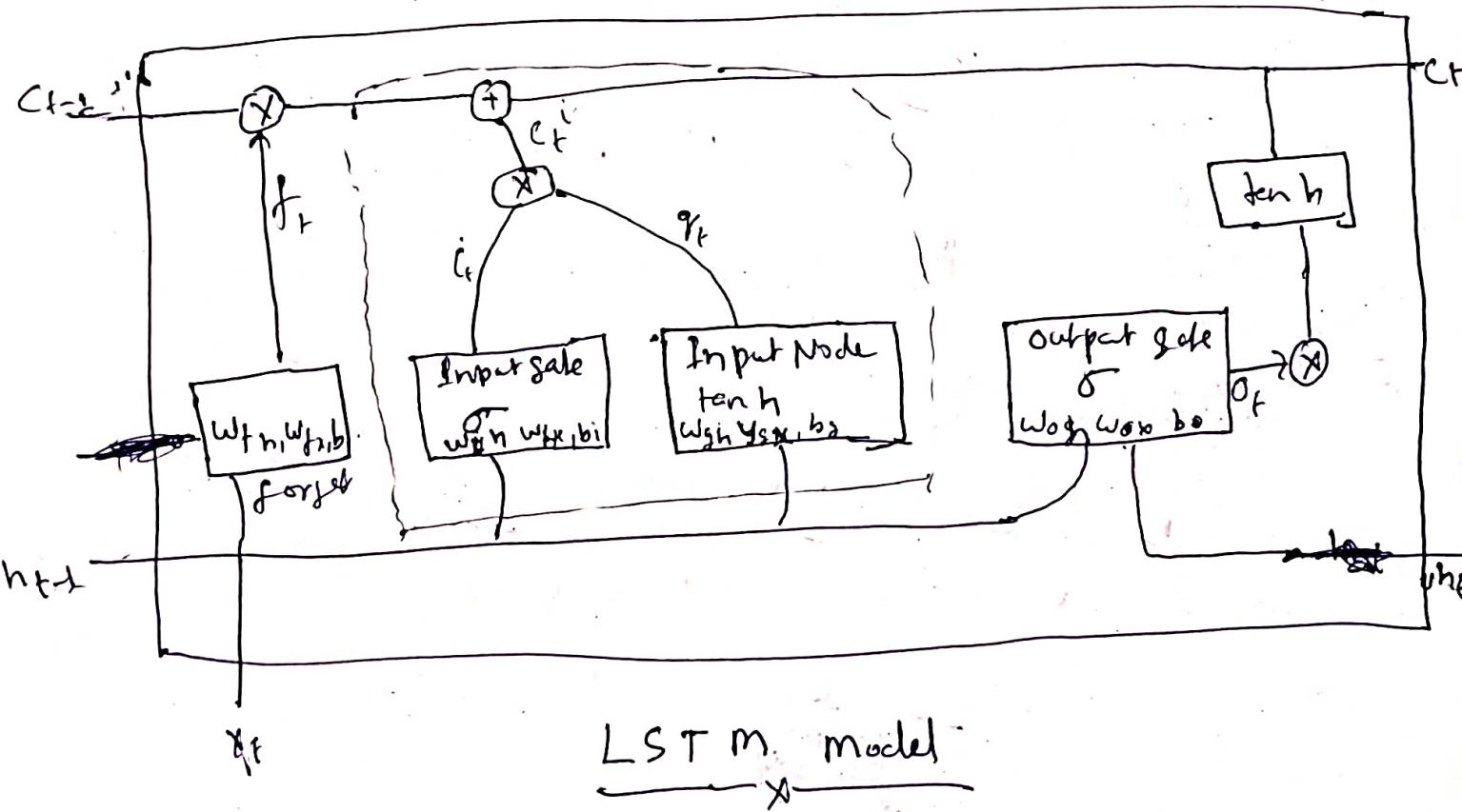
$$h_t = \tanh(C_t \cdot O_t)$$

$$w = \{ w_f \\ w_i \\ w_g \\ w_o \}$$

$$b = \{ b_f \\ b_i \\ b_g \\ b_o \}$$

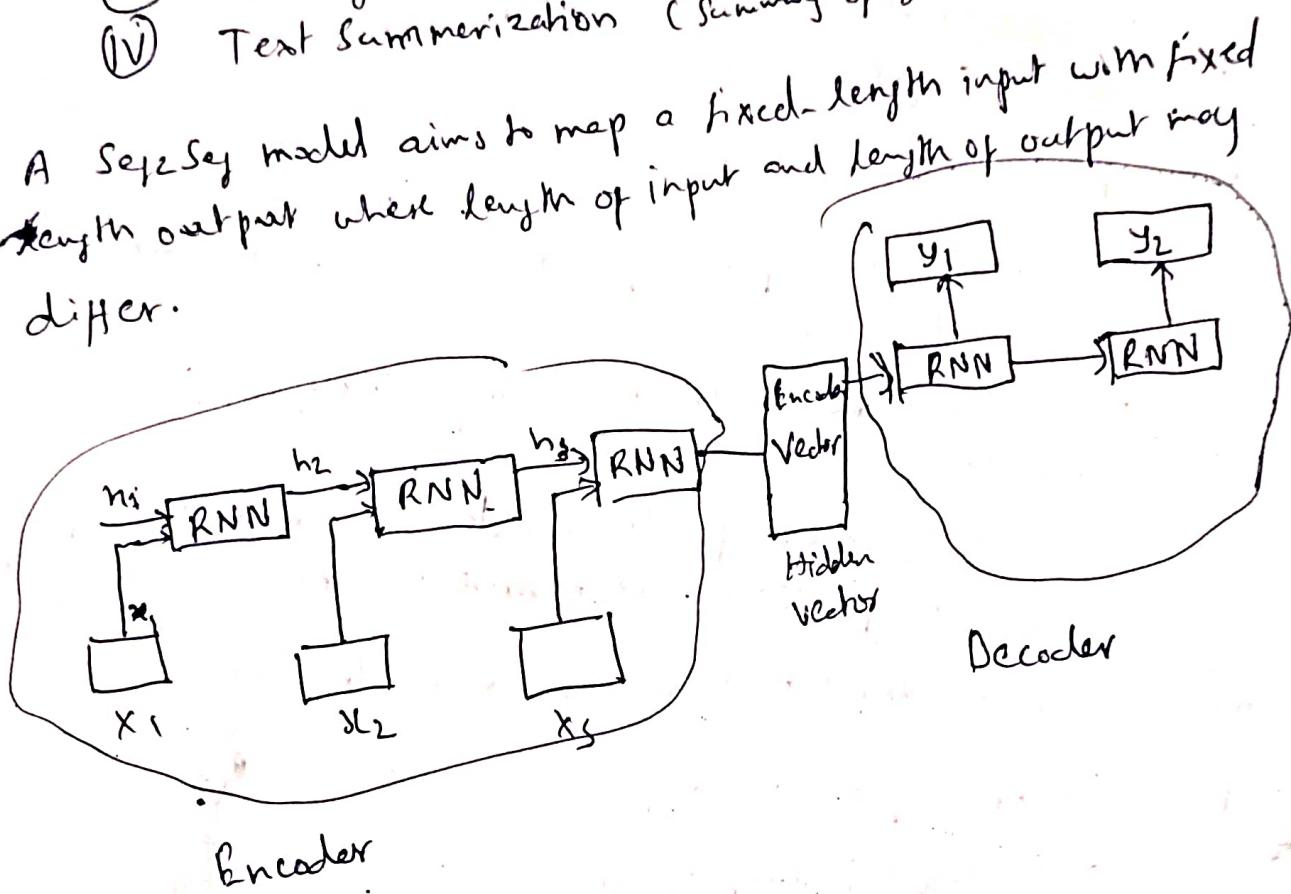
~~$$h_t = \{ w_{fh} \\ w_{ih} \\ w_{gh} \\ w_{oh} \}$$~~

$$g = \{ f_t \\ i_t \\ g_t \\ o_t \}$$



Encoder - Decoder - Sequence-to-sequence Architecture

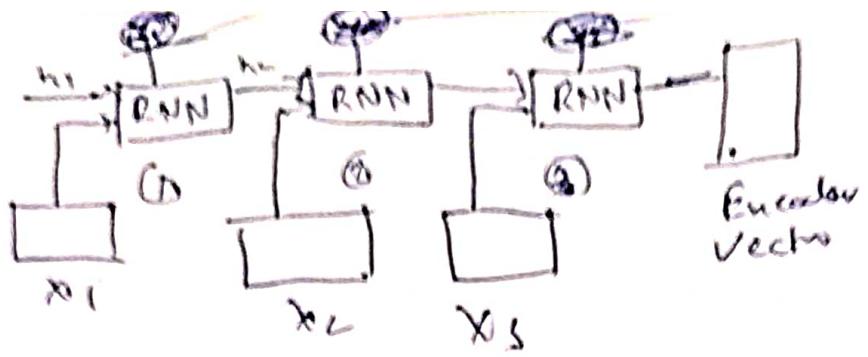
- Sequence-to-sequence (Seq2Seq) models is a special class of recurrent Neural Network architecture
- Seq2Seq models are used to solve complex language problem
 - (I) Machine translation
 - (II) Question Answering
 - (III) Creating chatbots
 - (IV) Text Summarization (Summary of given text)
- A Seq2Seq model aims to map a fixed-length input with fixed length output where length of input and length of output may differ.



- The encoder will convert the input sequence into a single dimensional vector (hidden vector)
- The decoder will convert the hidden vector into the output sequence.

Encoder

- Multiple RNN cells (LSTM cells for better performance) can be stacked together to form the encoder. RNN reads each input sequentially
- For every timestamp (each input) t , the hidden state (hidden vector) h is updated according to the input at that time stamp $x(t)$



Encoder

- After all the inputs are read by encoder model, the final hidden state of the model represents the context/summary of the whole input sequence.

Example:

- Consider the input sequence "I am a student" to be encoded.
- There will be totally 4 time steps (4 tokens) for the encoder model.
- At each time step, the hidden state h will be updated using the previous hidden state and the current input. h_0 can be considered as to be zero or randomly.
- Each layer outputs two things
 - ① Updated hidden states
 - ② Output for each stage.
- The output at each layer is rejected and only hidden states will be propagated to the next layer.

$$h_t = f(W^{hh} \cdot h_{t-1} + W^{hx} \cdot x_t)$$

W^{hh} → weight for recurrent neuron

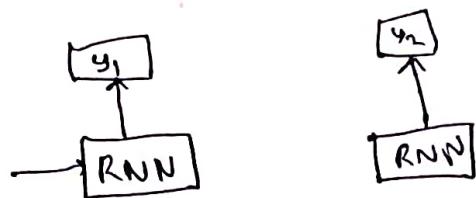
h_{t-1} → Previous state

W^{hx} → input weight neuron

x_t → Current input

- This encoder vector h_t aims to encapsulate the information for all input elements in order to help the decoder make accurate predictions.
- It acts as the initial hidden state of the decoder part of the model.

Decoder



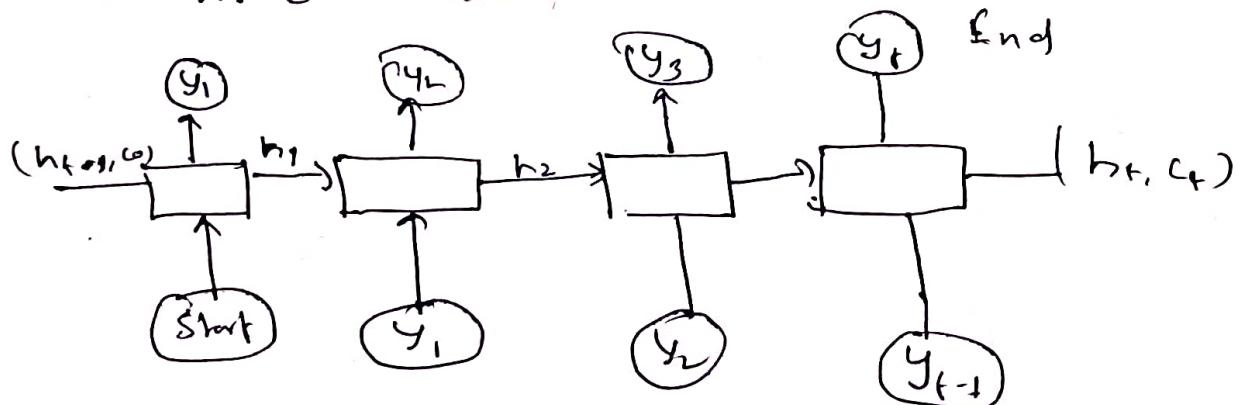
- A stack of several recurrent units where each unit predicts an output y_t at a time step t .
- The input for the decoder is the final hidden vector obtained at the end of encoder model.
- Each Recurrent Unit

Accept:

- a hidden state from the previous unit h_{t-1}
- Previous layer output y_{t-1}

Produces :

- An output as well as
- Its own hidden state



Any hidden state h_t is computed using

$$h_t = f(w^{(h)} h_{t-1})$$

- The output occurred at each time step of decoder is actual output
- The model will predict the output until the end symbol occurs

Output layer

- We use softmax activation function at the output layer
- It is used to produce the probability distribution from a vector of values with the target class of high probability

$$y_t = \text{Softmax}(w^s h_t)$$

- Loss is calculated on the predicted output from each time step and the errors are backpropagated through time in order to update the parameters of the network.
- Training the network over a longer period with sufficiently large amount of data results in pretty good predictions