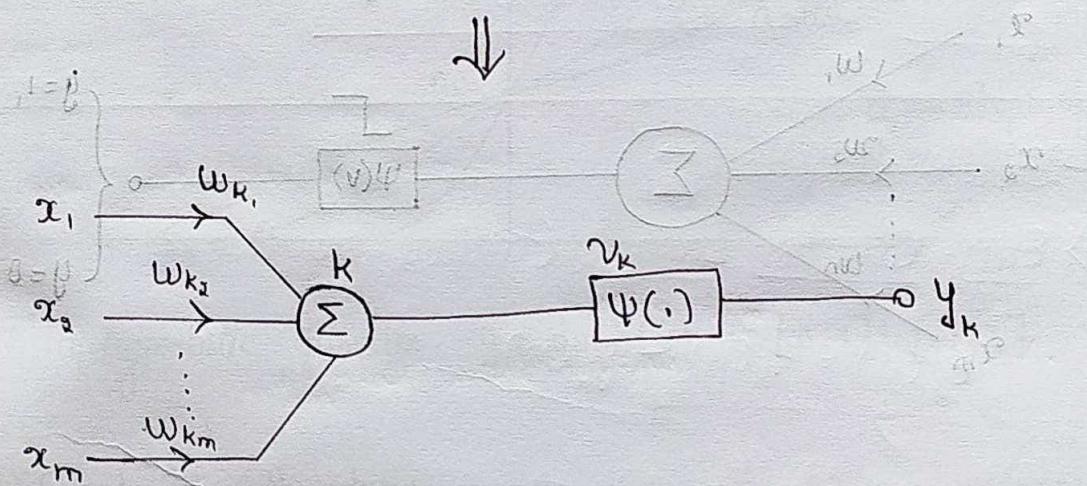
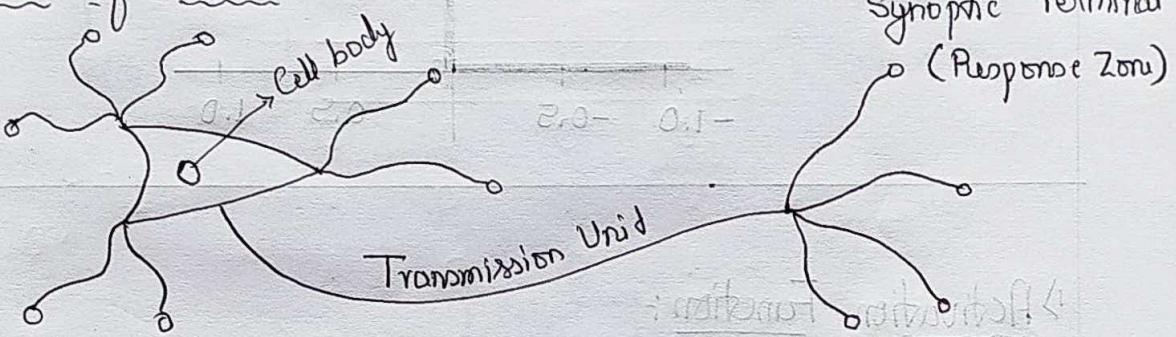


## Artificial Neural Network:

Benefits:

- \* Exploits non-linearity
- \* Input - Output Mapping
- \* Adaptability
- \* Evidential Response
- \* Ability for fault tolerance
- \* Uniformity in origin
- \* VLSI implementability

Model of Neuron:-



$$U_k = \sum_{j=1}^m x_j w_{kj}$$

$$V_k = U_k + b_k$$

$$U_k = \Psi(|A|)$$

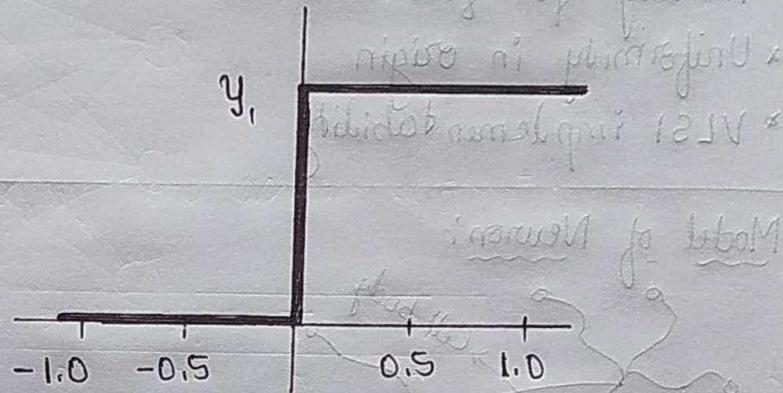
# Threshold Function (McCulloch & Pitts Model)

$$\Psi(\cdot) = \boxed{\text{L}}$$

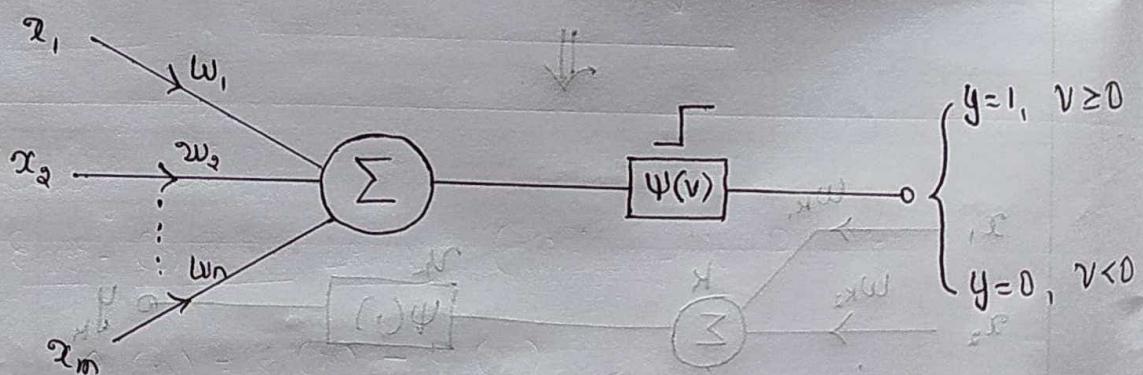
if  $V_K \geq 0, y_K = 1$

$V_K < 0, y_K = 0$

limiting range  
(and range)



## Activation Function:

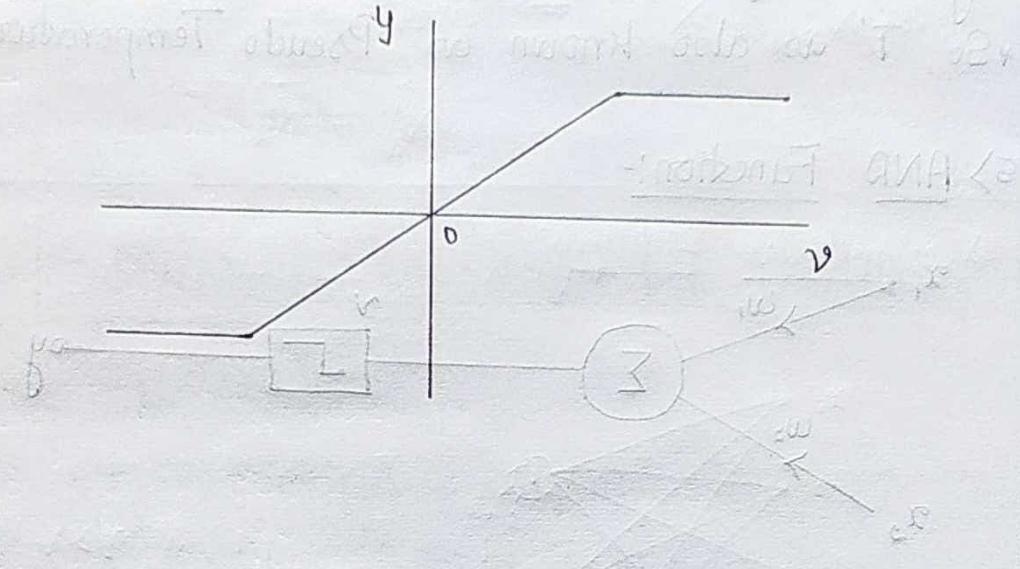


$$w_i x_i \sum_{i=1}^m = u$$

$$d + u = v$$

$$(d + u)\Psi \leq u$$

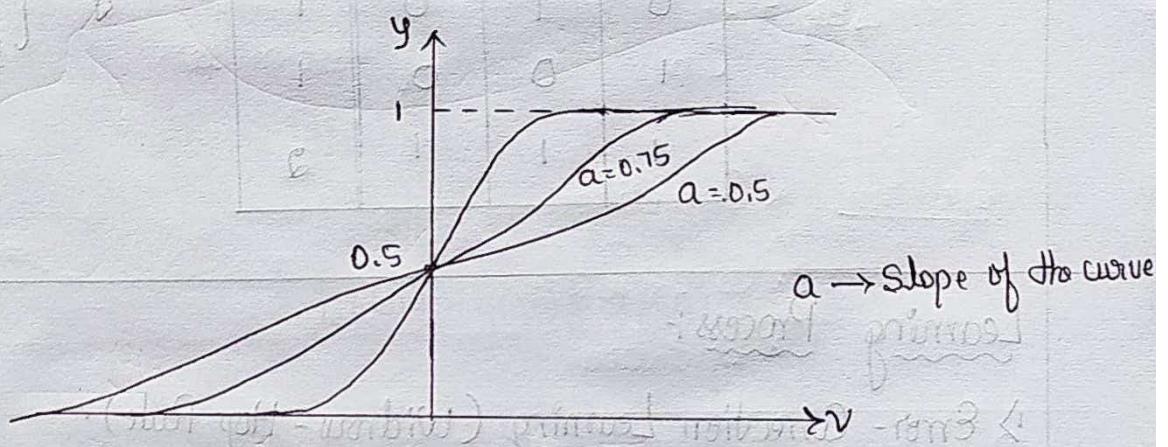
### 3) Linear Model:-



### 3) Sigmoid Function (S-Shaped) :- $\psi(v) = \frac{1}{1 + \exp(-av)}$

$$\psi(v) = \frac{1}{1 + \exp(-av)}$$

$v$	$\psi(v)$
0	0
1	0.5
0	0.5
-1	0



∴ Tanh hyperbolic function,  $y = \tanh(v) \Rightarrow [0, 1], [-1, 1]$

### 4) Stochastic function:

\* Uncertainty is introduced.

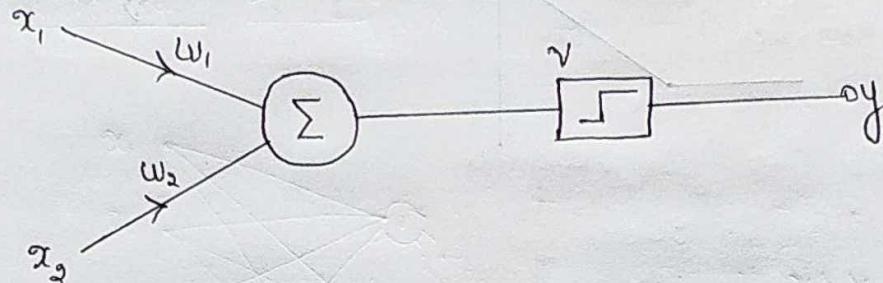
$$\therefore y = \psi(v) = \begin{cases} +1, & \text{with prob. } P(v) \\ -1, & \text{with prob. } 1 - P(v) \end{cases}$$

∴  $P(v) = \frac{1}{1 + \exp(-v/T)}$ ,  $T \rightarrow \text{Constant probability}$

$$\text{if } T=0, P(v)=1$$

- \* By increasing  $T$ , the noise is increases blow noise
  - \* So ' $T$ ' is also known as Pseudo Temperature

## S> AND Function:-



$$V = x_1 w_1 + x_2 w_2 \quad \text{and assume } w_1 = w_2 = 1$$

$x_1$	$x_2$	T	V
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	2

$$y = \begin{cases} 1, & v \geq 2 \\ 0, & v < 2 \end{cases}$$

## Learning Process :-

↳ Error - Correction Learning (Widrow - Hop Rule):

$$e_k(n) = d_k(n) - y_k(n), \text{ mit nachstehender Bedeutung.}$$

where, n: iteration steps

$k$  :  $k^{\text{th}}$  neuron

$t$ : target value

y : Actual value

## Minimum Entropy Function (Cost Function)

$$E(n) = \frac{1}{2} e_k^2(n)$$

$$\therefore \Delta w_i(n) = \eta e_k(n) \cdot x_i(n), \quad \eta \rightarrow \text{Learning Rate}$$

### Delta Rule:

$$w_j(n+1) = w_j(n) + \Delta w_j(n)$$

$$\therefore E^P(n) = \frac{1}{2} \sum_P e^P(n)$$

$$\therefore E^P(n) = \frac{1}{2} \sum [t(n) - y(n)]^2$$

Gradient,

$$\frac{\partial E}{\partial w_i} = \frac{\partial E}{\partial y} \cdot \frac{\partial y}{\partial w}$$

$$= [t(n) - y(n)] x_i(n)$$

$$= -e(n) \cdot x_i(n)$$

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

$$= \underline{\eta e(n) \cdot x_i(n)}$$

↳ known as Gradient Descent.

### Memory Based Learning:-

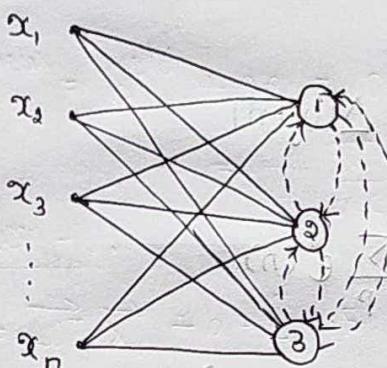
$$\{\vec{x}_i\}_{i=1}^N, \{d_i\}_{i=1}^N$$

$$(x_1, d_1), (x_2, d_2), \dots, (x_n, d_n)$$

$d(\vec{x}_{\text{Test}}, \vec{x}_i) \Rightarrow$  Gives minimum distance.

$$\vec{x}_N = \frac{\vec{x}_i}{d_N} = (n) \vec{z}$$

### 3) Competitive Learning:-



$$+ (a) w = (1-a) w, y = \begin{cases} 1, & \text{if } v > v_j \text{ for all } j \\ 0, & \text{otherwise} \end{cases}$$

$$\therefore \Delta w_j = \begin{cases} \eta(x_j - w_j), & \text{if neuron wins} \\ 0, & \text{if neuron loses} \end{cases}$$

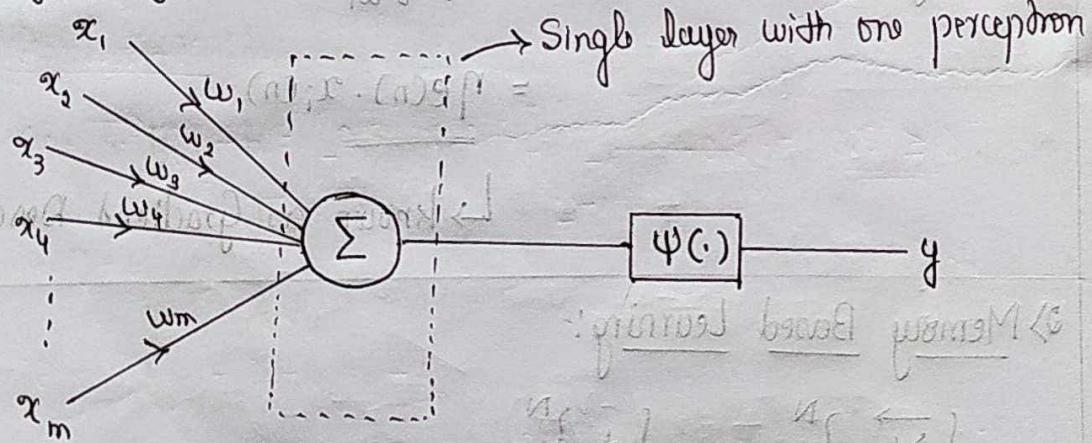
Operations:- \* Competition

\* Co-operation

\* Updation.

### Perceptron-

Single Layer Perceptron



$$\therefore v(n) = \sum_{i=1}^m w_i(n) \cdot x_i(n) = \vec{w}^T \vec{x}_i \text{ (dot product)}$$

- \* Works only when linearity Separable losses.
- \* Also called "Perceptron Convergence Algorithm"

Training data,

$$H_1 \in C_1, H_2 \in C_2, H = H_1 \cup H_2$$

For correct classification,

$$\vec{w}^T \vec{x} > 0, \text{ for every } \vec{x} \in C_1$$

$$\vec{w}^T \vec{x} \leq 0, \text{ for every } \vec{x} \in C_2$$

For wrong classification,

if  $\vec{w}^T \vec{x} \geq 0$  and  $\vec{x} \in C_2$ , so update the weight

$$\text{i.e. } \vec{w}(n+1) = \vec{w}(n) - \eta \vec{x}(n) \rightarrow \text{decrease}$$

if  $\vec{w}^T \vec{x} \leq 0$ , and  $\vec{x} \in C_1$ , then

$$\vec{w}(n+1) = \vec{w}(n) + \eta \vec{x}(n) \rightarrow \text{increase}$$

Let,

$$\vec{x}_1 = \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix}, \vec{x}_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \text{ initially, } \vec{w} = \begin{bmatrix} 0.5 \\ -1 \\ -0.5 \end{bmatrix}$$

$$d_1 = 0 \quad d_2 = 1 \quad b = 0.5$$

and  $y = \begin{cases} 1, & v \geq 0 \\ 0, & \text{otherwise} \end{cases}$  and assume  $\eta = 1$

Soln:-

Epoch-1:

Iteration-1:

$$v = \vec{w}^T \vec{x}_1 + b = [0.5 \ -1 \ -0.5] \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} + 0.5 = 0.5 \geq 0$$

$$\therefore y = 1$$

But  $d_1 = 0$  so mismatch with 'y' so calculate the error and update the weight

$$e = d - y = 0 - 1 = -1$$

$$w_{\text{new}} = w_{\text{old}} + \eta e \vec{x} = \begin{bmatrix} 0.5 \\ -1 \\ -0.5 \end{bmatrix} + (1)(-1) \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \\ 0.5 \end{bmatrix}$$

$$b_{\text{new}} = b_{\text{old}} + \eta e = 0.5 + (1)(-1) = -0.5$$

Iteration - 2:-

$$v = \vec{w}^T \vec{x}_2 + b = [-0.5 \ 0 \ 0.5] \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} + (-0.5) = -1.5 < 0$$

Follows the rule as  $v < 0 \Rightarrow y = 0$ .

But  $d_2 = 1 \neq 0$  is mismatch so update weight.

$$e = d - y = 1 - 0 = 1$$

$$w_{\text{new}} = w_{\text{old}} + \eta e \vec{x} = \begin{bmatrix} -0.5 \\ 0 \\ 0.5 \end{bmatrix} + (1)(1) \begin{bmatrix} 1 \\ 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1 \\ -0.5 \end{bmatrix}$$

$$b_{\text{new}} = b_{\text{old}} + \eta e = -0.5 + (1)(1) = 0.5$$

Epoch - 2:-

Iteration - 1:-

$$v = \vec{w}^T \vec{x}_1 + b = [0.5 \ 1 \ -0.5] \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} + 0.5 = 0.5 \geq 0$$

$\therefore y = 1, d_1 = 0$ .

$$e = d - y = 0 - 1 = -1$$

$$w_{\text{new}} = w_{\text{old}} + \eta e \vec{x} = \begin{bmatrix} 0.5 \\ 1 \\ -0.5 \end{bmatrix} + (1)(-1) \begin{bmatrix} 1 \\ -1 \\ -1 \end{bmatrix} = \begin{bmatrix} -0.5 \\ 0 \\ 0.5 \end{bmatrix}$$

$$b^{\text{new}} = 0.5 + (1)(-1) = -0.5$$

Iteration-2:

$$V = \vec{w}^T \vec{x} + b = [-0.5 \ 0.5] \begin{bmatrix} 1 \\ 1 \end{bmatrix} + (-0.5) = 0.5 \geq 0$$

$$\therefore y = 1, d_s = 1$$

No change.

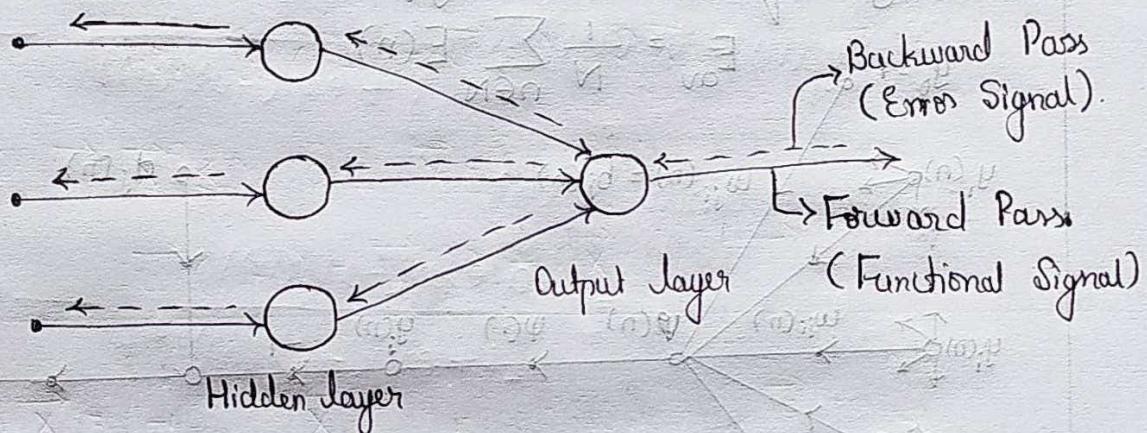
## Multi-layer Perceptron:

\* Multi-layer feed forward networks

→ Input layer (source node)

→ One or more hidden layer (computation node)

→ Output layer



Let,

\* Indices i, j, k refers to different neuron in network.

\* n is the iterations

\*  $\sum(n)$  refers to the instantaneous sum of error square at iteration n.

\*  $e_j(n)$  error signal at the output of neuron j at iteration n.

\*  $d_j(n)$  desired response for neuron j

\*  $y_i(n)$  function signal at the output neuron i

\*  $w_{ij}(n)$  synaptic weight connecting output of neuron i to j

\*  $v_j(n)$  activation potential neuron j

\*  $\psi_i(\cdot)$  activation function of neuron i

## Back-Propagation Algorithm:

Error signal at the output of neuron 'j' at iteration 'n' is,

$$e_j(n) = d_j(n) - y_j(n) \quad \text{--- (1)}$$

Instantaneous value of the error energy for neuron 'j'

$$\gamma_2 e_j^2(n)$$

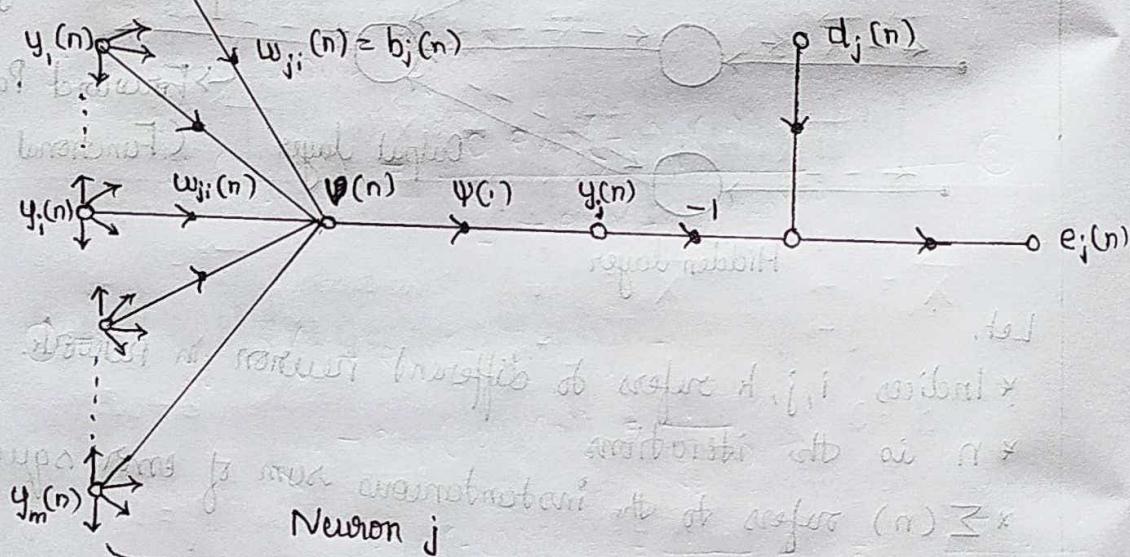
For total instantaneous error energy, sum over all output neurons,

$$E(n) = \frac{1}{q} \sum_{j \in C} e_j^2(n) \quad \text{--- (2)}$$

where 'C' is the all output neuron.

Let 'N' be the no. of examples in the training set. Then average squared error energy is,

$$E_{av} = \frac{1}{N} \sum_{n \in N} E(n)$$



$$\therefore v_j(n) = \sum_{i=0}^m w_{ji}(n) \psi_i(n) \quad \text{--- (3)}$$

$$y_j(n) = \psi_j(v_j(n)) \quad \text{--- (4)}$$

Using chain rule calculus,

$$\frac{\partial E(n)}{\partial w_{ji}(n)} = \frac{\partial E(n)}{\partial e_j(n)} \frac{\partial e_j(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial v_j(n)} \frac{\partial v_j(n)}{\partial w_{ji}(n)}$$

Differentiate eqn ① to ④ and put them in above eqn.

\* Diff' ② w.r.t.  $e_j(n)$

$$\frac{\partial E(n)}{\partial e_j(n)} = e_j(n)$$

\* Diff'n ③ w.r.t.  $y_j(n)$

$$\frac{\partial e_j(n)}{\partial y_j(n)} = -1$$

\* Diff'n ④ w.r.t.  $v_j(n)$

$$\frac{\partial y_j(n)}{\partial v_j(n)} = \psi'_j(v_j(n))$$

Diff'n ⑤ w.r.t.  $w_{ji}(n)$

$$\frac{\partial v_j(n)}{\partial w_{ji}(n)} = y_i(n)$$

Substitute these into the chain rule.

$$\textcircled{c} \quad \frac{\partial E(n)}{\partial w_{ji}(n)} = (-e_j(n) \psi'_j(v_j(n))) y_i(n) = ((\alpha)_j^T \psi) \cdot \psi$$

$\therefore$  The correction  $\Delta w_{ji}(n)$  applied to  $w_{ji}(n)$  is,

$$\Delta w_{ji}(n) = -\eta \frac{\partial E(n)}{\partial w_{ji}(n)} = -\eta \frac{(-e_j(n) \psi'_j(v_j(n))) y_i(n)}{((\alpha)_j^T \psi) \cdot \psi} = \eta \delta_j(n) y_i(n)$$

where  $\delta_j(n) = e_j(n) \psi'_j(v_j(n))$  is local gradient.

\* There are 2 possibilities

Case-1:- Neuron j is the output neuron  $\rightarrow$  direct solution.

Case-2:- Neuron j is the hidden neuron

\* When  $y_j$  is the hidden neuron we cannot get the desired result.  
We can redefine  $\delta_j(n)$  for hidden neuron.

BP Networks:-

$$\Delta w_{ji}(n) = \eta \delta_j(n) \cdot y_i(n)$$

↓

Local Gradient

Output Neuron

 $e_j(n) \psi'(v_j(n))$

Hidden Neuron

 $\psi(v_j(n)) \sum_k b_k(n) w_{kj}(n)$

i) Logistic Activation Function :-

$$y_j(n) = \psi_j(v_j(n)) = \frac{1}{1 + \exp(-\alpha v_j(n))} \quad \textcircled{1}$$

differentiate w.r.t  $v_j(n)$

$$\psi_j(v_j(n)) = \frac{\alpha \exp(-\alpha v_j(n))}{[1 + \exp(-\alpha v_j(n))]^2} \quad \textcircled{2}$$

$$\begin{aligned} \textcircled{1} \Rightarrow y_j(n)(1 + \exp(-\alpha v_j(n))) &= 1 \\ 1 + \exp(-\alpha v_j(n)) &= \frac{1}{y_j(n)} \\ \exp(-\alpha v_j(n)) &= \frac{1}{y_j(n)} - 1 = \frac{1 - y_j(n)}{y_j(n)} \end{aligned}$$

therefore local  $\alpha \psi'(v_j(n)) = \alpha y_j(n)(1 - y_j(n))$

$\therefore \textcircled{2}$  becomes,

$$\begin{aligned} \psi_j(v_j(n)) &= \alpha \frac{1 - y_j(n)}{y_j(n)} \cdot y_j(n) \\ &= \alpha y_j(n)[1 - y_j(n)] \quad \textcircled{3} \end{aligned}$$

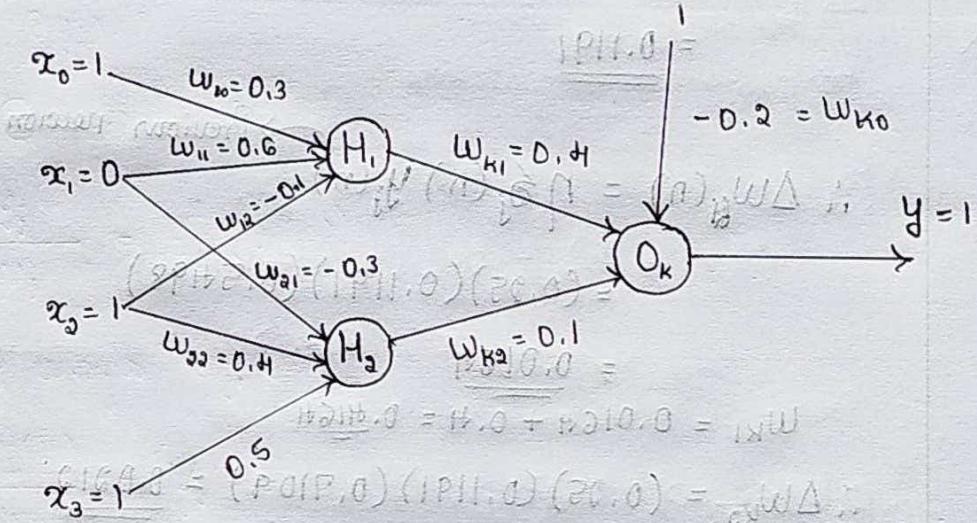
$\therefore \delta_j(n)$  for Output neuron is,

$$\delta_j(n) = e_j(n) \cdot a y_j(n) [1 - y_j(n)]$$

$\therefore \delta_j(n)$  for hidden neuron is,

$$\delta_j(n) = a y_j(n) [1 - y_j(n)] \sum_k \delta_k(n) w_{kj}(n)$$

Given,



and input pattern  $[0, 1]$ , target o/p = 10,  $\eta = 0.05$ ,  $a = 1$ .

Soln:-

For Hidden Layer

$$H_1, v_1 = (1)(0.3) + (0)(0.6) + (1)(-0.1) = \underline{0.2}$$

$$\Psi_1(n) = \frac{1}{1 + e(-v)} = \frac{1}{1 + e(-0.2)} = \underline{0.5498}$$

$$H_2, v_2 = (0)(-0.3) + (1)(0.4) + (1)(0.5) = \underline{0.9}$$

$$\Psi_2(n) = \frac{1}{1 + e(-0.9)} = \underline{0.7109}$$

Output k,

$$v_k = (0.5498)(0.4) + (0.7109)(0.1) = 0.2 = \underline{0.09101}$$

$$\therefore y_k = \frac{1}{1 + e(-0.09101)} = \underline{0.5227}$$

Here we need  $\delta_{\text{fp}} = y - 1$ , but we got  $y_k = 0.5507$ . So we need to update the weight at output neuron. So error,

$$e(n) = y - y_k < 1 - \underline{0.5927}$$

$$\begin{aligned}f_j(n) &= e_j(n) \alpha y_j(n) [1 - y_j(n)] \\&= [1 - 0.5997] (1) (0.5997) [1 - 0.5997] \\&= 0.1191\end{aligned}$$

$$\therefore \Delta w_{ij}(n) = \eta s_j(n) y_i(n)$$

$$= (0.05)(0.1191)(0.5498)$$

$$= 0.0164$$

$$\therefore \Delta w_{k_2} = (0.95)(0.1191)(0.7109) = \underline{\underline{0.0312}}$$

$$\therefore \Delta W_{K0} = (0.95)(0.1191)(*) = \underline{0.1098}$$

$$\therefore W_{k0} = 0.0398 + (-0.3) = \underline{\underline{-0.170}} \text{ joule}$$

$$\frac{27 \pi/2,0}{(6,0 \cdot 2,0 + 6,0 \cdot 2,0)} = \frac{1}{2} = \frac{\pm}{\pm} = f(x)W$$

$$1093.0 = 5.0 + (-1.0)(700P.O) + (4.0)(700) \quad \text{Satisfied}$$

Industries bring up oil wells so we have to move because

$$\left. \begin{array}{l} 0 < x \leq \frac{\mu}{\delta}, \quad x \\ \text{minimum}, \quad 0 \end{array} \right\} = (\bar{x})_{\text{upper}}$$

BEFLU → Regelfläche Flüssig Metall

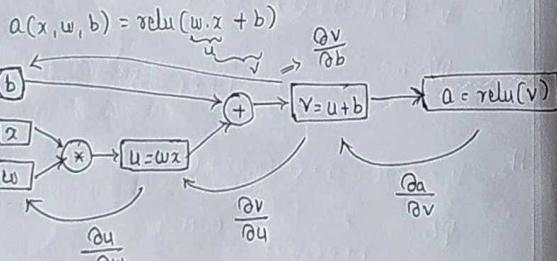
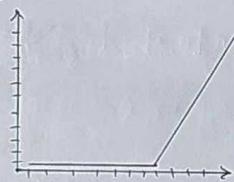
### Computation Graphs:-

- \* A type of graph can be used to represent mathematical expression
- \* In context of deep learning (PyTorch) we can think of neural networks as computational graph.
- \* Two types of calculation,
  - ↪ Forward computation
  - ↪ Backward computation
- \* Suppose we have a activation fn,

$$a(x, w, b) = \text{relu}(w \cdot x + b)$$

$$\text{relu}(x) = \begin{cases} x, & \text{if } x > 0 \\ 0, & \text{otherwise} \end{cases}$$

• ReLU → Rectified Linear Unit

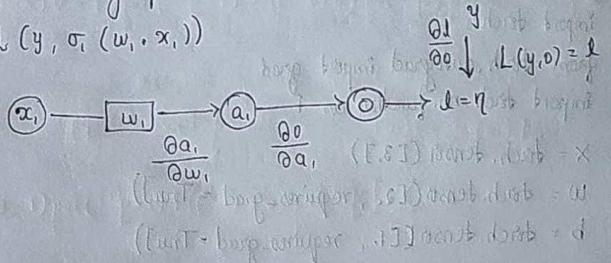


$$\therefore \frac{\partial a}{\partial b} = \frac{\partial v}{\partial b} \cdot \frac{\partial a}{\partial v}$$

$$\therefore \frac{\partial a}{\partial w} = \frac{\partial u}{\partial w} \cdot \frac{\partial a}{\partial u} = \frac{\partial u}{\partial w} \cdot \frac{\partial v}{\partial u} \cdot \frac{\partial a}{\partial v}$$

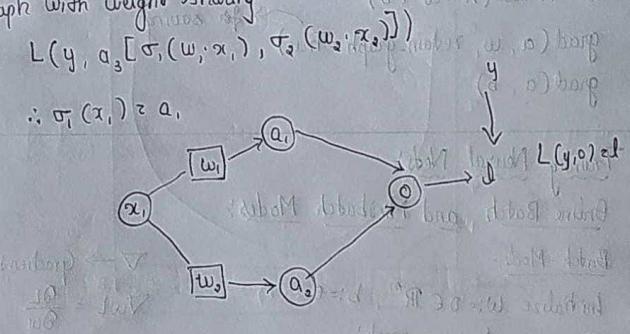
### Graph with single path,

$$L(y, \sigma_i(w_i \cdot x_i))$$

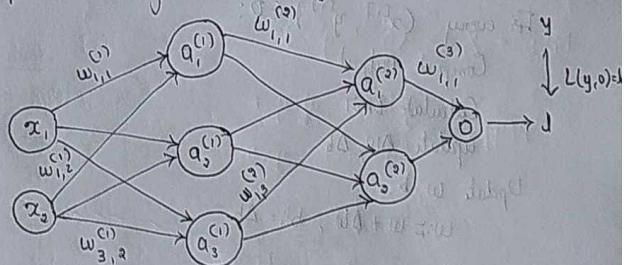


### Graph with weight sharing

$$L(y, a_3[\sigma_i(w_i \cdot x_i), \sigma_2(w_2 \cdot x_2)])$$



### Graph with Fully Connected Layer



$$\therefore \frac{\partial J}{\partial w_{1,1}} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial a_1^{(2)}} \cdot \frac{\partial a_1^{(2)}}{\partial a_1^{(1)}} \cdot \frac{\partial a_1^{(1)}}{\partial w_{1,1}}$$

$$+ \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial a_2^{(2)}} \cdot \frac{\partial a_2^{(2)}}{\partial a_2^{(1)}} \cdot \frac{\partial a_2^{(1)}}{\partial w_{1,1}}$$

```

import torch
from torch.autograd import grad
import torch.nn.functional as F

```

$x = \text{torch.} \text{.} \text{dense}([3,])$

$w = \text{torch.} \text{.} \text{dense}([2,], \text{requires\_grad}=\text{True})$

$b = \text{torch.} \text{.} \text{dense}([1,], \text{requires\_grad}=\text{True})$

$a = F. \text{relu}(x * w + b)$

grad(a, w, retain\_graph=True)  $\rightarrow$  for saving

grad(a, b)

### Training Neural Nets:

#### Online Batch, and Minibatch Modes:

Batch-Model:

Initialize  $w=0 \in \mathbb{R}^m$ ,  $b=0$

For every training epoch:

Initialize  $\Delta w=0$ ,  $\Delta b=0$

For every  $(x^{[i]}, y^{[i]}) \in D$ :

Compute Output

Calculate error

Update  $\Delta w$ ,  $\Delta b$

Update  $w, b$

$w := w + \Delta w$ ,  $b := \Delta b$

$\nabla \rightarrow \text{Gradient}$   
 $\nabla w L = \frac{\partial L}{\partial w}$

For every  $\{(x^{[i]}, y^{[i]}), \dots, (x^{[i+k]}, y^{[i+k]})\} \subset D$

Compute Output values off based with training set

Compute error  $\text{error} = (w)_{\text{true}} - (w)_{\text{cal}}$

Update  $\Delta w$ ,  $\Delta b$

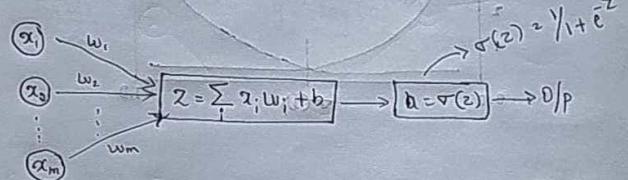
Update  $w, b: (w, b) + ((x)_{\text{true}} - (x)_{\text{cal}})^T [w]$

$w := w + \Delta w$ ,  $b := \Delta b$

### Linear Regression:

Activation function,  $\sigma(x) = x$

Output is real number,  $y \in \mathbb{R}$



### Likelihood Loss:

$$L(w) = P(y|x; w)$$

$$= \prod_{i=1}^n (y^{[i]} | x^{[i]}; w) \quad i=1, \dots, n$$

$$= \prod_{i=1}^n (\sigma(z^{[i]}))^{y^{[i]}} (1 - \sigma(z^{[i]}))^{1-y^{[i]}}$$

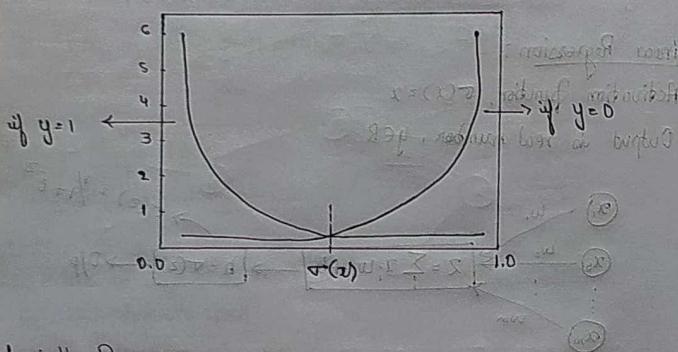
\* In practice it is easier to maximize the log of this equation to called log-likelihood loss.

i.e.  $J(w) = \log(L(w))$

$$= \sum_{i=1}^n [y^{[i]} \log(\sigma(x^{[i]})) + (1-y^{[i]}) \log(1-\sigma(x^{[i]}))]$$

To minimize this regard the equation, i.e  
 $L(w) = -L(w) \rightarrow \text{Binary Cross Entropy}$

$$= -\sum_{i=1}^n [y^{[i]} \log(\sigma(x^{[i]})) + (1-y^{[i]}) \log(1-\sigma(x^{[i]}))]$$



Logistic Regression:-

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (a^{[i]} - y^{[i]})^2$$

$$P(y|x) = \begin{cases} h(x), & \text{if } y=1 \\ 1-h(x), & \text{if } y=0 \end{cases} \quad (\omega : x | \omega) \neq (\omega)$$

$$P(y|x) = \frac{y}{1+y} \quad (\text{Rewrite})$$

$$P(y|x) = \frac{y}{1-y} \quad ((1-y)^{(1-y)})$$

Learning Rate:

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \frac{da}{dz} \frac{\partial z}{\partial w_i}$$

$$\therefore \frac{\partial L}{\partial a} = \frac{a-y}{a(1-a)}, \quad \frac{da}{dz} = \frac{e^{-z}}{(1+e^{-z})^2} = a(1-a), \quad \frac{\partial z}{\partial w_i} = x_i$$

Multi-Class Cross Entropy for k-different class labels

$$H_k(y) = \sum_{i=1}^n \sum_{k=1}^K -y_k^{[i]} \log(a_k^{[i]})$$

Softmax / Multi-class Classification

$$P(y=t | z_t^{[i]}) = \tau_{\text{softmax}}(z_t^{[i]}) =$$

Loss Function

$$L(w) = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^K -y_j^{[i]} \log(a_j^{[i]})$$

### Better ways to train Deep Neural Network:-

- ↳ Learning rate
- ↳ Activation function
- ↳ Weight updation
- ↳ Normalization
- ↳ Regularization

### Optimization:-

#### Batch Normalization:-

Step-1: Normalize the net input

$$\mu_j = \frac{1}{n} \sum_i z_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{n} \sum_i (z_j^{(i)} - \mu_j)^2$$

$$z_j^{(i)} = \frac{z_j^{(i)} - \mu_j}{\sigma_j}$$

$$z_j^{(i)} = \frac{z_j^{(i)} - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}$$

Step-2: Pre-Activation scaling

$$z_j^{(i)} = \frac{z_j^{(i)} - \mu_j}{\sigma_j}$$

$$\sigma_j^{(i)} = \gamma_j \cdot z_j^{(i)} + \beta_j \quad \begin{matrix} \rightarrow \text{Controls the mean} \\ \rightarrow \text{Controls the spread & scale.} \end{matrix}$$

#### Weight Initialization:- (Xavier Initialization)

Step-1: Initialize weight from Gaussian or uniform distribution

Step-2: Scale the weight proportional to the number of inputs to the layer.

Rationale:  $\rightarrow$  Variance of the sample linearly increases as the sample size increases.

$$w^{(d)} := w^{(d)} \cdot \sqrt{\frac{1}{m^{(d-1)}}} \quad \text{where } m \text{ is the no. of input units}$$

at previous layer. Now, w is the initial weight for the first layer.

### Optimizers:-

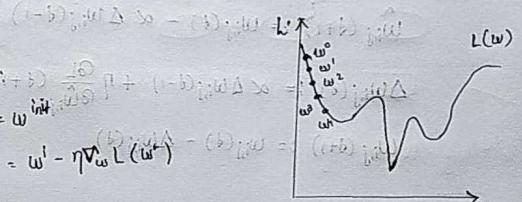
#### Challenges:-

Gradient Descent:

$w^0$

$w^1$

$w^{i+1}$



- \* Neural network loss ( $L(w)$ ) is not convex w.r.t. the network parameters  $w$ .
- \* There exist multiple local minima.
- \* Choosing the learning rate too low leads to very slow progress.

- \* Step cliffs can pose great challenges to optimization.
- \* Gradient clipping is a method to prevent step cliffs.
- \* At a Saddle point we have  $\nabla_w L(w) = 0$ , but we are not at a minimum.
- \* A flat region is called plateau where  $\nabla_w L(w) \approx 0 \Rightarrow$  Slow progress

### Training with Momentum: momentum, $\eta =?$ (from w) vanishing gradients

$$\Delta w_{i,j}(t) := \alpha \Delta w_{i,j}(t-1) + \eta \cdot \frac{\partial L}{\partial w_{i,j}}(t) \quad \text{P.D. = gradient}$$

$\therefore$  Weight update using the velocity vector  $v(t)$

$$w_{i,j}(t+1) := w_{i,j}(t) - \Delta w_{i,j}(t)$$

gradient

### Nesterov Momentum:-

- Look ahead and calculate the gradient w.r.t. predicted parameters
- $$\hat{w}_{ij}(t+1) := w_{ij}(t) - \alpha \Delta w_{ij}(t-1)$$
- $$w_{ij}(t+1) := w_{ij}(t) - \alpha \Delta w_{ij}(t-1) + \eta \frac{\partial L}{\partial w_{ij}}(t+1)$$
- $$\Delta w_{ij}(t) := \alpha \Delta w_{ij}(t-1) + \eta \frac{\partial L}{\partial w_{ij}}(t)$$

### Adaptive Learning:

- Define a local gain( $g$ ) for each weight (Initialize with  $g=1$ )

$$\Delta w_{ij} := \eta \cdot g_{ij} \cdot \frac{\partial L}{\partial w_{ij}}$$

- If gradient is consistent ( $\rightarrow$  temp. weight is still good)

$$g_{ij}(t) := g_{ij}(t-1) + \beta$$

$$\text{else } g_{ij}(t) := g_{ij}(t-1) \cdot (1-\beta)$$

### Adaptive Learning Rate via RMSProp:

$$\text{MeanSquare}(w_{ij}, t) := \beta \cdot \text{MeanSquare}(w_{ij}, t-1) + (1-\beta) \left( \frac{\partial L}{\partial w_{ij}} \right)^2$$

$$\text{where } \beta = 0.9 \text{ to } 0.999$$

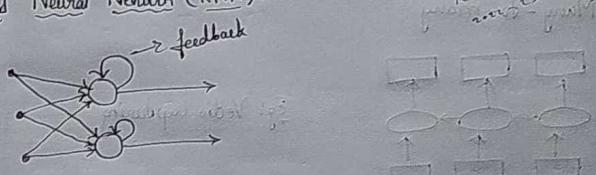
$$w_{ij}(t) := w_{ij}(t) - \eta_i \frac{\partial L}{\partial w_{ij}} / \sqrt{\text{MeanSquare}(w_{ij}, t) + \epsilon}$$

Adaptive

### Common Application of CNN:-

- Image Classification  $\rightarrow$  Object Detection  $\rightarrow$  Object Segmentation
- Image Synthesis  $\rightarrow$  Human Detection  $\rightarrow$  Face Recognition / Identification
- Facial expression  $\rightarrow$  Fatigue detection  $\rightarrow$  Lip reading
- Unmanned Aerial Vehicles (Drones)  $\rightarrow$  Biometrics
- Computer Aided Diagnosis

### Recurrent Neural Network (RNN):-



### Classical Approach:-

- Bag-of-word model

#### Raw training dataset

$x^{(1)} = \text{"The sun is shining"}$

$x^{(2)} = \text{"...."}$

$x^{(3)} = \text{"...."}$

$$y = [0, 1, 0]$$

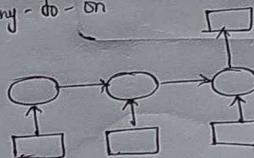
class label

dictionary	Vocabulary	Dating matrix
	"The": 0,	$\begin{bmatrix} 0 & 1 & 0 & 0 & \dots \end{bmatrix}$
	"and": 1,	$\begin{bmatrix} 0 & 0 & 1 & \dots \end{bmatrix}$
	"is": 2,	$\begin{bmatrix} 2 & 1 & 0 & 2 & \dots \end{bmatrix}$

$$y = [0, 1, 0]$$

### Different types of Sequence Modeling

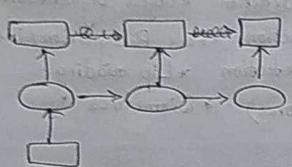
- Many-to-one



E.g.: Sentiment Analysis

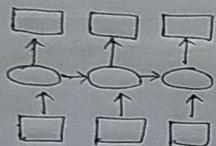
algorithm: RNN Encoder-Decoder  
decoding: beam search

⇒ One-to-many

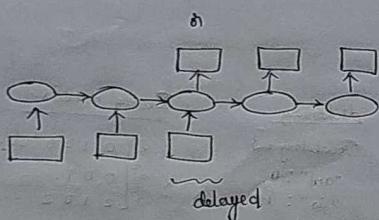


Eg: Text extracting from image

⇒ Many-to-many

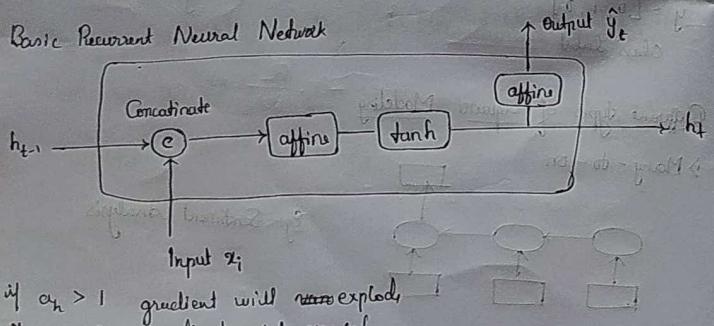


Eg: Video captioning



delayed

Basic Recurrent Neural Network



$\alpha_h > 1$  gradient will explode  
 $\alpha_h < 1$  gradient will vanish

Backpropagation through time:-

$$L = \sum_{t=1}^T L^{(t)} = \frac{\partial L^{(t)}}{\partial w_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial h^{(t)}} \cdot \left( \sum_{k=1}^T \frac{\partial h^{(k)}}{\partial h^{(t)}} \cdot \frac{\partial h^{(k)}}{\partial w_{hh}} \right)$$

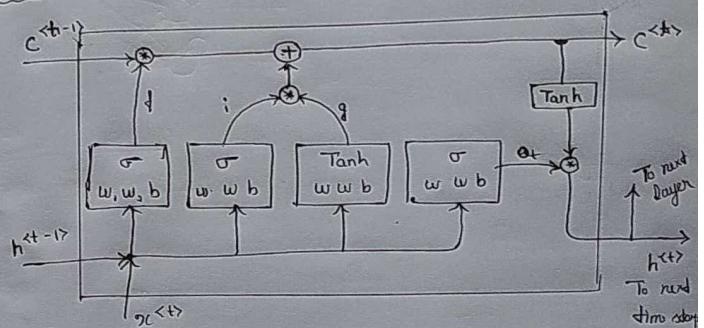
Solutions for Vanish / Exploding problem

↳ Gradient Clipping

↳ Truncated back propagation through time (TBPTT)

↳ Long short-term memory (LSTM)

LSTM:-



Forget gate,  $f_t = \sigma(w_f x^{(t)} + w_{fh} h^{(t-1)} + b_f)$

Input gate,  $i_t = \sigma(w_i x^{(t)} + w_{ih} h^{(t-1)} + b_i)$

Input node,  $g_t = \text{tanh}(w_g x^{(t)} + w_{gh} h^{(t-1)} + b_g)$

$$c^{(t)} = (c^{(t-1)} \odot f_t)$$

Output gate,  $o_t = \sigma(w_o x^{(t)} + w_{oh} h^{(t-1)} + b_o)$

$$\therefore h^{(t)} = o_t \odot \text{tanh}(c^{(t)})$$