

Deep Learning

Evaluation Policy:

Tutorials - 25

Mid Semester - 30

End Semester - 40

Assignment/Project - 30

Book: "Deep Learning" by Ian Goodfellow, Yoshua Bengio

No. of features = No. of nodes in input layer

No. of output nodes = No. of classes

$$o_j = f(w_{ij} \cdot x_j)$$

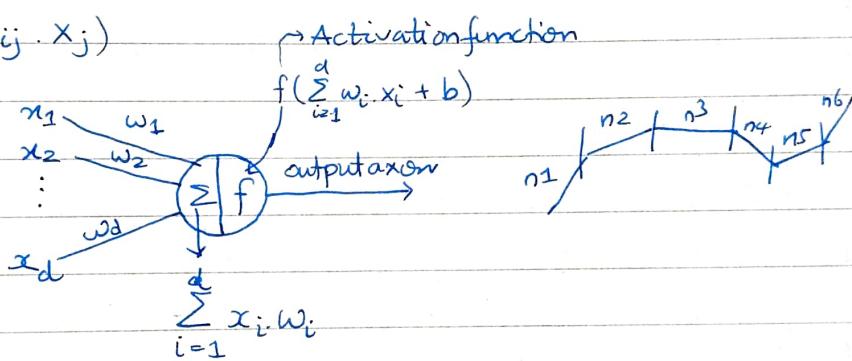


Fig.: Simple Mathematical Depiction of Neuron

$$g_i(x) = w^t x + w_0$$

$$g_i(x) = a^t y$$

$$g_i(x) = [w_1 \ w_2 \dots \ w_d \ w_0]^t [x_1 \ x_2 \ \dots \ x_d \ 1]$$

Neural Networks:

1) Neural Networks is a classifier

2) It is a supervised learning / training

3) (x_i, y_i) , x_i - i^{th} input feature, y_i - i^{th} class label
 ↳ target label / expected output

Two phases - train & test - good generalization

LDF:

- 1) The classes are linearly separable.
- 2) No probability distribution information about the class is required.
- 3) Intercept and orientation are two important parameters. $g(x) = g_1(x) - g_2(x)$
- 4) The weight vector w is orthogonal to the decision boundary.
- 5) One of the conditions to update the weight vector is called perception.
- 6) $g_i(x) = \sum_{i=1}^d w_i x_i + w_0 \Rightarrow g_i(x) \geq 0$ for correct classification
 random initially/ known \hookrightarrow Perception criterion function
 Uniform criterion

15 Neural Networks:

- 1) Changeable parameters : w
- 2) x_i are given as fixed inputs.
- 3) y_i is also fixed

20
$$g_i(x) = \sum_{i=1}^d w_i^T x_i + w_0 = w^T x$$

$$g_i(x) = a^T y > 0.$$
 (Append 1 to x , w_0 to w)
 \hookrightarrow uniform criterion function.

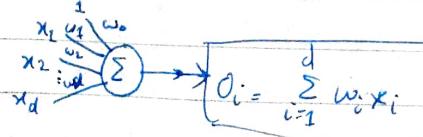
- 4) Perception criterion function

25 $g_i(x) > 0$

- 5) Relaxation criterion function

$g_i(x) > b$ \hookrightarrow bias.

- 6) To depict/ get a complex non-linear classes, you need to solve non-linearly separable classes.



Compare actual and expected output

$$E = \sum_{i=1}^N (A.O. - E.O.)^2 \quad \text{Min. } E \text{ by changing } w.$$

Cost function

Min. cost function $E \rightarrow$ correct class predicted.

7) $a[0] \Rightarrow$ randomly chosen weight vector.

$$8) a_{k+1} = a_k + \eta \sum_{y \in \text{Misclassified}} y$$

25/01/2024

Neural Networks (Cont'd)

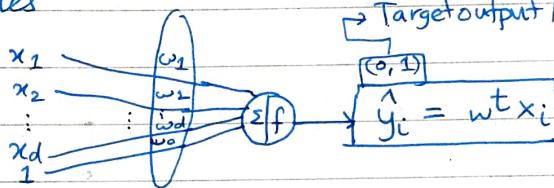
↑ no. of samples

$$i = 1, 2, \dots, N$$

It is supervised learning

$$(x_i, y_i)$$

↳ label



→ Target output / Expected output

y_i : Expected output (or) target output

\hat{y}_i : Actual output

x_i : i^{th} Sample Feature Vector

N : Total number of samples

$$E = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2, \text{ where } E \text{ is the sum of squared Error}$$

$$\frac{\partial E}{\partial w} = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i) \cdot \frac{\partial}{\partial w} (\hat{y}_i - y_i)$$

↓ target output

$$\frac{\partial E}{\partial w} = \sum_{i=1}^N (\hat{y}_i - y_i) \cdot \frac{\partial}{\partial w} (w^t \cdot x_i)$$

$$\frac{\partial E}{\partial w} = \sum_{i=1}^N (\hat{y}_i - y_i) \cdot x_i$$

Weight updation Rule:

$$W_{\text{new}} \leftarrow W - \eta \cdot \frac{\partial E}{\partial W}$$

5

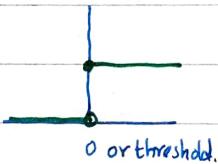
$$W_{\text{new}} \leftarrow W_{\text{old}} - \eta \sum_{i=1}^N (\hat{y}_i - y_i) \cdot X_i$$

Repeat until all samples are correctly classified.

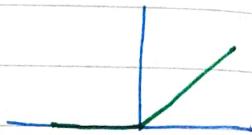
10 $a^T y > 0$: Uniform criterion function.

SLSO- Single Layer Single Output is used for simple perceptron

15 Step function:



ReLU function:



Training in Neural Network (Fill in the Blanks):

- 1) In this, the changeable parameter is weight vector (w)
- 2) In neural network, we should have some weight updation sub
- 3) The weight updation has to be done using training sample
- 4) Keep on updating weight, until all samples are correctly classified
- 5) Since it is supervised learning, we know expected / target output
- 6) The actual output predicted by NN is called y (predicted output)
- 7) The difference between A/P value and T/E value is error.
- 8) The convergence is until all samples are correctly classified
- 9) Single output NN takes care of two-class problem
- 10) But if multiple class problem is involved, then we should have neural network having a single layer Multi-class per
- 11) O/P N.L. function is called activation function.
- 12) Non-linear function used is step function.
- 13) The step function is used usually only in output layer.

- 14) The same function cannot be used in hidden layers as it is not differentiable because there is an abrupt change at $T=0$
- 15) Otherwise, sigmoid function is used at hidden layer.

30/1/2024

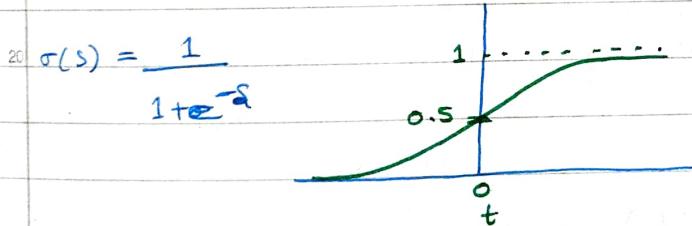
Weight updation rule

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \frac{\partial E}{\partial w}$$

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \sum_{i=1}^N (\hat{y}_i - y_i) \cdot x_i \quad \leftarrow \text{Batch optimization}$$

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta (\hat{y}_i - y_i) \cdot x_i \quad \leftarrow \text{Stochastic optimization}$$

The step function is not differentiable, so we can use the sigmoid activation function.



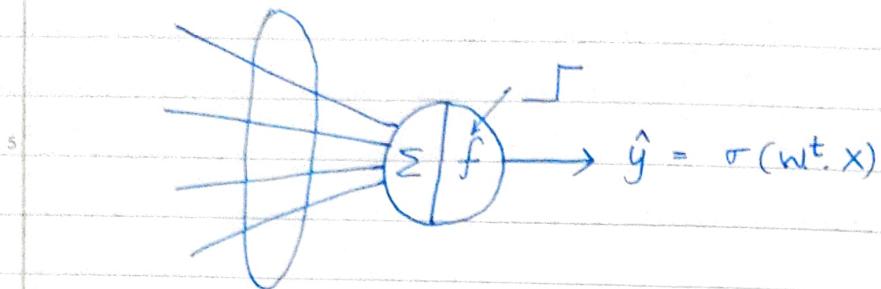
25) $\sigma(s) = \frac{1}{1+e^{-s}}, \quad s = w^T \cdot x$

when $s=0; \frac{1}{1+e^{-0}} = 0.5$

when $s \rightarrow \infty; \frac{1}{1+e^{-\infty}} \approx 1$

when $s \rightarrow -\infty; \frac{1}{1+e^{\infty}} \approx 0$

SLSO with sigmoid as non-linear function



The derivative of sigmoidal function is $\sigma(s)(1-\sigma(s))$

$$\sigma'(s) = \sigma(s)(1-\sigma(s))$$

\downarrow target value
 (x_i, y_i)

$$E = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

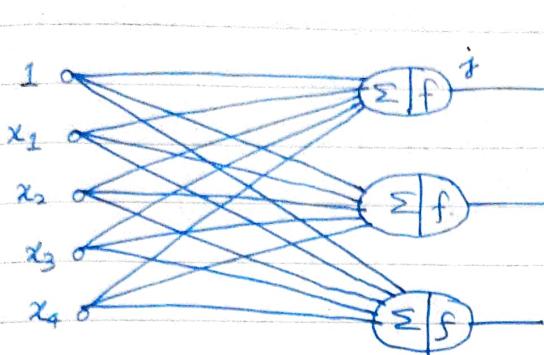
$$\frac{\partial E}{\partial w} = \frac{2}{2} \sum_{i=1}^N (\hat{y}_i - y_i) \cdot \frac{\partial}{\partial w} (\hat{y}_i - y_i)$$

$$\frac{\partial E}{\partial w} = \sum_{i=1}^N (\hat{y}_i - y_i) \cdot \frac{\partial}{\partial w} (\hat{y}_i)$$

$$\frac{\partial E}{\partial w} = \sum_{i=1}^N (\hat{y}_i - y_i) \cdot \frac{\partial}{\partial w} (\overbrace{\sigma(\hat{y}_i)}^{\sigma(\hat{y}_i)})$$

$$\frac{\partial E}{\partial w} = \sum_{i=1}^N (\hat{y}_i - y_i) \cdot x_i (\hat{y}_i (1 - \hat{y}_i))$$

SLMO (Single Layer Multiple Output)



$$\theta_j = f \left(\sum_{i=1}^d w_{ij} \cdot x_i \right)$$

↗ Sigmoidal function

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial \theta_j} \cdot \frac{\partial \theta_j}{\partial \theta_j} \cdot \frac{\partial \theta_j}{\partial w_{ij}}$$

$$\theta_j = \frac{1}{1 + e^{-\theta_j}}, \quad \theta_j = \sum_{i=1}^d w_{ij} \cdot x_i$$

θ_j - Predicted value / Actual Value

$(x_i, y_i), t_j$ - target output

x_i - Feature vector

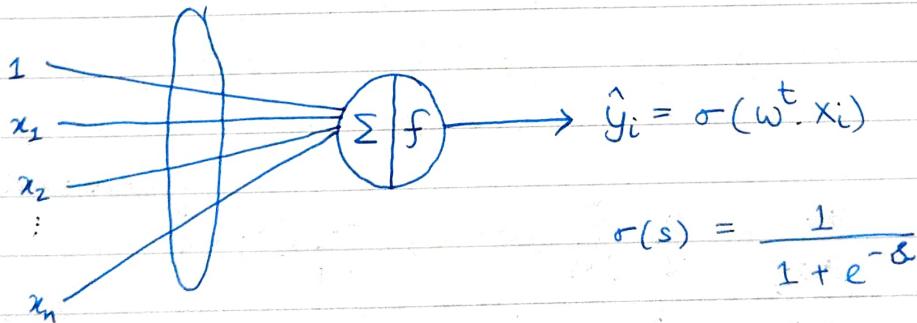
f - Sigmoid function

w_{ij} - Weight from i^{th} neuron (input layer) to j^{th} neuron

$$E = \sum_{j=1}^c (\theta_j - t_j)^2, \text{ where } c \text{ is the number of classes.}$$

31/1/2024

Recap (SLSO)



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

$$\sigma(w^t \cdot x) = \frac{1}{1 + e^{-w^t \cdot x}}; \quad \frac{\partial}{\partial s} (\sigma(s)) = \sigma(s)(1 - \sigma(s))$$

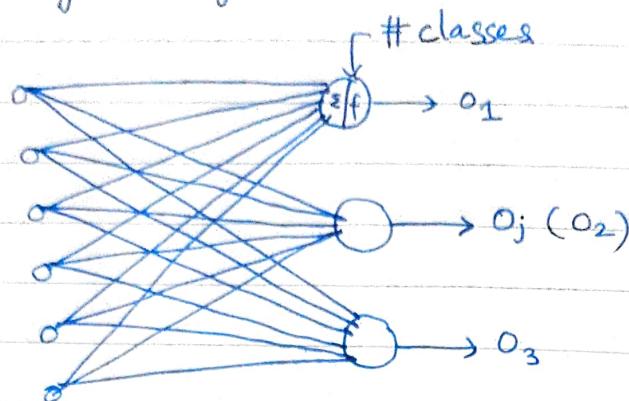
$$\sigma'(w^t \cdot x) = \sigma(w^t \cdot x) \cdot (1 - \sigma(w^t \cdot x)) \cdot x$$

Chain rule:

$$y = f(g(x))$$

$$y' = f'(g(x)) \cdot g'(x)$$

SLMO (Single Layer Multiple Output)



$$o_j = \sigma \left(\sum_{i=1}^d w_{ij} x_i \right), \quad \theta_j = \sum_{i=1}^d w_{ij} \cdot x_i, \quad o_j = \frac{1}{1 + e^{-\theta_j}}$$

ith input neuron jth output neuron

Non-linear function - activation function.

Hidden layer - No restriction on numbers of nodes

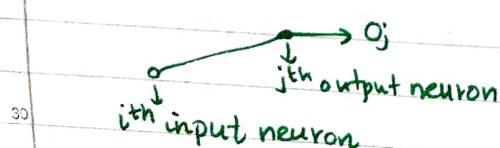
$i = 1, 2, \dots, d; j = 1, 2, \dots, c$, where d : dimension of F.V.

c : Number of classes

$$E = \frac{1}{2} \sum_{j=1}^c (o_j - t_j)^2, \text{ where } t_j \text{ or } y_j \text{ is called target output}$$

$$\frac{\partial E}{\partial w_{ij}} = \frac{2}{2} \sum_{j=1}^c (o_j - t_j) \cdot \frac{\partial}{\partial w} (o_j - t_j) \quad \text{constant}$$

$$\frac{\partial E}{\partial w_{ij}} = \sum_{j=1}^c (o_j - t_j) \cdot \frac{\partial}{\partial w} (o_j)$$



$$o_j = \sigma \left(\sum_{i=1}^d w_{ij} \cdot x_i \right)$$

$$o_j = \frac{1}{1 + e^{-\theta_j}}, \text{ where } \theta_j = \sum_{i=1}^d w_{ij} \cdot x_i$$

$$\frac{\partial E}{\partial w_{ij}} = \sum_{j=1}^c (o_j - t_j) \cdot \frac{\partial}{\partial w} (o_j)$$

$$\frac{\partial E}{\partial w_{ij}} = \sum_{j=1}^c (o_j - t_j) \cdot \frac{\partial o_i}{\partial \theta_j} \cdot \frac{\partial \theta_j}{\partial w_{ij}}$$

$$\frac{\partial E}{\partial w_{ij}} = \sum_{j=1}^c (o_j - t_j) \cdot o_j(1-o_j) \cdot x_i \quad (\because \sigma'(s) = \sigma(s)(1-\sigma(s)))$$

10 Weight updation rule :

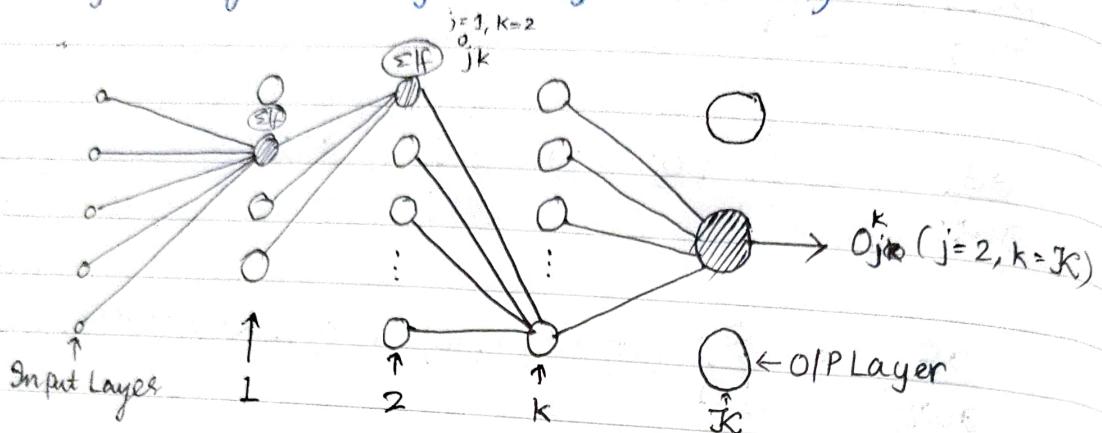
$$w_{new} \leftarrow w_{old} - \eta \frac{\partial E}{\partial w_{ij}}$$

$$15 \quad w_{new} \leftarrow w_{old} - \eta \sum_{j=1}^c (o_j - t_j) o_j(1-o_j) x_i$$

Repeat for 'n' iterations.

$$20 \quad \text{shape of } o_j = \begin{matrix} \text{shape of} \\ (w_{ij} \cdot x_i) \end{matrix} \quad (1 \times d+1) \\ (1 \times 1) \quad (1 \times d+1)$$

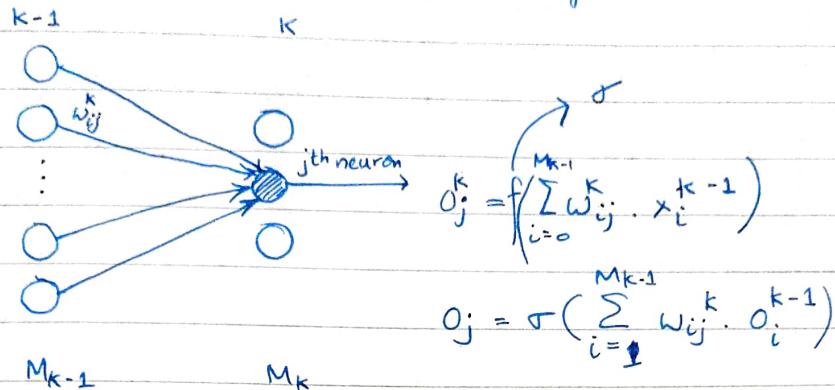
Multiple Layer Multiple Output (MLMO)



O_j^k - Output of j^{th} neuron at k^{th} layer.

MLMO (Multiple Layer Multiple output)

Neural Network (or) Multilayer Perceptron

Weight updation rule at 'output layer' $k \rightarrow k-1$.

(C) No. of classes.

 M_k : No. of nodes at k^{th} layer M_{k-1} : No. of nodes at $(k-1)^{\text{th}}$ layer o_j^k : Output at j^{th} neuron at k^{th} layer w_{ij}^k : Weight between i^{th} and j^{th} neuron to k^{th} layer t_j^k : target value. o_i^{k-1} : O/P obtained at i^{th} neuron in $(k-1)^{\text{th}}$ layer.

$$E = o_j^k - t_j^k$$

$$o_j^k = \frac{1}{1 + e^{-\theta_j}} ; \quad \theta_j^k = \sum_{i=1}^{M_{k-1}} w_{ij}^k \cdot o_i^{k-1}$$

$$E = \frac{1}{2} \sum_{j=1}^{M_k} (o_j^k - t_j^k)^2$$

$$\frac{\partial E}{\partial w_{ij}^k} = \frac{2}{2} \sum_{j=1}^{M_k} (o_j^k - t_j^k) \cdot \frac{\partial}{\partial w_{ij}^k} (o_j^k - t_j^k) \quad \text{fixed}$$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}^k} = \sum_{j=1}^{M_k} (o_j^k - t_j^k) \cdot \frac{\partial o_j^k}{\partial \theta_j^k} \cdot \frac{\partial \theta_j^k}{\partial w_{ij}^k}$$

$$\Rightarrow \frac{\partial E}{\partial w_{ij}^k} = \sum_{j=1}^{M_k} (\underbrace{o_j^k - t_j^k}_{\delta_j^k}) \cdot \underbrace{o_j^k (1-o_j^k)}_{s_j^k} \cdot o_i^{k-1}$$

$$\because \frac{\partial o_j^k}{\partial \theta_j^k} = \frac{\partial}{\partial \theta_j^k} \left(\frac{1}{1+e^{-\theta_j^k}} \right) = o_j^k (1-o_j^k)$$

$$\therefore \sigma'(s) = \sigma(s)(1-\sigma(s)), s = o_j^k$$

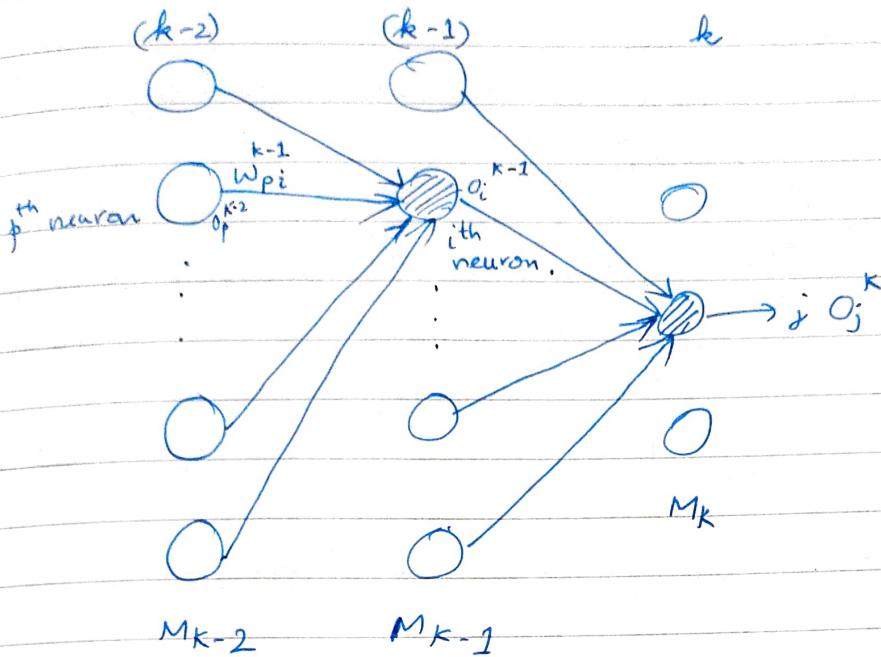
$$\frac{\partial \theta_i^k}{\partial w_{ij}^k} = \frac{\partial}{\partial w_{ij}^k} \left(\sum_{l=1}^{M_{k-1}} w_{lj}^k \cdot \underline{o_l^{k-1}} \right) = o_i^{k-1}$$

$$\frac{\partial E}{\partial w_{ij}^k} = \sum_{i=1}^{M_{k-1}} \sum_{j=1}^{M_k} \underbrace{\delta_j^k \cdot o_i^{k-1}}_{\text{Backprop. intuition.}}$$

The error obtained in j th neuron of k th layer is backpropagated to all i neurons in $(k-1)$ th layer connected to j th neuron of k th layer.

$$w_{ij}^k \text{ new} \leftarrow w_{ij}^k - \eta \delta_j^k \cdot o_i^{k-1}$$

Weight updation rule between hidden layer ($k-1$) to ($k-2$)



$$o_i^{k-1} = \frac{1}{1 + e^{-\theta_i^{k-1}}} ; \quad \theta_i^{k-1} = \sum_{p=1}^{M_{k-2}} w_{pi}^{k-1} \cdot o_p^{k-2}$$

~~$$E = \frac{1}{2} \left(\sum_{i=1}^{M_{k-2}} (o_i^{k-1} - t_i^{k-1})^2 \right)$$~~

$$E = \frac{1}{2} \sum_{j=1}^{M_k} (o_j^k - t_j^k)^2$$

$$\frac{\partial E}{\partial w_{pi}^{k-1}} = \frac{2}{2} \sum_{j=1}^{M_k} (o_j^k - t_j^k) \cdot \frac{\partial}{\partial w_{pi}^{k-1}} (o_j^k - t_j^k)$$

fixed.

$$= \sum_{j=1}^{M_k} (o_j^k - t_j^k) \cdot \frac{\partial o_j^k}{\partial \theta_j^k} \cdot \frac{\partial \theta_j^k}{\partial o_i^{k-1}} \cdot \frac{\partial o_i^{k-1}}{\partial \theta_i^{k-1}} \cdot \frac{\partial \theta_i^{k-1}}{\partial w_{pi}^{k-1}}$$

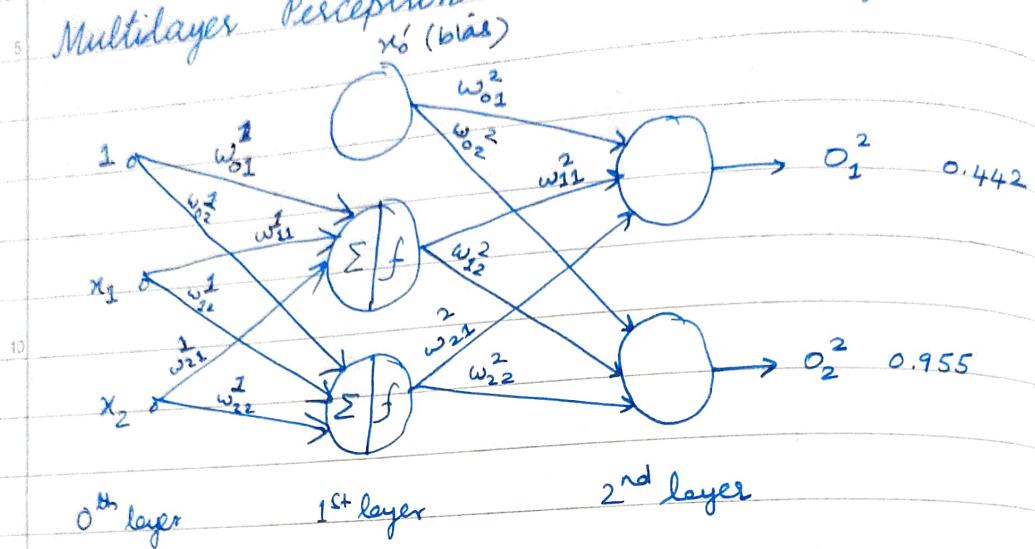
$$\frac{\partial E}{\partial w_{pi}^{k-1}} = \sum_{j=1}^{M_k} (o_j^k - t_j^k) \cdot \underbrace{o_j^k (1 - o_j^k)}_{s_j^k} \cdot w_{ji}^k \cdot \underbrace{o_i^{k-1} (1 - o_i^{k-1})}_{s_i^{k-1}} \cdot \underbrace{o_p^{k-2}}_{s_p^{k-1}}$$

Error from i^{th} neuron of k^{th} layer is passed to p^{th} neuron of $(k-2)^{th}$ layer.

$$w_{pi}^{k-1} = w_{pi}^{k-1} - \eta \delta_i^{k-1} \cdot o_p^{k-2}$$

2/2/2024

Multilayer Perception - Problem Solving



do

- ① Feed forward pass
- ② Back propagate pass (derivative of error)
- ③ Update the weight until 'E' is minimized.

[similar qn. for midsem]

$$\begin{bmatrix} w_{01}^1 & w_{11}^1 & w_{21}^1 \\ w_{02}^1 & w_{12}^1 & w_{22}^1 \end{bmatrix} \begin{bmatrix} w_{01}^2 & w_{11}^2 & w_{21}^2 \\ w_{02}^2 & w_{12}^2 & w_{22}^2 \end{bmatrix} \times = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} t = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 0.5 & 1.5 & 0.8 \\ 0.8 & 0.2 & -1.6 \end{bmatrix} \begin{bmatrix} 0.9 & -1.7 & 1.6 \\ 1.2 & 2.1 & -0.2 \end{bmatrix} x = \begin{bmatrix} 1 \\ 0.7 \\ 1.2 \end{bmatrix} t = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

- 1) Feed forward pass [layer '0' to layer '1']

$$\begin{bmatrix} 0.5 & 1.5 & 0.8 \\ 0.8 & 0.2 & -1.6 \end{bmatrix} \begin{bmatrix} 1 \\ 0.7 \\ 1.2 \end{bmatrix} = \begin{bmatrix} 2.51 \\ -0.98 \end{bmatrix}$$

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

$$\Rightarrow \begin{bmatrix} \frac{1}{1+e^{-2.51}} \\ \frac{1}{1+e^{-(0.98)}} \end{bmatrix} = \begin{bmatrix} 0.92 \\ 0.27 \end{bmatrix} \rightarrow \begin{array}{l} x_1^2/o_1^2 \\ x_2^2/o_2^2 \end{array}$$

Feed forward from layer 1 to 2.

$$\begin{bmatrix} 0.9 & -1.7 & 1.6 \\ 1.2 & 2.1 & -0.2 \end{bmatrix} \begin{bmatrix} 1 \\ 0.92 \\ 0.27 \end{bmatrix} = \begin{bmatrix} -0.232 \\ 3.057 \end{bmatrix}$$

$$\text{sigmoid} \left(\begin{bmatrix} -0.232 \\ 3.057 \end{bmatrix} \right) = \begin{bmatrix} 0.442 \\ 0.955 \end{bmatrix} \begin{array}{l} \xrightarrow{o_1^2} \\ \xrightarrow{o_2^2} \end{array}$$

Backpropagation from output layer(k) to previous layer(k-1)

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial E}{\partial w}$$

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \delta_j^k \cdot o_i^{k-1}$$

$$\delta_j^k = o_j^k (1-o_j^k) \cdot (o_j^k - t_j^k)$$

$$\text{Find } \delta_1^2 - o_1^2 \text{ and } \delta_2^2 - o_2^2$$

$$\delta_1^2 = o_1^2 (1-o_1^2) \cdot (o_1^2 - t_1^2)$$

$$= 0.442 (1-0.442) (0.442-1) = -0.138,$$

$$\boxed{\delta_1^2 = -0.14}$$

$$\delta_2^2 = \omega_2^2 (1 - \omega_2^2) (\omega_2^2 - t_2^2)$$

$$\delta_2^2 = 0.955 (1 - 0.955) (0.955 - 0)$$

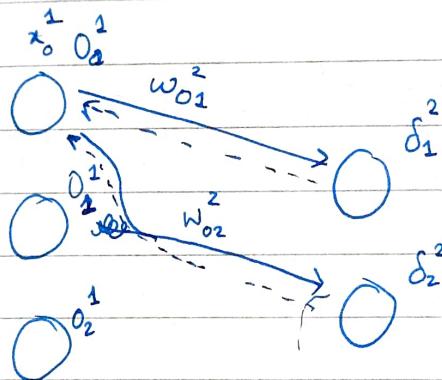
$$\boxed{\delta_2^2 = 0.045}$$

$$\delta_j^k = \omega_j^k (1 - \omega_j^k) (\omega_j^k - t_j^k)$$

$$\frac{\partial E}{\partial w} = \delta_j^k \cdot \omega_i^{k-1}$$

Formulae.

$$\frac{\partial E}{\partial w_{01}^2} = \delta_1^2 \cdot \omega_0^1 = -0.138 \times 1 = -0.138$$



Reference
Diagram

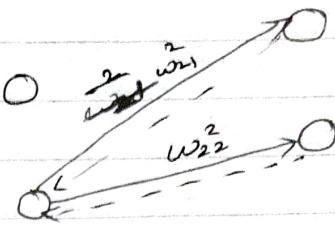
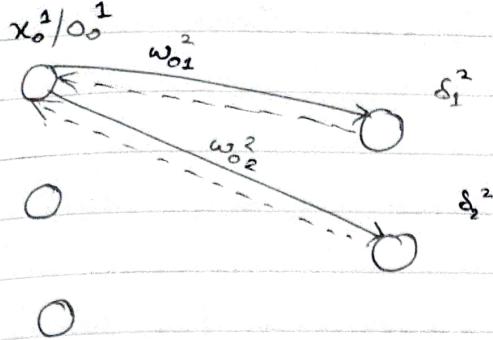
$$\frac{\partial E}{\partial w_{02}^2} = \delta_2^2 \cdot \omega_0^1 = 0.045 \times 1 = 0.045$$

$$\frac{\partial E}{\partial w_{11}^2} = \delta_1^2 \cdot \omega_1^1 = -0.138 \times 0.92 \approx -0.127$$

$$\frac{\partial E}{\partial w_{12}^2} = \delta_2^2 \cdot \omega_1^1 = 0.045 \times 0.92 = 0.0414$$

$$\frac{\partial E}{\partial w_{21}^2} = \delta_1^2 \cdot \omega_2^1 = -0.138 \times 0.27 = -0.037$$

$$\frac{\partial E}{\partial w_{22}^2} = \delta_2^2 \cdot \omega_2^1 = 0.045 \times 0.27 = 0.012$$



Update the respective weights :-

$$w_{\text{new}} \leftarrow w_{\text{old}} - \eta \cdot \frac{\partial E}{\partial w}, \quad \boxed{\eta = 0.5}$$

$$w_{01}^2 \leftarrow w_{01}^2 - \eta \cdot \frac{\partial E}{\partial w_{01}^2} = 0.9 - 0.5 \times -0.138 = 0.969$$

$$w_{11}^2 \leftarrow w_{11}^2 - \eta \cdot \frac{\partial E}{\partial w_{11}^2} = -1.7 - 0.5 \times -0.126 = -1.637$$

$$w_{02}^2 \leftarrow w_{02}^2 - \eta \cdot \frac{\partial E}{\partial w_{02}^2} = 1.2 - 0.5 \times 0.045 = 1.777$$

$$w_{12}^2 \leftarrow w_{12}^2 - \eta \cdot \frac{\partial E}{\partial w_{12}^2} = 2.1 - 0.5 \times 0.04 = 2.07$$

$$w_{21}^2 \leftarrow w_{21}^2 - \eta \cdot \frac{\partial E}{\partial w_{21}^2} = 1.6 - 0.5 \times -0.037 = 1.618$$

$$w_{22}^2 \leftarrow w_{22}^2 - \eta \cdot \frac{\partial E}{\partial w_{22}^2} = -0.2 - 0.5 \times 0.012 = -0.206$$

Next, we have to update the weights in layer 1.

$$\frac{\partial E}{\partial w_{01}^2}, \frac{\partial E}{\partial w_{11}^2}, \frac{\partial E}{\partial w_{21}^2}, \frac{\partial E}{\partial w_{02}^2}, \frac{\partial E}{\partial w_{12}^2}, \frac{\partial E}{\partial w_{22}^2}$$

$$-0.137 = w_{01}$$

$$-0.126 = w_{11}$$

$$-0.037 = w_{21}$$

$$0.3989 = w_{02}$$

$$0.0369 = w_{12}$$

$$-0.2088 = w_{22}$$

Algorithm for back propagation at any hidden layer 'k':

① For any hidden layer weight w_{ij}^k ,

$$\frac{\delta E}{\delta w_{ij}^k} = \delta_j^k \cdot o_i^{k-1} \quad \text{Error is back propagated to layer.}$$

$$\delta_i^k = \left(\sum_{j=1}^{M_{k+1}} \delta_j^{k+1} \cdot w_{ij}^{k+1} \right) o_i^k (1 - o_i^k)$$

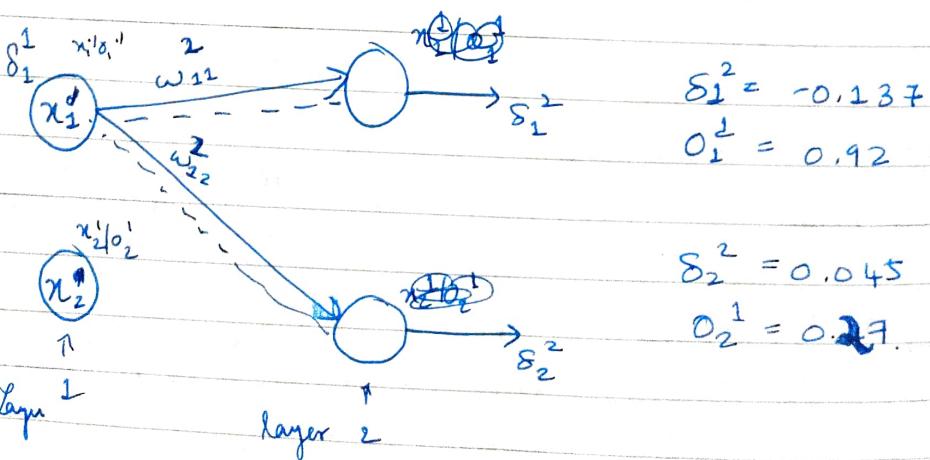
$$\delta_i^k = \sum_{j=1}^{M_{k+1}} \delta_j^{k+1} \cdot w_{ij}^{k+1}$$

③ Weight updation rule for any hidden layer

$$w_{ij}^k \leftarrow w_{ij}^k - \eta \frac{\partial E}{\partial w_{ij}^k}$$

$$\frac{\partial E}{\partial w_{ij}^k} = \delta_j^k \cdot o_i^{k-1}$$

Back propagation from $k-1 \rightarrow k-2$ (1 to 0)



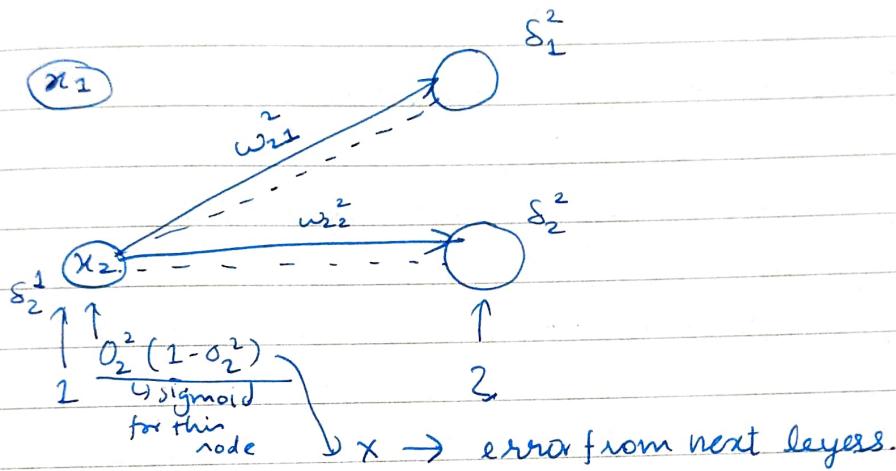
$$\delta_1^1 = o_1^1 (1 - o_1^1) \cdot [s_1^2 \cdot w_{11}^2 + s_2^2 \cdot w_{12}^2]$$

 $\delta_1^1 / 2$

First find error at 1st layer!

$$\delta_2^1 = 0.92 (1 - 0.92) [-0.137 \times -1.7 + 0.045 \times 2.1]$$

$$\boxed{\delta_1^1 = 0.0241.}$$



$$\delta_2^1 = 0.27 (1 - 0.27) [-\underbrace{0.137 \times 1.6}_{\delta_2^1} + \underbrace{0.045 \times -0.2}_{w_{22}^2}]$$

$$\boxed{\delta_2^1 = -0.044}$$

$$\frac{\partial E}{\partial w_{o1}^1} = \delta_1^1 \cdot x_1^0 = 0.024 \times \frac{1}{0.7} = 0.0168 \approx 0.017 = 0.024$$

$$\frac{\partial E}{\partial w_{o2}^1} = \delta_2^1 \cdot x_1^0 = -0.04 \times 1 = -0.04$$

$$\frac{\partial E}{\partial w_{11}^1} = \delta_1^1 \cdot x_1^0 = 0.024 \times 0.7 = 0.017$$

$$\frac{\partial E}{\partial w_{12}^1} = s_2^1 \cdot o_1^0 = s_2^1 \cdot x_1^0 = -0.04 \times 0.7 = -0.028$$

$$\frac{\partial E}{\partial w_{21}^1} = s_1^1 \cdot x_2^0 = 0.024 \times 1.2 = 0.0288$$

$$\frac{\partial E}{\partial w_{22}^1} = s_2^1 \cdot x_2^0 = -0.04 \times 1.2 = -0.048$$

10) $w_{01}^1 \leftarrow w_{01}^1 - \eta \cdot \frac{\partial E}{\partial w_{01}^1}$

Chptr 5 - PRML

$$w_{01}^1 \leftarrow 0.5 - 0.5 \times 0.024 = 0.488$$

15) $w_{02}^1 = w_{02}^1 - \eta \cdot \frac{\partial E}{\partial w_{02}^1} = 0.8 - 0.5 \times 0.04 = 0.82$

$$w_{21}^1 = 1.5 - 0.5 \times 0.017 = +1.4915$$

$$w_{12}^1 = 0.2 - 0.5 \times -0.028 = 0.214$$

20) $w_{22}^1 = 0.8 - 0.5 \times 0.0288 = -0.786$

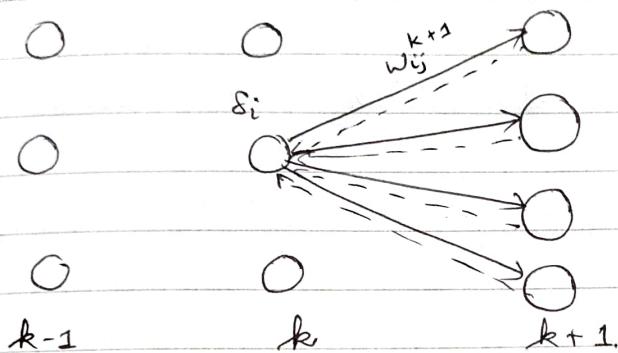
$$w_{22}^1 = -1.6 - 0.5 \times -0.048 = -1.594 / -1.576$$

Back propagation algorithm for any hidden layer 'k':

w_{ij}^k is the weight connecting for any hidden layer 'k'.

① calculate δ_i^k as follows:

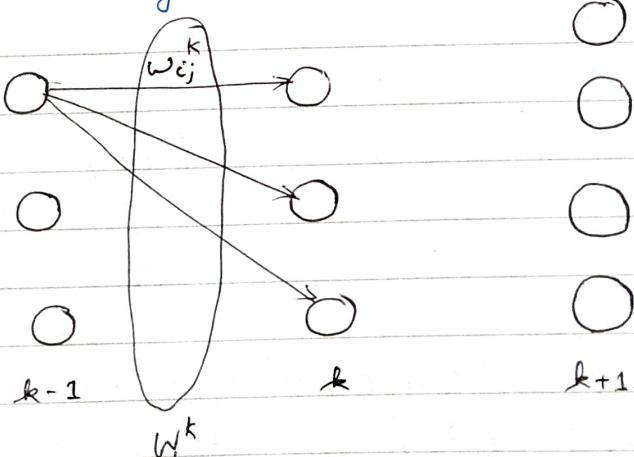
$$\delta_i^k = o_i^k (1 - o_i^k) \sum_{j=1}^{M_{k+1}} \delta_j \cdot w_{ij}^{k+1}$$



② Update the weight w_{ij}^k as follows:

$$w_{ij}^k \leftarrow w_{ij}^k - \eta \cdot \frac{\partial E}{\partial w_{ij}^k}$$

③ $k-1^{\text{th}}$ to k^{th} layer.



$$w_{ij}^k = w_{ij}^k - \eta \cdot \delta_j^k \cdot o_i^{k-1}$$

① Feed forward pass

$$f\left(\sum_{i=1}^d x_i \cdot w_i\right)$$

② Back propagation at output layer

$$\delta_i^k = (o_i^k - t_i^k) o_i^k (1 - o_i^k)$$

③ Back propagation at hidden layer 'k'

$$\delta_i^k = o_i^k (1 - o_i^k) \cdot \sum_{j=1}^{M_{k+1}} \delta_j^{k+1} w_{ij}^{k+1}$$

$$w_{ij}^k \leftarrow w_{ij}^k - \eta \cdot \frac{\partial E}{\partial w_{ij}^k}$$

$$\Rightarrow w_{ij}^k \leftarrow w_{ij}^k - \eta \cdot \delta_j^k \cdot o_i^{k-1} \quad i \rightarrow \text{left}, j \rightarrow \text{right layer} \\ (\text{notation}).$$

23) Simple NN for linear separable problem:

$\underbrace{x_1}_{F.V.} \quad \underbrace{x_2}_{F.V.} \quad \underbrace{y}_{\text{Target output}}$

0	0	0
0	1	0
1	0	0
1	1	1

$$d=2$$

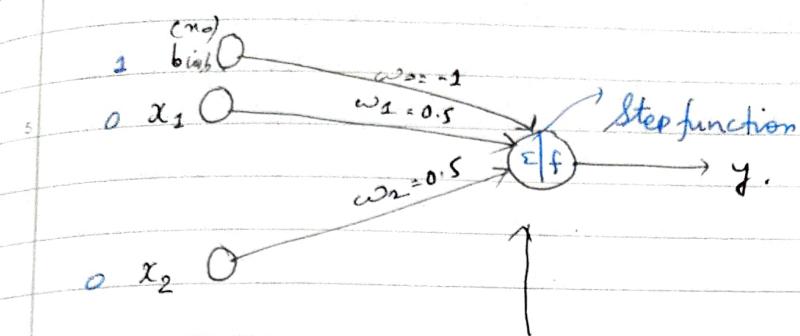
N=4, # samples.

samples

- ① # Neurons in IIP layer depends on no. of F.V. = 2 (here)
- ② # Neurons in OIP layer depends on no. of classes = 1 (here)
- ③ Decide the number of layers required.

Initialize weight vector

Decide non-linear function / activation function.

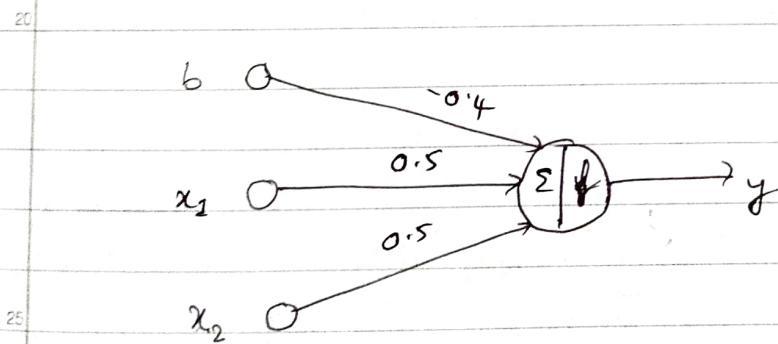


10 Input Layer (0th layer) Output layer.

Bias = 1 (Initially)

Draw perceptron for OR gate

x	y	$x \text{ OR } y$
0	0	0
0	1	1
1	0	1
1	1	1



X-NOR gate

$$\bar{x}_1 \bar{x}_2 + x_1 \bar{x}_2$$

$$\boxed{\bar{x}_1 \bar{x}_2 + x_1 \bar{x}_2}$$

$$x_1 \bar{x}_2 + \bar{x}_1 x_2$$

(Non-linearly separable)

x_1 x_2 y

0 0 0

0 1 1

1 0 1

1 1 0

$$h_1 = \boxed{x_1 + x_2}$$

OR

$$h_2 = \boxed{\bar{x}_1, \bar{x}_2}$$

NAND

$$h_1, h_2$$

AND

0 1 0

1 1 1

1 1 1

1 0 0

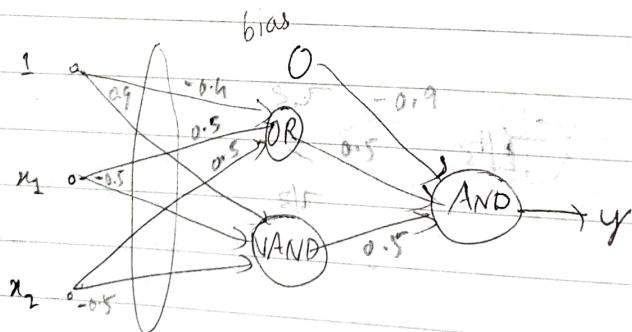
$w \cdot x$

$$\sum_{i=1}^d w_i \cdot \phi(x)$$

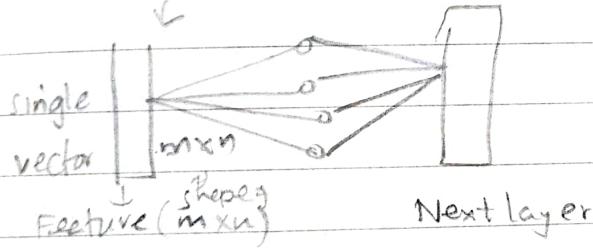
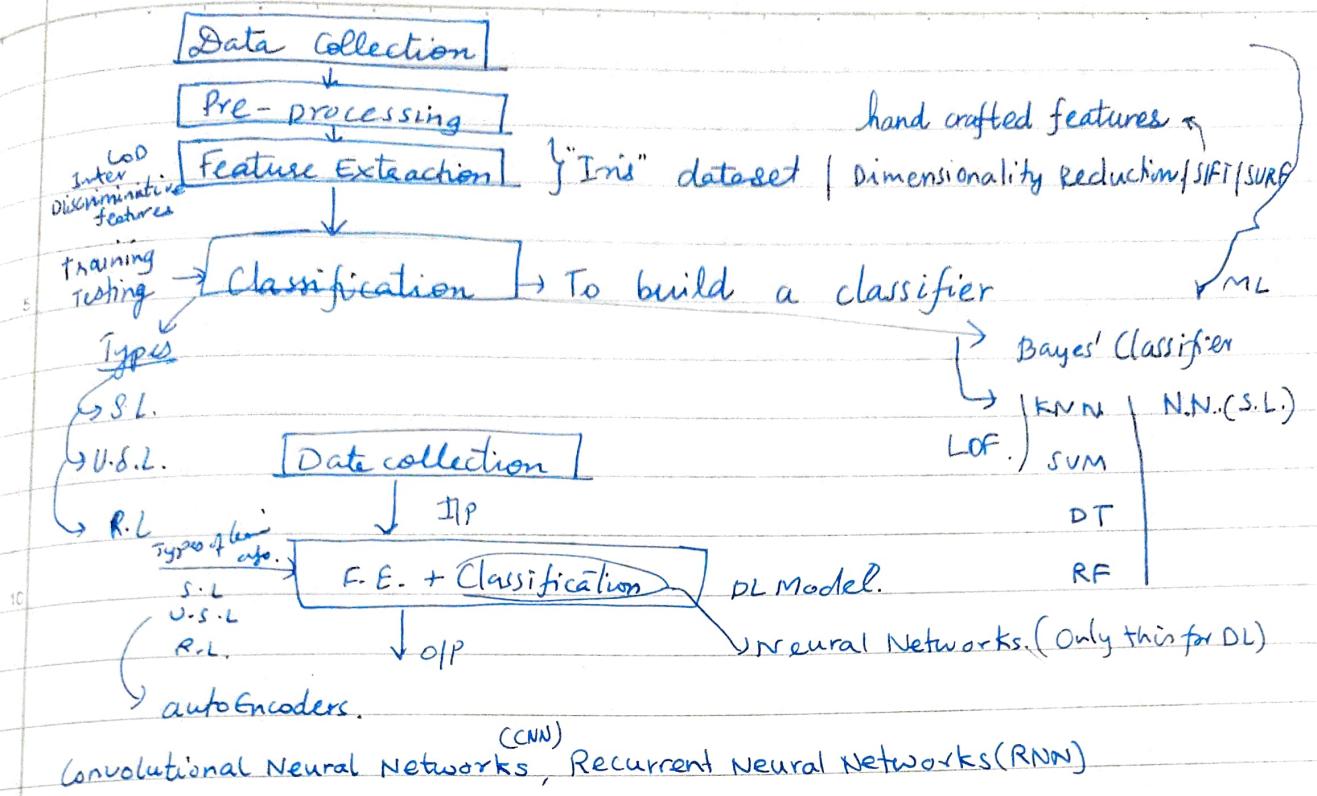
$w \cdot \phi(x)$

If samples are not linearly separable in low-dim. space is transformed to higher dim. space to make it linearly separable.

Draw neural network for the same. (H/W)



Find the weight vectors



G PLUS

All pixel intensity values - feature (Curse of Dimensionality).

ANN - Just extracts from given features. (Nothing new).

CNN does extraction of features.

25 ANN & CNN - Feed Forward NN.

Recurrent Neural Network