

Unit - 4

Introduction to machine learning

Part of AI

Training

- * Task to learn to classify inputs according to a finite set of classifications.
- * Learn from training data.
- * Classify new data such that it has not seen.
 $f(x) = y$.

Rate Learning.

- * Training which stores piece of training data & its classification.
- New data classified if it's stored in memory.

Version Space.

Given a set of training examples (positive & negative), the set of hypothesis that correctly map each of the training examples to its classification is called Version spaces.

~~Find 5 algorithm~~

| | Citations | size | in library | Price | editions | buy |
|---|-----------|--------|------------|------------|----------|-----|
| 1 | some | small | no | affordable | many | no |
| 2 | many | big | no | expensive | one | Yes |
| 3 | some | big | always | expensive | few | no |
| 4 | many | medium | no | expensive | many | Yes |
| 5 | many | small | no | affordable | many | Yes |

Applying find 5 algorithm to find maximally specific hypothesis

In this data set we have 5 attributes & the target variable is buy.

1. How many concepts are possible for this instance space.
2. How many hypothesis can be expressed by hypothesis language.
3. Apply Find 5 algorithm on give training set

Ans. 1 $2 \times 3 \times 2 \times 2 \times 3 = 72$

Ans. 2 There each attribute along with 2 more possibilities are considered ? and NULL.

$$4 \times 5 \times 4 \times 4 \times 5 = 1600$$

Ans : 3

step 1: Initially we initialize hypothesis to more specific value.

$$h_0 = \{ \phi, \phi, \phi, \phi, \phi, \phi \}$$

step 2: $x_1 = \{ \text{some}, \text{small}, \text{no}, \text{affordable}, \text{many} \}$ - No

$$h_1 = \{ \phi, \phi, \phi, \phi, \phi, \phi \}$$

$$x_2 = \{ \text{many}, \text{big}, \text{no}, \text{expensive}, \text{one} \} - \text{Yes}$$

$$h_2 = \{ \text{many}, \text{big}, \text{no}, \text{expensive}, \text{one} \}$$

It is a positive example, here we compare the attribute values in the example & the attribute values in the hypothesis.

If there is a match will keep the same value as it is. & if there is no match will replace the attribute value hypothesis with next most general values.

$$x_3 = \{ \text{some}, \text{big}, \text{always}, \text{expensive}, \text{few} \} - \text{No}$$

- ignore.

$$h_3 = \{ \text{many}, \text{big}, \text{no}, \text{expensive}, \text{one} \}.$$

h_3 remains same.

$$x_4 = \{ \text{many}, \text{medium}, \text{no}, \text{expensive}, \text{many} \} - \text{No}$$

$$h_4 = \{ \text{many}, \text{big}, \text{no}, \text{expensive}, \text{one} \}$$

$$x_5 = \{ \text{many}, \text{small}, \text{no}, \text{affordable}, \text{many} \} - \text{Yes}$$

$$h_5 = \{ \text{many}, ?, \text{no}, ?, ?, ? \}$$

Decision Tree Induction

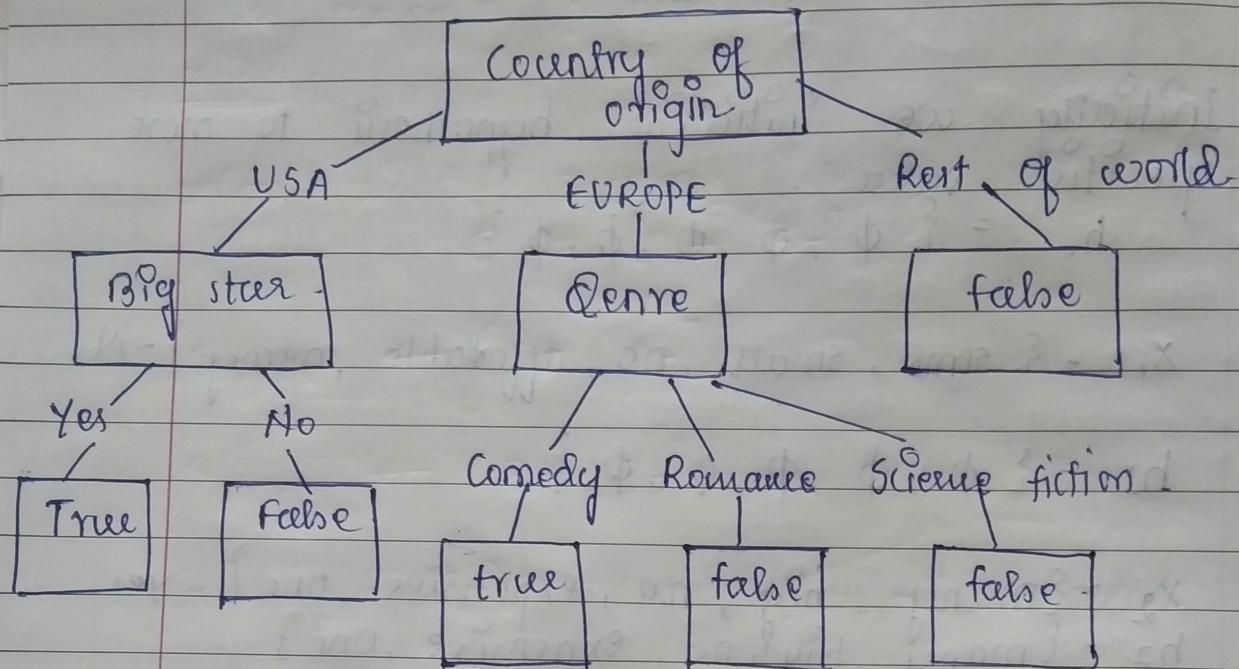


fig : Decision tree determining whether or not film will be by office success.

Entropy & Information gain

| Day | Outlook | Temp | Humidity | Wind | Play Tennis |
|-----------------|----------|------|----------|--------|-------------|
| D ₁ | Sunny | Hot | High | weak | No |
| D ₂ | Sunny | Hot | High | strong | No |
| D ₃ | Overcast | Hot | High | weak | Yes |
| D ₄ | Rain | mild | Normal | weak | Yes |
| D ₅ | Rain | cool | Normal | weak | Yes |
| D ₆ | Rain | Cool | Normal | strong | No |
| D ₇ | Overcast | Cool | High | strong | Yes |
| D ₈ | Sunny | mild | Normal | weak | No |
| D ₉ | Sunny | cool | Normal | weak | Yes |
| D ₁₀ | Rain | mild | Normal | weak | Yes |
| D ₁₁ | Sunny | mild | Normal | strong | Yes |
| D ₁₂ | overcast | mild | High | strong | Yes |
| D ₁₃ | overcast | Hot | Normal | weak | Yes |
| D ₁₄ | Rain | mild | High | strong | No |

Entropy : Measures the impurity of collection of examples.

S : Collection of training examples

P₊ : Proportion of +ve example in S

P₋ : Proportion of -ve " " S

$$\text{Entropy}(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

Characteristics of entropy - Some examples

$$\text{Entropy}([14+, 0-]) = -\frac{14}{14} \log_2 \left(\frac{14}{14}\right) - 0 \log_2 0 = 0$$

$$\text{Entropy}([9+, 5-]) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \left(\frac{5}{14}\right)$$

$$\begin{aligned} \text{Entropy}([7+, 7-]) &= -\frac{7}{14} \log_2 \left(\frac{7}{14}\right) - \frac{7}{14} \log_2 \left(\frac{7}{14}\right) \\ &= 1 \end{aligned}$$

Information Gain.

It measures the expected reduction in entropy caused by partitioning the examples according to an attribute.

Information gain (S, n)

$$\text{Gain}(S, n) = \text{Entropy}(S) - \sum_{v \in \text{values}(n)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

Example :

To compute information gain of attribute wind.

values (wind) = weak, strong
 $S = [9+, 5-]$

$S_{\text{weak}} \leftarrow [6+, 2-]$

$S_{\text{strong}} \leftarrow [3+, 3-]$

$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \sum_{v \in \{\text{weak}, \text{strong}\}} \frac{|S_v|}{|S|} \text{entropy}(S_v)$

$\text{Gain}(S, \text{wind}) = \text{Entropy}(S) - \sum_{v \in \{\text{weak}, \text{strong}\}} \frac{|S_v|}{|S|} \text{entropy}(S_v)$

$$= \text{entropy}(S) - \left(\frac{8}{14}\right) \text{entropy}(S_{\text{weak}}) - \left(\frac{6}{14}\right) \text{entropy}(S_{\text{strong}})$$

$$= (0.940) - \left(\frac{8}{14}\right) (0.811) - \left(\frac{6}{14}\right) (1.00)$$

$$= 0.048 \rightarrow \text{wind.}$$

Nearest Neighbour algorithm.

Eg: query $x = (\text{maths} = 6, \text{cs} = 8)$, $k = 3$

| | maths | cs | Result |
|---|-------|----|--------|
| 1 | 4 | 3 | Fail |
| 2 | 6 | 7 | Pass |
| 3 | 7 | 8 | Pass |
| 4 | 5 | 5 | Fail |
| 5 | 8 | 8 | Pass |

To compute we use euclidean distance

$$d = \sqrt{(x_{01} - x_{A1})^2 + (x_{02} - x_{A2})^2}$$

$$= \sqrt{(6-4)^2 + (8-3)^2} = \sqrt{29}$$

$$= 5.38$$

They store the training data & use to determine a classification for each new piece of data as it is encountered.

Learning Neural Networks

An artificial neural network is the network of simple processing nodes modeled on the human brain which achieves its power through connectivity between its neurons. Each neuron can either fire or not fire by combining millions of neurons together.

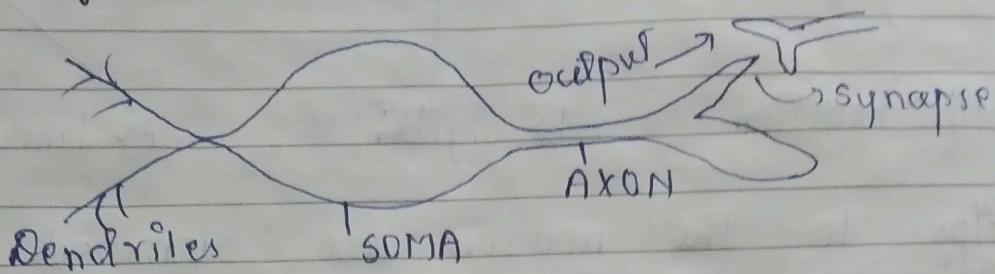
Supervised learning networks presented in three classified training data.

Unsupervised learning methods learn without any human intervention. This method is useful in situations where data needs to be classified / clustered into set of classifications but classifications are not known in advance.

A system that uses reinforcement learning a positive reinforcement given when it performs correctly & negative reinforcement is given when it performs incorrectly.

Ex: Robotic agent - learnt by reinforcement learning about how to pick up an object.

i) Neurons ii) Biological Neuron



human brain - neuron
connected to other neurons - synapse
SOMA

Neuron receives input from other neurons along its dendrites & when this IP signal exceeds a certain ~~feels~~ ^{threshold} then neuron fires a chemical reaction occurs called as action potential to be sent on to axon.

Artificial Neurons

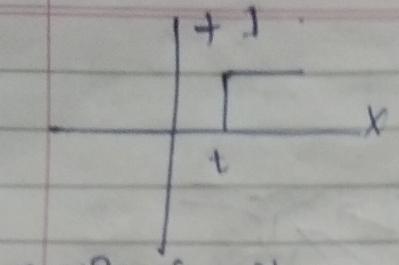
Each node in neural net → Receives a number of inputs.

Activation function are applied to the IP values.

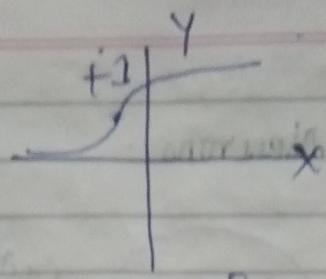
Activation level of neuron is IP value of neuron.

Behaviour of neuron.

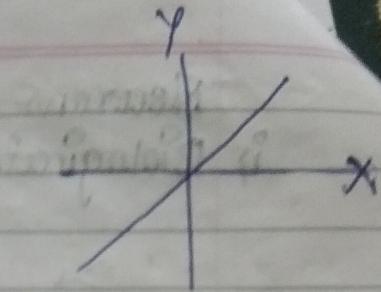
$$x = \sum_{i=1}^n w_i x_i$$



a) step function



b) sigmoid



c) linear function

Inputs to the neuron are summed & This sum is compared with the threshold if sum is greater than threshold the neuron fires & has an activation level of +1 otherwise it is inactive & has an activation level of zero.

$$X = \sum_{i=1}^n w_i x_i$$

Weights

$$w_1 = 0.8$$

$$w_2 = 0.4$$

Inputs

$$x_1 = 0.7$$

$$x_2 = 0.9$$

Summed weight

$$(0.8 \times 0.7) + 0.4 \times 0.9 = 0.92$$

Activation level.

$$Y = \begin{cases} +1 & \text{for } x > t \\ 0 & \text{for } x \leq t \end{cases}$$

Perceptions

JMP

Perception can have any number of inputs which are arranged into a grid. A grid can be used to represent an image or field of vision so perception can be used to perform some image classification or some recognition tasks.

Perception uses step functions.

- +] If weighted sum of i/p $x > \text{threshold } t$
-] If $x \leq t$.

$$x = \sum_{i=1}^n w_i x_i$$

$$y = \begin{cases} +1 & \text{for } x > t \\ 0 & \text{for } x \leq t \end{cases}$$

Function:

$$\text{step}(x) = \begin{cases} +1 & \text{for } x > t \\ 0 & \text{for } x \leq t \end{cases}$$

Activation function for perception.

$$y = \text{step} \left(\sum_{i=0}^n w_i x_i \right)$$

Learning process for perception.

Random weights are assigned between -0.5 & +0.5.

Training data is presented to the perception & it's o/p classification is observed. If o/p is incorrect the weights are adjusted to try to more closely classify this i/p.

If position of training data as negative, then weights need to be modify to increase the o/p for set of I/p's.

$$w_i \leftarrow w_i + \alpha x_i x_i$$

Eg: Logical OR function.

| Epoch | x_1 | x_2 | Expected Y | Actual Y | Error | w_1 | w_2 |
|-------|-------|-------|------------|----------|-------|-------|-------|
| 1 | 0 | 0 | 0 | 0 | 0 | -0.2 | 0.4 |
| 1 | 0 | 1 | 1 | 1 | 0 | -0.2 | 0.4 |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0.4 |
| 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0.4 |
| 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0.4 |
| 2 | 0 | 1 | 1 | 1 | 0 | 0 | 0.4 |
| 2 | 1 | 0 | 1 | 0 | 1 | 0.2 | 0.4 |
| 2 | 1 | 1 | 1 | 1 | 0 | 0.2 | 0.4 |

Logical OR Function for 2 I/p's t=0 & learning rate 0.2.

$$\alpha = 0.2$$

$$w_1 = -0.2 \quad w_2 = 0.4$$

$$x_1 = 0, x_2 = 0$$

∴ expected o/p is $x_1 \vee x_2 = 0$.

we apply formulae for Y.

$$Y = \text{step} \left(\sum_{i=0}^n w_i x_i \right)$$

$$= \text{step} (0x - 0.2) + (0x 0.4)$$

$$= 0$$

O/p Y is expected so weights do not change.

$$\omega_1 = -0.2$$

$$x_1 = 0$$

$$\omega_2 = 0.4$$

$$x_2 = 1$$

∴ expected o/p is $x_1 \vee x_2 = 1$.

we apply formula for Y

$$Y = \text{step} \left(\sum_{i=0}^n \omega_i x_i \right)$$

$$= \text{step} (0 \times (-0.2) + (0 \times (0.4))) \\ = 0.4$$

for $\omega_1 = -0.2 \quad \omega_2 = 0.4$

$$x_1 = 1 \quad x_2 = 0$$

∴ expected o/p is $x_1 \vee x_2 = 1$.

we apply formula for Y

$$Y = \text{step} \left(\sum_{i=0}^n \omega_i x_i \right)$$

$$= \text{step} (1 \times (-0.2) + (0 \times (0.4))) \\ = 0$$

This is incorrect $1 \vee 0 = 1$.

weights are adjusted.

so, assign new values to weights.

$$\omega_1 \leftarrow \omega_1 + (\alpha \times x_1 \times e)$$

$$\omega_1 = -0.2 + (0.2 \times 1 \times 1)$$

$$= -0.2 + 0.2$$

$$= 0$$

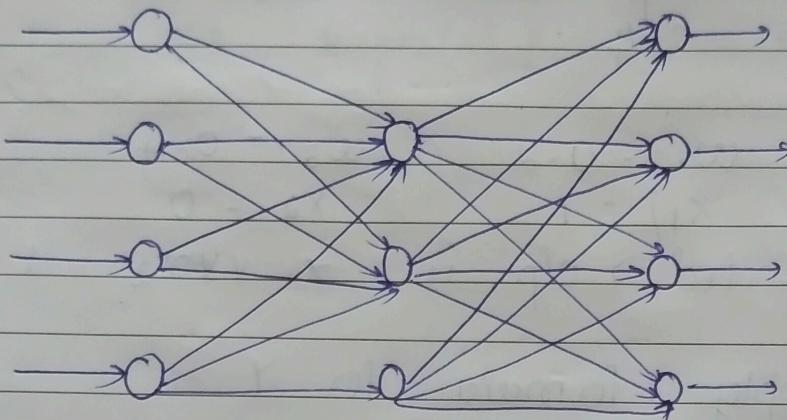
$$\omega_2 = 0.4$$

$$0.4 + (0.2 \times 0 \times 1)$$

for $x_1 = 1$ $x_2 = 1$

$$\begin{aligned}Y &= \text{step}(0x_1) + (0.4x_1) \\&= \text{step}(0) + (0.4) \\&= \text{step}(0.4) \\&= 1\end{aligned}$$

Multilayer Neural Network



Three layer feed forward neural net

Neural network consist of number of neurons that are connected together usually arranged in layers.

The above architectural is neural network also called as feed forward consist of 3 layers

First layer is Input layer.

In this layer it receives a single i/p signal;

If passes i/p signals on to nodes in the next layer which is called as middle layer.

The o/p of each node in this layer is passed through each node in final layer which is o/p layer.

IF carries final stage of processing & sends o/p sigles

Back propagation

multilayer n/w = each neuron has weights

when error is made, greater number of weights to be adjusted.

Recurrent Network

A recurrent network can have connection from that go backward from o/p node to i/p node

Recurrent n/w internal state can alter as set of i/p data presented to it & it can said to have a memory

Ex: Recurrent n/w could be used to predict the stock market price of particular stock based on all previous values.

Eg: Hopfield Networks

Form of recurrent NW.

Activation function used by Hopfield NW.

$$\text{sign}(x) = \begin{cases} +1 & \text{for } x > 0 \\ -1 & \text{for } x < 0 \end{cases}$$

weights of NW.

$$W = \sum_{i=1}^n x_i^o x_i^t - NI$$

Eg: Single layer Hopfield NW with 5 nodes & 3 training ips.

$$x_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad x_2 = \begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{bmatrix}, \quad x_3 = \begin{bmatrix} 1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

The 3 training states are used to predict the weight matrix will be stable state for the matrix.

We can test this by determining the output vector.

$$y_i = \text{Sign}(wx_i - \theta)$$

θ is threshold matrix which contains threshold for each of 5 inputs.

to

$$W = \sum_{i=1}^n X_i X_i^t - NI$$

$$= X_1 X_1^t + X_2 X_2^t + X_3 X_3^t - 3I$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \end{bmatrix}^t - I \begin{bmatrix} -1 & -1 & -1 & -1 & -1 \end{bmatrix}$$

$$+ \begin{bmatrix} 1 \\ -1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 & -1 & 1 & 1 & 1 \end{bmatrix}^t - 3I$$

$$= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} - 3I$$

$$= \begin{bmatrix} 3 & 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 & 3 \\ 3 & 3 & 3 & 3 & 3 \end{bmatrix} - \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 3 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 3 & 0 \\ 0 & 0 & 0 & 0 & 3 \end{bmatrix} - \begin{bmatrix} 0 & 1 & 3 & 3 & 1 \\ 1 & 0 & 1 & 1 & 3 \\ 3 & 1 & 0 & 3 & 1 \\ 3 & 1 & 3 & 0 & 1 \\ 1 & 3 & 1 & 1 & 6 \end{bmatrix}$$

$$Y_1 = \text{sign} \begin{bmatrix} 0 & 1 & 3 & 3 & 1 \\ 1 & 0 & 1 & 1 & 3 \\ 3 & 1 & 0 & 3 & 1 \\ 3 & 1 & 3 & 0 & 1 \\ 1 & 3 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \text{sign} \begin{bmatrix} 8 \\ 6 \\ 8 \\ 8 \\ 6 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = X_1$$

$$Y_2 = \text{sign} \begin{bmatrix} 0 & 1 & 3 & 3 & 1 & -1 & 0 \\ 1 & 0 & 1 & 1 & 3 & -1 & 0 \\ 3 & 1 & 0 & 3 & 1 & -1 & 0 \\ 3 & 1 & 3 & 0 & 1 & -1 & 0 \\ 1 & 3 & 1 & 1 & 0 & -1 & 0 \end{bmatrix}$$

$$= \text{sign} \begin{bmatrix} 8-8 \\ 6-6 \\ 8-8 \\ 2-8 \\ 4-6 \end{bmatrix}$$

$$\begin{bmatrix} -1 \\ -1 \\ -1 \\ -1 \\ -1 \end{bmatrix} = X_{21}$$

$$Y_4 = \text{sign} \begin{bmatrix} 0 & 1 & 3 & 3 & 1 \\ 1 & 0 & 1 & 1 & 3 \\ 3 & 1 & 0 & 3 & 1 \\ 3 & 1 & 3 & 0 & 1 \\ 1 & 3 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ -1 \\ 1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \text{sign} \begin{bmatrix} 2 \\ 2 \\ 2 \\ -4 \\ 2 \end{bmatrix}$$

$$= \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \end{bmatrix}$$

Not in one of the attractors we need to apply rule again.

$$Y_5 = \text{sign} \begin{bmatrix} 0 & 1 & 3 & 3 & 1 \\ 1 & 0 & 1 & 1 & 3 \\ 3 & 1 & 0 & 3 & 1 \\ 3 & 1 & 3 & 0 & 1 \\ 1 & 3 & 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ -1 \\ 1 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$= \text{sign} \begin{bmatrix} 2 \\ 4 \\ 2 \\ 8 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = X_1$$

The Hopfield net involves 3 stages.

stage 1: The net is trained to learn the attractor state.

This can be thought of as a storage / memorization stage.

stage 2: Involves testing in net by providing the attractor states as ip & checking that o/p are identical.

stage 3: Involves using the net which is required to retrieve the data from its memory.

Unsupervised learning network.

Kohonen maps.

- 1) Self organizing feature map.
- 2) Trained using unsupervised learning.
- 3) Used for clustering data.

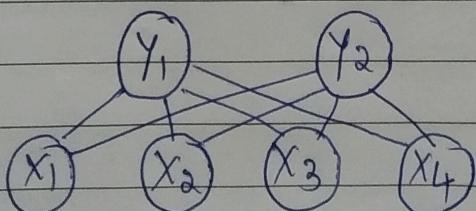
step 1: Initialize weights.

$$w_{ij}^{(0)}$$

no. of ip no. of cluster.

step 2: Compute the square of euclidean distance.

$$D(j) = \sum_{i=1}^n \sum_{j=1}^m (x_i - w_{ij}^{(0)})^2$$



$$m=2$$

$$j = 1 \text{ to } 2$$

step 3: Find winning index j , $D(j)$ is minimum

step 4: Calculate new weights

$$w_{ij}^{\text{new}} = w_{ij}^{\text{old}} + \alpha(x_i - w_{ij})$$

step 5: Update learning rate α using the formula
 $\alpha(t+1)$

step 6: After no. of epochs test for stopping condition of new

$$w_0 = 1.1234$$

$$w_0 = 0.1234$$

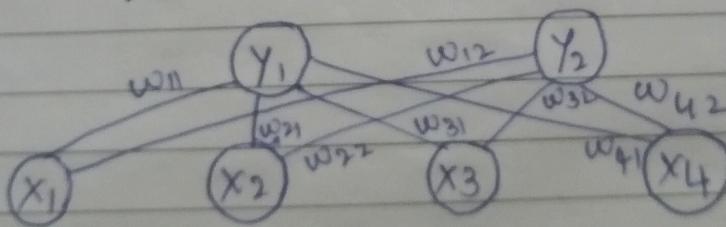
$$w_0 = 0.1233$$

Eg: $x \rightarrow [0 \ 0 \ 1 \ 1]$
 $y_1 \rightarrow [1 \ 0 \ 0 \ 0]$
 $y_2 \rightarrow [0 \ 1 \ 1 \ 0]$
 $u \rightarrow [0 \ 0 \ 0 \ 1]$

Number of clusters to form 2.

Assume

→ No. of I/p vectors $n = 4$
no. of clusters $m = 2$.



Initialize weight randomly b/w 0 & 1

$$w_{\cdot j} = \begin{bmatrix} 0.2 & 0.9 \\ 0.4 & 0.7 \\ 0.6 & 0.5 \\ 0.8 & 0.3 \\ w_{j1} & w_{j2} \end{bmatrix}$$

weights going for
cluster 1 → cluster 2

First: Input vector $x = [0 \ 0 \ 1 \ 1]$

Calculate euclidean distance

$$D(j) = \sum_{i=1}^n \sum_{j=1}^m (x_i - w_{ij})^2$$

$$\begin{aligned} D(1) &= (x_1 - w_{11})^2 + (x_2 - w_{21})^2 + (x_3 - w_{31})^2 + (x_4 - w_{41})^2 \\ &= (0 - 0.2)^2 + (0 - 0.4)^2 + (1 - 0.6)^2 + (1 - 0.8)^2 \\ &= 0.04 + 0.16 + 0.16 + 0.04 \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} D(2) &= (x_1 - w_{12})^2 + (x_2 - w_{22})^2 + (x_3 - w_{32})^2 + (x_4 - w_{42})^2 \\ &= (0 - 0.9)^2 + (0 - 0.7)^2 + (1 - 0.5)^2 + (1 - 0.3)^2 \\ &= 2.04 \end{aligned}$$

$$D_1 < D_2$$

$$\omega_{ij}^{(0)} = \omega_{ij}^{(0)} + \alpha (x_i^{(0)} - \omega_{ij}^{(0)}) .$$

$$\omega_{11} = \omega_{11} + 0.5 (0 - 0.2) = 0.1$$

$$\begin{aligned}\omega_{21} &= \omega_{21} + \alpha (x_2 - \omega_{21}) \\ &= 0.2 .\end{aligned}$$

$$\omega_{31} = 0.8$$

$$\omega_{41} = 0.9 .$$

$$\omega_{ij}^{(0)} = \begin{bmatrix} 0.1 & 0.9 \\ 0.2 & 0.7 \\ 0.8 & 0.5 \\ 0.9 & 0.3 \end{bmatrix}$$

Second P/p vector