

Neural Network

DATE: _____
PAGE: _____

UNIT 5

paradigms of Associative memory

Mathematic pattern

pattern term is used in context of neural networks to mean a set of activations across a pool of units (neurons)

ANNI combines biological principles with advanced statistics to solve problems in domains such as pattern recognition and game-play.

pattern →

Pattern is the quantitative description of an object, event.

Pattern can be classified as spatial or temporal pattern.

i) Spatial :- Example are pictures, weather, map, video, games

ii) Temporal :- Example are speech signal, ECG etc. These pattern involve the ordered sequence of data appeared in time.

* pattern classification:-

The main goal of pattern classification is to assign a physical object, event or phenomenon to one of the pre-specified

Good Write Classes

for example.

we classify various objects into different categories like living, non-living, thing etc.

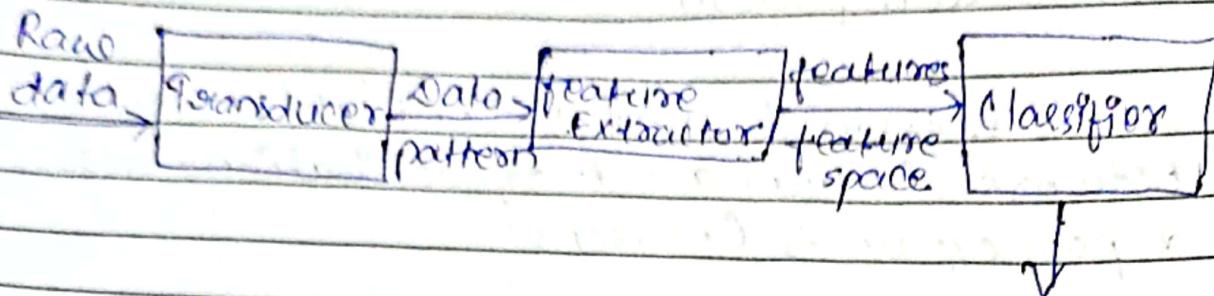


fig:- Block diagram of classification system.

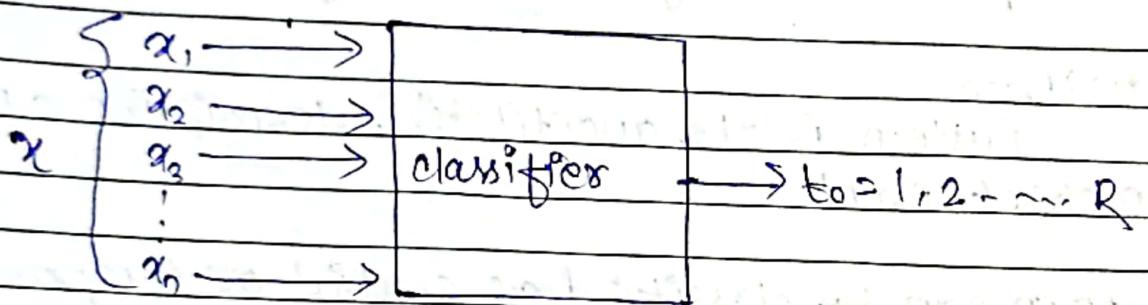


fig:- General diagram of pattern classifier.

→ Transducer generally converts one form of energy to another form of energy.

Here the input is raw data space
 $\alpha = [\alpha_1, \alpha_2, \dots, \alpha_n]$ if O/P is transducer
 is dia in pattern space that belong to
 a certain category that is

$$t_0 = 1, 2, \dots, R$$

converted data at the O/P can be compressed. Compressed data is called feature.

→ classifier at the end, convert or map the feature to the particular class.

General concept of Associative Memory

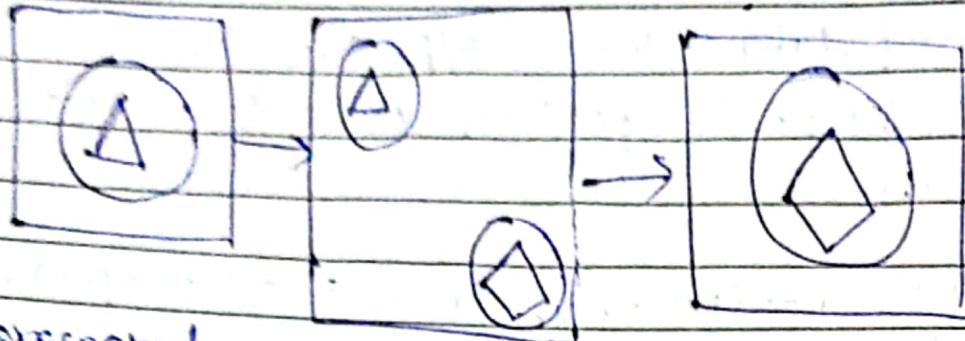
It is a storehouse of associated patterns which are encoded in some form.

- An associative memory is a content-addressable structure that maps a set of input patterns to a set of O/P patterns.
- content addressable & memory structure that enables the recollection of data based on the intensity of similarity between the input pattern and the pattern stored in the memory. This kind of memory is robust & fault-tolerant.

They are of two types

i) Auto-associative memory

It recovers a previously stored pattern that most closely relates to the current pattern. auto associative correlator.



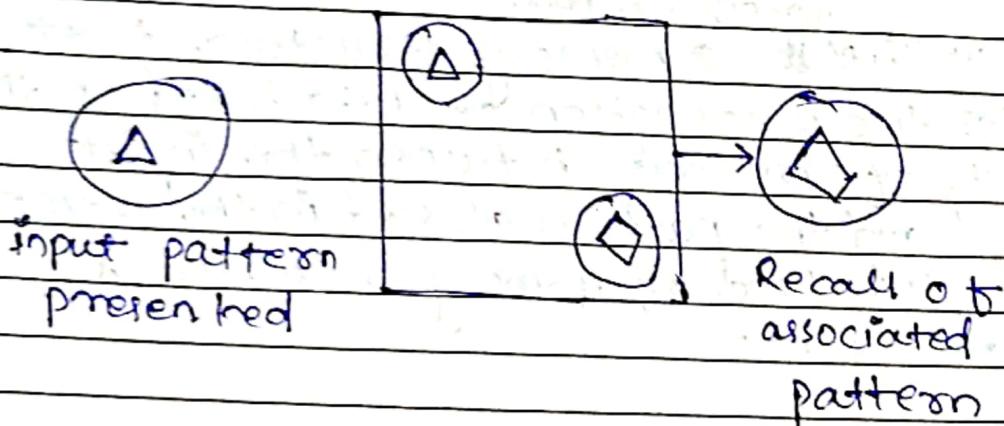
Presented
distorted
pattern.

Recall of
perfect
pattern.

Hetero-associative memory:-

The recovered pattern is generally different from the input pattern not only in type and format but also in content.

Also known as hetero-associative
correlators.

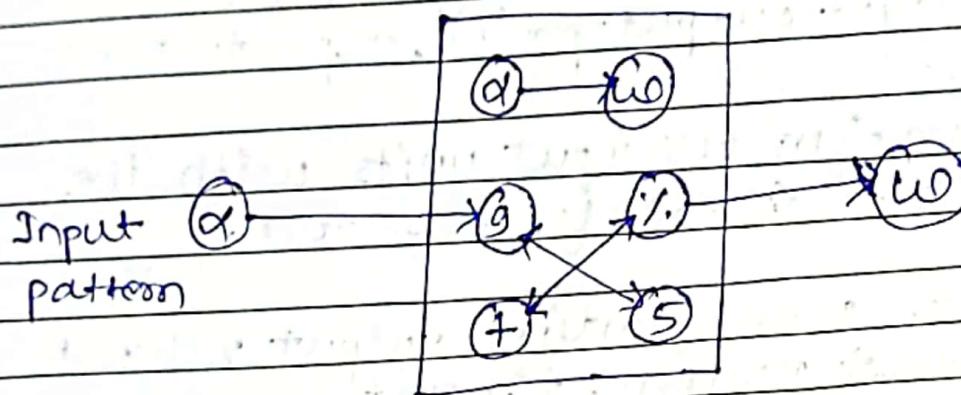


Neural Associative Memory are usually used to implement these memory model.

These models follows distinct architecture to memorize data.

Working of Associative Memory

Associative memory is a depository of associated pattern which in some form. If the depository is triggered with a pattern the associated pattern pair appear at the output. The input could be an exact or partial representation of a stored pattern.



If the memory is produced with an input pattern; may say α , the associated pattern w is recovered automatically

- Encoding or memorization.
- Errors & noise (performing highly parallel distributed computation)
- Performance Measure
- (Memory capacity & content-addressability mutually orthogonal)

Hebbian Learning:-

It was proposed by Donald O Hebb. It is one of the first and also easiest learning rule in the neural network used for pattern classification. It's a single layer NN having one input layer with n units and one output layer with m units.

single unit.

The rule works by updating the weights between neurons in the neural network for each training sample.

Hebbian Learning Rule Algol-

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 step 3 to 5.

3. Set activations for input units with the input vectors $x_i = s_i$ for $i = 1$ to n .

4. set the corresponding output value to the output neuron; i.e. $y = t$

5. Update weight and bias by applying Hebb rule for all $i = 1$ to n .

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

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

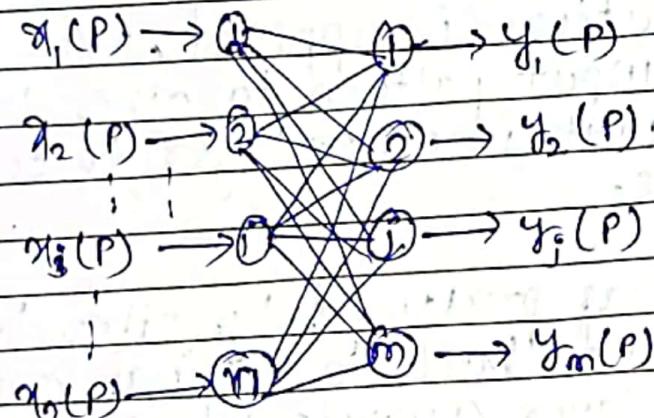
Practice Numerical for it.

Bidirectional Associative Memory

BAM is a hetero associative, content addressable memory consists of 2 layers. For an input pattern, it returns another pattern which is potentially of a different size.

BAM Architecture

When BAM accepts an input of n dimensional vector X from set A then the model recalls m -dimensional vector Y from set B. Similarly when Y is treated as input, the BAM recall X .



⑥ forward direction.

BAM model can perform both auto & hetero associative recall of stored information

Weight matrix can be calculated as

$$W_k = \mathbf{x}_k^T \mathbf{y}_k \quad (\text{to store a single associated patterns.})$$

$W = \sum_{k=1}^K (x_k, y_k)$ (to store simultaneously several associated pattern pair)

This process is generally called storage or encoding.

After encoding Nue can be used for decoding.

The decoding process is as follows

- ① Input layers can be applied, either on x layer or on the y layer depending on direction of propagation.
- ② When input pattern is applied, Nue will propagate the input pattern to other layer allowing the units in other layer to compute their O/P value.
- ③ Pattern that was produced by other layer is then propagated back to original layer & in original layer compute the O/P value.
- ④ This process is repeated until further propagation & computation do not result in a change in the associations of units in both layer where final pattern pair is one of the stored associated pattern pair.
- ⑤ Final pattern pair that will be produced in the Nue depend on initial pattern pair.

and connection weight matrix.

* Addition & Deletion of pattern pairs.

- Given a set of pattern pair (x_i, y_i) for $i = 1, 2, \dots, n$ & set of correlation matrix M .

- a new pair (x', y') can be added.
- an existing pair (x_j, y_j) can be deleted from memory model.

Addition

Add a new pair (x', y') to existing correlation matrix M , then new matrix

M_{new} is given by

$$M_{\text{new}} = X_1^T Y_1 + X_2^T Y_2 + \dots + X_n^T Y_n + x'^T y'$$

Deletion

Subtract the matrix corresponding to an existing pair (x_j, y_j) from matrix M then new Correlation Matrix is given by

$$M_{\text{new}} = M - (X_j^T Y_j)$$

- The addition & deletion of information is similar to the functioning of the system as a human memory exhibiting learning & forgetfulness.

Good Write

* Discrete BAM:-

We propagate an input pattern X to Y layer and Y layer compute the net input using

Input is calculated using

$$y_i = \sum_{j=1}^m x_j w_{ij} \quad i = 1, 2, \dots, m$$

O/P

$$y_i(t+1) = \begin{cases} S+1 & \text{if } y_i > 0 \\ y_i(t) & \text{if } y_i = 0 \\ -1 & \text{if } y_i < 0 \end{cases}$$

pattern produced on Y layer is then propagated back to X layer & X layer compute values.

Input

$$x_i = \sum_{j=1}^n y_j w_{ij}$$

O/P

$$x_i(t+1) = \begin{cases} S+1 & \text{if } x_i > 0 \\ x_i(t) & \text{if } x_i = 0 \\ -1 & \text{if } x_i < 0 \end{cases}$$

Energy function for OBAM

$$E = - \sum_{i=1}^m \sum_{j=1}^n a_i w_{ij} y_j$$

* continuous BAM

In this, the units use hyperbolic tangent O/P function. The units in X-layer have an extra external input I_i while units in Y-layer have an extra external input J_j .

$$\text{for } i = 1, 2, \dots, m \text{ & } j = 1, 2, \dots, m$$

These extra input lead to a modification in the computation of net input

$$x_i = \sum_{j=1}^m y_j w_{ij} + I_i$$

$$y_i = \sum_{j=1}^m a_i w_{ij} + J_i$$

Limitations of Bidirectional Associative Memory

- Storage capacity of the BAM: stored number of associations should not exceed the number of neurons.

- Incorrect convergence

Always the closest association may not be produced by BAM.

Back-propagation

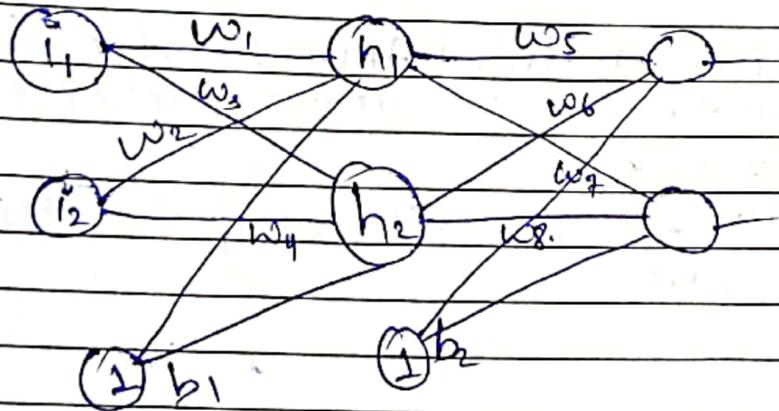
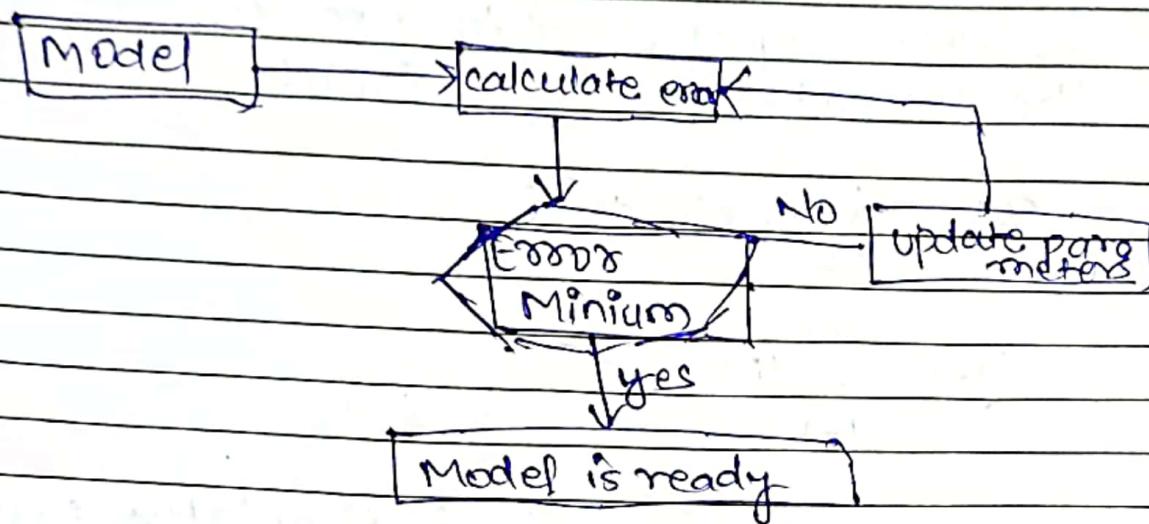
Backpropagation Network is multi-layer feed-forward ANN with diff. transfer function in ANN.

It is bind of gradient descent technique with backward error propagation.

It's a supervised learning and an implementation of delta learning rule.

gradient descent (differentiation)

gradient descent method is used to update the weight by reducing the error function



Credit Assignment Problem

A neural network faces the "credit assignment problem" in situations in which only incomplete performance evaluations are available. The credit assignment problem is that a network should assign credit or blame for its behaviour according to the contribution to the network performance.

If any distributed system or any distributed system network is there then it is a very common practice to give some credits or blame to each and every internal decision so that it can be reflected on the overall outcomes of the network.

→ Internal decisions are dependent on many actions. Among such, some actions are very prominent to get the outcome of internal decision. In this case credit assignment is divided into two parts:-

i) Temporal:- If we have given credit or blame for outcomes of the internal decision then this is called temporal credit assignment problem.

ii) Structural:-

If we have given credit or blame for actions of the internal decision, then it is called structural credit assignment.

temporal :- to get the selected outcomes of the internal decisions.

⇒ structural :- to get the selected actions of the internal decision.

the temporal & structural, now we

can think of some more information about

the data base management system in

order to have it as complete as

so as to obtain the following features of the

data base management system like as

(a) data base organization

and function of DBMS will be to store

the data in a structured form where as

the data can be easily retrieved, modified

and deleted from the data base system

so for example if we want to change

the data in the data base system then we

can do it by changing the data in the

data base system so that we can get

the required data in the required form

and we can also get the required data

in the required form by using the

function of the data base system

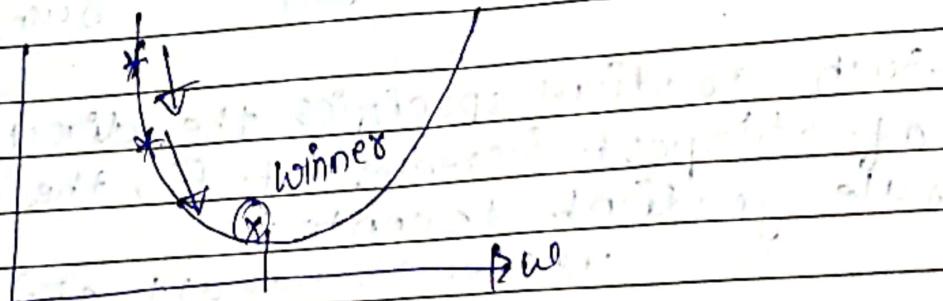
Generalized delta rule.

Linearly separable and non-linearly separable
perception rules fails for non-linearly
separable data. And here comes the
delta rule.

Here, delta rule tries to converge
the solution towards a best-fit approximation
to the target concept.

The key idea behind the delta rule is to
use gradient descent to search the hypothesis
space of possible weight vectors to find
the weights that best fit the training examples.

Gradient descent



$$O(\vec{x}) = \vec{w} \cdot \vec{x}$$

way to define error to derive a
weight learning rule for linear units.

$$E(\vec{w}) = \frac{1}{2} \sum (t_d - O_d)^2$$

t_d = desired output

O_d = actual output

\mathcal{D} is the set of training example

- * How can we calculate the direction of steepest descent along the error surface.

The direction can be found by computing the derivative of E with respect to the each component of the vector \vec{w} .

The vector derivative is called gradient of E with respect to (\vec{w})

$$\Delta E(\vec{w})$$

$$\Delta E(\vec{w}) = \left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Such gradient specifies the direction of steepest increase of E , the training rule gradient decent is

$$w_{(new)} = w_{(old)} + \Delta w$$

$$\text{where } \Delta w = -\eta \Delta E(\vec{w})$$

η = learning rate, it determines the step size in the gradient descent search.

$\Delta E(\vec{w})$ = derivative of error with respect to each component of weight vector.

The negative sign is present because we want to move the weight vector in the direction that decrease E(error)

$$\Delta E(\vec{w}) = \frac{dE}{d w_i} \quad \text{where } i \text{ is for respective weights}$$

We have to get the value of the given derivative term: $\frac{dE}{d w_i}$

$$E(\vec{w}) = \frac{1}{2} \sum dE$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \sum (t_d - o_d)^2$$

$$= \frac{1}{2} \sum \frac{\partial (t_d - o_d)^2}{\partial w_i}$$

$$= \frac{1}{2} \sum (t_d - o_d) \cdot \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= \frac{1}{2} \sum 2(t_d - o_d) \cdot \frac{d}{d w_i} (t_d - o_d)$$

$$= \sum (t_d - o_d) (-x_{id})$$

$$= - \sum (t_d - o_d) (\vec{x}_{id})$$

$$\Delta w_{ij} = n \sum_d (t_d - o_d) \cdot x_{id}$$

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

Back propagating Algorithm.

- To derive the equation for weight updation, stochastic gradient descent rule is used.
- For each training example d , every weight w_{ji} is updated by adding Δw_{ij}

$$w_{ij} \leftarrow w_{ij} + \Delta w_{ij}$$

$$\Delta w_{ij} = -n d E_j$$

$$\Delta w_{ij}$$

Firstly Error is calculated.

$$(Error) = E_d(\vec{w}) = \frac{1}{2} \sum (t_k - o_k)$$

Notations used:

x_{ji} = the i th input to unit j

w_{ji} = the weight associated with the i th input to unit j

$$\text{net}_j = \sum w_{ij} x_{ji}$$

o_j = the output computed by unit j .

t_j = the target output for unit j .

σ = the sigmoid function.

outputs = the set of units in the final layer of the network

$\text{Downstream}(j)$ = the set of units whose immediate inputs include the output of unit j .

To begin, notice that weight w_{ij} can influence the rest of the network only through net_j . Therefore applying chain rule as,

$$\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial w_{ij}} + \frac{\partial E_d}{\partial \text{net}_j} \times \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$$\text{time } \therefore \frac{\partial E_d}{\partial \text{net}_j} \times \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$$\Delta w_{ij} = -\eta \frac{\partial E_d}{\partial \text{net}_j} x_{ji}$$

$$\text{net}_j = \sum w_{ij} x_{ji}$$

o_j = the output computed by unit j .

t_j = the target output for unit j .

σ = the sigmoid function.

outputs = the set of units in the final layer of the network

Downstream (j) = the set of units whose immediate inputs include the output of unit j .

To begin, notice that weight w_{ij} can influence the rest of the network only through net_j . Therefore applying chain rule as,

$$\frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial \text{net}_j} \times \frac{\partial \text{net}_j}{\partial w_{ij}}$$

$$\therefore \frac{\partial E_d}{\partial w_{ij}} = \frac{\partial E_d}{\partial \text{net}_j} \times x_{ji}$$

$$\Delta w_{ij} = -n \frac{\partial E_d}{\partial \text{net}_j} x_{ji}$$

To derive a convenient expression for

$$\frac{\partial E_d}{\partial \theta_{netj}}$$

$$\frac{\partial \theta_{netj}}{\partial \theta_j}$$

We consider two cases in turn:

- Case 1, where unit j is an output unit for the network, and
- Case 2, where unit j is an internal unit of the network

Case 1:-

j as the output unit weight.

- Just as w_{ji} can influence the rest of the network only through θ_{netj} , θ_{netj} can influence the network through θ_j . Therefore, we can invoke the chain rule again to write,

$$\frac{\partial E_d}{\partial \theta_j} = \frac{\partial E_d}{\partial \theta_{netj}} \times \frac{\partial \theta_{netj}}{\partial \theta_j}$$

$$\therefore \frac{\partial \theta_{netj}}{\partial \theta_j} = \frac{\partial \theta_{netj}}{\partial \theta_j}$$

$$\frac{\partial E_d}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \sum_k (t_k - o_k)^2$$

$$= \frac{1}{2} (t_j - o_j) \frac{\partial (t_j - o_j)}{\partial \theta_j}$$

$$= -(t_j - o_j)$$

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial \sigma(\text{net}_j)}{\partial (\text{net}_j)}$$

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

$$= \sigma(\text{net}_j)(1 - \sigma(\text{net}_j))$$

$$= o_j(1 - o_j)$$

$$\boxed{\frac{\partial E_d}{\partial \text{net}_j} = -(t_j - o_j) o_j (1 - o_j)}$$

$$\Delta w_{ji} = \eta \frac{\partial E_d}{\partial w_{ji}}$$

$$= -\eta \frac{\partial E_d}{\partial \text{net}_j} x_{ji}$$

$$\delta_j = (t_j - o_j) o_j (1 - o_j)$$

$$\Delta w_{ji} = \eta \delta_j x_{ji}$$

CASE 2:-

Modify the weights at hidden layer.

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{\text{KE Downstream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \times \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{\text{KE Downstream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{\text{KE Downstream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial o_j} \times \frac{o_j}{\partial \text{net}_j}$$

$$= \sum_{\text{KE Downstream}} -\delta_k w_{kj} \frac{o_j}{\partial \text{net}_j}$$

$$= \sum_{\text{KE Downstream}} -\delta_k w_{kj} o_j (1 - o_j)$$

$$\textcircled{a} \quad \frac{\partial \text{net}_k}{\partial o_j} = \frac{\partial \delta_k w_{kj}}{\partial o_j} = \frac{\partial o_j w_{kj}}{\partial o_j}$$

$$\frac{\partial o_j}{\partial (\text{net}_j)} = o_j (1 - o_j)$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial \text{net}_j} x_{ji}$$

$$\Delta w_{ji} = -\eta o_j (1 - o_j) \leq \delta_k w_{kj} o_j$$

Kolmogorov Theorem

Chapman-Kolmogorov equation was derived independently by both the British mathematician Sydney Chapman and the Russian mathematician Andrey Kolmogorov.

Suppose

$x_n = n = 0, 1, 2, \dots$ is a homogeneous Markov chain then,

$$P_{ij}^{(m+n)} = \sum P_{ik}^m P_{kj}^n$$

P_{ij} = conditional probability that the Markov chain goes from state i to j in $m+n$ steps.

$P_{ik}^{(m)}$ = conditional probabilities that of reaching intermediate state k in m steps.

$P_{kj}^{(n)}$ = conditional probabilities of k reaching to j in steps n .

\Rightarrow Step transition probability:

$$\text{For } m+1, P_{ij}^{(m)} = P(X_{m+1} = j | X_0 = i)$$

That is

$$P(X_{m+1} = j | X_0 = i)$$

$$= \sum_{k=0}^K P(X_m = k | X_0 = i) P(X_{m+1} = j | X_m = k)$$

Proof

The (i, j) th entry of the matrix $P^{(m+n)}$ is given by,

$$P_{ij} = P(X_{m+n} = j | X_0 = i)$$

$$= \sum_k P(X_{m+n} = j, X_m = k | X_0 = i)$$

$$= \sum_k \frac{P(X_{m+n} = j, X_m = k, X_0 = i)}{P(X_0 = i)}$$

$$= \sum_k \frac{P(X_{m+n} = j, X_m = k, X_0 = i)}{P(X_m = k, X_0 = i)} \cdot \frac{P(X_m = k, X_0 = i)}{P(X_0 = i)}$$

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\begin{aligned} P(X_{m+n}=j | X_m=k, X_0=i) &= \sum_k P(X_{m+n}=j | X_m=k, X_0=i) \cdot P(X_m=k | X_0=i) \\ &= \sum_k P(X_{m+n}=j | X_m=k) \cdot P(X_m=k | X_0=i) \end{aligned}$$

by Markov property

Mathematically the markov property is stated as

$$P(X_{t+1}=j | X_t=i, X_{t-1}=i_{t-1}, \dots, X_0=i_0) = P(X_{t+1}=j | X_t=i)$$

for all $t=1, 2, 3, \dots$ and for all states.

The future depends only on the present, ~~not~~ not the past behaviour.

$$= \sum_k p_{kj}^{(n)} p_{ik}^{(m)}$$

$$= \sum_k p_{ik}^{(m)} p_{jk}^{(n)}$$

In the matrix notation, it can be written as,

$$P^{(m+n)} = P^{(m)} P^{(n)} \text{ or } P^{(m+n)} = P^{(m)} P^{(n)}$$

Its differential equation,

$$\dot{P}_{ij}^{(m+n)} = \frac{d}{dn} P_{ij}^{(m+n)}$$

$$= \sum_{k=1}^n P_{ik}^{(m)} \frac{d}{dn} P_{kj}^{(n)}$$

At $n=0$, we have,

$$P_{ij}^{(m)} = \sum_{k=1}^n P_{ik}^{(m)} a_{kj}$$

where $a_{kj} = \frac{d}{dn} P_{kj}^{(n)}$ at $n=0$

Or in matrix notation

$$P'(m) = P(m) A$$

$$\text{ie } \frac{d}{dn} P = PA$$

Let $\{X_n, n \geq 0\}$ be a homogeneous Markov chain with TPM (transition probability matrix) $P = (P_{ij})$ and n -step TPM $P^{(n)} = P^n$ where

$$P_{ij}^{(n)} = P(X_n = j | X_0 = i) \text{ if } P_{ij}^{(n)} = P_{ij}$$

Then the following properties hold:

$$(1) P^{(nm)} = P^n P^m$$

(2) $P^{(n)} = P^n$, i.e., n -steps TPM $P^{(n)}$ is equal to the n^{th} power of the one-step TPM P .

(2) we have to prove by induction

By definition.

we have,

$$P = P^{(1)}$$

thus the assertion $P^{(n)} = P^n$ is true
for $n=1$

Assume $P^{(m)} = P^m$ for some positive integer m .

Then by using result $P^{(m+n)} = P^m P^n$

(ie Chapman-Kolmogorov equation)

We have.

$$P^{(m+1)} = P^m (P^{(1)}) = P^m P = P^{m+1}$$

Thus the assertion $P^{(n)} = P^n$ is also true for $m+1$.

Hence proved.

The TPM of the Markov chain with three states 1, 2, 3 is

$$P = \begin{bmatrix} 0.1 & 0.5 & 0.4 \\ 0.6 & 0.2 & 0.2 \\ 0.3 & 0.4 & 0.3 \end{bmatrix}$$

Compute 2-step TPM

$$P^2 = P \cdot P_1$$

$$= \begin{bmatrix} 0.43 & 0.31 & 0.26 \\ 0.24 & 0.42 & 0.34 \\ 0.36 & 0.35 & 0.29 \end{bmatrix}$$

$$P = A \begin{bmatrix} A & B & C \\ B & C & A \\ C & A & B \end{bmatrix}$$

$$P^2 = P \cdot P$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 0.5 & 0.5 & 0 \\ 0 & 0.5 & 0.5 \end{bmatrix}$$

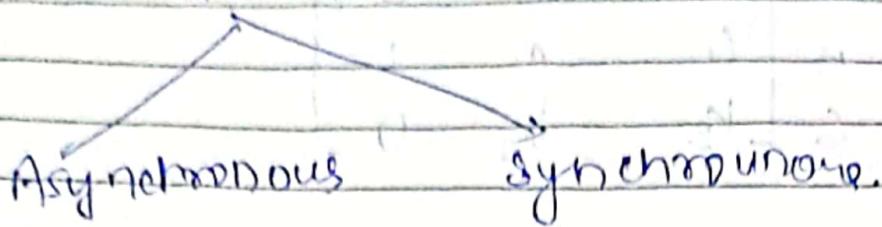
* Recurrent Autoassociative Memory:

In order to produce an improved association, there is a need for suppressing the output noise of memory output in linear associative memory. This can be done by recycling of output to input & by thresholding the o/p.

→ The repetitive process of recycling of the output to the input through the Nues can be performed by a Recurrent Neural Networks (RNN).

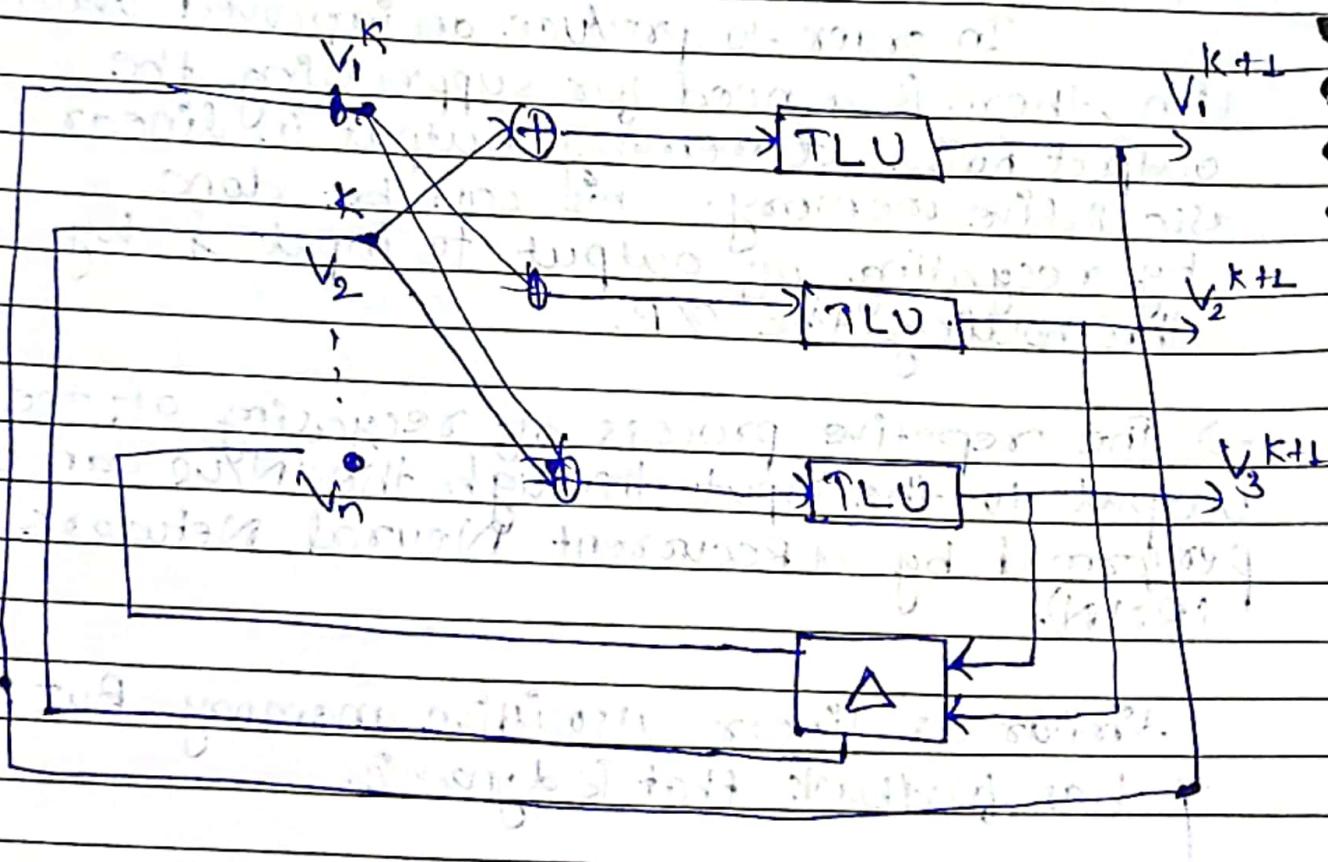
Similar as linear associative memory. But it has feedback that is dynamic.

2 modes.



- ① Under the asynchronous update mode, only one neuron is allowed to compute or change state at a time, and then all output are delayed by time μ . that is produced by unit delay element in the feed back loop.

- ② In synchronous mode.



Explanation of figure: dots represent the neuron TLU threshold logic unit & these are used because associative memories are updated in discrete time.

$\Delta \rightarrow$ unit delay in time.

Storage & Retrieval Algorithm: storage algorithm is called encoding & retrieval algorithm is called decoding process.

Given: p bipolar binary vector $\{S^1, S^2, \dots, S^p\}$ where S^m is $n \times 1$ for $m = 1, 2, \dots, p$

Initializing vector V^0 is $n \times 1$

Storage Algorithm:-

Step 1: weight matrix W is $(m \times n)$:
 $W \leftarrow 0, m \leftarrow 1$

Step 2: vector S^m is stored, therefore the storage algo for calculating the weight matrix is

$$W \leftarrow S^m S^{m(T)} - I$$

$$\text{or, } v_{0ij} = (1 - s_{ij}) \sum_{m=1}^p s_i^{(m)} s_j^{(m)}$$

$(s_{ij} \rightarrow \text{connector function})$

Step 3: If $m < p$ then $m \leftarrow m + 1$ &
go to step 2 otherwise go to
Step 4.

Step 4: Storage is complete, output vector
in.

Retrieval Algorithm:-

Step 1: Cycle counter k is initialized

$K \leftarrow 1$. Within the cycle counter
 i is initialized: $i \leftarrow 1$ & v^o is
initialized as $V \leftarrow V^o$.

Step 2: Integers $1, 2, \dots, n$ are arranged
in an ordered random sequence

$\alpha_1, \alpha_2, \dots, \alpha_n$ if randomized
sequence of updation is not used
then $\alpha_1 = 1, \alpha_2 = 2, \alpha_n = n$ or
 $\alpha_i = i$ at each update cycle.

Step 3: Neuron i is updated by computing
 V_i^{new} , α_i net = $\sum_{j=1}^n w_{ij} v_j$

Step 4: If $i < n$ then $i \leftarrow i + 1$ & go to
Step 3, otherwise go to
Step 5

Step 5: If $V_i^{\text{new}} \alpha_i = V_i$ for $i = 1, 2, \dots, n$

Hebb rule.

Algo:-

Step 1:- Initialize the weights for each input vectors $w_i = 0$ for $i = 0$ to n (number of neurons) and $b = 0$.

Step 2: for each input vector $s_i : t$ (target output pair), repeat steps 3 to 5.

Step 3: set activation for input unit in input vector $a_i = s_i$ for $i = 0$ to n

Step 4: Set the corresponding output value to the output neuron. $y = t$

Step 5: Update the weights & bias using Hebb rule

$$\Delta w = x_i y$$

x_i = Input vector

y = Output Value.

$$\Delta b = y$$

We will repeat the process for each training sample.

And function

Input	Output AND	w_1	w_2	B
x_1	x_2	B	y	$x_1 y_1$
-1	-1	1	-1	1
-1	1	1	-1	0
1	-1	1	-1	0
1	1	1	1	1

Initial weight = 0.5

x_1	x_2	B	y	$x_1 y_1$	$x_2 y_1$	Δb	w_1	w_2	B
1	1	1	1	1	1	1	1	1	1
1	-1	1	-1	-1	1	-1	0	2	0
-1	1	1	-1	1	-1	-1	1	1	-1
-1	-1	1	-1	1	1	-1	2	2	-2

Training with error back propagation

Initial weight = 0.5

Output after 100 iterations = 0.999

Output after 1000 iterations = 0.9999

Output after 10000 iterations = 0.99999

Output after 100000 iterations = 0.999999

Output after 1000000 iterations = 0.9999999

Training with error back propagation

Output after 1000000 iterations = 0.99999999

Good Write

Hopfield

Using Hebb Rule find weight required to perform the following classification of given input pattern '+' symbol represent the value 1 and empty squares indicates -1. Consider '1' belong to the member of class (so have target value 1) & '0' does not belong to the member of class (so has target value -1)

	+	+	+				+	+	+
		+					+	+	+
	+						+	+	+

'1'

'0'

SOL

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	Target	Targ.
1	1	1	1	-1	1	-1	1	1	1	1
0	1	1	1	1	1	-1	1	1	1	-1

Set the initial weight to 0.

i.e $w_1 = w_2 = w_3 = w_4 = w_5 = w_6 = w_7 = w_8 = w_b = 0$

and $b = 0$

Now, updating the weights

$$w_1(\text{new}) = w_1(\text{old}) + x_1 \cdot y = 1$$

$$w_2(\text{new}) = 0 + 1 = 1$$

$$w_3(\text{new}) = 0 + 1 = 1$$

$$w_8(\text{new}) = 0 \cdot 1 = 0$$

$$w_4(\text{new}) = 0 + -1 = -1$$

$$w_9(\text{new}) = 0 \cdot 1 = 0$$

$$w_5(\text{new}) = 0 + (+1) = 1$$

$$b = y = 1$$

$$w_6(\text{new}) = 0 + -1 = -1$$

so new weight vector is,

$$w = [1, 1, 1, -1, 1, -1, 1, 1, 1]$$

$$b = 1$$

Again

$$w_1 \text{ (new)} = 1 + (-1) = 0$$

$$w_2 \text{ (new)} = 1 + (1) = 0$$

$$w_3 \text{ (new)} = 1 + (-1) = 0$$

$$w_4 \text{ (new)} = -1 + (1) = 0$$

$$w_5 \text{ (new)} = 1 + (-1) = 0$$

$$w_6 \text{ (new)} = -1 + (-1) = -2$$

$$w_7 \text{ (new)} = 1 + (-1) = 0$$

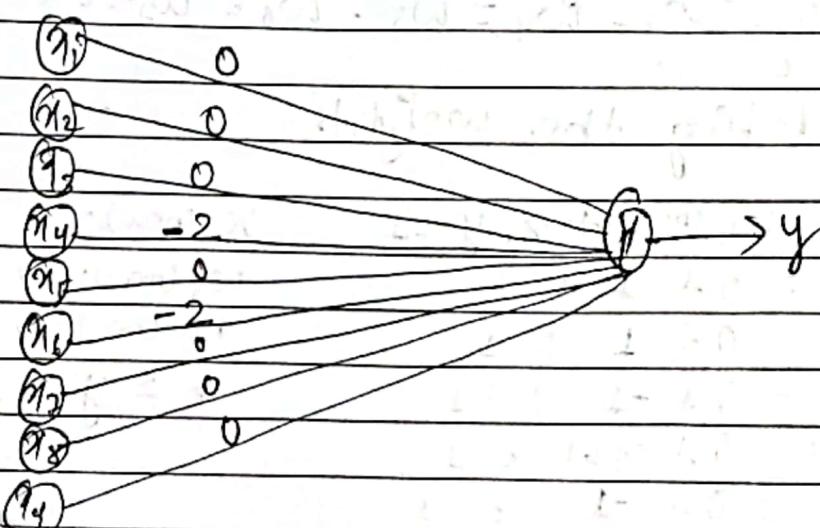
$$w_8 \text{ (new)} = 1 + (-1) = 0$$

$$w_9 \text{ (new)} = 1 + (-1) = 0$$

$$b = \Delta b + y$$

$$= 1 - 1$$

$$= 0$$



Good Write

Hopfield Network

The Hopfield Neural Networks, invented by Dr John J. Hopfield consists of one layer of n fully connected recurrent neurons.

Discrete Hopfield Network

It is fully interconnected neural network where each unit is connected to every other unit. It behaves in a discrete manner i.e. it gives finite distinct output generally of two types.

- Binary (0/1)
- Bipolar (-1/1)

The network has symmetrical weights with no self-connection

i.e.

$$w_{ij} = w_{ji} \text{ & } w_{ii} = 0$$

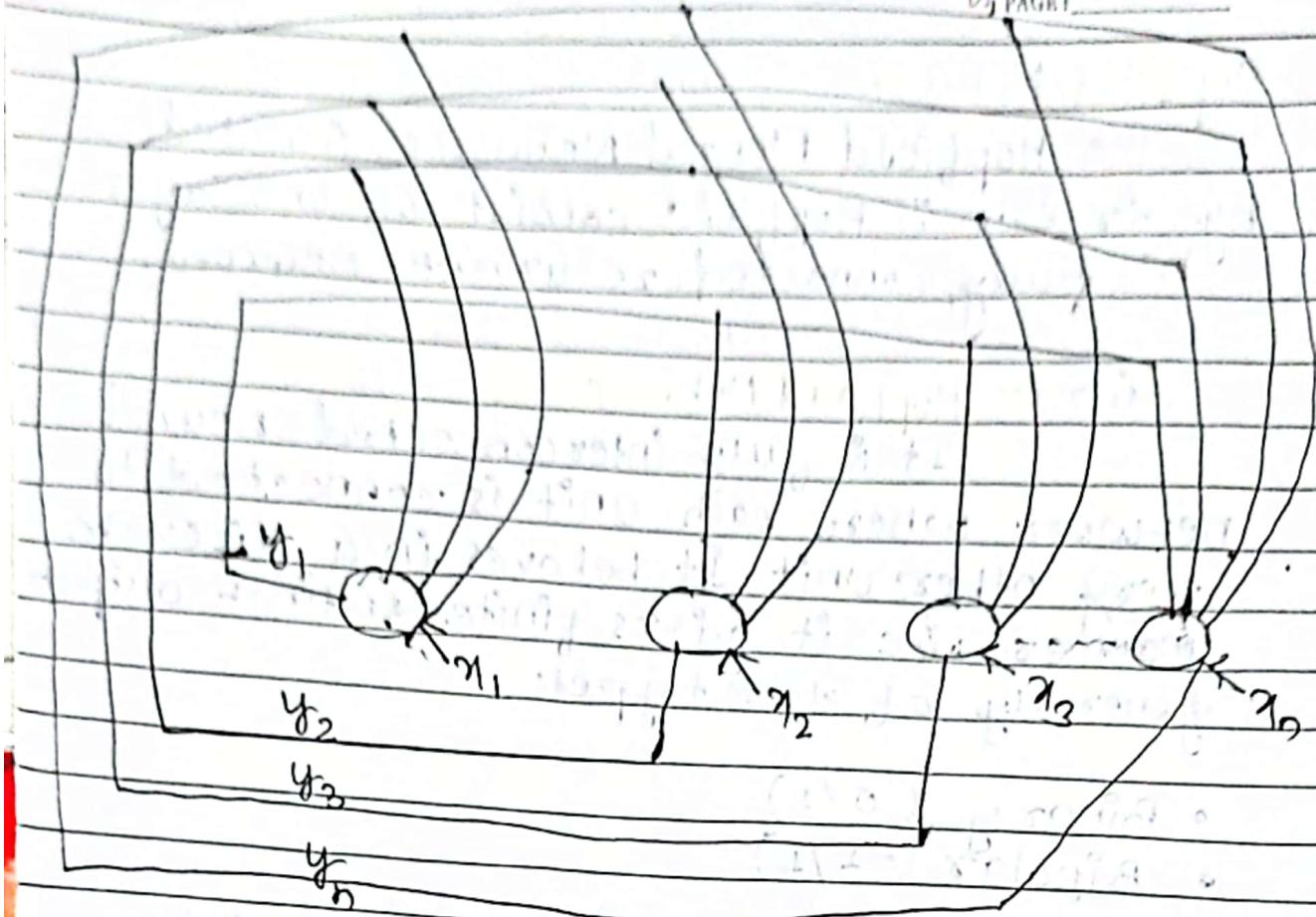
Structure & Architecture

- Each neuron has an inverting and a non-inverting output
- Being fully connected, the output of each neuron is an output to all other neurons but not self.

using w_{ij}, w_{ji} to find θ_i & θ_j for $i \neq j$

$$\theta_j = -\sum w_{ji} \cdot \theta_i + \theta_j$$

Good Write



Here $x_1, x_2, \dots, x_n \rightarrow$ Input to the n given neurons.

$y_1, y_2, y_3, \dots, y_n \rightarrow$ Output obtained from the n given neurons.

$w_{ij} \rightarrow$ weight associated with the connection between the i^{th} & j^{th} neuron.

Training Algorithm

for storing a set of input pattern $S(p)$ [$p = 1$ to P] where,

$$S(p) = S_1(p) - S_i(p)$$

Good Write

$\therefore S_n(p)$, the weight matrix is given by

for binary pattern

$$w_{ij} = \sum_{p=1}^P [2s_i(p)-1][2s_j(p)-1]$$

w_{ij} for all $i \neq j$

for bipolar

$$w_{ij} = \sum_{p=1}^P [s_i(p)s_j(p)]$$

where $w_{ii}=0$
for all $i=j$

(i.e. weights here have no self connections)

Testing Algorithm

Step 1 - Initialize weights (w_{ij}) to store patterns (using training algorithm)

Step 2 - For each input vector y_i , perform steps 3-8

Step 3 - Make initial activators of the network equal to the external input vector x .

$$y_i = x_i \quad (\text{for } i=1 \text{ to } n)$$

Step 4:- for each vector y_i , perform
step 5-8.

Step 5: calculate the total input of the
network y_{in} using the equation

$$y_{in} = x_i + \sum_j [y_j w_{ji}]$$

Step 6: Apply activation over the total
input to calculate the output as
per the equation below.

$$y_i = \begin{cases} 1 & \text{if } y_{in} > \theta_i \\ y_i & \text{if } y_{in} = \theta_i \\ 0 & \text{if } y_{in} < \theta_i \end{cases}$$

Where θ is normally taken as 0.

Step 7:- Now feedback the obtained
output y_i to all other units.
Thus the activation vectors are
updated.

Step 8:- Test the network for convergence.

Training

During training of the Discrete Hopfield network, the weight will be updated. Here we can take binary input as well bipolar inputs. The weight updation can be done by the following ways.

Here s_{ip}

where $S = \text{set of input pattern}$
 $p = 1 \text{ to } P$

For binary inputs:

$$w_{ij}^o = \sum_p [2s_i(p) - 1][2s_j(p) - 1] \quad \text{for } i \neq j$$

$$w_{ii}^o = 0$$

For bipolars:

$$w_{ij}^o = \sum_p [s_i(p)][s_j(p)] \quad \text{for } i \neq j$$

Working of the algorithm:

Step 1: Initialize all the weight that can be obtained from Hebbian principle

Step 2: For each input vectors y_i , repeat or perform steps from 3-7.

Step 3: Make the initial activators y_0 equal to the external input vectors, i.e. $y_0 = x_i$ for all $i=1$ to n .

Step 4: For all weight that need to be updated perform 5 to 7.

Step 5: Calculate the net input of all the set of input vectors.

$$y_{in} = x_i + \sum_j (w_{ji}x_j y_0)$$

Step 6: Activate all the y_0 using the below equation

$$y = \begin{cases} 1 & y_{in} \geq 0 \\ 0 & y_{in} < 0 \end{cases}$$

where $\theta =$

Step 7:- Feed back to obtained output

Step 8:- Test it for convergence.

Numerical.

PAGE

Consider the following problem. We are required to create discrete Hopfield net with bipolar representation of input vector as $[1 \ 1 \ 1 \ -1]$ or $[1 \ 1 \ 1 \ 0]$ is stored in the network. Test the hopfield net with missing entries in the first & second component of the stored vector (i.e. $[0 \ 0 \ 1 \ 0]$)

$$w_{ij} = \sum s_i^T c_j t(p)$$

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

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

and the weights with no self connection.

$$w_{ij} = \begin{bmatrix} 0 & 1 & 1 & -1 \\ 1 & 0 & 1 & -1 \\ 1 & 1 & 0 & -1 \\ -1 & -1 & -1 & 0 \end{bmatrix}$$

Step 3: As per the question the input vector x with missing entries $\gamma = [0 \ 0 \ 1 \ 0] \ ([\gamma_1, \gamma_2, \gamma_3, \gamma_4])$

- Make $y_i^0 = x = [0, 0, 1, 0]$

$$= (y_1, y_2, y_3, y_4)$$

Now

Choosing unit y_0 for updating its activation.

Take i th column of the weight matrix for calculation.

$$y_{ini} = x_i + \sum_{j=1}^4 [y_j w_{j,i}]$$

$$= 0 + [0 \ 0 \ 1 \ 0] \begin{bmatrix} 0 \\ 1 \\ 1 \\ -1 \end{bmatrix}$$

$$= 0 + 0 + 0 + 1 + 0$$

$$= 1$$

Applying activation
since $y_{ini} \geq 0$

$$\therefore y_1 = 1$$

giving feed back to other units,

$$y = [1 \ 0 \ 1 \ 0]$$

which is not equal to the input vector.

$$x = [1 \ 1 \ 0]$$

Hence no convergence.

Now, for next we have.

$$y_3 = [1 \ 0 \ 1 \ 0]$$

$$y_{in3} = y_3 + \sum_{j=1}^4 [y_j w_{j3}]$$

$$= 1 + [1 \ 0 \ 1 \ 0]^T \begin{bmatrix} 1 \\ 0 \\ -1 \\ 0 \end{bmatrix}$$

$$= 1 + (1+0+0+0)$$

$$= 2$$

Since $y_{in} > 0 \therefore y_3 = 1$

Now,

giving feed back to other units
we get $y = [1 \ 0 \ 1 \ 0]$

which is not equal to input vector

$$x = [1 \ 1 \ 1 \ 0]$$

$$\text{Hence no. } [0 \ 1 \ 0 \ 1] \neq x$$

$$y_{inj_4} = y_4 + \sum_{j=1}^4 [y_j \cdot w_{j4}]$$

$$= 0 + [1 \ 0 \ 1 \ 0] \begin{bmatrix} -1 \\ -1 \\ -1 \\ 0 \end{bmatrix}$$

$$= -2$$

which is less than 0.

$$\therefore y_4 = 0.0 \neq 0 \quad \text{not equal}$$

giving feed back to other.

$$\therefore y = [1 \ 0 \ 1 \ 0]$$

not equal to the $x = [1 \ 1 \ 1 \ 0]$

$$y_{inj_2} = y_2 + \sum_{j=L}^4 [y_j \cdot w_{j2}]$$

$$= 1 + [1 \ 0 \ 1 \ 0] \begin{bmatrix} 1 \\ 0 \\ 0 \\ -1 \end{bmatrix}$$

$$= 1 + (1+0+1+0) \\ = 3$$

Good Write $\therefore y_2 = 1$

giving feedback to the other units.
we get,

$$y_i^0 = [1 \ 1 \ 1 \ 0]$$

which is equal to input vector.

Hence convergence with vector x .

Continuous Hopfield Network

Here the time parameter is treated as a continuous variable. So, instead of getting binary/bipolar outputs we can obtain values that lie between 0 & 1.

- It can be used to solve constrained optimization and associative memory problems.

$$v_i^0 = g(u_i^0)$$

where

v_i^0 = output from the continuous

hopfield network

u_i^0 = internal activity of a node

- in continuous hopfield network

Energy function evaluation.

Energy function E_f

$$\text{activation energy} = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j w_{ij} -$$

$$\sum_{i=1}^n a_i y_i + \sum_{i=1}^n b_i y_i$$

It determines the stability of discrete Hopfield network

Characteristics

Condition

In a stable network, whenever the state of node changes, the above energy function will decrease.

$$\Delta E_f = E_f(y_i^{(k+1)}) - E_f(y_i^{(k)})$$

$$= \left(\sum_{j=1}^n w_{ij} y_j^{(k)} + a_i - b_i \right) (y_i^{(k+1)} - y_i^{(k)})$$

$$= (\text{net}_i) \Delta y_i$$

$$\Delta y_i = y_i^{(k+1)} - y_i^{(k)}$$

The change depends on the fact that only one unit can update its activation at a time.

continuous.

$$E_f = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j w_{ij} - \sum_{i=1}^n \alpha_i y_i + \frac{1}{\lambda}$$

$$\sum_{i=1}^n \sum_{j=1}^n w_{ij} g_{ji} \int_0^{y_i} a^{-1}(y) dy$$

Here λ is gain parameter and
 g_{ji} input conductance.

discrete.

$$E_f^k = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n y_i y_j w_{ij} - \sum_{i=1}^n \alpha_i y_i + \sum_{i=1}^n \theta_i y_i$$

In a stable condition, the energy also get decreases.

$$\begin{aligned} \Delta E_f^k &= E_f^k(y_i^{k+1}) - E_f^k(y_i^k) \\ &= - \left(\sum_{j=1}^n w_{ij} y_j^k + \alpha_i - \theta_i \right) (y_i^{k+1} - y_i^k) \\ &\approx -(\text{net}_i) \Delta y_i \end{aligned}$$

for continuous

$$E = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k y_i y_j w_{ij} - \sum_{i=1}^n x_i y_i + \text{initial } y_i \\ + \frac{1}{T} \sum_{i=1}^n \sum_{j=1}^k w_{ij} g_{ij} \int a^{-1}(y) dy$$

where a is a gain parameter and g_{ij} is input conductance.

$$E = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k w_{ij} y_i y_j - \sum_{i=1}^n x_i y_i + \\ \frac{1}{T} \sum_{i=1}^n \sum_{j=1}^k w_{ij} g_{ij} \int a^{-1}(y) dy$$

$$E = - \left(\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^k y_i y_j w_{ij} - \sum_{i=1}^n x_i y_i + \right. \\ \left. \sum_{i=1}^n \theta_i y_i \right)$$

① Storage (Learning)

weight matrix W is calculated between m pairs of patterns, and stored in synaptic weights of the networks.

$$W = \sum_{m=1}^M X_m Y_m^T$$

② Testing

$$Y_m = \text{sign}(W \cdot X_m)$$

$$X_m = \text{sign}(W \cdot Y_m)$$

③ Retrieval: for unknown vector X

(a corrupted or incomplete version of a pattern from set A or B) to the BAM and returns previously stored.

$$X \neq X_m$$

- Initialize the BAM

$$X(0) = X, p=0$$

- calculate the BAM output at iteration

$p:$

$$Y_p = \text{sign}(W^T X_p)$$

- Update the input vector

$$X(p+1) = \text{sign}(W Y_p)$$