

~~Section 8~~ ~~Topic 8~~ ~~NN~~
Ques-8 Define Single & Multi layer perceptron and their algorithm ~~with~~ ~~for~~ ~~any~~
Classification Models :- A Neural Network classifies a given object presented to it according to the off ~~function~~ & activation. For binary outputs, 1 corresponds to one class & 0 corresponds to the other.

2 types of classification models are discussed below:-

- 1) Single layer Perceptron
- 2) Multilayer Perceptron

The Perceptron :- The Perceptron is a program that learns concepts, i.e. it can learn to respond with true(1) or false(0) for inputs we present to it, by repeatedly "studying" examples presented on it.

Example for e.g., how do you teach a child to recognize what a chair is? You show him examples telling him "This is a chair, that one is not a chair" until the child learns the concept of what a chair is. This is exactly the idea behind the perceptron.

Single layer Perceptron :- A single layer Perceptron consists of an input layer & an output layer. An output unit will assume the value 1 if the sum of weighted inputs is greater than its threshold i.e.

$$\sum w_{ij}x_i > \theta_j$$

w_{ij} weight from i to j

$x_i \rightarrow$ input

$\theta_j \rightarrow$ threshold

There are 2 classes A & B. If $\sum w_{ij}x_i > \theta_j$ then object will be classified as class A otherwise

class B. Explain

Linear separability \rightarrow As discussed in notes of $\text{unit } 4$

K The Perceptron Algorithm

i) weight Initialization

set all weights & node thresholds to small random numbers.

ii) Calculation of Activation

The activation level O_j of an output unit

$$O_j = f_h(w_{ji}x_i - O_j)$$

f_h = Hard limit function

$$f_h(a) = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{Else} \end{cases}$$

iii) Weight Training

a) Adjust weight $w_{ji}(t+1) = w_{ji}(t) + \Delta w_{ji}$

$w_{ji}(t)$ → weight from unit i to unit j at time t

Δw_{ji} → weight adjustment

b) weight change may be computed by Delta learning rule:

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

δ_j → error at unit j

$$\delta_j = T_j - O_j$$

T_j → desired output activation

" " at output unit j

O_j → actual " " at output unit j

c) Repeat the iterations will converge.

* Limitations :- first, the output values of a perceptron

can take on only one of two values (T or F).

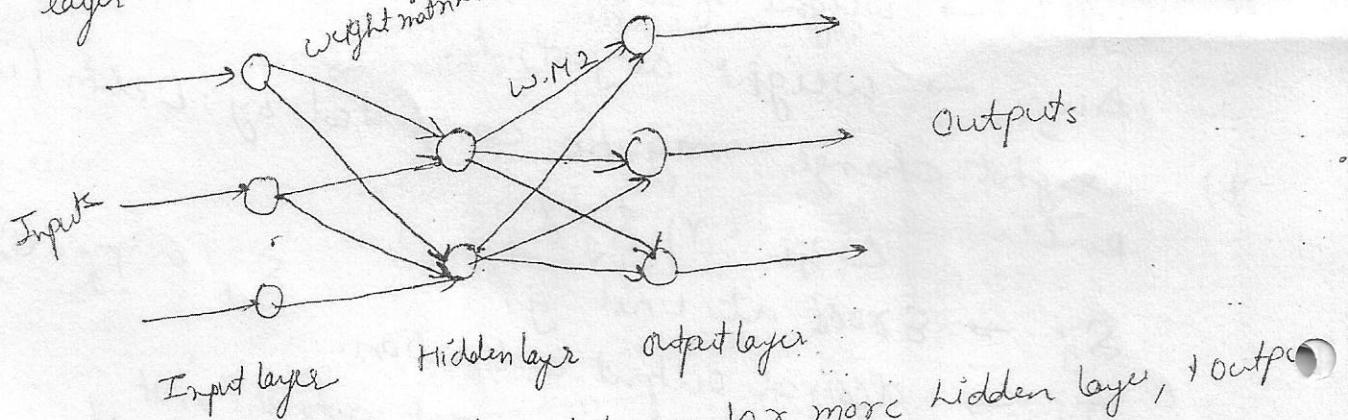
→ second perceptrons can only classify linearly separable sets of vectors.

→ If a straight line or plane can be drawn to separate the input vectors into their correct categories, the input vectors are linearly separable & the perceptron will find the solution. If the vectors are not linearly separable learning will never reach a point where all vectors are classified properly.

→ The most famous inability to solve problems with linearly nonseparable is boolean exclusive-or problem for these multilayer perceptrons is used.

* Multilayer Perceptron :- An MLP is a network of simple neurons called perceptrons. It was introduced by Rosenblatt in 1958.

→ It is a feedforward neural network with at least one hidden layer. It can deal with nonlinear classification problem.



Neuron layer → 1 input layer, 1 or more hidden layers, 1 output layer.

Input value type → Binary

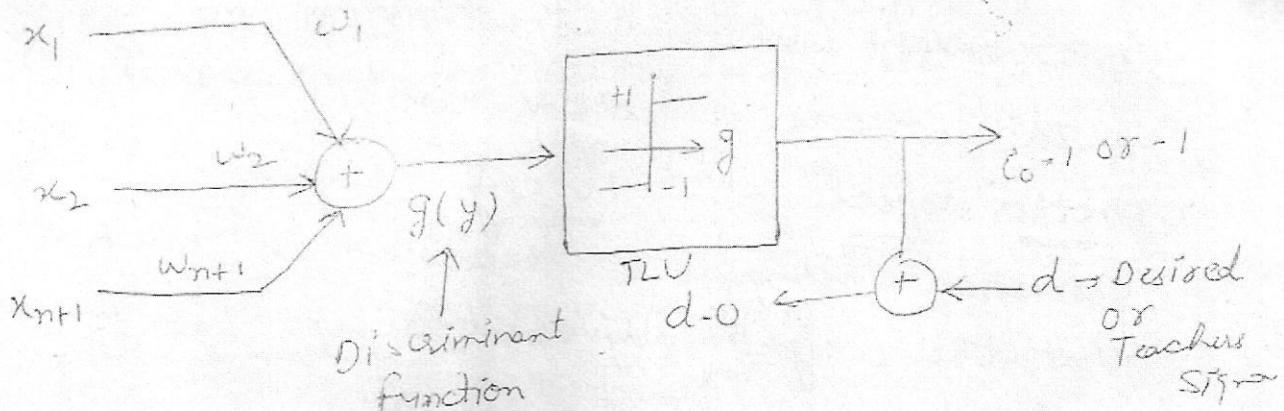
Activation func. → Sigmoid

Learning Method → Supervised

Learning algorithm → Delta learning Rule, Back propagation

B(2)

* Training & Classification: Using the discrete perceptron



Linear Dichotomizer using TLV Based Perception

TLV \rightarrow Threshold logical Unit.

Dichotomizer: Classifier that divides the pattern space into 2 classes (0 or 1 means 2). i.e x_1, x_2 only.

* Training Sequence: Sample pattern vectors x_1, x_2, \dots, x_n called training sequence i.e sequence in which ANN is trained. Response is provided by the teacher & it specifies the classification information for each input vector.

→ The N/w learns from the experience by comparing the targeted correct response with the actual response to classify structure & the weight is adjusted after each incorrect response based on error value generated.

→ In this discrete perceptron concept the formula for adjusted weights is

$$w' = w + Cy$$

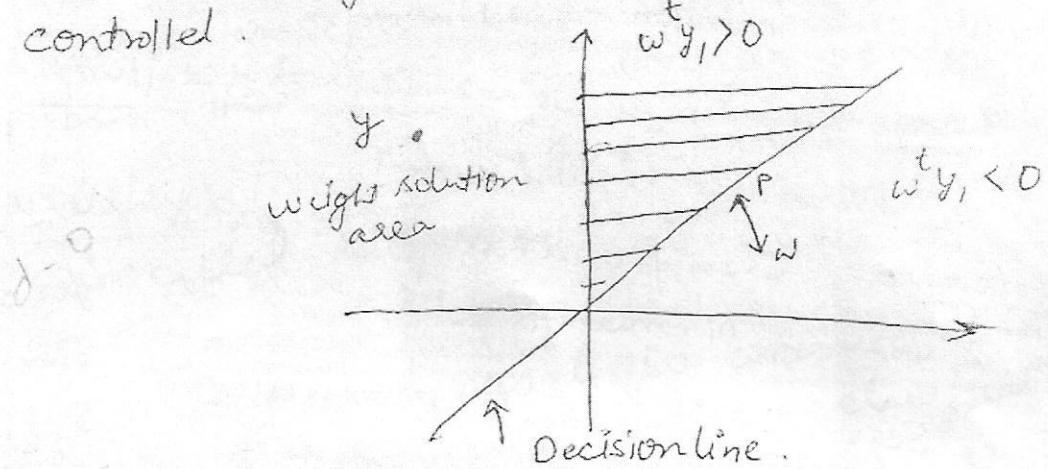
$C \rightarrow$ constant & $C > 0$ called correction increment

$+ve \rightarrow$ sign for pattern class 1, -ve \rightarrow for class 2

If a correct classification takes place under the rule, no adjustment of weight is made, the new weight (w') = old weight (w).

→ This is one aspect of weight adjustment.

Another aspect → Is by using the geometrical relation where correction increment C is chosen in such a way that the weight adjustment step size is meaningful controlled.



Let p be distance of a point w' from the plane $w^T y = 0$ in $(n+1)$ dimensional Euclidean space & it is calculated acc. to formula..

$$p = \frac{\pm w^T y}{\|y\|}$$

$p \rightarrow$ distance & is always a non-negative scalar.

∴ can be written as

$$p = \frac{w^T y}{\|y\|}$$

Now C must be selected such that corrected weight vector ω^2 based on $\omega' = \omega + Cy$ dislocated on the decision hyperplane $\omega'^T y = 0$ which is the decision hyperplane used for particular correction step

$$\begin{aligned} \omega'^T y &= 0 \\ (\omega' + Cy)^T y &= 0 \\ \therefore (\omega'^T + Cy^T) y &= 0 \\ \omega'^T y + Cy^T y &= 0 \\ Cy^T y &= -\omega'^T y \\ C &= -\frac{\omega'^T y}{y^T y} \end{aligned}$$

Since C is correction constant & is positive.

$$C = \frac{|\omega'^T y|}{y^T y}$$

* Single Discrete Perception Training Algorithm:

Given there are P training pairs.

$\{x_1, d_1, x_2, d_2, \dots, x_p, d_p\}$ where

x_i is $(n \times 1)$ & d_i is (1×1) , $i = 1, 2, \dots, P$

x_i is $(n \times 1)$ & d_i is (1×1) , $i = 1, 2, \dots, P$

Input vector $y_i = [x_i^T]^T$ for $i = 1, 2, \dots, P$

Step 1 $C > 0$ is chosen

Step 2 weights are initialized at ω at small random values & $K \leftarrow 1$, $P \leftarrow 1$, $E \leftarrow 0$ where K is training step.

Step 3: P is step counter within the training step or cycle.

Step 3 Training cycle starts.

Input is presented & output is computed.

$$y \leftarrow y_p, d \leftarrow d_p, o \leftarrow \text{sgn}(\omega^t y)$$

Step 4 Weights are updated

$$\omega \leftarrow \omega + \frac{1}{2} C(d - o)y$$

Step 5 Cycle error is computed

$$E \leftarrow \frac{1}{2}(d - o)^2 + E \quad \text{Previous cycle error}$$

Step 6 If $P < P$ then $P \leftarrow P+1$, $K \leftarrow K+1$ & go to step 3, otherwise go to step 7.

Step 7 Training is completed. If $E = 0$, terminate the training, output weight & K .

If $E > 0$, then $E \leftarrow 0$, $P \leftarrow 1$, new training step or iteration starts, starting from step 3.

* Single layer Continuous Perceptron Networks for linear separable classification (SCPTA).

Given P training pairs

$$\{x_1, d_1, x_2, d_2, \dots, x_P, d_P\}$$

Input vectors are used

$$y_i = \begin{bmatrix} n_i \\ 1 \end{bmatrix} \quad \text{for } i = 1, 2, \dots, P$$

where x_i is $(m \times 1)$, d_i is (1×1) , $i = 1, 2, \dots, P$

$P \rightarrow$ step counter

$K \rightarrow$ training step

1. choose $\eta > 0$, λ (steepness constant) = 1,

$$E_{max} \rightarrow 0$$

2. weight are initialized at w at small random value, counters & error are initialized.

$$k \leftarrow 1, p \leftarrow 1, E \leftarrow 0$$

3. Training starts \rightarrow input y is presented & o/p is computed i.e

$$y \leftarrow y_p, d \leftarrow d_p, o \leftarrow f(w^t y)$$

4. Weights are updated :

$$w^{k+1} \leftarrow w^k + \frac{1}{2} n(d^k - o^k)(1 - o^k)y$$

5. Error E is computed

$$E^{k+1} \leftarrow \frac{1}{2} (d^k - o^k)^2 + E^k$$

6. Error & If $p < P$, then $p \leftarrow p+1, k \rightarrow k+1$
& repeat step 3, otherwise go to step 7.

7. Training cycle is completed.

If $E < E_{max}$ terminate the session.

If $E \geq E_{max}$ then $E \leftarrow 0, p \leftarrow 1$ start new train

cycle & go to step 3.

Ques 9 Define the term \rightarrow Classification, Model and Decision Region.

Ans Classification Model, Features & Decision Region :-

1) Basic Meaning of Classifiers :- Assume that a set of 8 points a_0, \dots, a_7 in three dimensional space is as follows:

where

$$\therefore a_0 = (-1, -1, -1)$$

$$a_4 = (1, -1, -1)$$

$$a_1 = (-1, -1, 1)$$

$$a_5 = (1, -1, 1)$$

$$a_2 = (-1, 1, -1)$$

$$a_6 = (1, 1, -1)$$

$$a_3 = (-1, 1, 1)$$

$$a_7 = (1, 1, 1)$$

Elements of this set needs to be classified into two categories :-

(1) first category can be set of points with two or more ones.

(2) Second category contains all the remaining points that do not belong to first category.

∴ first category consists of a_3, a_5, a_6, a_7 & remaining points lie in second category.

This is the function of classifier that is to classify to divide the patterns in different categories
→ One of the most important application of ANN
pattern classification :

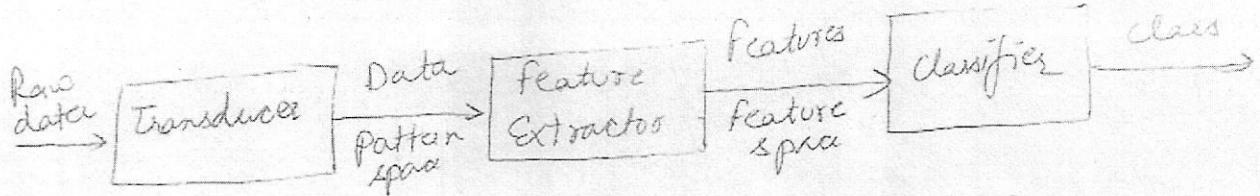
Pattern :- Pattern is the quantitative description of an object, event

Patterns can be classified as spatial or temporal patterns

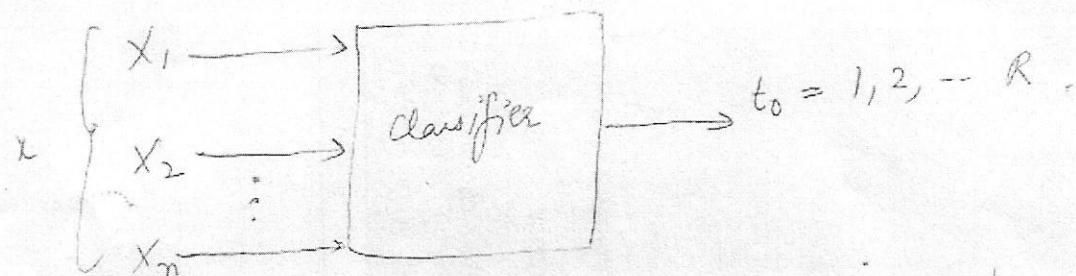
1) Spatial :- Ex. are pictures, whether maps, video go

2) Temporal :- Ex. are speech signal, ECG etc. These patterns involve the ordered sequence of data appear in time.

* Pattern Classification :- The main goal of pattern classification is to assign a physical object, event or phenomenon to one of the prespecified classes/categories for e.g. we classify various objects into different categories like living, non-living things etc.



a) Block diagram of classification sys



b) General diagram of pattern classifier

- b) General diagram of I

 - Transducer generally convert one form of energy to another form of energy. Here the input is raw data space $x = [x_1, x_2, \dots, x_n]$ & o/p is transducer is data in pattern space that belong to a certain category that is to $= 1, 2, \dots, R$. converted data at the o/p can be compressed. Compressed data is called features.
 - Classifier at the end, convert or map the feature to the particular class.

3.3 | LINEAR SEPARABILITY

The concept of linear separability in neural networks has been borrowed from perception in living beings (humans). In psychology, perception means mental organization and interpretation of sensory information. Perception is influenced by

- (1) Intensity and physical dimensions of the stimulus
- (2) Effects of preceding stimulation
- (3) Past experience
- (4) Attention factors such as readiness to respond to stimulus, motivation and emotional state.

Patterns are perceived in visual organisation according to :

- (1) The nearness to each other.
- (2) Similarity between them.

(3) The tendency to perceive complete figures.

(4) The ability to distinguish important figures from background.

Sensory organs merely gather the signals whereas the stimuli are actually perceived in the brain through :

- (1) *Perceptual constancy*. It is the tendency of brain to interpret one object in the same manner regardless of such variations as distance, angle of sight or brightness.
- (2) *Selective attention*. It is the tendency of the brain to focus on a limited number of stimuli and ignore those that are considered less important.
- (3) *Depth perception*. It is produced by a variety of visual cues indicating perspective and by a slight disparity in the images of an object on the two retinas.

Now we will see how this concept of linearity is extended to linear operators used extensively in neural networks. Our aim in training a neural network is that it will respond with the desired classification when presented with an input pattern that it was trained on. Alternatively, it should respond when presented with an input that is sufficiently similar to one of the training patterns. The desired response from the output is a 'yes' if the input pattern is a member of its class and a 'no' if it is not. A 'yes' response is represented by an output signal of 1, a 'no' by an output signal of -1 in case of bipolar signals. Since we want one of the two responses, the activation function may be taken to be a step function. The value of function is 1 if the net input is positive and -1 if the net input is negative. The induced local field of the neuron for m -dimensional space is

$$v = \sum_{i=1}^m w_i x_i + b \quad \dots(1)$$

where b is the externally applied bias.

Hence, the goal of the network becomes to classify the set of externally applied stimuli x_1, x_2, \dots, x_m into one of the two classes. The decision rule for the classification is to assign the point represented by the inputs x_1, x_2, \dots, x_n to class I if the output v is +1 and to class II if it is -1. One can easily visualize that the boundary between the region where $v > 0$ and the region where $v < 0$ which we call the *decision boundary*. It is determined by the relation

$$\sum_{i=1}^m w_i x_i + b = 0 \quad \dots(2)$$

Depending on the number of input units in the network, Eq. (2) can represent a line, a plane or a hyperplane.

If there are weights (and a bias) so that all the training input vectors for which the correct response is +1 lie on one side of the decision boundary and all the training input vectors for which the correct response is -1 lie on the other side of the decision boundary, we say that the problem is "*linearly separable*", i.e., separable by a linear combination of inputs.

Let us now consider the case of two input variables x_1 and x_2 for which the decision boundary takes the form of a straight line as shown in fig. (2). A point (x_1, x_2) that lies above the boundary line is assigned to class I and a point (x_1, x_2) that lies below the boundary line is assigned to class II. The choice of the sign for bias b determines which side of the separating line corresponds to a +1 response and which side to a -1 response.

Until now, we have considered problems with two input dimensions (x_1, x_2) . If there is only one input dimension (x) , then the two-class problem can be solved using a perceptron if and only if there is some value x_0 of x such all samples of one class occur for $x > x_0$ and all samples of other class occur for $x < x_0$ as shown in Fig. (3).

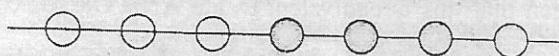


Fig. (3) Classification problem in one dimensional input space.

Similarly, if there are three input dimensions, a two class problem can be solved using a perceptron if and only if there is a plane that separates samples of different classes. Co-efficients of terms in the equation of the plane correspond to the weights of the perceptron.

For the perceptron to function properly, the two classes must be *linearly separable*. In other words, the patterns to be classified must be sufficiently separated from each other to ensure that the decision surface consists of a hyperplane. Refer to Fig. (4) which shows the case of a two-dimensional perceptron. In Fig. (4a), the two classes are sufficiently separated from each other so that we can draw a hyperplane (which is a straight line in this case) as the decision boundary.

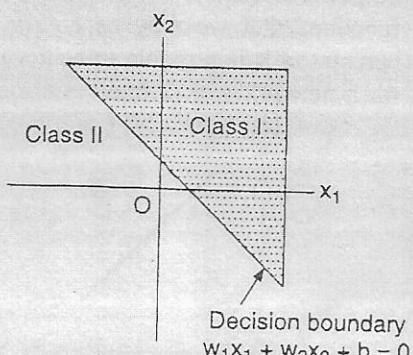


Fig. (2) Classification problem in two-dimensional input space.

Fig. (4) Classification problem in two-dimensional input space.

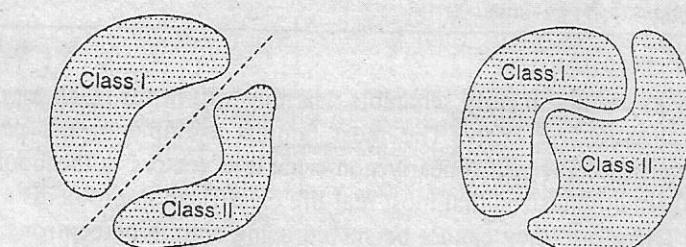


Fig. (4)

In Fig. (4b), they become non-linearly separable as no straight line can possibly separate samples belonging to the two classes. This type of situation is beyond the computing capability of perceptron. Most real-life classification problems are linearly non-separable. Hence they cannot be solved using perceptrons. This is the fundamental limitation of simple perceptron.

For simplicity sake, let us now consider a neural network with 2 inputs and 1 output. Since there are two possible responses for each input there are 2^4 different functions that we may be able to train a very simple network to perform. Using perceptron it is possible to achieve linear separability functions as shown in Fig. (5) for a network with 2 input and 1 output.

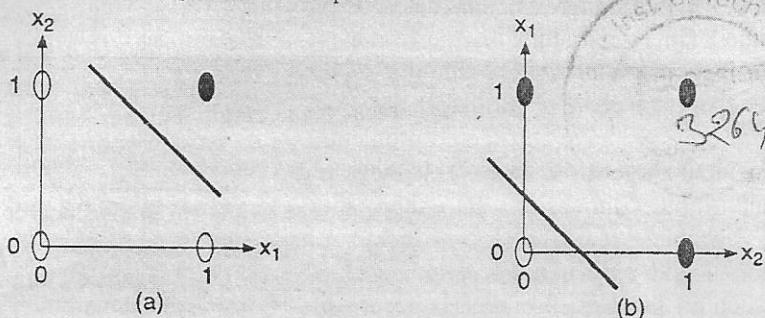


Fig. (5)

It can be seen that this is equivalent to the AND/OR logic gates as shown in Fig. (6).

x_1	x_2	Output
0	0	0
0	1	0
1	0	0
1	1	1

(a) AND logic gate

x_1	x_2	Output
0	0	0
0	1	1
1	0	1
1	1	1

(b) OR logic gate

Fig. (6)

Note : We have included a bias input in the examples-discussed. This is because if a bias weight is not included, the decision boundary will be forced to go through the origin. Sometimes, a third input component that is always 1 is also given (as in examples discussed after this article).

Limitation of Perceptron

If the classes are not linearly separable, learning will never reach a point where all inputs are classified properly. The most famous example of the perceptron's inability to solve problems with linearly non-separable vectors is the boolean XOR problem. It is again emphasized that most real-life classification problems are linearly non-separable. Therefore, they cannot be solved using simple perceptrons.

Example 1. Obtain the decision boundary for binary AND function using Hebb rule after first training pair. The given binary inputs and targets are :

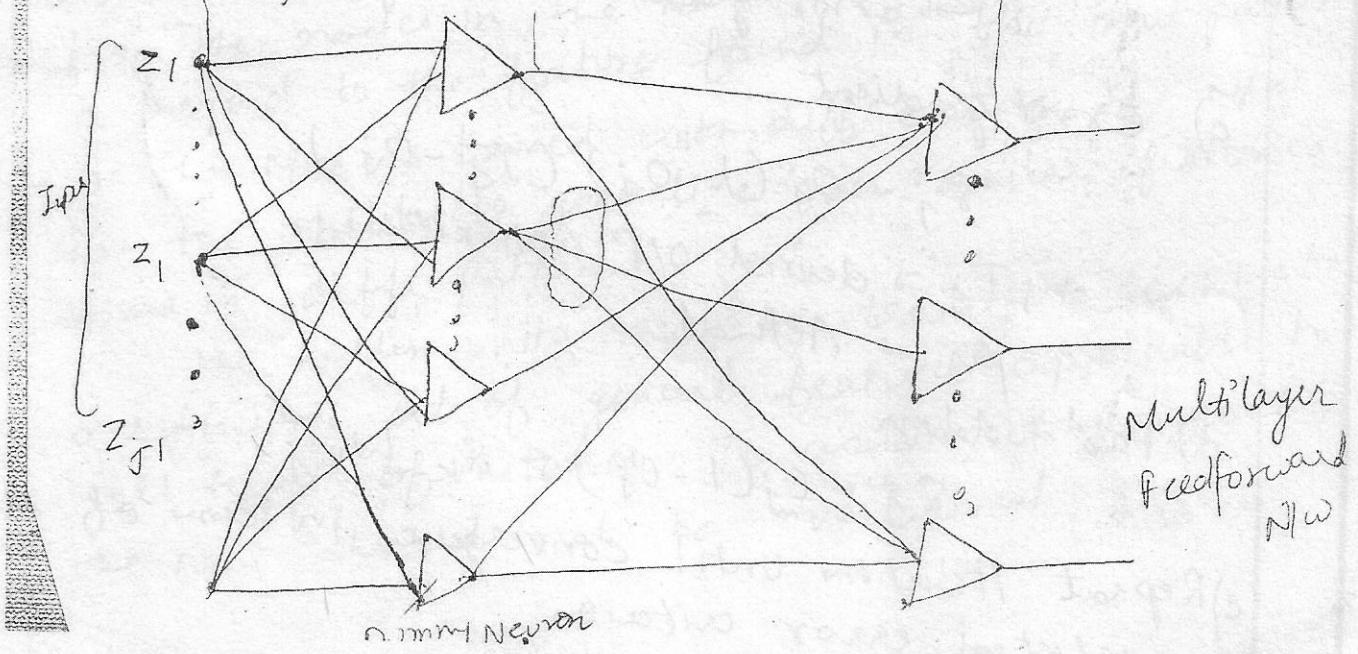
INPUT	TARGET
($x_1 \ x_2 \ 1$)	
(1 1 1)	1
(1 0 1)	0
(0 1 1)	0
(0 0 1)	0

* Neural Network learning :- The neural network has also been called the "connectionist". It contains a large no. of neuron like processing elements & a large no. of weighted connections b/w the elements. The weights on the connections encode the knowledge of a N/w. The various learning algorithms are:

1. Backpropagation (Generalized Delta Learning Rule)
2. Reinforcement learning
3. ART Networks

1. Backpropagation (Generalized Delta Learning Rule) :-

- Backpropagation Network is Multilayer feedforward N/w with diff. transfer function in ANN & a more powerful learning rule. This learning rule is called "Backpropagation" which is a kind of gradient descent technique with backward error propagation.
- It is a supervised learning method & is an implementation of Delta Rule. It requires a teacher that knows, or can calculate, the desired output for any given input layer j of neuron i & layer k of neuron j .



Backpropagation Algorithm

1) Weight Initialization - set all the weights & non threshold to small random nos.

2) Calculation of Activation

The activation level o_j of hidden & o/p unit is determined by

$$o_j = f(\sum w_{ji} o_i - \theta_j)$$

θ_j → threshold

f → sigmoid func.

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

a = weight & input.

3) Weight Training: a) Adjust weight by

$$w_{ji}(t+1) = w_{ji}(t) + \Delta w_{ji}$$

b) Weight change is computed by

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

η → trial independent learning rate.

δ → error gradient at unit j .

c) Error gradient

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

T_j → desired o/p activation

o_j → Actual

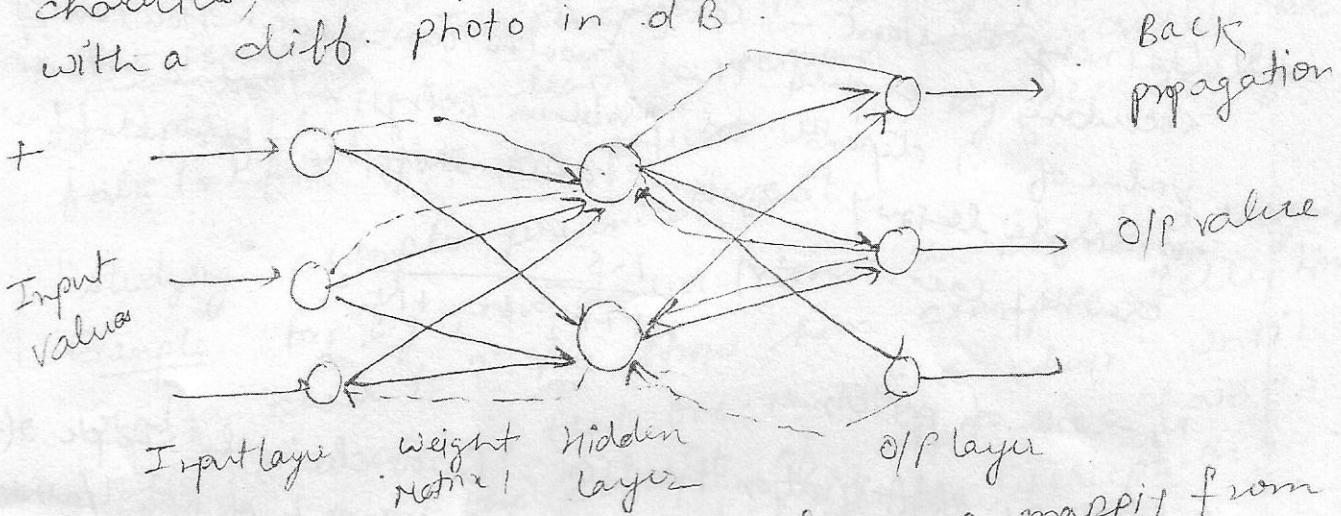
d) For hidden

$$\delta_j = o_j(1 - o_j) \sum \delta_k w_{kj}$$

e) Repeat iteration until convergence in terms of the selected error criterion

Some situations where a BP NN might be good idea:

- A large amount of input/output data is available, but you are not sure how to relate it to the o/p.
- The sol. to the problem may change over time, with the bounds of the given input & output parameters.
- Most common application is in image processing. Some eg would be: identifying hand-written characters, matching a photograph of a person's face with a diff. photo in dB.



- The backpropagation network learns a mapping from a set of input patterns to a set of o/p patterns. This NW can be designed & trained to accomplish a wide variety of mappings. This ability comes from the nodes in the hidden layer which learn to respond to the features found in the input pattern.
- As the NW is trained with diff. examples the NW has the ability to generalize over similar features found in diff. patterns.
- The hidden units must be trained to extract a sufficient set of general features applicable to all types of instances. To achieve this goal, the NW should not be overtrained.

Factors Affecting Backpropagation Training :-

(6)

- 1) Initialization of Weights - The initialization affects the complete solution. If all the weights start out with equal weight values & if the sol. requires that unequal weights be developed, the N/w may not be trained properly.
- There are 2 methods for training an ANN

- self-organizing ANN
- Back propagation ANN

- 2) Learning constant - Careful selection of step size (η) is often necessary to ensure the smooth convergence. However optimal value of η depends on problem being solved & there is no single learning constant value suitable for different training cases.

$$\eta = \frac{1.5}{N_1^2 + N_2^2 + \dots + N_n^2}$$

N_i → no. of patterns

- 3) Steepness of Activation function : Both choice & shape of the activation fun. strongly affect the speed of N/w learning.

$$f(\text{net}) = \frac{2}{(1 + \exp(-\gamma \text{net}))} - 1$$

$$\text{Then } f'(\text{net}) = 2 \gamma \exp(-\gamma \text{net}) / (1 + \exp(-\gamma \text{net}))^2$$

If reaches max. value of $\frac{1}{2} \gamma$ at $\text{net} = 0$
 Generally γ is kept at standard value of 1 & to control the learning speed using learning constant.

- 4) N/w size :- The N/w should have such a size as to capture the structure of the data & to model the problem. With some prior knowledge about the problem structure, one can form a good estimate of the proper N/w size.
- 5) Threshold value :- A neuron generates the O/P if the weighted sum of the input exceeds the threshold value.