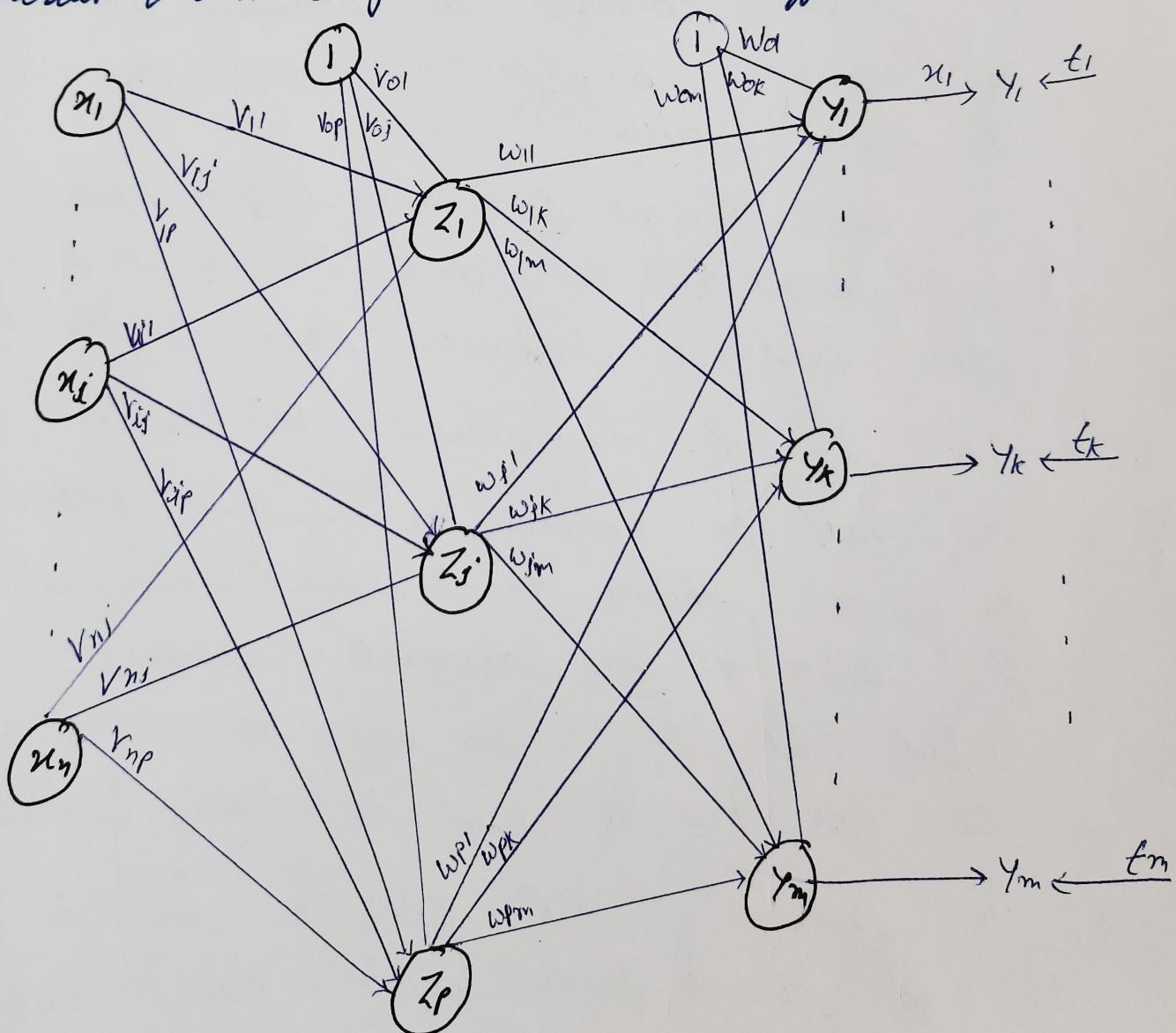


BACK PROPAGATION NETWORK

- * It is one of the most important developments in Neural Networks.
- * This learning algorithm is applied to multi-layer feed forward networks consisting of processing elements with continuous differentiable activation function.
- * The networks associated with back propagation learning algorithm are also called back propagation networks.
- * For a given set of training I/P O/P pairs, this algorithm provides a procedure for changing the weights in a BPN to classify the given I/P pattern correctly.
- * The basic concept for this weight update algo is gradient descent method. This is a method where error is propagated back to the hidden unit.
- * The aim of the neural N/w is to train the net to achieve a balance between the net's ability to respond (memorization) and its ability to give reasonable responses to the input that is similar but not identical to the one that is used in training (Generalization).
- * The training of the BPN is done in 3 stages:
 - 1) Feed forward of the input training pattern
 - 2) Calculation
 - 3) Back propagation of errors, and updation of weights.

Architecture of BPN:

- * A BPN N/w is a multilayer, feed-forward neural N/w consisting of an input layer, a hidden layer & an output layer.
- * The Neurons present in the hidden & O/P layers have bias which are connected from the unit whose activation is always 1.
- * The S/I/P are sent to the BPN as the output obtained from the net could be either binary (0,1) or bipolar (-1,+1)
- * The Activation function could be any function which increases monotonically and is also differentiable.



Terminology used in flow chart and in Training Algorithm are as follows

x = I/P training Vectors ($x_1 \dots x_i \dots x_n$)

t = layer output Vectors ($t_1 \dots t_k \dots t_m$)

α = Learning Rate

x_i = Input Unit i (since the I/P layer uses Identity Activation function, the I/P & O/P signals here are same)

v_{0j} = bias on j th hidden Unit

w_{0k} = bias on k th output unit

z_j = hidden unit j , The net Input to z_j is

$$z_{inj} = v_{0j} + \sum_{i=1}^n x_i \cdot v_{ij}$$

and the output is

$$z_j = f(z_{inj})$$

y_k = output unit k . The net Input to y_k is

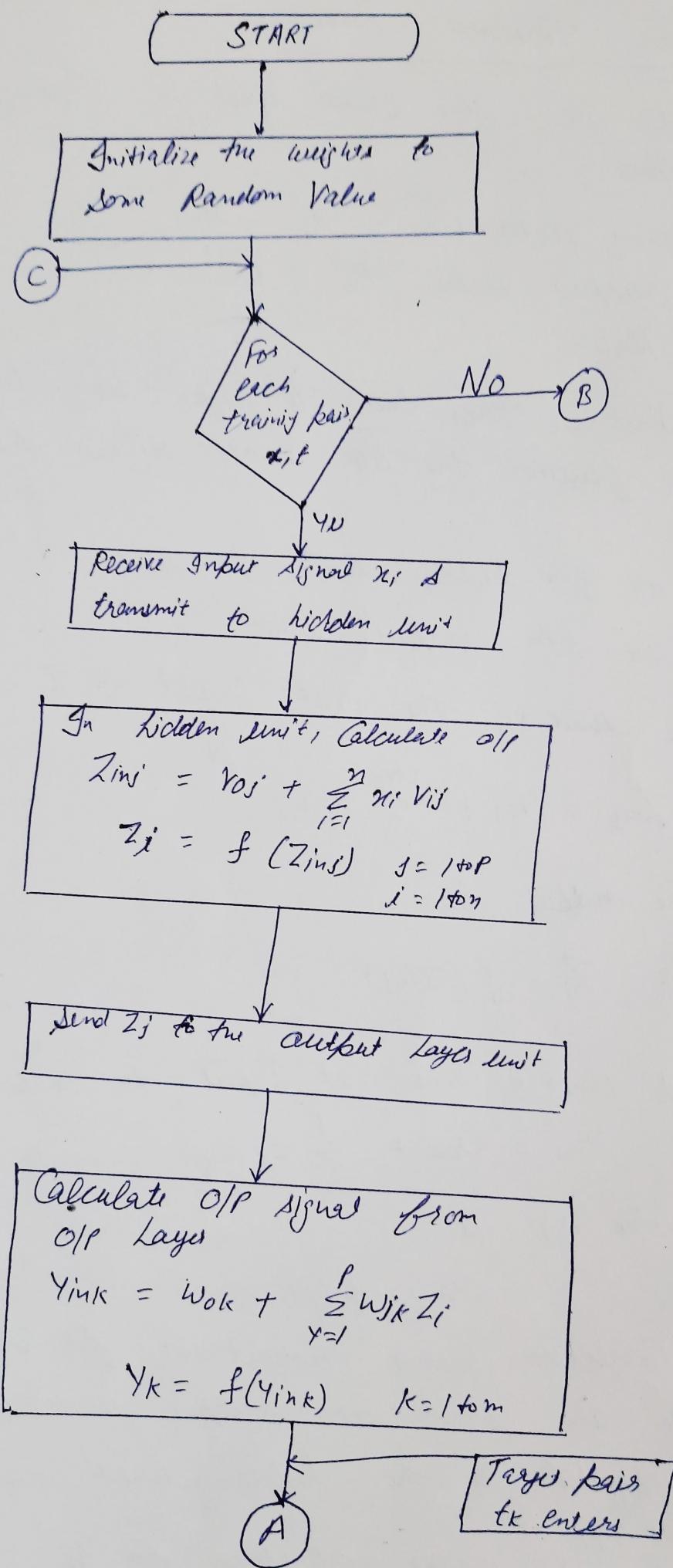
$$y_{in k} = w_{0k} + \sum_{j=1}^p z_j \cdot w_{jk}$$

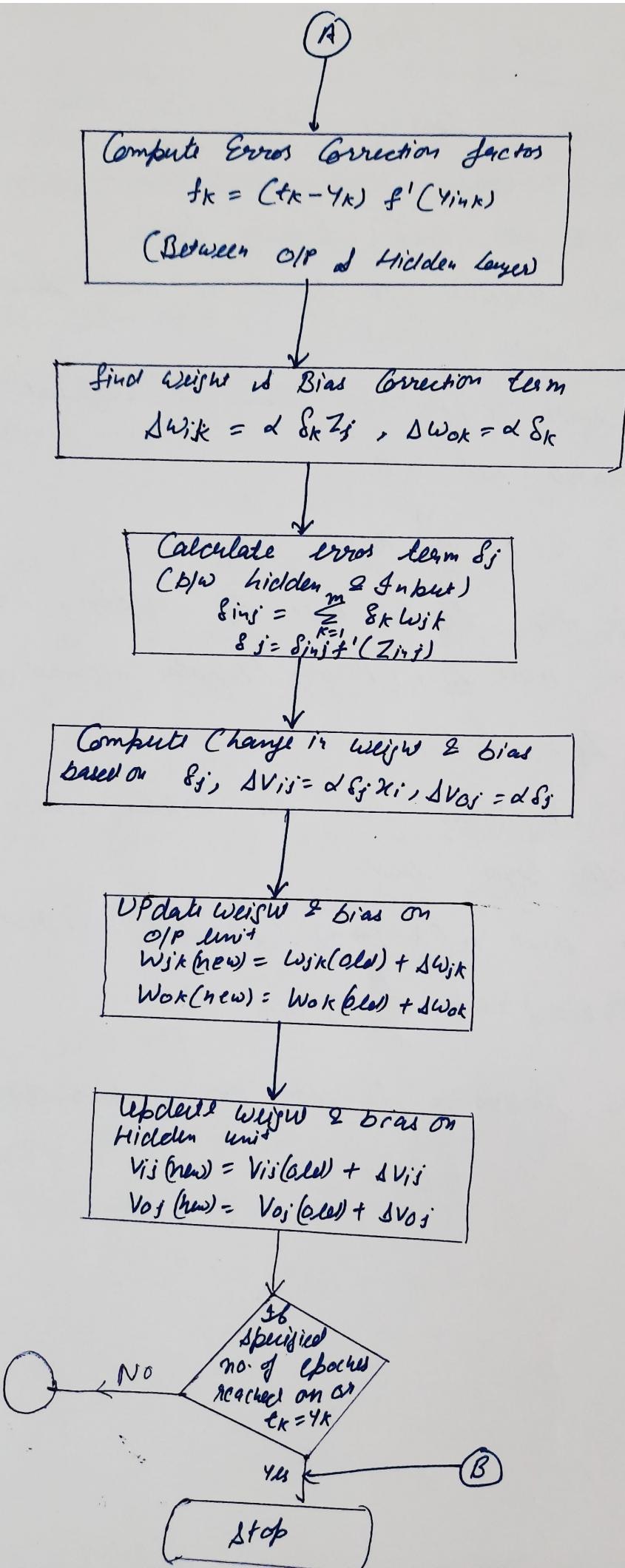
and the O/P is

$$y_k = f(y_{in k})$$

δ_k = error Correction weight adjustment for w_{jk} that is due to an error at output unit y_k which is back propagated to the hidden unit that feed into unit y_k

δ_j = error Correction weight adjustment for v_{ij} that is due to the back propagation of errors to the hidden unit z_j





TRAINING ALGORITHM FOR RPN

- Step 0: Initialize weights & learning rate (take some small random value)
- Step 1: perform step 2-9 when stopping condition is false
- Step 2: Perform steps 3-8 for each training pair
- Step 3: Each S/I/P unit receive S/I/P signal x_i and sends it to the hidden unit ($i=1 \text{ to } n$)
- Step 4: Each hidden unit Z_j ($j=1 \text{ to } p$) sums its weighted S/I/P signals to calculate net S/I/P

$$Z_{inj} = V_{0j} + \sum_{i=1}^n x_i v_{ij}$$

Calculate O/P of the hidden unit by applying its activation function over Z_{inj} (Binary / Bipolar sigmoid fns.)

$$Z_j = f(Z_{inj})$$

and send the O/P signal from the hidden unit to the S/I/P of O/P layer unit

- Step 5: for each O/P unit y_k ($k=1 \text{ to } m$) calculate the net Input

$$Y_{ink} = W_{0k} + \sum_{j=1}^p Z_j w_{jk}$$

and apply the Activation function to Compute O/P signal.

$$y_k = f(Y_{ink})$$

(4)

Back Propagation of errors (Phase II)

Step 6: Each output unit y_k ($k=1 \dots m$) receives a target pattern corresponding to the I/P training pattern & computes the error correction term.

$$\delta_k = (t_k - y_k) f'(y_{in,k})$$

The derivative $f'(y_{in,k})$ can be calculated as is.

On the basis of the calculated error correction term, update the change in weights & bias

$$\Delta w_{j,k} = \alpha \delta_k z_j \quad \Delta b_{ok} = \alpha \delta_k$$

Also send δ_k to the hidden layers backward

Step 7: Each hidden unit (z_j , $j=1 \dots p$) sums its delta I/P from the output units

$$\delta_{in,j} = \sum_{k=1}^m \delta_k w_{j,k}$$

The term $\delta_{in,j}$ gets multiplied with its derivative of $f(z_{in,j})$ to calculate the error term

$$\delta_j = \delta_{in,j} f'(z_{in,j})$$

The derivative $f'(z_{in,j})$ can be calculated depending on whether binary / bipolar sigmoidal function is used. On the basis of calculated δ_j , update the change in weight & bias:

$$\Delta v_{ij} = \alpha \delta_j x_i$$

$$\Delta v_{0j} = \alpha \delta_j$$

Weights & bias updation (Phase III):

Step 8: Each O/P unit (y_k , $k=1 \text{ to } m$) updates the bias & weights

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

$$w_{0k}(\text{new}) = w_{0k}(\text{old}) + \Delta w_{0k}$$

Each hidden unit (z_j , $j=1 \text{ to } p$) update it's bias & weights:

$$v_{ij}(\text{new}) = v_{ij}(\text{old}) + \Delta v_{ij}$$

$$v_{0j}(\text{new}) = v_{0j}(\text{old}) + \Delta v_{0j}$$

Step 9: Check for stopping condition. The stopping condition may be certain No of Epochs reached or when the Actual o/p equals the target o/p.

TESTING ALGO. IN BACK PROPAGATION NETWORK.

Step 0: Initialize the weights. The weights are taken from Training Algo.

Step 1: Perform steps 2-4 for each Input Vectors

Step 2: Set the Activation of Input unit x_i ($i=1 \text{ to } n$)

Step 3: Calculate the net input to hidden unit n and it's output for $j=1 \text{ to } p$

$$Z_{inj} = v_{0j} + \sum_{i=1}^n x_i v_{ij} \quad Z_j = f(Z_{inj})$$

Step 4: Now Compute the output of the output layer unit for $k=1 \text{ to } m$

$$Y_{ik} = w_{0k} + \sum_{j=1}^p Z_j w_{jk}$$

$$Y_k = f(Y_{ik})$$

Use Sigmoid Activation function for calculating the o/p.

FACTOR AFFECTING BACK PROPAGATION NETWORK

5

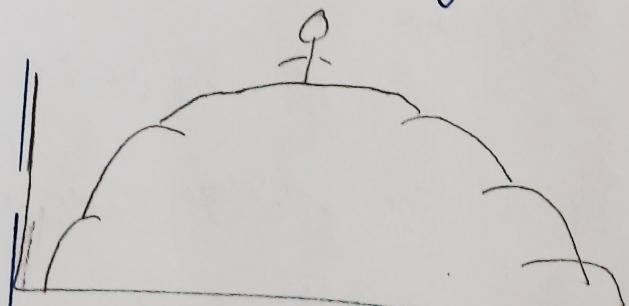
- * The Training of a BPN Based on the choice of Various Parameters
- * The Convergence of the BPN is Based on some Important Learning factors such as the Initial weight, The learning rate, The updation rule, The size & nature of Training set, & its Architecture.

1 Initial weights: * Initial weight are Initialized at small random value.

- * The choice of the initial weight determine how fast the N/w Converge
- * The Initial weight can not be very high because the Sigmoid Activation function used in BPN may get saturated from the beginning itself and the system may be stuck at a local minima or at a very flat plateau at the starting point.



Finding the lower point in a hilly landscape/ choose some small Random Value for weight Initially



Choose Large Random Value Initially

- * One Method of choosing the weights w_{ij} is choosing its Range

$$\left[\frac{-3}{\sqrt{\alpha_i}}, \frac{3}{\sqrt{\alpha_i}} \right]$$

Where α_i is the Number of processing element j that feed forward to processing element i .

- * The Initialization can also be done by a Method Called Nyugen-Widrow Initialization.
- * This type of Initialization leads to faster Convergence of Network.
- * This Method is used for Improving the learning ability of the hidden layers.
- * The Random Initialization of weights Connecting Input neurons to the hidden neuron is obtained by the equation

$$V_{ij} (\text{new}) = \gamma \frac{V_{ij} (\text{old})}{\| V_{ij} (\text{old}) \|}$$

V_{ij} is the Average weight calculated for all value of i

$$\gamma = 0.7(P)^{1/n}$$

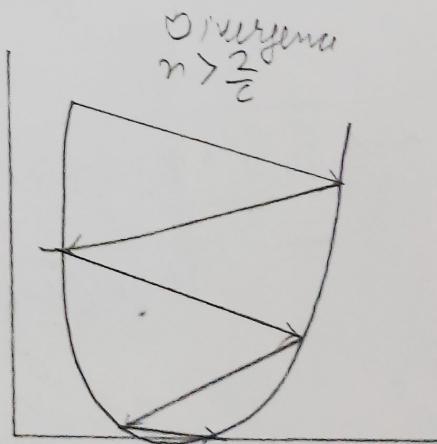
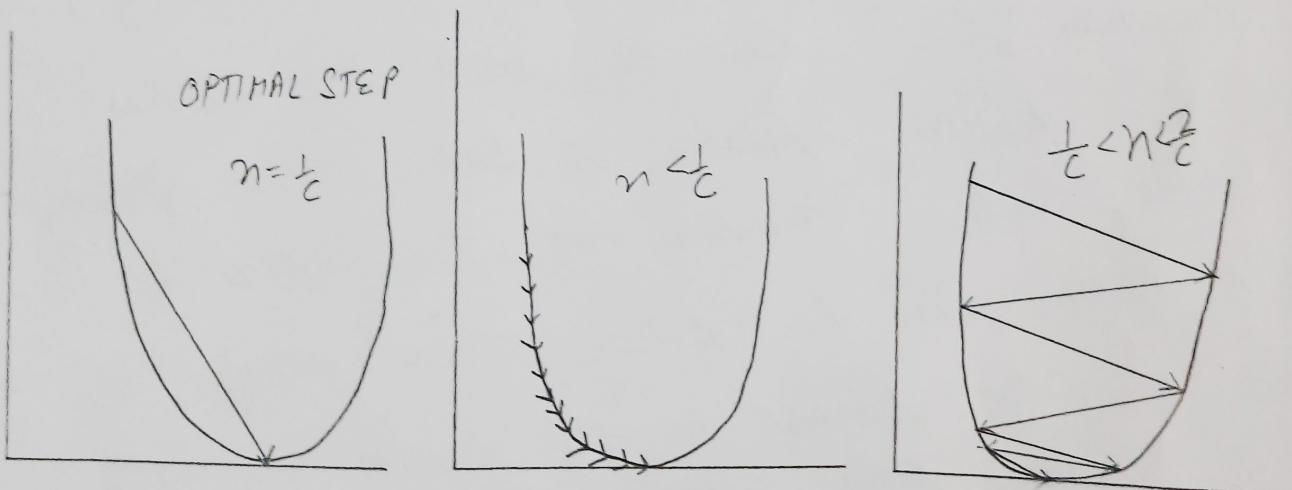
(n is the No of input neurons &

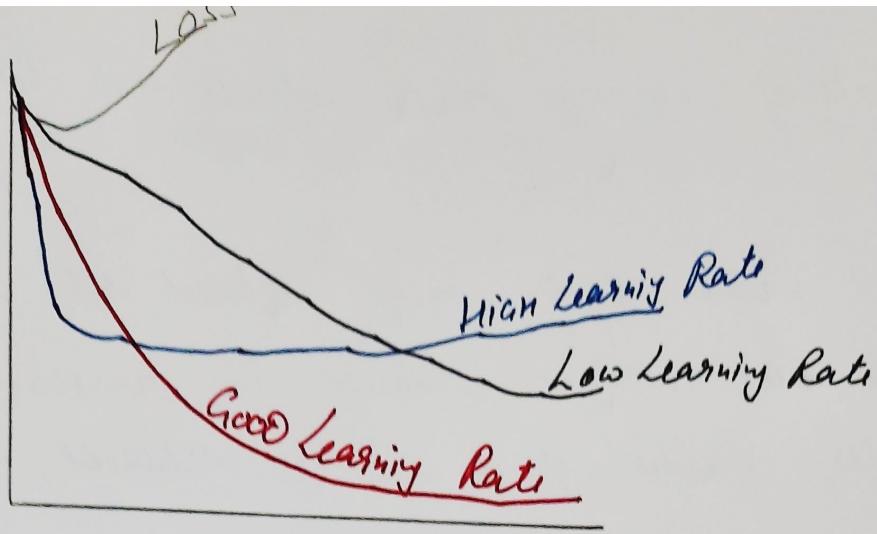
γ → scale factor P is the no of hidden neurons)

(6)

2 Learning Rate α : Learning rate affects the Convergence of the RPN.

- * A large Value of α may speed up the convergence but might result in overshooting the Minima.
- * Smaller Value of α may result undershoot.
- * The Range of α From 10^{-3} to 10 has been used successfully for several back propagation algorithm experiments.
- * Thus a large learning rate leads to rapid learning but there is an oscillation of weights. While the lower learning rate leads to slower learning.





3 Momentum factors : (n)

- * The Momentum factor is denoted by $n \in [0, 1]$ and the value of 0.9 is often used for the Momentum factor.
- * This approach is more useful when some training data are very different from the majority of data.
- * A Momentum factor can be used with either pattern by pattern updating or batch mode updating.
- * In Batch Mode, Momentum has the effect of complete averaging over the pattern. Even though the averaging is only partial in the pattern by pattern weight update.

$$w_{jk}(t+1) = w_{jk}(t) + \alpha \delta_k z_j + n [w_{jk}(t) - w_{jk}(t-1)]$$

$\overbrace{\qquad\qquad\qquad}^{\Delta w_{jk}(t+1)}$

and

$$v_{ij}(t+1) = v_{ij}(t) + \alpha \delta_i x_i + n [v_{ij}(t) - v_{ij}(t-1)]$$

$\overbrace{\qquad\qquad\qquad}^{\Delta v_{ij}(t+1)}$

Momentum also helps in fast convergence

4 GENERALIZATION

- * A Network is said to be Generalized when It Sensibly Interpolates with Input Networks that are new to the Network
- * When there are many trainable parameters for the given amount of training data, the Network learns but not Generalize & with small no of trainable parameters the N/w fails to learn the training data & perform very poorly on test data.
- * Solution to this problem is to monitor the errors on the test set & terminate when the errors increase
- * For Improving the ability of the Network to generalize from a training data set to a test data set.
- * It is desirable to make small changes in the I/p space of a pattern without changing the output components.
- * However Computationally, this method is very expensive
- * Small nets are preferred than large one

5 Number of Training data:

- * The training data should be sufficient & proper.
- * There exist a thumb rule, which states that the training data should cover the entire expected input space while training, training - vectors pair should be selected randomly from the set

- Assume, the Input space was being linearly separable into 'L' disjoint regions with their boundaries being part of Hyper plane.
- 'T' be the lower bound of the No of Training pattern
- Choosing T such that $T/L > 1$ will allow the network to discriminate pattern classes using fine piecewise hyperplane partitioning.

6 Number of Hidden Layer nodes:

- * The Number of Hidden units required for an application need to be determined separately.
- * The size of a hidden layer is usually determined experimentally.
- * If the network does not converge to a solution, it may need more hidden nodes on the other hand if the Net converges, the user may try a very few hidden nodes & then settle finally on a size based on overall system performance.

APPLICATIONS OF BACK PROPAGATION NETWORK

- * It is used in the field of speech Recognition
- * The neural network is trained to enunciate each letter of a word as a sentence
- * It is used in the field of character & face recognition.