

Perceptron Learning Rule

For the perceptron learning rule, the learning signal is the difference between the desired and actual neuron's response (Rosenblatt 1958). This example illustrates the perceptron learning rule of the network shown in Figure. The set of input training vectors is as follows where:

$$d_1 = -1, d_2 = -1, d_3 = 1, C = 0.1, \mathbf{W}^t = [1 \ -1 \ 0 \ 0.5]$$

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ -2 \\ 0 \\ -1 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1.5 \\ -0.5 \\ -1 \end{bmatrix}, \quad \mathbf{x}_3 = \begin{bmatrix} -1 \\ 1 \\ 0.5 \\ -1 \end{bmatrix}$$

Step 1 Input is \mathbf{x}_1 , desired output is d_1 :

$$net^1 = \mathbf{w}^{1t} \mathbf{x}_1 = [1 \ -1 \ 0 \ 0.5] \begin{bmatrix} 1 \\ -2 \\ 0 \\ -1 \end{bmatrix} = 2.5$$

Correction in this step is necessary since $d_1 \neq \text{sgn}(2.5)$. We thus obtain updated weight vector

$$\mathbf{w}^2 = \mathbf{w}^1 + 0.1(-1 - 1)\mathbf{x}_1$$

Plugging in numerical values we obtain

$$\mathbf{w}^2 = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0.5 \end{bmatrix} - 0.2 \begin{bmatrix} 1 \\ -2 \\ 0 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.8 \\ -0.6 \\ 0 \\ 0.7 \end{bmatrix}$$

Step 2 Input is \mathbf{x}_2 , desired output is d_2 . For the present weight vector \mathbf{w}^2 we compute the activation value net^2 as follows:

$$net^2 = \mathbf{w}^{2t} \mathbf{x}_2 = [0 \ 1.5 \ -0.5 \ -1] \begin{bmatrix} 0.8 \\ -0.6 \\ 0 \\ 0.7 \end{bmatrix} = -1.6$$

Correction is not performed in this step since $d_2 = \text{sgn}(-1.6)$

Step 3 Input is \mathbf{x}_3 , desired output is d_3 , present weight vector is \mathbf{w}^2 . Computing net^3 we obtain:

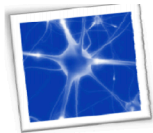
$$net^3 = \mathbf{w}^{2t} \mathbf{x}_3 = [-1 \ 1 \ 0.5 \ -1] \begin{bmatrix} 0.8 \\ -0.6 \\ 0 \\ 0.7 \end{bmatrix} = -2.1$$

Correction is necessary in this step since $d_3 \neq \text{sgn}(-2.1)$. The updated weight values are

$$\mathbf{w}^4 = \mathbf{w}^3 + 0.1(1 + 1)\mathbf{x}_3$$

or

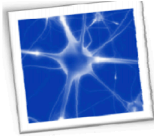
$$\mathbf{w}^4 = \begin{bmatrix} 0.8 \\ -0.6 \\ 0 \\ 0.7 \end{bmatrix} + 0.2 \begin{bmatrix} -1 \\ 1 \\ 0.5 \\ -1 \end{bmatrix} = \begin{bmatrix} 0.6 \\ -0.4 \\ 0.1 \\ 0.5 \end{bmatrix}$$



Example

- Consider the 2-dimensional training set $C_1 \cup C_2$,
- $C_1 = \{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$ with class label 1
- $C_2 = \{(1, -1, -1), (1, -1, 1), (1, 0, 1)\}$ with class label 0
- **Train a perceptron on $C_1 \cup C_2$**
- First pass is as follows

Input	Weight	Desired	Actual	Update?	New weight
(1, 1, 1)	(1, 0, 0)	1	1	No	(1, 0, 0)
(1, 1, -1)	(1, 0, 0)	1	1	No	(1, 0, 0)
(1, 0, -1)	(1, 0, 0)	1	1	No	(1, 0, 0)
(1, -1, -1)	(1, 0, 0)	0	1	Yes	(0, 1, 1)
(1, -1, 1)	(0, 1, 1)	0	0	No	(0, 1, 1)
(1, 0, 1)	(0, 1, 1)	0	1	Yes	(-1, 1, 0)



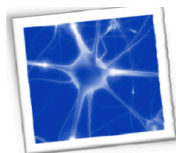
Example

C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$

C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$

Fill out this table sequentially (Second pass):

Input	Weight	Desired	Actual	Update?	New weight
(1, 1, 1)	(-1, 1, 0)	1	0	Yes	(0, 2, 1)
(1, 1, -1)	(0, 2, 1)	1	1	No	(0, 2, 1)
(1,0, -1)	(0, 2, 1)	1	0	Yes	(1, 2, 0)
(1,-1, -1)	(1, 2, 0)	0	0	No	(1, 2, 0)
(1,-1, 1)	(1, 2, 0)	0	0	No	(1, 2, 0)
(1, 0, 1)	(1, 2, 0)	0	1	Yes	(0, 2, -1)



Example

C1: $\{(1, 1, 1), (1, 1, -1), (1, 0, -1)\}$

C2: $\{(1, -1,-1), (1, -1,1), (1, 0,1)\}$

Fill out this table sequentially (Third pass):

Input	Weight	Desired	Actual	Update?	New weight
(1, 1, 1)	(0, 2, -1)	1	1	No	(0, 2, -1)
(1, 1, -1)	(0, 2, -1)	1	1	No	(0, 2, -1)
(1,0, -1)	(0, 2, -1)	1	1	No	(0, 2, -1)
(1,-1, -1)	(0, 2, -1)	0	0	No	(0, 2, -1)
(1,-1, 1)	(0, 2, -1)	0	0	No	(0, 2, -1)
(1, 0, 1)	(0, 2, -1)	0	0	No	(0, 2, -1)

At epoch 3 no weight changes.

⇒ stop execution of algorithm.

Final weight vector: (0, 2, -1).

⇒ decision hyperplane is $2x_1 - x_2 = 0$.

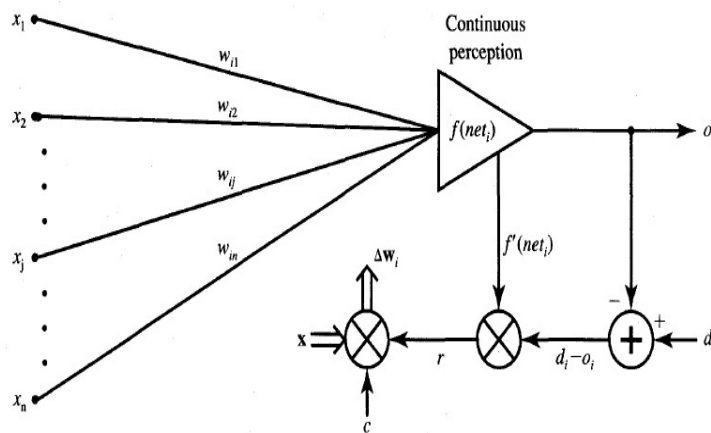
Delta Learning Rule

The delta learning rule is only valid for continuous activation functions as defined before, and in the supervised training mode. The adjustment of learning weight:

$$x_j$$

$$\Delta w_{ij} = [d_i - f(\mathbf{w}_i^t \mathbf{x})] f'(\mathbf{w}_i^t \mathbf{x})$$

for this rule is called delta and is defined as follows



This example discusses the delta learning rule as applied to the network shown in Figure. Training input vectors, desired responses, and initial weights are identical to those in Example. The delta learning requires

$$f(\text{net}) = \frac{2}{1 + e^{-\lambda \text{net}}} - 1$$

that the value $f(\text{net})$ be computed in each step.

$$f'(\text{net}) = \frac{1}{2}(1 - o^2)$$

Example/

$$d_1 = -1, d_2 = -1, d_3 = 1, C = 0.1, \lambda = 1$$

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ -2 \\ 0 \\ -1 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 0 \\ 1.5 \\ -0.5 \\ -1 \end{bmatrix}, \quad \mathbf{x}_3 = \begin{bmatrix} -1 \\ 1 \\ 0.5 \\ -1 \end{bmatrix}$$

Step 1 Input is vector \mathbf{x}_1 , initial weight vector is \mathbf{w}^1 :

$$\begin{aligned} net^1 &= \mathbf{w}^{1t} \mathbf{x}_1 = 2.5 \\ o^1 &= f(net^1) = 0.848 \\ f'(net^1) &= \frac{1}{2}[1 - (o^1)^2] = 0.140 \\ \mathbf{w}^2 &= c(d_1 - o^1)f'(net^1)\mathbf{x}_1 + \mathbf{w}^1 \\ &= [0.974 \quad -0.948 \quad 0 \quad 0.526]^t \end{aligned}$$

Step 2 Input is vector \mathbf{x}_2 , weight vector is \mathbf{w}^2 :

$$\begin{aligned} net^2 &= \mathbf{w}^{2t} \mathbf{x}_2 = -1.948 \\ o^2 &= f(net^2) = -0.75 \\ f'(net^2) &= \frac{1}{2}[1 - (o^2)^2] = 0.218 \\ \mathbf{w}^3 &= c(d_2 - o^2)f'(net^2)\mathbf{x}_2 + \mathbf{w}^2 \\ &= [0.974 \quad -0.956 \quad 0.002 \quad 0.531]^t \end{aligned}$$

Step 3 Input is x_3 , weight vector is \mathbf{w}^3 :

$$\begin{aligned} net^3 &= \mathbf{w}^{3t} \mathbf{x}_3 = -2.46 \\ o^3 &= f(net^3) = -0.842 \\ f'(net^3) &= \frac{1}{2}[1 - (o^3)^2] = 0.145 \\ \mathbf{w}^4 &= c(d_3 - o^3)f'(net^3)\mathbf{x}_3 + \mathbf{w}^3 \\ &= [0.947 \quad -0.929 \quad 0.016 \quad 0.505]^t \end{aligned}$$

Widrow-Hoff Learning Rule

The Widrow-Hoff learning rule (Widrow 1962) or Adaline is applicable for the supervised training of neural networks. It is independent of the activation function of neurons used since it minimizes the squared error between the desired output value d , and the neuron's activation value, $net = \mathbf{w}^t \mathbf{x}$. The learning signal for this rule is defined as follows:

$$r \triangleq d_i - \mathbf{w}_i^t \mathbf{x}$$

The weight vector increment under this learning rule is

$$\Delta \mathbf{w}_i = c(d_i - \mathbf{w}_i^t \mathbf{x}) \mathbf{x}$$

Correlation Learning Rule

By substituting $r = d_i$ into the general learning rule we obtain the correlation learning rule. The adjustments for the weight vector and the single weights, respectively are

$$\Delta \mathbf{w}_i = c d_i \mathbf{x}$$
$$\Delta w_{ij} = c d_i x_j, \quad \text{for } j = 1, 2, \dots, n$$

Multi-layer ANN and Backpropagation Learning Rule

All the aforementioned learning rules work with single neuron ANN but don't work with multilayer ANN. A multilayer perceptron (shown in Fig. below) is a feedforward neural network with one or more hidden layers. Typically, the network consists of an input layer of source neurons, at least one middle or hidden layer of computational neurons, and an output layer of computational neurons. The input signals are propagated in a forward direction on a layer-by-layer basis.

Each layer in a multilayer neural network has its own specific function. The input layer accepts input signals from the outside world and redistributes these signals to all neurons in the hidden layer. Actually, the input layer rarely includes computing neurons, and thus does not process input patterns. The output layer accepts output signals, or in other words a stimulus pattern, from the hidden layer and establishes the output pattern of the entire network.

Neurons in the hidden layer detect the features; the weights of the neurons represent the features hidden in the input patterns. These features are then used by the output layer in determining the output pattern. With one hidden layer, we can represent any continuous function of the input signals, and with two hidden layers even discontinuous functions can be represented.

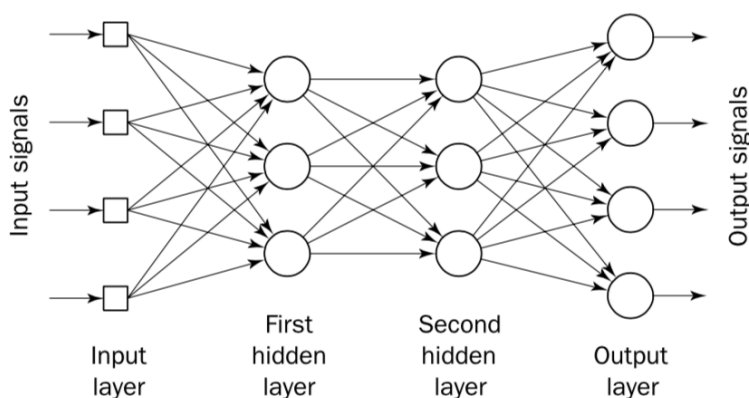


Figure: Multilayer perceptron with two hidden layers

The determination of the error is a recursive process which start with the o/p units and the error is back propagated to the I/p units. Therefore the rule is called error Back propagation (EBP) or simply Back Propagation (BP). The weight is changed exactly in the same form of the standard DR

$$\Delta w_{ij} = \xi \delta_j x_i$$

$$\Rightarrow w_{ij}(t+1) = w_{ij}(t) + \xi \delta_j x_i$$

There are two other equations that specify the error signal. If a unite is an o/p unit, the error signal is given by:-

$$\delta = (d_j - y_j) f_j(\text{net } j)$$

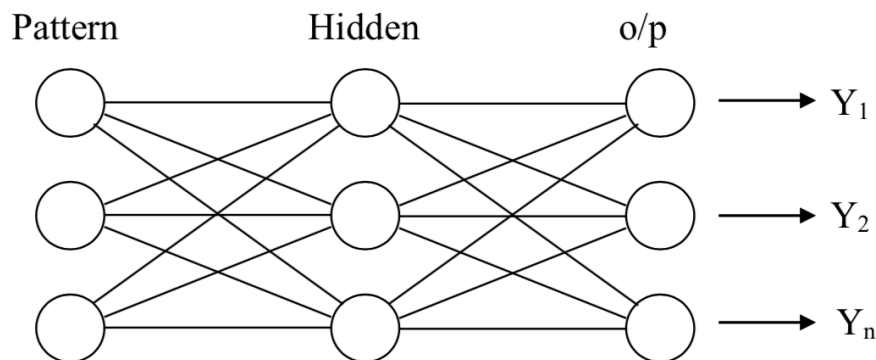
$$\text{Where } \text{net } j = \sum w_{ij} x_i + \theta$$

The GDR minimize the squares of the differences between the actual and the desired o/p values summed over the o/p unit and all pairs of I/p and o/p vectors. The rule minimize the overall error $E = \sum E_p$ by implementing a gradient descent in E: - where, $E_p = 1/2 \sum_j (d_j - y_j)^2$.

The BP consists of two phases:-

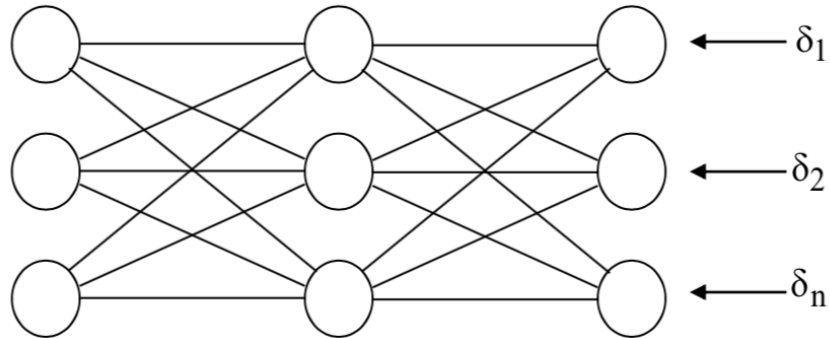
1- Forward Propagation:-

During the forward phase, the I/p is presented and propagated towards the o/p.



2- Backward Propagation:-

During the backward phase, the errors are formed at the o/p and propagated towards the I/p



3- Compute the error in the hidden layer.

$$\text{If } y = f(x) = \frac{1}{1 + e^{-x}}$$

$$f' = y(1 - y)$$

Equation is can rewrite as:-

$$\delta_j = y(1 - y)(d_j - y_j)$$

The error signal for hidden units for which there is no specified target (desired o/p) is determined recursively in terms of the error signals of the units to which it directly connects and the weights of those connections:-

That is

$$\delta_j = f'(\text{net}_j) \sum_k \delta_k w_{ik}$$

Or

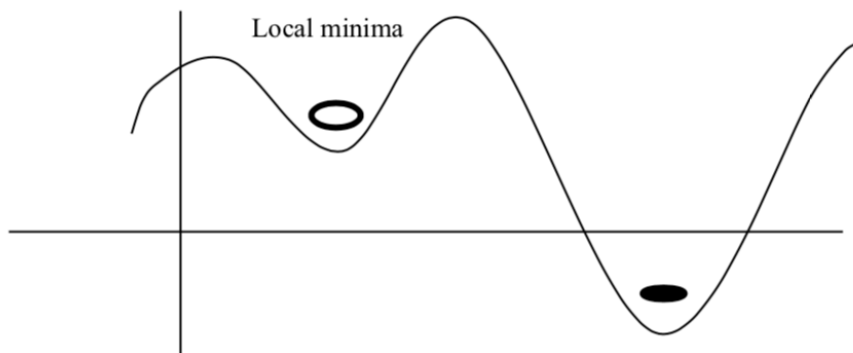
$$\delta_j = y_j(1 - y_j) \sum_k \delta_k w_{ik}$$

B.P learning is implemented when hidden units are embedded between input and output units.

Convergence

A quantitative measure of the learning is the :Root Mean Square (RMS) error which is calculated to reflect the "degree" of learning.

Generally, an RMS bellow (0.1) indicates that the net has learned its training set. Note that the net does not provide a yes /no response that is "correct" or "incorrect" since the net get closer to the target value incrementally with each step. It is possible to define a cut off point when the nets o/p is said to match the target values.



- Convergence is not always easy to achieve because sometimes the net gets stuck in a "Local minima" and stops learning algorithm.
- Convergence can be represented intuitively in terms of walking about mountains.

Momentum term

The choice of the learning rate plays important role in the stability of the process. It is possible to choose a learning rate as large as possible without leading to oscillations. This offers the most rapid learning. One way to increase the learning rate without leading to oscillations is to modify the GDR to include momentum term.

This can be achieved by the following rule:-

$$W_{ij}(t+1) = W_{ij}(t) + \zeta \delta_j x_i + \alpha (W_{ij}(t) - W_{ij}(t-1))$$

Where α ($0 < \alpha < 1$) is a constant which determines the effect of the past weight changes on the current direction of movement in weight space.

A "global minima" unfortunately it is possible to encounter a local minima, a valley that is not the lowest possible in the entire terrain. The net does not leave a local minima by the standard BP algorithm and special techniques should be used to get out of a local minima such as:-

- 1- Change the learning rate or the momentum term.
- 2- Change the no. of hidden units (10%).
- 3- Add small random value to the weights.
- 4- Start the learning again with different initial weights.

3.1.3.1 Back Propagation Training Algorithm

Training a network by back propagation involves three stages:-

- 1-the feed forward of the input training pattern
- 2-the back propagation of the associated error
- 3-the adjustment of the weights

let n = number of input units in input layer,

let p = number of hidden units in hidden layer

let m = number of output units in output layer

let V_{ij} be the weights between i/p layer and the hidden layer,

let W_{ij} be the weights between hidden layer and the output layer,

we refer to the i/p units as X_i , $i=1, 2, \dots, n$. and we refer to the hidden units as Z_j , $j=1, \dots, p$. and we refer to the o/p units as y_k , $k=1, \dots, m$.

δ_{1j} is the error in hidden layer,

δ_{2k} is the error in output layer,

ζ is the learning rate

α is the momentum coefficient (learning coefficient, $0.0 < \alpha < 1.0$,

y_k is the o/p of the net (o/p layer),

Z_j is the o/p of the hidden layer,

X_i is the o/p of the i/p layer.

η is the learning coefficient.

The algorithm is as following :-

Step 0 : initialize weights (set to small random value).

Step 1 : while stopping condition is false do steps 2-9

Step 2: for each training pair, do steps 3-8

Feed forward :-

Step 3:- Each i/p unit (X_i) receives i/p signal X_i & broad casts this signal to all units in the layer above (the hidden layer)

Step 4:- Each hidden unit (Z_j) sums its weighted i/p signals,

$$Z - \text{inj} = V_{aj} + \sum_{i=1}^n x_i v_{ij} \quad (V_{aj} \text{ is abias})$$

and applies its activation function to compute its output signal (the activation function is the binary sigmoid function),

$$Z_j f(Z - \text{inj}) = 1 / (1 + \exp - (Z - \text{inj}))$$

and sends this signal to all units in the layer above (the o/p layer).

Step 5:- Each output unit (Y_k)sums its weighted i/p signals,

$$y - \text{ink} = w_{ok} + \sum_{j=1}^p Z_j w_{jk} \quad (\text{where } w_{ok} \text{ is abias})$$

and applies its activation function to compute its output signal.

$$y_k = f(y - \text{ink}) = 1 / (1 + \exp - (y - \text{ink}))$$

back propagation of error:-

step 6 : Each output unit (y_k , $k= 1$ to m) receive a target pattern corresponding to the input training pattern, computes its error information term and calculates its weights correction term used to update W_{jk} later,

$$\delta_{2k} = y_k(1 - y_k) * (T_k - y_k),$$

where T_k is the target pattern & $k=1$ to m .

step 7 : Each hidden unit (Z_j , $j= 1$ to p) computes its error information term and calculates its weight correction term used to update V_{ij} later,

$$\delta_{1j} = Z_j * (1 - Z_j) * \sum_{k=1}^m \delta_{2k} W_{jk}$$

Update weights and bias :-

step 8: Each output unit (y_k , $k = 1$ to m) updates its bias and weights:

$$W_{jk}(\text{new}) = \eta * \delta_{2k} * Z_j + \alpha * [W_{jk}(\text{old})],$$

$j= 1$ to p

Each hidden unit (Z_j , $j= 1$ to p) update its bias and weights:

$$V_{ij}(\text{new}) = \eta * \delta_{1j} * X_i + \alpha [v_{ij}(\text{old})],$$

$I = 1$ to n

Step 9 : Test stopping condition.

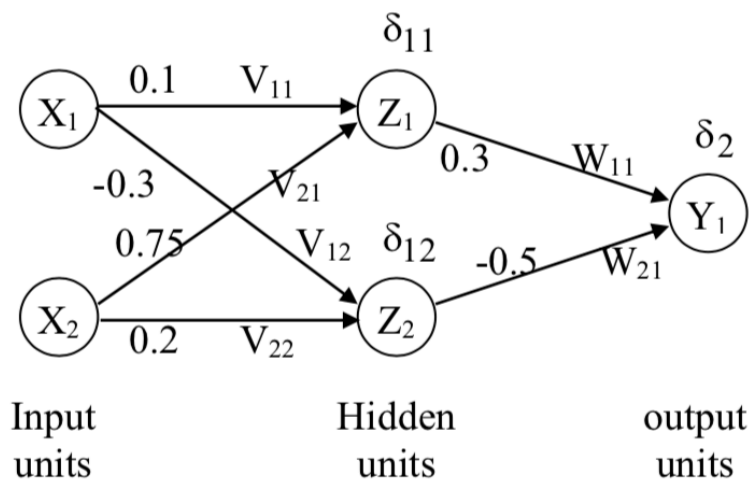
EX6

Suppose you have BP- ANN with 2-input , 2-hidden , 1-output nodes with sigmoid function and the following matrices weight, trace with 1-iteration.

$$V = \begin{bmatrix} 0.1 & -0.3 \\ 0.75 & 0.2 \end{bmatrix} \quad w = [0.3 \quad -0.5]$$

Where $\alpha = 0.9$, $\eta = 0.45$, $x = (1,0)$, and $T_k = 1$

Solution:-



1-Forward phase :-

$$Z - \text{in}1 = X_1 V_{11} + X_2 V_{21} = 1 * 0.1 + 0 * 0.75 = 0.1$$

$$Z - \text{in}2 = X_1 V_{12} + X_2 V_{22} = 1 * -0.3 + 0 * 0.2 = -0.3$$

$$Z_1 = f(Z - \text{in}1) = 1 / (1 + \exp - (Z - \text{in}1)) = 0.5$$

$$Z_2 = f(Z - \text{in}2) = 1 / (1 + \exp - (Z - \text{in}2)) = 0.426$$

$$y - \text{in}1 = Z_1 W_{11} + Z_2 W_{21}$$

$$= 0.5 * 0.3 + 0.426 * (-0.5) = -0.063$$

$$y_1 = f(y - \text{in}1) = 1 / (1 + \exp - (y - \text{in}1)) = 0.484$$

2-Backward phase :-

$$\delta_{2k} = y_k(1 - y_k) * (T_k - y_k)$$

$$\delta_{21} = 0.484(1 - 0.484) * (1 - 0.484)0.129$$

$$\delta_{1j} = Z_j * (1 - Z_j) * \sum_{k=1}^m \delta_{2k} W_{jk}$$

$$\begin{aligned} \delta_{11} &= Z_1(1 - Z_1) * (\delta_{21} W_{11}) \\ &= 0.5(1 - 0.5) * (0.129 * 0.3) = 0.0097 \end{aligned}$$

$$\begin{aligned} \delta_{12} &= Z_2(1 - Z_2) * (\delta_{21} W_{21}) \\ &= 0.426(1 - 0.426) * (0.129 * (-0.5)) = -0.015 \end{aligned}$$

3-Update weights:-

$$W_{jk}(\text{new}) = \eta * \delta_{2k} * Z_j + \alpha * [W_{jk}(\text{old})]$$

$$\begin{aligned} W_{11} &= \eta * \delta_{21} * Z_1 + \alpha * [W_{11}(\text{old})] \\ &= 0.45 * 0.129 * 0.5 + 0.9 * 0.3 = 0.299 \end{aligned}$$

$$\begin{aligned} W_{21} &= \eta * \delta_{21} * Z_2 + \alpha * [W_{21}(\text{old})] \\ &= 0.45 * 0.129 * 0.426 + 0.9 * -0.5 = -0.4253 \end{aligned}$$

$$V_{ij}(\text{new}) = \eta * \delta_{1j} * X_i + \alpha * [V_{ij}(\text{old})]$$

$$\begin{aligned} V_{11} &= \eta * \delta_{11} * X_1 + \alpha * [V_{11}(\text{old})] \\ &= 0.45 * 0.0097 * 1 + 0.9 * 0.1 = 0.0944 \end{aligned}$$

$$\begin{aligned} V_{12} &= \eta * \delta_{12} * X_1 + \alpha * [V_{12}(\text{old})] \\ &= 0.45 * 0.0158 * 1 + 0.9 * -0.3 = -0.2771 \end{aligned}$$

$$\begin{aligned} V_{21} &= \eta * \delta_{11} * X_2 + \alpha * [V_{21}(\text{old})] \\ &= 0.45 * 0.0097 * 0 + 0.9 * 0.75 = 0.675 \end{aligned}$$

$$\begin{aligned} V_{22} &= \eta * \delta_{12} * X_2 + \alpha * [V_{22}(\text{old})] \\ &= 0.45 * -0.0158 * 0 + 0.9 * 0.2 = 0.18 \end{aligned}$$

$$\therefore V = \begin{bmatrix} 0.0944 & -0.2771 \\ 0.675 & 0.18 \end{bmatrix} \quad W = \begin{bmatrix} 0.299 & -0.4253 \end{bmatrix}$$