

## Multi-layer Feedforward Neural Network (Back Propagation Algorithm)

- Error corrects the weights

$$Net = w_i^T x, \quad f(net) = \frac{1}{1+e^{-x}}$$

$$\frac{\partial ou}{\partial net} = f'(net) = out(1 - out)$$

- Training pattern  $(\bar{x}_p, \bar{t}_p)$

$$\bar{x}_p = \begin{bmatrix} x_{p1} \\ \vdots \\ x_{pn} \end{bmatrix}; \bar{t}_p = \begin{bmatrix} t_{p1} \\ \vdots \\ t_{pn} \end{bmatrix}$$

- For node in Hidden layer j, the net input is

$$Net_j = \sum_i w_{ij} o_i \quad \text{Where, } o_i = \text{Output of node } i$$

$w_{ij}$  = Connection weight  
from neuron i to neuron j

- The output of node j is

$$O_j = f(net_j) \quad (1)$$

Where f is the activation function

- For sigmoidal activation function

$$O_j = \frac{1}{1+e^{-net_j}}; \quad O_j = \frac{1}{1+e^{\left(\frac{-net_j - \theta_j}{\theta_o}\right)}}; \quad O_j = \frac{2}{1+e^{-net_j}} - 1$$

This is for Bipolar  
sigmoidal function

- The output of node k is

$$O_k = f(net_k) \quad (2)$$

Where f is the activation function

- For sigmoidal activation function

$$O_k = \frac{1}{1+e^{-net_k}}; \quad O_k = \frac{1}{1+e^{\left(\frac{-net_k - \theta_k}{\theta_o}\right)}}; \quad O_k = \frac{2}{1+e^{-net_k}} - 1$$

This is for Bipolar  
sigmoidal function

- For nodes in output layer k, the net input is

$$Net_k = \sum_j w_{jk} o_j \quad (3)$$

Where,  $O_j$  = Output of node j

$w_{jk}$  = Connection weight from neuron j to neuron k

- Compare the final output signals with a target signals then adjust weights between unit so that the difference between the final and target signals is minimized. The weights are changed to minimize the following cost function

$$E_p = \frac{1}{2} \sum_j (t_{pj} - o_{pj})^2 \quad (4)$$

Where,  $t_{pj}$  = target signal at unit j for input pattern p

$O_{pj}$  = output signal of unit j

- Error may be +ve or -ve
- Weights between neurons are determined as follows:

The **Steepest Decent Method** is utilized in order to minimize Eq. (4)

- Back Propagation Algorithm, attempts to minimize error by adjusting each weight of the network proportional to the derivative of the error with respect to that weight.
- The change in weights from unit ( i ) to ( j ) for input pattern ' p ' is given by

$$\Delta_p w_{ij} \propto - \left( \frac{\partial E_p}{\partial w_{ij}} \right)$$

- The change in proportional to gradient error

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

$$\Delta w = \eta (t - o) x \quad \leftarrow \text{input} \quad [ \text{Widrow-Hoff Rule} ]$$

$\eta \rightarrow$  Learning Rate

(modifies the rate of convergence)

Error +ve  $\rightarrow$  ADD

Error -ve  $\rightarrow$  SUBTRACT

- The change in weights from unit ( k ) to ( j )

$$\Delta w_{kj} \propto - \left( \frac{\partial E}{\partial w_{kj}} \right)$$

(-ve indicates minimization of error)

$$\Delta w_{kj} = -\eta \left( \frac{\partial E}{\partial w_{kj}} \right) \quad (5)$$

From Eq.(5)

$$\begin{aligned} - \left( \frac{\partial E}{\partial w_{kj}} \right) &= - \left( \frac{\partial E}{\partial w_{kj}} \right) \cdot \left( \frac{\partial net_k}{\partial net_k} \right) \\ &= - \left( \frac{\partial E}{\partial net_k} \right) \cdot \left( \frac{\partial net_k}{\partial w_{kj}} \right) \end{aligned} \quad (6)$$

$$E_p = \frac{1}{2} \sum_k (t_{pk} - O_{pk})^2$$

system error

$$E = \frac{1}{2} \sum_{p=1}^p \left[ \sum_k (t_{pk} - O_{pk})^2 \right]$$

$$\frac{\partial E}{\partial O_k} = -(t_k - O_k) \quad (7)$$

From Eq. (3)

$$net_k = w_k^t O_j$$

$$\frac{\partial net_k}{\partial w_{kj}} = O_j \quad (8)$$

From Eq. (6)

$$\delta_k = - \left( \frac{\partial E}{\partial net_k} \right) \quad (9)$$

$$= - \left( \frac{\partial E}{\partial net_k} \right) \cdot \left( \frac{\partial O_k}{\partial O_k} \right)$$

$$= - \left( \frac{\partial E}{\partial O_k} \right) \cdot \left( \frac{\partial O_k}{\partial net_k} \right)$$

From Eq. (6) and derived Eq. (2)

$$\delta_k = (t_k - O_k) \cdot f'(net_k) \quad (10)$$

Put Eq. (8) and (9) in Eq. (6)

$$= [(t_k - O_k) \cdot f'(net_k)] \cdot O_j$$

- Delta rule, it is valid for continuous activation function

Put Eq. (6) in Eq. (5)

$$\Delta w_{kj} = \eta [(t_k - O_k) \cdot f'(net_k)] \cdot O_j$$

$$\Delta w_{kj} = \eta \cdot \delta_k \cdot O_j$$

## Internal Layer

$$\Delta w_{ji} = -\eta \left( \frac{\partial E}{\partial w_{ji}} \right) \quad (11)$$

$$= -\eta \left( \frac{\partial E}{\partial w_{ji}} \right) \cdot \left( \frac{\partial net_j}{\partial net_j} \right)$$

$$= -\eta \left( \frac{\partial E}{\partial net_j} \right) \cdot \left( \frac{\partial net_j}{\partial w_{ji}} \right)$$

Similarly, as Eq. (8)

$$= -\eta \left( \frac{\partial E}{\partial net_j} \right) \cdot O_i$$

$$= -\eta \left( \frac{\partial E}{\partial net_j} \right) \cdot \left( \frac{\partial O_j}{\partial O_j} \right) \cdot O_i$$

$$= -\eta \left( \frac{\partial E}{\partial O_j} \right) \cdot \left( \frac{\partial O_j}{\partial net_j} \right) \cdot O_i$$

Similarly, as Eq. (10)

$$= -\eta \left( \frac{\partial E}{\partial O_j} \right) \cdot f'(net_j) \cdot O_i \quad (12)$$

From Eq. (12)

$$\delta_j = - \left( \frac{\partial E}{\partial O_j} \right) \cdot f'(net_j) \quad (13)$$

From Eq. (13)

$$- \left( \frac{\partial E}{\partial O_j} \right) = - \sum_k \left( \frac{\partial E}{\partial O_j} \right) \cdot \left( \frac{\partial net_k}{\partial net_k} \right)$$

$$= - \sum_k \left( \frac{\partial E}{\partial net_k} \right) \cdot \left( \frac{\partial net_k}{\partial O_j} \right)$$

Where,  $net_k = w_{kj} O_j$

$$= \sum_k \left( - \frac{\partial E}{\partial net_k} \right) \cdot \left( \frac{\partial \sum w_{kj} O_j}{\partial O_j} \right)$$

$$= \sum_k \left( - \frac{\partial E}{\partial net_k} \right) \cdot w_{kj}$$

$$= \sum_k (-(-\delta_k)) \cdot w_{kj} \quad (14)$$

From Eq. (9)

- Direction  $\delta_k$  has been reversed in this case, i.e.,  $\delta_k \rightarrow -\delta_k$

Put Eq. (14) in Eq. (13)

$$\delta_j = \sum_k \delta_k \cdot w_{kj} \cdot f'(net_j) \quad (15)$$

From Eq. (12), Eq. (13) and Eq. (15)

$$\Delta w_{ji} = \eta \left[ \sum_k \delta_k \cdot w_{kj} \cdot f'(net_j) \right] \cdot O_i$$

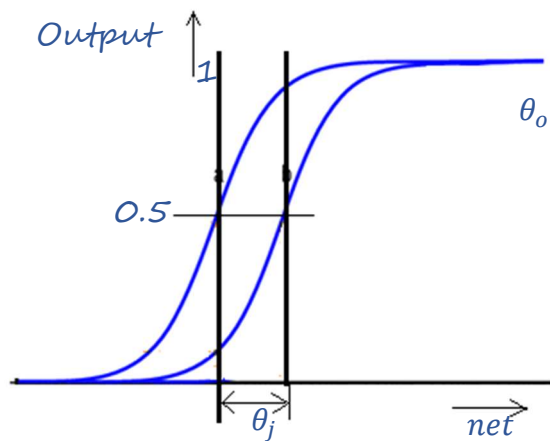
$$\Delta w_{ji} = \eta \left[ \sum_k \delta_k \cdot w_{kj} \cdot O_j (1 - O_j) \right] \cdot O_i$$

- BPA  $\rightarrow$  Propagating the output layer through the networks layer by layer, adjusting weight at each layer
- However, for hidden layers  $\delta$  must be generalised without benefit of a target vector.
- First  $\delta$  for each neuron in the output layer. It is used to adjust the weight feeding into the output layer, then it is propagated back through the same weights to generate a value. For  $\delta$ , for each neuron in the first hidden layer. These values of  $\delta$  are used, in turn to adjust the weight of this hidden layer and in a similar way, are propagate  $\delta$  back to all the preceding layers.
- Consider a single neuron in the hidden layer just before the output layer. In the forward pass, this neuron propagates in output value to neurons in the output layer through the inter connecting weights.
- During training these weights operate in reverse, passing the value of  $\delta$  from the output layer back to the hidden layer. Each of these weights multiplied by the  $\delta$  value of the neuron to which it connects in the output layer.
- The value of  $\delta$  needed for the hidden layer neuron is produced by summing all such products and multiplying by the derivative of the squashing function.

$$\delta_{pj} = out_{p,j}(1 - out_{p,j}) \left( \sum_q \delta_{q,k} \cdot w_{pq,k} \right)$$

## Comments or Observations

1. The sigmoidal function never becomes 1 or 0 unless function goes  $+\infty$  to  $-\infty$ . Therefore, the weights should be  $0 < w < 1$
2. Stopping values chosen randomly
3. If  $\eta$  is very large, it is oscillatory
4. Add a momentum term ( $\alpha$ ) such that it will not exchange too fast  
Large  $\alpha \rightarrow$  Convergence is slow  
$$\Delta w_{ji}(n+1) = \eta \delta_j O_i + \alpha \Delta w_{ji}(n)$$
5. Local minima
6.  $\theta_j$  could be viewed as equivalent weight



$\theta_j \rightarrow$  Threshold

$\theta_0 \rightarrow$  Slope

7. 
$$f'(net_j) = \frac{\partial O_j}{\partial net_j} = O_j(1 - O_j)$$