

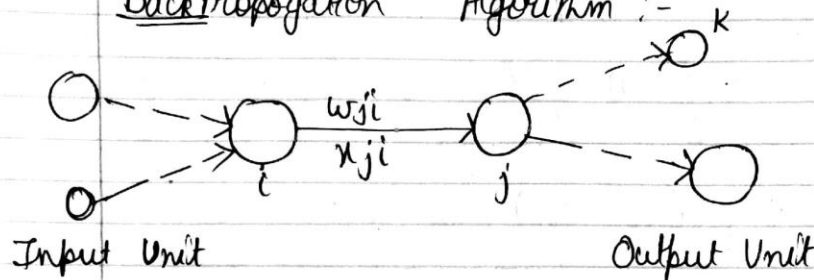
Ques 1.1

Q1.1

a) Tanh activation function:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Backpropagation Algorithm :-



Consider above network with two intermediate nodes 'i' and 'j'.

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

change in weight (Δw_{ji}) [w_{ji} is the weight from i^{th} node to j^{th} node].

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

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \cdot \frac{\partial \text{net}_j}{\partial w_{ji}}$$



$$\frac{\delta_{netj}}{\delta w_{ji}} = x_{ji}$$

For $\frac{\delta E_d}{\delta_{netj}}$, there are 2 cases:-

- (1) J is an Output unit
- (2) J is an hidden unit.

(1) J is an Output unit.

$$\frac{\delta E_d}{\delta_{netj}} = \frac{\delta E_d}{\delta o_j} \frac{\delta o_j}{\delta_{netj}}$$

Error (E_d) is

$$E_d(w) = \frac{1}{2} \sum_{R \in \text{Output}} (t_R - o_R)^2$$

Differentiate wrt o_j :

$$\frac{\delta E_d}{\delta o_j} = \frac{\delta}{\delta o_j} \left[\frac{1}{2} \sum_{R \in \text{Output}} (t_R - o_R)^2 \right]$$

$$= \frac{\delta}{\delta o_j} \left[\frac{1}{2} (t_j - o_j)^2 \right] = -(t_j - o_j)$$

for activation function tanh.

$$\frac{\delta o_j}{\delta_{netj}} = \frac{\delta}{\delta_{netj}} \left[\frac{e^x - e^{-x}}{e^x + e^{-x}} \right]$$

$$= \frac{(e^x + e^{-x})^2 - (e^x - e^{-x})^2}{(e^x + e^{-x})^2}$$

$$= 1 - \left(\frac{e^x - e^{-x}}{e^x + e^{-x}} \right)^2$$

$$\frac{\delta o_j}{\delta net_j} = 1 - \tanh^2(x) = 1 - o_j^2$$

$$\begin{aligned} \frac{\delta E_d}{\delta net_j} &= \frac{\delta E_d}{\delta o_j} \cdot \frac{\delta o_j}{\delta net_j} \\ &= -(t_j - o_j) (1 - o_j^2) \end{aligned}$$

$$\frac{\delta E_d}{\delta w_{ji}} = [-(t_j - o_j) (1 - o_j^2)] x_{ji}$$

$$\Delta w_{ji} = \eta (t_j - o_j) (1 - o_j^2) x_{ji}$$

$$\text{If } \delta_j = (t_j - o_j) (1 - o_j^2)$$

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

$$\frac{\delta E_d}{\delta net_j} = -\delta_j$$

Case 2: j is a hidden layer.

$$\frac{\delta E_d}{\delta net_j} = \sum_{k \in \text{downstream}(j)} \frac{\delta E_d}{\delta net_k} \cdot \frac{\delta net_k}{\delta net_j}$$

$$= \sum -\delta_k \frac{\delta_{netk}}{\delta_{netj}}$$

$$= \sum -\delta_k \frac{\delta_{netk}}{\delta o_j} \frac{\delta o_j}{\delta_{netj}}$$

$$= \sum -\delta_k w_{kj} (1 - o_j)^2$$

$$\delta_j = (1 - o_j)^2 \sum \delta_k w_{kj}$$

$$\Delta w_{ji} = \eta \frac{\delta E_d}{\delta w_{ji}} = \eta \delta_j x_{ji}$$

Relu Activation Function

$$\text{Relu}(x) = \max(0, x)$$

$$\frac{\delta E_d}{\delta w_{ji}} = \frac{\delta E_d}{\delta_{netj}} \cdot \frac{\delta_{netj}}{\delta w_{ji}}$$

Case 1: j is an output unit.

$$\frac{\delta E_d}{\delta_{netj}} = \frac{\delta E_d}{\delta o_j} \frac{\delta o_j}{\delta_{netj}}$$

Differentiate wrt o_j ;

$$\frac{\delta E_d}{\delta o_j} = -(t_j - o_j)$$

$$\frac{\delta o_j}{\delta_{netj}} = \delta(\text{Max}(0, x))$$

$$z = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x > 0 \end{cases}$$

at $x=0$ (function is not defined).

$$\frac{\delta E_d}{\delta w_{ji}} = [-(t_j - o_j)] x_{ji} \quad \text{for } x > 0$$

$$\Delta w_{ji} = \eta (t_j - o_j) (x_{ji})$$

$$\text{Substitute } t_j - o_j = \delta_j$$

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

$$\frac{\delta E_d}{\delta \text{net}_j} = -\delta_j$$

Case 2: j is a hidden unit.

$$\frac{\delta E_d}{\delta \text{net}_j} = \sum_{k \in \text{downstream}(j)} \frac{\delta E_d}{\delta \text{net}_k} \cdot \frac{\delta \text{net}_k}{\delta \text{net}_j}$$

$$= \sum_{k \in \text{downstream}(j)} -\delta_k \frac{\delta \text{net}_k}{\delta \text{net}_j}$$

$$= \sum_k -\delta_k \cdot \frac{\delta \text{net}_k}{\delta o_j} \frac{\delta o_j}{\delta \text{net}_j}$$

$$= \sum_k -\delta_k \frac{\delta \text{net}_k}{\delta o_j} \frac{\delta o_j}{\delta \text{net}_j}$$

$$= \sum_k -\delta_k \frac{\delta \text{net}_k}{\delta o_j} \quad \text{for } x > 0$$



$$= -\sum_k \delta_k w_{kj}$$

$$\Delta w_{ji} = \eta \delta_j' x_{ji} \quad \eta > 0.$$

Summarize

$$\frac{\delta f_d}{\delta w_{ij}} = \begin{cases} - \sum_k \delta_k w_{kj} & \text{for } o_j > 0 \\ 0 & o_j < 0 \end{cases}.$$



1.2 Gradient Descent

Single unit neuron with output O , can be defined as:

$$O = w_0 + w_1(x_1 + x_2) + \dots + w_n(x_n + x_n^2)$$

where x_1, x_2, \dots, x_n are inputs,
 w_1, w_2, \dots, w_n are the corresponding weights,
 and w_0 is the bias weight.

Gradient:

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

Training Rule,

$$\Delta \vec{w} = -\eta \Delta E[\vec{w}]$$

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i}$$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \left[\frac{1}{2d} \sum (t_d - o_d)^2 \right]$$

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

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

$$= \sum_d (t_d - o_d) \cdot \frac{\partial}{\partial w_i} (t_d - (w_0 + w_1(x_1 + x_1^2) + \dots + w_n(x_n + x_n^2)))$$

$$= \sum_d (t_d - o_d) - \frac{\delta}{\delta w_i} [\vec{w}_d (x_d + x_d^2)]$$

$$\frac{\delta \mathcal{E}}{\delta w_i} = - \sum_d (t_d - o_d) \cdot (x_d + x_d^2)$$

$$\Delta w = - \eta \frac{\delta \mathcal{E}}{\delta w_i}$$

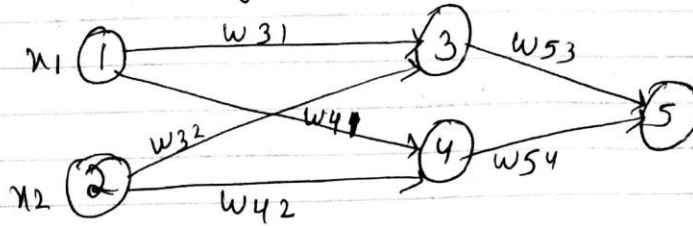
$$\Delta w = - \left[- \eta \sum_d (t_d - o_d) (x_d + x_d^2) \right]$$

$$\Delta w = \eta \sum_d (t_d - o_d) (x_d + x_d^2)$$



Ques 1.3

1.3 Comparing activation functions :-



| Node | Net | Output | Output |
|------|--------------------------------|-------------------|--------------------|
| 1 | x_1 | | |
| 2 | x_2 | | |
| 3 | $net3 = w_{31}x_1 + w_{32}x_2$ | | $x_3 = \eta(net3)$ |
| 4 | $net4 = w_{41}x_1 + w_{42}x_2$ | | $x_4 = \eta(net4)$ |
| 5 | $net5 = w_{53}x_3 + w_{54}x_4$ | | $x_5 = \eta(net5)$ |

$$x_5 = y_5 = \eta(net5)$$

$$= \eta(w_{53}x_3 + w_{54}x_4)$$

$$y_5 = h(w_{53} \eta(w_{31}x_1 + w_{32}x_2) + w_{54} \eta(w_{41}x_1 + w_{42}x_2))$$



(b) Input Layer to hidden Layer
Input 1 = $w^{(1)} x$

Output from hidden layer $\rightarrow \eta(w^{(1)} x)$

Input 2 = $w^{(2)} \cdot \eta(w^{(1)} \cdot x)$

Output $y_5 = \eta(w^{(2)} \cdot \eta(w^{(1)} x))$

$$\Rightarrow \eta \left([w_{53} \ w_{54}] \cdot \eta \left(\begin{bmatrix} w_{31} & w_{32} \\ w_{41} & w_{42} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \right) \right)$$



(c) Using Activation Function
 $hs(x) = \frac{1}{1+e^{-x}}$ (Sigmoid)

$$hs(x) = \frac{e^x}{1+e^x}$$

$$\Rightarrow hs(2x) = \frac{e^{2x}}{1+e^{2x}}$$

$$ht(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (\tanh)$$

$$= \frac{e^{2x} - 1}{e^{2x} + 1} \quad - (1)$$

$$2hs(2x) = \frac{2 \cdot e^{2x}}{e^{2x} + 1}$$

$$2hs(2x) - 1 = \frac{2e^{2x} - 1 - e^{2x}}{e^{2x} + 1}$$

$$2hs(2x) - 1 = \frac{e^{2x} - 1}{e^{2x} + 1} \quad - (2)$$

~~2hs(2x)~~ from (1) & (2)

$$ht(x) = 2hs(2x) - 1$$

The output of $\tanh(x)$ differs only by a linear transformation ($2x$) and a constant (-1).



Ques 1.4

$$\underline{1.4} \quad E(\vec{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{output}} (t_{kd} - o_{kd})^2 + \gamma \sum_{ij} w_{ji}^2$$

$$\Rightarrow \frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial n_{kj}} \cdot \frac{\partial n_{kj}}{\partial w_{ji}} + 2\gamma w_{ji}$$

$$= \frac{\partial E_d}{\partial n_{kj}} \cdot x_{ji} + 2\gamma w_{ji}$$

$$= \frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial n_{kj}} \cdot x_{ji} + 2\gamma w_{ji}$$

$$\frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial n_{kj}} = -(t_j - o_j) (1 - o_j) o_j$$

$$\frac{\partial E_d}{\partial w_{ji}} = -(t_j - o_j) (1 - o_j) o_j x_{ji} + 2\gamma w_{ji}$$

$$w_{ji}^{\text{new}} = w_{ji} + \eta (t_j - o_j) o_j (1 - o_j) x_{ji} - 2\eta \gamma w_{ji}$$

$$= (1 - 2\eta \gamma) w_{ji} + \eta (t_j - o_j) o_j (1 - o_j) x_{ji}$$

$$= (1 - 2\eta \gamma) w_{ji} + \eta s_j x_{ji} \quad (\text{for output layer})$$

for hidden layer

$$w_{ji}^{\text{new}} = w_{ji} (1 - 2\eta \gamma) + \eta s_j x_{ji}$$

here $s_j = o_j (1 - o_j)$

From the above, it can be interpreted that each rule can be implemented by multiplying $(1 - 2\eta \gamma)$ → weight update before Gradient descent update.

