

# 1. Backpropagation Algorithm,

for hidden layer,

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

for output layer,

$$\underline{f(x) = x \quad \Rightarrow \quad f'(x) = 1} \quad - (1)$$

Calculating derivatives of activation function  
for hidden & output layer,

for hidden layer,

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

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

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

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

$$= 1 - (\tanh(x))^2$$

$$\underline{\Rightarrow f'(x) = 1 - (f(x))^2} \quad - (2)$$

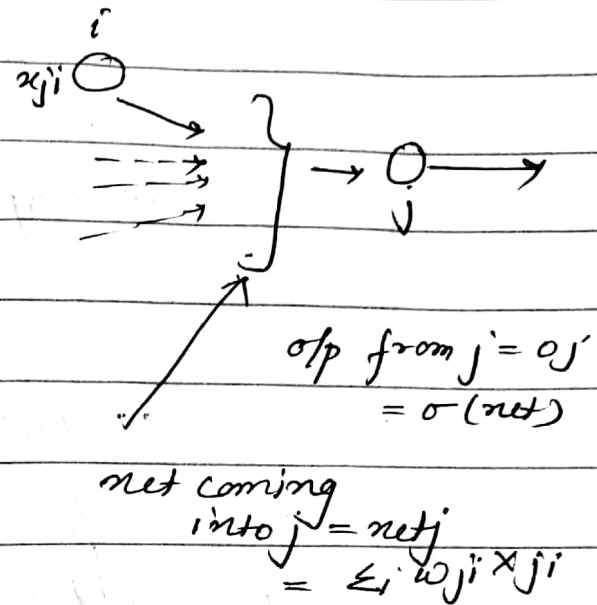
For output layer

$$f(x) = x$$

$$\Rightarrow f'(x) = 1$$

Algorithm,

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \times \frac{\partial net_j}{\partial w_{ji}}$$



$\frac{\partial E_d}{\partial net_j} \rightarrow 2 \text{ cases}$    
 ↗ hidden unit   
 ↘ output unit

Case 1, j is o/p unit

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \times \frac{\partial o_j}{\partial net_j} \rightarrow \text{Activation function}$$

$$\Rightarrow E_d = \frac{1}{2} \sum_{k \in \text{o/p}} (t_k - o_k)^2$$

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \left[ \frac{1}{2} \sum_{k \in \text{o/p}} (t_k - o_k)^2 \right] \quad (\text{from (1)})$$

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \left[ \frac{1}{2} (t_j - o_j)^2 \right] ; \quad \frac{\partial o_j}{\partial net_j} = 1$$

$$= -(t_j - o_j)$$

$$\Rightarrow -\delta_j = \frac{\partial E_d}{\partial \text{net}_j} = -(t_j - o_j) \times 1$$

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

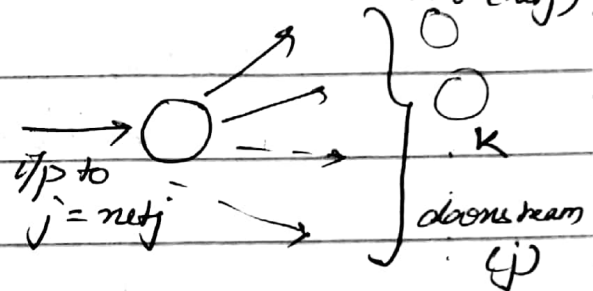
$$\Delta w_{ji} = \eta (t_j - o_j) x_{ji}$$

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

Case 2, j ∈ hidden unit

o/p from j = o\_j = σ(net\_j)

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in \text{down stream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j}$$



$$= \sum_{k \in \text{down stream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{k \in \text{down stream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial o_j} \times \frac{\partial o_j}{\partial \text{net}_j}$$

$$= \sum_{k \in \text{down stream}(j)} -\delta_k w_{kj} (1 - o_j^2)$$

Activation function

$$\delta_j = (1 - o_j^2) \sum_{k \in \text{down stream}} \delta_k w_{kj}$$

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

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

$$\begin{aligned} \text{Case 1 : } \delta_j' &= (t_j - o_j) \\ \text{Case 2 : } \delta_j &= (1 - o_j^2) \sum_{k \in \text{down stream}} \delta_k w_{kj} \end{aligned}$$

## 2. Gradient Descent - training rule

$$o = w_0 + w_1(x_1 + x_1^2) + \dots + w_n(x_n + x_n^2)^n$$

$x_1, x_2, \dots, x_n$  : Inputs

$w_1, w_2, \dots, w_n$  : weights

Activation function,  $f(x) = x$

Training rule,

$$\begin{aligned} \nabla E(\omega) &= \Delta \vec{\omega} = -\eta \nabla E(\omega) \\ \Delta \omega_i &= \eta \frac{\partial E}{\partial \omega_i} \end{aligned}$$

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

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

$$= \sum_d (t_d - o_d) \frac{\partial}{\partial \omega_i} (t_d - o_d)$$

$$= \sum_d (t_d - o_d) \frac{\partial}{\partial \omega_i} (t_d - \vec{\omega} \cdot \vec{x}_d)$$

$$\frac{\partial E}{\partial \omega_i} = \sum_d (t_d - o_d) (-(x_d + x_{id}^2))$$

$$\Delta w_i = \eta \sum_{d \in D} (t_d - o_d) (x_{id} + x_{id}^2)$$

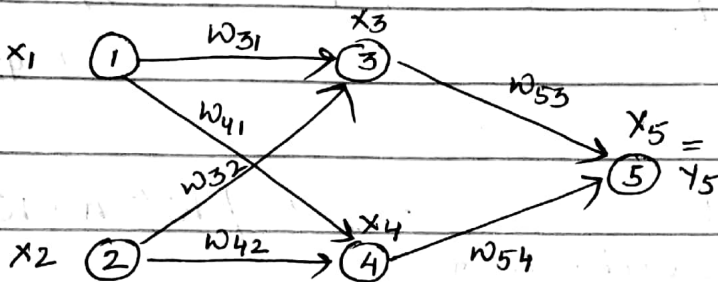
$\eta$  : learning rate,

$t_d$  : target output,

$o_d$  : Actual output,

$x_{id}$  = value of  $i$ th attribute of  $d$ th training eg.

3.



if  $\eta$

Activation function,  $f(x) = \sigma$   
hidden -  $h$

q)

Forward Pass

Node	Net	output-
1	$i_1$	$x_1 = i_1$
2	$i_2$	$x_2 = i_2$
3	$net_3 = w_{31}x_1 + w_{32}x_2$	$x_3 = \text{net}_3$ <del><math>h</math></del> $h(net_3)$
4	$net_4 = w_{41}x_1 + w_{42}x_2$	$x_4 = h(net_4)$
5	$net_5 = w_{53}x_3 + w_{54}x_4$	$x_5 = h(net_5)$

$$y_5 = h(net_5)$$

$$(x_5) = h(w_{53}x_3 + w_{54}x_4)$$

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

(b)  $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$

$$w^{(1)} = \begin{pmatrix} w_{31} & w_{32} \\ w_{41} & w_{42} \end{pmatrix}$$

$$w^{(2)} = [w_{53} \quad w_{54}]$$

from previous question,

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

output-  
in vector form  $= h[w^2 h(w^1 x)]$

---

(c)  $h_1(x) = \frac{1}{1+e^{-x}}$

$$h_2(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$h_2(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \times \frac{e^{-x}}{e^{-x}}$$

$$h_2(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$$

$$h_2(x) = h_1(x) [(1 - e^{-2x}) + 1 - 1]$$

$$h_2(x) = h_1(2x) \left[ \frac{2 - 1}{h_1(2x)} \right]$$

$$h_2(x) = 2h_1(2x) - 1$$

$$h_1(x) = \frac{1}{1 + e^{-x}} \Rightarrow h_1'(x) = h(x)(1 - h(x)) \leftarrow \text{sigmoid}$$

$$h_2(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \Rightarrow h_2'(x) = 1 - (h_2(x))^2$$

Substituting  $h_1(x)$  &  $h_2(x)$  to output functions,

for sigmoid,

$$s_j = (y_j - o_j) o_j (1 - o_j) \quad \text{--- (1)}$$

for tanhx,

$$s_j = (y_j - o_j) (1 - (o_j)^2)$$

$$s_j = (y_j - o_j) (1 + o_j) (1 - o_j) \quad \text{--- (2)}$$

Thus we observe that- the two output functions are same (1) & (2) with parameters differing only by linear transformations & constants.

4.

 $\sum \epsilon$ 

$$\epsilon(\vec{w}) = \frac{1}{2} \sum_{k \in D} \sum_{k \in \text{opp}} (t_{kd} - o_{kd})^2 + \gamma \sum_{i,j} w_{ji}^2$$

Gradient Descent update for  $\epsilon$

Update rule :

$$\Delta w_{ji} = -\eta \frac{\partial \epsilon(\vec{w})}{\partial w_{ji}}$$

$$w_{ji}^{\text{new}} \leftarrow w_{ji}^{\text{old}} + \Delta w_{ji}$$

$$\begin{aligned} \frac{\partial \epsilon(\vec{w})}{\partial w_{ji}} &= \frac{\partial}{\partial w_{ji}} \left[ \frac{1}{2} \sum_{k \in D} \sum_{k \in \text{opp}} (t_{kd} - o_{kd})^2 + \gamma \sum_{i,j} w_{ji}^2 \right] \\ &= - (t_{kd} - o_{kd}) (1 - o_{kd}) o_{kd} x_{ji} + 2 \gamma w_{ji} \end{aligned}$$

$$w_{ji} \leftarrow w_{ji} + \eta \delta_{kd} x_{ji} - 2\gamma w_{ji}$$

$\eta$  : learning rate

Since,

$$(t_{kd} - o_{kd}) (1 - o_{kd}) o_{kd} = \delta_{kd}$$

$$w_{ji} \leftarrow \eta \delta_{kd} x_{ji} + w_{ji} (1 - 2\gamma)$$

for hidden layer,

weight update,

$$w_{ji} \leftarrow \eta \delta_{kd} x_j + w_{ji} (1 - 2\gamma)$$

$$\text{Here, } \sum_j \delta_{kd} = o_{kd} (1 - o_{kd}) \sum_{k \in \text{dasm stream}} \delta_k w_{kj}$$



The above equation thus shows that the update rule can be implemented by multiplying each weight by some constant before standard Gradient-Descent-Update.