J.U.N.S. Rohith Reddy
Amandeep Singh

$$\frac{\partial t_d}{\partial \omega_{ji}} = \frac{\partial t_d}{\partial net_j} \times \frac{\partial net_j}{\partial \omega_{ji}} \qquad \text{where} \quad \frac{\partial net_j}{\partial \omega_{ji}} = x_{ji}$$

$$= \frac{\partial E_d}{\partial net_j} \times x_{ji} \quad -① \qquad \text{and} \quad net_j = \sum \omega_{ji} x_{ji}$$

$x_{ji}$ is the $i$th input to unit $j$

Two cases — ① unit $j$ is an output unit for the network

② unit $j$ is a hidden layer unit

Case 1 : Training rule for output unit weights

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \times \frac{\partial o_j}{\partial net_j} \quad -②$$

→) where $o_j$ is the output computed by unit $j$

Consider the first part of eq ②

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \, \text{outputs}} (t_k - o_k)^2$$

$$\frac{\partial}{\partial o_j} (t_k - o_k)^2 = 0 \quad \forall k \neq j$$

$$\therefore \frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} (t_j - o_j)^2 \qquad k : j$$

$$= \frac{1}{2} (2) (t_j - o_j) \frac{\partial}{\partial o_j} (t_j - o_j)$$

$$\frac{\partial E_d}{\partial o_j} = - (t_j - o_j) \quad -③$$

Now let the activation function used be tanh

$$O_j = \tanh(net_j)$$

Therefore
$$\frac{\partial O_j}{\partial net_j} = \frac{\partial}{\partial net_j}\left[\tanh(net_j)\right]$$

$$\boxed{\frac{d}{d\theta}\tanh\theta = 1 - \tanh^2\theta}$$

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

$$= 1 - O_j^2 \quad -④$$

From ②,③,④

$$\frac{\partial t_d}{\partial net_j} = -(t_j - O_j)(1 - O_j^2) \quad -⑤$$

Gradient descent rule for output units

$$\Delta w_{ji} = -\eta\,\frac{\partial t_d}{\partial O_{ji}}$$

from ① and ⑤, $\boxed{\Delta w_{ji} = \eta(t_j - O_j)(1 - O_j^2)\,x_{ji}} \quad -⑥$

Case 2: Training rule for hidden unit weights

Downstream$(j)$ refers to set of all units immediately downstream of unit $j$

$$\therefore \frac{\partial t_d}{\partial net_j} = \sum_{k \in downm(j)}^{} \frac{\partial t_d}{\partial net_k} \times \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in downt(j)} -\delta_k \cdot \frac{\partial net_k}{\partial net_j}$$

$$= \sum -\delta_k \times \frac{\partial net_k}{\partial O_j} \times \frac{\partial O_j}{\partial net_j}$$

$$= \sum_{k \in downnet(j)} -\delta_k\, w_{kj}\,\frac{\partial O_j}{\partial net_j}$$

from ④ $\dfrac{\partial O_j}{\partial net_j} = 1 - O_j^2$

$$\dfrac{\partial E_d}{\partial net_j} = \sum_{k \in down(j)} -S_k W_{kj} (1 - O_j)^2 \quad - ⑦$$

Let $S_j = -\dfrac{\partial E_d}{\partial net_j}$

then $S_j = (1 - O_j^2) \sum_{k \in down(j)} S_k W_{kj} \quad - ⑧$

$\Delta W_{ji} = \eta S_j x_{ji} \quad - ⑨$

Equations ⑧ and ⑨ ~~assist~~ shows the use of $tanh(x)$ activation function to calculate $\Delta W_{ji}$

## Using Relu as activation

Output from unit $j = O_j$

$O_j = Relu(net_j)$

$$O_j = \begin{cases} 0 & \text{for } net_j < 0 \\ net_j & \text{for } net_j \geq 0 \end{cases}$$

$$\dfrac{\partial O_j}{\partial net_j} = \begin{cases} 0 & \text{for } net_j < 0 \\ 1 & \text{for } net_j \geq 0 \end{cases}$$

Let $\dfrac{\partial O_j}{\partial net_j} = O_j'$

Two cases again

.

**Case 1 : Training rule for Output Unit Weights**

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

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

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial O_j} \times \frac{\partial O_j}{\partial net_j} \qquad \text{from ②}$$

$$\frac{\partial E_d}{\partial O_j} = -(t_j - O_j)$$

Now
$$\frac{\partial O_j}{\partial net_j} = \begin{cases} 0 & , \ net_j < 0 \\ 1 & , \ net_j \geq 0 \end{cases}$$

$$\frac{\partial O_j}{\partial net_j} = O_j'$$

$$\frac{\partial E_d}{\partial net_j} = -(t_j - O_j) O_j'$$

$$\therefore \Delta w_{ji} = \eta (t_j - O_j) O_j' x_{ji}$$

if $net_j < 0 \Rightarrow \Delta w_{ji} = 0 \quad \text{as } O_j' = 0$

if $net_j \geq 0 \Rightarrow \Delta w_{ji} = \eta (t_j - O_j) x_{ji} \quad (\text{as } O_j' = 1)$

**Case 2 : for hidden layer**

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in down(i)} -\delta_k w_{kj} \times \frac{\partial O_j}{\partial net_j}$$

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in downstream(i)} -\delta_k w_{kj} O_j'$$

$$\delta_j = -\frac{\partial E_d}{\partial net_j} = O_j' \sum_{k \in down(i)} \delta_k \cdot w_{kj}$$

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

If $net_j < 0$  (as $o_j' = 0$)

$$\Delta w_{ji} = 0$$

if $net_j \geq 0$  (as $o_j' = 1$)

$$\Delta w_{ji} = \eta \left( \sum_{k \in daonet(j)} \delta_k \cdot w_{kj} \right) x_{ji}$$

I.v.n.s. Rohith Reddy

Amandeep Singh

② $O = w_0 + w_1 (x_1 + x_1^2) + w_2 (x_2^2 + x_2) + \ldots + w_n (x_n + x_n^2)$

Error function $E(\bar{w}) = \frac{1}{2} \sum_{d \in D} (t_d - O_d)^2$

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

for term $w_0$

$$\frac{\partial E}{\partial w_0} = \frac{\partial}{\partial w_0} \frac{1}{2} \sum_{d \in D} (t_d - O_d)^2 = \frac{1}{2} \sum_{d \in D} \frac{\partial (t_d - O_d)^2}{\partial w_0}$$

$$= - \sum_{d \in D} (t_d - O_d)$$

Hence $\Delta w_0 = \eta \sum_{d \in D} (t_d - O_d)$

for $w_1, w_2 \ldots w_n$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{d \in D} (t_d - O_d)^2 = \frac{1}{2} \sum_{d \in D} \frac{\partial}{\partial w_i} (t_d - O_d)^2$$

$$= \frac{1}{2} \sum_{d \in D} 2 (t_d - O_d) \frac{\partial}{\partial w_i} (t_d - O_d) = \sum_{d \in D} (t_d - O_d)(-(x_{id} + x_{id}^2))$$

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

(B)

| Node | Net value | Output Value |
|---|---|---|
| 1 | $x_1$ | $x_1$ |
| 2 | $x_2$ | $x_2$ |
| 3 | $net_3 = \omega_{31}x_1 + \omega_{32}x_2$ | $x_3 = h(net_3)$ |
| 4 | $net_4 = \omega_{41}x_1 + \omega_{42}x_2$ | $x_4 = h(net_4)$ |
| 5 | $net_5 = \omega_{53}x_3 + \omega_{54}x_4$ | $x_5 = h(net_5)$ |

(a) $y_5 = h(net_5)$

$= h(net_5) = h(\omega_{53}x_3 + \omega_{54}x_4)$

$$= h\left(\omega_{53}\left(h(\omega_{31}x_1 + \omega_{32}x_2)\right) + \omega_{54}\left(h(\omega_{41}x_1 + \omega_{42}x_2)\right)\right)$$

(b) Output $= h\left[w^2 . h(w'x)\right]$

(c) $h_1(x) = \dfrac{1}{1+e^{-x}}$      $h_2(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

$h_1 = \dfrac{e^x}{e^x + 1}$      $h_2(x) = \dfrac{e^{2x} - 1}{e^{2x} + 1}$

$2h_1(2x) = \dfrac{2e^{2x}}{e^{2x} + 1}$     $\Rightarrow 2h_1(2x) - 1 = \dfrac{2e^{2x}}{e^{2x} + 1} - 1$

$2h_1(2x) - 1 = \dfrac{e^{2x} - 1}{e^{2x} + 1} = h_2(x)$

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

$\Rightarrow$ The output of both the functions are same, the only difference is that to get similar output as $h_2(x)$, $h_1(x)$ has to linear transform by multiplying

it with 2 and subtracting 1 from result

$\Rightarrow$ It shows that the $h_1(x)$ & $h_2(x)$ will generate same function

(4) Activation function $f(x) = \sigma(x) = \dfrac{1}{1+e^{-x}}$

$f' = \sigma(1-\sigma)$

Error function $E(\omega) = \dfrac{1}{2} \underset{d\in D}{\Sigma} \underset{k\in outputs}{\Sigma} (t_{kd} - O_{kd})^2 + \gamma \Sigma \omega_{ji}^2$

$\underbrace{\hspace{5cm}}_{Part\ 1}$  $\underbrace{\hspace{3cm}}_{part\ 2}$

We know $\Delta\omega_{ji} = -\eta \dfrac{\partial E_d}{\partial \omega_{ji}}$

$\Delta\omega_{ji} = -\eta\, \partial \dfrac{\dfrac{1}{2} \underset{d\in D}{\Sigma} \underset{k\in outputs}{\Sigma} (t_{kd}-O_{kd})^2}{\partial \omega_{ji}} - \dfrac{\eta\, \partial \gamma\, \Sigma\omega_{ji}^2}{\partial \omega_{ji}}$

Calculate part 1 $\begin{cases} \text{When } j \text{ in output layer} \dots \text{Case 1} \\ \text{When } j \text{ in hidden layer} \dots \text{Case 2} \end{cases}$

Case 1 : $\dfrac{\partial\, part\ 1}{\partial net_j} = \dfrac{\partial\, part\ 1}{\partial O_i} \times \dfrac{\partial O_i}{\partial net_j}$

$= -(t_i - O_j)\, O_j\, (1-O_j)$

$= -\delta_i \quad \text{where } \delta_i = (t_i - O_i)\, O_i\, (1-O_i)$

Case 2 : $\dfrac{\partial\, part\ 1}{\partial net_j} = \underset{k\in downstream(j)}{\Sigma} \dfrac{\partial\, part\ 1}{\partial net_k} \dfrac{\partial net_k}{\partial net_j}$

$= -\underset{k\in downstr(j)}{\Sigma} \delta_k \dfrac{\partial net_k}{\partial net_j}$

$$\frac{\partial Part\ 2}{\partial net\ i} = O_i\,(1-O_i)\sum_{k\in down(i)} \delta_k w_{ki}$$

for part 2 $\Rightarrow$ $\dfrac{\partial Part\ 2}{\partial w_{ji}} = 2\eta \sum_{i,j} w_{ji}$

$$\Delta w_{ji} = -\eta\,\delta_j x_{ji} - \eta\,2\eta\sum_{i,N} w_{ji}$$

where $j$ = Output layer then $\delta_j = (t_j - O_j)\,O_j\,(1-O_j)$

and $j$ = input layer then $\delta_j = O_j(1-O_j)\sum_{k\in down} \delta_k w_{kj}$

$$w_{ji} := w_{ji} + \Delta w_{ji}$$

$$w_{ji} = w_{ji} - \eta\,\delta_j x_{ji} - \eta\,2\eta\sum_{i,j} w_{ji}$$

$$= (1 - 2\eta\,\eta)\,w_{ji} + \eta\,\delta_j x_{ji}$$

$\Rightarrow$ $w_{ji} = \beta\,w_{ji} + \eta\,\delta_j x_{ji}$     where $\beta = (1-2\eta\eta)$

This shows we have to multiply $w_{ji}$ with constant $\beta$ before performing

Gradient descent