# CS 6375 ASSIGNMENT _____Neural Network_____

Names of students in your group:
Pallavi Pandey- PXP17009
Swapna Chintapalli- SXC180048

Number of free late days used: __0_____

Note: You are allowed a **total** of 4 free late days for the **entire semester**. You can use at most 2 for each assignment. After that, there will be a penalty of 10% for each late day.

Please list clearly all the sources/references that you have used in this assignment.

(i)

Q1. $\tanh(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}} = \dfrac{\sinh(x)}{\cosh(x)}$

$\dfrac{\partial(\tanh(x))}{\partial x} = \dfrac{\partial}{\partial x}\left(\dfrac{\sinh(x)}{\cosh(x)}\right)$

$= \dfrac{\frac{\partial}{\partial x}\sinh(x) \times \cosh(x) - \frac{\partial}{\partial x}\cosh(x) \times \sinh(x)}{\cosh^2(x)}$

$= \dfrac{\cosh^2(x) - \sin^2 h(x)}{\cosh^2(x)} = \dfrac{\cosh^2(x)}{\cos^2 h(x)} - \dfrac{\sin^2 h(x)}{\cosh^2(x)}$

$\boxed{\dfrac{\partial(\tanh(x))}{\partial x} = 1 - \tanh^2(x)}$ — ①

$E_d = \dfrac{1}{2}\sum_{k \in \text{output}}(t_k - O_k)^2$

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

$\eta$ is learning rate
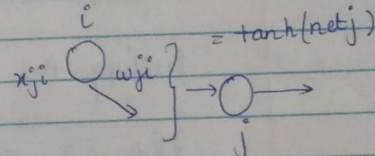
$x_{ji} \overset{i}{\bigcirc} w_{ji} \Big\} \to \bigcirc \to \qquad = \tanh(net_j)$
$\qquad\qquad\qquad\qquad j$

$\dfrac{\partial E_d}{\partial w_{ji}} = \dfrac{\partial E_d}{\partial net_j} \times \underbrace{\dfrac{\partial net_j}{\partial w_{ji}}}_{x_{ji}}$

$\boxed{\dfrac{\partial E_d}{\partial w_{ji}} = \dfrac{\partial E_d}{\partial net_j}\cdot x_{ji}}$ — ②

Case 1 :- When $j$ is an output unit

$\boxed{\dfrac{\partial E_d}{\partial net_j} = \dfrac{\partial E_d}{\partial O_j}\cdot \dfrac{\partial O_j}{\partial net_j}} \to ③$

②

$$\frac{\partial E_d}{\partial O_j} = \frac{\partial}{\partial O_j}\left(\frac{1}{2}\sum_{k \in Outputs}(t_k - O_k)^2\right) = \frac{\partial}{\partial O_j}\left[\frac{1}{2}(t_j - O_j)^2\right]$$

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

$$\boxed{\frac{\partial O_j}{\partial net_j} = 1 - O_j^2} \quad - ④$$

Putting the values of ③ and ④ in x

$$\boxed{\frac{\partial E_d}{\partial net_j} = -(t_j - O_j)(1 - O_j^2)} \rightarrow ⑤$$

Putting the value of ⑤ in ②

$$\frac{\partial E_d}{\partial w_{ji}} = -(t_j - O_j)(1 - O_j^2)x_{ji}$$

$$\Delta w_{ji} = -\eta\left(-(t_j - O_j)(1 - O_j^2)x_{ji}\right)$$
$$= \eta \underbrace{(t_j - O_j)(1 - O_j^2)}_{\delta_j} x_{ji}$$

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

Summary:
When $j$ is O/P unit $\Rightarrow \delta_j = (t_j - O_j)(1 - O_j^2)$
When $j$ is hidden unit $\rightarrow \delta_j = (1 - O_j^2)\sum \delta_n w_{kj}$
$w_{ij}$ now $= w_{ij}$ old $+ \Delta w_{ij}$
$$\boxed{\Delta w_{ij} = \eta \delta_j x_{ij}}$$

③

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

$$= \begin{cases} 1 & x > 0 \\ 0 & \text{otherwise} \end{cases}$$

$$\boxed{f_d = \frac{1}{2} \sum_{k \in \text{output}} (t_k - O_k)^2} \quad - \quad ①$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} \quad , \quad n - \text{learning rate}$$

As derived in part 1 a)

$$\boxed{\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \times X_{ji}} \quad - \quad ②$$

Case 1: when $j$ is O/P unit

$$\boxed{\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial O_j} \times \frac{\partial O_j}{\partial net_j}} \quad - \quad ③$$

$$\frac{\partial E_d}{\partial net_j} \quad \frac{\partial E_d}{\partial O_j} = \frac{\partial}{\partial O_j} \left[ \frac{1}{2} \sum_{k \in \text{output}} (t_k - O_k)^2 \right]$$

$$\boxed{\frac{\partial E_d}{\partial O_j} = - (t_j - O_j)} \quad - \quad ④$$

$$\frac{\partial O_j}{\partial net_j} = 1 \quad \text{when } net_j > 0 \atop = 0 \quad \text{otherwise} \Bigg\} \quad - \quad ⑤$$

$$\frac{\partial E_d}{\partial net_j} = - (t_j - O_j) \quad \text{when } net_j > 0$$
$$\qquad\qquad = 0 \quad \text{when } net_j \leq 0$$

④

$$\frac{\partial E_d}{\partial w_{ji}} = 0 \quad \text{when} \quad net_j \leq 0$$

$$= -(t_j - 0_j) x_{ji} \quad \text{when} \quad net_j > 0$$

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

$$\Delta w_{ji} = \underbrace{\eta (t_j - 0_j)}_{\delta_j} x_{ji} \quad \text{when} \quad net_j > 0$$

$$\Delta w_{ji} = 0 \quad \text{for} \quad net_j \leq 0$$

Case II    when $j$ is hidden unit

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in \text{downstream}} \frac{\partial E_d}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in \text{downstream}} -\delta_k \frac{\partial net_k}{\partial net_j}$$

$$= \sum_{k \in \text{downstream } j} -\delta_k \cdot \frac{\partial net_k}{\partial 0_j} \times \frac{\partial 0_j}{\partial net_j}$$

Using eq$^n$ 5

$$\frac{\partial E_d}{\partial net_j} = \sum -\delta_k w_{kj} \quad \text{for} \quad net_j > 0$$

$$= 0 \quad \text{Otherwise.}$$

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

$$= -\eta \underbrace{\left[\sum -\delta_k w_{kj}\right]}_{\delta_j} x_{ji} \quad \text{for} \quad net_j > 0$$

5

⑤

$$= 0 \qquad \text{for} \quad net_j \leq 0$$

## Summary

When $j$ is O/P unit
$$\delta_j = (t_j - 0_j), \quad net_j > 0$$
$$\delta_j = 0 \qquad , \quad net_j \leq 0$$

When $j$ is hidden layer
$$\delta_j = \sum \delta_k W_{kj}, \quad net_j > 0$$
$$\delta_j = 0, \quad net_j \leq 0$$

$$\boxed{\Delta W_{ij} = \eta \, \delta_i \, X_{ij}}$$

**Q1.2**

$$0 = w_0 + w_1(x_1 + x_1^2) + \cdots \quad w_n(x_n + x_n^2)$$
$$\underbrace{}_{bias}$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \qquad\qquad X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$E_d = \frac{1}{2} \sum_{k \in \text{output}} (t_k - 0_k)^2$$

$$\Delta w_{ji} = \eta \frac{\partial E_d}{\partial w_{ji}}, \quad \text{where} \quad \eta \text{ is the learning rate.}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \times \underbrace{\frac{\partial net_j}{\partial w_{ji}}}_{(X_j + X_j^2)}$$

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial 0_j} \times \underbrace{\frac{\partial 0_j}{\partial net_j}}_{1 \quad \text{using identity activation function}}$$

⑥

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

$$\Delta w_{ji} = -\eta \times \left[ -(t_j - o_j)(x_j + x_j^2) \right]$$

$$\boxed{\Delta w_{ji} = \eta (t_j \cdot o_j)(x_j + x_j^2)}$$

$$\Delta w_{ji} + w_j \, old = w_j \, new.$$

Q1.3   Output 1 = $x_1$
Output 2 = $x_2$

$net_3 = x_1 w_{31} + x_2 w_{32}$,   $o/p_3 = h(x_1 w_{31} + x_2 w_{32})$

$net_4 = x_1 w_{41} + x_2 w_{42}$,   $o/p_4 = h(x_1 w_{41} + x_2 w_{42})$

$net_5 = w_{53} \times h(x_1 w_{31} + x_2 w_{32}) + w_{54} \times h(x_1 w_{41} + x_2 w_{42})$

$$o/p_5 = h\left( w_{53} \cdot h(x_1 w_{31} + x_2 w_{32}) + w_{54} \cdot h(x_1 w_{41} + x_2 w_{42}) \right)$$

b)   $X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$$W(1) = \begin{bmatrix} w_{3,1} & w_{3,2} \\ w_{4,1} & w_{4,2} \end{bmatrix} \qquad W(2) = \begin{bmatrix} w_{5,3} & w_{5,4} \end{bmatrix}$$

$$W_1 \cdot x = \begin{bmatrix} w_{31} x_1 + w_{32} x_2 \\ w_{41} x_1 + w_{42} x_2 \end{bmatrix} \cdot \quad \text{o/p of hido}$$

O/P of hidden layer $X_2 = \begin{bmatrix} h(w_{31} x_1 + w_{32} x_2) \\ h(w_{41} x_1 + w_{42} x_2) \end{bmatrix}$

$$W_2 X_2 = \begin{bmatrix} w_{53} \cdot h(w_{31} x_1 + w_{32} x_2) + \\ w_{54} \cdot h(w_{41} x_1 + w_{42} x_2) \end{bmatrix}$$

⑦

O/P of $5(y_5) = \left[ h(w_{5,3} \cdot h(w_{31} X_1 + w_{32} X_2) \right.$

$\left. + w_{54} \cdot h(w_{41} X_1 + w_{42} X_2)) \right]$

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

$h_s(2x) = \dfrac{e^{2x}}{e^{2x} + 1}$

$h_t(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}} = \dfrac{e^{2x} - 1}{e^{2x} + 1}$ ——①

$2(h_s(2x)) - 1 = \dfrac{2e^{2x}}{e^{2x} + 1} - 1 = \dfrac{2e^{2x} - e^{2x} - 1}{e^{2x} + 1} = \dfrac{e^{2x} - 1}{e^{2x} + 1}$ ——②

So, $h_t(x) = 2(h_s(2x)) - 1$

Hence, the neural N/w created using the above two
activation function can generate the same function
with some diff in constants.

Q1.4

$E(\vec{w}) = \dfrac{1}{2} \sum\limits_{d \in D} \sum\limits_{k \in outputs} (t_{kd} - O_{kd})^2 + \gamma \sum\limits_{i,j} w_{ji}^2$

Assuming sigmoid activation

$\Delta w_{ji} = -\eta \dfrac{\partial E(\vec{w})}{\partial w_{ji}}$

$\dfrac{\partial E(\vec{w})}{\partial w_{ji}} = \dfrac{\partial}{\partial w_{ji}} \left[ \dfrac{1}{2} \sum\limits_{d \in D} \sum\limits_{k \in outputs} (t_{kd} - O_{kd})^2 + \gamma \sum\limits_{i,j} w_{ji}^2 \right]$

⑧

$$= \frac{1}{2} \times 2 \times -(t_j - 0_j)(0_j)(1-0_j) X_{ji} - 2\gamma w_{ji}$$

$$\Delta w_{ji} = \eta \underbrace{(t_j - 0_j)(0_j)(1-0_j) X_{ji}}_{\delta_j} - 2\gamma w_{ji}$$

$$\Delta w_{ji} = \eta \, \delta_j \, X_{ji} - 2\gamma w_{ji}$$

~~update~~ update rule

$$w_{ji} = \Delta w_{ji} + w_{ji} = \eta \delta_j X_{ji} - 2\gamma w_{ji} + w_{ji}$$

$$\boxed{w_{ji} = \eta \, \delta_j X_{ji} - (2\gamma - 1) w_{ji}}$$

$$\delta_j = (t_j - 0_j)(0_j)(1-0_j) \qquad \rightarrow \text{output layer}$$

$$\delta_j = 0_j(1-0_j) \sum_{k \in \text{downstream}} \delta_k \, w_{kj} \qquad \rightarrow \text{hidden layer}$$

Hence, it proves that update can be implemented by multiplying each weight by some constant before performing standard gradient descent.