

ML Assignment 3

Ans 1.1) a)

Weight for the training examples updates by this rule: $w_i = w_i + \Delta w$

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

wt. from i to j [input from node i to unit j]

E_d : error for each training example d

$$E_d(\vec{w}) = \frac{1}{2} \sum_{k \in O} (t_k - o_k)^2 \quad [k \in \text{output} \equiv k \in O]$$

\downarrow target output \downarrow output computed by unit 'k'

To implement stochastic gradient descent, we compute

$$\frac{\partial E_d}{\partial w_{ji}}$$

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

$$= \underbrace{\frac{\partial E_d}{\partial \text{net } j}}_{\text{computing this term firstly}} \cdot \left[\frac{\partial \text{net } j}{\partial w_{ji}} = \frac{\partial \sum w_{ji} x_{ji}}{\partial w_{ji}} \right] = x_{ji}$$

computing this term firstly

$$\frac{\partial E_d}{\partial \text{net } j} = \frac{\partial E_d}{\partial o_j} \frac{\partial o_j}{\partial \text{net } j}$$

Case 1) training rule for output unit weights

$$\frac{\partial E_d}{\partial \text{net}_j} = \frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j} \quad - (1)$$

solving this term first

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k=0} (t_k - o_k)^2$$

The derivatives $\frac{\partial}{\partial o_j} (t_k - o_k)^2$ will be zero for all outputs units k except when $k = j$

We drop the summation & set $k = j$

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

Now lets consider the next term

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial (\tanh(\text{net}_j))}{\partial \text{net}_j}$$

This is simply the derivative of tanh function

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

$$\frac{d}{dn} \tanh(n) = \frac{(e^n + e^{-n})(e^n + e^{-n}) - (e^n - e^{-n})(e^n - e^{-n})}{(e^n + e^{-n})^2}$$

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

$$\therefore \frac{\partial o_j}{\partial \text{net}_j} = 1 - (o_j)^2 \quad \text{--- (3)}$$

Putting (2), (3) in (1), we get

$$\frac{\partial E_d}{\partial \text{net}_j} = -(t_j - o_j)(1 - (o_j)^2) n_{ji}$$

$$\begin{aligned} \Delta w_{ji} &= -\eta \cdot [-(t_j - o_j)(1 - (o_j)^2) n_{ji}] \\ &= \eta \underbrace{(t_j - o_j)(1 - (o_j)^2)}_{f_j} n_{ji} \end{aligned}$$

$$\therefore \Delta w_{ji} = \eta f_j n_{ji}$$

$$\text{where } f_j = (t_j - o_j)(1 - (o_j)^2)$$

Case 2) training rule for hidden unit weights

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in \text{DS}(j)} \frac{\partial E_d}{\partial \text{net}_k} \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j} \quad [k \in \text{DS}(j) \equiv k \in \text{downstream}(j)]$$

$$= \sum_{k \in \text{DS}(j)} -f_k \cdot \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{k \in \text{DS}(j)} -f_k \cdot \underbrace{\frac{\partial \text{net}_k}{\partial o_j}}_{\downarrow} \cdot \underbrace{\frac{\partial o_j}{\partial \text{net}_j}}_{\downarrow}$$

$$= \sum_{k \in \text{DS}(j)} -f_k w_{kj} (1 - (o_j)^2)$$

$$\therefore \frac{\partial E_d}{\partial \text{net}_j} = -(1 - (o_j)^2) \sum_{k \in \text{DS}(j)} f_k w_{kj}$$

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}} = -\eta \frac{\partial E_d}{\partial \text{net}_j} \cdot n_{ji}$$

$$= -\eta \left[-(1 - (o_j)^2) \sum_{k \in \text{DS}(j)} f_k w_{kj} \right] n_{ji}$$

$$= \eta \underbrace{(1 - (o_j)^2) \sum_{k \in \text{DS}(j)} f_k w_{kj}}_{\downarrow f_j} (n_{ji})$$

$$\therefore \boxed{\Delta w_{ji} = \eta f_j n_{ji}} \quad \text{where } f_j = (1 - (o_j)^2) \sum_{k \in \text{DS}(j)} f_k w_{kj}$$

Ans 1) b) ReLu

$$E_d(\bar{w}) = \frac{1}{2} \sum_{k \in O} (t_k - o_k)^2$$

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

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial net_j} \eta_{ji} \quad (\text{from before})$$

Case 1) training rule for output unit weights

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial net} \cdot \frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

$$\frac{\partial E_d}{\partial o_j} = \frac{\partial}{\partial o_j} \frac{1}{2} \sum_{k \in O} (t_k - o_k)^2$$

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

$\frac{\partial o_j}{\partial net_j}$ is derivate of ReLu function

$$\text{Relu}(n) = \begin{cases} n & n > 0 \\ 0 & n \leq 0 \end{cases}$$

$$\frac{\partial (\text{Relu}(n))}{\partial n} = \begin{cases} 1 & n > 0 \\ 0 & n \leq 0 \end{cases}$$

$$\frac{\partial o_j}{\partial \text{net}_j} = \begin{cases} 1 & \text{net}_j > 0 \\ 0 & \text{net}_j \leq 0 \end{cases} \quad - \textcircled{4}$$

$$\therefore \frac{\partial E_d}{\partial \text{net}_j} = \begin{cases} -(t_j - o_j) & \text{net}_j > 0 \\ 0 & \text{net}_j \leq 0 \end{cases}$$

$$\frac{\partial E_d}{\partial w_{ji}} = \begin{cases} -(t_j - o_j) n_{ji} & \text{net}_j > 0 \\ 0 & \text{net}_j \leq 0 \end{cases}$$

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

$$\Delta w_{ji} = \begin{cases} \eta (t_j - o_j) n_{ji} & \text{net}_j > 0 \\ 0 & \text{net}_j \leq 0 \end{cases} = \eta f_k n_{ji}$$

$$0 \text{ where } f_j = (t_j - o_j)$$

Case 2) training rule for hidden unit weights

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in D S(j)} \frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{k \in D S(j)} -f_k \frac{\partial \text{net}_k}{\partial \text{net}_j}$$

$$= \sum_{k \in D S(j)} -f_k \frac{\partial \text{net}_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial \text{net}_j}$$

$$\frac{\partial E_d}{\partial \text{net}_j} = \begin{cases} \sum_{k \in D S(j)} -f_k w_{kj} & \text{net}_j > 0 \\ 0 & \text{net}_j \leq 0 \end{cases} \quad [\text{from } \textcircled{4}]$$

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

$$= -\eta \frac{\partial E_d}{\partial \text{net } j} \cdot n_{ji}$$

$$= \begin{cases} +\eta \sum_{k \in DS(j)} f_k w_{kj} & \text{net } j > 0 \\ 0 & \text{net } j \leq 0 \end{cases} \quad \text{where } f_j = \sum_{k \in DS(j)} f_k w_{kj}$$

$$\therefore \Delta W_{ji} = \begin{cases} \eta f_j n_{ji} & \text{net } j > 0 \\ 0 & \text{net } j \leq 0 \end{cases}$$

where $f_j = \sum_{k \in DS(j)} f_k w_{kj}$

Ans 1.2) $o = w_0 + w_1(x_1 + x_1^2) + \dots + w_n(x_n + x_n^2)$

x_i : inputs

w_i : weights

o : output

$$E_d(\bar{w}) = \frac{1}{2} \sum_{k \in O} (t_k - o_k)^2$$

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

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

← $(x_j + x_j^2)$

$$\frac{\partial E_d}{\partial net_j} = \frac{\partial E_d}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

→ $= f'(n) = 1$ [activation function is $f(n) = n$]

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

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

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

Ans 1.3) a) i: input o: output

$$i_3 = n_1 w_{31} + n_2 w_{32}$$

$$i_4 = n_1 w_{41} + n_2 w_{42}$$

$$o_3 = h(n_1 w_{31} + n_2 w_{32})$$

$$o_4 = h(n_1 w_{41} + n_2 w_{42})$$

$$i_5 = w_{53} o_3 + w_{54} o_4$$

$$i_5 = w_{53} [h(n_1 w_{31} + n_2 w_{32})] + w_{54} [h(n_1 w_{41} + n_2 w_{42})]$$

$$o_5 = h(i_5)$$

$$o_5 = y_5 = h(w_{53} \cdot h(n_1 w_{31} + n_2 w_{32}) + w_{54} \cdot h(n_1 w_{41} + n_2 w_{42}))$$

b)

$$X = \begin{bmatrix} n_1 \\ n_2 \end{bmatrix}$$

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

$$W^{(2)} = \begin{bmatrix} w_{53} & w_{54} \end{bmatrix}$$

$$W^{(1)} \cdot X = \begin{bmatrix} w_{31} n_1 + w_{32} n_2 \\ w_{41} n_1 + w_{42} n_2 \end{bmatrix}$$

$$\text{output for hidden layer} = \begin{bmatrix} h(w_{31} n_1 + w_{32} n_2) \\ h(w_{41} n_1 + w_{42} n_2) \end{bmatrix} = x_1$$

$$W^{(2)} x_1 = [w_{53} \cdot h(w_{31} n_1 + w_{32} n_2) + w_{54} \cdot h(w_{41} n_1 + w_{42} n_2)]$$

$$\text{output for node 5} = y_5 = h(W^{(2)} x_1)$$

$$= [h(w_{53} \cdot h(w_{31} n_1 + w_{32} n_2) + w_{54} \cdot h(w_{41} n_1 + w_{42} n_2))]$$

$$\text{Ans) 1.3). c) } h_s(n) = \frac{1}{1+e^{-n}} = s \text{ (let)} \quad - (1)$$

$$h_t(n) = \frac{e^n - e^{-n}}{e^n + e^{-n}} = t \text{ (let)} \quad - (2)$$

$$s = \frac{1}{1 + \frac{1}{e^n}} = \frac{e^n}{e^n + 1} \Rightarrow e^n = e^n \cdot s + s$$

$$e^n (1-s) = s$$

$$\therefore e^n = \frac{s}{1-s} \quad - (3)$$

Putting (3) in (2), we get

$$t = \frac{\frac{s}{1-s} - \left(\frac{1-s}{s}\right)}{\frac{s}{1-s} + \left(\frac{1-s}{s}\right)} = \frac{\frac{s^2 - (1-s)^2}{(1-s)s}}{\frac{s^2 + (1-s)^2}{(1-s)s}}$$

$$= \frac{s^2 - [1 + s^2 - 2s]}{s^2 + 1 + s^2 - 2s} = \frac{2s - 1}{2s^2 - 2s + 1}$$

$$\therefore h_t(n) = \frac{2h_s(n) - 1}{2(h_s(n))^2 - 2h_s(n) + 1}$$

Cont...

PTO

$$h_s(n) = \frac{l^n}{l^n + 1}$$

$$h_s(2n) = \frac{l^{2n}}{l^{2n} + 1} \quad \left[\text{by noticing power of 2 in } h_s(n) \right]$$

$$h_t(n) = \frac{l^{2n} - 1}{l^{2n} + 1}$$

By hit & trial we note that ' $2(h_s(2n)) - 1$ ' equals

$$\frac{2 \cdot \frac{l^{2n}}{l^{2n} + 1} - 1}{1} = \frac{2 \cdot \frac{l^{2n}}{l^{2n} + 1} - \frac{l^{2n} + 1}{l^{2n} + 1}}{1} = \frac{l^{2n} - 1}{l^{2n} + 1}$$

which equals $h_t(n)$

$$\therefore h_t(n) = 2 \cdot h_s(2n) - 1$$

Hence output generated is same, with the parameters differing only by linear transformations & constants.

$$\text{Ans 14)} E(\vec{w}) = \frac{1}{2} \sum_{k \in D} \sum_{k \in O} (t_{kd} - o_{kd})^2 + \gamma \sum_{i,j} w_{ji}^2$$

let's assume the activation function to be sigmoid.

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

Δw_{ji} : wt. from i to j
 η : learning rate
 E_d : error for each training example d

Case 1: for output unit weights

we already derived Δw_{ji} for $E(\vec{w})$ without the term $\gamma \sum_{i,j} w_{ji}^2$ in Q1.1) \therefore We would only derive for this term now & use the previous result.

$$= \frac{\partial}{\partial w_{ji}} \gamma \sum_{i,j} w_{ji}^2$$

$$= 2\gamma w_{ji}$$

$$\therefore \Delta w_{ji} = -\eta \left[-(t_j - o_j) o_j (1 - o_j) \eta_{ji} + 2\gamma w_{ji} \right]$$

$$= \eta \underbrace{(t_j - o_j) o_j (1 - o_j)}_{f_j} \eta_{ji} - 2\gamma w_{ji}$$

$$\boxed{\Delta w_{ji} = \eta f_j \eta_{ji} - 2\gamma w_{ji}}$$

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

$$= w_{ji} + \eta f_j \eta_{ji} - 2\gamma w_{ji}$$

$$\therefore \boxed{w_{ji} = \eta f_j \eta_{ji} - (2\gamma - 1) w_{ji}}$$

$$\text{where } f_j = (t_j - o_j) o_j (1 - o_j)$$

Case 2: for hidden unit weights

$$\Delta W_{ji} = -\eta \left[-o_j (1 - o_j) \sum_{k \in DS} f_k w_{kj} n_{ji} + 2\gamma w_{ji} \right]$$

$$= \eta \underbrace{o_j (1 - o_j) \sum_{k \in DS} f_k w_{kj}}_{f_j} n_{ji} - 2\gamma n_{ji}$$

$$\boxed{\Delta W_{ji} = \eta f_j n_{ji} - 2\gamma n_{ji}}$$

$$W_{ji}^{new} = W_{ji}^{old} + \Delta W_{ji}$$

$$= W_{ji} + \eta f_j n_{ji} - 2\gamma n_{ji}$$

$$\boxed{W_{ji} = \eta f_j n_{ji} - (2\gamma - 1) W_{ji}}$$

$$\text{where } f_j = o_j (1 - o_j) \sum_{k \in DS} f_k w_{kj}$$

In both the cases we get the same equation differing only by the value of the ' f_j ' term.

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