net id: rXS179030

Q.1.1 Consider the error on the outputs of training example $d$

$$E_d(w) = \frac{1}{2} \sum_{k \in d} (t_k - o_k)^2 \quad \text{where } k \in \text{output layer}$$

Let $w_{ji}$ be the weight between nodes $i$ and $j$

So     $\underset{i}{o} \xrightarrow{\ w_{ji}\ } \underset{j}{o}$

and let $x_{ji}$ be the input coming from $i$ to $j$.

Let $f_a$ be an activation function such as:

$$o_j = f_a(net_j) \quad \text{where} \quad net_j = \sum_i w_{ji} x_{ji}$$

Consider the update rule $w_{ji}^{new} = w_{ji}^{old} + \Delta w_{ji}$

where $\Delta w_{ji} = -\eta \dfrac{\partial E_d}{\partial w_{ji}}$. Let us compute $\dfrac{\partial E_d}{\partial w_{ji}}$.

By the chain rule of derivation we have:

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

We have that $\dfrac{\partial net_j}{\partial w_{ji}} = x_{ji}$

we have two cases for $\dfrac{\partial E_d}{\partial net_j}$

Case 1    $j$ is an output unit.

By the chain rule we have that

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

for $\dfrac{\partial E_d}{\partial o_j}$ we have $\dfrac{\partial E_d}{\partial o_j} = \dfrac{\partial}{\partial o_j}\left(\frac{1}{2}\sum_k (t_x - o_k)^2\right) = -(t_j - o_j)$

now $\dfrac{\partial o_j}{\partial net_j} = \dfrac{\partial f_a(net_j)}{\partial net_j} = f'_a(net_j)$

let    $f_a(x) = \tanh(x) = \dfrac{e^x - e^{-x}}{e^x + e^{-x}}$

Now by the division rule of derivation $\left(\frac{p}{q}\right)' = \frac{p'q - q'p}{q^2}$ ③

We have:

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

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

So $\dfrac{\partial o_j}{\partial net_j} = \dfrac{\partial f_a(net_j)}{\partial net_j} = f'_a(net_j) = 1 - \tanh^2(net_j)$

$$= 1 - o_j^2$$

Therefore if $f_a = \tanh$ and $j$ is an output layer.

$$\frac{\partial E_d}{\partial net_j} = -(t_j - o_j) \cdot (1 - o_j^2)$$

and $\dfrac{\partial E_d}{\partial w_{ji}} = -(t_j - o_j)(1 - o_j^2) X_{ji}$

So $\Delta w_{ji} = \eta (t_j - o_j)(1 - o_j^2) X_{ji}$ for $j$ an output layer
and $f_a = \tanh$.

now if $f_a(x) = Relu(x) = max(0,x)$

we can write $Relu(x) = \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases}$

then $Relu'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases} = \mathbb{1}_{(x>0)}(x)$ is the 0-1 indicator function over the set $x > 0$

where $\mathbb{1}_{x>0}(x)$

So if $f_a = Relu$ and $j$ is an output layer.

then $\dfrac{\partial o_j}{\partial net_j} = \dfrac{\partial f(net_j)}{\partial net_j} = f_a'(net_j) = Relu'(net_j)$

$= \mathbb{1}_{(x>0)}(net_j)$

So if $f_a = Relu$ and $j$ is an output layer

$\dfrac{\partial E_d}{\partial net_j} = -(t_j - o_j) \cdot \mathbb{1}_{(x>0)}(net_j)$

and $\dfrac{\partial E_d}{\partial w_{ji}} = -(t_j - o_j) \cdot \mathbb{1}_{(x>0)}(net_j) \cdot x_{ji}$

So $\Delta w_{ji} = \eta (t_j - o_j) \cdot \mathbb{1}_{x>0}(net_j) \, x_{ji}$ for $j$ an output layer and $f_a = Relu$.

Case 2: $j$ is a hidden unit

Let's denote $-\delta_j = -(t_j - o_j) f_a'(net_j) = \dfrac{\partial E_d}{\partial net_j}$

if $j$ is a hidden unit then

$$\frac{\partial E_d}{\partial net_j} = \sum_{k \in Downstream(j)} \frac{\partial E_d}{\partial net_k} \frac{\partial net_k}{\partial net_j}$$

$$= \sum -\delta_k \frac{\partial net_k}{\partial o_j} \cdot \frac{\partial o_j}{\partial net_j}$$

Now $\dfrac{\partial o_j}{\partial net_j} = f_a'(net_j)$

and $\dfrac{\partial net_k}{\partial o_j} = \dfrac{\partial (\sum w_{kj} o_{kj})}{\partial o_j} = w_{kj}$

So $\dfrac{\partial E_d}{\partial net_j} = f_a'(net_j) \sum -\delta_k w_{kj}$

So in the case $f_a = \tanh$, $f_a'(net_j) = 1 - o_j^2$

we have $\dfrac{\partial E_\ell}{\partial net_j} = (1 - o_j^2) \sum\limits_{k \in Downstream(j)} -\delta_k \cdot w_{kj}$

So if $j$ is a hidden layer and $f_a = \tanh$

$$\Delta w_{ji} = \eta \, (1 - o_j^2) \left( \sum\limits_{k \in Downstream(j)} \delta_k \, w_{kj} \right) x_{ji}$$

In the case $f_a = Relu$, $f_a'(net_j) = 1 \; (net_j)$
$(x > 0)$

we have $\dfrac{\partial E_\ell}{\partial net_j} = \underset{x > 0}{1 \; (net_j)} \cdot \sum\limits_{k \in Downstream(j)} -\delta_k \, w_{kj}$

So if $j$ is hidden layer and $f_a = Relu$

$$\Delta w_{ji} = \eta \cdot \underset{x > 0}{1 \; (net_j)} \left( \sum\limits_{k \in Downstream(j)} \delta_k \, w_{kj} \right) x_{ji}$$

Note: $\delta_k = (t_k - o_k)(1 - o_k^2)$ if $f_a = \tanh$

and $\delta_k = (t_k - o_k) \underset{x > 0}{1 \; (net_k)}$ if $f_a = Relu$

Q 1.2

Let $o = w_0 + w_1(x_1 + x_1^2) + \cdots + w_n(x_n + x_n^2)$

where $w_0$ is the bias weight and $\{x_i\}_{i=1,\ldots,n}$,

$\{w_i\}_{i=1}^{n}$ are the inputs and weights respectively.

Consider the identity activation function.

Let the Error be $E = \frac{1}{2} \sum_{k \in D} (t_k - o_k)^2$ where $D$ is the set of training examples.

and the update rule $w_i^{new} = w_i^{old} + \Delta w_i$

where $\Delta w_i = -\eta \dfrac{\partial E}{\partial w_i}$

Let us find $\dfrac{\partial E}{\partial w_i}$.

applying derivative of a finite sum

$$\frac{\partial E}{\partial w_i} = \frac{\partial \left( \frac{1}{2} \sum (t_k - o_k)^2 \right)}{\partial w_i} \overset{?}{=} \frac{1}{2} \sum_k \frac{\partial (t_k - o_k)^2}{\partial w_i}$$

now $\dfrac{\partial (t_k - o_k)^2}{\partial w_i} = \dfrac{\partial (t_k - o_k)^2}{\partial (t_k - o_k)} \cdot \dfrac{\partial (t_k - o_k)}{\partial w_i}$

by the chain rule

$$\frac{\partial (t_k - o_k)^2}{\partial (t_k - o_k)} = t_k - o_k$$

and $\dfrac{\partial (t_k - o_k)}{\partial w_i} = \dfrac{\partial \left[ t_k - \sum\limits_{j} w_j (x_{jk} + x_{jk}^2) \right]}{\partial w_i}$

$$= -\left( x_{ik} + x_{ik}^2 \right)$$

Here $x_{jk}$ denotes the $i^{th}$ component of the $k^{th}$ training example.

So $\dfrac{\partial E}{\partial w_i} = -\dfrac{1}{2} \cdot 2 \sum\limits_{k \in D} (t_k - o_k) \cdot (x_{ik} + x_{ik}^2)$

$$= -\sum\limits_{k \in D} (t_k - o_k)(x_{ik} + x_{ik}^2)$$

Therefore the gradient descent update rule would be

$$w_i^{new} = w_i^{old} + \Delta w_i$$

$$= w_i^{old} + \eta \sum\limits_{k \in D} (t_k - o_k)(x_{ik} + x_{ik}^2)$$
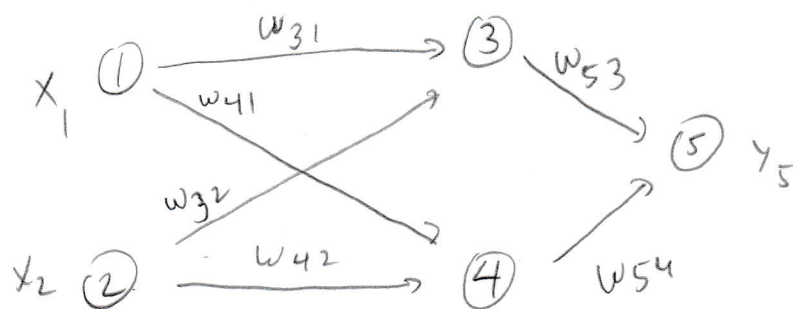
Q1.3

Consider

Input layer          Hidden layer



a)

Let $z_3$ and $z_4$ be:

$$z_3 = w_{31} x_1 + w_{32} x_2$$

$$z_4 = w_{41} x_1 + w_{42} x_2$$

then $x_3 = h(z_3)$

$x_4 = h(z_4)$

So   $z_5 = w_{53} x_3 + w_{54} x_4$

and  $y_5 = h(t_5)$

Putting all together we obtain

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

b) Let $X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$ $\quad W^1 = \begin{pmatrix} w_{31} & w_{32} \\ w_{41} & w_{42} \end{pmatrix}$ (10)

$$W^2 = (w_{53} \quad w_{54})$$

Following the notation from a)

we have:

$$\begin{pmatrix} z_3 \\ z_4 \end{pmatrix} = W^1 X$$

$$\begin{pmatrix} x_3 \\ x_4 \end{pmatrix} = h\begin{pmatrix} z_3 \\ z_4 \end{pmatrix}$$

then $z_5 = W^2 \begin{pmatrix} x_3 \\ x_4 \end{pmatrix}$

$$y_5 = h(z_5)$$

Thus: $y_5 = h\left[ W^2 h(W^1 X) \right]$

Note: We used the simplified notation $h(v)$ to represent the vector $\begin{pmatrix} h(v_1) \\ h(v_2) \\ \vdots \\ h(v_n) \end{pmatrix}$ where $v = \begin{pmatrix} v_1 \\ v_2 \\ \vdots \\ v_n \end{pmatrix}$

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

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

Let us try to find a relationship between $h_s$ and $h_t$

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

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

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

we can notice that $h_s(2x) = \frac{1}{1 + e^{-2x}}$

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

Therefore $h_s$ and $h_t$ generate functions with parameters just differing by linear transformations and constants.

Q1.4    Consider the error of the network    or

$$E(w) = \frac{1}{2} \sum_{d \in D} \sum_{k \in outputs} (t_{kd} - o_{kd})^2 + r \sum_{ij} w_{ji}^2$$

Consider the error or the outputs of the training example d.

$$E_d(w) = \frac{1}{2} \sum_{k \in d} (t_k - o_k)^2 + r \sum_{ij} w_{ji}^2$$

Let us derived the update rule:

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

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

Let us find $\dfrac{\partial E_d}{\partial w_{ji}}$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial \frac{1}{2} \sum_{k \in d} (t_k - o_k)^2}{\partial w_{ji}} + \frac{\partial r \sum w_{ji}^2}{\partial w_{ji}}$$

$$\frac{d r \sum_{ij} w_{ji}^2}{\partial w_{ji}} = 2 r \cdot w_{ji}$$

Now $\dfrac{\partial \frac{1}{2} \sum_{k \in d} (t_k - o_k)^2}{\partial w_{ji}}$ was obtained on Q.1.1

for the two cases ; an output layer and $j$ a hidden layer.

So if $j$ is an output layer.

$$\frac{\partial \frac{1}{2} \sum (t_k - o_k)^2}{\partial w_{ji}} = - (t_j - o_j) f'_a (net_j) x_{ij}$$

So $w_{ji}^{new} = w_{ji} - \eta \dfrac{\partial E_d}{\partial w_{ji}} = $

$= w_{ji} - \eta \left( - (t_j - o_j) f'_a (net_j) + 2 r \cdot w_{ji} \right)$

$= w_{ji} - 2 r \eta \, w_{ji} + \eta (t_j - o_j) f'(net_j)$

$= \left( 1 - 2 r \eta \right) w_{ji} + \eta (t_j - o_j) f'(net_j)$

Now if $j$ is a hidden layer we have that (14)

$$\frac{\partial \frac{1}{2} \sum_k (t_k - o_k)^2}{\partial w_{ji}} = f'_a(net_j)\left(\sum_{k \in Downstream(j)} -(t_j - o_j) f'_a(net_k) \cdot w_{kj}\right) X_{ji}$$

So

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

$$= w_{ji} - \eta \left[ -f'_a(net_j)\left(\sum_k (t_j - o_j) \cdot f'_a(net_k) w_{kj}\right) X_{ji} + 2r w_{ji}\right]$$

$$= w_{ji} - 2r\eta + \eta f'_a(net_j) \cdot \left(\sum_k (t_j - o_j) f'_a(net_k) w_{kj}\right) X_{ji}$$

$$= (1 - 2r\eta) w_{ji} + \eta f'_a(net_j)\left(\sum_k (t_j - o_j) f'_a(net_k) w_{kj}\right) X_{ji}$$

We can notice in either case $j$ an output or hidden layer we have obtained a weight update rule identical to the original backpropagation rule except that each weight is multiplied by the constant $1 - 2r\eta$