Part 1

1   Case 1 : when layer = output layer.

$$O = f(x) = x.$$

$$E(w) = \frac{1}{2} \sum_{k \in \text{outputs}} (t-o)^2$$

$$\nabla E(w) = \frac{\partial E(w)}{\partial w} = \frac{1}{2} \sum_{k \in \text{outputs}} \frac{\partial (t-o)^2}{\partial w}$$

$$= \frac{1}{2} \sum_{k \in \text{outputs}} 2(t-o)(-x_k) = \sum_{k} (t-o)(-x_k)$$

$$\Rightarrow \Delta w_{ji} = \eta \nabla E(w)$$

$$= \eta (t_i - o_i)(-x_{ji})$$

$$= -\eta (t-o) x_{ji}$$

$$\Rightarrow \delta_j = (t_j - o_j)$$

Case 2 : when layer = hidden layer.

$$O = \tanh(x)$$

$$net = x$$

$$E(w) = \frac{}{2}$$

For a hidden unit $h$

$$\delta_h = \frac{\partial E(w)}{\partial net_h} = \sum_{k \in \text{Downstream}(h)} \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial net_h}$$

$$= \sum_{k \in \text{Downstream}(h)} -\delta_k \cdot \frac{\partial net_k}{\partial net_h}$$

$$= \sum_{k \in \text{Downstream}(h)} -\delta_k \cdot \frac{\partial net_k}{\partial O_h} \cdot \frac{\partial O_h}{\partial net_h}$$

$$= \sum_{k \in \text{Downstream}(h)} -\delta_k \cdot \frac{\partial x_{kh}}{\partial O_h} \cdot (1 - O_h^2)$$

$$\equiv \sum_{k \in \text{Downstream}(h)} -\delta_k \cdot \partial O_{hh} \, w_{kh} \, (1 - O_h^2)$$

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

$$\Rightarrow \quad \delta_h = (1 - O_h^2) \sum_{k \in \text{Downstream}(h)} \delta_k \, w_{kh}$$

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

$$= -\eta \, \delta_h \cdot x_{ji}$$

Finally :

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

For output layer :

$$\delta_j = (t_j - O_j)$$

For hidden layer :

$$\delta_j = (1 - O_j^2) \sum_{k \in \text{downstream}(j)} \delta_k \, w_{kj}$$

2

linear activation function.

$$E_d(w) = \frac{1}{2} \sum_{d \in D} (t_d - O_d)^2 \quad d \in D$$

D = training data

d = training instance

$$\Delta w_j = -\eta \; \nabla E_d(w)$$

$$= -\eta \; \frac{\partial E_d(w)}{\partial w}$$

$$= -\eta \cdot \frac{1}{2} \sum_{d \in D} 2(t_d - O_d) \cdot \frac{\partial(t_d - O_d)}{\partial w_j}$$

$$= -\eta \cdot \frac{1}{2} \sum_{d \in D} (t_d - O_d) \cdot \left( -(x_j + x_j^2) \right)$$

$$\Rightarrow \quad \Delta w_j = \sum_{d \in D} \eta (t_d - O_d) (x_j + x_j^2)$$

$\eta$ = learning rate.

D = all training data

d = single training instance

$t_d$ = target output for dth instance

$O_d$ = observed output for dth instance

$x_{id}$ = value of ith attribute for dth training example.

final weight update rule :

$$w_j = w_j + \Delta w_j$$

## 3

### (a).

$$x_3 = w_{31} x_1 + w_{32} x_2 \quad, \quad O_3 = x_3$$

$$x_4 = w_{41} x_1 + w_{42} x_2 \quad, \quad O_4 = x_4$$

$$x_5 = w_{53} x_3 + w_{54} x_4 \quad, \quad O_5 = y_5 = x_5.$$

### (b)

$$X_H = W^{(1)} \times X = \begin{pmatrix} x_3 \\ x_4 \end{pmatrix}$$

$$X_0 = W^{(2)} \times X_H = \begin{pmatrix} x_5 \end{pmatrix}$$

$$y_5 = x_5.$$

### (c)

$$h_1(x) = \frac{1}{1+e^{-x}} \quad, \quad h_2(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$h_1 + h_1 e^{-x} = 1$$

$$e^{-x} = \frac{1-h_1}{h_1}$$

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

$$= \frac{h_1^2 - \left(1 + h_1^2 - 2h_1\right)}{h_1^2 + 1 + h_1^2 - 2h_1}$$

$$= \frac{2h_1 - 1}{2h_1^2 + 1 - 2h_1}$$

$$h_1 = \frac{e^x}{1+e^x}$$

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

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

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

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

$$= h_1(2x)\,(1+e^{2x}) - \frac{1}{1+e^{2x}}$$

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

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

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

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

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

$$= \frac{1}{2} \cdot \tanh\left(\frac{x}{2}\right)$$

Or $$\tanh\left(\frac{x}{2}\right) = 2\,h_1(x) - 1$$

or $$h_2(x) = 2\,h_1(2x) - 1$$

Since $h_2(x) \propto h_1(2x)$ ⇒ linear relation b/w tanh and sigmoid ⇒ changing to tanh will not affect a lot.

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

$$\Delta w_{ji} = -\eta \, \nabla E(w)$$

$$= -\eta \, \frac{\partial E(w)}{\partial w_{ji}}$$

$$\frac{\partial E(w)}{\partial w_{ji}} = \frac{1}{2} \sum_{d} \sum_{k} (t_{kd} - O_{kd}) \cdot (-x_{ji}) + 2\gamma \sum_{i,j} w_{ji}$$

$$\Delta w_{ji} = -\eta \left\{ \frac{1}{2} \sum_{d} \sum_{k} (t_{kd} - O_{kd})(-x_{ji}) + 2\gamma \sum_{i,j} w_{ji} \right\}$$

$$= \frac{1}{2} \eta \sum_{d} \cdot \sum_{k} (t_{kd} - O_{kd}) x_{ji} - 2\eta\gamma \sum_{i,j} w_{ji}$$

$$\Delta w_{ji} = \eta \sum_{d \in D} \sum_{k \in outputs} (t_{kd} - O_{kd}) x_{ji}$$

for backpropagation algorithm :

$$\frac{\partial E(w)}{\partial w_{ji}} = \sum_{k \in outputs} (t_{kd} - O_{kd}) \cdot$$

$$\sum_{d \in D} \delta_k (t \, x_{jk})$$

$$\frac{\partial E(w)}{\partial w_{ji}} = \frac{\partial}{\partial w} \left\{ \frac{1}{2} \sum_{d} \sum_{k} (t_{kd} - O_{kd})^2 \right\} + 2\gamma \sum_{i,j} w_{ji}$$

$$= \delta_k \cdot x_{ji} + 2\gamma \sum_{i,j} w_{ji}$$

$$\Delta w_{ji} = \eta \, \delta_k \cdot x_{ji} + 2\eta\gamma \sum_{i,j} w_{ji}$$

⇒ For weight updation

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

$$= w_{ji} + \eta \delta_k x_{ji} + 2\eta \gamma \sum_{j,i} w_{ji}$$

So, while updating the a single weight this can be approximated to,

$$w_{ji} = w_{ji} + \eta \delta_k x_{ji} + 2\eta \gamma w_{ji}$$

$$= (1 + 2\eta \gamma) w_{ji} + \eta \delta_k x_{ji}$$

$$w_{ji} = C w_{ji} + \eta \delta_k x_{ji}$$

⇒ In the traditional gradient descent, multiply the weight by a constant before updating.