

1.1

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = y \rightarrow \frac{dy}{dx} = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

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

$$\Delta w_{ji} = -\eta \frac{\partial E_d}{\partial w_{ji}}, \quad E_d(w) = \frac{1}{2} \sum_{k \in \text{outputs}} (t_k - o_k)^2$$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{\partial E_d}{\partial \text{net}_j} \frac{\partial \text{net}_j}{\partial w_{ji}} \quad \frac{\partial \text{net}_j}{\partial w_{ji}} = x_{ji}$$

Case I: j is an output

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

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

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

$$\rightarrow \Delta w_{ji} = \eta (t_j - o_j)(1 - o_j^2) x_{ji}$$

Case II: j is a hidden unit

$$\frac{\partial E_d}{\partial \text{net}_j} = \sum_{k \in \text{Downstream}(j)} \frac{\partial E_d}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial \text{net}_j} = \sum_{k \in \text{Downstream}(j)} -\delta_k \frac{\partial \text{net}_k}{\partial o_j} \frac{\partial o_j}{\partial \text{net}_j}$$

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

$$\rightarrow \Delta w_{ji} = \eta (1 - o_j^2) \sum_{k \in \text{Downstream}(j)} \delta_k w_{kj} x_{ji}$$

Hoang Nguyen

$$\frac{d \text{ReLU}(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & x < 0 \\ \text{undefined} & x = 0 \end{cases} = \text{ReLU}'(x)$$

Case 1: output layers $\rightarrow \Delta w_{ji} = \rho (t_j - o_j) \text{ReLU}'(x_i)$

Case 2: hidden layers $\rightarrow \Delta w_{ji} = \rho \text{ReLU}'(x_j) \sum_{k \in \text{downstream}} f_k w_{kj} x_{ij}$

$$\text{with } f_k = (t_k - o_k) \text{ReLU}'(x_k)$$

1. 2 Gradient Descent $O = w_0 + w_1(x_1 + x_1^2) + w_n(x_n + x_n^2)$

$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i} \left[\frac{1}{2} \sum_d (t_d - o_d)^2 \right], \text{ Gradient } \nabla E[\vec{w}] =$$

$$\left[\frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right]$$

Training rule: $\Delta \vec{w} = -\rho \nabla E[\vec{w}]$

$$\frac{\partial E}{\partial w_i} = \frac{1}{2} \sum_d 2(t_d - o_d) \frac{\partial}{\partial w_i} (t_d - o_d)$$

$$= \sum_d (t_d - o_d) \frac{\partial}{\partial w_i} (t_d - \vec{w} \cdot (\vec{x}_d + \vec{x}_d^2))$$

$$= -\sum_d (t_d - o_d) (x_{id} + x_{id}^2) \rightarrow \Delta w_i = \rho \sum_d (t_d - o_d) (x_{id} + x_{id}^2)$$

the final result for weight update

$$w_i^{\text{new}} = w_i^{\text{old}} + \rho \sum_d (t_d - o_d) (x_{id} + x_{id}^2)$$

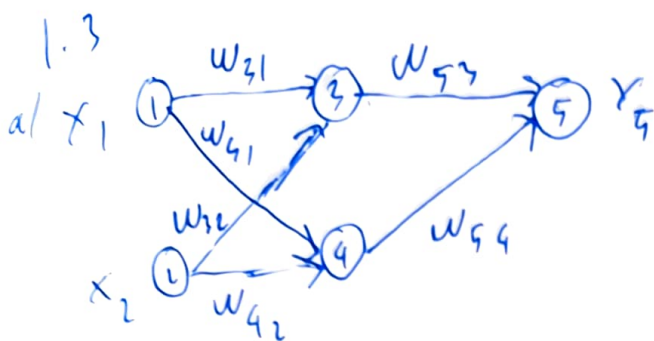
ρ is learning rate

D is the set of all training example

t_d is the target output for the d^{th} training example

o_d is the actual output for the d^{th} training example

x_{id} is the value of the i^{th} attribute for the d^{th} training example



$$\textcircled{3}: x_3 = h(w_{31} \cdot x_1 + w_{32} x_2)$$

$$x_4 = h(w_{41} x_1 + w_{42} x_2)$$

$$\rightarrow y_5 = h(w_{53} h(w_{31} x_1 + w_{32} x_2) + w_{54} h(w_{41} x_1 + w_{42} x_2))$$

b/ for $X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, $W^{(1)} = \begin{pmatrix} w_{31} & w_{32} \\ w_{41} & w_{42} \end{pmatrix}$ $W^{(2)} = (w_{53}, w_{54})$

Vector Format: $y_5 = h(W^{(2)} h(W^{(1)} \cdot X))$

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

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

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

$$O_t(x) = w_0 + w_h h_t(x) = w_0 + w_h [h_s(2x) - h_s(-2x)]$$

$$= w_0 + w_h h_s(2x) - w_h h_s(-2x)$$

$$O_s(x) = w'_0 + w_h h_s(x), \text{ assume } w_0 = 0 = w'_0 - w'_0$$

$$\rightarrow O_t(x) = w'_0 + w_h h_s(2x) - (w'_0 + w_h h_s(-2x))$$

$$= O_s(2x) - O_s(-2x)$$

$$\rightarrow O_t(x) = O_s(2x) - O_s(-2x)$$

\rightarrow We can see that $O_t(x)$ is exactly the same as the expression $O_s(x)$ with weights and biases adjusted, differing only by linear transformation and constant