



loss function: (Ordinary Least Squares) Back Propagation Algorithm

$$E(c_i) = \frac{1}{2} \| \mathbf{y}^{(i)} - \mathbf{o}^{(i)} \|^2 = \frac{1}{2} \sum_{k=1}^{n_L} (y_k^{(i)} - o_k^{(i)})^2$$

1> where $\mathbf{y}^{(i)}$ is the ground truth, $\mathbf{o}^{(i)}$ is the output of neural network.

2> The coefficient $\frac{1}{2}$ is unnecessary, which is for calculating easily in the future.

3> From the given picture of a neural network, let $n_L = 2$, $\mathbf{y}^{(i)} = (y_1^{(i)}, y_2^{(i)})^T$.

$$\text{Thus, } E(c_i) = \frac{1}{2} (y_1^{(i)} - a_1^{(3)})^2 + \frac{1}{2} (y_2^{(i)} - a_2^{(3)})^2$$

$$\text{If we go further into hidden layer, } E(c_i) = \frac{1}{2} (y_1^{(i)} - f(\underbrace{w_{11}^{(3)} a_1^{(2)} + w_{12}^{(3)} a_2^{(2)} + w_{13}^{(3)} a_3^{(2)} + b_1^{(3)}}_{z_1^{(3)}}))^2 + \frac{1}{2} (y_2^{(i)} - f(\underbrace{w_{21}^{(3)} a_1^{(2)} + w_{22}^{(3)} a_2^{(2)} + w_{23}^{(3)} a_3^{(2)} + b_2^{(3)}}_{z_2^{(3)}}))^2$$

If we go back further again into the input layer, $a_1^{(2)}, a_2^{(2)}, a_3^{(2)}$ can be replaced.

$$4> E_{\text{total}} = \frac{1}{N} \sum_{i=1}^N E(c_i)$$

The goal is to adjust the weights and ~~the~~ biases (W, b), decreasing the loss and finally acquiring the best weights and biases when the minimum of loss occurs.

5> Calculating the derivatives of weights in output layer by using derived chain rules.

$$\begin{aligned} \frac{\partial E}{\partial w_{11}^{(3)}} &= \frac{\partial \left[\frac{1}{2} (y_1^{(3)} - a_1^{(3)})^2 + \frac{1}{2} (y_2^{(3)} - a_2^{(3)})^2 \right]}{\partial w_{11}^{(3)}} \\ &= 2 \times \frac{1}{2} (y_1^{(3)} - a_1^{(3)}) \left(- \frac{\partial a_1^{(3)}}{\partial w_{11}^{(3)}} \right) + 0 \quad \because a_1^{(3)} = f(z_1^{(3)}) \\ &= - (y_1^{(3)} - a_1^{(3)}) f'(z_1^{(3)}) \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}} \Rightarrow = - (y_1^{(3)} - a_1^{(3)}) f'(z_1^{(3)}) \frac{\partial (w_{11}^{(3)} a_1^{(2)} + w_{12}^{(3)} a_2^{(2)} + w_{13}^{(3)} a_3^{(2)} + b_1^{(3)})}{\partial w_{11}^{(3)}} \\ &= - (y_1^{(3)} - a_1^{(3)}) f'(z_1^{(3)}) a_1^{(2)} \quad \textcircled{1} \quad \because \frac{\partial a_1^{(3)}}{\partial w_{11}^{(3)}} = f'(z_1^{(3)}) \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}} \end{aligned}$$

$$\text{Let } \delta_i^{(l)} = \frac{\partial E}{\partial z_i^{(l)}}, \text{ we will have: } \Rightarrow \because z_1^{(3)} = w_{11}^{(3)} a_1^{(2)} + w_{12}^{(3)} a_2^{(2)} + w_{13}^{(3)} a_3^{(2)} + b_1^{(3)} \\ \therefore \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}} = a_1^{(2)}$$

$$\frac{\partial E}{\partial w_{11}^{(3)}} = \frac{\partial E}{\partial z_1^{(3)}} \frac{\partial z_1^{(3)}}{\partial w_{11}^{(3)}} = \delta_1^{(3)} a_1^{(2)} \quad \textcircled{2}$$

$$\textcircled{1} \textcircled{2} \Rightarrow \delta_1^{(3)} = - (y_1^{(3)} - a_1^{(3)}) f'(z_1^{(3)}), \quad \delta_2^{(3)} = - (y_2^{(3)} - a_2^{(3)}) f'(z_2^{(3)})$$

$$\frac{\partial E}{\partial w_{12}^{(3)}} = \delta_1^{(3)} a_2^{(2)}, \quad \frac{\partial E}{\partial w_{13}^{(3)}} = \delta_1^{(3)} a_3^{(2)}$$

$$\frac{\partial E}{\partial w_{21}^{(3)}} = \delta_2^{(3)} a_1^{(2)}, \quad \frac{\partial E}{\partial w_{22}^{(3)}} = \delta_2^{(3)} a_2^{(2)}, \quad \frac{\partial E}{\partial w_{23}^{(3)}} = \delta_2^{(3)} a_3^{(2)}$$

Generally, assuming the number of layers is L , thus,

$$\delta_i^{(L)} = -(y_i - a_i^{(L)}) f'(z_i^{(L)}) \quad (1 \leq i \leq n_L)$$

L : the output layer

i : i^{th} neuron

n_L : the number of neurons in the L layer

$$\delta^{(L)} = -(y - a^{(L)}) \odot f'(z^{(L)})$$

$$\frac{\partial E}{\partial w_{ij}^{(L)}} = \delta_i^{(L)} a_j^{(L-1)} \quad (1 \leq i \leq n_L, 1 \leq j \leq n_{L-1}) \Leftrightarrow \nabla_{w^{(L)}} E = \delta^{(L)} (a^{(L-1)})^T$$

b) calculating the derivatives of weights in hidden layer.

$$\frac{\partial E}{\partial w_{ij}^{(L)}} = \frac{\partial E}{\partial z_i^{(L)}} \frac{\partial z_i^{(L)}}{\partial w_{ij}^{(L)}} = \delta_i^{(L)} a_j^{(L-1)} \quad (2 \leq L \leq L-1)$$

$$\delta_i^{(L)} = \frac{\partial E}{\partial z_i^{(L)}} = \sum_{j=1}^{n_{L+1}} \frac{\partial E}{\partial z_j^{(L+1)}} \frac{\partial z_j^{(L+1)}}{\partial z_i^{(L)}} = \sum_{j=1}^{n_{L+1}} \delta_j^{(L+1)} \frac{\partial z_j^{(L+1)}}{\partial z_i^{(L)}}$$

$$\therefore z_j^{(L+1)} = \sum_{i=1}^{n_L} w_{ji}^{(L+1)} a_i^{(L)} + b_j^{(L+1)} = \sum_{i=1}^{n_L} w_{ji}^{(L+1)} f(z_i^{(L)}) + b_j^{(L+1)}$$

$$\therefore \frac{\partial z_j^{(L+1)}}{\partial z_i^{(L)}} = \frac{\partial z_j^{(L+1)}}{\partial a_i^{(L)}} \frac{\partial a_i^{(L)}}{\partial z_i^{(L)}} = w_{ji}^{(L+1)} f'(z_i^{(L)})$$

$$\therefore \delta_i^{(L)} = \left(\sum_{j=1}^{n_{L+1}} \delta_j^{(L+1)} w_{ji}^{(L+1)} \right) f'(z_i^{(L)}) \Leftrightarrow \delta^{(L)} = (w^{(L+1)})^T \delta^{(L+1)} \odot f'(z^{(L)})$$

Thus, the $\delta^{(L+1)}$ in $L+1$ layer is calculated by $\delta^{(L)}$ in L layer, that is the reason why it's called Back Propagation.

T) ~~derivatives of~~ biases in output layer and hidden layer.

$$\frac{\partial E}{\partial b_i^{(L)}} = \frac{\partial E}{\partial z_i^{(L)}} \frac{\partial z_i^{(L)}}{\partial b_i^{(L)}} = \delta_i^{(L)} \Leftrightarrow \nabla_{b^{(L)}} E = \delta^{(L)}$$

8) Batch Gradient Descent:

$$w^{(L)} = w^{(L)} - \frac{\mu}{N} \sum_{i=1}^N \frac{\partial E_{(i)}}{\partial w^{(L)}}$$

$$b^{(L)} = b^{(L)} - \frac{\mu}{N} \sum_{i=1}^N \frac{\partial E_{(i)}}{\partial b^{(L)}}$$