



We will prove that the derivative of the cost function with respect to the weight $w_{kj}^{(l)}$ of the l^{th} hidden layer is $\frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial w_{kj}^{(l)}} = \alpha_j^{(l-1)} \delta_k^{(l)} = \alpha_j^{(l-1)} f'(z_k^{(l)}) \underbrace{\sum_m \delta_m^{(l+1)} w_{mk}^{(l+1)}}_{\delta_k^{(l)}}$, where f is the activation

function, i.e., $a_k^{(l)} = f(z_k^{(l)})$, and $\delta_m^{(l+1)}$ is the error propagated from layer $l + 1$, to layer l , i.e., $\delta_m^{(l+1)} = \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial z_m^{(l+1)}}$.

In the following, we will assume zero bias term for the sake of simplicity.

$$z_k^{(l)} = \sum_j w_{kj}^{(l)} \alpha_j^{(l-1)}, \quad z_m^{(l+1)} = \sum_k w_{mk}^{(l+1)} \alpha_k^{(l)}$$

$$\alpha_k^{(l)} = f(z_k^{(l)})$$

$$z_m^{(l+1)} = \sum_k w_{mk}^{(l+1)} \alpha_k^{(l)}$$

$$\begin{aligned} \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial w_{kj}^{(l)}} &= \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial z_k^{(l)}} \cdot \underbrace{\frac{\partial z_k^{(l)}}{\partial w_{kj}^{(l)}}}_{\alpha_j^{(l-1)}} \\ &= \frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial \alpha_k^{(l)}} \cdot \underbrace{\frac{\partial \alpha_k^{(l)}}{\partial z_k^{(l)}}}_{f'(z_k^{(l)})} \cdot \alpha_j^{(l-1)} \\ &= \left(\sum_m \underbrace{\frac{\partial J(\mathbf{W}, \mathbf{b})}{\partial z_m^{(l+1)}}}_{\delta_m^{(l+1)}} \frac{\partial z_m^{(l+1)}}{\partial \alpha_k^{(l)}} \right) \cdot f'(z_k^{(l)}) \cdot \alpha_j^{(l-1)} \\ &= \left(\sum_m \delta_m^{(l+1)} w_{mk}^{(l+1)} \right) f'(z_k^{(l)}) \cdot \alpha_j^{(l-1)} \end{aligned}$$