

1 Backpropagation

In this problem, you will derive the equations needed to calculate the gradient of the cost function in a feedforward neural network using backpropagation. This gradient is used within learning algorithms (such as gradient descent) that train the neural network.

Recall that the output from the j^{th} neuron in layer ℓ is given by

$$a_j^{(\ell)} = g_\ell(z_j^{(\ell)}),$$

where g_ℓ is a non-linear activation function and

$$z_j^{(\ell)} = \sum_{i=1}^{n_{\ell-1}} a_i^{(\ell-1)} W_{ij}^{(\ell)} + b_j^{(\ell)},$$

with $W_{ij}^{(\ell)}$ representing the weights and $b_j^{(\ell)}$ the biases of the network. Assume that the cost function C can be expressed as a sum over the dataset \mathcal{D} as

$$C = \frac{1}{|\mathcal{D}|} \sum_{\mathbf{x} \in \mathcal{D}} c(\mathbf{x}),$$

and define the notation

$$\delta_j^{(\ell)} = \frac{\partial c}{\partial z_j^{(\ell)}}.$$

a) Show that the quantity $\delta_j^{(\ell)}$ can be expressed as

$$\delta_j^{(L)} = \frac{\partial c}{\partial a_j^{(L)}} g_L'(z_j^{(L)}),$$

$$\delta_j^{(\ell)} = \sum_{k=1}^{n_{\ell+1}} \delta_k^{(\ell+1)} W_{kj}^{(\ell+1)T} g_\ell'(z_j^{(\ell)}) \quad (\text{for } \ell < L).$$

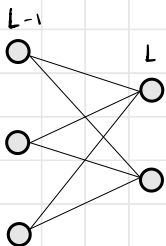
a)

$$\delta_j^{(L)} = \frac{\partial c}{\partial z_j^{(L)}} = \frac{\partial c}{\partial a_j^{(L)}} \cdot \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \quad \begin{array}{l} \text{chain rule} \\ \text{on the last layer} \end{array}$$

By definition, $a_j^{(L)} = g_L(z_j^{(L)}) \Rightarrow \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} = g_L'(z_j^{(L)})$

$$\Rightarrow \delta_j^{(L)} = \frac{\partial c}{\partial a_j^{(L)}} g_L'(z_j^{(L)})$$

Observe now layer $L-1$. Schematically, we have something like



Then,
$$\frac{\partial c}{\partial a_j^{(L-1)}} = \sum_{k=1}^{n_L} \frac{\partial z_k^{(L)}}{\partial a_j^{(L-1)}} \frac{\partial a_k^{(L)}}{\partial z_k^{(L)}} \frac{\partial c}{\partial a_k^{(L)}} \quad \begin{array}{l} \text{influenced by all} \\ \text{neurons in the previous layer} \end{array}$$

Although this was written for the second to last layer, the same would apply to layer $L-2$, just swapping $L \rightarrow L-1$. In general, we have

$$\frac{\partial C}{\partial a_j^{(L)}} = \sum_{k=1}^{n_{L+1}} w_{kj}^{(L+1)} \delta_k^{(L+1)}$$

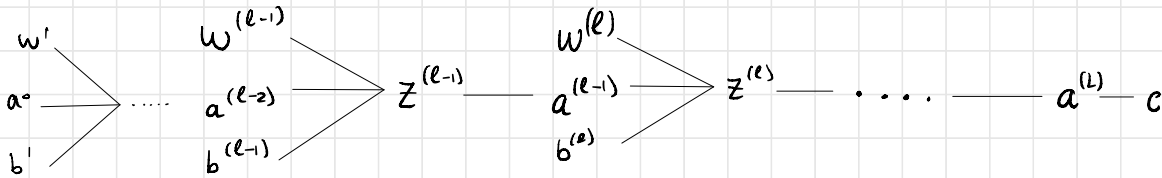
$$\Rightarrow \frac{\partial C}{\partial z_j^{(L)}} = \frac{\partial C}{\partial a_j^{(L)}} \left(\frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \right) = \sum_{k=1}^{n_{L+1}} w_{kj}^{(L+1)} \delta_k^{(L+1)}$$

$\left(\frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \right) = g'_L(z_j^{(L)})$

b) Show that the partial derivatives of the cost function with respect to the network's weights and biases can be calculated as

$$\frac{\partial C}{\partial w_{ij}^{(l)}} = a_i^{(l-1)} \delta_j^{(l)}, \quad \frac{\partial C}{\partial b_j^{(l)}} = \delta_j^{(l)}.$$

Functions are computed in chains:



Using the chain rule:

$$\frac{\partial C}{\partial w_{ij}^{(l)}} = \frac{\partial z_i^{(l)}}{\partial w_{ij}^{(l)}} \cdot \frac{\partial C}{\partial z_i^{(l)}} = a_i^{(l-1)} \cdot \delta_j^{(l)}$$

$\xleftarrow{\text{by def.}}$

$$\frac{\partial C}{\partial b_j^{(l)}} = \frac{\partial z_i^{(l)}}{\partial b_j^{(l)}} \cdot \frac{\partial C}{\partial z_i^{(l)}} = \delta_j^{(l)}$$

$\xleftarrow{\text{by def.}}$

$\frac{\partial z}{\partial b} = 1$ since $z = wa + b$