

2.5 Proof of the four fundamental equations (optional)

We'll now prove the four fundamental equations (BP1)–(BP4). All four are consequences of the chain rule from multivariable calculus. If you're comfortable with the chain rule, then I strongly encourage you to attempt the derivation yourself before reading on.

Let's begin with Equation (BP1), which gives an expression for the output error, δ^l . To prove this equation, recall that by definition

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}. \quad (2.14)$$

Applying the chain rule, we can re-express the partial derivative above in terms of partial derivatives with respect to the output activations,

$$\delta_j^l = \sum_k \frac{\partial C}{\partial a_k^l} \frac{\partial a_k^l}{\partial z_j^l}, \quad (2.15)$$

where the sum is over all neurons k in the output layer. Of course, the output activation a_k^l of the k -th neuron depends only on the weighted input z_j^l for the j -th neuron when $k = j$. And so $\partial a_k^l / \partial z_j^l$ vanishes when $k \neq j$. As a result we can simplify the previous equation to

$$\delta_j^l = \frac{\partial C}{\partial a_j^l} \frac{\partial a_j^l}{\partial z_j^l}. \quad (2.16)$$

Recalling that $a_j^l = \sigma(z_j^l)$ the second term on the right can be written as $\sigma'(z_j^l)$, and the equation becomes

$$\delta_j^l = \frac{\partial C}{\partial a_j^l} \sigma'(z_j^l), \quad (2.17)$$

which is just (BP1), in component form. Next, we'll prove (BP2), which gives an equation for the error δ^l in terms of the error in the next layer, δ^{l+1} . To do this, we want to rewrite $\delta_j^l = \partial C / \partial z_j^l$ in terms of $\delta_k^{l+1} = \partial C / \partial z_k^{l+1}$. We can do this using the chain rule,

$$\delta_j^l = \frac{\partial C}{\partial z_j^l} = \sum_k \frac{\partial C}{\partial z_k^{l+1}} \frac{\partial z_k^{l+1}}{\partial z_j^l} = \sum_k \frac{\partial z_k^{l+1}}{\partial z_j^l} \delta_k^{l+1}, \quad (2.18)$$

where in the last line we have interchanged the two terms on the right-hand side, and substituted the definition of δ_k^{l+1} . To evaluate the first term on the last line, note that

$$z_k^{l+1} = \sum_j w_{kj}^{l+1} a_j^l + b_k^{l+1} = \sum_j w_{kj}^{l+1} \sigma(z_j^l) + b_k^{l+1}. \quad (2.19)$$

Differentiating, we obtain

$$\frac{\partial z_k^{l+1}}{\partial z_j^l} = w_{kj}^{l+1} \sigma'(z_j^l). \quad (2.20)$$