

$L$  = number of layers in a network.

$l$  = a particular layer in the network:  $1 \leq l \leq L$

$a_i^l$  = value of the  $i$ th node in layer  $l$

$w_{ij}^l$  = The weight of the connection between  $a_j^{l-1}$  and  $a_i^l$

$b_i^l$  = bias of node  $i$  in the layer  $l$ .

$y_i$  = the expected outcome from a given input.

### Generating the output

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

← keeps our nodes between 0 & 1.

$$\text{let } z_i^l = w_{ij}^l a_j^{l-1} + b_i^l$$

← this will be useful later.

$$\Rightarrow a_i^l = \sigma(z_i^l)$$

### Cost function

$$C = \frac{1}{n_L} \sum_i^{n_L} (a_i^L - y_i)^2$$

$n_L$  = length of layer  $L$ .

$$\nabla C = \left( \frac{\partial C}{\partial w_{ij}^l}, \frac{\partial C}{\partial b_j^l} \right)$$

$$\forall l, 1 \leq l \leq L$$

Essentially one big vector to contain

the change in cost with respect to every ~~the~~ weight and bias in the network.

$$\frac{\partial C}{\partial w_{AB}^l} = \frac{\partial C}{\partial a^L} \frac{\partial a^L}{\partial a_A^l} \frac{\partial a_A^l}{\partial w_{AB}^l} = \frac{\textcircled{1}}{\partial a^L} \frac{\textcircled{2}}{\partial a^{l+1}} \frac{\textcircled{3}}{\partial a_A^l} \frac{\textcircled{4}}{\partial w_{AB}^l}$$

Note: I'm using A, B as fixed indices, so ~~this~~ ~~formula~~ only shows the cost w.r.t ONE specific weight.

Having got this formula, we can calculate each component separately.



$$\textcircled{1} \quad \frac{\partial \mathcal{L}}{\partial a_i^L} = \frac{1}{n_L} \sum \frac{\partial}{\partial a_i^L} (a_i^L - y_i)^2 = \frac{2}{n_L} (a_i^L - y_i)$$

↑ This is a vector. It helps a lot to visualise what it means in the context of the network.

Matrix

$$\frac{\partial a_i^L}{\partial a_j^{L+1}} = \underbrace{\frac{\partial a_i^L}{\partial a_i^{L-1}} \frac{\partial a_i^{L-1}}{\partial a_i^{L-2}} \dots \frac{\partial a_i^{L+2}}{\partial a_j^{L+1}}}_{\text{These are all matrices.}} = \left( \sigma'(\mathbf{Z}_i^L) \mathbf{W}_{ip}^L \right) \left( \sigma'(\mathbf{Z}_p^{L-1}) \mathbf{W}_{pm}^{L-1} \right) \dots \left( \sigma'(\mathbf{Z}_q^{L+2}) \mathbf{W}_{jq}^{L+2} \right)$$

$$\textcircled{2} \quad \frac{\partial a_m^{L+1}}{\partial a_p^L} = \frac{\partial}{\partial a_p^L} \left( \sigma(\mathbf{W}_{mp}^{L+1} a_p^L + b_m^{L+1}) \right) = \underbrace{\sigma'(\mathbf{Z}_m^{L+1}) \mathbf{W}_{mp}^{L+1}}_{\substack{\uparrow \text{ A matrix with } p \text{ columns} \\ \text{and } m \text{ rows.}}}$$

That's the worst one done. Phew! The indices are very tricky to code.

$$\textcircled{3} \quad \frac{\partial a_j^{L+1}}{\partial a_A^L} = \frac{\partial}{\partial a_A^L} \left( \sigma(\mathbf{W}_{jA}^{L+1} a_A^L + b_j^{L+1}) \right) = \left( \sigma'(\mathbf{Z}_j^{L+1}) \mathbf{W}_{jA}^{L+1} \right)$$

↑ vector. Again, it ~~is~~ ~~helpful~~ ~~to~~ ~~think~~ ~~of~~ ~~it~~ visually; How does each node in layer  $L+1$  change w.r.t the node  $a_A^L$ .

$$\textcircled{4} \quad \frac{\partial a_A^L}{\partial w_{AB}^L} = \frac{\partial}{\partial w_{AB}^L} \left( \sigma(\mathbf{W}_{AB}^L a_B^{L-1} + b_A^L) \right) = \sigma'(\mathbf{Z}_A^L) a_B^{L-1} \quad \leftarrow \text{Scalar.}$$

Quick Remark:  $a_A^L = \sigma(\mathbf{W}_{Ai}^L a_i^{L-1} + b_A^L)$  ~~for all  $i$~~   $0 \leq i \leq n_{L-1}$

but since,  $\frac{\partial}{\partial w_{AB}^L} \mathbf{W}_{Ai}^L = \begin{cases} 1 & \text{if } i=B \\ 0 & \text{otherwise} \end{cases}$

I've left it out.

Most of this, except eq (4), is the same for  $\frac{\partial \mathcal{L}}{\partial b^L}$ .