

$a_i^{(l)}$ = preactivation of unit i in layer l

$$a_i^{(l)} = \sum_{j=1}^{N_{L-2}} w_{ij}^{(l)} z_j^{(l-2)} + w_{io}^{(l)}$$

$$\frac{\partial a_i^{(l)}}{\partial w_{ij}^{(l)}} = z_j^{(l-2)}$$

$z_i^{(l)}$ = postactivation for unit i in layer l

$$z_i^{(l)} = \sigma(a_i^{(l)}) = \sigma\left(\sum_{j=1}^{N_{L-2}} w_{ij}^{(l)} z_j^{(l-1)} + w_{io}^{(l)}\right)$$

$$p(t=0|x) = y(x; \omega)$$

$$p(t=1|x) = 1 - y(x; \omega)$$

$$\begin{aligned} p(\{t(x^{(i)}) = t^{(i)}\}_{i=1}^N | x) &= \prod_{i=1}^N p(t(x^{(i)}) = t^{(i)} | x) \\ &= \prod_{i=1}^N y(x^{(i)}; \omega)^{1-t^{(i)}} (1 - y(x^{(i)}; \omega))^{t^{(i)}} \end{aligned}$$

→ taking the log

$$L(\omega) = - (1 - t^{(i)}) \log y(x^{(i)}; \omega) - t^{(i)} \log (1 - y(x^{(i)}; \omega))$$

1) We want

$$\frac{\partial L}{\partial w_{ij}^{(k)}} = \frac{\partial L}{\partial a_i^{(k)}} \frac{\partial a_i^{(k)}}{\partial w_{ij}^{(k)}}$$

$L \leftarrow a_i^{(k)}, \dots$
 $w_{ij}^{(k)}$

$\delta_i^{(k)}$ $z_j^{(k-1)}$

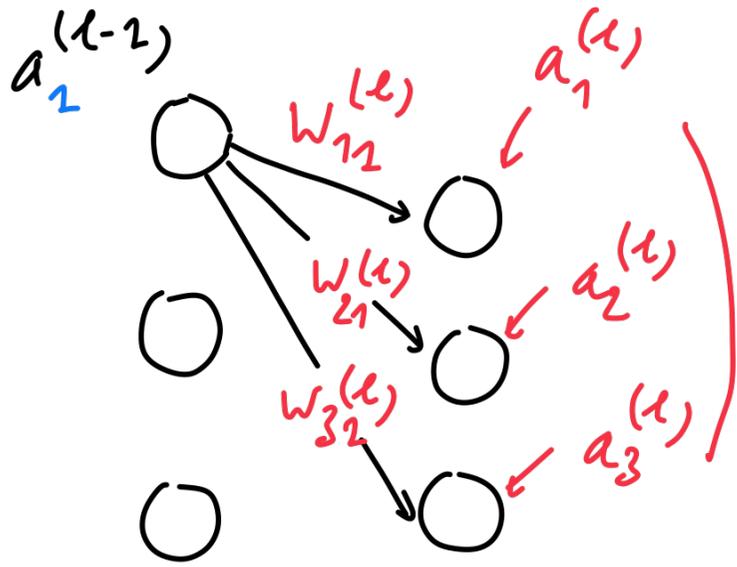
$$2) \delta_{out} = \frac{\partial L}{\partial a_{out}} \rightarrow \frac{\partial (-(1-t^{(k)}) \log \sigma(a_{out}) - t^{(k)} \log(1 - \sigma(a_{out})))}{\partial a_{out}}$$

$$= - (1-t^{(k)}) \frac{\sigma'}{\sigma(a_{out})} + t^{(k)} \frac{\sigma'}{1 - \sigma(a_{out})}$$

$\sigma'(a) = \sigma(a)(1 - \sigma(a))$

$$= - (1-t^{(k)}) (1 - \sigma(a_{out})) + t^{(k)} \sigma(a_{out})$$

$$\delta_{out} = \sigma(a_{out}) - (1-t^{(k)})$$



$$L(a_1^{(l)}, a_2^{(l)}, a_3^{(l)})$$

$$\frac{\partial L}{\partial a_1^{(l-2)}} \leftarrow L(a_1^{(l)}(a_1^{(l-2)}), a_2^{(l)}(a_2^{(l-2)}))$$

$$\frac{\partial L}{\partial a_i^{(l-2)}} = \sum_{j=2}^{N_L} \frac{\partial L}{\partial a_j^{(l)}} \cdot \frac{\partial a_j^{(l)}}{\partial a_i^{(l-2)}}$$

$$\delta_i^{(l-2)} = \sum_{j=2}^{N_L} \delta_j^{(l)}$$

$$\frac{\partial a_j^{(l)}}{\partial a_i^{(l-2)}} \rightarrow a_j^{(l)} = \sum_{i=1}^{N_{l-2}} W_{ji}^{(l)} \cdot \sigma(a_i^{(l-2)}) + W_{j0}^{(l)}$$

$$\frac{\partial a_j^{(l)}}{\partial a_i^{(l-2)}} = W_{ji}^{(l)} \sigma'(a_i^{(l-2)})$$

Backpropagation

1) \rightarrow Forward propagate $x^{(i)}$ through the network, and get all $a_i^{(l)}, z_i^{(l)}$

2) you get $\delta_{out} = \frac{\partial L}{\partial a_{out}} = \sigma(a_{out}) - (1 - t^{(i)})$

3) Backpropagate the δ_{out}

$$\delta_i^{(l-2)} = \sum_{j=1}^{N_L} \delta_j^{(l)} w_{ji}^{(l)} \cdot \sigma'(a_i^{(l-2)})$$

4) Get the gradient as $\frac{\partial L}{\partial w_{ij}^{(l)}} = \delta_i^{(l)} \cdot z_j^{(l-1)}$

$$\frac{\partial L}{\partial w_{i0}^{(l)}} = \delta_i^{(l)}$$

