

Q. Derive the backprop algorithm for the neural network + loss defined by

$$\tilde{J} = \frac{1}{N} \sum_{i=1}^N \sum_{k=1}^K -y_{ik} \log a_{ik}^{(3)} - (1-y_{ik}) \log (1-a_{ik}^{(3)}).$$

where $a_i^{(3)}$ is a k -dim vector of activations for the i th input. Omitting index i , recall that

$$a^{(3)} = \varphi(z^{(3)})$$

$$z^{(3)} = W^{(2)} a^{(2)} + b^{(2)}$$

$$a^{(2)} = \varphi(z^{(2)})$$

$$z^{(2)} = W^{(1)} a^{(1)} + b^{(1)}$$

$$a^{(1)} = X, \text{ input.}$$

We need to compute the gradients

$$\left(\frac{\partial \tilde{J}}{\partial W_{\alpha\beta}^{(2)}}, \frac{\partial \tilde{J}}{\partial b_{\alpha}^{(2)}}, \frac{\partial \tilde{J}}{\partial W_{\alpha\beta}^{(3)}}, \frac{\partial \tilde{J}}{\partial b_{\alpha}^{(3)}} \right) \text{ for every } \alpha, \beta.$$

For simplicity, ~~we~~ we first compute the gradients wrt

$$\tilde{J}^i = \sum_{k=1}^K -y_{ik} \log a_{ik}^{(3)} - (1-y_{ik}) \log (1-a_{ik}^{(3)}).$$

for a single input i . From this we can retrieve the gradients wrt \tilde{J} via

$$\frac{\partial \tilde{J}}{\partial W_{\alpha\beta}^{(2)}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \tilde{J}^i}{\partial W_{\alpha\beta}^{(2)}} + \lambda W_{\alpha\beta}^{(2)} \text{ etc. } (*)$$

(due to a.2d)

For simplicity we drop the index i for the computation of $\frac{\partial \hat{J}}{\partial W_{\alpha p}^{(2)}} = \frac{\partial \hat{J}^i}{\partial W_{\alpha p}^{(2)}}$ etc.

First compute

$$\eta_{\alpha}^{(3)} := \frac{\partial \hat{J}}{\partial a_{\alpha}^{(3)}} = -\frac{y_{\alpha}}{a_{\alpha}^{(3)}} + \frac{1-y_{\alpha}}{1-a_{\alpha}^{(3)}} \quad (\text{due to Q.2c}).$$

$\forall \alpha$

Then, compute via chain rule.

$$\begin{aligned} s_{\alpha}^{(3)} &:= \frac{\partial \hat{J}}{\partial z_{\alpha}^{(3)}} = \frac{\partial \hat{J}}{\partial a_{\alpha}^{(3)}} \frac{\partial a_{\alpha}^{(3)}}{\partial z_{\alpha}^{(3)}}. \\ &= \eta_{\alpha}^{(3)} \varphi(z_{\alpha}^{(3)}) (1 - \varphi(z_{\alpha}^{(3)})). \end{aligned}$$

(due to Q.2b)
 $\forall \alpha$
 \Rightarrow line 5 of Alg1.

In vectorial notation,

$$f^{(3)} = \eta^{(3)} \cdot \varphi(z^{(3)}) \cdot (1 - \varphi(z^{(3)})).$$

pointwise multiplications. \Rightarrow line 8 of Alg1.

We're ready to compute via multivar chain rule:

$$\begin{aligned} \frac{\partial \hat{J}}{\partial W_{\alpha\beta}^{(2)}} &= \sum_{\gamma} \frac{\partial \hat{J}}{\partial z_{\gamma}^{(3)}} \frac{\partial z_{\gamma}^{(3)}}{\partial W_{\alpha\beta}^{(2)}} \\ &= \sum_{\gamma} f'_{\gamma} \cdot \frac{\partial z_{\gamma}^{(3)}}{\partial W_{\alpha\beta}^{(2)}}. \end{aligned}$$

Note that $z_{\gamma}^{(3)} = \sum_{\tau} W_{\gamma\tau}^{(2)} a_{\tau}^{(2)} + b_{\gamma}^{(2)}$ and so

$$\begin{aligned} \frac{\partial z_{\gamma}^{(3)}}{\partial W_{\alpha\beta}^{(2)}} &= \sum_{\tau} \frac{\partial W_{\gamma\tau}^{(2)}}{\partial W_{\alpha\beta}^{(2)}} a_{\tau}^{(2)} \\ &= \sum_{\tau} \delta_{\alpha\tau} \delta_{\tau\beta} a_{\tau}^{(2)} \\ &= \delta_{\alpha\beta} a_{\beta}^{(2)}. \end{aligned}$$

$$\begin{aligned} &= \sum_{\gamma} f'_{\gamma} \delta_{\alpha\beta} a_{\beta}^{(2)}. \\ &= f'_{\alpha} a_{\beta}^{(2)}. \end{aligned}$$

In vectorial notation,

$$\frac{\partial W^{(3)}}{\partial W^{(2)}} = f^{(3)} a^{(2)T}, \text{ an outer product.}$$

← line 11 of Alg1.

and similarly,

$$\partial b_{\alpha}^{(3)} := \sum_{\gamma} \frac{\partial J}{\partial z_{\gamma}^{(3)}} \frac{\partial z_{\gamma}^{(3)}}{\partial b_{\alpha}^{(2)}}$$

$$= \sum_{\gamma} \cancel{\delta_{\gamma}} \delta_{\gamma}^{(3)} \cancel{\delta_{\alpha\gamma}}.$$

$$= \delta_{\alpha}^{(3)}.$$

← line 13 of Alg1.

Now we compute $\eta^{(2)}$ and $f^{(2)}$ to recursively obtain $\partial W^{(1)}$ and $\partial b^{(1)}$.

$$\begin{aligned} \eta_{\alpha}^{(2)} &:= \frac{\partial J}{\partial a_{\alpha}^{(2)}} = \sum_{\gamma} \frac{\partial J}{\partial z_{\gamma}^{(3)}} \frac{\partial z_{\gamma}^{(3)}}{\partial a_{\alpha}^{(2)}}. \\ &= \sum_{\gamma} \delta_{\gamma}^{(3)} \frac{\partial z_{\gamma}^{(3)}}{\partial a_{\alpha}^{(2)}}. \end{aligned}$$

$$\left[\begin{aligned} \text{Since } z_{\gamma}^{(3)} &= \sum_{\tau} W_{\gamma\tau}^{(2)} a_{\tau}^{(2)} + b_{\gamma}^{(2)}, \\ \frac{\partial z_{\gamma}^{(3)}}{\partial a_{\alpha}^{(2)}} &= \sum_{\tau} \cancel{W_{\gamma\tau}^{(2)}} \delta_{\alpha\tau} = W_{\gamma\alpha}^{(2)}. \\ &= \sum_{\tau} \delta_{\tau}^{(3)} W_{\tau\alpha}^{(2)} \end{aligned} \right]$$

In vectorial notations,

$$\eta^{(2)} = W^{\top} \delta^{(3)}.$$

← line 15 of Alg1.

It easily follows that

$$J^{(2)} = \eta^{(2)} - \varphi(z^{(2)}) \cdot \varphi'(z^{(2)}).$$

← line 17 of Alg 1.

And analogous computations lead to

$$\partial W^{(1)} \text{ and } \partial b^{(1)}.$$

We have thus derived the derivative wrt \tilde{J}
and from (*) follows the derivative wrt \tilde{J} . □

← $\frac{1}{N}$ factors at lines 11, 13,
line 26 of Alg 1.