

# MLE Ex 4

## Task 1

① We have:

$$\mathcal{L} = \mathcal{L}_0(x, y) = - \sum_{i=1}^c y_i \ln \hat{y}_i$$

so

$$\delta_{ij}^{(L)} = \frac{\partial \mathcal{L}}{\partial \tilde{z}_j^{(L)}}$$

$$= \frac{\partial \mathcal{L}}{\partial \tilde{z}_j^{(L)}} - \sum_{i=1}^c y_i \ln \hat{y}_i$$

$$= - \sum_{i=1}^c y_i \frac{\partial}{\partial \tilde{z}_j^{(L)}} \ln \hat{y}_i \quad (*)$$

Furthermore

$$\ln \hat{y}_i = \ln \left( \frac{\exp \tilde{z}_i^{(L)}}{\sum_{m=1}^c \exp \tilde{z}_m^{(L)}} \right)$$

$$= \tilde{z}_i^{(L)} - \ln \sum_{m=1}^c \exp \tilde{z}_m^{(L)}$$

$$\Rightarrow \frac{\partial \ln \hat{y}_i}{\partial \tilde{z}_j^{(L)}} = \delta_{ij} - \frac{1}{\sum_{m=1}^c \exp \tilde{z}_m^{(L)}} \frac{\partial}{\partial \tilde{z}_j^{(L)}} \sum_{m=1}^c \exp \tilde{z}_m^{(L)}$$

$$= \delta_{ij} - \frac{1}{\sum_{m=1}^c \exp \tilde{z}_m^{(L)}} \exp \tilde{z}_j^{(L)}$$

$$= \delta_{ij} - \hat{y}_j$$

Plugging into (\*):

$$\delta_j^{(L)} = - \sum_{i=1}^c y_i (\delta_{ij} - \hat{y}_j)$$

$$= \left( \sum_{i=1}^c y_i \hat{y}_j \right) - y_j$$

$$= \hat{y}_i - y_i \quad (\sum y_i = 1 \text{ \& classification task})$$

Thus, from  $\delta_i^{(L)} = \hat{y}_i - y_i$  we have  
 $\underline{\delta}^{(L)} = \hat{\underline{y}} - \underline{y}$ , as desired

②a) We have (for  $i \in [d_L]^{m \times k}$ ,  $j \in [d_L]$ )

$$(\nabla_{\underline{W}^{(L-1)}} L)_{ij} = \sum_{k=1}^c \frac{\partial L}{\partial \tilde{z}_k^{(L)}} \frac{\partial \tilde{z}_k^{(L)}}{\partial W_{ij}^{(L-1)}} \quad (*)$$

$\delta_k^{(L)}$

where

$$\begin{aligned} \tilde{z}_k^{(L)} &= (W^{(L-1)} Z^{(L-1)} + b^{(L-1)})_k \\ &= \left( \sum_{m=1}^{d_L} W_{km}^{(L-1)} Z_m^{(L-1)} \right) + b_k^{(L-1)} \\ \Rightarrow \frac{\partial \tilde{z}_k^{(L)}}{\partial W_{ij}^{(L-1)}} &= \sum_{m=1}^{d_L} \delta_{ik} \delta_{jm} Z_m^{(L-1)} \\ &= \delta_{ik} Z_j^{(L-1)} \end{aligned}$$

Plugging into (\*):

$$\begin{aligned} (\nabla_{\underline{W}^{(L-1)}} L)_{ij} &= \sum_{k=1}^c \delta_k^{(L)} \delta_{ik} Z_j^{(L-1)} \\ &= \delta_i^{(L)} Z_j^{(L-1)} \\ &= (\underline{\delta}^{(L)} (\underline{Z}^{(L-1)})^T)_{ij} \end{aligned}$$

So

$$\nabla_{\underline{W}^{(L-1)}} L = \underline{\delta}^{(L)} (\underline{Z}^{(L-1)})^T$$

As desired. (Can plug in  $\underline{\delta}^{(L)} = \hat{\underline{y}} - \underline{y}$ ).

b) Proceed similarly: (for  $i \in [d_L]$ )  $\leq [C]$

$$(\nabla_{\underline{b}^{(L-1)}} L)_i = \sum_{k=1}^C \underbrace{\frac{\partial L}{\partial \tilde{z}_k^{(L)}}}_{\delta_k^{(L)}} \frac{\partial \tilde{z}_k^{(L)}}{\partial b_i^{(L-1)}} \quad (*)$$

where

$$\tilde{z}_k^{(L)} = \left( \sum_{m=1}^{d_L} W_{km}^{(L-1)} z_m^{(L-1)} \right) + b_k^{(L-1)}$$

$$\Rightarrow \frac{\partial \tilde{z}_k^{(L)}}{\partial b_i^{(L-1)}} = \delta_{ik}$$

Plug into (\*):

$$\begin{aligned} (\nabla_{\underline{b}^{(L-1)}} L)_i &= \sum_{k=1}^C \delta_k^{(L)} \delta_{ik} \\ &= \delta_i^{(L)} \end{aligned}$$

$$\Rightarrow \nabla_{\underline{b}^{(L-1)}} L = \underline{\delta}^{(L)}$$

As desired. (Can plug in  $\underline{\delta}^{(L)} = \underline{\hat{y}} - \underline{y}$ ).

Dimensions of the gradients:

We computed the components  $\frac{\partial L}{\partial b_i^{(L-1)}} (\nabla_{\underline{b}^{(L-1)}} L)_i$  for  $i \in [d_{L-1}] = [C]$  and  $i \in [d_L]$ , and  $(\nabla_{\underline{b}^{(L-1)}} L)_i$  for  $i \in [d_L] = [C]$ , so necessarily they're of dimension  $C \times d_{L-1}$  &  $C$  respectively. This is because  $\nabla_a f$  is defined as a tensor with the same dimension as tensor  $a$ , for any  $a$  or  $f$ .

In this case, from the results,  $\dim(\nabla_{\underline{b}^{(L-1)}} L) = \dim \underline{\delta}^{(L)} = C$  &  $\dim(\nabla_{\underline{b}^{(L-1)}} L) = \dim \underline{\delta}^{(L)} (\underline{z}^{(L-1)})^T = \dim \underline{\delta}^{(L)} \times \dim \underline{z}^{(L-1)} = C \times d_{L-1}$ , so they match.

③. Since  $\varphi(x) = \frac{1}{2} \max(2, x)$ , we have applying the convention for  $x=0$

$$\varphi'(x) = \mathbb{1}_{\{x > 0\}}$$

Thus

$$\delta_i^{(L-1)} = [(W^{(L-1)})^T (\hat{y} - y) \odot \varphi'(\tilde{z}^{(L-1)})]_i$$

$$= [(W^{(L-1)})^T (\hat{y} - y)]_i \cdot [\varphi'(\tilde{z}^{(L-1)})]_i \quad (\text{Hadamard prod.})$$

$$= [(W^{(L-1)})^T (\hat{y} - y)]_i \varphi'(\tilde{z}_i^{(L-1)}) \quad (\text{elementwise application})$$

$$= [(W^{(L-1)})^T (\hat{y} - y)]_i \mathbb{1}_{\{\tilde{z}_i^{(L-1)} > 0\}}$$

$$\text{Thus, } \tilde{z}_i^{(L-1)} \leq 0 \Rightarrow \mathbb{1}_{\{\tilde{z}_i^{(L-1)} > 0\}} = 0 \Rightarrow \delta_i^{(L-1)} = 0$$

$\varphi$  isn't diff'ble at 0 (at 0, the left derivative is 0 while the right derivative is 1), so we apply the convention  $\varphi'(0) = 0$ . As a consequence, the  $\tilde{z}_i^{(L-1)} = 0$  case is the same as the  $\tilde{z}_i^{(L-1)} \leq 0$  case.

Recall from the hint that

$$\nabla_{b^{(L)}} L = \underline{\delta}^{(L-1)}, \quad \nabla_{W^{(L)}} L = \underline{\delta}^{(L-1)} \underline{z}^{(L)T}$$

$$\nabla_{b^{(L-1)}} L = \underline{\delta}^{(L)}, \quad \nabla_{W^{(L-1)}} L = \underline{\delta}^{(L-1)} (\underline{z}^{(L)})^T$$

Thus  $(\nabla_{b^{(L-1)}} L)_i = 0$  and  $(\nabla_{W^{(L-1)}} L)_{ij} = 0 \quad \forall i \in [d_{L-1}]$  so all the weights & the bias <sup>parameters</sup> for neuron  $i$  in layer  $L-1$  will not update under a gradient step.

How this propagates backwards through the prior layers is more difficult to comment on:

$$\delta^{(L-2)} = (W^{(L-2)})^T \delta^{(L-1)} \odot \varphi'(\underline{z}^{(L-2)})$$

is the decisive equation since the error signals are (very) closely related to the parameter derivatives. However, the effect of  $\delta_j^{(L-1)} = 0$  on the product  $(W^{(L-2)})^T \delta_j^{(L-1)}$  isn't enormously significant. Prior components can still have non-zero derivative (they affect  $L$ , just not on a path through neuron  $j$  of layer  $(L-1)$ ).