# Homework 1

## Solution 1.2

Change in notation:

$$\text{Output of } Linear_1 : z^{(1)} \rightarrow s^{(1)}.$$
$$\text{Output of } f : z^{(2)} \rightarrow z^{(1)}.$$
$$\text{Output of } Linear_2 : z^{(3)} \rightarrow s^{(2)}.$$
$$\text{Output of } g : \text{Remains the same, } \hat{y}.$$

## Solution a)

1. `torch.nn.Linear` : $\text{Linear}(x) = Wx + b$.
2. `torch.nn.ReLU` : $\text{ReLU}(x) = \max(0, x)$.
3. `torch.nn.Linear` : $\text{Linear}(x) = Wx + b$.
4. `torch.nn.ReLU` : $\text{ReLU}(x) = \max(0, x)$.
5. `torch.nn.MSELoss` : $l_{\text{MSE}}(\hat{y}, y) = ||\hat{y} - y||^2$.

## Solution b)

Strictly using $x, y, W^{(1)}, W^{(2)}, b^{(1)}, b^{(2)}$:

| Layer | Input | Output |
|---|---|---|
| $Linear_1$ | $x$ | $W^{(1)}x + b^{(1)}$ |
| $f$ | $W^{(1)}x + b^{(1)}$ | $\text{ReLU}(W^{(1)}x + b^{(1)})$ |
| $Linear_2$ | $\text{ReLU}(W^{(1)}x + b^{(1)})$ | $W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)}$ |
| $g$ | $W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)}$ | $\text{I}(W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)})$ |
| $Loss$ | $\hat{y}, \text{I}(W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)})$ | $(\hat{y} - \text{I}(W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)}))(\hat{y} - \text{I}(W^{(2)}\text{ReLU}(W^{(1)}x + b^{(1)}) + b^{(2)}))^T$ |

Using intermediate variables:

| Layer | Input | Output |
|---|---|---|
| $Linear_1$ | $x$ | $s^{(1)} = W^{(1)}x + b^{(1)}$ |
| f | $s^{(1)}$ | $z^{(1)} = \text{ReLU}(s^{(1)})$ |
| $Linear_2$ | $z^{(1)}$ | $s^{(2)} = W^{(2)}z^{(1)} + b^{(2)}$ |
| g | $s^{(2)}$ | $\hat{y} = \text{I}(s^{(2)})$ |
| Loss | $\hat{y}, y$ | $\ell_{\text{MSE}} = (\hat{y} - y)(\hat{y} - y)^T$ |

Using components:

| Layer | Input | Output |
|---|---|---|
| $Linear_1$ | $x_j$ | $s_i^{(1)} = \sum_j W_{ij}^{(1)} x_j + b_i^{(1)}$ |
| f | $s_i^{(1)}$ | $z_i^{(1)} = \text{ReLU}(s_i^{(1)})$ |
| $Linear_2$ | $z_i^{(1)}$ | $s_k^{(2)} = \sum_i W_{ki}^{(2)} z_i^{(1)} + b_k^{(2)}$ |

| Layer | Input | Output |
|-------|-------|--------|
| g | $s_k^{(2)}$ | $y_k = g(s_k^{(2)})$ |
| Loss | $\hat{y}_k, y_k$ | $\ell_{\mathrm{MSE}} = \sum_k (\hat{y}_k - y_k)(\hat{y}_k - y_k)$ |

## Solution c)

### Dimensions

Following [numerator layout](#):

$$\boldsymbol{x} : d_{\boldsymbol{x}} \times 1.$$
$$\boldsymbol{s}^{(1)} : d_{\boldsymbol{s}^{(1)}} \times 1.$$
$$\boldsymbol{z}^{(1)} : d_{\boldsymbol{z}^{(1)}} \times 1.$$
$$\boldsymbol{s}^{(2)} : d_{\boldsymbol{s}^{(2)}} \times 1.$$
$$\hat{\boldsymbol{y}} : d_{\hat{\boldsymbol{y}}} \times 1.$$
$$W^{(1)} : d_{\boldsymbol{s}^{(1)}} \times d_{\boldsymbol{x}}.$$
$$W^{(2)} : d_{\boldsymbol{s}^{(2)}} \times d_{\boldsymbol{z}^{(1)}}.$$
$$b^{(1)} : d_{\boldsymbol{s}^{(1)}} \times 1.$$
$$b^{(2)} : d_{\boldsymbol{s}^{(2)}} \times 1.$$
$$\frac{\partial \ell}{\partial W^{(2)}} : d_{\boldsymbol{z}^{(1)}} \times d_{\boldsymbol{s}^{(2)}}.$$
$$\frac{\partial \ell}{\partial W^{(1)}} : d_{\boldsymbol{x}} \times d_{\boldsymbol{s}^{(1)}}.$$
$$\frac{\partial \ell}{\partial b^{(2)}} : 1 \times d_{\boldsymbol{s}^{(2)}}.$$
$$\frac{\partial \ell}{\partial b^{(1)}} : 1 \times d_{\boldsymbol{s}^{(1)}}.$$

Where:

$$d_{\boldsymbol{s}^{(1)}} = d_{\boldsymbol{z}^{(1)}}.$$
$$d_{\boldsymbol{s}^{(2)}} = d_{\hat{\boldsymbol{y}}}.$$

## Gradient of $W^{(2)}$

Using chain rule and tensor notation:

$$\frac{\partial \ell}{\partial W_{ij}^{(2)}} = \sum_{k,l} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \frac{\partial s_l^{(2)}}{\partial W_{ij}^{(2)}}.$$
$$= \sum_{k,l} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \frac{\partial}{\partial W_{ij}^{(2)}} \left( \sum_m W_{lm}^{(2)} z_m^{(1)} + b_l^{(2)} \right).$$
$$= \sum_{k,l} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} z_m^{(1)} \delta_{il} \delta_{jm}.$$
$$= \sum_k \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_i^{(2)}} z_j^{(1)}.$$
$$= \delta_i^{(2)} z_j^{(1)}.$$

In matrix form:

$$\frac{\partial \ell}{\partial W^{(2)}} = \begin{pmatrix} \frac{\partial \ell}{\partial W_{00}^{(2)}} & \frac{\partial \ell}{\partial W_{10}^{(2)}} & \cdots & \frac{\partial L}{\partial W_{d_{s(2)}0}^{(2)}} \\ \frac{\partial \ell}{\partial W_{01}^{(2)}} & \frac{\partial \ell}{\partial W_{11}^{(2)}} & \cdots & \frac{\partial \ell}{\partial W_{d_{s(2)}1}^{(2)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \ell}{\partial W_{0d_{z(1)}}^{(2)}} & \frac{\partial \ell}{\partial W_{1d_{z(1)}}^{(2)}} & \cdots & \frac{\partial \ell}{\partial W_{d_{s(2)}d_{z(1)}}^{(2)}} \end{pmatrix}.$$

$$= \begin{pmatrix} z_0^{(1)} \\ \vdots \\ z_{d_{z(1)}}^{(1)} \end{pmatrix} \begin{pmatrix} \frac{\partial \ell}{\partial \hat{y}_0} & \cdots & \frac{\partial \ell}{\partial \hat{y}_{d_{\hat{y}}}} \end{pmatrix} \begin{pmatrix} \frac{\partial \hat{y}_0}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_0}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_0}{\partial s_{d_{s(2)}}^{(2)}} \\ \frac{\partial \hat{y}_1}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_1}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_1}{\partial s_{d_{s(2)}}^{(2)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_{d_{s(2)}}^{(2)}} \end{pmatrix} = z^{(1)} \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}}.$$

$$= \begin{pmatrix} z_0^{(1)} \\ \vdots \\ z_{d_{z(1)}}^{(1)} \end{pmatrix} \begin{pmatrix} \delta_0^{(2)} & \cdots & \delta_{d_{s(2)}}^{(2)} \end{pmatrix} = z^{(1)} [\delta^{(2)}]^T.$$

This results are for any activation function and any loss, in our case:

$$\frac{\partial \hat{y}_k}{\partial s_i^{(2)}} = \frac{\partial}{\partial s_i^{(2)}} g\left(s_k^{(2)}\right) = \delta_{ki}.$$

$$\frac{\partial \hat{y}}{\partial s^{(2)}} = I_{d_{\hat{y}} \times d_{s(2)}}.$$

And for the loss:

$$\frac{\partial \ell}{\partial \hat{y}_k} = \frac{\partial}{\partial \hat{y}_k} \left[ \sum_i (\hat{y}_i - y_i)^2 \right].$$
$$= \sum_i 2(\hat{y}_i - y_i)\delta_{ik}.$$
$$= 2(\hat{y}_k - y_k).$$

$$\frac{\partial L}{\partial \hat{y}} = 2(\hat{y} - y)^{\mathrm{T}}.$$

Inserting that in the formula for $\delta^{(2)}$:

$$\delta_i^{(2)} = 2(\hat{y}_i - y_i).$$

$$\delta^{(2)} = 2(\hat{y} - y) = 2 \begin{pmatrix} \hat{y}_0 - y_0 \\ \vdots \\ \hat{y}_{d_{s(2)}} - y_{d_{s(2)}} \end{pmatrix}.$$

And for $\frac{\partial \ell}{\partial W^{(2)}}$:

$$\frac{\partial \ell}{\partial W_{ij}^{(2)}} = 2(\hat{y}_i - y_i)z_j^{(1)}.$$

$$\frac{\partial \ell}{\partial W^{(2)}} = 2 \begin{pmatrix} z_0^{(1)} \\ \vdots \\ z_{d_{z^{(1)}}}^{(1)} \end{pmatrix} \begin{pmatrix} \hat{y}_0 - y_0 & \cdots & \hat{y}_{d_{s^{(2)}}} - y_{d_{s^{(2)}}} \end{pmatrix}.$$

$$= 2z^{(1)}(\hat{y} - y)^T$$

## Gradient of $b^{(2)}$

Using chain rule and components and having into account the previous results for $W^{(2)}$:

$$\frac{\partial \ell}{\partial b_i^{(2)}} = \sum_{k,l} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \frac{\partial s_l^{(2)}}{\partial b_i^{(2)}}.$$

$$= \sum_{k,l} \frac{\partial l}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \frac{\partial}{\partial b_i^{(2)}} \left( \sum_m W_{lm}^{(2)} z_m^{(1)} + b_l^{(2)} \right).$$

$$= \sum_{k,l} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \delta_{il}.$$

$$= \sum_k \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_i^{(2)}}.$$

$$= \delta_i^{(2)}.$$

$$= 2(\hat{y}_i - y_i).$$

In matrix form:

$$\frac{\partial L}{\partial b^{(2)}} = \begin{pmatrix} \frac{\partial \ell}{\partial b_0^{(2)}} & \frac{\partial \ell}{\partial b_1^{(2)}} & \cdots \end{pmatrix}$$

$$= \begin{pmatrix} \frac{\partial \ell}{\partial \hat{y}_0} & \cdots & \frac{\partial \ell}{\partial \hat{y}_{d_{\hat{y}}}} \end{pmatrix} \begin{pmatrix} \frac{\partial \hat{y}_0}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_0}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_0}{\partial s_{d_{s^{(2)}}}^{(2)}} \\ \frac{\partial \hat{y}_1}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_1}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_1}{\partial s_{d_{s^{(2)}}}^{(2)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_0^{(2)}} & \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_1^{(2)}} & \cdots & \frac{\partial \hat{y}_{d_{\hat{y}}}}{\partial s_{d_{s^{(2)}}}^{(2)}} \end{pmatrix}.$$

$$= \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}}.$$

$$= [\delta^{(2)}]^T.$$

$$= 2(\hat{y} - y)^T.$$

$$= 2 \begin{pmatrix} \hat{y}_0 - y_0 & \cdots & \hat{y}_{d_{s^{(2)}}} - y_{d_{s^{(2)}}} \end{pmatrix}.$$

## Gradient of $W^{(1)}$

Using chain rule and tensor notation:

$$\frac{\partial \ell}{\partial W_{ij}^{(1)}} = \sum_{k,l,m,n} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s_l^{(2)}} \frac{\partial s_l^{(2)}}{\partial z_m^{(1)}} \frac{\partial z_m^{(1)}}{\partial s_n^{(1)}} \frac{\partial s_n^{(1)}}{\partial W_{ij}^{(1)}}.$$

$$= \sum_{l,m,n} \delta_l^{(2)} \frac{\partial s_l^{(2)}}{\partial z_m^{(1)}} \frac{\partial z_m^{(1)}}{\partial s_n^{(1)}} \frac{\partial s_n^{(1)}}{\partial W_{ij}^{(1)}}.$$

$$= \sum_{n} \delta_n^{(1)} \frac{\partial s_n^{(1)}}{\partial W_{ij}^{(1)}}.$$

$$= \sum_{n} \delta_i^{(1)} x_j.$$

$$\frac{\partial \ell}{\partial W^{(1)}} = \boldsymbol{x} \left[ \boldsymbol{\delta^{(1)}} \right]^T.$$

Where $\boldsymbol{\delta^{(L=1)}}$ are the so called "errors" for the linear layer $L=1$. Then, we can compute $\frac{\partial \ell}{\partial W^{(1)}}$ in terms of the jacobians:

$$\frac{\partial \ell}{\partial W^{(1)}} = \boldsymbol{x} \left[ \boldsymbol{\delta^{(1)}} \right]^T.$$

$$\boldsymbol{\delta^{(1)}} = \left[ \frac{\partial \boldsymbol{s^{(2)}}}{\partial \boldsymbol{z^{(1)}}} \frac{\partial \boldsymbol{z^{(1)}}}{\partial \boldsymbol{s^{(1)}}} \right]^T \rightarrow \boldsymbol{x} \left[ \left[ \frac{\partial \boldsymbol{s^{(2)}}}{\partial \boldsymbol{z^{(1)}}} \frac{\partial \boldsymbol{z^{(1)}}}{\partial \boldsymbol{s^{(1)}}} \right]^T \boldsymbol{\delta^{(2)}} \right]^T.$$

$$\boldsymbol{\delta^{(2)}} = \left[ \frac{\partial \ell}{\partial \hat{\boldsymbol{y}}} \frac{\partial \hat{\boldsymbol{y}}}{\partial \boldsymbol{s^{(2)}}} \right]^T \rightarrow \boldsymbol{x} \frac{\partial \ell}{\partial \hat{\boldsymbol{y}}} \frac{\partial \hat{\boldsymbol{y}}}{\partial \boldsymbol{s^{(2)}}} \frac{\partial \boldsymbol{s^{(2)}}}{\partial \boldsymbol{z^{(1)}}} \frac{\partial \boldsymbol{z^{(1)}}}{\partial \boldsymbol{s^{(1)}}}.$$

## Gradient of $W^{(L)}$

The errors are easily generalizable:

$$\delta_i^{(L)} = \sum_{p,q} \delta_p^{(L+1)} \frac{\partial s_p^{(L+1)}}{\partial z_q^{(L)}} \frac{\partial z_q^{(L)}}{\partial s_i^{(L)}}.$$

In matrix form:

$$\boldsymbol{\delta^{(L)}} = \left[ \left( \delta_0^{(L+1)} \quad \cdots \quad \delta_{d_s(2)}^{(L+1)} \right) \begin{pmatrix} \frac{\partial s_0^{(L+1)}}{\partial z_0^{(L)}} & \frac{\partial s_0^{(L+1)}}{\partial z_1^{(L)}} & \cdots & \frac{\partial s_0^{(L+1)}}{\partial z_{d_z(L)}^{(L)}} \\ \frac{\partial s_1^{(L+1)}}{\partial z_0^{(L)}} & \frac{\partial s_1^{(L+1)}}{\partial z_1^{(L)}} & \cdots & \frac{\partial s_1^{(L+1)}}{\partial z_{d_z(L)}^{(L)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial s_{d_s(L+1)}^{(L+1)}}{\partial z_0^{(L)}} & \frac{\partial s_{d_s(L+1)}^{(L+1)}}{\partial z_1^{(L)}} & \cdots & \frac{\partial s_{d_s(L+1)}^{(L+1)}}{\partial z_{d_z(L)}^{(L)}} \end{pmatrix} \begin{pmatrix} \frac{\partial z_0^{(L)}}{\partial s_0^{(L)}} & \frac{\partial z_0^{(L)}}{\partial s_1^{(L)}} & \cdots & \frac{\partial z_0^{(L)}}{\partial s_{d_s(L)}^{(L)}} \\ \frac{\partial z_1^{(L)}}{\partial s_0^{(L)}} & \frac{\partial z_1^{(L)}}{\partial s_1^{(L)}} & \cdots & \frac{\partial z_1^{(L)}}{\partial s_{d_s(L)}^{(L)}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_{d_z(L)}^{(L)}}{\partial s_0^{(L)}} & \frac{\partial z_{d_z(L)}^{(L)}}{\partial s_1^{(L)}} & \cdots & \frac{\partial z_{d_z(L)}^{(L)}}{\partial s_{d_s(L)}^{(L)}} \end{pmatrix}^T \right].$$

$$= \left[ \left[ \boldsymbol{\delta^{(L+1)}} \right]^T \frac{\partial \boldsymbol{s^{(L+1)}}}{\partial \boldsymbol{z^{(L)}}} \frac{\partial \boldsymbol{z^{(L)}}}{\partial \boldsymbol{s^{(L)}}} \right]^T.$$

$$= \left[ \frac{\partial \boldsymbol{s^{(L+1)}}}{\partial \boldsymbol{z^{(L)}}} \frac{\partial \boldsymbol{z^{(L)}}}{\partial \boldsymbol{s^{(L)}}} \right]^T \boldsymbol{\delta^{(L+1)}}.$$

Now, let's compute $\frac{\partial \boldsymbol{s^{(L+1)}}}{\partial \boldsymbol{z^{(L)}}}$ for a linear layer:

$$\frac{\partial s_i^{(L+1)}}{\partial z_j^{(L)}} = \frac{\partial}{\partial z_j^{(L)}}\left(\sum_k W_{ik}^{(L+1)} z_k^{(L)} + b_i^{(L+1)}\right).$$

$$\frac{\partial s_i^{(L+1)}}{\partial z_j^{(L)}} = W_{ij}^{(L+1)}.$$

$$\frac{\partial s^{(L+1)}}{\partial z^{(L)}} = W^{(L+1)}.$$

Taking into account the previous expressions, we can compute the gradient for any linear layer and any activation function:

$$\frac{\partial \ell}{\partial W^{(L)}} = z^{(L-1)}\left[\delta^{(L)}\right]^T.$$

$$\delta^{(L)} = \left[W^{(L+1)}\frac{\partial z^{(L)}}{\partial s^{(L)}}\right]^T \delta^{(L+1)}.$$

$$z^{(0)} = x.$$

$$\delta^{(L_{\max})} = \left[\frac{\partial \ell}{\partial \hat{y}}\frac{\partial \hat{y}}{\partial s^{(L_{\max})}}\right]^T.$$

In regard to $\frac{\partial z^{(L)}}{\partial s^{(L)}}$, we can compute it for $g = I(\cdot)$ and $f = \text{ReLU}(\cdot)$ (one of the most common cases):

$$f = \text{ReLU}(\cdot) \rightarrow \frac{\partial z_i^{(L)}}{\partial s_j^{(L)}} = \max(0, \text{sign}(s_j^L))\delta_{ij}.$$

$$\frac{\partial z^{(L)}}{\partial s^{(L)}} = I_{z^{(L)} \times s^{(L)}}^{+s^{(L)}} = \begin{pmatrix} \max(0, \text{sign}(s_0^{(L)})) & 0 & \ldots & 0 \\ 0 & \max(0, \text{sign}(s_1^{(L)})) & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \max(0, \text{sign}(s_{d_{s^{(L)}}}^{(L)})) \end{pmatrix}$$

$$g = I(\cdot) \rightarrow \frac{\partial z_i^{(L)}}{\partial s_j^{(L)}} = \delta_{ij}.$$

$$\frac{\partial z^{(L)}}{\partial s^{(L)}} = I_{z^{(L)} \times s^{(L)}} = \begin{pmatrix} 1 & 0 & \ldots & 0 \\ 0 & 1 & \ldots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \ldots & 1 \end{pmatrix}.$$

Then, the "errors" for any linear layer are given by:

$$\delta_i^{(L)} = \sum_{p,q} \delta_p^{(L+1)} \frac{\partial z_q^{(L)}}{\partial s_i^{(L)}} \frac{\partial}{\partial z_q^{(L)}}\left(\sum_l W_{pl}^{(L+1)} z_l^{(L)} + b_p^{(L+1)}\right).$$

$$= \sum_{p,q} \delta_p^{(L+1)} W_{pq}^{(L+1)} \frac{\partial z_q^{(L)}}{\partial s_i^{(L)}}.$$

$$\delta^{(L)} = \left[W^{L+1}\frac{\partial z^{(L)}}{\partial s^{(L)}}\right]^T \delta^{(L+1)}.$$

For a $\text{ReLU}(\cdot)$:

$$\delta^{(L)} = \left[W^{L+1} I_{z^{(L)} \times s^{(L)}}^{+s^{(L)}}\right]^T \delta^{(L+1)}.$$

As an example, let's particularize our computation of $\frac{\partial \ell}{\partial W^{(1)}}$:

$$\frac{\partial L}{\partial W^{(1)}} = \boldsymbol{x}\left[\boldsymbol{\delta^{(1)}}\right]^T.$$

$$= \boldsymbol{x}\left[\left[W^{(2)} I^{+\boldsymbol{s}^{(1)}}_{\boldsymbol{z}^{(1)} \times \boldsymbol{s}^{(1)}}\right]^T \boldsymbol{\delta^{(2)}}\right]^T.$$

$$= 2\boldsymbol{x}(\hat{\boldsymbol{y}} - \boldsymbol{y})^T W^{(2)} I^{+\boldsymbol{s}^{(1)}}_{\boldsymbol{z}^{(1)} \times \boldsymbol{s}^{(1)}}.$$

## Gradient of $b^{(1)}$

Following the same idea:

$$\frac{\partial \ell}{\partial b^{(1)}_i} = \sum_{k,l,m,n} \frac{\partial \ell}{\partial \hat{y}_k} \frac{\partial \hat{y}_k}{\partial s^{(2)}_l} \frac{\partial s^{(2)}_l}{\partial z^{(1)}_m} \frac{\partial z^{(1)}_m}{\partial s^{(1)}_n} \frac{\partial s^{(1)}_n}{\partial b^{(1)}_i}.$$

$$= \sum_n \delta^{(1)}_n \frac{\partial s^{(1)}_n}{\partial b^{(1)}_i}.$$

$$= \delta^{(1)}_i.$$

$$\frac{\partial \ell}{\partial b^{(1)}} = \left[\boldsymbol{\delta^{(1)}}\right]^T.$$

In terms of the jacobians:

$$\frac{\partial \ell}{\partial b^{(1)}} = \frac{\partial \ell}{\partial \hat{\boldsymbol{y}}} \frac{\partial \hat{\boldsymbol{y}}}{\partial \boldsymbol{s}^{(2)}} \frac{\partial \boldsymbol{s}^{(2)}}{\partial \boldsymbol{z}^{(1)}} \frac{\partial \boldsymbol{z}^{(1)}}{\partial \boldsymbol{s}^{(1)}}.$$

## Gradient of $b^{(L)}$

For any linear layer, the gradient respect to the bias is:

$$\frac{\partial \ell}{\partial b^{(L)}_i} = \delta^{(L)}_i.$$

$$\frac{\partial \ell}{\partial b^{(L)}} = \left[\boldsymbol{\delta^{(L)}}\right]^T.$$

Where $\boldsymbol{\delta^{(L_{\max})}}$ is given in the previous section. Let's particularize for our special case:

$$\frac{\partial L}{\partial b^{(1)}} = 2(\hat{\boldsymbol{y}} - \boldsymbol{y})^T W^{(2)} I^{+\boldsymbol{s}^{(1)}}_{\boldsymbol{z}^{(1)} \times \boldsymbol{s}^{(1)}}.$$

## Summary

Shapes:

$$\boldsymbol{s}^{(L)} = 1 \times d_{\boldsymbol{s}^{(L)}}.$$

$$\boldsymbol{z}^{(L)} = 1 \times d_{\boldsymbol{z}^{(L)}}.$$

$$W^{(L)} : d_{\boldsymbol{s}^{(L)}} \times d_{\boldsymbol{z}^{(L-1)}}.$$

$$b^{(L)} : d_{\boldsymbol{s}^{(L)}} \times 1.$$

$$\frac{\partial \ell}{\partial W^{(L)}} : d_{\boldsymbol{z}^{(L-1)}} \times d_{\boldsymbol{s}^{(L)}}.$$

$$\frac{\partial \ell}{\partial b^{(L)}} : 1 \times d_{\boldsymbol{s}^{(L)}}.$$

Where:

$$d_{z^{(0)}} = d_x.$$

$$d_{z^{(L_{\max})}} = d_{\hat{y}} = d_y.$$

Backpropagation for a stack of linear layers in matrix form:

$$\frac{\partial \ell}{\partial W^{(L)}} = z^{(L-1)} \left[\boldsymbol{\delta}^{(L)}\right]^T.$$

$$\boldsymbol{\delta}^{(L)} = \left[W^{(L+1)} \frac{\partial z^{(L)}}{\partial s^{(L)}}\right]^T \boldsymbol{\delta}^{(L+1)}.$$

$$\frac{\partial \ell}{\partial b^{(L)}} = \left[\boldsymbol{\delta}^{(L)}\right]^T.$$

$$z^{(0)} = x.$$

$$\boldsymbol{\delta}^{(L_{\max})} = \left[\frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(L_{\max})}}\right]^T.$$

$\frac{\partial z^{(L)}}{\partial s^{(L)}}$ for a $\mathrm{ReLU}(\cdot)$:

$$\frac{\partial z^{(L)}}{\partial s^{(L)}} = I_{z^{(L)} \times s^{(L)}}^{+s^{(L)}} = \begin{pmatrix} \max(0, \mathrm{sign}(s_0^L)) & 0 & \ldots & 0 \\ 0 & \max(0, \mathrm{sign}(s_1^L)) & \ldots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \ldots & \max(0, \mathrm{sign}(s_{d_{s^{(L)}}}^L)) \end{pmatrix}.$$

| Parameter | Gradient | Gradient shape |
|---|---|---|
| $W^{(1)}$ | $x \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}} \frac{\partial s^{(2)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial s^{(1)}} = 2x(\hat{y}-y)^T W^{(2)} I_{z^{(1)} \times s^{(1)}}^{+s^{(1)}}.$ | $d_x \times d_{s^{(1)}}.$ |
| $b^{(1)}$ | $\frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}} \frac{\partial s^{(2)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial s^{(1)}} = 2(\hat{y}-y)^T W^{(2)} I_{z^{(1)} \times s^{(1)}}^{+s^{(1)}}.$ | $1 \times d_{s^{(1)}}.$ |
| $W^{(2)}$ | $z^{(1)} \frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}} = 2z^{(1)}(\hat{y}-y)^T.$ | $d_{z^{(1)}} \times d_{s^{(2)}}.$ |
| $b^{(2)}$ | $\frac{\partial \ell}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial s^{(2)}} = 2(\hat{y}-y)^T.$ | $1 \times d_{s^{(2)}}.$ |

## Solution d)

With the change in notation:

$$\frac{\partial z^{(2)}}{\partial z^{(1)}} \to \frac{\partial z^{(1)}}{\partial s^{(1)}}.$$

$$\frac{\partial \hat{y}}{\partial z^{(3)}} \to \frac{\partial \hat{y}}{\partial s^{(2)}}.$$

$\frac{\partial z^{(1)}}{\partial s^{(1)}}$:

$$f = \mathrm{ReLU}(\cdot) \to \frac{\partial z_i^{(1)}}{\partial s_j^{(1)}} = \frac{\partial}{\partial s_j^{(1)}} \mathrm{ReLU}(s_i^{(1)}) = \max(0, \mathrm{sign}(s_j^{(1)})) \delta_{ij}.$$

$$\frac{\partial z^{(1)}}{\partial s^{(1)}} = I_{d_{z^{(1)}} \times s^{(1)}}^{+s^{(1)}} = \begin{pmatrix} \max(0, \mathrm{sign}(s_0^{(1)})) & 0 & \ldots & 0 \\ 0 & \max(0, \mathrm{sign}(s_1^{(1)})) & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & \max(0, \mathrm{sign}(s_{d_{s^{(1)}}}^{(1)})) \end{pmatrix}.$$

$\frac{\partial \hat{y}}{\partial s^{(2)}}$:

$$g = I(\cdot) \rightarrow \frac{\partial \hat{y}_i}{\partial s_j^{(2)}} = \frac{\partial}{\partial s_j^{(2)}} I(s_i^{(2)}) = \delta_{ij}, \quad i = 0, \ldots, d_{\hat{y}}, \; j = 0, \ldots, d_{s^{(2)}}.$$

$$\frac{\partial \hat{y}}{\partial s^{(2)}} = I_{d_{\hat{y}} \times s^{(2)}} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{pmatrix}.$$

$\frac{\partial \ell}{\partial \hat{y}}$ :

$$\frac{\partial \ell}{\partial \hat{y}_i} = \frac{\partial}{\partial \hat{y}_i} \left[ \sum_j (\hat{y}_j - y_j)^2 \right].$$
$$= \sum_j 2(\hat{y}_j - y_j) \delta_{ij}.$$
$$= 2(\hat{y}_i - y_i).$$

$$\frac{\partial \ell}{\partial \hat{y}} = 2(\hat{y} - y)^{\mathrm{T}}.$$
$$= 2 \begin{pmatrix} \hat{y}_0 - y_0 & \cdots & \hat{y}_{d_{\hat{y}}} - y_{d_y} \end{pmatrix}$$

# Solution 1.3

### Solution a)

In the case of **b)** the loss function (the replacement is done in the table with intermediate variables only) :

| Layer | Input | Output |
|---|---|---|
| Linear$_1$ | $x$ | $s^{(1)} = W^{(1)} x + b^{(1)}$ |
| $\sigma$ | $s^{(1)}$ | $z^{(1)} = \sigma(s^{(1)})$ |
| Linear$_2$ | $z^{(1)}$ | $s^{(2)} = W^{(2)} z^{(1)} + b^{(2)}$ |
| $\sigma$ | $s^{(2)}$ | $\hat{y} = \sigma(s^{(2)})$ |
| Loss | $\hat{y}, y$ | $\ell_{\mathrm{MSE}} = (\hat{y} - y)(\hat{y} - y)^T$ |

In the case of **c)** the jacobians $\frac{\partial z^{(1)}}{\partial s^{(1)}}$ and $\frac{\partial \hat{y}}{\partial s^{(2)}}$.

In the case of **d)** , we need to compute the derivatives so we can see the components explicitly, the derivative of $\sigma$ is:

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

Then, $\frac{\partial z^{(1)}}{\partial s^{(1)}}$ :

$$f = \sigma(\cdot) \rightarrow \frac{\partial z_i^{(1)}}{\partial s_j^{(1)}} = \frac{\partial}{\partial s_j^{(1)}} \sigma(s_i^{(1)}) = \sigma(s_i^{(1)})(1 - \sigma(s_i^{(1)})) \delta_{ij}.$$

$$\frac{\partial z^{(1)}}{\partial s^{(1)}} = \begin{pmatrix} \sigma(s_0^{(1)})(1 - \sigma(s_0^{(1)})) & 0 & \cdots & 0 \\ 0 & \sigma(s_1^{(1)})(1 - \sigma(s_1^{(1)})) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma(s_{d_{s^{(1)}}}^{(1)})(1 - \sigma(s_{d_{s^{(1)}}}^{(1)})) \end{pmatrix}.$$

$\frac{\partial \hat{y}}{\partial s^{(2)}}$:

$$g = \sigma(\cdot) \rightarrow \frac{\partial \hat{y}_i}{\partial s_j^{(2)}} = \frac{\partial}{\partial s_j^{(2)}} \sigma(s_i^{(2)}) = \sigma(s_i^{(2)})(1 - \sigma(s_i^{(2)}))\delta_{ij}.$$

$$\frac{\partial \hat{y}}{\partial s^{(2)}} = \begin{pmatrix} \sigma(s_0^{(2)})(1 - \sigma(s_0^{(2)})) & 0 & \cdots & 0 \\ 0 & \sigma(s_1^{(2)})(1 - \sigma(s_1^{(2)})) & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma(s_{d_{s^{(2)}}}^{(2)})(1 - \sigma(s_{d_{s^{(2)}}}^{(2)})) \end{pmatrix}.$$

$\frac{\partial \ell}{\partial \hat{y}}$ remains the same.

## Solution b)

In the equations of **b)** only the loss function, $\ell_{\text{MSE}} \rightarrow \ell_{\text{BCE}}$

| Layer | Input | Output |
|---|---|---|
| Linear$_1$ | $x$ | $s^{(1)} = W^{(1)}x + b^{(1)}$ |
| $\sigma$ | $s^{(1)}$ | $z^{(1)} = \sigma(s^{(1)})$ |
| Linear$_2$ | $z^{(1)}$ | $s^{(2)} = W^{(2)}z^{(1)} + b^{(2)}$ |
| $\sigma$ | $s^{(2)}$ | $\hat{y} = \sigma(s^{(2)})$ |
| Loss | $\hat{y}, y$ | $\ell_{\text{BCE}} = -\frac{1}{K}\left[y^T \log(\hat{y}) + (1 - y)^T \log(1 - \hat{y})\right]$ |

In the equations of **c)** the derivative $\frac{\partial \ell}{\partial \hat{y}}$.

In the equations of **d)**, since the derivative $\frac{\partial \ell}{\partial \hat{y}}$ changes, so do its components, let's compute them and write the matrix representation:

$$\ell_{\text{BCE}} = -\frac{1}{K} \sum_j \left[y_j \log(\hat{y}_j) + (1 - y_j) \log(1 - \hat{y}_j)\right].$$

$$\frac{\partial \ell_{\text{BCE}}}{\partial \hat{y}_i} = \frac{1}{K} \frac{\hat{y}_i - y_i}{\hat{y}_i(1 - \hat{y}_i)}.$$

$$\frac{\partial \ell_{\text{BCE}}}{\partial \hat{y}} = \frac{1}{K}\left(\begin{matrix} \frac{\hat{y}_0 - y_0}{\hat{y}_0(1 - \hat{y}_0)} & \frac{\hat{y}_1 - y_1}{\hat{y}_1(1 - \hat{y}_1)} & \cdots & \frac{\hat{y}_{d_{\hat{y}}} - y_{d_{\hat{y}}}}{\hat{y}_{d_{\hat{y}}}(1 - \hat{y}_{d_{\hat{y}}})} \end{matrix}\right).$$

## Solution c)

Because the the calculation and the calculation of the gradient is faster and $\text{ReLU}$ is good avoiding gradient vanishing.