# A Numpy Neural Network

Saffet Gökçen Şen

May 31, 2020

1

A neural network with 3 hidden layers is to be constructed. The gradient of the loss function with respect to the weights and biases of the hidden units are derived. In a separate Jupyter notebook, the forward propagation and backward propagation are implemented using numpy to fit the network to a dataset. The output data and the input data are created randomly. The activation functions of the hidden units are all relu except the last hidden layer.

The dimension of the input layer is $d_{in}$. It means that there are $d_{in}$ features. The dimensions of the first, second and third layers are $d_1$, $d_2$ and $d_{out}$ respectively. Hence, there are $d_{out}$ outputs. Let the batch size be $n$. Then, the input to the first hidden layer is a matrix $\mathbf{x}$ of dimension $n \times d_{in}$. The weight matrix of the first hidden layer is $\mathbf{w}_1$ of dimension $d_{in} \times d_1$. The bias matrix of the first hidden layer is $\mathbf{b}_1$ of dimension $1 \times d_1$. The output of the first hidden layer is $\mathbf{h}_1$ given by

$$\mathbf{h}_1 = \mathbf{x}.\mathbf{w}_1 + \mathbf{1}_{n \times 1}.\mathbf{b}_1 \tag{1}$$

$$\mathbf{h}_{1_{\text{relu}}} = \text{Relu}\,(\mathbf{h}_1) \tag{2}$$

$\mathbf{1}_{n \times 1}$ is a matrix of ones with dimension n $\times$ 1. Let $\mathbf{1}_{n \times 1}.\mathbf{b}_1$ be denoted by $\mathbf{B}_1$. $\mathbf{B}_1$, $\mathbf{h}_1$ and $\mathbf{h}_{1_{\text{relu}}}$ are of dimension $n \times d_1$. The weight matrix of the second hidden layer is $\mathbf{w}_2$ of dimension $d_1 \times d_2$. The bias matrix of the second hidden layer is $\mathbf{b}_2$ of dimension $1 \times d_2$. The output of the second hidden layer is $\mathbf{h}_2$ given by

$$\mathbf{h}_2 = \mathbf{h}_{1_{\text{relu}}}.\mathbf{w}_2 + \mathbf{1}_{n \times 1}.\mathbf{b}_2 \tag{3}$$

$$\mathbf{h}_{2_{\text{relu}}} = \text{Relu}\,(\mathbf{h}_2) \tag{4}$$

Let $\mathbf{1}_{n \times 1}.\mathbf{b}_2$ be denoted by $\mathbf{B}_2$. $\mathbf{B}_2$, $\mathbf{h}_2$ and $\mathbf{h}_{2_{\text{relu}}}$ are of dimension $n \times d_2$. The weight matrix of the third hidden layer is $\mathbf{w}_3$ of dimension $d_2 \times d_{\text{out}}$. The bias matrix of the third hidden layer is $\mathbf{b}_3$ of dimension $1 \times d_{\text{out}}$. The output of the

third hidden layer is $y_{\text{pred}}$ given by

$$\mathbf{y}_{\text{pred}} = \mathbf{h}_{2_{\text{relu}}}.\mathbf{w}_3 + \mathbf{1}_{\text{n}\times1}.\mathbf{b}_3 \tag{5}$$

Let $\mathbf{1}_{\text{n}\times1}.\mathbf{b}_3$ be denoted by $\mathbf{B}_3$. $\mathbf{B}_3$ and $\mathbf{y}_{\text{pred}}$ are of dimension $n \times d_{\text{out}}$. The loss function is the mean squared error function. Hence,

$$\text{L} = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \left( y_{\text{pred}_{ij}} - y_{ij} \right)^2 \tag{6}$$

The derivative of the loss function L with respect to the matrix $\mathbf{w}_3$ is a matrix of dimension $1 \times (\text{d}_2.\text{d}_{\text{out}})$. The element in the column k.l of this matrix is given by

$$\left( \frac{\partial \text{L}}{\partial \mathbf{w}_3} \right)_{1(\text{k.l})} = \frac{\partial L}{\partial w_{3_{\text{kl}}}} = \frac{\partial}{\partial w_{3_{\text{kl}}}} \left[ \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \left( y_{\text{pred}_{ij}} - y_{ij} \right)^2 \right] = \tag{7}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \frac{\partial}{\partial w_{3_{\text{kl}}}} \left[ \left( y_{\text{pred}_{ij}} - y_{ij} \right)^2 \right] = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2 \left( y_{\text{pred}_{ij}} - y_{ij} \right) \frac{\partial y_{\text{pred}_{ij}}}{\partial w_{3_{\text{kl}}}} \tag{8}$$

$y_{\text{pred}_{ij}}$ is equal to the following:

$$y_{\text{pred}_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{\text{mj}}} + B_{3_{ij}} \Rightarrow \tag{9}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{3_{\text{kl}}}} = h_{2_{\text{relu}_{im}}} I\left(\text{m=k \& j=l}\right) \tag{10}$$

Let the derivation for $\left( \frac{\partial \text{L}}{\partial \mathbf{w}_3} \right)_{1(\text{k.l})}$ be completed:

$$\left( \frac{\partial \text{L}}{\partial \mathbf{w}_3} \right)_{1(\text{k.l})} = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2 \left( y_{\text{pred}_{ij}} - y_{ij} \right) h_{2_{\text{relu}_{im}}} I\left(\text{m=k \& j=l}\right) \Rightarrow \tag{11}$$

3

$$\left(\frac{\partial L}{\partial \mathbf{w}_3}\right)_{1(k.l)} = \sum_{i=1}^{n} 2\left(y_{\text{pred}_{il}} - y_{il}\right) h_{2_{\text{relu}_{ik}}} \Rightarrow \tag{12}$$

$$\left(\frac{\partial L}{\partial \mathbf{w}_3}\right)_{1(k.l)} = \left[\mathbf{h}_{2_{\text{relu}}}^{T}.2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right)\right]_{kl} \tag{13}$$

If $\frac{\partial L}{\partial \mathbf{w}_3}$ is put into the form of a matrix of dimension $d_2 \times d_{\text{out}}$, then it's equal to $\mathbf{h}_{2_{\text{relu}}}^{T}.2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right)$.

Now, let the derivative of the loss with respect to the $\mathbf{b}_3$ matrix be derived. It is a matrix of dimension $1 \times d_{\text{out}}$. The element in its $k^{\text{th}}$ column is given as follows:

$$\left(\frac{\partial L}{\partial \mathbf{b}_3}\right)_{1k} = \frac{\partial L}{\partial b_{3_k}} = \frac{\partial}{\partial b_{3_k}}\left[\sum_{i=1}^{n}\sum_{j=1}^{d_{\text{out}}}\left(y_{\text{pred}_{ij}} - y_{ij}\right)^2\right] = \tag{14}$$

$$\sum_{i=1}^{n}\sum_{j=1}^{d_{\text{out}}}\frac{\partial}{\partial b_{3_k}}\left[\left(y_{\text{pred}_{ij}} - y_{ij}\right)^2\right] = \sum_{i=1}^{n}\sum_{j=1}^{d_{\text{out}}}2\left(y_{\text{pred}_{ij}} - y_{ij}\right)\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{3_k}} \tag{15}$$

$$y_{\text{pred}_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{mj}} + B_{3_{ij}} \Rightarrow \frac{\partial y_{\text{pred}_{ij}}}{\partial b_{3_k}} = \frac{\partial B_{3_{ij}}}{\partial b_{3_k}} = \frac{\partial b_{3_j}}{\partial b_{3_k}} \Rightarrow \tag{16}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{3_k}} = I\,(\text{j=k}) \tag{17}$$

Hence:

$$\left(\frac{\partial L}{\partial \mathbf{b}_3}\right)_{1k} = \sum_{i=1}^{n}\sum_{j=1}^{d_{\text{out}}}2\left(y_{\text{pred}_{ij}} - y_{ij}\right)\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{3_k}} \Rightarrow \tag{18}$$

$$\left(\frac{\partial L}{\partial \mathbf{b}_3}\right)_{1k} = \sum_{i=1}^{n}\sum_{j=1}^{d_{\text{out}}}2\left(y_{\text{pred}_{ij}} - y_{ij}\right)I\,(\text{j=k}) \Rightarrow \tag{19}$$

$$\left(\frac{\partial L}{\partial \mathbf{b}_3}\right)_{1k} = \sum_{i=1}^{n}2\left(y_{\text{pred}_{ik}} - y_{ik}\right) \tag{20}$$

The $k^{\text{th}}$ column of $\frac{\partial L}{\partial \mathbf{b}_3}$ is equal to the sum of the elements of the $k^{\text{th}}$ column of $2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right)$.

The derivative of L with respect to the matrix $\mathbf{w}_2$ is to be derived. It is a

matrix of dimension $1 \times (d_1.d_2)$. The element in the column k.l of this matrix is given by

$$\left(\frac{\partial L}{\partial \mathbf{w}_2}\right)_{1(k.l)} = \frac{\partial L}{\partial w_{2_{kl}}} = \frac{\partial}{\partial w_{2_{kl}}} \left[\sum_{i=1}^{n}\sum_{j=1}^{d_{out}} \left(y_{\text{pred}_{ij}} - y_{ij}\right)^2\right] = \tag{21}$$

$$\sum_{i=1}^{n}\sum_{j=1}^{d_{out}} \frac{\partial}{\partial w_{2_{kl}}} \left[\left(y_{\text{pred}_{ij}} - y_{ij}\right)^2\right] = \sum_{i=1}^{n}\sum_{j=1}^{d_{out}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right)\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{2_{kl}}} \tag{22}$$

$$y_{\text{pred}_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{mj}} + B_{3_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{mj}} + B_{3_{ij}} \Rightarrow \tag{23}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{2_{kl}}} = \sum_{m=1}^{d_2} \frac{\partial h_{2_{\text{relu}_{im}}}}{\partial w_{2_{kl}}} w_{3_{mj}} = \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial w_{2_{kl}}} w_{3_{mj}} = \tag{24}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial h_{2_{im}}}{\partial w_{2_{kl}}} w_{3_{mj}} = \tag{25}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial}{\partial w_{2_{kl}}} \left(\sum_{p=1}^{d_1} h_{1_{\text{relu}_{ip}}} w_{2_{pm}} + B_{2_{im}}\right) w_{3_{mj}} = \tag{26}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} h_{1_{\text{relu}_{ip}}} I\left(\text{p=k \& m=l}\right) w_{3_{mj}} \Rightarrow \tag{27}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{2_{kl}}} = \frac{\partial \text{Relu}\left(h_{2_{il}}\right)}{\partial h_{2_{il}}} h_{1_{\text{relu}_{ik}}} w_{3_{lj}} \tag{28}$$

Going on with the derivation:

$$\left(\frac{\partial L}{\partial \mathbf{w}_2}\right)_{1(k.l)} = \sum_{i=1}^{n}\sum_{j=1}^{d_{out}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right)\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{2_{kl}}} = \tag{29}$$

$$\sum_{i=1}^{n}\sum_{j=1}^{d_{out}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right)\frac{\partial \text{Relu}\left(h_{2_{il}}\right)}{\partial h_{2_{il}}} h_{1_{\text{relu}_{ik}}} w_{3_{lj}} = \tag{30}$$

$$\sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{2_{il}}\right)}{\partial h_{2_{il}}} h_{1_{\text{relu}_{ik}}} \left( \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) w_{3_{lj}} \right) \Rightarrow \tag{31}$$

$$\left( \frac{\partial \text{L}}{\partial \mathbf{w}_2} \right)_{1(k.l)} = \sum_{i=1}^{n} h_{1_{\text{relu}_{ik}}} \frac{\partial \text{Relu}\left(h_{2_{il}}\right)}{\partial h_{2_{il}}} \left[ 2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T \right]_{il} \tag{32}$$

Let the following definition be made for the matrix $\mathbf{A}$:

$$A_{il} = \frac{\partial \text{Relu}\left(h_{2_{il}}\right)}{\partial h_{2_{il}}} = \begin{cases} 1 & \text{if } h_{2_{il}} \geq 0 \\[2mm] 0 & \text{if } h_{2_{il}} < 0 \end{cases} \tag{33}$$

If $\frac{\partial \text{L}}{\partial \mathbf{w}_2}$ is put into a matrix of dimension $d_1 \times d_2$, then it is equal to

$$\mathbf{h}_{1_{\text{relu}}}^T . \left( \mathbf{A} \circ \left( 2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T \right) \right) \tag{34}$$

where $\circ$ denotes the Hadamard or Schur product of two matrices.

The derivative of L with respect to the bias matrix $\mathbf{b}_2$ is going to be obtained. It is a matrix of $1 \times d_2$. The element in its $k^{\text{th}}$ column is given as follows:

$$\left( \frac{\partial \text{L}}{\partial \mathbf{b}_2} \right)_{1k} = \frac{\partial L}{\partial b_{2_k}} = \frac{\partial}{\partial b_{2_k}} \left[ \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \left(y_{\text{pred}_{ij}} - y_{ij}\right)^2 \right] = \tag{35}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \frac{\partial}{\partial b_{2_k}} \left[ \left(y_{\text{pred}_{ij}} - y_{ij}\right)^2 \right] = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) \frac{\partial y_{\text{pred}_{ij}}}{\partial b_{2_k}} \tag{36}$$

$$y_{\text{pred}_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{mj}} + B_{3_{ij}} \Rightarrow \tag{37}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{2_k}} = \sum_{m=1}^{d_2} \frac{\partial h_{2_{\text{relu}_{im}}}}{\partial b_{2_k}} w_{3_{mj}} = \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial b_{2_k}} w_{3_{mj}} = \tag{38}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial h_{2_{im}}}{\partial b_{2_k}} w_{3_{mj}} = \tag{39}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} \frac{\partial}{\partial b_{2_k}} \left( \sum_{p=1}^{d_1} h_{1_{\text{relu}_{ip}}} w_{2_{pm}} + B_{2_{\text{im}}} \right) w_{3_{mj}} = \tag{40}$$

$$\sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} \frac{\partial b_{2_m}}{\partial b_{2_k}} w_{3_{mj}} = \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} I\left(\text{m=k}\right) w_{3_{mj}} \Rightarrow \tag{41}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{2_k}} = \frac{\partial \text{Relu}\left(h_{2_{\text{ik}}}\right)}{\partial h_{2_{\text{ik}}}} w_{3_{kj}} \tag{42}$$

Going on with the derivation:

$$\left( \frac{\partial \text{L}}{\partial \mathbf{b}_2} \right)_{1k} = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left( y_{\text{pred}_{ij}} - y_{ij} \right) \frac{\partial y_{\text{pred}_{ij}}}{\partial b_{2_k}} = \tag{43}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left( y_{\text{pred}_{ij}} - y_{ij} \right) \frac{\partial \text{Relu}\left(h_{2_{\text{ik}}}\right)}{\partial h_{2_{\text{ik}}}} w_{3_{kj}} = \tag{44}$$

$$\sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{2_{\text{ik}}}\right)}{\partial h_{2_{\text{ik}}}} \left( \sum_{j=1}^{d_{\text{out}}} 2\left( y_{\text{pred}_{ij}} - y_{ij} \right) w_{3_{kj}} \right) \Rightarrow \tag{45}$$

$$\left( \frac{\partial \text{L}}{\partial \mathbf{b}_2} \right)_{1k} = \sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{2_{\text{ik}}}\right)}{\partial h_{2_{\text{ik}}}} \left[ 2\left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right]_{ik} \tag{46}$$

The $k^{\text{th}}$ column of $\frac{\partial \text{L}}{\partial \mathbf{b}_2}$ is equal to the sum of the elements in the $k^{\text{th}}$ column of $\mathbf{A} \circ \left[ 2\left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right]$.

The derivative of the loss L is to be derived with respect to the weight matrix $\mathbf{w}_1$. This derivative will be of dimension $1 \times \left( d_{\text{in}} . d_1 \right)$.

$$\left( \frac{\partial \text{L}}{\partial \mathbf{w}_1} \right)_{1(k.l)} = \frac{\partial L}{\partial w_{1_{kl}}} = \frac{\partial}{\partial w_{1_{kl}}} \left[ \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \left( y_{\text{pred}_{ij}} - y_{ij} \right)^2 \right] = \tag{47}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \frac{\partial}{\partial w_{1_{kl}}} \left[ \left( y_{\text{pred}_{ij}} - y_{ij} \right)^2 \right] = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left( y_{\text{pred}_{ij}} - y_{ij} \right) \frac{\partial y_{\text{pred}_{ij}}}{\partial w_{1_{kl}}} \tag{48}$$

$$y_{\text{pred}_{ij}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{im}}} w_{3_{mj}} + B_{3_{ij}} = \sum_{m=1}^{d_2} \text{Relu}\left(h_{2_{im}}\right) w_{3_{mj}} + B_{3_{ij}} \Rightarrow \tag{49}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{1_{kl}}} = \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial h_{2_{im}}}{\partial w_{1_{kl}}} w_{3_{mj}} = \tag{50}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial}{\partial w_{1_{kl}}} \left(\sum_{p=1}^{d_1} h_{1_{\text{relu}_{ip}}} w_{2_{pm}} + B_{2_{im}}\right) = \tag{51}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} \frac{\partial h_{1_{ip}}}{\partial w_{1_{kl}}} w_{2_{pm}} = \tag{52}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} \frac{\partial}{\partial w_{1_{kl}}} \left(\sum_{q=1}^{d_{\text{in}}} x_{iq} w_{1_{qp}} + B_{1_{ip}}\right) w_{2_{pm}} = \tag{53}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} x_{iq} I\left(\text{q=k \& p=l}\right) w_{2_{pm}} \Rightarrow \tag{54}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial w_{1_{kl}}} = \sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{il}}\right)}{\partial h_{1_{il}}} x_{ik} w_{2_{lm}} \tag{55}$$

Going on with the derivation of $\frac{\partial L}{\partial \mathbf{w}_1}$:

$$\left(\frac{\partial \text{L}}{\partial \mathbf{w}_1}\right)_{1(k.l)} = \frac{\partial L}{\partial w_{1_{kl}}} = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) \frac{\partial y_{\text{pred}_{ij}}}{\partial w_{1_{kl}}} = \tag{56}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) \left(\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{il}}\right)}{\partial h_{1_{il}}} x_{ik} w_{2_{lm}}\right) = \tag{57}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \sum_{m=1}^{d_2} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{il}}\right)}{\partial h_{1_{il}}} x_{ik} w_{2_{lm}} = \tag{58}$$

$$\sum_{i=1}^{n} \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{il}}\right)}{\partial h_{1_{il}}} x_{ik} w_{2_{lm}} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) w_{3_{mj}} = \tag{59}$$

$$\sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{1_{\text{il}}}\right)}{\partial h_{1_{\text{il}}}} x_{\text{ik}} \sum_{m=1}^{d_2} w_{2_{\text{lm}}} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)_{\text{im}} = \tag{60}$$

$$\sum_{i=1}^{n} x_{\text{ik}} \frac{\partial \text{Relu}\left(h_{1_{\text{il}}}\right)}{\partial h_{1_{\text{il}}}} \left(\left(\mathbf{A} \circ \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)\right).\mathbf{w}_2^T\right)_{\text{il}} = \tag{61}$$

$$\left(\frac{\partial \text{L}}{\partial \mathbf{w}_1}\right)_{1(\text{k.l})} = \left(\mathbf{x}^T.\left(\mathbf{C} \circ \left(\left(\mathbf{A} \circ \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)\right).\mathbf{w}_2^T\right)\right)\right)_{\text{kl}} \tag{62}$$

where the matrix $\mathbf{C}$ is defined as:

$$C_{\text{il}} = \frac{\partial \text{Relu}\left(h_{1_{\text{il}}}\right)}{\partial h_{1_{\text{il}}}} = \begin{cases} 1 & \text{if } h_{1_{\text{il}}} \geq 0 \\ 0 & \text{if } h_{1_{\text{il}}} < 0 \end{cases} \tag{63}$$

If $\frac{\partial \text{L}}{\partial \mathbf{w}_1}$ is put into the form of a matrix with dimension $d_{\text{in}} \times d_1$, then this matrix is given by $\mathbf{x}^T.\left(\mathbf{C} \circ \left(\left(\mathbf{A} \circ \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)\right).\mathbf{w}_2^T\right)\right)$.

It is now the turn of the derivative of L with respect to the bias matrix $\mathbf{b}_1$. This derivative is a matrix with dimension $1 \times d_1$:

$$\left(\frac{\partial \text{L}}{\partial \mathbf{b}_1}\right)_{1\text{k}} = \frac{\partial L}{\partial b_{1_{\text{k}}}} = \frac{\partial}{\partial b_{1_{\text{k}}}} \left[\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \left(y_{\text{pred}_{\text{ij}}} - y_{\text{ij}}\right)^2\right] = \tag{64}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \frac{\partial}{\partial b_{1_{\text{k}}}} \left[\left(y_{\text{pred}_{\text{ij}}} - y_{\text{ij}}\right)^2\right] = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{\text{ij}}} - y_{\text{ij}}\right) \frac{\partial y_{\text{pred}_{\text{ij}}}}{\partial b_{1_{\text{k}}}} \tag{65}$$

$$y_{\text{pred}_{\text{ij}}} = \sum_{m=1}^{d_2} h_{2_{\text{relu}_{\text{im}}}} w_{3_{\text{mj}}} + B_{3_{\text{ij}}} = \sum_{m=1}^{d_2} \text{Relu}\left(h_{2_{\text{im}}}\right) w_{3_{\text{mj}}} + B_{3_{\text{ij}}} \Rightarrow \tag{66}$$

$$\frac{\partial y_{\text{pred}_{\text{ij}}}}{\partial b_{1_{\text{kl}}}} = \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} \frac{\partial h_{2_{\text{im}}}}{\partial b_{1_{\text{k}}}} w_{3_{\text{mj}}} = \tag{67}$$

$$\sum_{m=1}^{d_2} w_{3_{\text{mj}}} \frac{\partial \text{Relu}\left(h_{2_{\text{im}}}\right)}{\partial h_{2_{\text{im}}}} \frac{\partial}{\partial b_{1_{\text{k}}}} \left(\sum_{p=1}^{d_1} h_{1_{\text{relu}_{\text{ip}}}} w_{2_{\text{pm}}} + B_{2_{\text{im}}}\right) = \tag{68}$$

9

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} \frac{\partial h_{1_{ip}}}{\partial b_{1_k}} w_{2_{pm}} = \tag{69}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} \frac{\partial}{\partial b_{1_k}} \left(\sum_{q=1}^{d_{in}} x_{iq} w_{1_{qp}} + B_{1_{ip}}\right) w_{2_{pm}} = \tag{70}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} \frac{\partial b_{1_p}}{\partial b_{1_k}} w_{2_{pm}} = \tag{71}$$

$$\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \sum_{p=1}^{d_1} \frac{\partial \text{Relu}\left(h_{1_{ip}}\right)}{\partial h_{1_{ip}}} I\left(\text{p=k}\right) w_{2_{pm}} \Rightarrow \tag{72}$$

$$\frac{\partial y_{\text{pred}_{ij}}}{\partial b_{1_k}} = \sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} w_{2_{km}} \tag{73}$$

Going on with the derivation of $\frac{\partial L}{\partial \mathbf{b}_1}$:

$$\left(\frac{\partial \text{L}}{\partial \mathbf{b}_1}\right)_{1k} = \frac{\partial L}{\partial b_{1_k}} = \sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) \frac{\partial y_{\text{pred}_{ij}}}{\partial b_{1_k}} = \tag{74}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) \left(\sum_{m=1}^{d_2} w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} w_{2_{km}}\right) = \tag{75}$$

$$\sum_{i=1}^{n} \sum_{j=1}^{d_{\text{out}}} \sum_{m=1}^{d_2} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) w_{3_{mj}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} w_{2_{km}} = \tag{76}$$

$$\sum_{i=1}^{n} \sum_{m=1}^{d_2} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} w_{2_{km}} \sum_{j=1}^{d_{\text{out}}} 2\left(y_{\text{pred}_{ij}} - y_{ij}\right) w_{3_{mj}} = \tag{77}$$

$$\sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} \sum_{m=1}^{d_2} w_{2_{km}} \frac{\partial \text{Relu}\left(h_{2_{im}}\right)}{\partial h_{2_{im}}} \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)_{im} \Rightarrow \tag{78}$$

$$\left(\frac{\partial \text{L}}{\partial \mathbf{b}_1}\right)_{1k} = \sum_{i=1}^{n} \frac{\partial \text{Relu}\left(h_{1_{ik}}\right)}{\partial h_{1_{ik}}} \left(\left(\mathbf{A} \circ \left(2\left(\mathbf{y}_{\text{pred}} - \mathbf{y}\right).\mathbf{w}_3^T\right)\right).\mathbf{w}_2^T\right)_{ik} \tag{79}$$

The k$^{\text{th}}$ column of $\frac{\partial L}{\partial \mathbf{b}_1}$ is equal to the sum of the elements in the k$^{\text{th}}$ column of

10

$$\mathbf{C} \circ \left( \left( \mathbf{A} \circ \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right) \right) . \mathbf{w}_2^T \right).$$

The gradient results can be summarized as follows:

$$\frac{\partial L}{\partial \mathbf{w}_3} \to \mathbf{h}_{2_{\text{relu}}}^T . 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) \tag{80}$$

$$\left( \frac{\partial L}{\partial \mathbf{b}_3} \right)_k = \sum_{i=1}^{n} \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) \right)_{ik} \tag{81}$$

$$\frac{\partial L}{\partial \mathbf{w}_2} \to \mathbf{h}_{1_{\text{relu}}}^T . \left( \mathbf{A} \circ \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right) \right) \tag{82}$$

$$\left( \frac{\partial L}{\partial \mathbf{b}_2} \right)_k = \sum_{i=1}^{n} \left( \mathbf{A} \circ \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right) \right)_{ik} \tag{83}$$

$$\frac{\partial L}{\partial \mathbf{w}_1} \to \mathbf{x}^T . \left( \mathbf{C} \circ \left( \left( \mathbf{A} \circ \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right) \right) . \mathbf{w}_2^T \right) \right) \tag{84}$$

$$\left( \frac{\partial L}{\partial \mathbf{b}_1} \right)_k = \sum_{i=1}^{n} \left( \mathbf{C} \circ \left( \left( \mathbf{A} \circ \left( 2 \left( \mathbf{y}_{\text{pred}} - \mathbf{y} \right) . \mathbf{w}_3^T \right) \right) . \mathbf{w}_2^T \right) \right)_{ik} \tag{85}$$

The numerical implementation of the above gradients, the updates of the weight matrices and the bias matrices are carried out in a separate Jupyter notebook.