

# Neural Network implementation

Martin Lazo

April 2022

## 1 Notation

$\odot$	Hadamard (element-wise) product
$\mathbf{v}$	Vector, vectors are denoted in lower case and bold
$\mathbf{M}$	Matrix, matrices are denoted in upper case and bold
$\mathbf{v}_i$	$i$ th element of a vector $\mathbf{v}$
$\mathbf{M}_{i,j}$	Element in the $i$ th row and $j$ th column of a matrix $\mathbf{M}$
$\mathbf{M}_{i,}$	$i$ th row of a matrix $\mathbf{M}$
$\mathbf{M}_{,i}$	$i$ th column of a matrix $\mathbf{M}$
$\mathbf{o}_i$	Vector which preserves only the $i$ th element of a vector as a scalar

$$\mathbf{o}_3 \mathbf{v} = \begin{pmatrix} 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = c$$

$\mathbf{O}_i$	Square matrix which preserves only the $i$ th element of a vector
----------------	---

$$\mathbf{O}_3 \mathbf{v} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ c \end{pmatrix}$$

Hidden layers are indexed from 0 to  $n$  with  $n$  being the output layer and 0 being the input layer.

$\mathbf{X}$	Features
$\mathbf{Y}$	Labels
$\mathbf{x}$	Features data point
$\mathbf{y}$	Labels data point
$\mathbf{b}^m$	Bias vector of the $m$ th layer
$\mathbf{W}^m$	Weights matrix of the $m$ th layer
$\phi^m$	Activation function of the $m$ th layer, applied element-wise
$\phi^{m'}$	Derivative of the activation function of the $m$ th layer

## 2 Backpropagation calculus

Neural network is a function  $\hat{y} : \mathbb{R}^a \rightarrow \mathbb{R}^b$  where  $a$  and  $b$  are positive integers. Loss function  $\ell_2$  for a set  $S = \{(\mathbf{y}_n, \mathbf{x}_n) : 1 \leq n \leq N\}$  :

$$\ell_2 = \frac{1}{N} \sum_{\mathbf{y}, \mathbf{x} \in S} (\hat{y}(\mathbf{x}) - \mathbf{y})^T (\hat{y}(\mathbf{x}) - \mathbf{y})$$

The aim is to minimize  $\ell_2$  with respect to weights and biases. This is done by iteratively updating weights using gradient descent:

$$(\mathbf{W}_{i,j}^m)_{n+1} = (\mathbf{W}_{i,j}^m)_n - \alpha \frac{\partial \ell_2}{\partial (\mathbf{W}_{i,j}^m)_n}$$

$$(\mathbf{b}_i^m)_{n+1} = (\mathbf{b}_i^m)_n - \alpha \frac{\partial \ell_2}{\partial (\mathbf{b}_i^m)_n}$$

where  $\alpha$  is the learning rate.

$$\frac{\partial \ell_2}{\partial \mathbf{W}_{i,j}^m} = \frac{2}{N} \sum_{\mathbf{y}, \mathbf{x} \in S} (\hat{y}(\mathbf{x}) - \mathbf{y})^T \frac{\partial \hat{y}}{\partial \mathbf{W}_{i,j}^m}(\mathbf{x})$$

$$\frac{\partial \ell_2}{\partial \mathbf{b}_i^m} = \frac{2}{N} \sum_{\mathbf{y}, \mathbf{x} \in S} (\hat{y}(\mathbf{x}) - \mathbf{y})^T \frac{\partial \hat{y}}{\partial \mathbf{b}_i^m}(\mathbf{x})$$

Neural network output for an input  $\mathbf{x}$ :

$$\hat{y}(\mathbf{x}) = \phi^n(\mathbf{W}^n \phi^{n-1}(\dots \phi^0(\mathbf{W}^0 \mathbf{x} + \mathbf{b}^0) \dots) + \mathbf{b}^n)$$

Partial derivatives w.r.t. weights:

$$\frac{\partial \hat{y}}{\partial \mathbf{W}_{i,j}^m}(\mathbf{x}) = \phi^{n'}(\dots) \odot \dots \odot (\mathbf{W}^{m+2}(\phi^{m+1'}(\dots) \odot (\mathbf{W}^{m+1}(\mathbf{O}_i \phi^{m'}(\dots) \mathbf{o}_j \phi^{m-1}(\dots))))))$$

$$\frac{\partial \hat{y}}{\partial \mathbf{W}_{i,j}^0}(\mathbf{x}) = \phi^{n'}(\dots) \odot \dots \odot (\mathbf{W}^2(\phi^{1'}(\dots) \odot (\mathbf{W}^1(\mathbf{O}_i \phi^{0'}(\dots) \mathbf{x}_j))))$$

$$\frac{\partial \hat{y}}{\partial \mathbf{W}_{i,j}^n}(\mathbf{x}) = \mathbf{O}_i \phi^{n'}(\dots) \mathbf{o}_j \phi^{n-1}(\dots)$$

Partial derivatives w.r.t. biases:

$$\frac{\partial \hat{y}}{\partial \mathbf{b}_i^m}(\mathbf{x}) = \phi^{n'}(\dots) \odot \dots \odot (\mathbf{W}^{m+2}(\phi^{m+1'}(\dots) \odot (\mathbf{W}^{m+1}(\mathbf{O}_i \phi^{m'}(\dots))))))$$

$$\frac{\partial \hat{y}}{\partial \mathbf{b}_i^0}(\mathbf{x}) = \phi^{n'}(\dots) \odot \dots \odot (\mathbf{W}^2(\phi^{1'}(\dots) \odot (\mathbf{W}^1(\mathbf{O}_i \phi^{0'}(\dots))))))$$

$$\frac{\partial \hat{y}}{\partial \mathbf{b}_i^n}(\mathbf{x}) = \mathbf{O}_i \phi^{n'}(\dots)$$