# Equations in the matrix form
# for implementing
# simple artificial neural networks

**Ajit R. Jadhav**, Ph.D.

https://JadhavResearch.wordpress.com

---

**Abstract**

This document presents the basic equations in reference to which artificial neural networks are designed and implemented. The scope is restricted to the simpler feedforward networks, including those having hidden layers. Convolutional and recurrent networks are out of the scope.

Equations are often initially noted using an index-basd notation for the typical element. However, all the equations are eventually cast in the direct matrix form, using a consistent set of notation. Some of the minor aspects of notation were invented to make the presentation as simple and direct as possible.

The presentation here regards a layer as the basic unit. The term "layer" is understood in the same sense in which APIs of modern libraries like TensorFlow-Keras 2.x take it. The presentation here is detailed enough that neural networks with hidden layers could be implemented, starting from the scratch.

---

# Preface

### Raison d'être

I wrote this document mainly for myself, to straighten out the different notations and formulae used in different sources and contexts.

In particular, I wanted to have a document that better matches the design themes used in today's libraries (like TensorFlow-Keras 2.x) than the description in the text-books.

For instance, in many sources, the input layer is presented as consisting of both a fully connected layer and its corresponding activation layer. However, for flexibility, libraries like TF-Keras 2.x treat them as separate layers.

Also, some sources uniformly treat the input of any layer as $\vec{X}$ and output of any layer as activation, $\vec{a}$, but such usage overloads the term "activation". Confusions also creep in because different conventions exist: treating the bias by expanding the input vector with $1$ and the weights matrix with $w_0$; the "to-from" vs "from-to" convention for the weights matrix, etc.

I wanted to have a consistent notation that dealt with all such issues with a uniform, matrix-based notation that came as close to the `numpy ndarray` interface as possible.

### Level of coverage

The scope here is restricted to the simplest ANNs, including the simplest DL networks. Convolutional neural networks and recurrent neural networks are out of the scope.

Yet, this document wouldn't make for a good tutorial for a complete beginner; it is likely to confuse him more than explaining anything to him.

So, if you are completely new to ANNs, it is advisable to go through sources like Nielsen's online book [1] to learn the theory of ANNs. Mazur's fully worked out example of the back-propagation algorithm [2] should also prove to be very helpful, before returning back to this document.

If you already know ANNs, and don't want to see equations in the fully expanded forms—or, plain dislike the notation used here—then a good reference, roughly at the same level as this document, is the set of write-ups/notes by Mallya [3].

### Feedback

Any feedback, especially that regarding errors, typos, inconsistencies in notation, suggestions for improvements, etc., will be thankfully received.

email: aj175tp@yahoo.co.in.

**How to cite this document**

Cite this document as:

<span style="color:red">**Insert citation info. at the time of the Release version**</span>          <span style="color:red">**TBD**</span>

ARJ
Pune, India

# Contents

# 1  Layers as building blocks of ANNs

## 1.1  Data

The data for training the network consists a table having $\mu$ number of rows.

Each row has one tuple for the independent variables, $\{X\}$, and another tuple for the dependent variable $\{Y\}$. The independent variables may be regarded as "features" of the dataset, and the dependent variables form the "targets" of the dataset (which act as the "ground truth" for training of the ANN).

In this document, the input vector $\{X\}$ is assumed to have $M$ number of elements, and the target vector $\{Y\}$ is assumed to have $N$ number of elements.

In case of the MNIST dataset, the tuple for the independent variables consists of a flattened vector of floating point values for the pixel intensities; its size is $M = 28 \times 28 = 784$. The tuple for the dependent variables consists of a one-hot encoded vector of size $N = 10$; it indicates the class of digits to which the input picture belongs.

## 1.2  Layers involved

In this document, the phrase "simplest ANN" is taken to mean an artificial neural network that has *no hidden layer*.

Accordingly, the simplest ANN may be regarded as composed from three layers which appear in the following sequence:

- $L_1$: A fully connected layer: It contains the weights and biases, i.e., the parameters of the model

- $L_2$: An activation layer: It determines how each neuron fires

- $L_3$: A loss layer: It determines the measure for the error between the parameters- and activation-based predictions, and the ground truth contained in the target data.

Many textbooks present the material as if $L_1$ and $L_2$ were a single layer.

# 2  Matrix-notation for vectors, matrices, and operations on them

Before turning to the ANNs, it is important to develop a consistent notation for denoting vectors, matrices and their transposes, and the operations involving them—algebraic operations as well as differentiations. Only those cases are considered which are directly relevant to the equations of ANNs.

## 2.1   Basic notation for vectors, matrices, and indices for the typical elements

### Vectors are always column-matrices

We always represent a *vector* using a *column* matrix, and never using a row matrix. A column matrix is denoted using the braces (curly brackets); e.g. $\{V\}$.

We always represent the transpose of a vector, i.e. a row-vector, using the floor symbols, as in $\lfloor V^T \rfloor = \{V\}^T$.

We always represent a matrix (whether a square matrix or a general-rectangular matrix) using square brackets, as in $[A]$.

To avoid confusion, in this article, we never use braces, floor symbols, or square brackets for grouping of terms in an equation. We only use parentheses "( )" for this purpose, even if they get nested.

### Sizes appear in subscripts, sometimes

Sometimes, the dimensions (i.e. the "size" or the number of rows and columns) are explicitly noted using square-brackets in the subscripts, as in

$$\{V\}_{[R \times 1]}, \qquad \text{or} \qquad [W]_{[R \times C]},$$

where $R$ is the number of rows and $C$ is the number of columns.

### Transposed matrices and vectors:  basic form

The transpose of a rectangular matrix is denoted as:

$$[W^T]_{[C \times R]} = \left([W]_{[R \times C]}\right)^T$$

where $R$ and $C$ refer to the number of rows and columns of the original, untransposed matrix. Similarly, a row-matrix is denoted as

$$\lfloor v^T \rfloor_{[1 \times C]} = \left(\{v\}_{[R \times 1]}\right)^T ; \qquad C = R.$$

### Notation for vectors using the typical elements

Vectors are sometimes also denoted using a notation involving its typical symbol and its shape (whether a column- or a row-matrix):

$$\{X\} \overset{\text{def}}{=} \begin{Bmatrix} x_1 \\ x_2 \\ \vdots \\ x_r \\ \vdots \\ x_M \end{Bmatrix} \overset{\text{denoted}}{=} \{x_r\}$$

7

and

$$\lfloor A \rfloor \stackrel{\text{def}}{=} \lfloor a_1 \quad a_2 \quad \ldots \quad a_c \quad \ldots \quad a_N \rfloor \stackrel{\text{denoted}}{=} \lfloor a_c \rfloor$$

**Notation for matrices using the typical elements**

Similarly, matrices are sometimes denoted using a notation that involves two indices for its typical symbol, put inside the rectangular matrix symbol. Thus:

$$
[W]_{[M \times N]} \stackrel{\text{def}}{=}
\begin{bmatrix}
w_{11} & w_{12} & \cdots & w_{1c} & \cdots & w_{1N} \\
w_{21} & w_{22} & \cdots & w_{2c} & \cdots & w_{2N} \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
w_{r1} & w_{r2} & \cdots & w_{rc} & \cdots & w_{rN} \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
w_{M1} & w_{M2} & \cdots & w_{Mc} & \cdots & w_{MN}
\end{bmatrix}_{[M \times N]}
$$

$$\stackrel{\text{denoted}}{=} \left[ w_{rc} \right]_{[M \times N]}$$

where $M$ is the number of rows and $N$ is the number of columns. The transpose of this matrix would be denoted as:

$$\left[ W^T \right]_{[N \times M]} \stackrel{\text{def}}{=} \left( [W]_{[M \times N]} \right)^T$$

$$\stackrel{\text{denoted}}{=}
\begin{bmatrix}
wt_{11} & wt_{12} & \cdots & wt_{1c} & \cdots & wt_{1M} \\
wt_{21} & wt_{22} & \cdots & wt_{2c} & \cdots & wt_{2M} \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
wt_{r1} & wt_{r2} & \cdots & wt_{rc} & \cdots & wt_{rM} \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
wt_{N1} & wt_{N2} & \cdots & wt_{Nc} & \cdots & wt_{NM}
\end{bmatrix}_{[N \times M]}
$$

$$\stackrel{\text{denoted}}{=} \left[ wt_{rc} \right]_{[N \times M]}.$$

Note carefully that the typical element is here listed as $wt$ and not as $w$. The number of rows now is $N$ and the number of columns $M$. Since $wt$ is the typical element and not $w$, the first subscript refers to the row of *transposed* matrix and the second subscript to the column of the *transposed* matrix.

Sometimes, we also note the transpose of $[W]$ using the indices of the *original*

*matrix*, as in:

$$\left[\mathrm{W^T}\right]_{[\mathrm{N \times M}]} \overset{\text{def}}{=} \left(\ \left[\mathrm{W}\right]_{[\mathrm{M \times N}]}\ \right)^{\mathrm{T}}$$

$$\overset{\text{denoted}}{=} \begin{bmatrix} \mathrm{w_{11}} & \mathrm{w_{21}} & \cdots & \mathrm{w_{c1}} & \cdots & \mathrm{w_{N1}} \\ \mathrm{w_{12}} & \mathrm{w_{22}} & \cdots & \mathrm{w_{c2}} & \cdots & \mathrm{w_{N2}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathrm{w_{1r}} & \mathrm{w_{2r}} & \cdots & \mathrm{w_{cr}} & \cdots & \mathrm{w_{Nr}} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \mathrm{w_{1M}} & \mathrm{w_{2M}} & \cdots & \mathrm{w_{cM}} & \cdots & \mathrm{w_{NM}} \end{bmatrix}_{[\mathrm{N \times M}]}$$

$$\overset{\text{denoted}}{=} \left[\mathrm{w_{cr}}\right]_{[\mathrm{N \times M}]}.$$

Note carefully that the typical element listed in the last equation was as $\mathrm{w}$. However, the notation still indicates a *transpose*, because the *order of the indices is now flipped*. The first subscript now refers to the *column* of the *original* matrix, and the second subscript now refers to the *row* of the *original* matrix. Thus $\mathrm{w_{cr}}$.

### Summary: Typical elements-based notation

In short, in the typical element-based notation, we denote a matrix and its transpose as:

$$\left[\mathrm{W}\right] \overset{\text{denoted}}{=} \left[\mathrm{w_{rc}}\right]\ ; \qquad \mathrm{r} \in [1, \mathrm{M}], \mathrm{c} \in [1, \mathrm{N}]$$

$$\left[\mathrm{W^T}\right] \overset{\text{denoted}}{=} \left[\mathrm{wt_{rc}}\right]\ ; \qquad \mathrm{r} \in [1, \mathrm{N}], \mathrm{c} \in [1, \mathrm{M}]$$

$$\overset{\text{denoted}}{=} \left[\mathrm{w_{cr}}\right]\ ; \qquad \mathrm{r} \in [1, \mathrm{M}], \mathrm{c} \in [1, \mathrm{N}]$$

Similarly, a vector and its transpose are denoted as:

$$\{\mathrm{V}\} \overset{\text{denoted}}{=} \{\mathrm{v_r}\}\ ; \qquad \mathrm{r} \in [1, \mathrm{M}]$$

$$\lfloor \mathrm{V^T} \rfloor \overset{\text{denoted}}{=} \lfloor \mathrm{vt_r} \rfloor\ ; \qquad \mathrm{r} \in [1, \mathrm{M}]$$

$$\overset{\text{denoted}}{=} \lfloor \mathrm{v_c} \rfloor\ ; \qquad \mathrm{c} \in [1, \mathrm{N}]\ ; \qquad \text{but for a vector, } \mathrm{N} = \mathrm{M}$$

## 2.2 A rectangular matrix as one column of many rows, or as one row of many columns

A rectangular matrix can also be seen as a column matrix having row-matrices for its elements. It can also be seen as a row matrix having column-matrices for its elements. We adopt the following notation to simplify noting these representa-

9

tions:

$$[W]_{[M \times N]} = \left\{ \ \lfloor w_{r*} \rfloor_{[1 \times N]} \ \right\}_{[M \times 1]} \quad = \quad \left\{ \begin{array}{c} \lfloor w_{1*} \rfloor \\[6pt] \lfloor w_{2*} \rfloor \\ \vdots \\ \lfloor w_{r*} \rfloor \\ \vdots \\ \lfloor w_{M*} \rfloor \end{array} \right\}$$

where $*$ denotes expanding over the entire range of that index.

Similarly,

$$[W]_{[M \times N]} = \left\lfloor \ \{w_{*c}\}_{[M \times 1]} \ \right\rfloor_{[1 \times N]}$$
$$= \left\lfloor \ \{w_{*1}\} \quad \{w_{*2}\} \quad \ldots \quad \{w_{*c}\} \quad \ldots \quad \{w_{*N}\} \ \right\rfloor .$$

This notation denotes a matrix using its typical elements.

## 3   Matrix-notation for algebraic operations between vectors and matrices

### 3.1   Initialization of vectors and matrices with numbers

When all elements of a vector or a matrix are initialized to some number, say 0.0, we use the following notation:

$$\{A\} \leftarrow \{0.0\} .$$

where it is understood that $\{0.0\}$ is a vector of all zeros, and of the same size as that of $\{A\}$.

Similarly, for a matrix:

$$[W] \leftarrow [0.0] .$$

Initialization with random values is noted as:

$$\{A\} \leftarrow \text{random values}$$

where it is understood that each element is initialized to a different random value drawn from a certain distribution (normal or uniform) over a certain numerical range, e.g., $[-1.0, 1.0]$.

## 3.2 Summing all the elements of a vector

Let $\{A\}$ be a vector. The sum of all its elements can be noted, in the index notation, as:

$$S \;=\; \sum_{r=1}^{M} a_r$$

We note this operation, in our notation, as:

$$S \;\overset{\text{denoted}}{=}\; \Sigma\left(\; \{A\}\; \right) \;\overset{\text{denoted}}{=}\; \Sigma_r\left(\; \{a_r\}\; \right)$$

Similarly, for summing all elements of a row-vector $\lfloor B \rfloor$

$$S \;=\; \sum_{c=1}^{N} b_c,$$

we use the notations:

$$S \;\overset{\text{denoted}}{=}\; \Sigma\left(\; \lfloor B \rfloor\; \right) \;\overset{\text{denoted}}{=}\; \Sigma_c\left(\; \lfloor b_c \rfloor\; \right)$$

## 3.3 Multiplication of a vector by a scalar

We denote the multiplication of each element of a vector with a scalar as:

$$s\,\{A\} \;\overset{\text{denoted}}{=}\; \{s\,a_r\}$$

Similarly,

$$s\,\lfloor B \rfloor \;\overset{\text{denoted}}{=}\; \lfloor s\,a_c \rfloor$$

Similar operations can be defined for element-wise divisions, as also for element-wise additions and subtractions.

## 3.4 Functions acting element-wise on one vector to produce another vector

Let $\{X\}$ and $\{Y\}$ be two vectors of equal lengths.

If each $r$-th element of $\{Y\}$, $y_r$, can be obtained by applying the same scalar function $\mathcal{F}$ on the corresponding elements of $\{X\}$, $x_r$, i.e., if

$$y_r \;=\; \mathcal{F}_r\left(\; x_r\; \right) \qquad \forall\; 1 \leq r \leq M,$$

then we denote this fact as:

$$\{y_r\} \;=\; \mathcal{F}_r\left(\; \{x_r\}\; \right)$$
$$\{Y\} \;=\; \mathcal{F}_r\left(\; \{X\}\; \right),$$

where it is understood that there can be a set of different row-wise acting functions $\mathcal{F}_r$.

Similarly, we can define a function relationship between two transposed vectors:

$$\lfloor y_c \rfloor \stackrel{\text{def}}{=} \mathcal{F}_c \left( \lfloor x_c \rfloor \right)$$

$$\lfloor Y \rfloor \stackrel{\text{def}}{=} \mathcal{F}_c \left( \lfloor X \rfloor \right)$$

Further, we can also define functions that send some column-matrix to a row-matrix, and vice-versa, as:

$$\lfloor y_c \rfloor = \tau_{r \to c} \left( \{ x_r \} \right)$$

$$\{ y_r \} = \tau_{c \to r} \left( \lfloor x_c \rfloor \right)$$

and

$$\lfloor Y \rfloor = \tau_{r \to c} \left( \{ X \} \right)$$

$$\{ Y \} = \tau_{c \to r} \left( \lfloor X \rfloor \right)$$

Seen in this light, the transpose operation is a special function that converts a column-matrix to a row-matrix while keeping the same elements, and vice versa:

$$\lfloor X^T \rfloor \stackrel{\text{def}}{=} \mathcal{T}_{r \to c} \left( \{ X \} \right) \stackrel{\text{def}}{=} \left( \{ X \} \right)^T$$

$$\{ X \} \stackrel{\text{def}}{=} \mathcal{T}_{c \to r} \left( \{ X^T \} \right) \stackrel{\text{def}}{=} \left( \{ X^T \} \right)^T$$

### 3.5 Algebraic binary operators between two vectors

We use the following symbols for operations on vectors (i.e. column matrices):

- Element-wise addition of $\{ X \}$ and $\{ Y \}$:

$$\{ A \} = \{ X \} \oplus \{ Y \}$$

$$\{ a_r \} = \{ x_r \} \oplus \{ y_r \}$$

$$\stackrel{\text{def}}{=} \{ x_r + y_r \}$$

Similarly for the following operators.

- Element-wise subtraction of $\{ Y \}$ from $\{ X \}$ is denoted as:

$$\{ S \} = \{ X \} \ominus \{ Y \} \qquad \Leftrightarrow \qquad \{ s_r \} \stackrel{\text{def}}{=} \{ x_r - y_r \}$$

- The Hadamard product: Element-wise multiplication of $\{ X \}$ and $\{ Y \}$ is denoted as:

$$\{ H \} = \{ X \} \circledast \{ Y \} \qquad \Leftrightarrow \qquad \{ h_r \} \stackrel{\text{def}}{=} \{ a_r b_r \}$$

12

- Element-wise division of $\{A\}$ by $\{B\}$ is denoted as:

$$\{D\} \;=\; \{A\} \oslash \{B\} \qquad \Leftrightarrow \qquad \{d_r\} \;\overset{\text{def}}{=}\; \left\{\frac{a_r}{b_r}\right\}$$

The last two operations are not found defined in the usual undergraduate courses on vector algebra. However, the `ndarrays` from the Python library `numpy` come with intuitive operator overloadings for all these four operations.

## 3.6  The outer- and inner-products in the vector and matrix notations

It is to be realized that when writing equations using an explicit matrix form, there is no need to have separate symbols for the outer- or inner-products. To appreciate this point, consider the following vector equations:

$$\vec{A} \cdot \vec{B} \;=\; s \qquad\qquad \text{(Inner, i.e. dot, product)}$$
$$\vec{A} \otimes \vec{B} \;=\; \overline{\overline{\sigma}} \qquad\qquad \text{(Outer product)}$$

where $s$ is a single number (a scalar), and the double-overbars denotes a tensor.

Notice that the symbols for the operators, viz. $\cdot$ and $\otimes$, explicitly appeared in the above equations. If the same equations are expressed using the matrix notation, appropriate transposes have to be taken, but there is no need to insert the explicit symbols for the inner or outer products. The rules of the transpose-taking and order-based matrix-multiplications together take care of the appropriate semantics. Thus:

$$\lfloor A^T \rfloor_{[1 \times M]} \{B\}_{[M \times 1]} \;=\; s \qquad\qquad \text{(Inner, i.e. dot, product)}$$
$$\{A\}_{[M \times 1]} \lfloor B^T \rfloor_{[1 \times N]} \;=\; [\sigma]_{[M \times N]} \qquad\qquad \text{(Outer product)}$$

Further, using our typical element-based notation, the same equations will now be put as:

$$\lfloor a_c \rfloor_{[1 \times M]} \{b_r\}_{[M \times 1]} \;=\; \sum_{i=1}^{M} a_i b_i \;=\; c \qquad \text{(Inner, i.e. dot, product)}$$
$$\{a_r\}_{[M \times 1]} \lfloor b_c \rfloor_{[1 \times N]} \;=\; [t_{rc}]_{[M \times N]} \;; t_{rc} = a_r b_c \qquad \text{(Outer product)}$$

In this document, we use mostly the matrix notation. Hence, the $\otimes$ and the $\cdot$ operators will not appear in the equations as stated. In contrast, we will often use the four element-wise operators noted earlier.

To avoid confusion, make sure to note again the difference between the $\otimes$ operator (the outer product) which will *not* appear in our matrix-based equations, and the $\circledast$ operator (the Hadamard product) which will.

### 3.7 Algebraic binary operators between the same rows of a vector and a matrix

If the number of rows of a vector and a matrix is the same, then, all the afore-mentioned element-wise operations can be conducted, in a row-by-row manner, between the two.

For instance, the element-wise multiplication of the elements of a column vector and the rows of a matrix can be denoted as:

$$[P]_{[M \times N]} = \{V\}_{[M \times 1]} \circledast [D]_{[M \times N]},$$

which means:

$$
\begin{aligned}
[p_{rc}]_{[M \times N]} &= \{v_r\} \circledast [d_{rc}] \\
&= \{v_r\} \circledast \{\lfloor d_{r*} \rfloor\} \\
&= \{v_r\}_{[M \times 1]} \circledast \left\{ \lfloor d_{r1} \quad d_{r2} \quad \ldots \quad d_{rN} \rfloor_{[1 \times N]} \right\}_{[M \times 1]} \\
&= \left\{ \lfloor v_r d_{r1} \quad v_r d_{r2} \quad \ldots \quad v_r d_{rN} \rfloor_{[1 \times N]} \right\}_{[M \times 1]} \\
&= \begin{bmatrix} v_1 d_{11} & v_1 d_{12} & \ldots & v_1 d_{1N} \\ v_2 d_{21} & v_2 d_{22} & \ldots & v_2 d_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ v_M d_{M1} & v_r d_{M2} & \ldots & v_M d_{MN} \end{bmatrix}_{[M \times N]}.
\end{aligned}
$$

Of course, a similar operation can be defined also for a row-vector and a matrix, provided that the number of columns of the row-vector is equal to the number of columns of the matrix:

$$[P]_{[M \times N]} = \lfloor V \rfloor_{[1 \times N]} \circledast [D]_{[M \times N]},$$

because the number of columns N is the same for both the operands. In the typical element-notation:

$$
\begin{aligned}
[p_{rc}]_{[M \times N]} &= \lfloor v_c \rfloor_{[1 \times N]} \circledast [d_{rc}]_{[M \times N]} \\
&= \lfloor v_c \rfloor_{[1 \times N]} \circledast \left\lfloor \{d_{*c}\}_{[M \times 1]} \right\rfloor_{[1 \times N]} \\
&= \begin{bmatrix} v_1 d_{11} & v_2 d_{12} & \ldots & v_N d_{1N} \\ v_1 d_{21} & v_2 d_{22} & \ldots & v_N d_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ v_1 d_{M1} & v_2 d_{M2} & \ldots & v_N d_{MN} \end{bmatrix}
\end{aligned}
$$

Other binary operators can be defined on similar lines.

### 3.8 Converting a vector into its diagonal matrix form

**Our notation**

Let's take a concrete example. Consider

$$\{A\}_{[3 \times 1]} \overset{\text{def}}{=} \begin{Bmatrix} a_1 \\ a_2 \\ a_3 \end{Bmatrix}_{[3 \times 1]}$$

Its corresponding diagonal matrix is given by:

$$\begin{bmatrix} a_1 & 0 & 0 \\ 0 & a_2 & 0 \\ 0 & 0 & a_3 \end{bmatrix}_{[3 \times 3]} .$$

In our notation, we denote the above matrix as:

$$\begin{bmatrix} \mathbf{diag}\,\{A\} \end{bmatrix}_{[3 \times 3]}$$

**Identity matrix**

The identity matrix is defined as a matrix with all its diagonal elements 1, and all other elements 0:

$$[\mathcal{I}] \overset{\text{def}}{=} \begin{bmatrix} 1 & 0 & \ldots & 0 \\ 0 & 1 & \ldots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \ldots & 1 \end{bmatrix}_{[N \times N]}$$

The identity matrix may be seen as embedding the operation of converting the identity vector to its diagonal matrix form.

**Diagonal form for an arbitrary vector**

Let $\{A\}$ be a $[N \times 1]$-sized column vector. Then, perform:

$$[O]_{[N \times N]} \overset{\text{def}}{=} \{A\}_{[N \times 1]} \otimes \lfloor 1 \rfloor_{[1 \times N]} \tag{1}$$

where, remember, $\otimes$ represents the outer product. Realize, the result of the outer product here looks like:

$$[O]_{[N \times N]} = \begin{bmatrix} a_1 & a_1 & \ldots & a_1 \\ a_2 & a_2 & \ldots & a_2 \\ \vdots & \vdots & \ddots & \vdots \\ a_N & a_N & \ldots & a_N \end{bmatrix}_{[N \times N]}$$

Now, the diagonal form for $\{A\}$ can be easily obtained as:

$$\left[\mathbf{diag}\{A\}\right]_{[N \times N]} = \left[O\right]_{[N \times N]} \circledast \left[\mathcal{I}\right]_{[N \times N]}$$

$$= \left( \{A\}_{[N \times 1]} \otimes \lfloor 1 \rfloor_{[1 \times N]} \right) \circledast \left[\mathcal{I}\right]_{[N \times N]}$$

where $\circledast$ implies multiplications of the corresponding elements.

The diagonal matrix-form of the identity vector $\lfloor 1 \rfloor$ is the identity matrix. It *can* be seen as row-wise or column-wise additions of unit vectors to the zero-matrix. Similarly, the vector $\{A\}$ could also be put into the diagonal matrix form using row-wise matrix additions its component-vectors. However, we will not pursue these matters because our focus here is on ANNs.

## 4 Matrix-notation for element-wise differentiations of vectors and matrices

**The importance of the indices of the typical element**

The most "unintuitive" part of our notation appears in denoting differentiations. The best strategy here is to

- first focus on the typical elements (like the value at r-th row and c-th column),

- then express the desired differentiations among these typical elements,

- then put the typical differential expression inside the symbol that indicates the shape of the resulting matrix (a column matrix, a row matrix, or a rectangular matrix).

The subsection sub-sections will make clear what exactly we mean.

**The notation for functions relating vectors and matrices**

In general, a given matrix $\left[B\right]$ could be the result of applying a *different* function $f_{ij}$ to each element of another matrix $\left[A\right]$. Thus,

$$\left[b_{ij}\right] = \left[f_{ij}\left(a_{ij}\right)\right],$$

16

which expands to

$$
\begin{bmatrix} b_{11} & b_{12} & \dots & b_{1N} \\ b_{21} & b_{22} & \dots & b_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ b_{M1} & b_{M2} & \dots & b_{MN} \end{bmatrix}
$$
$$
= \begin{bmatrix} f_{11}\left(a_{11}\right) & f_{12}\left(a_{12}\right) & \dots & f_{1N}\left(a_{1N}\right) \\ f_{21}\left(a_{21}\right) & f_{22}\left(a_{22}\right) & \dots & f_{2N}\left(a_{2N}\right) \\ \vdots & \vdots & \ddots & \vdots \\ f_{M1}\left(a_{M1}\right) & f_{M2}\left(a_{M2}\right) & \dots & f_{MN}\left(a_{MN}\right) \end{bmatrix}.
$$

In this document, we denote even the above kind of a very general functional relationship using only the simple notation:

$$
[B]_{[M \times N]} = \mathcal{F}\left( [A]_{[M \times N]} \right).
$$

Even if a single symbol $\mathcal{F}$ appears on the outside of the parentheses, it is understood that there in fact are $M \times N$ number of different fuctions acting here.

So, when it comes to equations involving functions of vectors and matrices, particular attention will have to be paid to the context in which they arise.

## 4.1 Differentiation of a single variable with respect to the rows of a vector

Suppose that

$$
\mathcal{E} = \mathcal{L}\left( \{A\}_{[N \times 1]}, \{Y\}_{[N \times 1]} \right)
$$

Now, if the single variable $\mathcal{E}$ is differentiated with respect to each element of $\{A\}$, then we use the following notation to capture the resulting column of the respective differential quantities:

$$
\left\{ \begin{array}{c} \partial\mathcal{E}/\partial a_1 \\ \partial\mathcal{E}/\partial a_2 \\ \vdots \\ \partial\mathcal{E}/\partial a_r \\ \vdots \\ \partial\mathcal{E}/\partial a_N \end{array} \right\}_{[N \times 1]} \overset{\text{denoted}}{=} \left\{ \frac{\partial\mathcal{E}}{\partial a_r} \right\}_{[N \times 1]}
$$

Notice that there is no subscript on $\mathcal{E}$ because it is a single variable, not a vector.

## 4.2 Row-wise differentiation of a vector with respect to another vector

Suppose that:

$$
\{B\}_{[N \times 1]} = \mathcal{F}\left( \{A\}_{[N \times 1]}, \dots \right).
$$

where ... indicates some more, unspecified, independent variables. In terms of typical elements,

$$\{b_r\} = \mathcal{F}\left(\ \{a_r\}\ , \dots\ \right).$$

Now, if we partially differentiate each $b_r$ element with respect to the corresponding $a_r$ element, then we denote the resulting set in our matrix notation as:

$$\begin{Bmatrix} \partial b_1/\partial a_1 \\ \partial b_2/\partial a_2 \\ \vdots \\ \partial b_r/\partial a_r \\ \vdots \\ \partial b_N/\partial a_N \end{Bmatrix} \overset{\text{denoted}}{=} \left\{\frac{\partial b_r}{\partial a_r}\right\}_{[N \times 1]}$$

## 4.3 Column-wise differentiation of a column-matrix with respect to another column-matrix

Similar results can be had also for transposed vectors. For instance, if

$$\lfloor B \rfloor_{[1 \times N]} = \mathcal{F}\left(\ \lfloor A \rfloor_{[1 \times N]}, \dots\ \right),$$

i.e., in terms of typical elements, if,

$$\lfloor b_r \rfloor = \mathcal{F}\left(\ \lfloor a_r \rfloor\ , \dots\ \right),$$

then,

$$\left\lfloor \frac{\partial b_1}{\partial a_1} \quad \frac{\partial b_2}{\partial a_2} \quad \dots \quad \frac{\partial b_r}{\partial a_r} \quad \dots \quad \frac{\partial b_N}{\partial a_N} \right\rfloor_{[1 \times N]} \overset{\text{denoted}}{=} \left\lfloor \frac{\partial b_r}{\partial a_r} \right\rfloor_{[1 \times N]}$$

## 4.4 Differentiation of each element of a vector with respect to all the elements of another vector

Suppose that each element $a_r$ of a vector $\{A\}$ is a function of *all* the elements of another vector $\{Z\}$. Such a relationship is best expressed by considering the transpose of the second vector. Thus, in terms of typical elements,

$$a_r = \mathcal{F}\left(\ \lfloor z_1 \quad z_2 \quad \dots \quad z_c \quad \dots \quad z_N \rfloor\ , \dots\ \right).$$

Now, suppose we differentiate each $a_r$ with respect not just the corresponding $z_r$, but *all* the rest of them too, taken separately. Then, they naturally form an

$[N \times N]$ matrix:

$$
\begin{bmatrix}
\partial a_1/\partial z_1 & \partial a_1/\partial z_2 & \ldots & \partial a_1/\partial z_c & \ldots & \partial a_1/\partial z_N \\
\partial a_2/\partial z_1 & \partial a_2/\partial z_2 & \ldots & \partial a_2/\partial z_c & \ldots & \partial a_2/\partial z_N \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
\partial a_r/\partial z_1 & \partial a_r/\partial z_2 & \ldots & \color{red}{\partial a_r/\partial z_c} & \ldots & \partial a_r/\partial z_N \\
\vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\
\partial a_N/\partial z_1 & \partial a_N/\partial z_2 & \ldots & \partial a_N/\partial z_c & \ldots & \partial a_N/\partial z_N
\end{bmatrix}_{[N \times N]}
\tag{2}
$$

$$
\overset{\text{denoted}}{=} \quad \left[ \frac{\partial a_r}{\partial z_c} \right] \qquad \overset{\text{denoted}}{=} \quad \mathbf{J}_{rc},
$$

where $\mathbf{J}_{rc}$ stands for the Jacobian matrix. The pnemonic for the $\left[ \bullet_{rc} \right]$ order for the indices is this:

Hold the the row $_r$ fixed, and unroll all the $_c$ columns!. Then repeat for the next rows.

Note that, some times, in literature, the Jacobian matrix is defined as the transpose of the matrix we gave. Be careful about what index occupies what place.

A further generalization of the Jacobian matrix for the case when $\{B\}$ has $N$ elements, $\{A\}$ has $M$ elements, and $M \neq N$, is possible. In such a case, the Jacobian matrix has $[M \times N]$ size. However, we will not need it in this document.

### 4.5 Differentiation of each element of a vector with respect to the corresponding row of a matrix

Suppose that an $[N \times 1]$-sized vector $\{Z\}$ is a function of a matrix $\left[ W^T \right]$ of size $[N \times M]$. (Yes, there does exist another matrix $[W]$ of size $[M \times N]$. However, for the fully connected layer, the function of interest is better captured with respect to its transpose.) Thus, suppose

$$
\{Z\}_{[N \times 1]} \;=\; \mathcal{F} \left( \; \left[ W^T \right]_{[N \times M]} , \ldots \right).
$$

Now, express $\left[ W^T \right]$ as a column matrix of rows:

$$
\{Z\}_{[N \times 1]} \;=\; \mathcal{F} \left( \; \left\{ \lfloor W_r^T \rfloor_{[1 \times M]} \right\}_{[N \times 1]} , \ldots \right)
$$

where $\lfloor W_r^T \rfloor$ denotes the entire $r$-th row of the rectangular matrix $\left[ W^T \right]$.

In terms of typical elements,

$$
z_r \;=\; \mathcal{F} \left( \; \lfloor wt_{r*} \rfloor , \ldots \right)
$$

where $\lfloor wt_{r*} \rfloor$ denotes the entire $r$-th row of the rectangular matrix $\left[ W^T \right]$.

**Differentiation with respect to a single typical row**

Now, consider the differentiation of $z_r$ with respect to each element of the r-th row of $\left[W^T\right]$. Thus, the quantities of interest are:

$$\left[\frac{\partial z_r}{\partial wt_{r1}} \quad \frac{\partial a_r}{\partial wt_{r2}} \quad \cdots \quad \frac{\partial z_r}{\partial wt_{rc}} \quad \cdots \quad \frac{\partial z_r}{\partial wt_{rM}}\right].$$

Note how the differential quantities naturally form a row-matrix. Using the $*$-based notation introduced earlier, we can condense down the above long expression to:

$$\left[\frac{\partial z_r}{\partial wt_{r1}} \quad \frac{\partial z_r}{\partial wt_{r2}} \quad \cdots \quad \frac{\partial z_r}{\partial wt_{rc}} \quad \cdots \quad \frac{\partial z_r}{\partial wt_{rM}}\right] \overset{\text{denoted}}{=} \left[\frac{\partial z_r}{\partial wt_{r*}}\right]_{[1 \times M]},$$

where it is understood that the expression refers to a single row-vector of a rectangular matrix.

**Differentiation with respect to a matrix**

Repeating the above procedure for every row of $\left[W^T\right]$, we can say that we have arrived at an expression for the partial differentiation of a vector with respect to a matrix:

$$\begin{bmatrix} \partial z_1/\partial wt_{11} & \partial z_1/\partial wt_{12} & \ldots & \partial z_1/\partial wt_{1c} & \ldots & \partial z_1/\partial wt_{1M} \\ \partial z_2/\partial a_{21} & \partial z_2/\partial wt_{22} & \ldots & \partial z_2/\partial a_{2c} & \ldots & \partial z_2/\partial wt_{2M} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \partial z_r/\partial wt_{r1} & \partial z_r/\partial wt_{r2} & \ldots & \partial z_r/\partial wt_{rc} & \ldots & \partial z_r/\partial wt_{rM} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \partial z_N/\partial wt_{N1} & \partial z_N/\partial wt_{N2} & \ldots & \partial z_N/\partial wt_{Nc} & \ldots & \partial z_N/\partial wt_{NM} \end{bmatrix} \tag{3}$$

$$\overset{\text{denoted}}{=} \left[\frac{\partial z_r}{\partial wt_{rc}}\right]_{[N \times M]}$$

The logic behind adopting such unusual notation is simply this much: The actual operations happen, at the most basic level, only with respect to sets of linear equations. If the operation involved is differentiation, that form of matrix notation should be adopted which reduces the cognitive load stating the equations.

## 5 Training a network using the mini-batch method: main steps

The description here refers to the mini-batch gradient descent. We will take the example of the MNIST dataset.

Suppose that the number of independent variables is $M$. In the MNIST example, $M = 28 \times 28 = 784$.

Suppose that the number of dependent variables is $N$. In the MNIST example, $N = 10$.

Suppose that there are $\mu$ number of rows in the training dataset. In the MNIST example, $\mu = 60,000$.

An epoch refers to processing all the rows of the training dataset (here, $60,000$ of them) once. Training is conducted for many such epochs, may be up to a few hundred, until the errors drop down sufficiently. We will not look into the criteria to stop training.

Now, the processing flow is the following.

## 5.1  Initialize the parameters of the model

### Convention for denoting the weights matrix

In this document, we adopt the convention that the weights matrix for the fully connected layer is arranged in the "from–to" order. Thus, if you take the typical element $w_{ij}$, then it refers to the weight of the edge *from* i-th *input* neuron *to* the j-th *output* neuron.

With this convention, all our matrix equations come to refer to only the transpose of the weight matrix, i.e., to $\left[W^T\right]$, which can be seen as being listed in the "to–from" order. Thus:

$$\{Z\} \;=\; \left[W^T\right]\{X\} \;\oplus\; \{B\}$$

In index notation:

$$z_i \;=\; \sum_{j=1}^{M} \left(\, wt_{ij}\, x_j \,\right) + b_j, \tag{4}$$

where the index i runs over all rows of $\{Z\}$ and the index j runs over all columns of $\left[W^T\right]$ (or all rows of $\left[W\right]$) and also over all rows of $\{X\}$.

In matrix notation using the typical elements:

$$\begin{aligned}
\{z_i\} \;&=\; \left[wt_{ij}\right]\{x_j\} \;\oplus\; \{b_j\} \tag{5}\\
&=\; \left[w_{ji}\right]\{x_j\} \;\oplus\; \{b_j\}
\end{aligned}$$

We follow the "from–to" convention for two reasons: (i) the ordering is in the natural forward sequence, i.e., from input to output, and not from output to input, and (ii) the equations in the matrix form remain in direct correspondence with those in the vector-tensor form, as in:

$$\vec{Z} \;=\; \overline{\overline{W}} \cdot \vec{X} + \vec{B}$$

### Initialization

First of all, before any training begins, initialize the parameters (the weights and biases) with some random values.

$$\left[\mathrm{W^T}\right]_{[\mathrm{N \times M}]} \leftarrow \left[\text{random values}\right] \qquad \text{range: } -1.0 \text{ to } 1.0$$
$$\left[\mathrm{B}\right]_{[\mathrm{N \times 1}]} \leftarrow \left\{\text{random values}\right\} \qquad \text{range: } -1.0 \text{ to } 1.0$$

During training, these values will get *updated* at the end of each mini-batch to reflect learning, but they will never be *initialized* at any step during the training.

Training is conducted continuously for a number of epochs (say, between a few tens to a few hundreds).

## 5.2   Processing during each epoch

Processing during each epoch consists of the following steps.

### Randomly select rows for all the mini-batches of an epoch

At the beginning of each epoch, row-wise shuffle the dataset.

Then divide the datasets into a number of mini-batches. For instance, in case of the MNIST dataset, mini-batches of $2,000$ rows might be created. Denote the number of rows in a mini-batch as $\nu$.

### Initialize the accumulated gradients

Let $\left[\mathcal{D}_{\mathcal{E},\mathrm{W^T}}\right]$ and $\left\{\mathcal{D}_{\mathcal{E},\mathrm{B}}\right\}$ denote the gradients accumulated at the end of processing of a mini-batch.

At the beginning of each mini-batch, initialize these matrices to zero:

$$\left[\mathcal{D}_{\mathcal{E},\mathrm{W^T}}\right] \leftarrow \left[0.0\right]$$
$$\left\{\mathcal{D}_{\mathcal{E},\mathrm{B}}\right\} \leftarrow \left\{0.0\right\}$$

### Process all the rows of a mini-batch

For each row of the current mini-batch, make one forward pass, and one backward pass. Gradients of errors with respect to parameters will be obtained at the end of the backward pass. Accumulate these gradients for all rows of the current mini-batch:

$$\left[\mathcal{D}_{\mathcal{E},\mathrm{W^T}}\right] \leftarrow \left[\mathcal{D}_{\mathcal{E},\mathrm{W^T}}\right] \oplus \left[\mathrm{D}_{\mathcal{E},\mathrm{W^T}}\right]_{\mathrm{r}}$$
$$\left\{\mathcal{D}_{\mathcal{E},\mathrm{B}}\right\} \leftarrow \left\{\mathcal{D}_{\mathcal{E},\mathrm{B}}\right\} \oplus \left\{\mathrm{D}_{\mathcal{E},\mathrm{B}}\right\}_{\mathrm{r}}$$

where the quantities in black, i.e. $\left[\mathrm{D}_{\mathcal{E},\mathrm{W^T}}\right]_{\mathrm{r}}$ and $\left\{\mathrm{D}_{\mathcal{E},\mathrm{B}}\right\}$, represent the gradients computed at the end of the r-th row, and the quantities in blue represent the accumulated gradients.

**Update the weights and biases**

At the end of a mini-batch, update all the weights and biases, using the accumulated gradients:

$$\left[\mathrm{W^T}\right]_{[\mathrm{M \times N}]} \quad \leftarrow \quad \left[\mathrm{W^T}\right]_{[\mathrm{M \times N}]} \quad \ominus \quad \rho\left[\mathrm{D}_{\mathcal{E},\mathrm{W^T}}\right]_{[\mathrm{N \times M}]}$$

$$\left[\mathrm{B}\right]_{[\mathrm{N \times 1}]} \quad \leftarrow \quad \left[\mathrm{B}\right]_{[\mathrm{N \times 1}]} \quad \ominus \quad \rho\left\{\mathcal{D}_{\mathcal{E},\mathrm{B}}\right\}$$

where $\rho$ is the learning rate, a hyper-parameter of the training process (i.e., it is not a parameter denoting learning). The minus sign appears because the gradient by definition points "uphill", whereas the objective to reach the global minimum, not the maximum.

## 5.3   Variations for processing the rows of a dataset

The stochastic gradient descent and the batch-gradient descent are two other variations used during the training phase.

Further, instead of accumulating gradients as we have done here, it is possible to increase performance using certain matrix-based techniques. Consult sources like Nielsen [1] for such variations.

## 6   Forward pass through each layer

In this section and the next, we will develop the equations used in the forward and backward passes through all the three layers (the fully connected, activation, and loss layers). At this point, it is assumed that the reader is familiar with the ideas of forward and backward passes and mini-batches.

The forward pass is conducted for each r-th row of a mini-batch (which has in all $\nu$ number of rows).

Remember, the sizes of the vectors:

- The input vector $\left\{\mathrm{X}\right\}$ has the size of $[\mathrm{M \times 1}]$.

- The output vector $\left\{\mathrm{Y}\right\}$ has the size of $[1 \times \mathrm{N}]$.

- Therefore, in a simple ANN without any hidden layer, all the following vectors have the size of $[1 \times \mathrm{N}]$:

    - The bias vector of the fully connected layer, $\left\{\mathrm{B}\right\}$.

    - The output of the fully connected layer $\left\{\mathrm{Z}\right\}$, which also is the input of the activation layer.

    - The output of the activation layer $\left\{\mathrm{A}\right\}$, which also is one of the inputs of the loss layer (the other input for the loss layer being, the target

vector $\{Y\}$).

Also, remember, we define the weights matrix using the "from input $\{X\}$ to $\{Z\}$" convention. Hence,

- The size of the weights matrix, $[W]$, is $[M \times N]$, whereas,

- the size of its transpose, $[W^T]$ is $[N \times M]$.

## 6.1 Fully connected layer

$$\{z_i\}_r \;\leftarrow\; [W^T]\{x_i\}_r + \{b_i\} \;; \qquad \forall\, i \in [1, N]$$

where the subscript $r$ denotes the index of the current input data-row in the current mini-batch. The parameters, i.e., the weights and biases (printed in the blue color) remain the same for all rows of a mini-batch.

## 6.2 Activation layer

A detailed treatment in respect of the various activation functions is given in a later section. Here, we note the activation equation in generic terms.

**When each element $a_i$ of the activation vector is a function of only the row-wise corresponding $z_i$ elements**

$$\{a_i\}_r \;\leftarrow\; \mathcal{A}\left(\; \{z_i\}_r \;\right) \;; \qquad \forall\, i \in [1, N]$$

The exact matrix operation depends on the activation function $\mathcal{A}$ chosen. Typical examples include the sigmoid, $\tanh$ or ReLU functions for classification problems, and the identify function (i.e. no processing) for regression problems.

**When each element $a_i$ of the activation vector is a function of all the $z_k$ elements**

This case arises when we use the softmax activation function. Each $r$-th element of the activation vector produced by the softmax layer, $a_r$, is a function of not only the row-wise corresponding element $z_r$ but also all the other $z_k$ elements of the $\{Z\}$ vector:

$$\{a_j\}_r \;\leftarrow\; \mathcal{A}\left(\; [z_{i \to j}]_r \;\right) \;; \qquad \forall\, i \in [1, N],\; \forall\, k \in [1, N]$$

## 6.3 Loss layer

$$\mathcal{E}_r \;\leftarrow\; \mathcal{L}\left(\; \{a_i\}_r,\; \{y_i\}_r \;\right)$$

The exact matrix operation depends on the loss function $\mathcal{L}$ chosen. Typical examples include the SSE (sum of squared errors) for regression, and the log-loss i.e. the cross-entropy loss for classification. A detailed treatment of these loss functions is given in a later section.

Realize that $\mathcal{E}_r$ is a single number.

# 7 Backward pass through each layer

The backward pass is conducted once per row of the mini-batch.

## 7.1 Loss layer

Backward pass through the loss layer:

$$\left\{ D_{\mathcal{E},A} \right\}_{[N \times 1]} \leftarrow \left\{ \frac{\partial \mathcal{E}}{\partial a_i} \right\}_{[N \times 1]}.$$

The exact expression for the right hand-side will depend on the particular loss layer actually being used.

All the vectors here are of the dimensions $[N \times 1]$.

## 7.2 Activation layer

### Case 1: Each $a_i$ is a function of the corresponding $z_i$ only

In such a case, a backward pass through the activation layer can be written as:

$$\left\{ D_{A,Z} \right\}_{[N \times 1]} \leftarrow \left\{ \frac{\partial a_i}{\partial z_i} \right\}_{[N \times 1]}$$

The exact expression for the right hand-side will depend on the particular activation layer actually being used.

Then, by chain rule,

$$
\begin{aligned}
\left\{ \frac{\partial \mathcal{E}}{\partial z_i} \right\}_{[N \times 1]} &= \left\{ \left( \frac{\partial \mathcal{E}}{\partial a_i} \right) \left( \frac{\partial a_i}{\partial z_i} \right) \right\}_{[N \times 1]} \\
&= \left\{ \frac{\partial \mathcal{E}}{\partial a_i} \right\}_{[N \times 1]} \circledast \left\{ \frac{\partial a_i}{\partial z_i} \right\}_{[N \times 1]} ; \qquad \forall \, i \in [1, N].
\end{aligned}
$$

In other words,

$$\left\{ D_{\mathcal{E},Z} \right\}_{[N \times 1]} \leftarrow \left\{ D_{\mathcal{E},A} \right\}_{[N \times 1]} \circledast \left\{ D_{A,Z} \right\}_{[N \times 1]}$$

All the vectors here are of the dimensions $[N \times 1]$.

**Case 2: Each $a_i$ is a function of all $z_j$'s i.e., of $z_1, z_2, \ldots, z_N$**

$$a_i = \mathcal{F}\left( \begin{bmatrix} z_1 & z_2 & \ldots & z_j & \ldots & z_N \end{bmatrix} \right) ; \qquad \forall\, j \in [1, N].$$

The Jacobian matrix involved here is:

$$\left[ D_{A,z} \right]_{[N_A \times N_z]} \stackrel{\text{def}}{=} \mathbf{J}_{ij} = \left[ \frac{\partial a_i}{\partial z_j} \right]_{[N_A \times N_z]},$$

where $N_z = N_A = N$, but we have noted them separately only for convenience in reading, in understanding the vector from where a given range comes. See eq. (2) for the full expansion of the Jacobian matrix.

Now, by chain rule,

$$\left\{ \frac{\partial \mathcal{E}}{\partial z_i} \right\}_{[N_z \times 1]} = \left\{ \sum_{j=1}^{N_A} \left( \frac{\partial \mathcal{E}}{\partial a_j} \right) \left( \frac{\partial a_j}{\partial z_i} \right) \right\}_{[N_z \times 1]}$$

$$= \left[ \frac{\partial a_j}{\partial z_i} \right]_{[N_z \times N_A]} \left\{ \frac{\partial \mathcal{E}}{\partial a_j} \right\}_{[N_A \times 1]}$$

$$= \left( \mathbf{J}_{ij} \right)^{\mathrm{T}}_{[N_z \times N_A]} \left\{ \frac{\partial \mathcal{E}}{\partial a_j} \right\}_{[N_A \times 1]},$$

Thus, using our matrix notation:

$$\left\{ D_{\mathcal{E},z} \right\}_{[N_z \times 1]} \leftarrow \left( \left[ D_{A,z} \right] \right)^{\mathrm{T}}_{[N_z \times N_A]} \left\{ D_{\mathcal{E},A} \right\}_{[N_A \times 1]}. \tag{6}$$

## 7.3 Fully connected layer

**Derivation**

The equation for forward pass through the fully connected layer is eq. (5), which we reproduce below with color-coding:

$$\left\{ z_i \right\} \leftarrow \left[ wt_{ij} \right] \left\{ x_j \right\} \oplus \left\{ b_j \right\}.$$

Of course, at this point in the backward pass,

$$\mathcal{E} = \mathcal{F}\left( \left\{ Z \right\} \right)$$

So, for the backward pass, by chain rule,

$$\left[ \frac{\partial \mathcal{E}}{\partial wt_{ij}} \right]_{[N \times M]} = \left\{ \frac{\partial \mathcal{E}}{\partial z_i} \right\}_{[N \times 1]} \circledast \left[ \frac{\partial z_i}{\partial wt_{ij}} \right]_{[N \times M]}$$

$$= \left\{ \frac{\partial \mathcal{E}}{\partial z_i} \right\}_{[N \times 1]} \circledast \left\{ \left\lfloor \frac{\partial z_i}{\partial wt_{i*}} \right\rfloor_{[1 \times M]} \right\}_{[N \times 1]}$$

Now, given eq. (5), if we conduct the partial differentiations for the $z_r$ with respect to the r-th row of $\left[W^T\right]$, the result turns out to be:

$$\left\lfloor \frac{\partial z_i}{\partial wt_{i*}} \right\rfloor_{[1 \times M]} = \left\lfloor xt_i \right\rfloor_{[1 \times M]}$$

where $xt_i$ is the i-th element of $\left\lfloor X^T \right\rfloor$, i.e., of the transpose of $\{X\}$. Therefore,

$$\left[ \frac{\partial \mathcal{E}}{\partial wt_{ij}} \right]_{[N \times M]} = \left\{ \frac{\partial \mathcal{E}}{\partial z_i} \right\}_{[N \times 1]} \circledast \left\{ \left\lfloor xt_i \right\rfloor_{[1 \times M]} \right\}_{[N \times 1]}$$

**Gradients of errors with respect to weights**

Let

$$\left[ D_{Z,W^T} \right]_{[N \times M]} \overset{\text{denoted}}{=} \left[ \frac{\partial z_i}{\partial wt_{ij}} \right]_{[N \times M]}.$$

For full expansion of $\left[ \dfrac{\partial z_i}{\partial wt_{ij}} \right]$, see eq. (3).

In the light of the immediately preceding paragraph, the gradients of $\{Z\}$ with respect to $\left[W^T\right]$ turn out to be N number of copies of the row-vectors $\left\lfloor X^T \right\rfloor$ stacked in a column:

$$\left[ D_{Z,W^T} \right]_{[N \times M]} \leftarrow \left\{ \begin{array}{c} \left\lfloor X^T \right\rfloor_{[1 \times M]} \\[6pt] \left\lfloor X^T \right\rfloor_{[1 \times M]} \\[6pt] \vdots \\[6pt] \left\lfloor X^T \right\rfloor_{[1 \times M]} \end{array} \right\}_{[N \times 1]}.$$

Then

$$\left[ D_{\mathcal{E},W^T} \right]_{[N \times M]} \leftarrow \left\{ D_{\mathcal{E},Z} \right\}_{[N \times 1]} \circledast \left[ D_{Z,W^T} \right]_{[N \times M]}$$

**Gradients of errors with respect to biases**

Let

$$\left\{ D_{Z,B} \right\} \overset{\text{denoted}}{=} \left\{ \frac{\partial z_i}{\partial b_i} \right\}.$$

Then, it should be the case that

$$\left\{ D_{\mathcal{E},B} \right\} \leftarrow \left\{ D_{\mathcal{E},Z} \right\} \circledast \left\{ D_{Z,B} \right\}$$

However, given eq. (5), it so turns out that

$$\left\{\frac{\partial z_i}{\partial b_i}\right\} = \{1.0\}.$$

Therefore,

$$\{D_{\mathcal{E},B}\} \leftarrow \{D_{\mathcal{E},Z}\}$$

**Gradients of errors with respect to input**

In case of networks with hidden layers, there are many pairs of (fully connected layer + activation layers). The output of the activation layer part of the previous pair forms the input for the next hidden pair's fully connected layer part. Hence, errors have to backpropagated to the activation layer part of the previous pair.

Thus, gradients have to be propagated not only to the weights and biases of the current fully connected layer, but also to its input.

Let the index j run over $\{X\}$ and i over $\{Z\}$. Then, reproducing eq. (4), we have

$$z_i = b_i + \sum_{j=1}^{M} \left( wt_{ij} \, x_j \right).$$

Therefore, for each pair of i and j,

$$\frac{\partial z_i}{\partial x_j} = wt_{ij} = w_{ji}$$

So, in our matrix notation:

$$\left[\frac{\partial z_i}{\partial x_j}\right]_{[N \times M]} = \left[W^T\right]_{[N \times M]}$$

Then, by chain rule, for each typical element,

$$\frac{\partial \mathcal{E}}{\partial x_j} = \frac{\partial z_i}{\partial x_j}\frac{\partial \mathcal{E}}{\partial z_i}$$

So, for all the elements,

$$\left\{\frac{\partial \mathcal{E}}{\partial x_j}\right\}_{[M \times 1]} = \left( \left[\frac{\partial z_i}{\partial x_j}\right]_{[N \times M]} \right)^T \left\{\frac{\partial \mathcal{E}}{\partial z_i}\right\}_{[N \times 1]}$$

$$= \left( [wt_{ij}]_{[N \times M]} \right)^T \left\{\frac{\partial \mathcal{E}}{\partial z_j}\right\}_{[N \times 1]}$$

$$= \left( [W^T]_{[N \times M]} \right)^T \left\{\frac{\partial \mathcal{E}}{\partial z_j}\right\}_{[N \times 1]}$$

$$= [W]_{[M \times N]} \left\{\frac{\partial \mathcal{E}}{\partial z_j}\right\}_{[N \times 1]}$$

Hence, using our matrix notation:

$$\{D_{\mathcal{E},X}\} \leftarrow [W]\{D_{\mathcal{E},Z}\}$$

28

## 8 Activation functions and their gradients

A discussion of the relative advantages and disadvantages of the different activation functions used, is out of the scope of this document.

### 8.1 ReLU

**Forward pass: Activation function**

The ReLU (rectified linear unit) is defined as:

$$a_i = \begin{cases} 0.0 & \forall \ z_i < 0.0 \\ z_i & \forall \ z_i \geq 0.0 \end{cases}$$

See [4] for a discussion on the fastest way to implement this method in `numpy`. For instance, the multiplication method could be put, using our matrix notation, as:

$$\{t\} \ \leftarrow \ \{Z\} > \{0.0\}$$
$$\{A\} \ \leftarrow \ \{Z\} \circledast \{t\}$$

where $>$ denotes the element-wise binary comparison operator, and its result, the truth-hood vector $\{t\}$, has either 1.0 or 0.0 as its elements.

**Backward pass: Gradient of $\{A\}$ with respect to $\{Z\}$**

$$\frac{\partial a_i}{\partial z_i} = \begin{cases} 0.0 & \forall \ z_i < 0.0 \\ 1.0 & \forall \ z_i \geq 0.0 \end{cases}$$

See the preceding paragraph for examples of how the case-wise assignment of element-values can be realized in matrix notation (or `numpy`).

We thus get the vector $\{D_{A,Z}\}$ for the ReLU function.

### 8.2 Sigmoid

**Forward pass: Activation function**

The sigmoid function is defined as:

$$a_i \ \overset{\text{def}}{=} \ \sigma(z_i) \ = \ \frac{1}{1 + e^{-z_i}}$$

where i goes over the rows of the respective vectors. The domain of the sigmoid function is $[-\infty, \infty]$ and its range is $(0, 1.0)$ for any finite input.

Hence, in our matrix notation, the sigmoid function becomes:

$$\{A\} \ \leftarrow \ \{1.0\} \oslash \left( \ \{1.0\} \ \oplus \ \exp\left( \ \{-1.0\} \otimes \{Z\} \ \right) \ \right),$$

where it is understood that the exp function is applied element-wise.

**Backward pass: Gradient of $\{A\}$ with respect to $\{Z\}$**

The partial derivative of $a_i$ with respect to $z_i$ can be worked out as:

$$\frac{\partial a_i}{\partial z_i} = \frac{e^{-z_i}}{\left(1 + e^{-z_i}\right)^2} = a_i\left(1 - a_i\right)$$

Hence, in our matrix notation,

$$\left\{\frac{\partial a_i}{\partial z_i}\right\} \overset{\text{denoted}}{=\!=} \{D_{A,Z}\} \leftarrow \{A\} \circledast \left(\{1.0\} \ominus \{A\}\right)$$

### 8.3 tanh

**Forward pass: Activation function**

The tanh function is defined as:

$$a_i \overset{\text{def}}{=} \tanh\left(z_i\right) = \frac{1 - e^{-2z_i}}{1 + e^{-2z_i}} = \frac{2}{1 + e^{-2z_i}} - 1$$

where $i$ goes over the rows of the respective vectors. The domain of the tanh function is $[-\infty, \infty]$ and its range is $(-1.0, 1.0)$ for any finite input.

Hence, in our matrix notation, the tanh function becomes:

$$\{A\} \leftarrow \left(\{2.0\} \oslash \left(\{1.0\} \oplus \exp\left(\{-2.0\} \circledast \{Z\}\right)\right)\right) \ominus \{1.0\},$$

where it is understood that the $\exp$ function is applied element-wise.

**Backward pass: Gradient of $\{A\}$ with respect to $\{Z\}$**

The partial derivative of $a_i$ with respect to $z_i$ can be derived as:

$$\frac{\partial a_i}{\partial z_i} = 1 - a_i^2$$

Hence, in our matrix notation,

$$\left\{\frac{\partial a_i}{\partial z_i}\right\} \overset{\text{denoted}}{=\!=} \{D_{A,Z}\} \leftarrow \{1.0\} \ominus \left(\{a_i\} \circledast \{a_i\}\right)$$

**Remarks on tanh vs. sigmoid**

Andrew Ng reportedly [5] says that using tanh is almost always preferable to using the sigmoid.

### 8.4 Softmax

The distinctive feature of the softmax function is that the $a_i$ values it produces are such that they all add up to 1.0. Thus:

$$\sum_{i=1}^{N} a_i = 1.0.$$

Hence, we can think of the softmax function as directly providing probabilities for different output classes in the multi-class classification problem that uses logistic regression.

It can be shown [1] that the softmax function is optimal for the maximum-likelihood estimation of the model's parameters (weights and biases).

**Forward pass: Activation**

The softmax function is defined as:

$$a_i \ \overset{\text{def}}{=} \ \frac{e^{z_i}}{\sum\limits_{j=1}^{N} e^{z_j}} \qquad \forall\, i \in [1, N],\ \forall\, j \in [1, N].$$

The domain of the softmax function is $[-\infty, \infty]$ and its range is $(0.0, 1.0)$ for any finite input.

Using our matrix notation,

$$\{Q\} \ \leftarrow \ \exp\big(\ \{Z\}\ \big)$$
$$s \ \leftarrow \ \Sigma_r\big(\ \{Q\}\ \big)$$
$$\{A\} \ \leftarrow \ \{Q\} \oslash \{s\}$$

**Backward pass: Gradients of $\{A\}$ with respect to $\{Z\}$**

Since each $a_i$ is a function of all $z_j$'s, there are $N$ number of partial derivatives for each i-th row, and so, the partial derivatives form an $N \times N$ matrix, i.e. the square-Jacobian.

The specific steps are the following:

$$a_i \ \overset{\text{def}}{=} \ \frac{e^{z_i}}{s} \qquad \forall\, i \in [1, N]$$

$$s \ \overset{\text{def}}{=} \ \sum_{j=1}^{N} e^{z_j} \qquad \forall\, j \in [1, N]$$

Given this definition, when $i \neq j$, then by product and quotient rules of differentiation,

$$\begin{aligned}
\frac{\partial a_i}{\partial z_j} &= e^{z_i} \frac{\partial}{\partial z_j}\left(\frac{1}{s}\right) \ + \ \left(\frac{1}{s}\right)\frac{\partial e^{z_i}}{\partial z_j} \\
&= e^{z_i}\left(\ -\frac{1}{s^2}\, e^{z_j}\ \right) \ + \ \left(\frac{1}{s}\right)(0) \\
&= -\frac{e^{z_i}}{s}\frac{e^{z_j}}{s} \\
&= -a_i\, a_j \\
&= a_i\left(\ 0 - a_j\ \right) \\
&= a_i\left(\ \delta_{ij}\ -\ a_j\ \right),
\end{aligned}$$

31

where $\delta_{ij}$ is Kronecker's delta, defined as:

$$\delta_{ij} \overset{\text{def}}{=} \begin{cases} 1 & \forall\ i = j \\ 0 & \forall\ i \neq j \end{cases} .$$

When $i = j$, by product and quotient rules of differentiation,

$$\begin{aligned}
\frac{\partial a_i}{\partial z_i} &= e^{z_i} \frac{\partial}{\partial z_i}\left(\frac{1}{s}\right) + \left(\frac{1}{s}\right) \frac{\partial e^{z_i}}{\partial z_i} \\
&= e^{z_i}\left(-\frac{1}{s^2} e^{z_i}\right) + \frac{1}{s} e^{z_i} \\
&= -\frac{e^{z_i}}{s}\frac{e^{z_i}}{s} + \frac{e^{z_i}}{s} \\
&= -a_i\, a_i + a_i \\
&= a_i\,(\, 1 - a_i\,) \\
&= a_i\,(\, \delta_{ij} - a_i\,) \\
&= a_i\,(\, \delta_{ij} - a_j\,) \qquad \because i = j
\end{aligned}$$

To express the above component-wise equations to the matrix form, we use eq. (1), which is repeated here:

$$\left[O\right]_{[N \times N]} \overset{\text{def}}{=} \{A\}_{[N \times 1]} \otimes \lfloor 1 \rfloor_{[1 \times N]}$$

To repeat,

$$\left[O\right]_{[N \times N]} = \begin{bmatrix} a_1 & a_1 & \ldots & a_1 \\ a_2 & a_2 & \ldots & a_2 \\ \vdots & \vdots & \ddots & \vdots \\ a_N & a_N & \ldots & a_N \end{bmatrix}_{[N \times N]},$$

and so,

$$\left[O^{\mathrm{T}}\right]_{[N \times N]} \overset{\text{def}}{=} \left(\, \left[O\right]_{[N \times N]}\,\right)^{\mathrm{T}}$$

$$= \begin{bmatrix} a_1 & a_2 & \ldots & a_N \\ a_1 & a_2 & \ldots & a_N \\ \vdots & \vdots & \ddots & \vdots \\ a_1 & a_2 & \ldots & a_N \end{bmatrix}_{[N \times N]}$$

Then, the necessary equation is given by:

$$\left[D_{A,Z}\right]_{[N_A \times N_Z]} \overset{\text{def}}{=} \left[\frac{\partial a_i}{\partial z_j}\right] \tag{7}$$

$$= [O] \circledast \left(\, [\mathcal{I}]_{[N \times N]} \ominus [O^{\mathrm{T}}]\,\right)$$

Note carefully the use of the outer product operator $\otimes$ and the element-wise multiplication operator $\circledast$.

We can then substitute eq. (7) into eq. (6).

### 8.5 Identity

In the forward pass, the identify function simply passes $\{Z\}$ unmodified as $\{A\}$. It thus is a mere placeholder. Its only utility is so that the linear regression problems (simple, polynomial, or multivariate) can be handled within the framework of ANNs. The loss function used in such a case is the SSE (sum of squared errors).

In the backward pass, the identity function simply passes the gradient of the loss vector $\{\mathcal{E}\}$ with respect to $\{Z\}$ unmodified.

## 9  Loss functions and their gradients

The terms "loss", "error", and "cost" are used interchangeably in the context of gradient descent algorithm for artificial neural networks.

### 9.1  Euclidean loss (SSE)

Note that SSE (or HSSE) produces a convex loss-function surface only when the activation function is linear. Thus, SSE is used only for regression problems, not for classification problems.

**Forward pass: Loss function**

The half-sum of squared error for the i-th element is given by:

$$\mathcal{E} \stackrel{\text{def}}{=} \frac{1}{2} \sum_{i=1}^{M} (y_i - a_i)^2$$

Therefore, in our matrix notation,

$$\{d\}_{[N \times 1]} \leftarrow \{Y\}_{[N \times 1]} \ominus \{A\}_{[N \times 1]}$$
$$\{\epsilon\}_{[N \times 1]} \leftarrow \{0.5\}_{[N \times 1]} \circledast \{d\}_{[N \times 1]} \circledast \{d\}_{[N \times 1]}$$
$$\mathcal{E} \leftarrow \Sigma_r \left( \{\epsilon\}_{[N \times 1]} \right)$$

**Backward pass: Gradient of errors with respect to activations**

Given the above-mentioned loss function, the gradient can be worked out as:

$$\frac{\partial \epsilon_i}{\partial a_i} = \left(\frac{1}{2}\right)(2)(y_i - a_i) = y_i - a_i = d_i$$

Therefore, in our matrix notation,

$$\{D_{\mathcal{E},A}\}_{[N \times 1]} \leftarrow \{Y\}_{[N \times 1]} \ominus \{A\}_{[N \times 1]}.$$

### 9.2 Log loss/Cross entropy loss

The log-loss function is used when the activation function is the sigmoid, as in the simple binary classification problems (i.e. the logistic regression).

**Forward pass: Loss function**

$$\mathcal{E} \overset{\text{def}}{=} - \sum_{i=1}^{N} (y_i \log a_i).$$

So, define

$$\epsilon_i \overset{\text{def}}{=} (y_i \log a_i)$$

Then, in our matrix notation,

$$\{\epsilon\}_{[N \times 1]} \leftarrow \{Y\}_{[N \times 1]} \circledast \log \left( \{A\}_{[N \times 1]} \right),$$

where it is understood that the log function is applied element-wise. Then,

$$\mathcal{E} \leftarrow - \Sigma_r \left( \{\epsilon\}_{[N \times 1]} \right),$$

**Backward pass: Gradient of errors with respect to activations**

Given the above-mentioned loss function, the gradient can be worked out as:

$$\frac{\partial \epsilon_i}{\partial a_i} = - \frac{y_i}{a_i}$$

Therefore, in our matrix notation,

$$\{D_{\mathcal{E},A}\}_{[N \times 1]} \leftarrow \left( \{-1.0\} \circledast \{Y\} \right) \oslash \{A\}$$

## 10   Combining layers for simpler tasks

The three layers presented so far form the basic building blocks to construct different kinds of feedforward networks that are designed to perform various tasks. Given below is a brief indication of some of the various combinations that are possible.

### 10.1   Linear regressions

Performing linear regressions with ANNs is, really speaking, only an issue of terminology.

Linear regression involves only the fully connected layer and the loss layer.

**Simple linear regression**

The design to follow is:

- Fully connected layer of size $M = 1$ for the input and $N = 1$ for its output

- No activation layer (or one with the identity function)

- Loss layer using the SSE loss function.

**Multivariable linear regression**

The design to follow is:

- Fully connected layer of size $M > 1$ for the input and $N = 1$ for its output

- No activation layer (or one with the identity function)

- Loss layer using the SSE loss function.

**Multivariate linear regression**

The design to follow is:

- Fully connected layer of the size $M \geq 1$ for the input and of the size $N > 1$ for its output

- No activation layer (or one with the identity function)

- Loss layer using the SSE loss function.

**Polynomial regression**

The model here is:

$$\hat{y} \;=\; \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2$$

Realize that this is an instance of *linear* regressions, and hence, in principle, it can be handled in ANNs.

To make it work in the framework of ANNs, perform some pre-processing by defining new independent variables:

$$\chi_1 \overset{\text{def}}{=} x_1$$
$$\chi_2 \overset{\text{def}}{=} x_2$$
$$\chi_3 \overset{\text{def}}{=} x_1 x_2$$
$$\chi_4 \overset{\text{def}}{=} x_1^2$$
$$\chi_5 \overset{\text{def}}{=} x_2^2$$

Now, treat $\{\chi_i\}$ as the input vector for any of the previously mentioned linear regressions.

## 10.2 Nonlinear regression having $\beta^n$ terms

It is not possible to directly realize an ANN that corresponds to the model:

$$\hat{y} = \beta_0 + \beta_1^2 x$$

The reason is that the fully connected layer, as presented, performs only linear combinations. Therefore, implementing a nonlinear model of this kind would involve using a special-purpose layer in place of the fully connected layer. The loss layer can be used as is.

## 10.3 Regression for classification

This, again, is an issue only of terminology.

### Logistic regression (binary classification)

We take the term "logistic regression" to mean binary classification using the sigmoid function. Accordingly, the design to follow is:

- Fully connected layer of size $M = 1$ for the input, and $N = 1$ for its output

- Activation layer of size $N = 1$ using the sigmoid function

- Loss layer of size $N = 1$ using the binary log-loss function.

### Multinomial regression (multi-class classification)

The design to follow is:

- Fully connected layer of size $M \geq 1$ for the input, and $N > 1$ for its output

- Activation layer of size $N$ using the sigmoid function

- Loss layer of size $N$ using the multi-class log-loss function.

# 11 Simple deep learning networks

The term "deep learning" is typically used for denoting advanced architectures like the convolutional neural networks (ConvNNs or CNNS), and the recurrent neural networks (RNNs). The layers presented in this document are not sufficient to build such advanced architectures. Mind you, they do form parts of CNNs and RNNs too. However, only these layers, by themselves, are not sufficient.

However, technically speaking, any neural network that has hidden layers also is called a deep network.

Using the terminology of layers we introduced, a hidden layer of a deep network consists of a *pair* of a fully connected and an activation layer.

Accordingly, the design to follow is:

- Input stage ("layer")

  - One fully connected layer of the size $M$

  - Its corresponding activation layer of some arbitrary size $H_1^I$

- One or more hidden stages ("layers"), *each* consisting of

  - One fully connected layer with size of its input vector the same as the activation vector of the previous stage (which, for the first hidden stage is $H_1^I$)

  - One activation layer of some arbitrary size $H_1^O$.

  When multiple hidden layers are used, the size of the input for the next hidden layer $H_{n+1}^I$ is, of course, the same as the size of the output for the previous hidden layer output $H_n^O$.

- Output stage ("layer")

  - Just one loss layer. It takes two inputs: (i) the output of the last hidden layer $H_{last}^O$, and (ii) the target variables vector $\{Y\}$

During mini-batch processing, for each $r$-th row of a mini-batch, all the layers are traversed through during the forward and backward passes. The gradients are accumulated in each fully connected layer, until the mini-batch gets over. The weights and biases in each fully connected layer get updated at the end of the mini-batch.

## References

[1] Nielsen, Michael (2015) "Neural networks and deep learning", Determination Press. https://neuralnetworksanddeeplearning.com/

[2] Mazur, Matt (2015) "A step by step backpropagation example". https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/

[3] Mallya, Arun (2017) "Writeups/Notes". https://arunmallya.github.io

[4] Online discussion (2015) "How to implement the ReLU function in Numpy", StackOverFlow, https://stackoverflow.com/questions/32109319/how-to-implement-the-relu-function-in-numpy.

[5] Online discussion (2018) "Why is tanh almost always better than sigmoid as an activation function?", Stats.StackExchange, https://stats.stackexchange.com/questions/330559/why-is-tanh-almost-always-better-than-sigmoid-as-an-activation-function