

INSTITUTO SUPERIOR TÉCNICO

DEEP LEARNING

---

## Homework 1

---

*Group 28:*

Luís FREIRE D'ANDRADE (*ist194179*)

Rita DUARTE (*ist195531*)

December 23, 2022



## Work Distribution

The questions 1.2.a), 1.2.b), 2.2 and 2.3 were assigned to Luís D'Andrade.

The questions 1.1.a), 1.1.b) and 2.1 were assigned to Rita Duarte.

Both the colleagues worked together to solve the 3.1, 3.2 and 3.4 questions.

## References

- Theoretical slides
- Jupyter notebooks from practical 2, 3 and 5. And also annotations taken from the practical classes
- Piazza forum of the Deep Learning subject
- Aparentamentos Álgebra Linear IST

## Question 1

### 1.1.a)

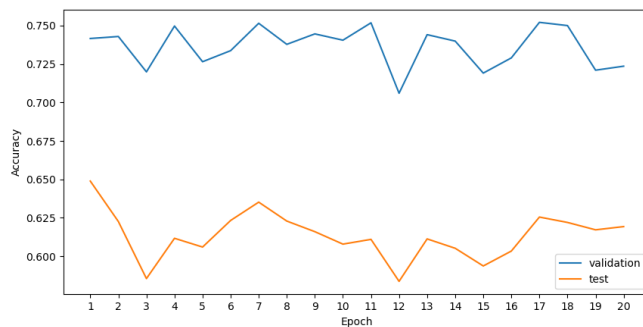


Figure 1: Plotting of the accuracy for each of the 20 trained epochs by the implemented perceptron

### 1.1.b)

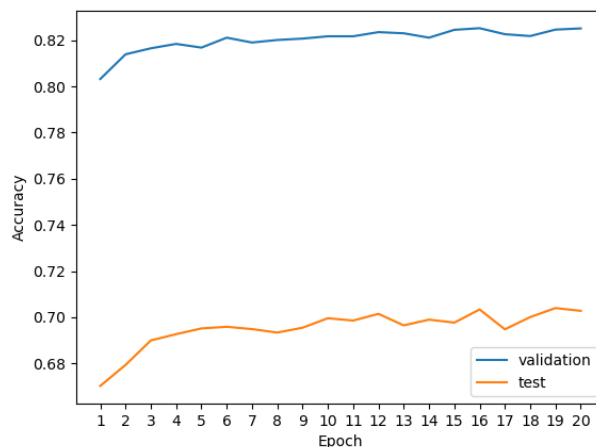


Figure 2: Plotting of the accuracy for each of the 20 trained epochs by the implemented logistic regression with learning rate of 0.001

### 1.2.a)

Unlike the the Simple Single Layer Perceptron (SLP) implemented above, a Multi Layer Perceptron (MLP) contains one or more hidden layers (apart from one input and one output layer). While in a SLP the neuron must have a linear activation function that imposes a threshold, neurons in a MLP can use any arbitrary (including non-linear) activation function.

This conveys that a MLP can learn non-linear target functions, i.e. learn a non-linear relationship between the inputs and outputs. If a neural network consists of only linear activation functions (like a SLP), it can only model a linear relationship between the inputs and outputs, which makes it not as expressive or useful in a lot of applications.

1.2.b)

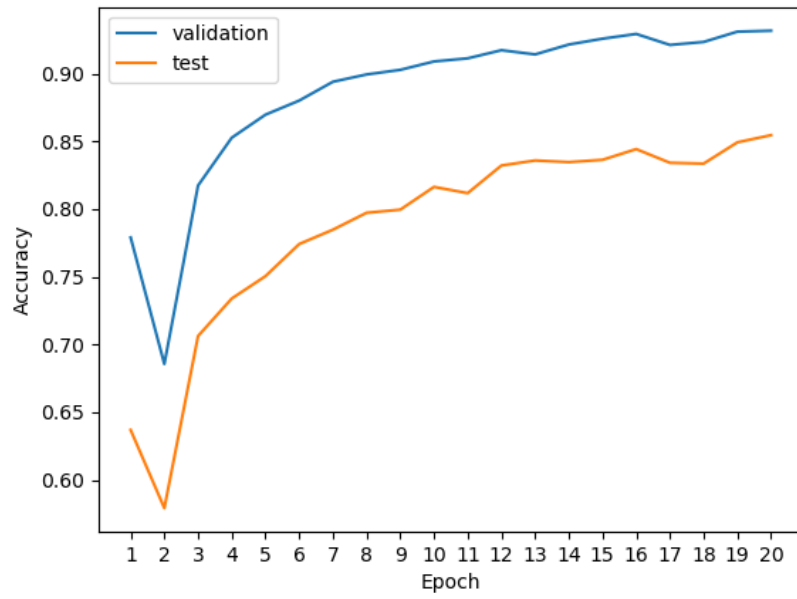


Figure 3: Plotting of the accuracy for each of the 20 trained epochs by a multi-layer perceptron with a single hidden layer

Number of epochs	20
Hidden Units	200
Hidden Layers	Activation function <b>ReLU</b>
Output layer	Multinomial logistic loss
Training algorithm	Stochastic gradient descent
Learning rate	0.001
Bias	Zero vectors
Weights	$w_{i,j} \sim \mathcal{N}(\mu, \sigma^2)$ with $\mu = 0.1$ $\sigma^2 = 0.1^2$

## Question 2

### 2.1.

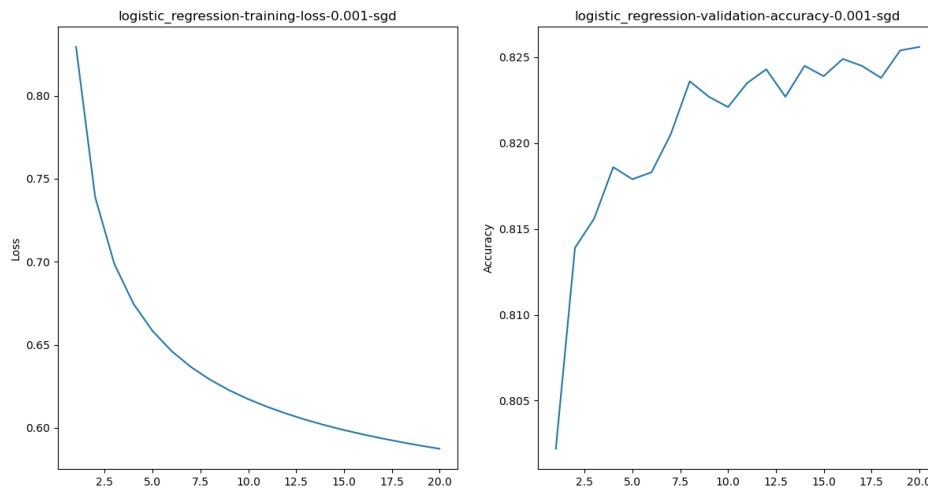


Figure 4: Plotting of the accuracy and training loss for each of the 20 trained epochs by a linear model with logistic regression

Final test accuracy: 0.7019

Number of epochs	20
Hidden Units	200
Batch size	1
Optimization algorithm	Stochastic gradient descent
Learning rate	0.001

## 2.2.

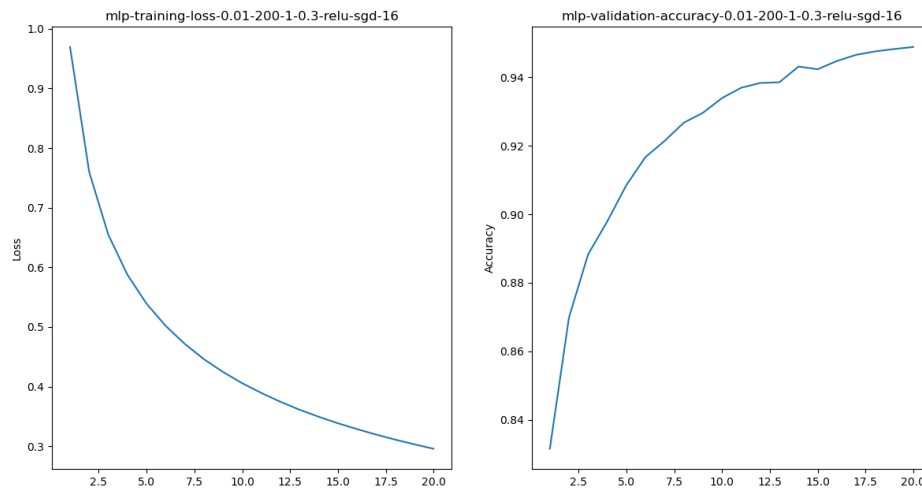


Figure 5: Plotting of the accuracy and training loss for each of the 20 trained epochs by feed-forward neural network using a dropout regularization

Final test accuracy:0.8818

Number of epochs	20
Number of Layers	1
Learning Rate	0.01
Hidden Size	200
Dropout	0.3
Batch Size	16
Activation	ReLU
Optimizer	SGD

### 2.3.

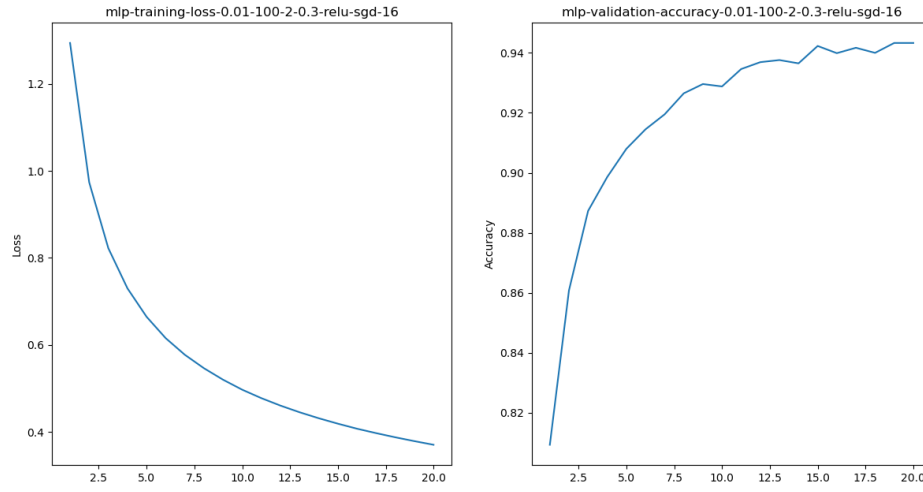


Figure 6: Plotting of the accuracy and training loss for each of the 20 trained epochs by feed-forward neural network using a dropout regularization

Final test accuracy: 0.8724

Number of epochs	20
Number of Layers	2
Learning Rate	0.01
Hidden Size	100
Dropout	0.3
Batch Size	16
Activation	ReLU
Optimizer	SGD

### Question 3

#### 3.1.

For  $h$  to be a linear transformation of  $\Phi(x)$  we need to find a feature transformation  $\Phi(x)$  and a matrix  $A_\Theta$  such that  $h = A_\Theta \Phi(x)$ .

Since  $h = g(Wx), x \in \mathbb{R}^D$  and the activation function is  $g(z) = z^2$  then:

$$h = g(Wx) \Leftrightarrow h = (Wx)^2$$

Also  $\Theta = (W, v) \in \mathbb{R}^{k \times D} \times \mathbb{R}^k$  hence:

$$W = \begin{bmatrix} W_{11} & \dots & W_{1D} \\ \vdots & \dots & \vdots \\ W_{k1} & \dots & W_{kD} \end{bmatrix}$$

$$\text{and } x = \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_D \end{bmatrix}$$

Now,

$$h = (Wx)^2 \Leftrightarrow h = \left( \begin{bmatrix} W_{11} & \dots & W_{1D} \\ \vdots & \dots & \vdots \\ W_{k1} & \dots & W_{kD} \end{bmatrix} \times \begin{bmatrix} x_1 \\ \vdots \\ \vdots \\ x_D \end{bmatrix} \right)^2 \Leftrightarrow \left( \begin{bmatrix} \sum_{i=1}^D W_{1i} \times x_i \\ \vdots \\ \vdots \\ \sum_{i=1}^D W_{ki} \times x_i \end{bmatrix} \right)^2$$

From the Square of the Sum:

$$(\sum_{i=1}^N a_i)^2 = \sum_{i=1}^N a_i^2 + 2 \times \sum_{j=1}^N \sum_{i=1}^{j-1} a_i a_j$$

We can expand the last matrix into,

$$\Leftrightarrow \begin{bmatrix} \sum_{i=1}^D (W_{1i} x_i)^2 + 2 \times \sum_{j=1}^D \sum_{i=1}^{j-1} (W_{1i} x_i)(W_{1j} x_j) \\ \vdots \\ \vdots \\ \sum_{i=1}^D (W_{ki} x_i)^2 + 2 \times \sum_{j=1}^D \sum_{i=1}^{j-1} (W_{ki} x_i)(W_{kj} x_j) \end{bmatrix}$$

In the description of the exercise 3.1),  $\Phi : \mathbb{R}^D \rightarrow \mathbb{R}^{\frac{D(D+1)}{2}}$  and  $A \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$ .

As stated before  $h = A_\Theta \Phi(x)$  and the product of the matrix  $A_\Theta$  with size  $k \times \frac{D(D+1)}{2}$  by the vector  $\Phi(x)$  with size  $\frac{D(D+1)}{2} \times 1$  will result in a vector of size  $k \times 1$ :

$$h = A_\Theta \Phi(x) = \begin{bmatrix} a_{11} & \dots & a_{1 \frac{D(D+1)}{2}} \\ \vdots & \ddots & \vdots \\ a_{k1} & \dots & a_{k \frac{D(D+1)}{2}} \end{bmatrix} \begin{bmatrix} b_{11} \\ \vdots \\ b_{k1} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^{\frac{D(D+1)}{2}} a_{1i} b_{1i} \\ \vdots \\ \sum_{i=1}^{\frac{D(D+1)}{2}} a_{ki} b_{ki} \end{bmatrix}$$

By only looking at line  $k$  of the matrix  $h$ ,

$$\sum_{i=1}^D (W_{ki} x_i)^2 + 2 \times \sum_{j=1}^D \sum_{i=1}^{j-1} (W_{ki} x_i)(W_{kj} x_j) \text{ needs to be written as } \sum_{i=1}^{\frac{D(D+1)}{2}} a_{ki} b_{ki}$$



$$\begin{aligned}
& \sum_{i=1}^D (W_{ki}x_i)^2 + 2 \times \sum_{j=1}^D \sum_{i=1}^{j-1} (W_{ki}x_i)(W_{kj}x_j) \Leftrightarrow \\
& \Leftrightarrow \sum_{i=1}^D W_{ki}^2 x_i^2 + 2 \times \sum_{j=1}^D \sum_{i=1}^{j-1} (W_{ki}x_i)(W_{kj}x_j) \Leftrightarrow \\
& \Leftrightarrow (W_{k1}^2 x_1^2 + W_{k2}^2 x_2^2 + \dots + W_{kD}^2 x_D^2) + (2W_{k1}W_{k2}x_1x_2 + 2W_{k1}W_{k3}x_1x_3 + \dots + \\
& \quad 2W_{k(j-1)}W_{kD}x_{(j-1)}x_D)
\end{aligned}$$

For each  $n \in [1, D]$ , by putting

$x_n$  in evidence

$W_{kn}$  in evidence

$$\begin{aligned}
& x_n^2 \times (x_n x_{n+1} \times x_n x_{n+2} \times x_n x_{n+3} \times \dots \times x_n x_D) \\
& \text{and} \\
& W_{in}^2 \times 2W_{in}W_{i(n+1)} \times 2W_{in}W_{i(n+2)} \times 2W_{in}W_{i(n+3)} \times \dots \times 2W_{in}W_{iD}
\end{aligned}$$

By analysing the behaviour of  $h$  in line  $= k$  we can detect a pattern and isolate the  $x$  from the  $W$  in order to write  $h$  as the product of matrices  $A_\Theta$  and  $\Phi(x)$ ,

$$\begin{aligned}
h &= [A_\Theta] [\Phi(x)] \Leftrightarrow \\
& \Leftrightarrow h = \begin{bmatrix} W_{11}^2 & \dots & W_{1D}^2 & 2W_{11}W_{12} & \dots & 2W_{1(D-1)}W_{1D} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ W_{k1}^2 & \dots & W_{kD}^2 & 2W_{k1}W_{k2} & \dots & 2W_{k(D-1)}W_{kD} \end{bmatrix} \begin{bmatrix} x_1^2 \\ \vdots \\ \vdots \\ x_D^2 \\ x_1x_2 \\ \vdots \\ \vdots \\ x_{D-1}x_D \end{bmatrix}
\end{aligned}$$

Hence the mapping  $\Phi$  and the matrix  $A_\Theta$  are given by:

$$\begin{aligned}
\Phi(x) &= [x_1^2 x_1 x_2 \dots x_1 x_D \quad x_2^2 x_2 x_3 \dots x_2 x_D \quad \dots \quad x_D^2]^T \\
A_\Theta &= [W_{i1}^2 2W_{i1}W_{i2} \dots 2W_{i1}W_{iD} \quad W_{i2}^2 2W_{i2}W_{i3} \dots 2W_{i2}W_{iD} \quad \dots \quad W_{iD}^2]
\end{aligned}$$

### 3.2.

In the description of exercise 3.2) there is information relevant for the solving of the problem such as:

1.  $\hat{y} = v^T h \in \mathbb{R}$  hence  $\hat{y} \in \mathbb{R}$
2.  $\Theta = (W, v) \in \mathbb{R}^{K \times D} \times \mathbb{R}^K$  hence  $v \in \mathbb{R}^k$ .
3.  $h \in \mathbb{R}^K$
4.  $A_\Theta \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$
5.  $\Phi(x) \in \mathbb{R}^{\frac{D(D+1)}{2}}$

Based on the previous claim  $h = A_\Theta \Phi(x)$ ,

$$\hat{y} = v^T h = v^T A_\Theta \Phi(x)$$

The hypothesis is we can write  $\hat{y}(x; c_\Theta) = c_\Theta^T \Phi(x)$ ,

$$\hat{y} = v^T A_\Theta \Phi(x) = c_\Theta^T \Phi(x),$$

where  $v^T$  is a vector size  $1 \times k$ ,  $A_\Theta$  is a matrix size  $k \times \frac{D(D+1)}{2}$ :  
 $c_\Theta^T$  is size  $1 \times \frac{D(D+1)}{2}$ .

$$c_\Theta^T = v^T A_\Theta \Leftrightarrow c_\Theta^T = [v_1 \quad \dots \quad v_k] \begin{bmatrix} - & A_{\Theta 1} & - \\ - & \dots & - \\ - & A_{\Theta k} & - \end{bmatrix} \Leftrightarrow$$

Having into consideration the following transpose properties for matrices:

$(A^T)^T = A$ , let A be a matrix.

$(A \times B)^T = B^T \times A^T$ , let A and B be matrices.

Then:

$$\begin{aligned} \Leftrightarrow (c_\Theta^T)^T &= \left( [v_1 \quad \dots \quad v_k] \begin{bmatrix} - & A_{\Theta 1} & - \\ - & \dots & - \\ - & A_{\Theta k} & - \end{bmatrix} \right)^T \Leftrightarrow \\ \Leftrightarrow c_\Theta &= [v_1 \quad \dots \quad v_k]^T \begin{bmatrix} - & A_{\Theta 1} & - \\ - & \dots & - \\ - & A_{\Theta k} & - \end{bmatrix}^T \Leftrightarrow \\ \Leftrightarrow c_\Theta &= \begin{bmatrix} | & \dots & | \\ A_{\Theta 1} & \dots & A_{\Theta k} \\ | & \dots & | \end{bmatrix} \begin{bmatrix} v_1 \\ \dots \\ v_k \end{bmatrix}, \end{aligned}$$

$c_\Theta$  corresponds to the product of a matrix  $A_\Theta^T$  of size  $\frac{D(D+1)}{2} \times k$  and a vector  $v$  of size  $k \times 1$ .

Therefore  $c_\Theta$  is size  $\frac{D(D+1)}{2} \times 1$  and we can write the matrix as:

$$c_{\Theta i} = \sum_{j=1}^k A_{\Theta ij} v_j, \text{ } i \text{ of size } \frac{D(D+1)}{2}$$

Therefore it is proven that,

$$\hat{y} = (v^T A_\Theta) \Phi(x) \Leftrightarrow \hat{y} = c_\Theta^T \Phi(x), \text{ for some } c_\Theta \in \mathbb{R}^{K \times \frac{D(D+1)}{2}}$$

Linear regression, as confirmed, can be used to learn this model.

Nevertheless, this is still a non-linear model in terms of the original parameters  $\Theta$  because its activation function is non-linear:  $g(z) = z^2$ .

### 3.4.

Since  $\hat{y} = c_\Theta^T \Phi(x)$  is linear, it's possible to minimize the squared loss,  $L(c_\Theta; \mathcal{D})$ , and find a closed form solution for  $\hat{c}_\Theta$ . To do so we need to solve,

$$\hat{c}_\Theta = \arg \min_{c_\Theta} \frac{1}{2} \sum_{n=1}^N (\hat{y}(x_n; c_\Theta) - y_n)^2$$

The factor  $\frac{1}{2}$  is irrelevant and the subtraction can be commuted, as it's squared,

$$\hat{c}_\Theta = \arg \min_{c_\Theta} \sum_{n=1}^N (y_n - \hat{y}(x_n; c_\Theta))^2 = \arg \min_{c_\Theta} \sum_{n=1}^N (y_n - c_\Theta^T \Phi(x_n))^2$$

Recalling the Euclidean/Frobenius norm formula,

$$\sum_{n=1}^N (y_n - c_\Theta^T \Phi(x_n))^2 = \|y - X c_\Theta\|^2,$$

$$\text{with } X = \begin{bmatrix} \Phi(x_1)^T \\ \dots \\ \Phi(x_N)^T \end{bmatrix} \in \mathbb{R}^{N \times \frac{D(D+1)}{2}}, \quad y = \begin{bmatrix} y_1 \\ \dots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

Writing the optimality condition, when the gradient equals zero,

$$\begin{aligned} 0 &= \nabla_{c_\Theta} \|y - X c_\Theta\|^2 \\ &= \nabla_{c_\Theta} (c_\Theta^T X^T X c_\Theta - 2 c_\Theta^T X^T y + \|y\|^2) \\ &= 2 X^T X c_\Theta - 2 X^T y \end{aligned}$$

Which results in,

$$\hat{c}_\Theta = (X^T X)^{-1} X^T y$$

Feedforward neural network training is a special case of function minimisation, and many of its optimization algorithms can easily get stuck in non-global minima, which makes global minimization usually intractable. However, in our case, we're dealing with a squared loss (a convex optimization problem) and a polynomial activation function. Local minimizers of a convex optimization problem are also global minimizers, which makes convex formulations a preferable special case.