



Universitat Autònoma de Barcelona

Degree Thesis

---

# IMPROVEMENTS OF DETERMINISTIC PROCESSES THROUGH NEURAL NETWORKS

---

Author:  
Gerard Lahuerta Martín

Supervisor:  
Dr. Lluís Alsedà Soler

A THESIS SUBMITTED IN FULFILMENT OF THE  
REQUIREMENTS FOR THE DEGREE OF COMPUTATIONAL  
MATHEMATICS AND DATA ANALYTICS IN THE

SCIENCE FACULTY

June 2024



## Declaration of Authorship

## Abstract

## Acknowledgements

I would like to express my gratitude to my family for their support throughout my thesis, degree, and my entire life.

I am especially grateful to my mother, Maria Montserrat, for her encouragement and guidance in facing the challenges of this thesis and my career.

Thanks to all the peoples that treat me as I was part of their family.

I also appreciate the encouragement and support of my colleagues.

Finally, I want to give a special mention to Dr. Lluís Alsedà for guiding me through the thesis and being an inspiration.

Thank you all for everything.

# Contents

<b>Declaration of Authorship</b>	<b>2</b>
<b>Abstract</b>	<b>3</b>
<b>Acknowledgements</b>	<b>4</b>
<b>Preface</b>	<b>6</b>
<b>Introducción</b>	<b>7</b>
<b>1 Neural Networks</b>	<b>8</b>
1.1 Multilayer Perceptron and Perceptron neuron . . . . .	8
1.2 Why Neural Networks works . . . . .	9
1.3 Cost Function . . . . .	10
1.4 BackPropagation . . . . .	11

## Preface

## Introducción



# 1 Neural Networks

Before starting programming and testing a Neural Network is necessary to understand how Neural Networks work.

A neural Network is made of individual and independent elements connected between them, passing and managing the information through the network formed.

In this thesis we will focus on one of the simplest networks, a multilayer perceptron, to test the different methods of optimization.

## 1.1 Multilayer Perceptron and Perceptron neuron

One of the simple Neural Networks to analyse and program is the Multilayer Perceptron.

It was first proposed by Frank Rosenblat<sup>1</sup> in 1958 (nevertheless its approach did not learn either produce accurate results).

This Neural Network is formed by elements (the artificial neurons) called *Perceptrons* (which gives its name to this network). This neuron is formed by input, weight and activation functions.

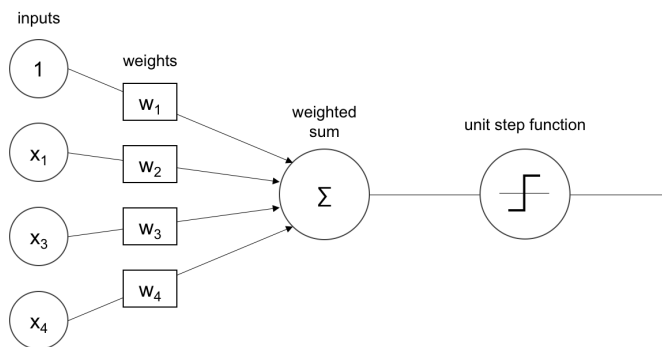


Figure 1: Schema of the Perceptron neuron

The number of inputs that the Perceptron receives is variable but at least needs one input and a bias (represented in the picture as the input with value 1).

All the inputs are escalated by a factor (the weights of the Perceptron) and then summed together

Finally, the result of this sum is introduced in the activation function which returns the output of the neuron.

The traditional activation function used in the Multilayer Perceptron is the Sigmoid:

$$f(x) = \frac{1}{1 + e^{-w \cdot x}}, \text{ where: } x, w \in \mathbb{R}^n$$

The objective of the activation function is to obtain how relevant the result of the weighted sum is.

The use of the sigmoid function as the activation function allows to represent the output of the perception as a probability (because the image of the sigmoid function is  $[0, 1]$ ).

<sup>1</sup>Frank Rosenblat, psychologist and father of deep learning, check its [biography](#).

The Multilayer Perceptron is divided into layers with neurons, we can identify three types of layers:

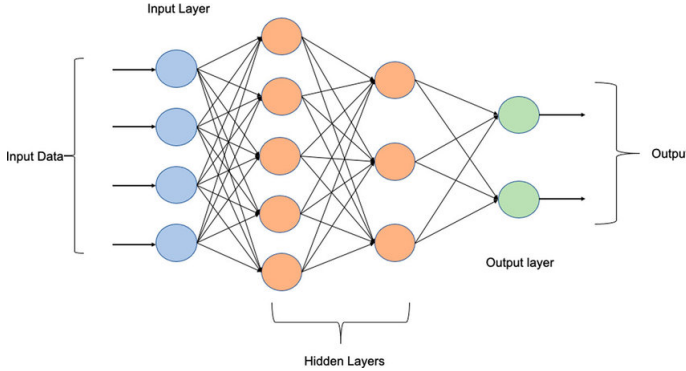


Figure 2: Schema of the Multilayer perceptron

- **Input Layer:** The initial set of neurons of the Multilayer Perceptron.
- **Output Layer:** The final set of neurons of the Multilayer Perceptron.
- **Hidden Layers:** The set of neurons (in layers) between the input and output layers.

The number of hidden layers depends on the problem and the number of neurons in each layer could not be the same.

Moreover, mention some vocabulary used in the Thesis in relation to the Neural Networks architectures:

- **Fully connected:** Architecture with all the neurons connected between layers.
- **Deep:** Architecture where the number of layers is huge to perform complex regression or classification tasks.

Therefore, we can interpret a Multilayer Perceptron as a weighted sum of activation functions.

But, why a weighted sum of functions can perform complex tasks like classification or regression?

## 1.2 Why Neural Networks works

Let's imagine that we want to perform regression (for example) on a set of points that exhibit a certain trend, such as a polynomial trend.

In this scenario, it's important to note that any function can be approximated by an infinite sum of sigmoids:

$$f(x) \approx \sum_{i=0}^N a_i \sigma(w_i^T \cdot x + b_i) \quad (1)$$

In this context, the function  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  is a continuous function to be approximated,  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is the sigmoid function,  $w_i$  represents the weights of the variable  $x \in \mathbb{R}^n$ ,  $b_i$  represents the bias, and the parameter  $a_i$  is a vector of size  $(m, 1)$  containing the weights of the sigmoid function.

This approximation is known as the **Universal Approximation Theorem**, and it enables the Multilayer Perceptron to perform both regression and classification tasks.

### 1.3 Cost Function

The most relevant application of the **Universal Approximation Theorem** in the field of Neural Networks is that it allows us to approximate cost functions, which are functions that evaluate the errors made by the model.

This approximation enables the model to learn by minimizing the error.

The most used cost functions are:

1. **Mean Squared Error:**

The Mean Squared Error (*MSE*) is used in regression tasks. It estimates the error of the model. Its expression is:

$$MSE = \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad (2)$$

Where  $\hat{y}_i := f(x_i)$  represents the predicted value of the dependent variable  $y_i$  to be estimated, and  $n$  is the number of samples in the dataset, where  $n \in \mathbb{N}$ .

2. **Cross-Entropy Loss:**

The Cross-Entropy Loss (*CEL*) is used in multi-class classification tasks. It estimates the accuracy of the model. Its expression is:

$$CEL = -\frac{1}{n} \sum_{i=1}^n (y_i \cdot \log(\hat{y}_i)), \text{ where } \hat{y}_i := f(x_i), n \in \mathbb{N} \quad (3)$$

Where  $\hat{y}_i$  is the predicted probability of the  $i$ -th data classified as the class  $y_i$  and  $n \in \mathbb{N}$  the number of samples.

Despite the existence of these two cost functions, there are more complex functions that are widely used, such as the **Mean Absolute Error**.

This minimization process (despite the cost function) is used iteratively to determine the values of the parameters in Neural Networks, regardless of the type of task or network.

The use of iterative optimizers to obtain the best parameters for a Neural Network arises due to the complexity of finding the coefficients that minimize the cost function.

Moreover, in some cases, it's possible that such coefficients may not exist or cannot be obtained analytically.

This is why it's important to use differentiable cost functions in training neural networks, as they allow us to leverage the gradient of the cost function to obtain the optimal set of parameters. However, we will discuss this topic in more detail shortly.

## 1.4 BackPropagation

The importance of using differentiable cost functions lies in their compatibility with the BackPropagation training algorithm. This algorithm relies on the derivatives of the cost function to compute the gradients for each neuron, which are then used in gradient descent optimization.

One widely used example, often employed for educational purposes to practice and understand the algorithm, is the following network:

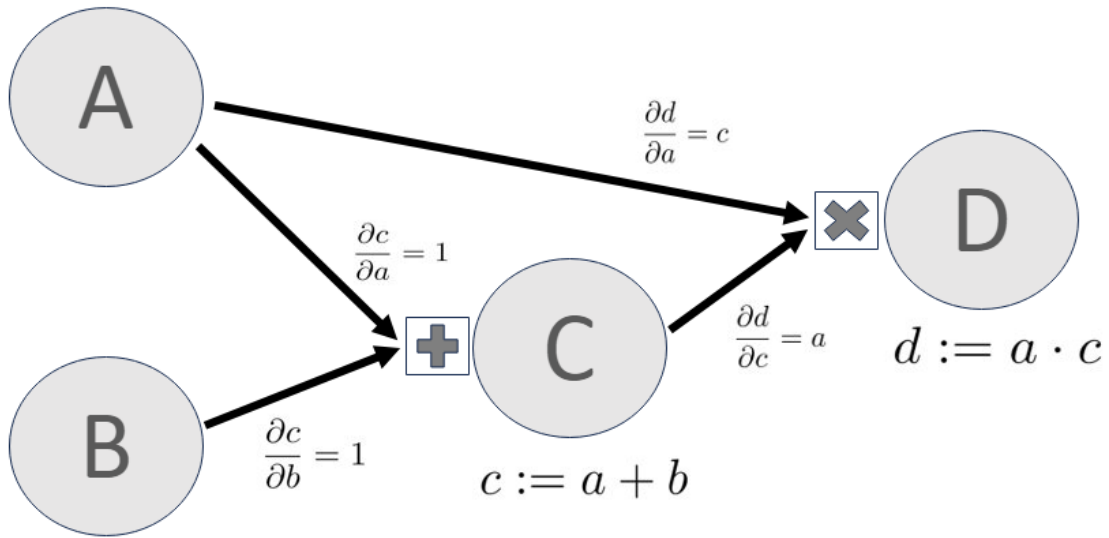


Figure 3: Example of BackPropagation algorithm in simple graph

In the example shown in Figure 3, if we want to modify the value of the parameter  $a$  based on the output  $d$ , we can apply the chain rule of derivatives, as demonstrated in this case:

$$\frac{\partial D}{\partial a} = \frac{\partial d}{\partial a} + \frac{\partial d}{\partial c} \cdot \frac{\partial c}{\partial a} = c + a = 2a + b$$

This algorithm solves one of the biggest problems in early neural network studies related to their development and training. In practice, when a neural network is used, regardless of the task, the only information available is the input and the output.

This aspect renders the usage of neural networks with more than 2 layers (with more than 1 hidden layer) unviable.

The solution to this problem was the **Back-Propagation** algorithm, which "propagates the effects backwards", through each neuron based on its relevance to the obtained output.

This solution was first introduced by Frank Rosenblatt in 1962 for training an *MLP* (Multi-Layer Perceptron). It was subsequently adopted and used in the training process of various types of neural networks due to its effectiveness, allowing the network to optimize the parameters of each neuron to minimize the error, as represented by the cost function, through **gradient descent**.

## References

- [1] WIKIPEDIA, *Multilayer perceptron*,  
[https://en.wikipedia.org/wiki/Multilayer\\_perceptron](https://en.wikipedia.org/wiki/Multilayer_perceptron)
- [2] WIKIPEDIA, *Perceptrón*,  
<https://es.wikipedia.org/wiki/Perceptr%C3%B3n>
- [3] WIKIPEDIA, *Activation function*,  
[https://en.wikipedia.org/wiki/Activation\\_function](https://en.wikipedia.org/wiki/Activation_function)
- [4] WIKIPEDIA, *Frank Rosenblatt*,  
[https://es.wikipedia.org/wiki/Frank\\_Rosenblatt](https://es.wikipedia.org/wiki/Frank_Rosenblatt)
- [5] WIKIPEDIA, *Activating Function*,  
[https://en.wikipedia.org/wiki/Activating\\_function](https://en.wikipedia.org/wiki/Activating_function)
- [6] WIKIPEDIA, *Artificial neural network*,  
[https://en.wikipedia.org/wiki/Artificial\\_neural\\_network](https://en.wikipedia.org/wiki/Artificial_neural_network)
- [7] WIKIPEDIA, *Neural network*,  
[https://en.wikipedia.org/wiki/Neural\\_network](https://en.wikipedia.org/wiki/Neural_network)
- [8] WIKIPEDIA, *Universal Approximation Theorem*,  
[https://en.wikipedia.org/wiki/Universal\\_approximation\\_theorem](https://en.wikipedia.org/wiki/Universal_approximation_theorem)
- [9] WIKIPEDIA, *Mean Squared Error*,  
[https://en.wikipedia.org/wiki/Mean\\_squared\\_error](https://en.wikipedia.org/wiki/Mean_squared_error)
- [10] WIKIPEDIA, *Mean Absolute Error*,  
[https://en.wikipedia.org/wiki/Mean\\_absolute\\_error](https://en.wikipedia.org/wiki/Mean_absolute_error)
- [11] WIKIPEDIA, *Cross-Entropy*,  
<https://en.wikipedia.org/wiki/Cross-entropy>
- [12] WIKIPEDIA, *BackPropagation*,  
<https://en.wikipedia.org/wiki/Backpropagation>
- [13] DANTZIG, G.B. y P. WOLFE, «Decomposition principle for linear programs», *Operations Research*, **8**, págs. 101–111, 1960.