



Universitat Autònoma de Barcelona

Degree Thesis

IMPROVEMENTS OF DETERMINISTIC PROCESSES THROUGH NEURAL NETWORKS

Author:
Gerard Lahuerta Martín

Supervisor:
Dr. Lluís Alseda 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

I, *Gerard Lahuerta Martín*, declare that this thesis titled, "**Improvements of Deterministic Processes through Neural Networks**" and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a bachelor degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed: Gerard.

Date: 21-06-2024

Abstract

This study in the context of neural networks and computational cost gives a new perspective of a possible method that reduces time-cost of deterministic processes.

Its main objective is to sample enough evidence for which deterministic processes can be modelled with neural networks, allowing in cases where this deterministic processes are computationally expensive come by an approximated value within a fixed time.

This study is performed using the scientific method of: *Observation, Hypothesis, Experimentation, Conclusions*.

It develops some working hypotheses that could not be proven but if future investigations verify it, these working hypotheses gives a better understanding of neural networks and how their architecture can be planned.

The results are not as relevant as should be to ensure the viability of the Working Hypothesis proposit but they give a starting point for future works in the area.

Key words: *neural networks, multilayer perceptron, deterministic process, Universal Approximation Theorem, distribution of roman numbers.*

Acknowledgements

I want to thank my family for their support throughout my thesis, degree, and 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 everyone who treated me as if 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

Introduction	2
1 Neural Networks	3
1.1 Multilayer Perceptron and Perceptron neuron	3
1.2 Why Neural Networks work	4
2 Neural Networks topology	5
2.1 Hypotesis and Working Hypothesis	5
2.2 Study of deterministic processes	6
2.2.1 Study of the <i>sin</i> function	6
2.2.2 Study of the <i>parabolic</i> function	8
2.2.3 Study of the <i>squared root</i> function	9
2.3 Viability of approximation discontinuous functions	9
2.3.1 Study of Jump discontinuity	10
2.3.2 Study of Asimptotyc discontinuity	10
2.4 Observations of the results	11
3 Arabic to Roman numbers	13
3.1 Function to approximate	13
3.2 Preprocess of data and reduction of variables	14
3.3 Analisis of the dataset and creation of the neural network	15
3.4 Efficiency and time-cost results	17
4 Conclusions	18

List of Figures

1	Schema of the Perceptron neuron	3
2	Schema of the Multilayer perceptron	3
3	Results of different architectures approximating the function <i>sin</i>	7
4	Architectures that improves the modelisation of <i>sin</i> function	7
5	Results of different architectures approximating the functions x^2 and x^3	8
6	Results of different architectures approximating the function \sqrt{x}	9
7	Architecture 1 – 2 – 2 – 1 approximation to a jump discontinuity function	10
8	Architecture 1 – 5 – 5 – 1 approximation to $\tan x$ function	11
9	Distribution of the associated Arabic number to the Roman natural number	15
10	Working Hypotesis architecture of the neural network that approximates the Roman number distribution	15
11	Test of different neural networks to approach the distribution of Roman numbers	16
12	Improvement of sub-network capabilities to approach the roman numbers distribution	16
13	Final architecture proposed to approach the distribution of roman numbers	17

List of Tables

1	Analysis of cost and efficiency about the different architectures develop	17
---	---	----

Preface

This report is an ambitious project that has been developing for months and has been meditated even more.

Every result obtained remarks the viability of neural network development and how it can be implemented to solve problems that have high computational costs.

All the resources used can be consulted in the [GitHub repository](#).

Introduction

This work is a starting point for other investigations about neural networks in deterministic processes and how they can be developed.

It is known that a fixed series of rules which associates two groups can be explained as an application between these sets. This application is unknown and in some cases cannot be described at all. Nevertheless, since the deterministic process exists, it is possible to create some registers to train a neural network.

This neural network, thanks to the **Universal Approximation Theorem**, can mimic the application, with enough time and data. Therefore, it can approach the results of the deterministic process.

In cases where it is not possible to know how much time a deterministic process could take to obtain a solution, neural networks can perform an approximated solution within a fixed number of calculations.

In this study, observations and a *Working Hypothesis* about neural networks will be tested: It is possible to recreate a neural network with the ability to mimic any deterministic process?

To solve this hypothesis, it is necessary to test the limits of the **Universal Approximation Theorem**.

Therefore, it is compulsory to test the results and observations throughout a practical problem. While solving this problem enough evidence will be recapped to ensure that the proposed *Working Hypothesis* can be advantageous.

1 Neural Networks

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 is the Multilayer Perceptron¹.

It was first proposed by Frank Rosenblat² in 1958 (nevertheless its approach did not learn either produce accurate results).

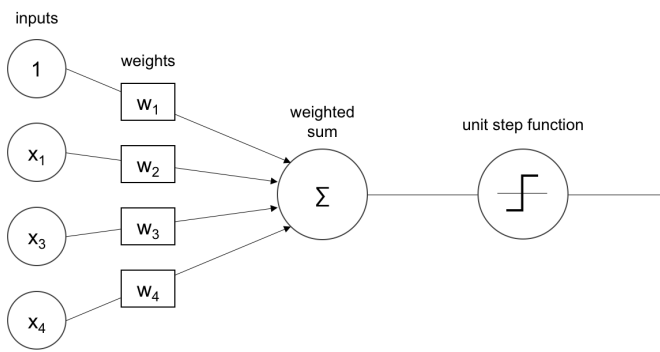


Figure 1: Schema of the Perceptron neuron

The Neural Network is made by individual neurons called *Perceptrons*. The neurons are split into input, weights and activation functions.

Nevertheless, the most important part of the neuron that determines significantly the capabilities of the neuron, is 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 Multilayer Perceptron Topology can be split into layers of three types:

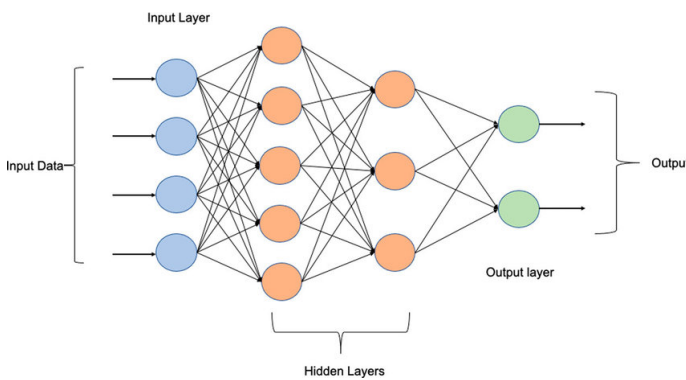


Figure 2: Schema of the Multilayer perceptron

- **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.

¹To obtain more information about Multilayer Perceptron functionality check the following [link](#)

²Frank Rosenblat, psychologist and father of deep learning, check its [biography](#).

The training process uses the **BackPropagation**³ method. The requirements of *backpropagation* are the *learning rate* (a parameter that determines how much the network will be modified while applying the method) and the *optimization method* (an algorithm to reduce the network's *error function*⁴).

The training process of neural networks is divided into *epochs*, the number of iterations where *backpropagation* is applied with a number of *batches*. The *batch* (or *batch size*) is the number of samples used in every *epoch*, and the *mini batch* (or *mini batch size*) is the number of samples used to apply the *optimization method* (using this number of samples to create a mean modification of the parameters).

In this study, the *mini batch* will be 1 and the *batch size* will be the number of samples; that is to say, every *epoch* will apply *backpropagation* for every sample.

Moreover, the neural network will be implemented in **C language** to improve efficiency, with all methods (*backpropagation*, *optimization method*, *error function*, etc.) developed by hand to ensure minimal computational cost.

1.2 Why Neural Networks work

It is important to remark that what gives us the capability to model any system is that we can express every problem as a function (no matter if it is a classification task, probability function, prediction, or regression function).

The association between a task and a function allow us to apply the **Universal Approximation Theorem**.

• **Definition:** Universal Approximation Theorem.

For any function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, with $m, n \in \mathbb{N}$, and a subset $D \subset \mathbb{R}^n$ where f is continuous at all D , $\exists \{(w_i, b_i, c_i)\}_{i=0}^k$ that:

$$f(\vec{x}) - \lim_{k \rightarrow \infty} \sum_{i=0}^k c_i \sigma(w_i^T \cdot \vec{x} + b_i) = 0$$

Where $w_i \in \mathbb{R}^n$, $b_i \in \mathbb{R}$, $c_i \in \mathbb{R}^m$, $\vec{x} \in D$ and σ the sigmoid function.

The parameters $\{w_i, b_i, c_i\}$ are associated (respectively) with the weights, bias and scale factor of the i -th neuron.

Neural Networks essentially use this theorem with a limited number of sigmoids (equal to de number of neurons), consequently an error of approximation is given.

$$f(\vec{x}) - \sum_{i=0}^k c_i \sigma(w_i^T \cdot \vec{x} + b_i) = \epsilon, k \in \mathbb{N}$$

Modeling systems the information of the function shape/tendency is limited or none, making it only possible to study the problem using data values.

This cases are handle using the *square-norm* metric (or *Mean Squared Error*).

³The BackPropagation method and how it works can be consulted in the following [link](#)

⁴The error function is introduced and explained on page 4.

• **Definition:** Square-norm:

Being $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ function followed by the system. Being $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$, with $m, n \in \mathbb{N}$, as $g(\vec{x}) = \sum_{i=0}^k c_i \sigma(w_i^T \cdot \vec{x} + b_i)$, $k \in \mathbb{N}$.

Given a dataset $B = \{(\vec{x}_i, \vec{y}_i)\}_{i=0}^N$, with $N \in \mathbb{N}$, where $\forall(\vec{x}_i, \vec{y}_i) \in B, \vec{y}_i = f(\vec{x}_i)$.

$$\Delta^2 = \frac{1}{N} \sum_{i=0}^N (\vec{y}_i - g(\vec{x}_i))^2$$

Decreasing the error Δ implies reducing the error ϵ because as Δ decrease the model improves its approximation.

As more neurons are added to the Network, the approximation and the training time increases. This increase in training time may not be worth as it can represent a minimum improvement of the model.

Moreover, the Backpropagation⁵ method used to train Neural Networks have a complexity that increases the processing time exponentially to the number of neurons.

Therefore, trying to obtain the ideal topology for solving a problem is a difficult task that in many cases resides in trial and error.

2 Neural Networks topology

The main objective of the study is to obtain a more efficient process using neural networks than the traditional ways that have been used to date.

As it has been explained before, there is not an understandable process through which to obtain an efficient neural network.

While approaching this problem, an idea that could partially solve the lack of understanding of what the neural network is doing, and deduce the optimal architecture, came to mind.

2.1 Hypotesis and Working Hypothesis

As it has been said before, when approaching this problem an intuitive idea appeared.

-
- **Working Hypothesis:** If there is enough information about the system to model, there exists an efficient neural network whose architecture can be obtained using rules.

The idea proposed is influenced by the following observations:

1. Every deterministic process is a combination of restrictions that can be represented as a function.
2. Every function can be split into domains its behaviour can be categorized as periodical, irrational, or polynomial.

⁵The Backpropagation method is an algorithm to modify the internal parameters of a Neural network by using the chain rule, more information in the following [link](#).

3. The periodical, polynomial and irrational functions are continuous in "all"⁶ domain.
4. All continuous functions can be approximated using the **Universal Approximation Theorem**.

Moreover, if a function does not follow one of these types of functions, it can be split by more intervals until each of them can be approximated by one of these types or even can be approximated using *Fourier Series* or *Taylor series* allowing it to be modelled by periodical and polynomial functions.

Therefore, this Working Hypothesis can be verified by reducing it to corroborate the following points:

1. Each category of functions mentioned can be modelled using a particular architecture of a neural network.
2. There is a determined minimum number of data registers where the neural network can be trained.

To verify the Working Hypothesis, efficiency and precision studies are required.

2.2 Study of deterministic processes

To start, let's assume the following axioms to simplify the study and give some feedback to make the correlations to verify or dismiss the Working Hypothesis.

1. All periodic functions can be approximated as a combination of sin functions.
2. All polynomial functions can be approximated as a combination of x^n functions, where $n \in \{\mathbb{N} \cup \{0\}\}$.
3. All irrational functions⁷ can be approximated as a combination of \sqrt{x} function.

All points can be proven using *Taylor Series* and *Fourier Series*. The study will be reduced to taking patterns of the architectures used to approximate the functions $\sin x$, \sqrt{x} , x^2 .

This paper uses $I - H - N - O$ nomenclature to represent the architecture of the neural network, where the letters I, H, N, O represent the number of inputs, the hidden layers, the number of neurons for each hidden layer and the number of outputs respectively.

For instance, the architecture $1 - 3 - 4 - 1$ is a network with 1 input, 3 hidden layers with 4 neurons each one and 1 output.

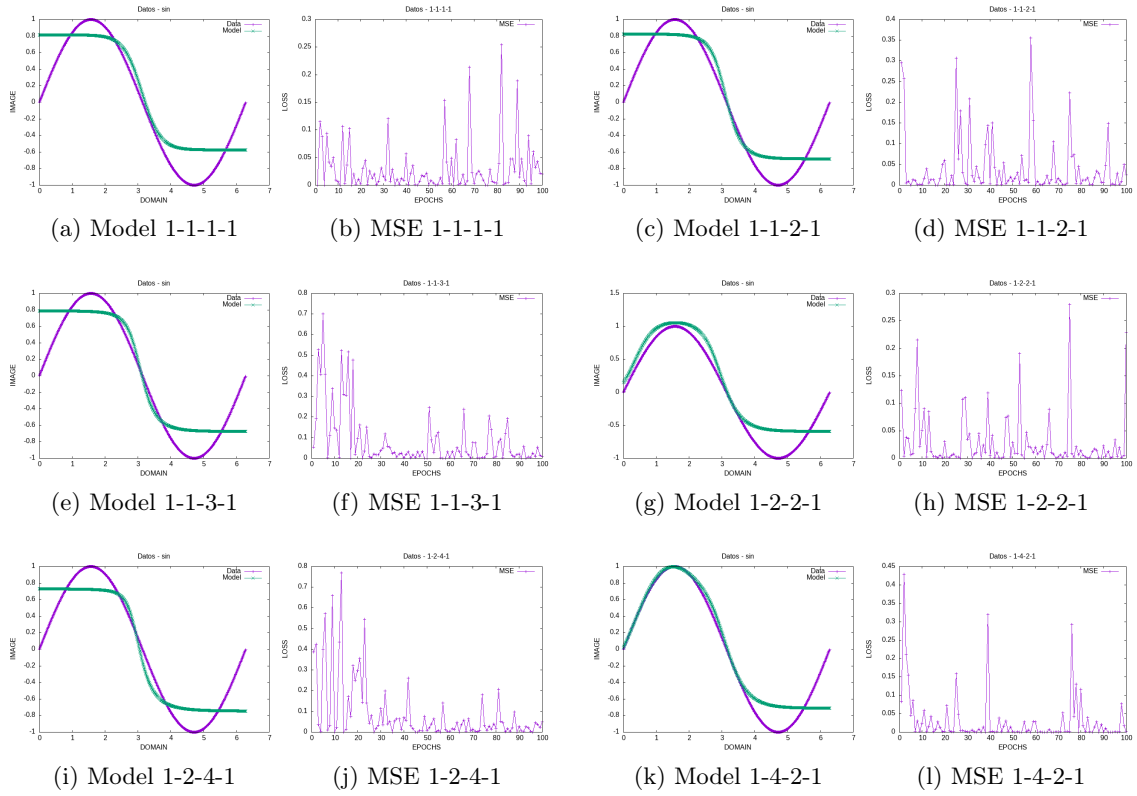
2.2.1 Study of the *sin* function

To test if there exists an architecture that can approximate efficiently the sin function in the real plane, the approximation of the function is studied in the interval $[0, 2\pi]$ using a fixed number of iterations with different architectures.

The results shown in Figure 3 were obtained by fixing 100 epochs and training the networks using the *Stochastic Gradient Descent* method with a learning rate of 0.1 trying to minimize the *Square-norm* error (*MSE*).

⁶In case of irrational functions such as \sqrt{x} , it can be interpreted as $f(x) = \begin{cases} 0, & x < 0 \\ \sqrt{(x)}, & x \geq 0 \end{cases}$

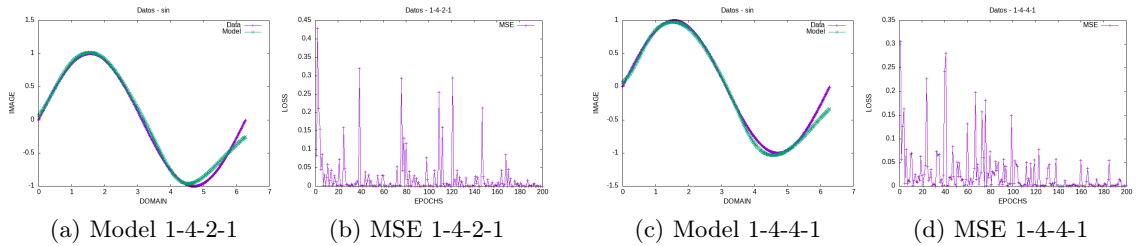
⁷In this study, *irrational functions* define the functions that present an n-th root.

Figure 3: Results of different architectures approximating the function \sin

As can be seen in the Figure, almost all architectures do not approximate the \sin function properly, not counting the architectures $1-2-2-1$ and $1-4-2-1$ which approximates the function relatively correctly.

If the experiment is repeated changing the epoch number (instead of 100 change to 200) the results are almost similar (with more precision in some parts but no improvement of modelation).

However, the architectures $1-4-4-1$ and $1-4-2-1$ give an augment of precision significantly relevant.

Figure 4: Architectures that improves the modelisation of \sin function

The approximation done by these architectures is so significant that allows us to observe and conclude that the model that approximates the best the \sin function is $1-2-4-1$.

The following conclusion is obtained because this is an architecture that efficiently approximates the function it aims to "learn"⁸.

Moreover, it can be an evidence of the following idea:

- **Suggested Hypotesis:** To approximate any periodical function in a real plane a minimum of 2 neurons is needed in each of the K hidden layers, where K is the sum of the relative extremes and points of inflexion of the function.

This suggestion is obtained by the intuition and observation that every layer of the neural network could be modelling one interval (delimited by the extremes of the function and inflexion points) and performing that approximation seems to need at least 2 sigmoids.

2.2.2 Study of the *parabolic* function

The same experiment has been done with the function x^2 (also with function x^3 to have more results to compare) with fixed epochs, learning rate and domain. Specifically, the values of epochs, learning rate and domain used are, respectively, 100, 0.1 and $[-2, 2]$.

The results obtained for polynomial functions are shown in 5.

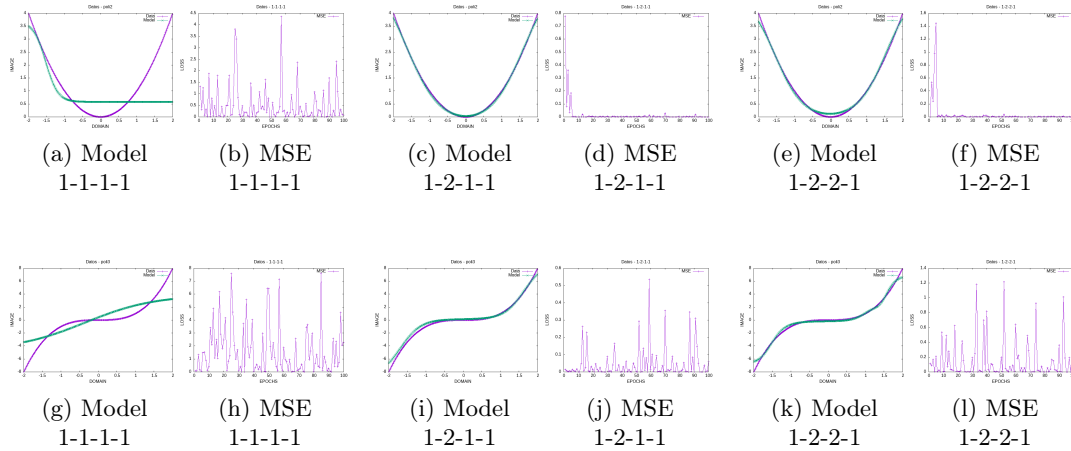


Figure 5: Results of different architectures approximating the functions x^2 and x^3

It shows how the most efficient architecture, which gives the better performance in modelling and precision, is the $1 - 2 - 1 - 1$.

This architecture was tested before with the periodical functions but without positive results. However, in this case is the best, which was one of the axioms proposed, each type of function has a diferent optimal architecture.

Therefore, a last study will be done to obtain the architecture needed to model the irrational functions.

⁸Learn in this term refers as if the training process allows to reduce the loss (also named error) near 0 with a significant decrease throughout the epochs.

2.2.3 Study of the *squared root* function

Finally, reproducing the experiment for the function \sqrt{x} , the results obtained are shown in 6.

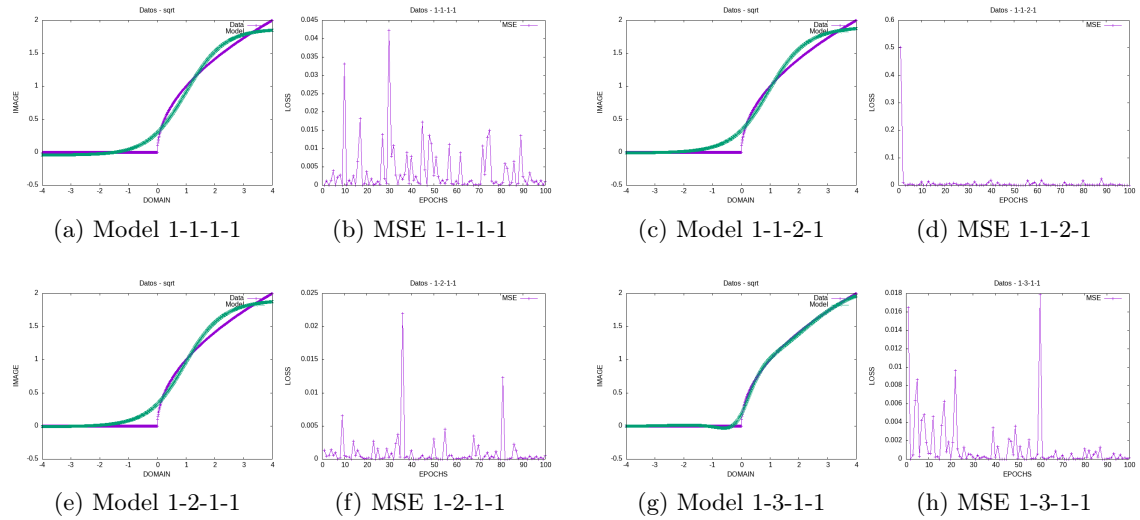


Figure 6: Results of different architectures approximating the function \sqrt{x}

In this study, the only parameter that was changed with respect to the previous experiments was the domain (which was increased to the interval $[-4, 4]$).

As it can be observed, the architecture that approximates the best the function is $1 - 3 - 1 - 1$.

Furthermore, it should also be commented that the intuition of the inclusion of irrational functions such as the \sqrt{x} in the axioms of the Working Hypothesis was made to use the property of having a function similar to *ReLU*⁹.

2.3 Viability of approximation discontinuous functions

It is important to keep in mind that splitting functions into domains can develop in discontinuities which do not allow the use of the **Universal Approximation Theorem**.

Nevertheless, the application of neural networks in tasks such as classification or regression in many areas, with no information of the function that is approximating, using activation functions different from the *sigmoid* that give (apparently) correct results give the idea that possibly the theorem could be expanded and (in some cases) have no need of the continuity restriction.

The widely use of *ReLU* function in neural networks is because it is easier to program and with enough neurons, its effects (depite it should not be possible) give a similar or even better results.

⁹The *ReLU* function is widely used as an activation function which is defined as $f(x) = \begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$

- **Suggested Hypothesis:** If a function has discontinuities different from asymptotic, the **Universal Approximation Theorem** can be used to give an approximation with error 0 when the number of sigmoid tends to infinity.

This Suggested Hypothesis can be sustained by the ability of sigmoid functions to approximate jump discontinuities. However, an experiment with discontinuous functions has been done to give evidence for this hypothesis.

2.3.1 Study of Jump discontinuity

In this experiment, a simple architecture of $1 - 2 - 2 - 1$ with a high amount of epochs (concretely 10000) and fixed learning rate of 0.1 have been used to approximate the function:

$$f(x) = \begin{cases} 1, & x \leq 1 \\ 0, & x \in (1, 2] \\ -1, & x > 2 \end{cases}$$

The results of approximating the function in the interval $[0, 3]$ are shown in 7.

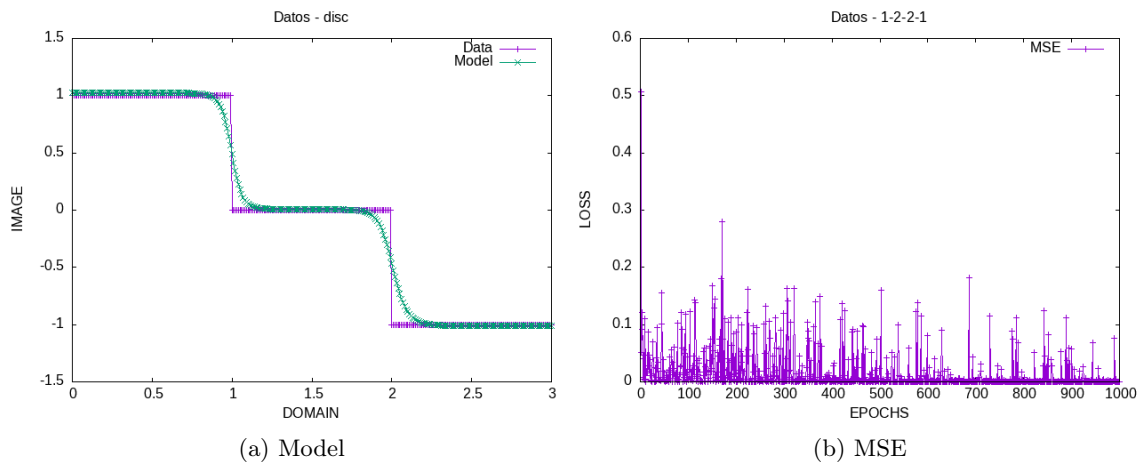


Figure 7: Architecture $1 - 2 - 2 - 1$ approximation to a jump discontinuity function

As can be seen in the figure 7, using a limited model with enough number of epochs the approximation is quite acceptable. Therefore the use of *sigmoid* functions in approximate split functions with jump discontinuities can be used at least in some cases.

2.3.2 Study of Asimptotyc discontinuity

However, the application of the **Universal Approximation Theorem** in discontinuities such as asymptotic can cause problems that avert the modelation using sigmoids.

Repeating the experiment done with the jump discontinuity function, but with a more powerful architecture, to approximate the function $\tan x$ in the interval $[-\frac{\pi}{2}, \frac{\pi}{2}]$ give us an idea of this problem.

The results obtained are shown in 8.

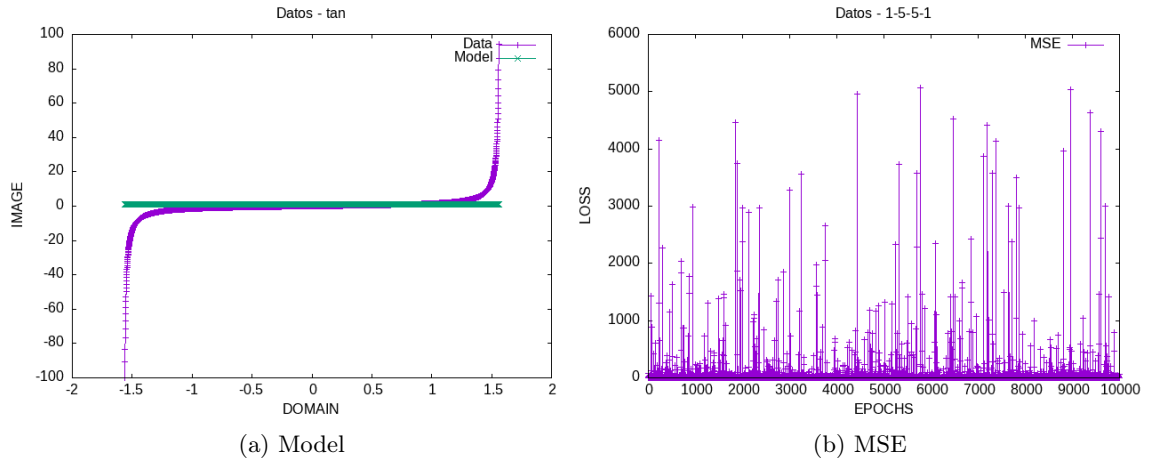


Figure 8: Architecture 1 – 5 – 5 – 1 approximation to $\tan x$ function

As can be observed the neural network is limited and cannot approximate properly the function. However, the use of other activation functions such as the *ReLU* mentioned before can give a better performance.

In the case of using the *ReLU* function, it gives another problem that can be difficult to handle, the **Universal Approximation Theorem** does not ensure that the resulting function approximates the system with a limited error.

Moreover, the use of the *BackPropagation* method and its dependency with the *Newtons Method* (the *Steepest Gradient Descend*) to obtain zeros of the *cost function* in the training process, does not ensure that the solution it arrives is the best or, even though, if it is correct in all cases (and we can be blind in situations were the network was not tested).

2.4 Observations of the results

To sum up, all the results and conclusions obtained in the studies we can conclude if there is enough evidence for the initial working hypothesis.

-
- **Working Hypothesis:** If there is enough information about the system to model, there exists an efficient neural network whose architecture can be obtained using rules.
-

To verify the viability of the Working Hypothesis needed to accomplish the following conditions:

1. All periodic functions can be approximated as a combination of sin functions.

The first condition can be proven using *Fourier series*

However, it is important to mention that a periodical function such as $\tan x$ can not be approximated due to the asymptotes.

2. All polynomial functions can be approximated as a combination of x^n functions, where $n \in \{\mathbb{N} \cup \{0\}\}$.

Such as the first condition, the second one can be proven using *Taylor Series*, not having any inconvenient.

3. All irrational functions can be approximated as a combination of the irrational functions \sqrt{x} .

Point 3 can be explained by trying to approximate a function in a certain interval. Therefore it can be obtained the following deduction:

- **Deduction:** Approximation of irrational functions.

Being $f_n : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$ where $f_n(x) = \sqrt[n]{x}$, $\forall n \in \mathbb{N}/n > 1$ and $g(x) = \sqrt{x}$.

$$\exists c, d \in \mathbb{R} / \forall x \in [a, b], f_n(x) - g(c \cdot x + d) < \epsilon$$

Where $\epsilon \in \mathbb{R}^+$ is the error and $\{c, d\}$ are parameters that adjust g to approximate f_n in the interval $[a, b] \cup \mathbb{R}$.

Nevertheless, mention that when n (the degree of the root) is even the interval needs to be in the set \mathbb{R}^+ . Moreover, this deduction can handle situations where the set of functions f_n are defined as $f_n(x) = \alpha \sqrt[n]{x} + \beta$ where the parameters $\alpha, \beta \in \mathbb{R}$. This deduction can also be demonstrated using tangency:

- **Proof:** Approximation of irrational functions.

Being $f_n : \mathbb{R} \rightarrow \mathbb{R}$, $g : \mathbb{R} \rightarrow \mathbb{R}$ where $f_n(x) = \sqrt[n]{x}$, $\forall n \in \mathbb{N}/n > 1$ and $g(x) = \sqrt{x}$.

Being $P = (p_0, p_1) \in [a, b] \subset \mathbb{R}/p_0 = \frac{\alpha + \beta}{2}, p_1 = f_n(p_0)$,

$\exists c, d \in \mathbb{R}/g(c \cdot p_0 + d) = p_1$. Therefore,

$$\exists (\alpha, \beta) \subset [a, b]/p_0 = \frac{\alpha + \beta}{2} \text{ where } \forall x \in (\alpha, \beta), |f_n(x) - g(c \cdot x + d)| < \epsilon$$

$\epsilon \in \mathbb{R}^+$ the error.

These demonstrations could be used for every function, not limiting them only to irrational functions.

4. The mentioned functions can be modelled using a particular minimal and distinguishable architecture.

The studies about the different architectures of the different functions ($x^2, \sin x, \sqrt{x}$) show that there exists a distinguishable minimum architecture that approximates the best of each function.

5. There is a determined minimum number of data to train the neural network.

That point is controversial¹⁰. In all the studies the number of data available were different with at least 600 registers. However, the relevant part about this is that it really needs a deterministic process that generates the dataset.

The minimal information that is required is the deterministic process intrinsically. If there is no deterministic process the information needed will not be enough despite how many data it could recap.

This is the consequence of the *backpropagation* method and its dependency in the *Steepest Gradient Descend*. If there is no "correct" model to compare and interpret the functionality of the network it is not possible to "extract theoretically the best architecture".

Hence, there is no "minimum number of samples" despite there is necessary existence of a (deterministic) *non-black box*¹¹ method.

Concluding the observations, the studies give enough evidence of the Working Hypothesis. To give more evidence of that proof, an experimental process will be done to extract the best minimal architecture that approximates a deterministic process.

Moreover, an efficiency test will be also performed to ensure that neural networks in collaboration with the information that gives the deterministic process improve the computation time needed.

3 Arabic to Roman numbers

The idea about that topic is the following: any natural number (normally written in the Arabic form) can be expressed as a Roman number. This process is a sequence of fixed rules that if you follow them, give you the result (it is a deterministic process).

As it is a deterministic process ensures the existence of a function that models the process, but this function is unknown.

The main objective is, with the observation and the Working Hypothesis done before, to create a neural network that approximates an unknown function.

3.1 Function to approximate

The unknown function that is going to be approximated is the composition of two functions:

1. Function that represents each natural number into its roman number expression (f_1).
2. Function that associates every expression of roman number to a natural number (f_2).

The usage of this two functions in composition, $f_2(f_1)$, is unknown but it ensures that can be represented in the plain:

$$\mathbb{N} \xrightarrow{f_1} \{C_n\} \xrightarrow{f_2} \mathbb{N}, \quad f_2(f_1) : \mathbb{N} \longrightarrow \mathbb{N}$$

Where $\{C_n\}$ is the set of roman numbers with n letters ($n \in \mathbb{N}$).

Therefore the input and output number can be different between them.

¹⁰It is controversial because the number of samples depends on the range of the interval and the back-propagation requires to perform good adjustments

¹¹A *black box* model is a process where the internal operations and functionality of it can not be interpreted partially or at all, giving the aspect of being a non-understandable or/and verified model.

3.2 Preprocess of data and reduction of variables

The main problem when tracking with Roman numbers is that is a sequence of different characters that can not be handled easily. Therefore a simplification was performed.

-
- **Application:** Interpretation of Roman numbers.

Being the set $C_n = \{c_0, \dots, c_{n-1}\}$ where $\forall c_i \in C_n$ is a roman number.

Let's define that each $\forall c \in C_n$ is formed by two characters: $c = (a, b)$, where b is a letter of the Roman number system, the set $\{I, V, X, C, \dots\}$, and a the element in the roman number system just one positions before which is b or "I" or nothing (\emptyset).

Using this assumption we can write every number, for example, the element $c = (I, V) \equiv 4$ and the sequence $C_3 = \{(\emptyset, X), (\emptyset, X), (I, V)\} \equiv 24$.

With this interpretation and definition of the Roman numbers set and its properties, another application can be done.

-
- **Application:** Reduction of liberty degrees function.

Using the previous definition of the *Roman numbers set*, it is possible to create a function $f : \{C_n\} \rightarrow \mathbb{Z}, \forall n \in \mathbb{N}$.

Before defining this function it essential to assign a number system to each letter of the Roman number system (for example its position in the alphabet starting with $A \equiv 1$ and defining $\emptyset \equiv 0$).

This method allows us to represent the example sequence used before as a set of numbers:

$$C_3 = \{(\emptyset, X), (\emptyset, X), (I, V)\} \equiv \{(0, 24), (0, 24), (9, 22)\} \equiv 24$$

Now we can create the function $f : C_n \rightarrow \mathbb{Z}, \forall n \in \mathbb{N}$ as:

$$f(C_n) = \sum_{i=0}^{n-1} b_i - a_i, \text{ where } c_i = (a_i, b_i)$$

This application allows to represent each Roman number in one Arabic number (but does not ensure that this process satisfies the injection property).

The importance about this process is that we can represent each natural number expressed in Arabic with another number (in the integers) in Arabic which is equivalent to, at least one, roman number.

In consequence, it can be represented graphically in \mathbb{R}^2 . This representation is shown in 9.

The function distribution that appears is the result of applying the function defined to the deterministic rules of Roman numbers, therefore there is no knowledge of which is the function behind this process.

However is an ideal case to use the experimental observation and try to obtain with the dataset given by the deterministic method a function that approximates the distribution of the dataset.

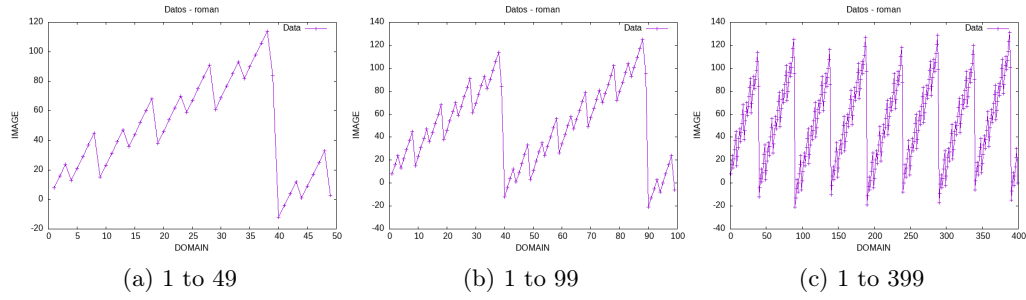


Figure 9: Distribution of the associated Arabic number to the Roman natural number

3.3 Analisis of the dataset and creation of the neural network

The shape of the function shown in 9 tends to mimic the *Sawtooth wave* but, in this case, it is not exactly periodical. This consequence (not being periodical) creates some problems because the application of the architecture developed to approximate the periodical functions can not be applied correctly.

Despite this problem, it is possible to use the architecture developed for periodical functions (tending to approach the *Fourier Serie* of the *Sawtooth wave*) in what seems to be the period of the wave $T = 49$.

Therefore, the neural neural network should have the architecture shown in 10.

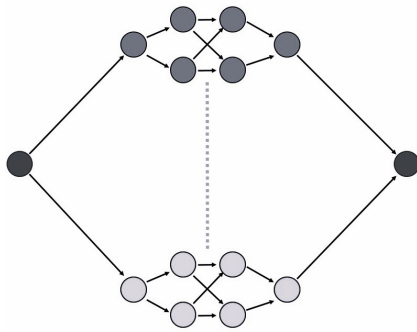


Figure 10: Working Hypotesis architecture of the neural network that approximates the Roman number distribution

This neural network will be composed of two neurons (one input and one output) that send/receive the information to/from sub-neural networks with the architecture obtained to modulate the sin function.

The neural network will be trained with only the first subset of the dataset (number between 0 and 49) and if a bigger number is introduced it only needs to be reduced as the remainder of the division of this number and the period T .

The important part of the network is the number of sub-neural networks, the learning rate and the number of epochs that need to approximate properly the distribution.

Unfortunately, the number of epochs and learning rate can not be deduced until it has been tested several times.

However, the number of sub-networks could be deduced as the amount needed to approximate the *Fourier Serie* the *Sawtooth wave* with enough precision.

To test and deduce the optimal number of sub-networks and the value of the learning rate, two similar approaches have been done, as can be seen in 11

The conclusion obtained from this test is that the simple architecture of the sin function of $1 - 2 - 2 - 1$ is not enough and could be improved by using the most robust one of $1 - 2 - 4 - 1$.

Moreover, the variable learning rate (decreasing its value while the epochs increase) adjusts better the extremis of the distribution. However, the number of epochs tested (in the case of figure 11 1000) seem not to be enough.

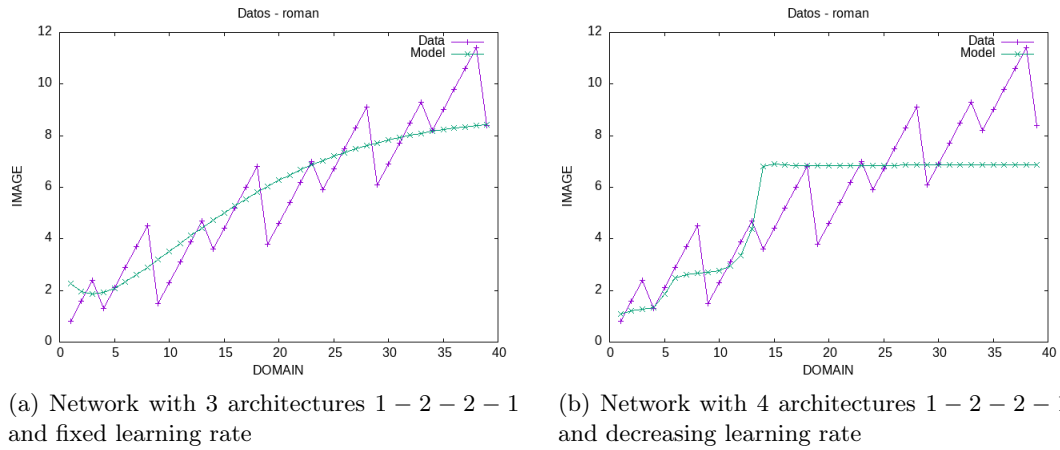


Figure 11: Test of different neural networks to approach the distribution of Roman numbers

Remark that the scale of the distribution has been decreased by one order of magnitude to reduce the variance between values and help to focus the Backpropagation when training (reducing the gradient vector to avoid erratic moves through the function cost).

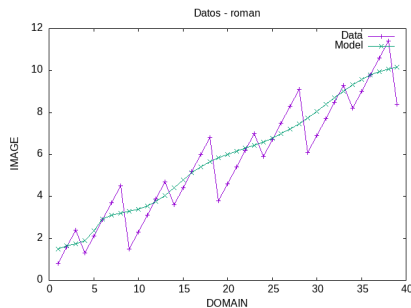


Figure 12: Improvement of sub-network capabilities to approach the roman numbers distribution

If the test is repeated with 4 sub-networks with architectures $1 - 2 - 4 - 1$, introducing the learning rate reduction by one order of magnitude each time the number of epochs done increases its order of magnitude, and with 10000 epochs the results obtained perform a better approach, as can be seen in 12.

From the figure, it is possible to deduce that the critical number of sub-networks has not been obtained yet. Moreover, observing the figure another conclusion can be made; the number of epochs is not enough.

Nevertheless, it is possible that the model is more complex than it should be.

Studying the distribution and behaviour of the Roman numbers, it can be deduced that it has 3 subperiods. The first period occurs when the number ends in 4, the second when the number ends in 9 and the third period is when the number starts with 4 or 9. Using this information can be considered that the exact amount of subnetworks of the type $1 - 2 - 4 - 1$ to perform the best relation with precision-efficiency is 3 (one for each period).

This problem could not be solved using the Working Hypotesis and gives problems when planning the correct architecture for modelling a process.

Therefore, as it is impossible to perform a comparison between the results of the methods due to the poor precision of the models, an analysis of training cost has to be done to verify if it is practical to make improvements in the process to allow the models to train with more capabilities (such as parallelization using GPU, improvements of the neural network code or more data an training time).

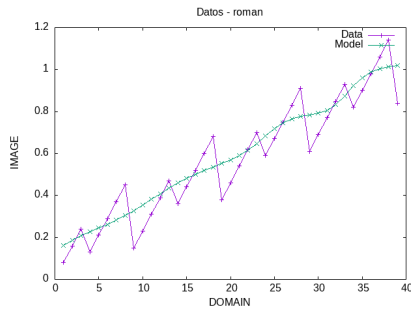


Figure 13: Final architecture proposed to approach the distribution of roman numbers

In this approach, the number of epochs has been increased to 10^6 epochs and includes 3 sub-networks with architectures $1 - 2 - 4 - 1$. The results obtained in 13 are not as good as was expected but give the best precision at the moment. However, until now there is not enough precision to consider any of these architectures a good candidate to approximate the distribution of the Roman numbers. Notwithstanding, all this unexpected non-approximation of the function could be caused by the unknown correct values of necessary epochs and/or learning rate.

3.4 Efficiency and time-cost results

To sum up, for all four tests we can reflect results¹² in the table 1, where the *Training time* units is seconds.

Sub-networks	Architecture	Epochs	Learning rate	Training time	Mean MSE
3	$1 - 2 - 2 - 1$	10^3	0.01	1.03"	4.09
4	$1 - 2 - 2 - 1$	10^3	$\{10^{-2}, 10^{-3}\}$	1.70	2.96
4	$1 - 2 - 4 - 1$	10^5	$\{10^{-(2+i)}\}_{i=0}^2$	121.04	0.87
3	$1 - 2 - 4 - 1$	10^6	$\{10^{-(3+2i)}\}_{i=0}^3$	617.55	0.14

Table 1: Analysis of cost and efficiency about the different architectures develop

It can be observed that the training models used are quite good taking into account that the epochs are "relatively low". Also, mention that the "size" of the architecture is not as big as normally MLPs are done to approach this type of problem.

Therefore it can be considered a shred of evidence that the Working Hypothesis made as a set of "considerations" to develop neural networks could be verified in other dimensions and consolidated as a theorem if it is possible to prove its properties.

In the worst scenario, the study has shown that knowing the process behind a model can be used as a sustain to create neural network experts in solving this process (which could perform a fixed amount of operations to develop a task which, in its deterministic process of calculation, will not be able to determine or even have the solution in a reasonable amount of time).

¹²The results, studies and figures that have been shown/mentioned in the report where using a Raspberry PI 4 to ensure that the computational cost and time process can be completed without any problems of different capabilities.

4 Conclusions

All the objectives initially planned have been achieved:

1. Develop neural networks using C language coding.
2. Propose a possible improvement of how neural networks can be designed.
3. Verify that The **Universal Approximation Theorem** definition can be expanded and used in cases that theoretically it should not be able to.
4. Propose a neural network using the obtained Working Hypotesis that approximates a deterministic process.

In spite of the fact that the most challenging objective, being able to improve the time process for a deterministic process using a neural network, was not achieved, during the study, it has been shown that there exists enough evidence that this challenge can be done with a better performance and analysis of the deterministic processes.

It is unfair not to mention that the deterministic process selected to verify the improvements of applying neural networks was not as easy as it was thought and with another deterministic process this verification objective could have been achieved much easily.

Also, emphasize that the different Working Hypoteses done during the study can be an initial point for future studies.

Encouraging everyone to take it as a starting point to develop better the Working and Suggested Hypothesis and correct any possible mistakes that could have been done and not taken into consideration.

References

- [1] WIKIPEDIA, *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
- [4] WIKIPEDIA, *Frank Rosenblatt*,
https://es.wikipedia.org/wiki/Frank_Rosenblatt
- [5] WIKIPEDIA, *Activating Function*,
https://en.wikipedia.org/wiki/Activating_function
- [6] WIKIPEDIA, *Artificial neural network*,
https://en.wikipedia.org/wiki/Artificial_neural_network
- [7] WIKIPEDIA, *Neural network*,
https://en.wikipedia.org/wiki/Neural_network
- [8] WIKIPEDIA, *Universal Approximation Theorem*,
https://en.wikipedia.org/wiki/Universal_approximation_theorem
- [9] WIKIPEDIA, *Mean Squared Error*,
https://en.wikipedia.org/wiki/Mean_squared_error
- [10] WIKIPEDIA, *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.