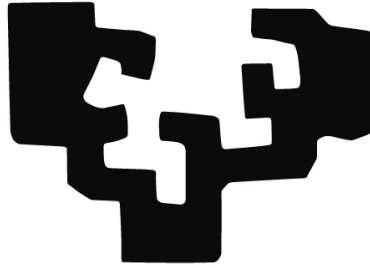


eman ta zabal zazu



Universidad
del País Vasco

Euskal Herriko
Unibertsitatea

Máster Universitario en Modelización e Investigación
Matemática, Estadística y Computación 2019/2020

Trabajo Fin de Máster

On the use of Neural Networks to solve Differential Equations

Alberto García Molina

Tutor/es

Carlos Gorria Corres

Lugar y fecha de presentación prevista

12 de Octubre del 2020

Abstract

English.

Artificial neural networks are parametric models, generally adjusted to solve regression and classification problem. For a long time, a question has laid around regarding the possibility of using these types of models to approximate the solutions of initial and boundary value problems, as a means for numerical integration. Recent improvements in deep-learning have made this approach much attainable, and integration methods based on training (fitting) artificial neural networks have begin to spring, motivated mostly by their mesh-free nature and scalability to high dimensions. In this work, we go all the way from the most basic elements, such as the definition of artificial neural networks and well-posedness of the problems, to solving several linear and quasi-linear PDEs using this approach. Throughout this work we explain general theory concerning artificial neural networks, including topics such as vanishing gradients, non-convex optimization or regularization, and we adapt them to better suite the initial and boundary value problems nature. Some of the original contributions in this work include: an analysis of the vanishing gradient problem with respect to the input derivatives, a custom regularization technique based on the network's parameters derivatives, and a method to rescale the subgradients of the multi-objective of the loss function used to optimize the network.

Spanish.

Las redes neuronales son modelos paramétricos generalmente usados para resolver problemas de regresiones y clasificación. Durante bastante tiempo ha rondado la pregunta de si es posible usar este tipo de modelos para aproximar soluciones de problemas de valores iniciales y de contorno, como un medio de integración numérica. Los cambios recientes en deep-learning han hecho este enfoque más viable, y métodos basados en entrenar (ajustar) redes neuronales han empezado a surgir motivados por su no necesidad de un mallado y su buena escalabilidad a altas dimensiones. En este trabajo, vamos desde los elementos más básicos, como la definición de una red neuronal o la buena definición de los problemas, hasta ser capaces de resolver diversas EDPs lineales y casi-lineales. A lo largo del trabajo explicamos la teoría general relacionada con redes neuronales, que incluyen tópicos como los problemas de desvanecimiento de gradientes (vanishing gradient), optimización no-convexa y técnicas de regularización, y los adaptamos a la naturaleza de los problemas de valores iniciales y de contorno. Algunas de las contribuciones originales de este trabajo incluyen: un análisis del desvanecimiento de gradientes con respecto a las variables de entrada, una técnica de regularización customizada basada en las derivadas de los parámetros de la red neuronal, y un método para rescalar los subgradients de la función de coste multi-objectivo usada para optimizar la red neuronal.

Acknowledgements

To my advisor Carlos Gorria Corres, for his advice, and to my family and friends who have given me their support in all these months.

Preamble

The structure of this work is divided into 5 chapters and 2 annexes.

Chapter 0 starts by giving an initial pragmatic overview of multi-linear algebra. Its purpose is to give anyone foreign to this subject a working knowledge of tensors: defining their notation and how to operate with them. Tensors will be extensively used throughout Chapter 2 when describing artificial neural networks.

Chapter 1 contains the actual introduction to problem at hand. Here we will be exploring the motivations for using artificial neural networks to numerically integrate initial/boundary value problems. On top of this, we will also be listing the differential operators that will be used, describe the general conditions under which we will be guaranteeing well-posedness, and examine state of the art.

Chapter 2 will layout the theoretical framework of artificial neural networks. It will be covering the everything necessary to define and train a deep learning model from ground zero. The topics covered in this section include: definition and design choices, establishment of an objective (loss) function and non-convex optimization, and the use of regularization techniques. Although these topics are general to deep-learning, throughout this whole chapter we have adapted them, where necessary, to fit the subject of this work.

Chapter 3 is the experimental part of this work. The first three sections contain the discussion on some practical issues, namely, the programming, approximating capacities of artificial neural networks and training multi-objective functions. Following these sections, lie the experiments and simulations of this work. Here we put into practice all the previous knowledge that we have build up to numerically integrate some instances of initial/boundary value problems. On each instance we benchmark and discuss the results for several set-ups based on the different architectures and training options seen up to this point.

Chapter 4 has the final conclusions to this work. An analysis on the limitations and the advantages of this technique with respect to others, as a way to approximate solutions of differential equations, is made. Also, based on the experience from this work, we suggest possible lines of work and open related questions, which can be consider for further work.

Annexes A & B include: a linear algebra perspective of some expressions in Chapter 2 for further clarity, and the code, respectively.

Contents

Abstract	I
Preamble	III
List of Figures	V
Table Index	VI
0 Overview of Multi-linear Algebra	1
0.1 What is a tensor?	1
0.2 Tensor Operations and Summation Convention	2
0.3 Linear Algebra as Multi-linear Algebra	3
0.4 Derivatives of Vector Functions and Tensors	4
0.5 The Chain Rule in Tensor Notation	5
1 Introduction	6
1.1 Posing the Problem	9
1.2 Relevant Literature	12
2 Artificial Neural Networks Framework	14
2.1 What are Artificial Neural Networks?	14
2.2 From Numerical Integration to Deep-Learning	17
2.3 Derivatives: Back Propagation and Gradient Issues	18
2.3.1 Derivatives Behaviour (Vanishing and Exploding Gradients)	19
2.4 Optimizers	25
2.4.1 First Order Methods	26
2.4.2 Second Order Methods	31
2.5 Activation Functions and Parameter Initialization	35
2.5.1 Parameter Initialization	38
2.6 Regularization	39
2.6.1 Noise-based Regularizations	40
2.6.2 Restriction-based Regularizations	42
2.6.3 Other Regularizations	45
3 Case Studies and Simulations	46
3.1 Coding Artificial Neural Networks	46
3.2 Approximating a Function	47
3.3 Training with Multi-Objective Loss Functions	51
3.4 Model Simulation	55
3.4.1 Model 1: The 1D Divergence Operator	55
3.4.2 Model 2: The 2D Divergence Operator	57
3.4.3 Model 3: The 2D Laplacian Operator	61
3.4.4 Model 4: The 1D Advection Operator	62

3.4.5	Model 5: The 2D Clairaut Operator	64
3.4.6	Model 6: The 2D Burgers Operator	65
4	Conclusions	68
4.1	Author's Final Thoughts	69
4.2	Further Work	69
A	Linear Algebra Formulation of 2.3.1	70
B	The Code	72
B.1	imports Cell	73
B.2	auxiliryPlotting Class	73
B.3	myDataSets Class	79
B.4	problemInstance Class	82
B.5	secondOrderOptimizers Class	86
B.6	myLayer Class	87
B.7	myModel Class	89
B.8	execution Cell	100
	Bibliography	102

List of Figures

2.1	Perceptron scheme.	14
2.2	A directed graph which could be a possible representation of the architecture or an artificial neural network. Nodes are artificial neurons and edges indicate which neurons feed into each other.	15
2.3	General scheme of a perceptron based fully-connected feed-forward artificial neural network.	15
2.4	Computational graph of example (2.6).	19
2.5	Computational graph (derivatives) of example (2.6). In green the flow of nodes required to compute $\partial f(x, y)/\partial x$	19
2.6	Example model: A 2-3-4-2 artificial neural network.	20
2.7	Main activation functions and their first order derivatives.	36
2.8	Combination of sigmoid functions.	37
2.9	Secondary activation functions and their first order derivatives.	37
2.10	Example of overfitting of a model.	40
2.11	Example of a model adding noisy input.	41
3.1	Comparison for different activation functions training performance for a [3,4,1]-ANN, with Adam $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_1 = 0.999$. Log10 scale.	48
3.2	Comparison for different first order optimizers training performance for a [3,4,1]-ANN, with sigmoid activations. Lower image in log10 scale.	49
3.3	Training performance of a [3,4,1]-ANN with sigmoid activations, to fit (3.1), using BFGS and L-BFGS.	50

3.4	Training performance of a [3,4,1]-ANN with sigmoid activations, to fit (3.1), using Adam with $\eta = 0.01$	51
3.5	Example of possible multi-objective functions. Component and total representation.	52
3.6	Example of possible multi-objective functions. Adjusted factors.	53
3.7	Training performance of 3 models trained for a [1,5,5,1]-ANN scheme, with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 3000 epochs. (3.5)	56
3.8	Final results. Best performing trained model (tanh) for (3.7) against the exact solution.	57
3.9	Result of a [1,10,10,1]-ANN model and tanh activations, trained with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 12000 epochs. Left plot: model against exact solution. Right plot MSE error of the model, for each point in the domain.	58
3.10	Comparison of different regularization techniques in training performance of 3 models trained for a [1,10,10,1]-ANN scheme, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs. (3.10)	60
3.11	Comparison of different regularization techniques in training performance of 3 models trained for a [1,40,40,1]-ANN scheme, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs. (3.10)	60
3.12	Final results of the best performing trained model ([1,40,40,1]-ANN, trained with the custom regularization (2.58)) for (3.7) against the exact solution. . .	60
3.13	Results and performance of the model trained for (3.11).	62
3.14	Positive and negative sign solutions of 3.13.	63
3.15	Results and performance of the model trained for (3.16).	64
3.16	Results and performance of the model trained for (3.17).	65
3.17	Results and performance of the model trained for (3.18).	67

Table Index

1.1	Comparison between FEMs and the Artificial Neural Network Methods.	8
1.2	List of differential operators used in Chapter 3.	9
2.1	List of main activation functions.	36
2.2	List of secondary activation functions.	37
3.1	Results of 3 models trained for a [1,5,5,1]-ANN scheme, with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 3000 epochs. (3.5)	56
3.2	Results of 6 models with different architectures, trained for (3.10), using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs and different regularization techniques.	59

Chapter 0

Overview of Multi-linear Algebra

Generally, when working in the context of artificial neural networks, the framework linear algebra is more than enough to describe the elements and operations taking place. Even when dealing with convolutional networks, which may involve operations on 3 dimensional arrays of objects, one can be decompose everything into vectors matrices, matrix multiplications and element-wise products. Thus, many times, when in the context of artificial neural networks, any explicit reference to multi-linear algebra or the tensor nature of such objects is disregarded.

In this work, however, we will be taking the multi-linear algebra approach. There are two main draws for doing this, i.e. generalizing vectors and matrices to tensors:

- First, the tensor notation is very powerful. This notation serves two purposes: it allows us to represent operations between tensors in a very compact way and it also helps to keep track of dimensions at any time.
- Second, multi-linear algebra provides a simple and natural framework to characterize high order derivatives of multidimensional objects such as vectors or matrices, which is a particularity of this work. In this framework derivatives and the chain rule are really easy to interpret as they visually take the form of the one dimensional case.

For the next part of this chapter we will be covering the basics of tensors. However, since the objective of this work is not to discuss multi-linear algebra, and the only purpose of this chapter is to serve as an entry point to the concepts and the notation of tensors, we will be taking a hands-on informal approach. This means that, there will be no formal definitions and every concept will be explained through an example. For a proper introduction with due rigour one can refer to chapters 2 to 4 in [1].

0.1 What is a tensor?

Perhaps the simplest way to define a tensor is as an element in a tensor space, which is nothing else than a direct product of vector spaces and dual vector spaces. So, for example, lets imagine a random tensor T in the following tensor space:

$$T \in \mathbb{R}^{4^*} \otimes \mathbb{R}^{2^*} \otimes \mathbb{R}^3, \quad (1)$$

then T is of the form $T = v \otimes w \otimes z$, where $v \in \mathbb{R}^{4^*}$, $w \in \mathbb{R}^{2^*}$, $z \in \mathbb{R}^3$. Observe that T is uniquely defined in the tensor space by $4 \times 2 \times 3 = 24$ scalar components (the individual coordinates of v , w and z , having fixed a base in each (dual) vector space).

The previous is essentially a definition of a tensor, but in practice we want to describe a tensor not by a direct product of vectors but by a set of scalar coordinates, the same way we do with a vector space. This is achieved by defining a tensor base. So, given a base for each of the (dual) vector spaces in the tensor space; for the previous example $\{e_1, e_2, e_3, e_4\}_{\mathbb{R}^{4^*}}$, $\{\hat{e}_1, \hat{e}_2\}_{\mathbb{R}^{2^*}}$, $\{\tilde{e}^1, \tilde{e}^2, \tilde{e}^3\}_{\mathbb{R}^3}$; we can intuitively build a base for the tensor space as follows:

$$\{e_i \otimes \hat{e}_j \otimes \tilde{e}^k \mid i = 1, 2, 3, 4, j = 1, 2, k = 1, 2, 3\}_{\mathbb{R}^{4^*} \otimes \mathbb{R}^{2^*} \otimes \mathbb{R}^3} \quad (2)$$

Now we can describe any tensor in the tensor space through its coordinates in the tensor base, the same way we do with vectors in vector spaces. Thus, T can be represented through its 24 scalar coordinates of the tensor base as:

$$T \equiv \sum_{i,j,k} T_k^{i,j} e_i \otimes \hat{e}_j \otimes \tilde{e}^k, \quad (3)$$

where the $T_k^{i,j}$'s are the coordinates. This latest characterization of a tensor as coordinates of a tensor base is the main way of defining tensors formally, and the one that can be found in almost every source, including [1], being any other characterization equivalent to this.

From (3) we can derive the abstract notation for tensors. Looking back at (3) we realize that once the tensor base has been selected, the explicit reference for the summation over all indices and the tensor base can be implicitly assumed, thus we can omit it. So in the end we end up writing the tensor T as $T_k^{i,j}$, where where i , j and k become free indexes which are general place-holders that can take values in $i = 1, 2, 3, 4$, $j = 1, 2$ and $k = 1, 2, 3$. It is important to note that while $T_k^{i,j}$ denotes a tensor in abstract notation, once we replace the free indexes with actual values, for example $T_1^{2,1} \in \mathbb{R}$ we denote a component. In fact, every time we replace a free index with an actual value, we fix the base element associated with that index, hence we obtain a new tensor in a lower dimension tensor space. For instance:

$$T_k^{2,j} \equiv \sum_{j,k} T_k^{2,j} e_2 \otimes \hat{e}_j \otimes \tilde{e}^k \in \mathbb{R}^{2*} \otimes \mathbb{R}^3. \quad (4)$$

Just for the sake of completeness, we will give some important facts and ideas about tensor spaces that will help understand tensors better:

- In principle, since the tensor space is a direct product, the order in which the (dual) vector spaces appear in the product does not affect the nature of the space. The only convention is that the dual vector spaces must appear before the vector spaces; so in practice $\mathbb{R}^{4*} \otimes \mathbb{R}^{2*} \otimes \mathbb{R}^3$ can describe the same tensors as $\mathbb{R}^{2*} \otimes \mathbb{R}^{4*} \otimes \mathbb{R}^3$. However, in terms of notation the order is important as the indexes must appear in the same order as the (dual) vector spaces in the direct product, by convention. Given the previous example, it is said that the tensors in the first tensor space $T_k^{i,j}$ and the tensors $T_k^{j,i}$ in the second tensor space are related by index juggling.
- A tensor space is said to be of dimension (r, s) , if it is the direct product of r dual vector spaces and s vector spaces. Hence, in the example (1), the tensor space is of dimension $(2, 1)$. Note that one of the perks of the abstract notation is that we can know a tensor's dimension easily by simply counting the number of free super-indexes (which yields r) and the number of free sub-indexes (which yields s).

0.2 Tensor Operations and Summation Convention

There are a few operations that can be performed on tensors. For example, if two tensors exist in the same tensor space, $a_k^{i,j}, b_k^{i,j} \in \mathbb{T}$ we can add them by simply adding the coordinates that correspond to the same element of the base. This is:

$$c_j^i = a_k^{i,j} + b_k^{i,j} \equiv \sum_{i,j,k} (a_k^{i,j} + b_k^{i,j}) e_i \otimes \hat{e}_j \otimes \tilde{e}^k. \quad (5)$$

Two tensors $a_k^i \in \mathbb{T}$ and $b_{k'}^{i',j'} \in \mathbb{T}'$, which can exist in very different tensor spaces can always be multiplied in what is known as the tensor product. The resulting tensor exists in the tensor space formed by the direct product combination of the original tensor spaces, $\mathbb{T} \otimes \mathbb{T}'$, the new base becomes the direct product of the original tensor spaces bases, and the new components are the product of the originals. For this example:

$$c_{k,k'}^{i,i',j'} = a_k^i \cdot b_{k'}^{i',j'} \equiv \sum_{i,k} \sum_{i',j',k'} \left(a_k^i \cdot b_{k'}^{i',j'} \right) e_i \otimes \tilde{e}^k \otimes e_{i'} \otimes \hat{e}_{j'} \otimes \tilde{e}^{k'} \in \mathbb{T} \otimes \mathbb{T}'. \quad (6)$$

There is another “way” in which we can multiply tensors. Lets imagine that we have two tensors $a_j^i \in \mathbb{R}^{2^*} \otimes \mathbb{R}^3$ and $b_{k'}^{j'} \in \mathbb{R}^3 \otimes \mathbb{R}^4$. Observe that in this particular case, the first tensor space contains the vector space \mathbb{R}^3 , and the second tensor space contains its dual counterpart, the dual \mathbb{R}^{3^*} . Since \mathbb{R}^{3^*} is the space of linear applications $\mathbb{R}^3 \rightarrow \mathbb{R}$, instead of taking the direct product $\hat{e}_j \otimes \hat{e}^{j'}$ we can apply one to the other $\hat{e}_j(\hat{e}^{j'}) = \delta_j^{j'}$, where $\delta_j^{j'}$ is the Kronecker’s delta, and we are assuming the basis are orthonormal. The rest of the elements of the basis will behave as in the tensor product. This “way” of multiplying tensors is called a contraction, and under this rules the multiplication looks like:

$$\begin{aligned} c_{k'}^i &= a_j^i \cdot b_{k'}^{j'} \equiv \sum_i \sum_{k'} \sum_{j,j'} \hat{e}_j(\hat{e}^{j'}) \left(a_j^i \cdot b_{k'}^{j'} \right) e_i \otimes \tilde{e}^{k'} \\ &= \sum_i \sum_{k'} \sum_{j,j'} \delta_j^{j'} \left(a_j^i \cdot b_{k'}^{j'} \right) e_i \otimes \tilde{e}^{k'} \\ &= \sum_i \sum_{k'} \sum_j \left(a_j^i \cdot b_{k'}^j \right) e_i \otimes \tilde{e}^{k'} \in \mathbb{R}^{2^*} \otimes \mathbb{R}^4. \end{aligned} \quad (7)$$

To this point we have defined how to sum tensors, and two ways two multiply them: the tensor product and the contraction (when compatible). However, it is confusing that we have been using the same notation to denote the tensor product and the contraction. This is where the Einstein Summation Convention comes in. Note that in (7) we have denoted the contraction operation by $a_j^i \cdot b_{k'}^{j'}$, which contains the free indexes j and j' , but looking at the final expression, because we are applying the Kronecker’s delta, one of the indexes has become irrelevant. Hence, in the product of two tensors, when the same free index is repeated, the element of the basis for that index gets contracted (Einstein Summation Convention). This means that from now on $a_j^i \cdot b_{k'}^{j'}$ will be a tensor product as in (6), and $a_j^i \cdot b_{k'}^j$ will be a contraction on the j -th index, defined as:

$$c_{k'}^i = a_j^i \cdot b_{k'}^j \equiv \sum_i \sum_{k'} \sum_j \left(a_j^i \cdot b_{k'}^j \right) e_i \otimes \tilde{e}^{k'}, \quad (8)$$

and we will call both tensor multiplications. In a practical sense the Einstein summation convention provides some kind of cancellation intuition. We can imagine as if repeating indexes up and down cancel each other in the resulting tensor, $c_{k'}^i = a_{i,j}^i \cdot b_{k'}^j$.

As a final note, we can also define the transposition of a tensor as $(a_k^{i,j})^\top = a_{i,j}^k$. In essence the transposition operator acts in each of the vector spaces that compose the tensor space.

0.3 Linear Algebra as Multi-linear Algebra

Multi-linear algebra is a generalization on linear algebra. Thus we can express all objects in linear algebra through tensors and any linear algebra operators as tensor operators. In this section we will draw a parallel between one and the other for the most basic concepts.

The most simple objects that can be described as tensors are vector and dual vectors. A vector in a vector space (a.k.a contravariant vector), which is represented in linear algebra as column array, is a (0,1) dimensional tensor. Similarly, a dual vector in a dual space (a.k.a dual vector, covariant vector, covector or one-form), which is represented in linear algebra as row array, is a (1,0) dimensional tensor. Next, there is an example of a vector and a dual vector resp., in the left written in tensor form and in the right in its linear algebra representation:

$$v^i \equiv \begin{pmatrix} v^1 \\ v^2 \\ v^3 \end{pmatrix}, \quad v_j \equiv (v_1 \ v_2 \ v_3). \quad (9)$$

A matrix in its tensor form simply is a (1,1) dimensional tensor. If we think of a matrix as a linear application, i.e. a vector of one-forms (a vector of applications from $\mathbb{R}^n \rightarrow \mathbb{R}$), it is natural that its tensor representation would be a_j^i . Note that $a^{i,j}$ and $a_{i,j}$ would also be able to describe the same number of elements, but they define objects of different nature than a matrix. An example of a matrix $a_j^i \in \mathbb{R}^{3*} \otimes \mathbb{R}^3$ in tensor and linear algebra representation would be the following:

$$a_j^i = \begin{pmatrix} a_j^1 \\ a_j^2 \\ a_j^3 \end{pmatrix} \equiv \begin{pmatrix} (a_1^1 \ a_2^1 \ a_3^1) \\ (a_1^2 \ a_2^2 \ a_3^2) \\ (a_1^3 \ a_2^3 \ a_3^3) \end{pmatrix} = \begin{pmatrix} a_1^1 & a_2^1 & a_3^1 \\ a_1^2 & a_2^2 & a_3^2 \\ a_1^3 & a_2^3 & a_3^3 \end{pmatrix}. \quad (10)$$

A special case is the identity matrix I_j^i . Observe that its components satisfy that $I_j^i = 1$ when $i = j$, and $I_j^i = 0$ when $i \neq j$. This is actually the definition of the Kronecker's delta, and indeed, it is the Kronecker's delta that is the tensor generalization of the identity, and as a consequence of the identity matrix in linear algebra. Hence, from now on $I_j^i = \delta_j^i$.

Finally, regarding operations we can represent: the scalar product of two vectors as a transposition and a contraction, $w^i \cdot (v^i)^\top$; the outer product of two vectors as transposition and a tensor product, $w^i \cdot (v^j)^\top$; a matrix-vector multiplication as a contraction, $a_j^i \cdot v^j$; and a matrix-matrix multiplication as a contraction, $a_j^i \cdot b_k^j$.

0.4 Derivatives of Vector Functions and Tensors

In this section we will discuss the need of tensors as a tool to deal with high order derivatives of vector functions through an example. Lets suppose we have a vector function f which takes as an input a vector in \mathbb{R}^3 and outputs a vector in \mathbb{R}^2 , through a linear algebra representation this would be:

$$f : \quad \mathbb{R}^3 \quad \longrightarrow \quad \mathbb{R}^2$$

$$x = \begin{pmatrix} x^1 \\ x^2 \\ x^3 \end{pmatrix} \longrightarrow \begin{pmatrix} f^1(x) \\ f^2(x) \end{pmatrix}. \quad (11)$$

Note that the application does not have to be lineal, as the components f^1 and f^2 can be non-linear.

Now we would like to get the partial derivatives with respect to the components x^1, x^2, x^3 , of the input vector. These partial derivatives define the Jacobian matrix, which is well known in vector calculus:

$$J = \begin{pmatrix} \frac{\partial f^1(x)}{\partial x^1} & \frac{\partial f^1(x)}{\partial x^2} & \frac{\partial f^1(x)}{\partial x^3} \\ \frac{\partial f^2(x)}{\partial x^1} & \frac{\partial f^2(x)}{\partial x^2} & \frac{\partial f^2(x)}{\partial x^3} \end{pmatrix}. \quad (12)$$

Considering further higher order derivatives into account is already a problem. The function (11) is a one dimensional object, and its first order derivatives (12) already require a two dimensional object to be represented. Therefore its second order derivatives will require a three dimensional object (a matrix whose elements are dual vectors), to be represented. This is impractical and it is where tensors become really useful.

It is possible to describe the vector function as a tensor, $f^k(x^i) \in \mathbb{R}^2$, and take advantage of multi-linear algebra to represent its partial derivatives. However, first we need to know how to differentiate in the context of tensors. This is part of what is known as Ricci calculus or tensor calculus (a subdomain of differential geometry), which has several ways of defining derivatives on tensors. Nevertheless, in practice, since we will only make use of partial derivatives (which are done with respect to a single component, ignoring the rest), the only thing we need to do is to multiply the following dual vector (one-form) through a tensor product:

$$\frac{\partial}{\partial x^i} = \left(\frac{\partial}{\partial x^1} \quad \frac{\partial}{\partial x^2} \quad \frac{\partial}{\partial x^3} \right) \in \mathbb{R}^{3*}. \quad (13)$$

Applying the operator to $f^k(x^i)$ once we obtain:

$$\frac{\partial}{\partial x^{i'}} \cdot f^k(x^i) \in \mathbb{R}^{3*} \otimes \mathbb{R}^2, \quad (14)$$

which has, in fact, the dimensions of a matrix as shown in the previous section. Indeed, this is the Jacobian matrix, and if we were to replace the free indices k and i' to obtain every component and write them in matrix form, we would exactly recover (12).

Applying the operator to $f^k(x^i)$ multiple times we obtain its higher order derivatives as:

$$\frac{\partial^2}{\partial x^{i'} \partial x^{i''}} \cdot f^k(x^i) \equiv \frac{\partial}{\partial x^{i'}} \cdot \frac{\partial}{\partial x^{i''}} \cdot f^k(x^i) \in \mathbb{R}^{3*} \otimes \mathbb{R}^{3*} \otimes \mathbb{R}^2, \quad (15)$$

$$\frac{\partial^3}{\partial x^{i'} \partial x^{i''} \partial x^{i'''}} \cdot f^k(x^i) \equiv \frac{\partial}{\partial x^{i'}} \cdot \frac{\partial}{\partial x^{i''}} \cdot \frac{\partial}{\partial x^{i'''}} \cdot f^k(x^i) \in \mathbb{R}^{3*} \otimes \mathbb{R}^{3*} \otimes \mathbb{R}^{3*} \otimes \mathbb{R}^2. \quad (16)$$

0.5 The Chain Rule in Tensor Notation

As we have seen in (14-16), derivatives using the Einstein Summation Convention take a really simple form, which resembles in appearance to the one variable calculus framework. Another case in which this happens is when applying the chain rule. For instance, given $(f^k \circ g^l)(x^i) = f^k(g^l(x^i)) = f^k(y^l) : \mathbb{R}^2 \rightarrow \mathbb{R}^3 \rightarrow \mathbb{R}^2$ where $f^k : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ and $g^l : \mathbb{R}^2 \rightarrow \mathbb{R}^3$ the chain rule is applied just by making use of contractions as follows:

$$\frac{\partial}{\partial x^{i'}} \cdot f^k(g^l(x^i)) = \frac{\partial f^k(g^l(x^i))}{\partial g^l(x^i)} \cdot \frac{\partial g^l(x^i)}{\partial x^{i'}} = \frac{\partial f^k}{\partial y^l} \cdot \frac{\partial y^l}{\partial x^{i'}} \in \mathbb{R}^{2*} \otimes \mathbb{R}^2 \quad (17)$$

Chapter 1

Introduction

Many phenomenons in physics, finance and other branches of science are modelled through equations or systems of equations involving rates of change, i.e. derivatives. These systems are generally classified as: systems of ordinary differential equations or ODEs, when every equation depends on the derivatives of a single variable (usually called time); or systems of partial differential equations or PDEs, when some equation in the system depends on the derivatives of more than one variable.

Finding the solutions to these systems of differential equations is of great importance to make predictions based on these models. Nevertheless, there are very few cases in which one can find an analytic expression for these solutions. Usually this involves having a very specific type of equations and making some kind of assumption on the form of solution. This is the case of methods such as the integrating factor to solve linear ODEs; or the method of characteristics to solve certain types of PDEs, whereby we assume that there is a parametric form to the solutions of the system. Therefore, in many situations, the only way to gain insights from the solutions of these systems of differential equations is to compute them numerically (known as numerical integration).

There are a wide variety of methods to integrate numerically ODEs and PDEs. In general, the main methods follow the following ideas:

- In the case of a system of ODEs we can think of its solutions as functions of the time variable, or simply as trajectories in time. Hence, starting at an initial point (the initial condition), the idea is to numerically compute the trajectory moving along the tangent in discrete intervals (steps) of time. The two principal methods that use this idea are the Runge-Kutta and the Taylor Series methods; where the first “averages” the tangent of intermediate substeps of time, and the second uses the Taylor theorem locally up to some order, to get better approximations of the tangent for larger steps of time.
- When it comes to PDEs the approach of moving through the trajectory becomes useless, as the solutions are no more one dimensional curves, but at least two dimensional or higher surfaces for which there is no unique way to move through them starting at any given point. For these types of systems the most widely used methods are the ones in the Galerkin projection family, whose general oversimplified idea consist in using a base of functionals $\{\phi_k\}_{k=0}^{\infty}$ in an adequate function space (a Sobolev space) to approximate the solution via a linear combination of base elements $u(x) = \sum_{k=0}^{\infty} c_k \phi_k(x)$. Replacing the expression $u(x)$ into the differential equations and projecting into the function space $\langle \cdot, \phi_k \rangle$, a linear system of equations is obtained which has to be solved to obtain the coefficients c_k . Actually, in this process the main computational expense usually comes with the resolution of the linear system, which has complexity $O(m^3)$, where m is the number of variables, i.e. the number coefficients c_k . This is why selecting a good base is very important, since a good choice can lead to a sparser matrix for the linear system after the projection which can greatly reduce the number of operations.

The two principal methods in the Galerkin family are the spectral method and the finite elements method (FEM), the main difference of which relies in the choice of the base. While the first uses a truncated base of functions which are globally defined in the solution's domain, like a truncated Fourier base $\{e^{2\pi k/L}\}_{k=0}^N$; the second uses a finite base $\{\phi_k\}_{k=0}^N$ of piece-wise locally defined “tent-like” functions, each associated to a node of a mesh of the solution's domain $\{x_k\}_{k=0}^N$, in a way that $\phi_k(x_j) = \delta_j^k$. This last type of methods that use a mesh approach are referred to as collocation methods.

Each of these methods have their own pros and cons. The spectral methods are typically used in very specific equations where a particular base is known to work well, for example the Fourier base on the Laplace equation, or in general, where solutions are known to be periodic. One good quality of these methods is that they are very stable and enjoy very good error properties, however, they are difficult to generalize. Applying an ill-thought choice of base to an equation, can lead to projections $\langle \cdot, \phi_k \rangle = \int_{\Omega} \cdot \phi_k(x) dx$ with very hard integrals to compute and a dense linear system matrix to solve for c_k . In contrast, FEMs are very generalizable as they use piece-wise polynomial and locally defined functions as a base, which make projection integrals $\langle \cdot, \phi_k \rangle$ relatively easy to compute and the linear system matrix to obtain the c_k s is very sparse given that there is very little overlapping in the base function's compact support. As a downside, its error and stability are strongly tied to the number of nodes (usually large) in the mesh, which cannot be trivially generated either, and requires some computational complexity.

These are some of the most popular methods, but there are many more. For PDEs in particular, some other notable mentions are: the finite differences methods, which makes use of differences between nodes to approximate derivatives; the pseudo-spectral methods, which are similar to the spectrals, but it uses quadrature formulas for the projection integrals; and the Monte Carlo methods, which consist in simulating and averaging copies of the system.

In general, there is no single best method, each has its own trade-offs. Nonetheless, in most instances for PDEs (especially in engineering), the preferred method of numerical integration is the FEM due to its versatility. It avoids any complications of having to choose a particular basis of functions, the mesh can be adapted easily to irregular solution domains, the projection integrals are easy to compute, the linear system matrix to obtain the coefficients c_k is sparse, the stability is good enough for many equations, and the error behaviour is well known and can be reduced refining (adding points to) the mesh. However, as good as this sounds, this method incurs in what is known as a dimensionality curse (impracticality when scaling to larger dimensions). For instance, suppose we want to solve some PDE within a certain error threshold in various dimensions. Solving the 1D version on a segment would require a n node mesh; the 2D version on a square would require a n^2 node mesh; the 3D version on a cube would require a n^3 node mesh; and so on. Hence, to maintain the same error threshold in a system with d dimensions requires a mesh with n^d nodes. If we recall, solving the linear system for the coefficients c_k had computational complexity of $O(m^3)$, with m being the number of nodes in the mesh (same as functions in the basis in the FEMs). Putting all this together means that in order to keep the error under a threshold as the dimension increases the complexity scales as $O((n^d)^3) \equiv O(n^{3d})$. Thus, the conclusion is that: *to keep the error within a certain threshold, the computational complexity of the FEM at least grows exponentially with the dimension of the equations*, not even taking into account the increase in computations for the mesh or the projections. Already a system in 4 dimensions starts to be impractical to solve with FEMs. (From now on, 4 or more dimensions is high-dimensional).

High dimensional systems are not very common in physics, but arise in many fields such as sociology and economics. For example, if we were to consider option pricing in finance, assuming the market parameters constant (to not incur in a stochastic problem), the system would be modelled after a PDE which has at least as many variables and dimensions as stocks in the portfolio as well as the time, which is generally a large number [2]. In cases such as the one we have just exposed, FEMs are impractical and Monte Carlo methods are used [3], but still have some stability limitations. For this reason in recent times, with the many improvements in artificial neural network, new machine learning methods have resurged as potential candidates to deal with these kinds of high dimensional problems. The main idea is based in using the good qualities of artificial neural networks as function approximators.

An artificial neural network is just a complex parametrized function $\mathcal{N}(x; \theta)$, which uses modular architecture based on the concept of neurons, has a structure optimized for computer processing, and makes use of non-linear optimization algorithms to train its parameters to fit some model (we will cover this in Chapter 2). Basing ourselves in the previous simplified definition, the deep-learning approach should be straight forward, simply put: the method will approximate the exact solution by taking an artificial neural network, replacing it into the differential equation and using an optimization algorithm to train its parameters so that the equation is satisfied; all while making use of deep-learning strategies to speed up the process. In [4] this type of methods is referred to as “Deep Galerkin Method”, the reason being that: both methodologies revolve around approximating the exact solutions of a differential equation via a parametrized function, either a linear combination of base functions or an artificial neural network; and both involve replacing this approximation into the differential equation and solving an inverse problem to find its coefficients or parameters. However, the artificial neuron strategy differ much in nature and lacks many of the elements of the methods in the Galerkin family as it does not: take into account the idea of weak formulation (which we have not explained here for simplicity); use linear combination of base functions and projections; and the resulting inverse problem does not lead to solving a linear system of equations in favour all in favour of a pure non-convex optimization. In fact, it is because of these differences that this machine learning approach should be, in theory, able to scale well with dimension, since in using non-convex optimization, all dimensions are trained at the same time, which should not increase much computational cost. On the other hand, one of the main problems is that the error is not bound by an order and is unpredictably subject to optimization and training particularities. The following table summarizes all the above:

Galerkin Methods (FEM)	Artificial Neural Network Methods
Approximates the solution with a base of linear functions.	Approximates the solution with an artificial neural network.
Requires computing some integrals (or quadratures) and solving a linear system.	Requires solving a non-convex optimization problem.
Error order and stability properties known.	Error and stability unknown and depends on the specifics of the optimization.
The complexity scales exponentially with the dimension.	Generalizes well to higher dimensions with just a few more neurons.

Table 1.1: Comparison between FEMs and the Artificial Neural Network Methods.

In this work, we will be using deep-learning techniques and methodologies to try and solve some instances of differential equation. The objective will be to analyse the viability and capabilities of these methods. Although the main interest for this methods is in integrating PDEs (as the existing ODE integration methods are already very efficient), we will start in a progressive way, by studying its application on ODEs (which can be seen as a particular case of PDEs). Then we will scale up the complexity of the operators until we are able to solve some low-dimensional PDEs. Higher dimensional equations will be out of scope since the aim is to illustrate the feasibility, and strengths-weaknesses of this strategy for which two dimensions will be enough.

1.1 Posing the Problem

In its most general form, a system of differential equations with solution in the real space may be represented as follows:

$$\mathcal{L}[u(x)] = f(x), \quad x \in \Omega \subseteq \mathbb{R}^n, \quad (1.1)$$

where $\Omega \subseteq \mathbb{R}^n$ is a compact manifold, $\mathcal{L}[\cdot]$ is a differential operator, $f(x) : \Omega \rightarrow \mathbb{R}^m$ is the external force, and $u : \Omega \rightarrow \mathbb{R}^m$ is a solution of the system. Note that (1.1) may represent either a system of ODEs or PDEs depending of the differential operator. The list of operators which we will be solving in Chapter 3 are:

Name	Expression
Identity Operator	$\mathcal{L}[u(x)] = u(x)$
1D Divergence Operator	$\mathcal{L}[u(x)] = \frac{\partial u(x)}{\partial x}$
2D Divergence Operator	$\mathcal{L}[u(x, y)] = \frac{\partial u(x, y)}{\partial x} + \frac{\partial u(x, y)}{\partial y}$
2D Laplacian Operator	$\mathcal{L}[u(x, y)] = \frac{\partial^2 u(x, y)}{\partial x^2} + \frac{\partial^2 u(x, y)}{\partial y^2}$
1D Advection Operator	$\mathcal{L}[u(x)] = u(x) \cdot \frac{\partial u(x)}{\partial x}$
2D Clairaut Operator	$\mathcal{L}[u(x, y)] = x \cdot \frac{\partial u(x, y)}{\partial x} + y \cdot \frac{\partial u(x, y)}{\partial y}$
2D Burgers Operator	$\begin{aligned} \mathcal{L}[\mathbf{u}(x, y)] &= \mathcal{L}[(u_x(x, y), u_y(x, y))] \\ &= \left(u_x(x, y) \cdot \frac{\partial u_x(x, y)}{\partial x} + u_y(x, y) \cdot \frac{\partial u_x(x, y)}{\partial y}, \right. \\ &\quad \left. u_x(x, y) \cdot \frac{\partial u_y(x, y)}{\partial x} + u_y(x, y) \cdot \frac{\partial u_y(x, y)}{\partial y} \right) \end{aligned}$

Table 1.2: List of differential operators used in Chapter 3.

Recall that at the start of this section, in defining (1.1) we indicated that u was “a” solution to the system of differential equations. In fact, there are usually many solutions or none may even exist. To ensure existence and uniqueness of the solution we need to impose some additional conditions to (1.1), namely initial conditions on ODEs and boundary conditions on PDEs. The most common set of these conditions are:

$$\text{Cauchy (ODE):} \quad u(x_0) = u_0 \quad (1.2)$$

$$\text{Dirichlet (PDE):} \quad u(x) = g(x), \quad x \in \Gamma \equiv \partial\Omega \quad (1.3)$$

$$\text{Neumann (PDE):} \quad \frac{\partial u(x)}{\partial n(x)} = g(x), \quad x \in \Gamma \equiv \partial\Omega \quad (1.4)$$

$$\text{Cauchy (PDE):} \quad u(x) = g_1(x) \wedge \frac{\partial u(x)}{\partial n(x)} = g_2(x), \quad x \in \Gamma \equiv \partial\Omega \quad (1.5)$$

where Γ or $\partial\Omega$ (depending on the convention) is the border of the domain Ω and $n(x)$ is the normal vector at the point $x \in \Omega$. Descriptively, Cauchy initial conditions fix the solution value at a certain point; Dirichlet border conditions, fix the solution values at the border of the domain; Neumann border condition, fix the flow coming in and out of the domain; and finally, Cauchy border conditions are a mix of Dirichlet and Neumann conditions. [5]

Before proceeding, one observation has to be made on ODEs. Given a single ODE of n -th order (the highest derivative in the equation has order n) with $n > 1$, it is common practice to transform the equation into a system of first order ODEs by simply introducing the following set of $n - 1$ equations $u_1 = u'$, ..., $u_{n-1} = u'_{n-2} = u^{(n-1)}$, and using them to replace any derivatives of order higher than one in the original equation. This means that any given n -th order ODE is equivalent to a system of n first order ODEs; thus finding a solution to the original n -th order equation, $u(x)$, is equivalent to finding an extended manifold solution in the corresponding system of first order equations, $\mathbf{u}(x) = (u, u_1, \dots, u_{n-1})(x)$, which includes its derivatives. The (1.2) definition of Cauchy initial conditions is based on this last paradigm where we consider systems of first order ODEs. Hence, when considering a n -th order ODE in its original form, the equivalent of fixing an initial point on the manifold solution $\mathbf{u}(x_0) = (u_0, u_{1,0}, \dots, u_{n-1,0})$ is to fix the value of its first $n - 1$ derivatives, and on those premises these conditions should be written as $u(x_0) = u_0, \dots, u^{(n-1)}(x_0) = u_{n-1,0}$.

Summarizing, we will be considering systems of differential equations (1.1) in combination some initial/boundary conditions (1.2-1.5), mainly Cauchy conditions, to formulate what are known as an initial/boundary value problems. The main objective is to formulate a “well-posed” problem: a set of basic properties which is required to apply any numerical integration successfully. A system of differential equations is said to be well-posed in the Hadamard sense [5], if it holds the following three conditions:

- A solution exists.
- The solution is unique.
- The is stable, i.e. it changes continuously with small variations of its initial conditions, boundary conditions and external force.

Proving that a given problem is well-posed is a really tricky matter. There are very few general results and many of the proofs are case specific: they may apply to certain types of differential operators (for example linear or Poisson operators), require a certain type of boundary conditions and impose several degrees of regularity.

As shown in Table 1.2, in this work we will be using very simple differential operators, all of them linear or quasi-linear. Also, the external force terms will always be an analytic functions (actually polynomials) and the initial/boundary conditions will be for the most part of Cauchy type. Under these specific conditions the Cauchy-Kovalevskaya theorem guarantees the existence of an unique analytic solution to the expression in both the ODE and PDE cases.

Nonetheless, this theorem has its limitations:

- First, it is a local theorem, although this can be remedied if all the terms are analytic everywhere to form a global version by “stitching” the local solutions in several local neighbourhoods to form a cover of the domain and build a global solution. Since the solution has to be unique in the intersection of the neighbourhoods the global solution has to be unique.
- Second, the proof is very dependant on the analyticity of the coefficients in the operator and external force as its proof relies on the methods of majorants. The sketch of this proof goes as follows [6, 7]: first we assume that the solution can be written as power series in some neighbourhood $U \subseteq \Omega$, the coefficients of which are obtained from the initial conditions and replacing the power series into the differential equation. Then we attempt to find some power series that majorates the solution power series, the definition of which is that $\sum_k a_k x^k$ majorates $\sum_k b_k x^k$ iff $|a_k| < b_k$. Finally, we use the property that states that if a series is majorated by a series that converges, so does that series. If the majorating series to the solution power series is adequately chosen and converges, so does the solution power series which converges to the local unique analytic solution.

When the differential equation is linear or quasi-linear, for every $x_0 \in U \subseteq \Omega$, there is always a systemic change of variables $h : U \rightarrow V$ such that $h(x_0) = 0$. Then, on this new domain V , we can always construct a power series that converges to 0 with radius of convergence $\rho = 1$, and majorates the power series solution $u(h(x)) : V \rightarrow \mathbb{R}^n$. This makes the proof independent of the differential operator as long as it is linear or quasi-linear. Note that, this proof is constructive as the power series solution satisfies the differential operator and initial/boundary conditions, and it converges on a neighbourhood $h^{-1}(D_0(1))$ of x_0 . Therefore, it is a valid local solution, and likewise its uniqueness is proven from a similar argument.

The assumptions of analyticity and constructiveness of the proof in this theorem implies that the theorem is proving that there is a unique analytic solution. This is much different than claiming that the only solution is analytic. Hence, the initial/boundary problem could still have other non-analytic solutions. (Actually, in the case of ODEs, the Picard–Lindelöf theorem guarantees general uniqueness over other solutions, so the analytic one is the only one; but there are no similar results for PDEs.)

Despite the two potential limitations in applying the Cauchy-Kovalevskaya theorem that we have just seen, this will be enough for the artificial neural network to approximate the analytic solution of the problem. The reason for this assumption is that the artificial neural networks will be composed of analytic functions (almost everywhere), thus we expect them to fit preferentially that solution. From now on, there will be no further discussions about the well-posedness of the initial/boundary problems that we will attempt to solve in this work, the Cauchy-Kovalevskaya theorem will always apply.

1.2 Relevant Literature

The approach for solving differential equation systems dates relatively “old”; at least, we have found an article [8], dating back to 1994. Although this article uses a graph-like structure acknowledged as an artificial neural network to solve ODEs, it applies a FEM type of “tend-like” activation functions and does not rely in “training” in the modern sense, i.e. defining a non-convex optimization problem, opting instead for some kind of Galerkin method hybrid. Throughout this article there are some references of some papers which use some kind of mean square error and non-convex optimization (the most standard approach nowadays), but the author regards them as computationally expensive. This shows that the state of the field of deep-learning back then did not allow for these strategies to be viable candidates to integrate differential equations.

Moving to more recent times, articles explicitly integrating ODEs with artificial neural networks are hard to come by, since as explained before, there are very efficient methods already to integrate ODEs, and the main interest is in PDEs. A related case that we found very interesting and worth mentioning is [9], which uses a reverse approach. Instead of training an artificial neural network to integrate an ODE, it uses ODE numerical integrators to train artificial neural networks.

With regards to PDEs, [4] is a very complete work. It defines a loss by the discretization over a random collocation of points, of the error of the artificial neural network with respect to the boundary value problem (the same idea we will be using to define a loss in 2.2). Then, it goes to solving very high dimensional free boundary PDEs (for American options), and boundary problems (for the Hamilton-Jacobi-Bellman). Two interesting features in this work are that: in the article is called Monte Carlo method for fast computation of second derivatives, which is a type of synthetic gradient; and proof to restricted version of a universal approximation theorem for the solution of PDEs. A synthetic gradient [10] is usually used for very large networks or very large amounts of data. Instead of computing the exact derivatives of the loss function with respect to the parameters required to minimize the loss, the derivatives are drawn from a distribution which is updated for every step of the training. This technique trades off not having the exact derivatives, with less computational cost and possibility of asynchronous training. The authors of [4] use this Monte Carlo method to avoid the expensive cost of a second order automatic differentiation for very high dimensions. In this work, we will not be considering this technique since, unlike [4] which integrates PDEs of up to 200 dimensions, we only integrate PDEs of up to 2 dimensions (not really high dimensions) like almost all of the other papers that we will review next do. However, the use of a synthetic gradient is a recurrent theme in papers dealing with very high dimensional systems.

Other paper focused on high dimensions are [11] and [12]. Both are similar in that they do not consider deterministic PDEs but BSDEs (Backward Stochastic Differential Equations) such as the Allen-Cahn (physics) or Black-Scholes (economics) equations, and they consider systems of up to 100 dimensions.

Closer to the line of work of this project are [13, 14, 15, 16, 17, 18]. The outlines of these work are quite similar: they simulate PDEs of up to 2 dimensions and do some kind of error analysis. Some the characteristics of [15] is that it analyses the effect in the error of the mesh and number of hidden nodes, and in [16, 17] the method is compared to a standard FEM method. On the more interesting side of things lie [13] and [18].

[18] follows up on the architecture of [4], which uses a special kind of feed-forward artificial neural network. In a regular feed-forward artificial neural network the neurons are divided into sequential layers, then the outputs of a layer strictly get fed as input to the next layer (we will see this in section 2.1). However, [4] used an architecture where the outputs of a layer feed all of its successive layers. This seem to yield good results although there is no comparison to other type of architectures.

[13] is focused in, instead of using the artificial neural network approach for its capabilities to integrate systems in high dimensions, in using its mesh-free nature to integrate over irregular domains. In this paper artificial neural networks are trained to fit the advection and diffusion operators for collocations of very irregular domains. It also applies a very original idea which is to consider the approximated solution as $\hat{u}(x) = g(x) + D(x) \cdot \tilde{u}(x)$, where $g(x)$ is the boundary condition, $D(x)$ is a distance function to the border such that $D(x) = 0$ iff $x \in \Gamma$, and $\tilde{u}(x)$ is a regular artificial neural network. This way $\hat{u}(x)$ always satisfies the boundary conditions by construction and it is only required to train the model to fit differential equation, thus one can pre-compute the distance from the collocation to the border since it does not change throughout the training, and focus on a single objective loss function. For this work we reckoned that this idea would only work well with Neumann or Dirichlet boundary conditions, but with Cauchy boundary conditions which include both at the same time.

Other papers that relate artificial neural networks to PDEs are: [19], which is a version of [9] relating PDEs to the dynamics of non-convex optimization in convolutional networks; and [20] which draw the same relation between the dynamics of the optimization of general artificial neural networks and PDEs, using statistical physics techniques. Also [21, 22, 23] are a series of papers by the same author which train artificial neural networks with experimental data from physics to learn the underlying behaviours modelled by PDEs.

Finally, we want to remark that the need for the derivatives of artificial neural networks with respect to inputs, which seem to be something that would not appear in simple regression or classification problems, thus only related to this topic, has been used in other contexts. [24, 25] are examples of this, both make use of information about the derivatives as regularization techniques and to speed up training in problems with no imperative use for them.

Chapter 2

Artificial Neural Networks Framework

In this chapter we will be covering from scratch everything about artificial neural networks that we will be using to solve differential equations in the next chapter. We will start by defining what an artificial neuron and artificial neural network are; explain the architecture of a fully-connected feed-forward neural network; their possible activation functions and initializations; show how to assign a loss function to train the model to fit the initial/boundary value problem; discuss the pros and cons of the main optimizer options available to train the model; and examine diverse regularization techniques which help improve training results.

2.1 What are Artificial Neural Networks?

Artificial neural networks are ensembles of units called artificial neurons. There are many designs for these artificial neurons, and by combining and arranging them in different ways we can create networks with very different behaviours and results.

In this work, the only type of artificial neuron that we will be using is known as the perceptron, which is probably the simplest and the most widely used. Figure 2.1 shows the basic scheme of a perceptron. A perceptron takes in n inputs, which we can view as coordinates of a vector x_n ; it combines them linearly multiplying weights and adding a bias; and then applies a (mostly) non-linear function a to the linear combination.

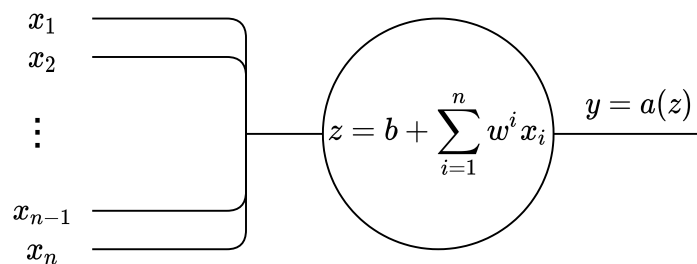


Figure 2.1: Perceptron scheme.

Some other popular design are convolutional neurons, which are used in image recognition; and memory cells, which are used in data with time dependencies such as video processing.

The way in which we combine artificial neurons to form artificial neural networks is such that the outputs of group of neurons become the input of another group of neurons. In this light, one can think of an artificial neural network as a directed graph with entry and exit edges, where each of the nodes correspond to a neuron and the directed edges indicate which neurons outputs feed another neuron as inputs. This is why many times the neurons in a network are referred to as nodes.

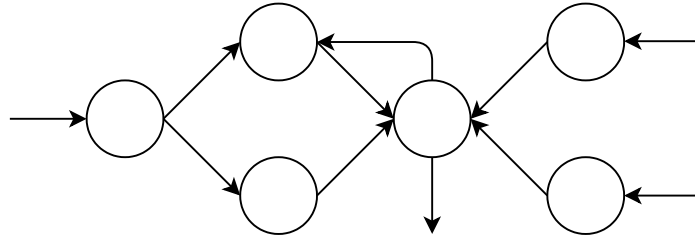


Figure 2.2: A directed graph which could be a possible representation of the architecture of an artificial neural network. Nodes are artificial neurons and edges indicate which neurons feed into each other.

One approach for organizing these neurons is using a feed-forward scheme. By this scheme we divide the neurons into sequential layers (groups). Then the outputs of a layer can only become inputs of the next layer. Using the graph characterization, this would correspond to artificial neural networks defined by directed graphs without cycles. The main advantage of this scheme is it allows for the use standard back-propagation algorithm (which we will be explaining in the next sections) to train the parameters in the neurons to fit a certain model.

A particular case of feed-forward artificial neural networks is the fully-connected. This happens when all the neurons in a layer are connected to all in the neurons in the next layer.

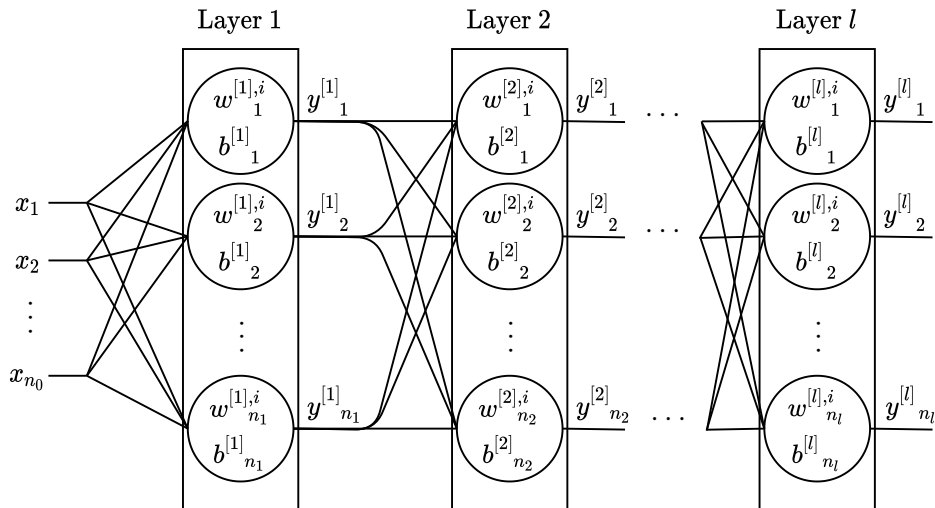


Figure 2.3: General scheme of a perceptron based fully-connected feed-forward artificial neural network.

Throughout this work we will be exclusively using perceptron fully-connected feed-forward artificial neural networks with different number of layers and different number of nodes per layer to approximate the solutions of differential equations. Figure 2.3 shows the standard graph representation of the structure of such neural networks based on the concepts explained up to this point. Note that the edges are not directed, as it is unnecessary, since by standard convention the flow of the neurons goes from left to right.

An artificial neural network provides a frame to define complex parametric functions in a modular way as a composition of simpler operations encapsulated in artificial neurons (we will insist in this point the next sections).

The following observations are a way to better understand artificial neural networks:

- A single neuron can form a neural network. In that case, if the activation function is identity, and we adjust the parameters of the neuron so the outputs fit a continuous dataset, we perform a linear regression. Indeed w^i and b are simply the slope and the intercept. Similarly, if the activation function is a sigmoid, and we adjust the parameters of the neuron so the outputs fit a binary dataset, we perform a logistic regression.
- In more complex deep neural networks, we kind of expand the same ideas as in the single neuron. In general, in artificial neural networks, what happens in regression problems is that we fit the parameters to make the hyper-surface defined by network structure get as close as possible to the sample data; and in classification problems we fit the parameters as much as possible to match the underlying marginal distribution of each category with the network structure.

Finally we will end this introductory section with some nomenclature for the rest of the work:

- A feed-forward neural network is said to be deep when it contains more than one layer.
- In deep neural networks, layers are classified as input, hidden and output layers. The input layer is the one that simply takes in the inputs and does no transformations; the output layer is the last layer of the sequence of layers and its outputs are the output results of all the whole network; and the hidden layers are all that lie between the input and output layers. Based on Figure 2.3, the input vector corresponds to the input layer (layer 0, even though it is not explicitly referenced as that), layers 1 to $l - 1$ would be the hidden layers and layer l would be the output layer.
- A feed-forward fully-connected network is defined by its number of layers, the number of neurons (nodes) in each layer, and its types of neurons. A neuron is defined by its weights, bias and activation function. Therefore, from now on we will use the following standard nomenclature, which contains all the elements that we need, to completely define our networks:

ℓ	layers indexes,
n_ℓ, m_ℓ	neurons indexes for layer ℓ ,
$w_{n_\ell}^{[\ell] m_{\ell-1}}$	weights of the neuron n_ℓ in the layer ℓ (parameters),
$b_{n_\ell}^{[\ell]}$	bias of the neuron n_ℓ in the layer ℓ (parameters),
$a_{n_\ell}^{[\ell]}$	activation function of the neuron n_ℓ in the layer ℓ ,
$z_{n_\ell}^{[\ell]}$	result of the linear combination of the neuron n_ℓ in the layer ℓ ,
$y_{n_\ell}^{[\ell]}$	result of applying the activation function of the neuron n_ℓ in the layer ℓ .

By this notation we will consider the input vector as layer 0 for which no transformations are performed, therefore, $x_i = z_i^{[0]} = y_i^{[0]}$. Also, the output layer (the n -th layer) never has an activation function, therefore $\hat{u}_i = z_i^{[n]} = y_i^{[n]}$.

2.2 From Numerical Integration to Deep-Learning

The intended use of the artificial neural networks in this work is for them to approximate the solution of some initial/boundary value problem. In order to achieve this, the parameters of the network should be tweaked to minimize some measure representing how well the network satisfies the differential operator and initial/boundary conditions.

Given the artificial neural network approximation of the solution, $\hat{u}(x; w, b)$, which depends parametrically on the set of all weights w and all biases b , we can define a positively defined loss or cost function, $L(w, b)$, that quantifies the degree of satisfaction of the network to the problem, in the following terms:

$$L(w, b) = L_1(w, b) + L_2(w, b) + R(w, b), \quad (2.1)$$

where $L_1(w, b)$ is the loss term measuring how well the neural network approximates the differential operator (1.1), $L_2(w, b)$ is the loss term measuring how well the neural network approximates the initial/boundary conditions (1.2-1.5), and $R(w, b)$ is the regularization term of the loss which is a term that will help to stabilize and improve the convergence in the optimization (this terms will be covered in a later section). In particular, using Cauchy boundary conditions, which are the most complex of all, the terms would be:

$$L_1(w, b) = \left\| \mathcal{L}[\hat{u}(x; w, b)] - f(x) \right\|_{\Omega, 2} = \int_{\Omega} \left(\mathcal{L}[u(x; w, b)] - f(x) \right)^2 dx, \quad (2.2)$$

$$\begin{aligned} L_2(w, b) &= \left\| \hat{u}(x; w, b) - g_1(x) \right\|_{\Gamma, 2} + \left\| \frac{\partial \hat{u}(x; w, b)}{\partial n(x)} - g_2(x) \right\|_{\Gamma, 2} \\ &= \int_{\Gamma} \left(\hat{u}(x; w, b) - g_1(x) \right)^2 dx + \int_{\Gamma} \left(\frac{\partial \hat{u}(x; w, b)}{\partial n(x)} - g_2(x) \right)^2 dx, \end{aligned} \quad (2.3)$$

with $\|\cdot\|_{\Omega, 2}$ the norm of the $L^2(\Omega)$ Hilbert space (which is the space of square integrable functions in Ω), and the same concept applies to $\|\cdot\|_{\Gamma, 2}$. However, in practice, as the integrals in (2.2-2.3) are virtually impractical to compute, instead of using the $\|\cdot\|_{\Omega, 2}$ and $\|\cdot\|_{\Gamma, 2}$ norms, a discrete approximation is used. This is achieved by taking a random collocation of N_{Ω} points in Ω and N_{Γ} points in Γ , which can be obtained by using a Monte Carlo hit-and-miss approach, and discretizing as $\int_{\Omega} \rightarrow 1/N_{\Omega} \sum_{N_{\Omega}}$ and $\int_{\Gamma} \rightarrow 1/N_{\Gamma} \sum_{N_{\Gamma}}$. Thus the actual loss terms become:

$$L_1(w, b) \approx \frac{1}{N_{\Omega}} \sum_{i \in N_{\Omega}} \left(\mathcal{L}[u(x_i; w, b)] - f(x_i) \right)^2, \quad (2.4)$$

$$L_2(w, b) \approx \frac{1}{N_{\Gamma}} \sum_{i \in N_{\Gamma}} \left(\hat{u}(x_i; w, b) - g_1(x_i) \right)^2 + \frac{1}{N_{\Gamma}} \sum_{i \in N_{\Gamma}} \left(\frac{\partial \hat{u}(x_i; w, b)}{\partial n(x_i)} - g_2(x_i) \right)^2. \quad (2.5)$$

Observe that we have transformed the continuous loss functions into the MSE (mean squared error) on a random collocation of points of the domain and the border. Now, on these premises, the problem has changed in nature, from a numerical integration problem, to an almost purely deep-learning regression type of problem. Other discretizations using the absolute error or the Huber error would have yielded equally valid approximations.

The next steps are straightforward, train the artificial neural network as done in any other deep-learning regression type of problem. This entails, taking the loss $L(w, b)$ which defines a hyper-surface in the parameter space, and given a certain initialization of parameters (w_0, b_0) , use a gradient-based optimization technique to minimize the loss function which is equivalent to minimizing the error of the approximation.

2.3 Derivatives: Back Propagation and Gradient Issues

As we have just stated we would like to use some kind of gradient-based optimization technique to optimize $L(w, b)$. The idea of this kind of techniques is to start at some point and move that point in steps along the direction of the gradient until we reach the minimum or sufficiently low value of the loss function. This strategy implies that in every step we are required to compute the gradient $\nabla_{(w,b)}L(w_0, b_0)$ at the point we are in.

Using a feed-forward scheme allows for the application of “standard” back-propagation algorithms to compute the derivatives of these gradients, which other schemes do not allow, or require of many modifications. In actuality, back-propagation is the name that receives in deep-learning a backward accumulation automatic differentiation algorithm (or autodiff), which is a type of dynamic programming algorithm used to compute derivatives of composed functions efficiently in computers. [26]

We will illustrate how a backward autodiff works with the following example, imagine we wanted to calculate in a computer the derivative in (x_0, y_0) of:

$$f(x, y) = \sin(e^x + xe^y). \tag{2.6}$$

What an autodiff algorithm does is decompose progressively the function into primitive operations (sums, subtractions, multiplications, divisions and basic functions: exponentials, sinus, cosines...). A computer can take any of this primitive operations and compute their values and the values of its derivatives at any point with arithmetic precision. For (2.6), the decomposition would be:

$$\begin{aligned}
 z_1 = \sin(z_2) &\longrightarrow \frac{\partial z_1}{\partial z_2} = \cos(z_2), & z_5 = x &\longrightarrow \frac{\partial z_5}{\partial x} = 1 \\
 z_2 = z_3 + z_4 &\longrightarrow \frac{\partial z_2}{\partial z_3} = 1, \frac{\partial z_2}{\partial z_4} = 1 & z_6 = x &\longrightarrow \frac{\partial z_6}{\partial x} = 1 \\
 z_3 = e^{z_5} &\longrightarrow \frac{\partial z_3}{\partial z_5} = e^{z_5}, & z_7 = e^{z_8} &\longrightarrow \frac{\partial z_7}{\partial z_8} = e^{z_8} \\
 z_4 = z_6 \cdot z_7 &\longrightarrow \frac{\partial z_4}{\partial z_6} = z_7, \frac{\partial z_4}{\partial z_7} = z_6, & z_8 = y &\longrightarrow \frac{\partial z_8}{\partial y} = 1.
 \end{aligned} \tag{2.7}$$

From a implementation stance, this decomposition is used to build what is known as a computational graph, which is a directed graph such that, its nodes store the computed numerical values of the primitive operations z_i and their derivatives at the point of calculation (x_0, y_0) , and its directed edges store the precedence of all operations. Figures (2.4-2.4) show the computational graph of (2.6). In Tensorflow v2.3 which is the deep-learning framework that we will be using in the next chapter to build artificial neural networks, computational graphs are natively implemented in a class called *tape*, which automatically builds these graphs.

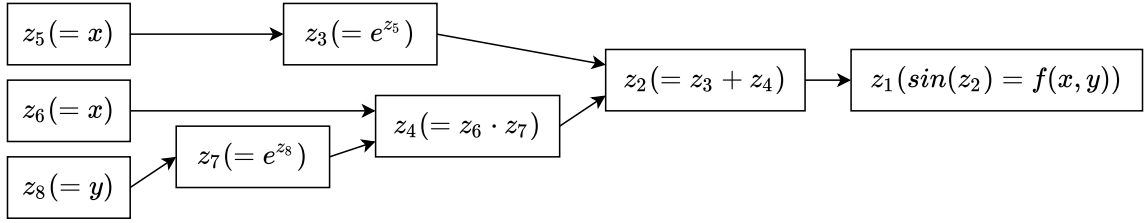


Figure 2.4: Computational graph of example (2.6).

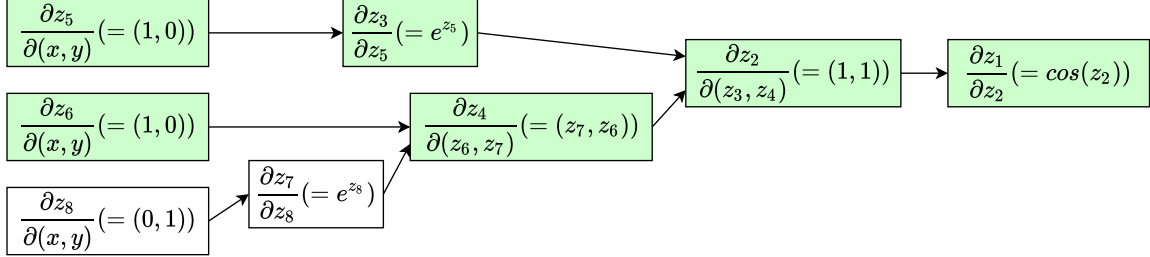


Figure 2.5: Computational graph (derivatives) of example (2.6). In green the flow of nodes required to compute $\partial f(x, y)/\partial x$

The computation of the derivatives at the point (x_0, y_0) is done by multiplying and adding (when more than one edge enters a node) the values of derivatives (actually, the Jacobians) at the nodes following the directions of the edges. This algorithm in its backward version is exactly the same as successively applying the chain rule in order, from the most fundamental operation to the most composed one, in a systematic way. Its forward version applies the chain rule in the other order. For (2.6) the algorithm would be the same as using the chain rule as follows:

$$\begin{aligned}
 \frac{\partial f(x, y)}{\partial(x, y)} &= \left(\left(\frac{\partial z_1}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_3} \right) \cdot \frac{\partial z_3}{\partial z_5} \right) \cdot \frac{\partial z_5}{\partial(x, y)} \\
 &+ \left(\left(\frac{\partial z_1}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_4} \right) \cdot \frac{\partial z_4}{\partial z_6} \right) \cdot \frac{\partial z_6}{\partial(x, y)} \\
 &+ \left(\left(\left(\frac{\partial z_1}{\partial z_2} \cdot \frac{\partial z_2}{\partial z_4} \right) \cdot \frac{\partial z_4}{\partial z_7} \right) \cdot \frac{\partial z_7}{\partial z_8} \right) \cdot \frac{\partial z_8}{\partial(x, y)}.
 \end{aligned} \tag{2.8}$$

The dynamic programming element to this algorithm is very helpful as it allows to use the same computation graph of a given point to compute with no additional cost the derivatives with respect to intermediate operations and parameters (which will be extremely useful). It also provides better errors and stability than numerical differentiation which approximates derivatives by differences. This is also different from symbolic differentiation. [27]

2.3.1 Derivatives Behaviour (Vanishing and Exploding Gradients)

In this subsection we will compute the derivatives of an example artificial neural network and look at how they behave. Although by using the back propagation algorithm we are able to numerically compute derivatives up to any order with great precision by successively applying the algorithm to its outputted derivatives, this does not imply that the derivatives we compute (in a structural sense) are well suited to train the model. Differences in magnitude and sensitivity to variations of the parameters can make training almost impossible.

When computing the gradient of the loss function with respect to the parameters to train the artificial neural network, we obtain the following:

$$\nabla_{(w,b)} L(w,b) := \frac{\partial L(w,b)}{\partial(w,b)} = \frac{\partial L_1(w,b)}{\partial(w,b)} + \frac{\partial L_2(w,b)}{\partial(w,b)} + \frac{\partial R(w,b)}{\partial(w,b)}, \quad (2.9)$$

where, in particular,

$$\frac{\partial L_1(w,b)}{\partial(w,b)} = \frac{\partial \mathcal{L}[\hat{u}(x;w,b)]}{\partial(w,b)} - \frac{\partial f(x)}{\partial(w,b)}. \quad (2.10)$$

and thus, we observe from the term in L_1 that, in order to obtain the gradient, we will require the derivatives of up to first order for (w,b) , and up to the highest order in the differential operator \mathcal{L} for (x) , which in this work can be of up to order 2 from the Laplacian operator.

Next, we will use an example artificial neural network consisting of 4 layers: an input layer of 2 neurons, two hidden layers of 3 and 4 neurons respectively, and an output layer of 2 neurons (see Figure 2.6). Using this example for greater clarity, we will first characterize its elements in detail, which will serve as a consolidation example of the first section of this chapter. Then, we will be applying the chain rule to express some derivatives of up to order 2 in variables and order 1 in parameters, and we will discuss their potential issues when used in gradient based training, namely the issues with vanishing and exploding gradients. This whole subsection will be written in the tensor formulation introduced in Chapter 0. In Appendix A, we reproduce the results in this subsection that make sense in linear algebra notation, for didactic reasons.

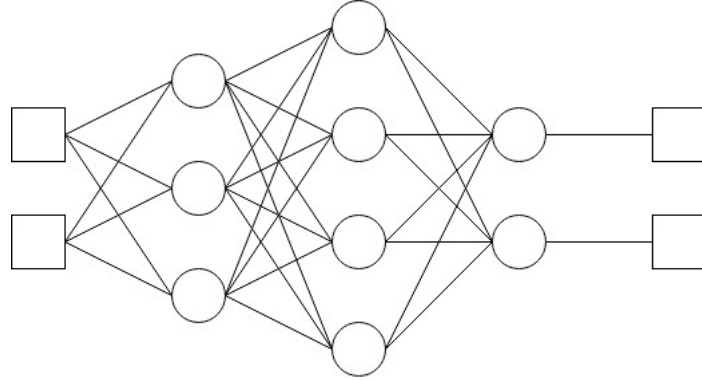


Figure 2.6: Example model: A 2-3-4-2 artificial neural network.

Given our example artificial neural network (Figure 2.6), the elements that characterize it (variables, parameters and activation functions), in the most general way for perceptron type neurons, are:

$$\begin{aligned} n_0, m_0 &= 1, 2, & n_1, m_1 &= 1, 2, 3, & n_2, m_2 &= 1, 2, 3, 4, & n_3, m_3 &= 1, 2, \\ w_{n_1}^{[1] m_0}, w_{n_2}^{[2] m_1}, w_{n_3}^{[3] m_2}, & z_{n_0}^{[0]} = x_{n_0}, z_{n_1}^{[1]}, z_{n_2}^{[2]}, z_{n_3}^{[3]} = \hat{u}_{n_3}, & (2.11) \\ b_{n_1}^{[1]}, b_{n_2}^{[1]}, b_{n_3}^{[1]}, & y_{n_0}^{[0]} = x_{n_0}, y_{n_1}^{[1]}, y_{n_2}^{[2]}, y_{n_3}^{[3]} = \hat{u}_{n_3}, \\ a_{n_1}^{[1]}, a_{n_2}^{[2]} : \mathbb{R} &\longrightarrow \mathbb{R}. \end{aligned}$$

Recall that each of this elements is a tensor in itself.

In particular, for each of the layers ℓ , the output variables $y^{[\ell]}_{n_\ell}$ and $z^{[\ell]}_{n_\ell}$ of each neuron are (where remember that from now on everything is written in Einstein convention notation):

$$\begin{aligned} z^{[1]}_{n_1} &= b^{[1]}_{n_1} + w^{[1]}_{n_1}{}^{n_0} \cdot y^{[0]}_{n_0}, & y^{[1]}_{n_1} &= a^{[1]}_{n_1}(z^{[1]}_{n_1}), \\ z^{[2]}_{n_2} &= b^{[2]}_{n_2} + w^{[2]}_{n_2}{}^{n_1} \cdot y^{[1]}_{n_1}, & y^{[2]}_{n_2} &= a^{[2]}_{n_2}(z^{[2]}_{n_2}), \\ z^{[3]}_{n_3} &= b^{[3]}_{n_3} + w^{[3]}_{n_3}{}^{n_2} \cdot y^{[2]}_{n_2}, & y^{[3]}_{n_3} &= z^{[3]}_{n_3}. \end{aligned} \quad (2.12)$$

If we were compose all these output variables we can see that this instance of an artificial neural network actually represents the composite function:

$$y^{[3]}_{n_3} = b^{[3]}_{n_3} + w^{[3]}_{n_3}{}^{n_2} \cdot a^{[2]}_{n_2} \left(b^{[2]}_{n_2} + w^{[2]}_{n_2}{}^{n_1} \cdot a^{[1]}_{n_1} \left(b^{[1]}_{n_1} + w^{[1]}_{n_1}{}^{n_0} \cdot y^{[0]}_{n_0} \right) \right), \quad (2.13)$$

or equivalently,

$$\hat{u}_{n_3} = b^{[3]}_{n_3} + w^{[3]}_{n_3}{}^{n_2} \cdot a^{[2]}_{n_2} \left(b^{[2]}_{n_2} + w^{[2]}_{n_2}{}^{n_1} \cdot a^{[1]}_{n_1} \left(b^{[1]}_{n_1} + w^{[1]}_{n_1}{}^{n_0} \cdot x_{n_0} \right) \right). \quad (2.14)$$

Now, we will compute some derivatives applying the chain rule:

Derivatives of order 1 in x :

$$\begin{aligned} \frac{\partial \hat{u}_{n_3}}{\partial x_{m_0}} &:= \frac{\partial y^{[3]}_{n_3}}{\partial z^{[0]}_{n_0}} \\ &= \frac{\partial y^{[3]}_{n_3}}{\partial z^{[3]}_{m_3}} \cdot \frac{\partial z^{[3]}_{m_3}}{\partial y^{[2]}_{n_2}} \cdot \frac{\partial y^{[2]}_{n_2}}{\partial z^{[2]}_{m_2}} \cdot \frac{\partial z^{[2]}_{m_2}}{\partial y^{[1]}_{n_1}} \cdot \frac{\partial y^{[1]}_{n_1}}{\partial z^{[1]}_{m_1}} \cdot \frac{\partial z^{[1]}_{m_1}}{\partial y^{[0]}_{n_0}} \cdot \frac{\partial y^{[0]}_{n_0}}{\partial z^{[0]}_{m_0}} \\ &= \delta_{n_3}^{m_3} \cdot w^{[3]}_{n_3}{}^{n_2} \cdot D^{n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot \delta_{n_2}^{m_2} \cdot w^{[2]}_{n_2}{}^{n_1} \cdot D^{n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{m_1} \cdot w^{[1]}_{n_1}{}^{n_0} \cdot \delta_{n_0}^{m_0} \\ &= w^{[3]}_{n_3}{}^{n_2} \cdot D^{n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2}{}^{n_1} \cdot D^{n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1}{}^{m_0}. \end{aligned} \quad (2.15)$$

Derivatives of order 2 in x :

$$\begin{aligned} \frac{\partial^2 \hat{u}_{n_3}}{\partial x_{\bar{m}_0} \partial x_{m_0}} &:= \frac{\partial}{\partial z^{[0]}_{\bar{m}_0}} \left(\frac{\partial \hat{u}_{n_3}}{\partial x_{m_0}} \right) \\ &= w^{[3]}_{n_3}{}^{n_2} \cdot \frac{\partial}{\partial z^{[0]}_{\bar{m}_0}} D^{n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2}{}^{n_1} \cdot D^{n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1}{}^{m_0} \\ &\quad + w^{[3]}_{n_3}{}^{n_2} \cdot D^{n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2}{}^{n_1} \cdot \frac{\partial}{\partial z^{[0]}_{\bar{m}_0}} D^{n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1}{}^{m_0} \\ &= w^{[3]}_{n_3}{}^{n_2} \cdot D^{n_2, n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \\ &\quad \cdot \left(w^{[2]}_{n_2}{}^{\bar{n}_1} \cdot D^{\bar{n}_1} \left(a^{[1]}_{\bar{n}_1}(z^{[1]}_{\bar{n}_1}) \right) \cdot w^{[1]}_{\bar{n}_1}{}^{\bar{m}_0} \right) \cdot \left(w^{[2]}_{n_2}{}^{n_1} \cdot D^{n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1}{}^{m_0} \right) \\ &\quad + w^{[3]}_{n_3}{}^{n_2} \cdot D^{n_2} \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2}{}^{n_1} \cdot D^{n_1, n_1} \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1}{}^{\bar{m}_0} \cdot w^{[1]}_{n_1}{}^{m_0}. \end{aligned} \quad (2.16)$$

Derivatives of order 0 in x and order 1 in w :

$$\begin{aligned}
\frac{\partial \hat{u}_{n_3}}{\partial w^{[3]}_{m_3} n_2} &= \delta_{n_3}^{m_3} \cdot y^{[2]}_{n_2}, \\
\frac{\partial \hat{u}_{n_3}}{\partial w^{[2]}_{m_2} n_1} &= w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot \delta_{n_2}^{m_2} \cdot y^{[1]}_{n_1}, \\
\frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0} &:= \frac{\partial y^{[3]}_{n_3}}{\partial w^{[1]}_{m_1} n_0} \\
&= \frac{\partial y^{[3]}_{n_3}}{\partial z^{[3]}_{m_3}} \cdot \frac{\partial z^{[3]}_{m_3}}{\partial y^{[2]}_{n_2}} \cdot \frac{\partial y^{[2]}_{n_2}}{\partial z^{[2]}_{m_2}} \cdot \frac{\partial z^{[2]}_{m_2}}{\partial y^{[1]}_{n_1}} \cdot \frac{\partial y^{[1]}_{n_1}}{\partial z^{[1]}_{m_1}} \cdot \frac{\partial z^{[1]}_{m_1}}{\partial w^{[1]}_{m_1} n_0} \\
&= \delta_{n_3}^{m_3} \cdot w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot \delta_{n_2}^{m_2} \cdot w^{[2]}_{m_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{m_1} \cdot y^{[0]}_{n_0} \\
&= w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{m_1} \cdot y^{[0]}_{n_0}.
\end{aligned} \tag{2.17}$$

Derivatives of order 0 in x and order 1 in b :

$$\begin{aligned}
\frac{\partial \hat{u}_{n_3}}{\partial b^{[3]}_{m_3}} &= \delta_{n_3}^{m_3}, \\
\frac{\partial \hat{u}_{n_3}}{\partial b^{[2]}_{m_2}} &= w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot \delta_{n_2}^{m_2}, \\
\frac{\partial \hat{u}_{n_3}}{\partial b^{[1]}_{m_1}} &:= \frac{\partial y^{[3]}_{n_3}}{\partial b^{[1]}_{m_1}} \\
&= \frac{\partial y^{[3]}_{n_3}}{\partial z^{[3]}_{m_3}} \cdot \frac{\partial z^{[3]}_{m_3}}{\partial y^{[2]}_{n_2}} \cdot \frac{\partial y^{[2]}_{n_2}}{\partial z^{[2]}_{m_2}} \cdot \frac{\partial z^{[2]}_{m_2}}{\partial y^{[1]}_{n_1}} \cdot \frac{\partial y^{[1]}_{n_1}}{\partial z^{[1]}_{m_1}} \cdot \frac{\partial z^{[1]}_{m_1}}{\partial b^{[1]}_{m_1}} \\
&= \delta_{n_3}^{m_3} \cdot w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot \delta_{n_2}^{m_2} \cdot w^{[2]}_{n_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{m_1} \\
&= w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{m_1}.
\end{aligned} \tag{2.18}$$

Derivatives of order 1 in x and order 1 in $w^{[1]}$:

$$\begin{aligned}
\frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} \bar{n}_0 \partial x_{m_0}} &:= \frac{\partial}{\partial w^{[1]}_{m_1} \bar{n}_0} \left(\frac{\partial \hat{u}_{n_3}}{\partial x_{m_0}} \right) \\
&= w^{[3]}_{n_3} \cdot D^{n_2, n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \\
&\quad \cdot \left(w^{[2]}_{n_2} \bar{n}_1 \cdot D^{\bar{n}_1} \left(a^{[1]}_{\bar{n}_1} (z^{[1]}_{\bar{n}_1}) \right) \cdot \delta_{\bar{n}_1}^{\bar{m}_1} \cdot y^{[0]}_{\bar{n}_0} \right) \cdot \left(w^{[2]}_{n_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot w^{[1]}_{n_1} \right) \\
&\quad + w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2} \cdot D^{n_1, n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_1}^{\bar{m}_1} \cdot y^{[0]}_{\bar{n}_0} \cdot w^{[1]}_{n_1} \\
&\quad + w^{[3]}_{n_3} \cdot D^{n_2} \left(a^{[2]}_{n_2} (z^{[2]}_{n_2}) \right) \cdot w^{[2]}_{n_2} \cdot D^{n_1} \left(a^{[1]}_{n_1} (z^{[1]}_{n_1}) \right) \cdot \delta_{n_0}^{m_0} \cdot \delta_{n_1}^{\bar{m}_1}.
\end{aligned} \tag{2.19}$$

The vanishing and exploding gradient problems are a recurrent issue in the deep-learning field, and in fact, it was the main reason for which, although the concept artificial neural networks was well known since the 1950s, it was not until the early 2000s when the effect of this issues were “mitigated” with new techniques, that the field of deep-learning exploded became what it is today. The solution up until now has been to design improved optimizers, use regularization techniques, increase computer power and make use of parallel computation, and create new architectures that preserve data structures (such as convolutional networks in image processing). In the following sections we will see the ideas for the typical optimizers and regularizations that are used to avoid the “classical” vanishing and exploding gradients.

From the derivatives with respect to w and b in (2.17) and (2.18), these gradient problems can be straight forward explained. Observe that in the derivatives of the parameters, as we go from the deepest (closest to the input) layer to the shallowest layers (closest to the output) of the network, the number of terms and derivatives of the activation functions increase. Thus, when the weights w become small (or large) with respect to 1, or the activation function has a small (or large) slope in the range of the linear combination of the neurons z , both of which are extremely common, then as we consider the derivatives of deeper parameters, these become the product of an increasing number of small (or large) terms. In conclusion, as we consider deeper layers, the derivatives used to train the parameters become vanishingly small or explosively large with respect to the shallower layers, making the optimization of these deeper parameters almost impossible. Given our problem (only in w , but equivalent in b):

$$\begin{aligned}
\text{“Calssic” Vanishing Gradient: } & \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[3]} \frac{n_2}{m_3}} \right| \gg \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[2]} \frac{n_1}{m_2}} \right| \gg \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]} \frac{n_0}{m_1}} \right|. \\
\text{“Calssic” Exploding Gradient: } & \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[3]} \frac{n_2}{m_3}} \right| \ll \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[2]} \frac{n_1}{m_2}} \right| \ll \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]} \frac{n_0}{m_1}} \right|.
\end{aligned} \tag{2.20}$$

At this point, we will refer to the vanishing and exploding gradient problems that we have just explained as “classical” because it is the one that always takes place in any deep-learning classification and regression problem. However, by the very nature of this problem, being that the loss function contains derivatives of artificial neural network in x , we incur in a second type of vanishing and exploding gradient problems which is harder to see. If the classical vanishing/explode type implies that the derivatives of the artificial neural network \hat{u} with respect to the parameters w and b increase/decrease as we consider parameters in deeper layers; in this second type, for higher order derivatives of \hat{u} in x , their derivatives with respect to the same parameter of the same layer also vanishes or explodes as we consider higher derivatives in x . Given our problem (only in w , but equivalent in b):

$$\begin{aligned}
\text{“Second type” Vanishing Gradient: } & \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[l]} \frac{n_2}{m_3}} \right| \gg \left| \frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[l]} \frac{n_1}{m_2} \partial x_{\bar{m}_0}} \right| \gg \dots \\
\text{“Second type” Exploding Gradient: } & \left| \frac{\partial \hat{u}_{n_3}}{\partial w^{[l]} \frac{n_2}{m_3}} \right| \ll \left| \frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[l]} \frac{n_1}{m_2} \partial x_{\bar{m}_0}} \right| \ll \dots
\end{aligned} \tag{2.21}$$

Showing this second type of gradient problem a bit less obvious than the classic one. To be able to give an intuition of the problem we will consider the case in which all the activation functions are exponential $a^{[\ell]}_{n_\ell}(x) = e^x$. In such case the derivatives of the activation function are simply:

$$\begin{aligned} D^{n_\ell} \left(a^{[\ell]}_{n_\ell}(z^{[\ell]}_{n_\ell}) \right) &= a^{[\ell]}_{n_\ell}(z^{[\ell]}_{n_\ell}) \cdot \mathbb{1}^{n_\ell}, \\ D^{n_\ell, n_\ell} \left(a^{[\ell]}_{n_\ell}(z^{[\ell]}_{n_\ell}) \right) &= a^{[\ell]}_{n_\ell}(z^{[\ell]}_{n_\ell}) \cdot \mathbb{1}^{n_\ell, n_\ell}, \end{aligned} \quad (2.22)$$

where $\mathbb{1}^{n_\ell}$ is the tensor whose every component is 1, and $\mathbb{1}^{n_\ell, n_\ell} = \mathbb{1}^{n_\ell} \cdot \mathbb{1}^{n_\ell}$. Then, for an exponential activation function, the derivatives (2.17) and (2.19) of the weights $w^{[1]}$ in the first layer become:

$$\frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0} = w^{[3]}_{n_3} \cdot \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot \mathbb{1}^{n_2} \cdot w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot \delta_{n_1}^{m_1} \cdot y^{[0]}_{n_0}, \quad (2.23)$$

$$\begin{aligned} \frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0 \partial x_{m_0}} &= w^{[3]}_{n_3} \cdot \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot \mathbb{1}^{n_2, n_2} \\ &\cdot \left(w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot \delta_{n_1}^{m_1} \cdot y^{[0]}_{n_0} \right) \cdot \left(w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot w^{[1]}_{n_1}^{m_0} \right) \\ &+ w^{[3]}_{n_3} \cdot \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot \mathbb{1}^{n_2} \cdot w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1, n_1} \cdot \delta_{n_1}^{m_1} \cdot y^{[0]}_{n_0} \cdot w^{[1]}_{n_1}^{m_0} \\ &+ w^{[3]}_{n_3} \cdot \left(a^{[2]}_{n_2}(z^{[2]}_{n_2}) \right) \cdot \mathbb{1}^{n_2} \cdot w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot \delta_{n_0}^{m_0} \cdot \delta_{n_1}^{m_1}, \end{aligned} \quad (2.24)$$

Since the derivative of the exponential is essentially itself, and the activation function tensor is made of copies of exponentials for which differentiating means multiplying $\mathbb{1}^{n_\ell}$ tensors, both expressions (2.23) and (2.24) are written in the same terms, hence are easy to compare. As the tensor operations (sums and products) are commutative; and given a tensor $T_k^{i,j}$, which lets say has non-zero components for simplicity, we can define a “pseudo-inverse” to the contraction $(T_k^{i,j})^{-1} = [T^{-1}]_{i,j}^k$ (element-wise) such that it holds $T_k^{i,j} \cdot [T^{-1}]_{i,j}^k = i \cdot j \cdot k$ (if there are zero values, we would have we would have to fix an inverse element for the zero components and discount the zeroes from the count $i \cdot j \cdot k$; here we will assume there are no zeroes for clarity); then we can easily replace (2.23) in (2.24) yielding:

$$\begin{aligned} \frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0 \partial x_{m_0}} &= \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0} \cdot \left(\mathbb{1}^{n_2} \cdot w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot w^{[1]}_{n_1}^{m_0} \right) \\ &+ \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0} \cdot \mathbb{1}^{n_1} \cdot w^{[1]}_{n_1}^{m_0} + \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0} \cdot \delta_{n_0}^{m_0} \cdot (\bar{n}_0)^{-1} \cdot \left(y^{[0]}_{n_0} \right)^{-1}, \end{aligned} \quad (2.25)$$

and grouping,

$$\begin{aligned} \frac{\partial^2 \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0 \partial x_{m_0}} &= \left(\mathbb{1}^{n_2} \cdot w^{[2]}_{n_2} \cdot \left(a^{[1]}_{n_1}(z^{[1]}_{n_1}) \right) \cdot \mathbb{1}^{n_1} \cdot w^{[1]}_{n_1}^{m_0} \right. \\ &\left. + \mathbb{1}^{n_1} \cdot w^{[1]}_{n_1}^{m_0} + \delta_{n_0}^{m_0} \cdot (\bar{n}_0)^{-1} \cdot \left(y^{[0]}_{n_0} \right)^{-1} \right) \cdot \frac{\partial \hat{u}_{n_3}}{\partial w^{[1]}_{m_1} n_0}. \end{aligned} \quad (2.26)$$

From (2.26) we see that, in the case of exponential activation functions, we can write the derivatives of a given weight and input in terms of a lower order derivative in terms of the input (this can be done for any order of the input). Now, regarding the factor in parenthesis in (2.26), assuming that the input is normalized and we do not shuffle data: the third term behaves as a very large constant ($\gg 1$) which increases with the input dimension; when the weights are small, the second term is the largest, thus since the weights are small, the whole factor is small ($\ll 1$); and when the weights are large, it is the first term that dominates, making the whole factor very large ($\gg 1$). This creates the vanishing/exploding gradient effect described previously in (2.21). Also, this same analysis can be carried out with respect to any other weight or bias parameter making a few changes, but in case of considering other activation functions, although the end conclusion is the same and can be empirically visualized, the analytic study becomes much harder.

This issue is important when considering operators with that contain derivatives of different order or products of derivatives, which is always ignored. For example, the laplacian operator which contain only additions of second order derivatives does not have this problem as the effect of no derivative vastly dominates the effect of another when deriving over the parameters, but the burger's operator presents it.

2.4 Optimizers

Recall that in section 2.2 we transformed the problem of approximating the solution of an initial/boundary problem into a non-convex optimization problem, whereby we had to find the minimum (or a sufficiently small value) of a loss function $L(w, b)$, defined by (2.1), (2.4) and (2.5). On that section we went on to anticipate that, in order to do that, we would be using a gradient based optimization technique (optimizer), which required the computation of the loss derivatives with respect to the parameters, i.e. the gradient $\nabla_{(w,b)}L(w, b)$. This had led to section 2.3 where we explained the algorithm of back propagation to compute such derivatives and the discussions of their potential problems, namely the vanishing/exploding gradients. Now everything is set up and it is finally time to get into the process of actually making the parameters of the artificial neural network approximate the solution (optimizing the loss function), which in deep-learning jargon is known as the training process.

The following discussion will be devoted to explaining the design of some of the most important gradient based optimizers used in deep-learning, and the ones we will be using in this work. These optimizers are the so called methods of steepest descent or methods of gradient descent, which are a family of methods used to solve general non-linear (convex or non-convex) unrestricted optimization problems. Intuitively, the idea behind these methods relies on thinking of the loss function as a hyper-surface $L(w, b) : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$, where $w \in \mathbb{R}^n$ and $b \in \mathbb{R}^m$. Then, starting at some (w_0, b_0) , initial point, the method goes on to calculate new points which should reduce the loss function value by moving within a certain rate, η , named the *learning rate*, in the direction of $\nabla_{(w,b)}L(w, b)$. The typical analogy for this idea is thinking of it as having a ball (initial point), and letting it roll downhill along the slope (the direction of the gradient) until it reaches the bottom.

As simple as these methods look conceptually, in practice it is not that easy to reach the minimum. If we were to apply one of these methods to a linear or quadratic bowl loss function $((ax + b)^2, a > 0)$, we are guaranteed that the gradient at any point would always point in the direction of the only existing minimum, thus given adequate learning rates, these methods would have perfect convergence. However, with almost every other loss function, the direction of steepest descent (gradient) does not necessarily point to the global minimum. Moreover, if the problem is non-convex, as all the ones we will be considering in this work, we are almost guaranteed that there are many local minimums, and the direction of steepest descent may lead the method to a local minimum and not the global one.

Another tricky issue for these methods is the presence of saddle or “saddle-like” regions of the loss function. These are regions for which we have very small derivatives of the gradient in certain directions, and very large in others. Visualizing these regions in the loss function, they resemble to, and thus are often called, “valleys”. What happens in these areas is that, in the ball analogy, the ball starts oscillating up and down along the valley’s walls (directions of large value derivatives) but is unable to make any progress across the valley (directions of low value derivatives). When using these steepest descent, this “saddle-like” region effect, as well as the effect of not being unable to escape a local minimum, is often reflected in the method when the point and loss function start oscillating between the same two very similar values.

These are the main three problems with steepest descent: the gradient not pointing in the direction of the global minimum; getting trapped in a local minimum; and stagnating when passing through “saddle-like” regions of the loss function. In order to avoid or mitigate these issues as much as possible, there are also three measures that can be applied: choosing a good initialization (starting point); applying some regularization technique, which somewhat has the effect of smoothing the loss function; and adjusting “properly” the learning rate at each step. In the next sections we will be looking at the initialization (which is tightly related to the selection of activation function), and the regularization techniques. For the rest of this section we will see different designs of steepest descent methods which adjust the learning rates for every step based on different ideas. We will divide these designs into first order if they require only the gradient, and second order if they also require estimates of the curvature.

For an extensive qualitative survey on gradient based methods [28] has a good coverage; in particular, in Table I and Table II there is a very complete comparison among first and higher order methods respectively. Other non gradient based methods are quite rare, for instance, in [29] a bio-inspired approach is used: a population of artificial neural networks is generated using different weights and architectures (hyper-parameters); the networks get tested and ranked by complexity and performance; then, a new population is generated based on the best performing networks with small alterations; and the process gets repeated.

2.4.1 First Order Methods

As we have already explained, these methods only depend on the gradient. The idea behind being so many variations is to have the method correct its learning rate by keeping some kind of memory of the gradients at previous points (steps) to improve convergence [30]. Next, we will discuss this methods grouping them in the following categories, from least to most refined:

- Vanilla (No Learning Rate Correction)
- Momentum Learning Rate Correction
- Component Learning Rate Adaptation
- Momentum + Component Learning Rate Adaptation

Vanilla (No Learning Rate Correction)

This group is the simplest and easiest to implement. It is actually the plain idea we have just explained, thus at every new step $t + 1$ we update the previous point with the formula:

$$(w_{t+1}, b_{t+1}) \longrightarrow (w_t, b_t) - \eta \nabla_{(w,b)} L(w_k, b_k). \quad (2.27)$$

Generally, if the batches of input data (here the random collocation of points in Ω) is large, computing $\nabla_{(w,b)} L(w_t, b_t)$ in every step can be computationally very expensive. Recall that the factors in the loss function are of the form $1/N \sum_{i=1}^N (\dots)$, which means that if N is large, in every step we have to compute a large sum of sub-gradients, $1/N \sum_{i=1}^N \nabla_{(w,b)} (\dots)$, which can be costly. The fix is to take what is known as an stochastic or on-line approach, this is, to divide the input data into partitions, $\{x_i\}_{i \in M_1}, \dots, \{x_i\}_{i \in M_k}$, with $M_1 + \dots + M_k$, and compute the gradient in every step just for the data of one of those partitions. If the input data is well randomized into the partitions, and the partitions are rotated consistently at every step, the differences in the gradient from not using the whole batch should be evened out throughout the many steps. When partitions of more than one data input are used the method is called **mini-batch gradient descent**, when the partitions contain a single data input the method is called **stochastic gradient descent** (SGD), and when the full batch is used the method is simply called **gradient descent** (GD). Often times, no distinction is made between mini-batch and stochastic, and both get referred to as stochastic gradient descent.

This stochastic approach can be applied to all the variations that we will be seeing next. In this work, however, we will not be sampling very large input data batches and the artificial neural networks will not be very large either, so we will always take full-batch approaches.

In terms of usage one would think that these version being the least refined would also be the least used, but it is far from the truth. It is true that the learning rate, η , has to be manually adapted in every step, which requires much try-and-error experimentation. This is done by setting a **learning schedule**, which is the set of instructions on how to vary the learning rate (for example, one could be as follows: for the first 1000 steps use $\eta = 0.001$, then every 1000 steps reduce $\eta/10$). Nonetheless, in recent times, and specially for artificial neural networks with large amounts of parameters (something that happens in very deep perceptron neural networks, or in convolutional networks by design), there have been many papers that claim plain SGD (or at most the momentum we will be seeing next) can outclass any other variation that we will see here. The strategy in these papers is to use an unusually large learning rate, which means moving too far in the direction of the gradient and straying from the optimal path of minimal loss value, in order to create an **annealing** effect [31]. This annealing effect is a direct parallel from its homonym in metallurgy. Using these very “long jumps” allows for greater mobility for the point we are at in the method, giving it the capacity to get over “walls” and explore the loss hyper-surface to get into a better region, before switching to the regular small learning rate strategy used to achieve convergence. This works the same way as heating a metal to allow for greater mobility of its molecules, and then letting them settle by cooling the metal. Adding noise to the gradient has been for a long time an extremely successful regularization technique in deep-learning problems following the same principle of annealing of adding some exploration component. However, this concept takes it further, the objective being achieving **superconvergence**, which happens when entering in a very good region where the method suffers a drastic drop in its loss value, and convergence can be obtained many orders of magnitude faster than with a standard approach. In [32] successive cycles of short and long learning rates are used to obtain superconvergence, and [33] develop an adapted version called **SGD with Entropy** following these same ideas.

As we will see the idea behind the next variations is to speed up the method by auto-adjusting the initial learning rate at every step. This implies less tweaking of the learning rates as the method will reduce it natively when the loss is worsening to stay in the right track, and increase it when the loss is improving to go faster. This also make these variations incompatible with superconvergence, as in the first step where the loss worsens, the method will immediately damp the learning rate.

Momentum Learning Rate Correction

Adding momentum to correct the learning rate in GD is very old and one of the first improvements on GD, the idea being based on keeping the inertia. In the ball analogy, if a ball is located at certain point but was carrying some velocity in a some direction, at that point it would not stop cold and resume its movement following the steepest descent. The ball will combine its previous inertia with the movement defined by the slope it is in. This is the idea behind **classical momentum** (CM), where the effective change v_{t+1} at the step $t + 1$ is not only given by the gradient at that point, but also by a certain proportion μ by the effective change of the previous step v_t :

$$\begin{aligned} v_{t+1} &\longrightarrow \mu v_t - \eta \nabla_{(w,b)} L(w_t, b_t), \\ (w_{t+1}, b_{t+1}) &\longrightarrow (w_t, b_t) + v_{t+1}. \end{aligned} \tag{2.28}$$

Intuitively, if we were are moving in a very consistent direction through the loss hypersurface, the inertial term from the previous step v_t adds to the gradient making larger jumps in that direction. Conversely if the direction suddenly changes, v_t dampens the jump as we might have overstepped into a bad area by taking to large of a jump in the previous step. A second variation of this idea is the **Nesterov's accelerated gradient** (NAG), which tend to yield better results than CM. The difference is that in NAV we look at the gradient not at the point we are in, but in a point projected ahead as if we had done a second jump in the previous step:

$$\begin{aligned} v_{t+1} &\longrightarrow \mu v_t - \eta \nabla_{(w,b)} L((w_t, b_t) + \mu v_t), \\ (w_{t+1}, b_{t+1}) &\longrightarrow (w_t, b_t) + v_{t+1}. \end{aligned} \tag{2.29}$$

This subsection is based on the article [34].

Component Learning Rate Adaptation

While momentum uses the information of the previous gradient to speed up or slow down the method when there is consistent or changing behaviour, it does little to move forward in saddle regions. Recall that this kind of regions occur when some derivatives are several order of magnitude larger than others, i.e. some components of the gradient are much larger than others, which can be caused by vanishing or exploding gradient problems. In these cases we cannot increase the global learning rate to take longer jumps in the flatter directions because this would also make the jumps longer in the steeper directions, which require shorter steps to not stray from the convergence path. Also, momentum cannot help either, as it only adds up on the previous gradient, which is still small for the flatter directions. The solution is to rescale the learning rate for each component in the gradient individually based on previous gradients. So, if there have been directions which have had consistently small derivatives, we want to take larger jumps just in those directions, and conversely for directions which have had consistently large derivatives, we want to make smaller jumps to not to overstep out of the convergence path.

The first method that we are going to review is the **Adaptative gradient Algorithm** (AdaGrad). In its original paper [35], the method is presented as follows:

$$\begin{aligned} G_t &= \sum_{\tau=1}^t (\nabla_{(w,b)} L(w_\tau, b_\tau)) \cdot ((\nabla L_{(w,b)}(w_\tau, b_\tau))^\top) \in R^{n+m \times n+m}, \\ (w_{t+1}, b_{t+1}) &\longrightarrow (w_t, b_t) - \eta (\text{diag}(G_t) + \varepsilon Id)^{-1/2} \nabla_{(w,b)} L(w_\tau, b_\tau). \end{aligned} \quad (2.30)$$

where G_t is the cumulative matrix of products of the past gradients, $\text{diag}(G_t)$ is the diagonal of such matrix, Id corresponds to the identity matrix, and ε is a small constant to avoid diving by zero. As the matrix G_t can be computed accumulatively and only its diagonal elements are used, we suggest rewriting the method in the following vectorized way:

$$\begin{aligned} \mathcal{G}_t &\longrightarrow \mathcal{G}_{t-1} + \nabla_{(w,b)} L(w_t, b_t) \odot \nabla_{(w,b)} L(w_t, b_t) \in R^{n+m}, \\ (w_{t+1}, b_{t+1}) &\longrightarrow (w_t, b_t) - \eta (\mathcal{G}_t + \varepsilon \mathbf{1})^{-1/2} \odot \nabla_{(w,b)} L(w_\tau, b_\tau). \end{aligned} \quad (2.31)$$

where products, inverses and root are all element-wise, and $\mathbf{1}$ is the vector consisting of all ones. Observe that if some direction of the gradient has consistently small derivatives, the cumulative value of \mathcal{G}_t will be small, and thus dividing by the square root of that value will increase the learning rate of direction (the inverse happens for components with large derivatives). This is sort of approximating the curvature in the principal directions by the values of its past gradients. However, this cumulative nature is this method's main problem, as we are constantly accumulating positive values, \mathcal{G}_t becomes increasingly large at each step, and since we are constantly dividing the learning rate by it, the methods halts the progress and is unable to scape local minima as time passes.

An improvement of AdaGrad comes with **AdaDelta**, [36], which mitigates the effect of the strong decay in learning rates of AdaGrad. Instead of using the accumulated information of all the squared previous gradients, it uses an exponential decay moving average of the square values of the gradient. This is instead of \mathcal{G}_t , it uses $E[(\nabla_{(w,b)} L(w_t, b_t))^2]$:

$$\begin{aligned} E[(\nabla_{(w,b)} L(w_t, b_t))^2] &= \rho E[(\nabla_{(w,b)} L(w_{t-1}, b_{t-1}))^2] \\ &+ (1 - \rho) \nabla_{(w,b)} L(w_t, b_t) \odot \nabla_{(w,b)} L(w_t, b_t), \end{aligned} \quad (2.32)$$

where ρ is the decay rate of the moving average. On top of this change, the method also adds a second idea. Since dividing by the square root of (2.32) is sort of a very brute approximation of dividing by the local curvature, in an attempt to resemble a second order Newton method, a term approximating the slope is multiplied. This term is a moving average of the squares of previous increments, $E[(\Delta(w_t, b_t))^2]$ which uses the same decay rate as before:

$$\begin{aligned} E[(\Delta(w_t, b_t))^2] &= \rho E[(\Delta(w_{t-1}, b_{t-1}))^2] \\ &+ (1 - \rho) \Delta(w_t, b_t) \odot \Delta(w_t, b_t), \end{aligned} \quad (2.33)$$

then the final algorithm at each step t works as:

$$\begin{aligned} &\text{Compute (2.32),} \\ \Delta(w_t, b_t) &\longrightarrow \frac{\sqrt{E[(\Delta(w_t, b_t))^2] + \varepsilon \mathbf{1}}}{\sqrt{E[(\nabla_{(w,b)} L(w_t, b_t))^2] + \varepsilon \mathbf{1}}}, \\ &\text{Compute (2.33),} \\ (w_t, b_t) &\longrightarrow (w_{t-1}, b_{t-1}) + \Delta(w_t, b_t). \end{aligned} \quad (2.34)$$

As a general note on the equations posed for this method, all the operations in (2.32-2.34) have been element wise. Finally, observe that the increment $\Delta(w_t, b_t)$ in (2.34) has in its numerator a term that approximates the slope and in its denominator a term that approximates the curvature, which tries to replicate a structure $H(f)^{-1} \cdot \nabla f$ of a Newton method.

A parallel developed, very popular and much simpler method than AdaDelta to solve the fast damping of AdaGrad is **RMSProp**. This is an unpublished method proposed in a Coursera course by Geoffrey Hinton, in lecture 6.5. [37]. This method was thought as an adaptation of the **RProp** which is a method originally designed for full-batches, to be able to account for mini-batches. This RProp method does not take into account the magnitude of the derivatives in the gradient, and instead, only takes into consideration the sign of the derivatives. Each direction learning rate is increased slightly every time the sign of its corresponding derivative is preserved, and drastically decreased whenever the sign of the derivative changes, everything within a certain threshold. When working with mini-batches this method can have many problems, as some sub-gradient may change in sign for some derivative due to the characteristics of that specific mini-batch, and not because the method has entered into region with a different behaviour. For instance, if the last 9 out of 10 derivatives in a direction have been positive and the only one has been negative, we do not want to drastically reduce its learning rate. To fix this resilience RMSProp uses the following moving average:

$$\begin{aligned}
 E[(\nabla_{(w,b)}L(w_t, b_t))^2] &= 0.9 E[(\nabla_{(w,b)}L(w_{t-1}, b_{t-1}))^2] \\
 &\quad + 0.1 \nabla_{(w,b)}L(w_t, b_t) \odot \nabla_{(w,b)}L(w_t, b_t), \tag{2.35} \\
 (w_{t+1}, b_{t+1}) &\longrightarrow (w_t, b_t) - \eta (E[(\nabla_{(w,b)}L(w_t, b_t))^2] + \epsilon)^{-1/2} \odot \nabla_{(w,b)}L(w_t, b_t),
 \end{aligned}$$

where again all the operations are element-wise. Note that the directional adaptation of the gradient at the current step t is introduced with a factor of 0.1. This gives robustness to the method as it requires persistence in the change of a sign through several steps to change the behaviour of the method. RMSProp also works better than RProp with full-batch, due to this robustness.

Momentum + Component Learning Rate Adaptation

This last type of methods combine the ideas of momentum and component learning rate adaptation. Recall that momentum introduced information about the slope by preserving some of the gradient of the last step, and component adaptation rescaled the components of the gradient in each direction dividing by the square root of the square of the gradient, which is some kind of approximation of the curvature in the principal directions corresponding to the elements in the diagonal of the Hessian, and tell us about the variation of the slope and the directions that we can go faster. Combining slope and curvature to get some sort of first order Newton method has already been done AdaDelta, however, as the information of the slope came from previous increments (already corrected gradients) and not strictly from previous gradients (the definition of momentum), we refrained from including it in this section.

Mainly the first method that embraced this approach is **Adam**, [38] (2014), not taking into account AdaDelta, (2012). Adam at the present time (2020) is one of the best result yielding first order methods, and it has become the almost de facto optimizer in deep-learning applications. It combines the versatility of pure classical momentum (to scape local minima) and component learning rate adaptation (to escape saddle-points), and it is quite fast.

The method design is as follows:

$$\begin{aligned}
m_t &\longrightarrow \frac{1}{1 - (\beta_1)^t} \left(\beta_1 m_{t-1} + (1 - \beta_1) \nabla_{(w,b)} L(w_t, b_t) \right), \\
v_t &\longrightarrow \frac{1}{1 - (\beta_2)^t} \left(\beta_2 v_{t-1} + (1 - \beta_2) \nabla_{(w,b)} L(w_t, b_t) \odot \nabla_{(w,b)} L(w_t, b_t) \right), \\
(w_t, b_t) &\longrightarrow (w_{t-1}, b_{t-1}) - \eta \frac{m_t}{\sqrt{v_t + \varepsilon}},
\end{aligned} \tag{2.36}$$

where all operations (sums, products, roots and inverses) are element wise, m_t is called the first order momentum in t and β_1 is its decay rate, and v_t is called the second order momentum in t and β_2 is its decay rate.

Observe that inside the big parenthesis of m_t , we have a decaying average of the gradient, which resembles the classical momentum as defined in (2.28), and the big parenthesis of v_t is exactly the same as the decaying average principal direction curvature approximation in the AdaDelta (2.32). Each of the momentums is given an exponential rescale factor in the form of the term $1/(1 - (\beta)^t)$, which tends to 1 as t increases. Hence, since $0 < \beta_2 < \beta_1 < 1$, in the beginning m_t dominates over v_t giving some sort of annealing effect by prioritizing momentum over curvature in the first steps. Finally the steps are updated as in AdaDelta following a Newton-like approach.

Other notable variations of Adam are: **AMSGrad** [39], which tries to fix the convergence problem of Adam in some instances (however, it is argued that the very specific instances that AMSGrad fixed do not really occur real problems, thus it is sometimes regarded as more complex and noisier version of Adam); **Nadam** [40] which uses Nesterov’s momentum instead of classical momentum (from which the N in its name comes from); and **AdamW**, which tries to incorporate a Tikhonov regularization (which we will see in the Regularization section) inside the optimizer, instead of adding it to the loss function. As an additional comment, very recently a new type of more sophisticated first order methods which do not even require specifying a learning rate have appeared yielding apparently better results than Adam and its variations, one such method is **YellowFin** [41].

2.4.2 Second Order Methods

The previous first order methods yield good results in relatively small artificial neural networks (a few layers deep). They are not very computationally intensive and have linear convergence (which is often times enough), all at the expense of fixing a hyper-parameter, namely the learning rate. In particular, AdaDelta and Adam have proven to work really well against vanishing/exploding gradient and **sparse gradient** problems. Sparse gradients are a “kind” of vanishing gradients which happens when the dataset is sparse, i.e. there are rare features that occur in very few data points. Hence, if we recall that given the loss function form, the gradient is actually a sum of sub-gradients each associated to an individual data point, $1/N \sum_{i=1}^N \nabla_{(w,b)}(\dots)$, then the contributions to the gradient to fit these rare features are small in comparison to more common features, as there are fewer points and sub-gradients that can add to the sum. In that case it is said that there is a weak signal for that feature, and in practice this means that the parameters associated with that feature have smaller derivatives, creating saddle regions as the vanishing gradient problem does.

Nevertheless, as well as these first order methods work in many small problems with a somewhat homogeneous dataset, there are two related motives occurring in more complex problems that may require the consideration of higher order methods:

- The first, motive is computational cost: As we consider larger artificial neural networks, the number of parameters scales up, and the smaller number of steps required with the quadratic convergence (or almost quadratic) of second order methods start to become a computational advantage to the simpler but larger amount of steps required with linear convergence of first order methods.
- The second motive, very related to the first, is the high slope variation: As the number of parameters increase or the input dataset becomes more noisy, the loss function hyper-surface starts becoming more “bumpy”, meaning that in using first order methods the strides in the direction of the gradient have to be shorter to account for its variation, i.e. the learning rate has to be reduced. Recall that AdaDelta and Adam corrected the learning rate based on some sort of approximation of diagonal of the Hessian. Therefore, when the Hessian increases (which happen when the number of parameters increase), the effect of elements outside of diagonal aggregate to become more relevant, and the methods lose part of their effectiveness.

Out of all the second order methods, the classic **optimization Newton method** is the principal one. This method relies on the Taylor expansion up to second order to approximate the function to be optimized by a quadratic function, in a local neighbourhood or region of confidence of a point (w_0, b_0) . In our case, the loss function can be approximated by:

$$L((w_0, b_0) + p) \approx L(w_0, b_0) + \nabla_{(w,b)} L(w_0, b_0)^\top p + \frac{1}{2} p^\top H(w_0, b_0) p, \quad (2.37)$$

where p is an increment within the region of confidence and $H(w_0, b_0)$ is the Hessian matrix in (w_0, b_0) . Then, as (2.37) is a quadratic function of p , it should have a unique minimum p_0 , thus deriving the expression (2.37) with respect to $p \in \mathbb{R}^{n \times m}$, the minimum p_0 must satisfy:

$$\begin{aligned} \frac{\partial}{\partial p} \left(L((w_0, b_0) + p) \right) &\approx \frac{\partial}{\partial p} \left(L(w_0, b_0) + \nabla_{(w,b)} L(w_0, b_0)^\top p + \frac{1}{2} p^\top H(w_0, b_0) p \right), \\ \frac{\partial}{\partial p} \left(L((w_0, b_0) + p) \right) &\approx \nabla_{(w,b)} L(w_0, b_0) + H(w_0, b_0) p, \\ 0 = \frac{\partial}{\partial p} \left(L((w_0, b_0) + p_0) \right) &\approx \nabla_{(w,b)} L(w_0, b_0) + H(w_0, b_0) p_0, \\ p_0 &\approx -H^{-1}(w_0, b_0) \nabla_{(w,b)} L(w_0, b_0). \end{aligned} \quad (2.38)$$

Therefore the optimization Newton method, being at (w_t, b_t) in step t , computes a new step point by approximating the original function in a confidence region around (w_t, b_t) by a quadratic function using Taylor’s theorem, then finds the increment p_t that minimizes that quadratic approximation of the original function, and moves using that increment. In summary, the optimization Newton method update rule is:

$$\begin{aligned} p_k &\rightarrow -H^{-1}(w_t, b_t) \nabla_{(w,b)} L(w_0, b_0), \\ (w_{t+1}, b_{t+1}) &\rightarrow (w_k, b_t) + p_t. \end{aligned} \quad (2.39)$$

The increment p_t is known as *search direction* these types of methods.

Note that no hyper-parameters are required, and second order convergence is guaranteed by Taylor's theorem. As a major drawback, the method involves computing the Hessian matrix and inverting it, which is an extremely impractical and computationally expensive task, even if the number of parameters is just moderately large. Thus, there are a series of methods that modify this optimization Newton method to use approximations instead of the whole inverse of the Hessian, but preserve many of the good properties of the original. As a trade-off for their decrease in computational complexity, these methods lose their quadratic convergence, but they still get a much better than linear convergence, usually referred to as super-linear convergence, which outclass any first order method's convergence.

Quasi-Newton Method

In the Quasi-Newton family each step update uses the same idea as in the Newton method, with the a small variation. Instead of computing and using the Hessian matrix $H(w_t, b_t)$, we use an approximation matrix B_t which we have to update in every step. This means that, in essence, all the reasoning and derivation for the update rule are completely analogous to that of (2.37-2.38) with the only difference being writing B_t instead of $H(w_t, b_t)$. The final update rule will in fact have the same blueprint as the Newton's,

$$\begin{aligned} p_k &\rightarrow -B_t^{-1} \nabla_{(w,b)} L(w_0, b_0), \\ (w_{t+1}, b_{t+1}) &\rightarrow (w_k, b_t) + p_t, \end{aligned} \tag{2.40}$$

with a few additions (two in particular). The first is that, in every step $H(w_t, b_t)$ is being replaced by B_t , which is a matrix that changes with the curvature, but does not necessarily have to be an exact approximation of the Hessian matrix (for instance, it could be a scaled down version or be displaced). This means that we can trust the search direction p_t for its direction but not for its magnitude, thus, we will require a learning rate α_t to scale p_t at every step. There are two ways to compute α_t at each step, namely *inexact lines search* and *trust regions*. Coincidentally, the Quasi-Newton methods that we will be seeing next use inexact lines search, and the truncated Newton methods of the next subsection use trust region. In particular, the inexact lines search that we will use is the satisfaction of *Wolfe conditions* which is given by the following set of inequalities:

$$\begin{aligned} L((w_t, b_t) + \alpha_k p_t) &\leq L(w_t, b_t) + c_1 \alpha_t \nabla L(w_t, b_t)^\top p_t, \\ \nabla L((w_t, b_t) + \alpha_k p_t)^\top p_k &\geq c_2 \nabla L((w_t, b_t))^\top p_t, \end{aligned} \tag{2.41}$$

with $0 < c_1 < c_2 < 1$. Using Wolfe conditions, α_k is progressively decreased until the inequalities (2.41) are satisfied. This guarantees the learning rate holds sufficient decrease in curvature conditions.

The second issue is how to compute the matrices B_t at every step. In principle, two general requirements are demanded, that help calculate the matrix: it has to be symmetric like the Hessian and it must satisfy the *secant equation* (or *Quasi-Newton equation*). This secant equation can be obtained by differentiating in terms of the increment variable p for the quadratic approximation (2.37 with B_t) for a given step t :

$$\begin{aligned} \frac{\partial}{\partial p} \left(L((w_t, b_t) + p) \right) &\approx \frac{\partial}{\partial p} \left(L(w_t, b_t) + \nabla_{(w,b)} L(w_t, b_t)^\top p + \frac{1}{2} p^\top B_t p \right), \\ \frac{\partial L((w_t, b_t) + p)}{\partial((w_t, b_t) + p)} \cdot \frac{\partial((w_t, b_t) + p)}{\partial p} &\approx \nabla_{(w,b)} L(w_t, b_t) + B_t p, \end{aligned} \tag{2.42}$$

$$\nabla_{(w,b)}L((w_t, b_t) + p) \cdot (0 + Id) \approx \nabla_{(w,b)}L(w_t, b_t) + B_t p,$$

then, if we take the value of p for the actual increment $s_k = \alpha_k p_k$, we expect $B_t \approx B_{t+1}$, which yields:

$$\begin{aligned} \nabla_{(w,b)}L((w_t, b_t) + s_k) &\approx \nabla_{(w,b)}L(w_t, b_t) + B_{t+1} s_k, \\ \nabla_{(w,b)}L(w_{t+1}, b_{t+1}) - \nabla_{(w,b)}L(w_t, b_t) &\approx B_{t+1} s_k, \\ y_k &\approx B_{t+1} s_k, \end{aligned} \tag{2.43}$$

where $y_k = \nabla_{(w,b)}L(w_{t+1}, b_{t+1}) - \nabla_{(w,b)}L(w_t, b_t)$. The last expression of (2.43) is what is known as the secant equation.

Therefore the general update rule of a Quasi-Newton method would be:

1. Obtain the search direction solving $p_k = -B_t^{-1} \nabla_{(w,b)}L(w_0, b_0)$.
2. Use Wolfe conditions (or any other approach) to find α_k .
3. Compute the next gradient $\nabla_{(w,b)}L((w_t, b_t) + \alpha_k p_k)$.
4. Compute the next matrix B_{t+1} imposing symmetry and satisfying the secant equation.

Every Quasi-Newton method varies in their approaches to compute α_t , and speciality B_t , as are many ways to construct symmetric matrices that satisfy the secant equation. Ideally, we want to use a expression of B_t which is easy to compute and easy to invert. The first of these methods was the **Davidon–Fletcher–Powell** method (DFP), but in current times the **Broyden–Fletcher–Goldfarb–Shanno** method (BFGS) and its Limited Memory version (which require less storage memory) are the most widely used methods in this family.

The BFGS impose the extra condition that B_t has to be positive defined, on top of the symmetry and satisfaction of the secant equation, and use Wolfe conditions. The following expression is the one that satisfies the three conditions:

$$B_{t+1} = B_t - \frac{B_t s_t s_t^\top B_t}{s_t^\top B_t s_t} + \frac{y_t y_t^\top}{y_t^\top s_t}, \tag{2.44}$$

which allows for the Sherman–Morrison formula to invert the matrix recurrently:

$$B_{t+1}^{-1} = \left(Id - \frac{s_t y_t^\top}{y_t^\top s_t} \right) B_t^{-1} \left(Id - \frac{y_t s_t^\top}{y_t^\top s_t} \right) + \frac{s_t s_t^\top}{y_t^\top s_t}. \tag{2.45}$$

Section based on [42, 43]. [44] has a really recommendable coverage on, line search conditions, trust regions and Quasi-Newton methods, far beyond this work. Also, [45] proposes a mini-batch adaptation of the L-BFGS method.

Truncated Newton Method

While Quasi-Newton methods such as BFGS have been adapted and used for large deep-learning applications successfully, not many other types second order methods have historically been used. In 2010 though, Martens [46] proposed an adaptation of a type of methods called the truncated Newton methods or Hessian-Free methods, which has been fairly replicated. The set-up for this method also revolves around trying to fixing the main problem of the optimization Newton method. Hence it tries to find the search direction p_t by solving the system in (2.38) without having to invert the Hessian matrix.

Rearranging the last expression of (2.38), the system to find the search direction in step t , can be written as:

$$H(w_t, b_t) p_t \approx \nabla_{(w,b)} L(w_t, b_t). \quad (2.46)$$

Here, the point of this method is that we do not have to compute the Hessian or any approximation matrix at all. Given any vector v , its product by the Hessian matrix can be thought of a directional derivative of the gradient as a function. Thus, we can compute the product $H(w_t, b_t) v$ by evaluating an extra derivative:

$$H(w_t, b_t) v = \lim_{\varepsilon \rightarrow 0} \frac{\nabla_{(w,b)} L((w_t, b_t) + \varepsilon v) - \nabla_{(w,b)} L(w_t, b_t)}{\varepsilon}, \quad (2.47)$$

where the limit can actually be replaced by some finite difference taking a small ε or by any other method. Then, we can use the Conjugate Gradient (CG) algorithm to solve the linear system (2.46) combined with (2.47) to obtain the search direction p_t . In summary the steps for step t would be:

1. Randomly initialize a search direction $p_{t,0}$
2. Apply the CG algorithm to solve the linear system of (2.46). Since with regards to the Hessian (which is the matrix of the linear system), the CG only requires the products $p_{t,i}^\top H(w_t, b_t) p_{t,i}$, use the finite differences of the gradient from (2.47) to compute this product without needing the Hessian.
3. Once the CG has converged (or is close enough) to the solution p_t , update the point $(w_{t+1}, b_{t+1}) = (w_t, b_t) + p_t$ and repeat the process.

This last method is the only one in this section that we have not used, as the implementation required much more technical work than the rest, since there is no package from which to adapt it. However, we still wanted it to be explained as a promising second order alternative to the LBFGS and BFGS.

2.5 Activation Functions and Parameter Initialization

In this section we will be introducing the most commonly used activation functions in deep-learning. We have avoided this issue up to this point because the choice of activation functions is crucial in mitigating the gradient problems and speeding up the optimization process, which we have now covered. Also, in practice, all the activation functions of an artificial neural network are chosen to be the same (except for the output layer), so from now on we will assume this is the case.

Although any function can do theoretically as long as it is continuous and one time piecewise differentiable, there are some other that can prove beneficial. For example, using bounded or dissipative functions can add a component of localized training, making the output of the activation function in certain neurons mostly dominant with respect to others for some range of input values, thus the parameters of those neurons become descriptive of the characteristics of that range of input values. As well as this, functions with derivatives close in magnitude to their primitives, or derivatives that have a close to linear behaviour in some region are desirable, as this helps mitigate the gradient problems and improve convergence through the initialization of the parameters (we will look this last part in the next subsection).

The following table shows the main classic activation functions used in deep-learning (holding many of the properties stated above), in order, from most used (top) to the least used (bottom).

Name	Value	First Derivative
Sigmoid:	$a(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$	$a'(x) = \sigma(x)(1 - \sigma(x))$
ReLU	$a(x) = \max(x, 0)$	$a'(x) = \Theta(x)$
Hyperbolic Tangent:	$a(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$a'(x) = \frac{1}{\cosh^2(x)} = \frac{4}{(e^x + e^{-x})^2}$
Exponential:	$a(x) = e^x$	$a'(x) = e^x$

Table 2.1: List of main activation functions.

About the nomenclature in table above: the sigmoid function is almost always represented as $\sigma(x)$, and $\Theta(x)$ stands for the Heaviside step function, which is 0 if the input is negative and 1 if the input is positive. Also due to the specifics of this work, where the loss function contains derivatives of artificial neural network, we would like to impose the following reasonable condition. If n is the degree of the differential system linear operator (this is, the highest order derivative of the operator), then at least the first $n + 1$ derivatives of the activation function should be different from zero. This condition makes sense as we would like to have some contribution (or signal) to model the higher order derivatives and not to have them vanish. Therefore, from the previous list the ReLU (Rectified Lineal Unit) function is mostly unsuitable for this work, since its second order derivative is strictly zero. We can see this effect clearly from (2.16), where we expressed the Hessian tensor (second order derivatives tensor) of our example network. If we are to consider the tensor arrangement of activation functions with ReLUs, then $D^{n_\ell, n_\ell}(a^{[\ell]}_{n_\ell}(z^{[\ell]}_{n_\ell})) = 0^{n_\ell, n_\ell}$, is the tensor of zeros, and since these terms appear multiplying in every factor of (2.16), all the second order derivatives of the artificial network with respect to the inputs are zero.

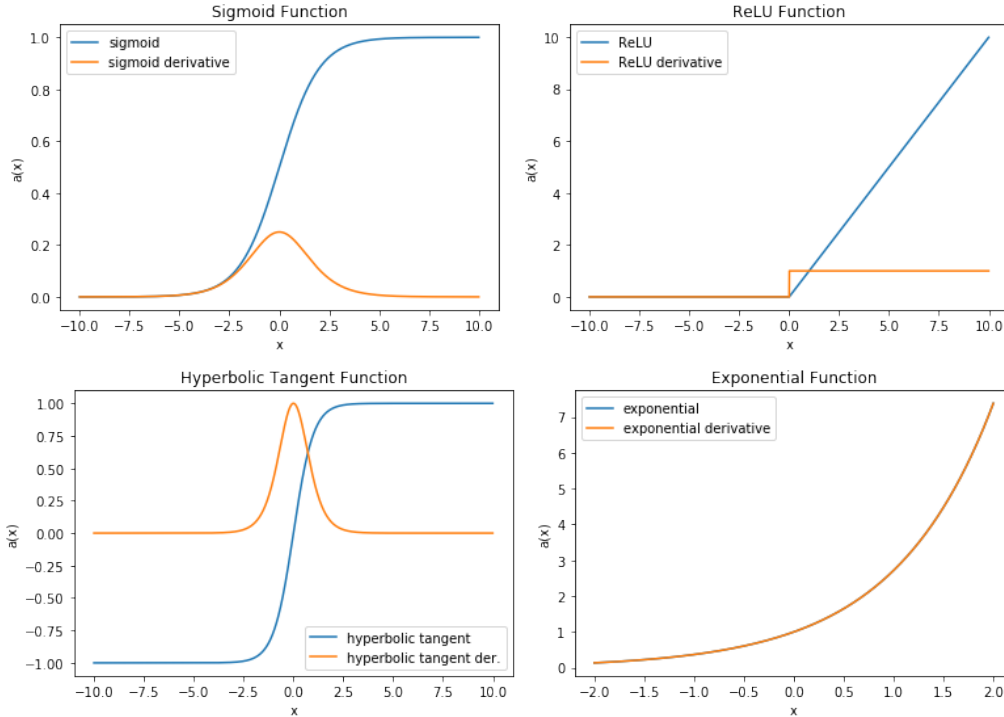


Figure 2.7: Main activation functions and their first order derivatives.

Observe that all the classic activation functions are monotonically increasing and bounded: the exponential and ReLU in the positive values, the sigmoid between 0 and 1 and the hyperbolic tangent between -1 and 1. In general, using activation functions with these two characteristics allows for some sort of localized pulse-like behaviour, which can be helpful in modelling traits over very specific regions. For instance, lets assume we have two neurons that use sigmoids, both share the same one input x , and having weights and biases $w_1 = -3$, $w_2 = 2$, $b_1 = -2$, $b_2 = 10$. Figure 2.8, shows the result of imputing x to the neurons and adding them up in the next layer. Note that the end result is some kind of continuous version of a square pulse function, which in this case differentiates between the region $(-2.5, 7.5)$ and the rest of the real line.

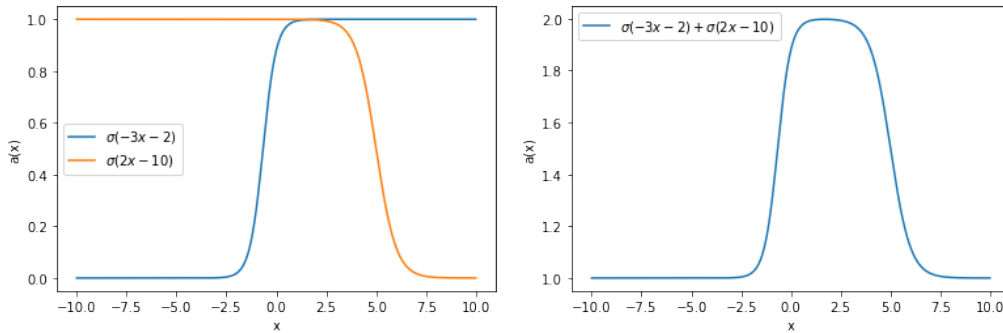


Figure 2.8: Combination of sigmoid functions.

Through time many secondary variations of these functions have appeared trying to fix some of their inconveniences. Two popular ones that we will consider are the Swish function [47], which also appear under the name of Sigmoid-Weighted Linear Unit [48], and the Softplus function [49]. The first offers a smooth continuous variation of the ReLU function using the sigmoid in its formula (thus has no vanishing higher order derivatives), and the second fixes the acute growth of the exponential function.

Name	Value	First Derivative
Swish:	$a(x) = x \cdot \sigma(x)$	$a'(x) = \sigma(x) + x \sigma(x) (1 - \sigma(x))$
Softplus	$a(x) = \log(e^x + 1)$	$a'(x) = \frac{e^x}{e^x + 1}$

Table 2.2: List of secondary activation functions.

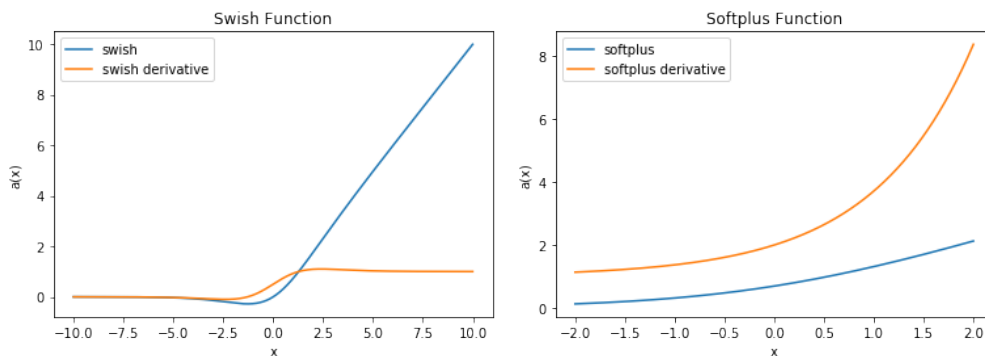


Figure 2.9: Secondary activation functions and their first order derivatives.

2.5.1 Parameter Initialization

An issue which depends directly on the activation function choice is the parameter initialization, necessary as the starting point to perform the iterative optimization methods we saw in the previous section. Actually, the optimization process is very sensitive to the initial point, and slight variations can even lead to very different results.

Each activation function has its own set of different initialization arrangements (many named after the authors that proposed them). However, all of them share the same the basic principle, which is to set the parameters so that, at least for the first inputs, the activation functions behave in their linear regime. For example, if we were to choose exponentials as activation functions, $a(x) = e^x \approx 1 + x + x^2/2 + \dots$, we would like the parameters to be close to 0, so the inputs to each function are also close to 0 (or small) for each activation function behave like $1 + x$ in the first steps. The reason behind using this idea (with a few refinements that we will be covering in the next paragraph) is because, in practice, solving exactly the non-convex optimization problem is virtually impossible, and many times even finding a decent enough result becomes a matter of luck with trial and error (using different methods and initialization points). By starting in a point at a region where activation functions behave linearly, we guarantee that the slopes (the loss derivatives with respect to the parameters) will not vanish or explode abruptly from the start and convergence will be smooth (at least in the beginning).

There are two extra considerations that we have to add to the previous idea of starting in the linear regime.

- First we have to break the symmetry of the parameters. This is means that we cannot initialize all the parameters in a layer with the same values. Recall the derivatives computed in section 2.3.1. Giving the same values to the parameters in a layer implies that the derivatives with respect to the parameters for that layer are all the same, thus at every step of the optimization the parameters also increment by the same amount, and stay equal along the whole process. To avoid this symmetry which effectively impede the ability of the parameters of a layer to adapt to different traits of the dataset, we add some variance by drawing the parameters from a distribution, usually a uniform or a normal. For instance, in the exponential case, instead of initializing the parameters with 0 (which is actually the only value for which the parameters cannot be trained), we would draw them from something like $\mathcal{N}(0, 0.1)$ or $\mathcal{U}(-0.1, 0.1)$.
- Second, we have to balance the variance of the weights in every neuron, so that the neurons outputs all stay within a certain magnitude. The two conditions to impose are for the linear combination of the neurons during evaluation to have similar variance (2.48), and for each loss derivative with respect to the parameters of a neuron in during back-propagation to also be similar in variance (2.49).

$$\forall \ell, \ell', n, n', \quad \text{Var}(z_n^{[\ell]}) = \text{Var}(z_n^{[\ell']}), \quad (2.48)$$

$$\forall \ell, \ell', n, n', \quad \text{Var}\left(\frac{\partial L}{\partial w_n^{[\ell]}}\right) = \text{Var}\left(\frac{\partial L}{\partial w_n^{[\ell']}}\right). \quad (2.49)$$

Assuming independence with respect to inputs and zero mean (2.48-2.49) become:

$$\forall \ell, \ell', n, n', \quad n_\ell \text{Var}(w_n^{[\ell]}) = n_{\ell'} \text{Var}(w_n^{[\ell']}), \quad (2.50)$$

$$\forall \ell, \ell', n, n', \quad n_{\ell+1} \text{Var}(w_n^{[\ell]}) = n_{\ell'+1} \text{Var}(w_n^{[\ell']}). \quad (2.51)$$

Averaging the previous conditions yield:

$$\forall \ell, n, \quad \text{Var}(w^{[\ell]}_n) = \frac{2}{n_\ell + n_{\ell+1}}, \quad (2.52)$$

where n_ℓ is the number of inputs of the neuron (called the fan-in), and $n_{\ell+1}$ is the number of outputs to the neuron (called the fan-out), which in fully-connected artificial neural networks is equivalent to the number of neurons in the current and the next layer respectively. Note that in practice we are simply scaling by the size of the network.

Getting these three ideas together (linear regime, variation and fan-in/out scaling) lead to the Xavier/Glorot [50] and He [51] initializations, which are the main ones we will be using throughout this work. The former uses straight forward (2.52), and the latter makes a correction for linear units (which are activation functions which drop or almost drop their negative values, for example the ReLU) by dropping the fan-out term. These initializations in their uniform and normal distribution versions therefore are:

$$\text{Xavier Uniform:} \quad \mathcal{U} \left(-\sqrt{\frac{6}{n_\ell + n_{\ell+1}}}, \sqrt{\frac{6}{n_\ell + n_{\ell+1}}} \right), \quad (2.53)$$

$$\text{Xavier Normal:} \quad \mathcal{N} \left(0, \sqrt{\frac{2}{n_\ell + n_{\ell+1}}} \right), \quad (2.54)$$

$$\text{He Uniform:} \quad \mathcal{U} \left(-\sqrt{\frac{6}{n_\ell}}, \sqrt{\frac{6}{n_\ell}} \right), \quad (2.55)$$

$$\text{He Normal:} \quad \mathcal{N} \left(0, \sqrt{\frac{2}{n_\ell}} \right), \quad (2.56)$$

For the sigmoid and hyperbolic tangent activation functions we will use the Xavier/Glorot initialization, and for the ReLU, exponential, swish and softplus activation functions we will be using the He initialization.

2.6 Regularization

At this point we have all the tools necessary to build an artificial neural network and to apply an optimization method to train it. However, there is still the issue of selecting and tuning the hyper-parameters of our inverse problem (recall that these are problems in which we have sample data and we want to fit the parameters of certain model to fit the data), which is hard. The issue here is that, given any model with certain hyper-parameters defining it, in our case the number of layers and neurons, there are two opposite situations that may arise. Either the number of hyper-parameters are too few and the model is too simple to fit the data (underfitting), or the number of hyper-parameters are too many and the model has extra degrees of freedom that do not necessarily fit the data (overfitting).

Underfitting is a problem which is relatively simple to spot, it is impossible to accurately fit the model to the data for as much as we train the model, and it has its only solution in increasing the number of hyper-parameters or the complexity of the model. For overfitting, detecting and solving to the problem is not as easy. To mitigate overfitting, as well as sometimes fixing the ill-posedness of a problem, regularization techniques are commonly used.

An overfitted model can be harder to train, this is because in the optimization process, the gradient that updates the parameters contains contributions of the extra degrees of freedom of the model. This can make the optimization process much noisier and susceptible to over-represent outliers. Also, overfitted models tend to generalize data poorly, i.e. their ability to predict on new data decreases, as any data not strictly in the training set would not simply be the extrapolation of the relations among the training points, but would also contain the contributions of the degrees of freedom. Regularization techniques help reduce the number degrees of freedom of the models.

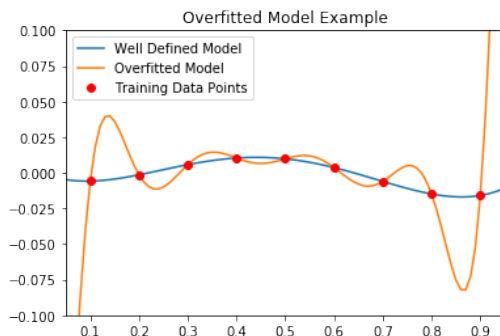


Figure 2.10: Example of overfitting of a model.

The previous Figure 2.10 show an example of an overfitted model. Both the blue line and orange line are polynomial models that fit the training data (red points). The orange model (which has more coefficients than the blue model) is overfitted though, since, most prominently at the extreme points, it generates a bumpy behaviour completely unrelated to that of the training data, related to the unnecessary extra parameters that we have added.

In practice finding the exact right number of hyper-parameters required in a model is an almost impossible task by the sheer amount of possibilities. Besides a trial and error strategy would impractical to compare models due to training being a computationally costly process sensitive to the optimizer and initial conditions. Hence effectively, since there is no general rule that can be followed, the hyper-parameters are often chosen within a reasonable rough margin (mostly based on the results of the first few steps of the optimization), guaranteeing some overfitting. Then regularization techniques are used to clamp down on the extra degrees of freedom of the model. This is much more viable approach than training an almost exponentially increasing amount of models with different numbers of hyper-parameters to narrow down the right number which does not underfit or overfit the data.

Each of the following subsections will be dedicated to a different regularization technique. We will improvise two categories to group these techniques based on the main general principles behind them, namely noise-based regularizations and restriction-based regularizations. In this work we will only be using restriction-based regularizations though.

2.6.1 Noise-based Regularizations

Behind the noise-based regularizations lies the idea that adding an stochastic component throughout the training process to add some (small) variance into model. To much variance can lead to a chaotic model (undesirable), but adding a small variance during training can be very beneficial as it can somewhat be seen as employing the extra degrees of freedom to account for the extra variability of the model.

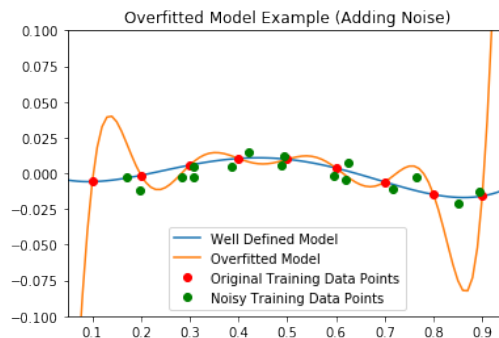


Figure 2.11: Example of a model adding noisy input.

Intuitive the concept in its most trivial form can be seen in Figure 2.11, which is nothing more than Figure 2.10 to which we have added some extra noisy points (the original points plus some extra noise). It now becomes apparent that when fitting the model the end result would be closer to the well-defined model (blue line) than to the overfitted model (orange line), since now the model has to also account for the green dots for which the blue line has a smaller error, specially at the extremes. Thus, if the variance is small the extra degrees of freedom are spent to ensure that small variations in the data do not yield overwhelmingly large changes in the model.

Now that we have explained how noise works, we will be looking at how to introduce it into the model for the training phase. The most obvious way is to introduce variance into the model is by using the **direct approach**, this is add noise to the input data, like for example, $x_i + \mathcal{N}(\mu, \sigma)$. A much smarter way to apply this concept is making use of some invariant to generate new data, in what is called **data augmentation**. This happens a lot in object recognition, whereby a car in a picture is car independently of the image being rotated 90° or the car appearing in the center or a corner of the picture, thus we can rotate or shift the pictures to generate new inputs.

There are much more sophisticated way to introduce variance into the model. One of these is using **noisy neurons**, which implies that, only during the training phase, we add noise to the output of each neuron, this is $y^{[l]}_{n_\ell} + \mathcal{N}(\mu, \sigma)$. Another one is the **dropout** technique [52], which only during the training phase uses a probability to suppress the output of a given neuron. Hence, for example, for every neuron at each step we would draw a number from a uniform distribution $p \sim \mathcal{U}(0, 1)$, and if $0.9 < p < 1$ we would set its output to 0 (this would be a dropout of 10%). At last we recall that the annealing effect explained in section 2.4.1 for the vanilla stochastic gradient descent can also be considered some sort of noise-based regularization technique.

One of the main issues with noise is that, we want to introduce some small variance throughout the model during training, but we want this variance to be controlled and small along the process, for the end result not to be a chaotic model fitting only noise. By this we mean that, we do not want the effect of the variance introduced in the early layers (the ones closest to the input) to explode in the following layers. We want the variance contribution to remain small as it gets imputed into its next layers. The principal idea to assure this, as well as providing other very good properties, are the **normalization** techniques which has become a staple in many deep-learning, namely batch and layer [53] normalization. It consists in normalizing either input batch or the outputs of the neurons of every layer, respectively.

In general for small models (specially in their number of layers) with well initialized parameters do not require normalization because in those dimensions the noise will most likely not scale up. One inconvenience of normalizations is that it correlates the gradients. Recall from (2.4-2.5), that the loss function is a summation over each of the individual losses at the points of the collocation, which makes the derivatives of the parameters with respect to the loss a sum of uncorrelated derivatives. Normalization entangles these derivatives which makes the update gradient for the optimizer correlated with respect to collocation points, thus, the gradient computation at every step of the training becomes less sparse and more computationally expensive. Additionally, normalization does not work well with dropout.

The reason we have disregarded noise-based in favour of restriction-based regularizations, although both are mutually compatible, is because understanding the behaviour of noise in a context where we are considering the derivatives of the artificial neural network is very risky and becomes exponentially more complex with the order of the derivatives. Also, in the case of normalization (which we have tested for this work), where we do not actually introduce variance but limit its effects, the extra increase in computational cost, caused by the correlation of the gradient, builds up on the already unavoidable computation of higher order order derivatives of the neural network required in this work, making the training process many times slower and impractical. On the contrary, the techniques that we have categorized as restriction-based are mostly (soft or hard) binds on the parameters. Thus, their application do not interfere with the computation of L_1 and L_2 , and so, their effect is applies afterwards.

2.6.2 Restriction-based Regularizations

On the other side of the spectrum lie what we have named as restriction-based regularization techniques. These are a set of soft or hard constrains on the parameters, added in the loss function or are applied independently. Hence, any extra degrees of freedom in the model may be invested in fulfilling these constrains.

The most common of these all are the **weight penalties** (actually it would be parameter penalties). This regularization techniques rely on an extra term, which is added to the loss function, and impose some preference in the parameters (this would correspond to the placeholder term R introduced in (2.1)). The most notable of these regularizations is the popular Tikhonov regularization which implies adding the following term:

$$R = \lambda (\|w\|_2^2 + \|b\|_2^2) := \lambda \sum_{\ell, n_\ell, m_{\ell-1}} \left(\|w_{n_\ell}^{[\ell] m_{\ell-1}}\|_2^2 + \|b_{n_\ell}^{[\ell]}\|_2^2 \right), \quad (2.57)$$

where λ is the regularization scale factor. Technically, this is same Tikhonov regularization as the one used in the least squares method for linear regression. One possible interpretation for it is that we impose a “minimum energy” model, i.e. we are looking for a model where the bias and specially the slopes are the smallest possible (yielding a “flatter model”), which happens naturally as we minimize the term R in the loss function. Another common reading is that the additive contribution of $\nabla_{(w,b)} R$ to the total update gradient introduces some sort of dampening force which penalizes the optimization method when moving in directions where the loss function, $L(w, b)$, is less smooth (which happen with larger values of w and b). Going back to Figure 2.10, the well-defined model would satisfy best the term (2.57).

Furthermore, the reason why these are called weight penalties and not parameter penalties, is because in most cases this regularization only affects the slopes (weights), and $\|b\|_2^2$, the bias terms in (2.57), are dropped. For this work though, we are considering the bias term as we believe it provides some sort of centring effect on the neuron outputs which (in the particular problem instances chosen for this work) help in speeding up the training. However, in a general context, this could prove tricky and is generally undesirable, specially in instances where the solution to the initial/boundary problem we want to approximate have many different localized features, i.e. the solution is very bumpy (which will not be case in this work). The reason for this, is that biases, although not strictly necessary in an artificial neural network (an only weights network is absolutely functional), when applied help to optimize the differentiation among neurons in the same layer. For example, in a neuron using sigmoid activation functions, given two inputs from two different regions and combining them with the weights, suppose we obtain values 0.5 and 1.5. Then, computing the activation, giving the neuron $b = 0$, yields a difference in output between the two inputs of ~ 0.19 , and for $b = 2$, it yields a difference in output between the two inputs of ~ 0.05 . This is apparent from Figure 2.7 as we see that the maximal slope is centred around zero. Hence, using a bias to shift the product of inputs and weights in a neuron can lead to an increase or decrease of the difference in outputs among values in different regions, an effect that applying a Tikhonov term which pushes $b \rightarrow 0$ can even negate, as contrary to the multiplicative contribution of the weights in the neuron, biases have an additive one requiring much larger values to have a significant effect. A second interpretation can be drawn by explanation given in Figure 2.8, whereby we argued that two increasing functions could be combined to form a sort of dissipative square pulse function. By reducing the biases, we reduce the amplitude of the plateaus (width of the windows) of these arrangements, which reduces some the specificity that can be achieved for certain regions in the model. Finally, we can still argue that the loss of some localized specialization in the neurons due to the term $\|b\|_2^2$ should still not pose a problem in the adaptability of the model, as it would simply make the contributions of a the neurons in a layer for a region more overlapping, why should this be a concern and undesirable? Although this last statement is true, a generally desired feature for a good artificial neural network is for it to be a *sparse artificial neural network*, i.e. that for any input given to the network almost all of the signal is contributed by just few neurons (or in other words any input only require passing through few relevant neurons neurons, and not all of them have to be active at the same time). Because of the objective of this work, which is proving it is possible to approximate solutions of initial/boundary problems by artificial neural networks, we would rather have the extra regularization effects of the term to obtaining a sparser network.

Another custom weight penalty that we have devised and seems to work rather well in this work is the following:

$$\begin{aligned}
R &= \lambda \left(\left\| \frac{\partial \hat{u}(x)}{\partial w} - \frac{\partial}{\partial w} \frac{\partial \hat{u}(x)}{\partial x} \right\|_2^2 + \left\| \frac{\partial \hat{u}(x)}{\partial b} - \frac{\partial}{\partial b} \frac{\partial \hat{u}(x)}{\partial x} \right\|_2^2 \right) \\
&:= \lambda \sum_{\ell, n_\ell, m_{\ell-1}} \left(\left\| \frac{\partial \hat{u}(x)}{\partial w^{[\ell]} \frac{m_{\ell-1}}{n_\ell}} - \frac{\partial}{\partial w^{[\ell]} \frac{m_{\ell-1}}{n_\ell}} \frac{\partial \hat{u}(x)}{\partial x} \right\|_2^2 + \left\| \frac{\partial \hat{u}(x)}{\partial b^{[\ell]} \frac{n_\ell}{n_\ell}} - \frac{\partial}{\partial b^{[\ell]} \frac{n_\ell}{n_\ell}} \frac{\partial \hat{u}(x)}{\partial x} \right\|_2^2 \right), \quad (2.58)
\end{aligned}$$

with its idea being, instead of using the extra degrees of freedom to obtain the solution with minimal slopes, to find a solution whose derivatives with respect to the parameters and with respect to parameters and inputs are similar in magnitude. This turns out to work quite well as we will see in the next section, and in fact, it requires no extra computations as all the derivatives involved in (2.58) are automatically calculated when computing $\nabla_{(w,b)} L_1$ (the loss of the differential operator).

A topic that we have not covered yet is the setting of the regularization coefficient λ . This has to be done manually and it has to be revised and updated at every step of the training process. We want the model to minimize its errors $L_1 + L_2$ (main objective) over fitting the term R (secondary objective). Typically, to establish priority as with any other multi-objective function we use the coefficient λ , to limit the magnitude in which the term regularization R contributes the total loss without becoming irrelevant. One possible reasonable demand would be to ask for the regularization term magnitude to be between 10% and 20% of the main term magnitude, or in other words $0.1 \cdot (|L_1| + |L_2|) \lesssim |R| \lesssim 0.2 \cdot (|L_1| + |L_2|)$, and adjust λ every time the criterion is not met. A much logical approach at first sight would seem to be, adapt λ as a function of the magnitudes of the terms, for example in the Tikhonov case, we could always make $\lambda = 0.1(|L_1| + |L_2|) / (||w||_2^2 + ||b||_2^2)$, for R to always be 10% of the other two terms. However, bear in mind that every time we modify λ we are in fact changing the total loss function $L(w, b)$, which is detrimental for the robustness and convergence of the optimization process. Thus, it is best to use a threshold that allows the optimizer to train on a fixed hyper-surface for many steps until a correction has to be done, than to have the optimizer move on a hyper-surface that changes in every step, more so considering that all but vanilla SGD use some kind of memory from previous gradients, which becomes irrelevant if the hyper-surface has changed. Nevertheless, in section 3.3 in the next chapter, we will propose an alternative approach to deal with this issue as part of a larger framework to deal with multi-objective loss functions, which seems to perform much better than this classical threshold strategy and achieve faster convergence.

An alternative to weight penalties are the **weight constrains** (or parameter constrains), which are sets of inequalities that can be applied to the parameters, either by component value like $a < |w_{n_\ell}^{[\ell]}| < b$, or by node or layer norm $a < ||w_{n_\ell}^{[\ell]}||_2^2 < b$. The way to apply these inequalities is usually by *clipping*, this is for example, if we update a parameter in a given training step and surpass the upper bound b , then parameter is set to b . This effectively stalls the training of parameters that have become too large (usually very dominant effects) or prevents from vanishing parameters that have become too small (usually very negligible effects), forcing a more even distribution in the relevance of the parameters. A more common clipping practice is **gradient clipping**, which is applied on the gradients with respect to the parameters used to update the parameters on an upper bound. Hence, this limits the effect of any exploding derivatives case, as any extremely large derivative which would mean an extremely large update, whereby using gradient clipping would be instantly reduced to a maximum reasonable range. In the training of all the models in this work we have implemented a component upper bound weights and bias clipping 10^3 , which is reasonable enough for the small scale of the models, and an upper bound gradient clipping by norm of all the parameters in the layer of 1.

Finally, we propose and will implement the following idea for a regularization, which can be drawn from the context of this work. Since want to train artificial neural networks to satisfy differential equations, and thus approximate their unique exact solutions, any conservation law satisfied by the exact solutions must also be approximately satisfied by the artificial neural network. Hence we can add **conservation laws** to the loss function the same way we did with the other weight penalty regularization, by simply replacing the conservation laws into the placeholder R (and optionally adding some regularization constant). Actually we could argue that, since the network must also satisfy the conservation law as closely as possible as the exact solution does, the term is not actually a regularization but a legitimate extra term which speeds up training, and not an extra condition which help select some specific model among the many that approximate the solution (essentially a regularization).

Not many differential equation have known conservation laws though, and thus, are hard to come by. Besides, in cases where conservation laws are known, usually adding an external force to the equation (something we will be doing in this work) invalidates such laws. Sometimes, although it is rare, this can be accounted for by deriving again the conservation law with the external force, and this can lead to the original law with some extra terms, like an the integral over the domain of the external force if the domain is bounded.

2.6.3 Other Regularizations

In this subsection we will look at two very common practices that might as well fall in the category of regularization. The first is known as **pre-training**, which consists of, instead of initializing a new artificial neural network to approximate an initial/boundary problem, we would use an already existing one as a starting point, with the hopes that this network is already closer to the desired outcome. All the attempts in this work to use pre-training with artificial neural network trained to only fit the initial/border data, to only fit the domain, or to fit the differential equation dropping any of its terms, have either had the same performance as using no pre-training, or worsened. The most plausible explanation might be that being the loss function multi-objective, it is best to keep a balanced agreement between the two parts from the start, rather than starting by fitting either L_1 or L_2 , as the region in the parameter space we can fall in during these one term optimizations might be useless or even detrimental to the other term, thus making worse the combined optimization.

Second is **early stop**, which is not only always applicable, but useful in many ways. Data used to fit artificial neural networks (or any model by that), is usually split into two groups, the *training data* and the *validation data*. The training data is used to fit the model (it is the data imputed in the loss function during the optimization), and validation data is used as a control mechanism to prevent overfitting of the model. Therefore every certain number of iterations in the process (for example 1000 iteration), we evaluate the loss with respect to the validation set, and if this validation set loss has worsened with respect to its previous evaluation, then we stop the training (this is early stop). The principle here is that the training process is blind to the validation data (not used), however the model should still fit this data as part of its capabilities to generalize beyond the training points. Similarly to the workings of noise, if the validation loss gets worst, the behaviour beyond the training points becomes undesirable, and thus, the model is overfitting. When this happens we can simply stop the training completely, or it might be a sign that the learning rate of the optimizer is too large and we have to reduce it, we can try introduce some extra regularization or adapt the regularization coefficient λ to correct the model, or we can **generate a new batch** which is also a regularization technique, and resume the training. Hence, early stopping is very useful not only as a regularization technique, but it gives a cue to rectify the training of the model when the process is stalled. Throughout this work we use validation intervals (we check the validation loss) every 1000 steps of training.

As a final note, we want to address the reason why we do not check the validation loss at every step. The first motive is because first order optimizers are not always smooth, i.e. the loss function can be decreasing but in an oscillating (conjugate) manner (specially around valleys), thus in a very short span early stop could confuse one of these fluctuations where the is at local maximum with a stop criterion. Still, for second order methods which use line-search that guarantees that there is always a decrease in loss (or they simply stop), the reason is that it is computationally more expensive to evaluate an extra loss at each step, and a few more steps from the early stop criterion will not substantially change the model.

Chapter 3

Case Studies and Simulations

In this chapter we will finally be training artificial neural networks to approximate the solutions of some instances of initial/boundary, starting by the most simple case and building up to more complex operators.

The layout of this chapter will be fairly consistent. Besides the first three sections, dedicated to the general implementation related topics practical to this work: the coding framework, function approximation capabilities of artificial neural networks, and adaptation to multi-objective function training; each of the remaining sections follows the same structure of posing a problem instance, training, and result analysis, with different operators. All of the operators used in this chapter have already been detailed in Table 1.2, and as mentioned in the introduction, we will be using only Cauchy initial/boundary conditions. With regards to the external forces, we will be selecting them ad-hoc in every problem so that the solution is a simple known polynomial. This way we can benchmark the artificial neural network results against the exact solution with ease.

3.1 Coding Artificial Neural Networks

From an implementation standpoint, deep-learning model training require the computation of many operations, specially linear combinations (tensor operations). Recall that an artificial neural network neuron is composed of a linear combination of the outputs of the previous layer, and an application of a non-linear activation functions. In terms of activation functions, little can be done to improve performance, but the many sums and products of the linear combinations are susceptible to high parallelization, as they are mostly independent among neurons, low in computational cost and high in number. Therefore, in order to speed up learning, instead of using CPUs, which at the time of this work have up to 8/16 cores, i.e. processing units and maximum amount of operations that can be perform in parallel, we can make use of the already existing GPU hardware. GPUs are optimized for image processing, a process which rely heavily in matrix multiplication. Thus contrary to CPUs composed of a few powerful cores, GPU are built using a large number of lower end cores, which at the time of this work can be in average of 120 cores. By parallelizing the linear combination operations to the many cores of a GPU, we can reduce the training time of an artificial neural network manifold, especially in larger networks. Nowadays, a new piece of hardware specially designed for deep-learning training has irrupted called TPUs (Tensor Processing Units). This hardware contains an even larger number of cores and its architecture is ad-hoc designed to parallelize tensor operations, improving on the capabilities of GPUs.

In order to manage and distribute the flow of operations to make the most use of GPUs and TPUs, there are several developed software solutions. We will briefly give a basic understanding on the most prominent high/medium/low level options.

On the lowest level, almost exclusively, lies the APIs named CUDA, developed by GPU maker NVidia. This API allows for direct control and distribution of operations of the cores in a GPU/TPU. However, from a practical perspective, unless we want to really customize and micromanage the resources in our GPU/TPUs, this level of control is too much. Thus, there are several middle-level libraries used in deep-learning that automatically handle these tasks, the most popular ones being TensorFlow, developed by Google, and PyTorch, which is open source (both running on CUDA). The way this libraries work is by implementing their own class for multi-dimensional objects (like arrays), and every time tensor operations are performed on these objects, they use CUDA under the hood to distribute its operations into the GPUs and/or TPUs cores automatically. This simplifies the work by allows us to concentrate on programming the mathematical framework for the models without having to deal with the management of the parallelization tasks. Lastly, on top of these libraries, there are also higher level ones, such as Keras library which hinges on TensorFlow. These build on the multi-dimensional class to further implement classes for layers, models, optimizers, training, and more, with many options, creating a structure that allows to build and train a model in a very simple and encapsulated manner.

The code for this work has been written using Python's version of TensorFlow 2.3. Some of the principal classes of the Keras library that implement the layers, models and optimizers, have been imported but only serve as a structure, since they were not applicable to the special formulation of this work, they had to be completely overwritten. Moreover, the execution has been done through Jupyter Notebooks in the Google Colab cloud environment which offers a free Nvidia K80/T4 GPU. For the code, address to Appendix B.

3.2 Approximating a Function

Here we will be studying the approximating capabilities of an artificial neural network to model a function. This can be considered, in the context of this work, as simplest case of differential equation possible, the trivial case of the identity operator, whereby the artificial neural network should be adjusted to satisfy:

$$\mathcal{L}[\hat{u}(x)] = f(x) \Rightarrow \hat{u}(x) = f(x), \quad (3.1)$$

which is equivalent to simply having the artificial neural network model the external force function. Being this operator of order zero, initial/boundary conditions are irrelevant, and thus, the loss function to optimize is:

$$L(w, b) = L_1(w, b) + R = \frac{1}{N_\Omega} \sum_{i \in N_\Omega} \left(\hat{u}(x_i; w, b) - f(x_i) \right)^2 + R. \quad (3.2)$$

The (external force) function that we will be approximating in this section will be the polynomial $f(x) = x(x - 1)$. Using this instance as an example, we will compare how well different optimizers and activation functions work at the task of modelling functions, as well as explain some of the behaviours of training. Here, the basic metric to assess performance is the relation loss function - iterations, this represents how well the model fits the solution at every step. As the loss function can have very steep decreases in value, we will be using logarithmic scales for better representations. Moreover, for every model we will be showing a plot of the end result compared to the real solution, and in evaluation future evaluations, where it actually applies, we will also be decomposing the total loss into its components L_1 and L_2 .

At last, we will be using an artificial neural network with an input layer with 1 neuron, two hidden layer with 3 and 4 neurons each, and an output layer with 1 neuron; which we will call a [3,4,1]-ANN, to approximate (3.1). First, we will start by comparing, how different activation functions work for the same network layout with different choices of activation functions. For this purpose we will use an Adam optimizer fixed to $\eta = 0.01$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$, and we will see the performance for the first 10000 iterations without any adjustments. The only regularization applied will be a parameter upper bound of $10e^3$ and a gradient clipping by layer norm of 1, which as explained in the regularization section of this work, will be the standard. Initialization from here on are done as detailed in 2.5.1, using the normal distribution versions.

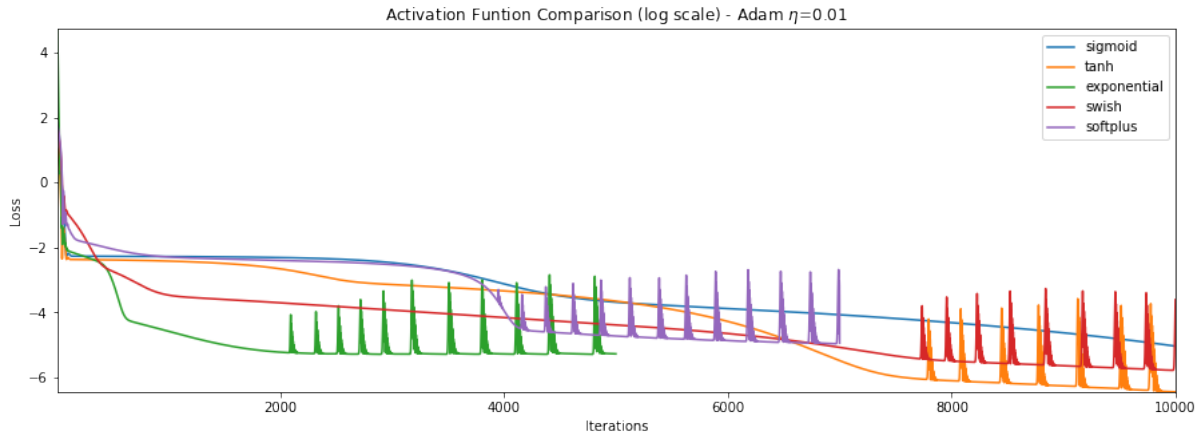


Figure 3.1: Comparison for different activation functions training performance for a [3,4,1]-ANN, with Adam $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$. Log10 scale.

From Figure 3.1 we can observe that by the end of the training, for many of this activation functions, the loss value stagnates and an oscillating behaviour starts to appear. This is in some sources called *saturation*, meaning that the model is unable to learn more. Fundamentally this is intrinsic to the model because we are approximating functions which may (and actually have) a very different analytic structure from the parametric model we are using. Hence, the same way it happens when we use a Taylor series expansion, where we have to truncate at some order to obtain a finite model obtaining an error, here we will also have an intrinsic minimal error to the model. However, being this a non-convex optimization problem, we do not know if these saturations correspond to reaching the intrinsic error of the model (global minimum), or if it corresponds to a local minimum or a valley. When this happens, if we have implemented early stop in the training loop, the process will stop (which happened for the exponential and softplus activations in Figure 3.1). Then, we can choose to strengthen the regularization (not very effective), or to reduce the learning rate in the optimizer or use a second order one in the hopes it is valley and we can scape it. If we use a second order method (in this work L-BFGS), and the method stops, we can be almost completely sure that at loss is in some minimum and we will not be able to scape it. This is because the stop criterion with line search is not finding any ratio in the gradient direction that can actually decrease the loss function (and line search looks for this ratio with exponential decay), thus almost guaranteeing we are in a minimum. Here is where luck of non-convex optimization comes into place, as a different initialization, or instance of the same initialization, or a simpler or more complex network layout, or an apparently worst performing optimizer can lead to a different optimization path through the loss hyper-surface, leading to a better fitting model.

In this benchmark we have used quite a minimal model to ensure it is not too overfitted (regularization can only fix some overfitting) and the loss function is quite smooth, and we have used a fairly robust optimizer. Therefore, we can assume with some confidence that the saturation corresponds at least to some minimum close to the global one. This lead us to extrapolate as a general criterion that, the activation functions that saturate the latest and at lower values, i.e. sigmoid, hyperbolic tangent and swish, are preferable to the exponential or softplus activation functions, and thus, we will prioritize the in the upcoming simulations (which does not mean that for some particular instance an exponential or softplus activation could outperform the others).

As a second part of this benchmarking, we will compare possibilities of the other crucial choice in training, the optimizers. We will be using the same set-up as before, but this time instead of fixing the optimizer, we will be fixing the activation function to be sigmoids.

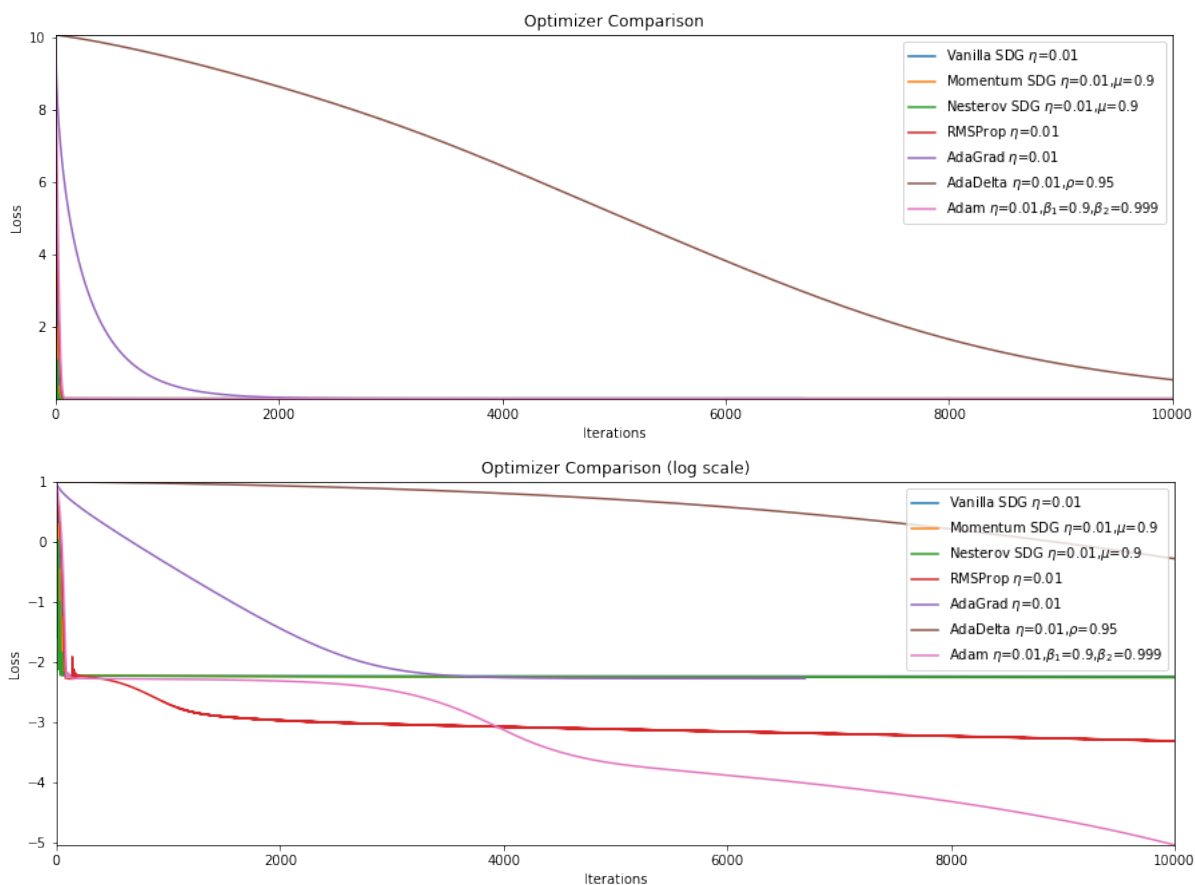


Figure 3.2: Comparison for different first order optimizers training performance for a $[3,4,1]$ -ANN, with sigmoid activations. Lower image in \log_{10} scale.

From these first 10000 iterations, all using a learning rate of $\eta = 0.01$ (the rest of the hyperparameters in the optimizers can be drawn from the legend in Figure 3.2), we can see many of the behaviours expected from section 2.4.1. Vanilla, classic and Nesterov momentum SGD, all had a very similar performance reaching an almost-saturation around the same loss value, which means that at that point we should have reduced manually the learning rate. Although, it is hard to discern from Figure 3.2, Nesterov momentum had the fastest initial decrease until reaching the state of almost-saturation, followed by classic momentum and vanilla SGD, as expected.

As for the optimizers that adjust the learning component-wise, AdaGrad and AdaDelta performed poorly. In the AdaGrad case, which completely stopped learning in iteration 6696, it matches well with what we had explained in the previous chapter, the exponential accumulation of previous gradients in $G(t)$ from (2.30) vanished the update gradient. For the AdaDelta case, the method is steady but starts slow, which occurs because of the quotient of (2.34). Since the moving average of gradient and the increment $\Delta(w_t, b_t)$ are similar, the quotient that yield the actual increment ~ 1 , making the starting convergence small.

On the other side of the component-wise adapting optimizers lay the RMSProp and Adam methods. These will be our preferred methods, especially Adam, since they yield very good results and have great autonomy (require little adjusting of the learning rate).

Concluding this discussion on optimizers, in the next Figure 3.4, we show the performance of the BFGS and L-BFGS, which are second order optimizers we use in this work. As it should be the behaviours are the exact same, since both, the limited memory version and the original methods are the same method, the only change being that the original method keeps the required matrices and vectors in memory, and the limited memory version stores only vectors by keeping the already multiplied matrix vector products. The line search will have 10 attempts to reduce the coefficient α . Comparing the performance of these methods to the first order ones, we see that the gains are much smaller (the figure does not use logarithmic scale) and each step takes much more time. While this is true, it is also true that the method is much more steady than the first order ones, and can further make gains in regions where first order methods get stuck. Hence, we will use these methods, especially L-BFGS, to unblock the training when first order methods stuck (mostly because of valleys).

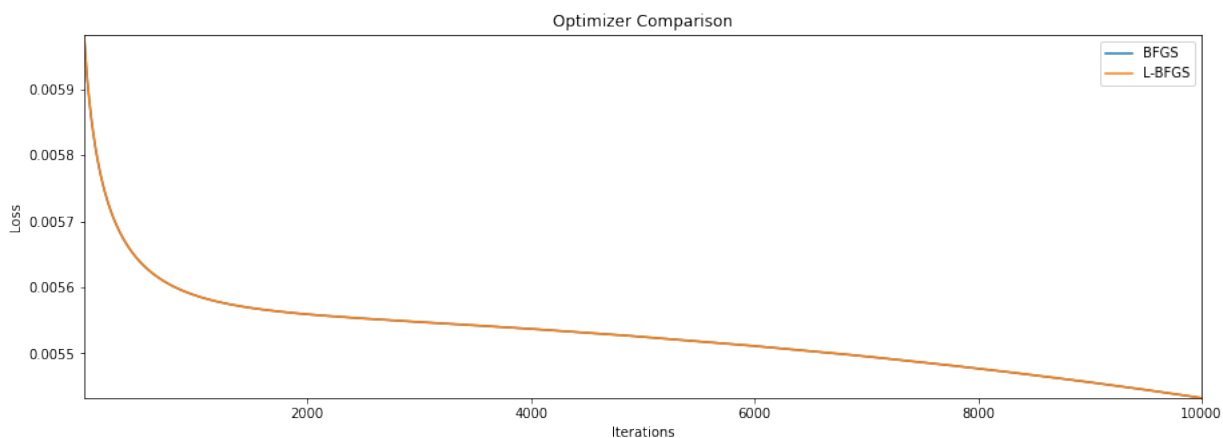


Figure 3.3: Training performance of a [3,4,1]-ANN with sigmoid activations, to fit (3.1), using BFGS and L-BFGS.

Finally we will show the end result of fitting the [3,4,1]-ANN model with sigmoids, using the Adam optimizer in the same conditions as before, for executions for 15000 iterations, with a training set of 10000 training points randomly collocated using a uniform distribution on the [0,1] segment. In the next figure will represent the model against the exact solution and the loss function for the iterations of training. We can see that model and the exact solution overlap very well, with a final total loss achieved of $L = L_{sol} = 4.331806e - 07$.

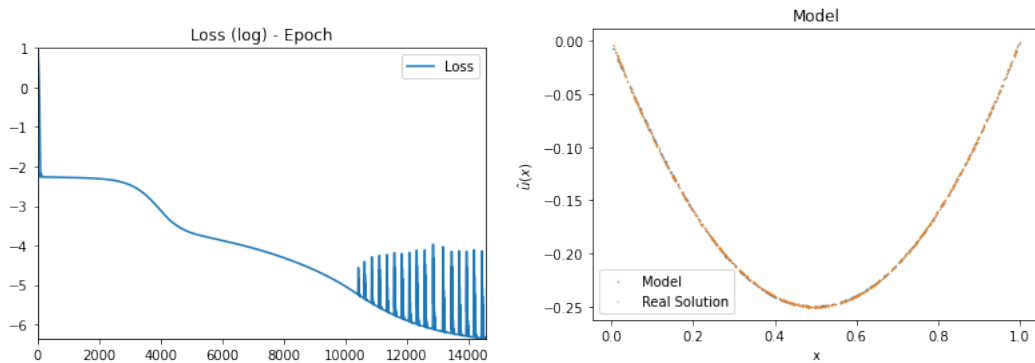


Figure 3.4: Training performance of a $[3,4,1]$ -ANN with sigmoid activations, to fit (3.1), using Adam with $\eta = 0.01$.

3.3 Training with Multi-Objective Loss Functions

In this section we will be discussing some of the issues that arise when training multi-objective functions, and we will propose a method to adjust the gradients that will stabilize the training process.

One of the main inconveniences in this work is that our target loss function (2.1) is composed of two very different components: L_1 satisfying the differential operator and L_2 satisfying the initial/boundary conditions. However, in reality, we actually want to approximate the exact unique analytic solution to the initial/boundary problem, $u(x)$. Therefore, this means that ideally, we would want to be minimizing:

$$L_{sol}(w, b) \approx \frac{1}{N_{\Omega} + N_{\Gamma}} \sum_{x_i \in N_{\Omega}, N_{\Gamma}} \left(u(x_i) - \hat{u}(x_i; w, b) \right)^2. \quad (3.3)$$

Since the exact solution, $u(x)$, is generally unknown, as this is the whole purpose of using a numerical integration method, $L_{sol}(w, b)$ cannot be computed. The question thus is, how well does the target loss $L(w, b)$ (2.1) and its terms correlate with the ideal loss $L_{sol}(w, b)$ (3.3)?

Throughout this work we have introduced ad-hoc initial/boundary conditions, as well as external forces, to modify the models so that the unique solution is known and has a simple form, which means that $L_{sol}(w, b)$ could always be computed. Observing how the models in the next sections behaves, we can conclude that the relation between the target and ideal losses is not strictly correlational. Indeed, although $L(w, b)$ and $L_{sol}(w, b)$ have their only global minimum at 0, which corresponds to the model being the exact solution $\hat{u} = u$ (the desired end result and the one we would surely obtain optimizing both if the problem was purely convex), both losses describe completely different hyper-surfaces. In fact, we could argue that $L(w, b)$ is a poor choice of a loss function, since $L_1(w, b) = 0$ for any solution of the differential operator satisfying any other set of initial/boundary conditions, and $L_2(w, b) = 0$ for any function that satisfies the initial/boundary conditions. This implies that there is an infinite number of potential local minima where the optimization can fall in, or at least by the standards of this work as we expect some continuity in with respect to the operators and conditions, extremely flat valleys.

Imagine, we would want to approximate the solution of $du/dx = 2x - 1$ with $u(0) = 0$, which is nothing more than $u(x) = x^2 - x$, with a model of the type $\hat{u} = ax^2 + bx + c$. The losses corresponding to the parameter values $a = 1.5$, $b = -2$ and $c = 0.05$, for the point 0.5 of the domain take values $L_1 = 0.25$ and $L_{sol} = 0.105625$, and for the second point 0.1 of the domain take values $L_1 = 0.81$ and $L_{sol} = 0.002025$. Through this basic example we have seen that, not even in this simplified model, the loss functions at two different points are positively correlated, as in the first case $L_1 \downarrow L_{sol} \uparrow$ and in the second case $L_1 \uparrow L_{sol} \downarrow$. On top of this, we see that the magnitudes of both losses are different, and furthermore this proportion does not remain constant during training. This means that the scenario where in the first stages of training, a decrease from 10^{-1} to 10^{-4} in L_1 represents a 10^{-1} to 10^{-2} decrease in L_{sol} for the domain points, and the later on in training, a decrease from 10^{-4} to 10^{-6} in L_1 represents a 10^{-1} to 10^{-5} decrease in L_{sol} for the domain points, is very plausible. This poses one of main problems of this method for integrating differential equations compared to others, such as FEM, which is: in a general situation where the exact solution is unknown, it is not possible to know the real error of the approximated artificial neural network solution, \hat{u} . (The same arguments are valid with respect between L_2 and L_{sol} , and L_1 and L_2).

The fact that the terms L_1 and L_2 behave so differently between them, and during training, makes simultaneous training of these two objectives quite hard. For this work, we always chose the domain in which we want to approximate the solutions to be $\Omega = [0, 1]$ in the one dimensional case, and $\Omega = [0, 1] \times [0, 1]$ in the two dimensional case. Under this regimes, the two main behaviours of the losses after the first steps of training, when the optimization has stabilized, are that $L \ll L_{sol}$, and $\nabla_{(w,b)} L_2(w, b) \gg \nabla_{(w,b)} L_1(w, b)$ always (at least in this work). Both are expected characteristics of the training: the first is related to the explanations of the previous paragraph as having similar information about the derivatives is not the same as being the same function; and the second occurs due to the initial/boundary conditions always having information of one derivative less than the differential operator, thus the vanishing gradient problem (2.21) with respect to input derivatives explained in 2.3.1.

Having $\nabla_{(w,b)} L_2(w, b) \gg \nabla_{(w,b)} L_1(w, b)$ is a real problem. In fact if we were to optimize the model with what tools we have discussed until now, we would not be able to model anything past the divergence operator, as the artificial neural networks would fit the initial/border conditions with great precision and nothing else. At this point the optimization gets stuck in a very flat valley where the gradient is basically too weak and noisy that not even L-BFGS can find a minimum loss improvement in that direction. We will illustrate this problem with the following example shown in Figure 3.5.

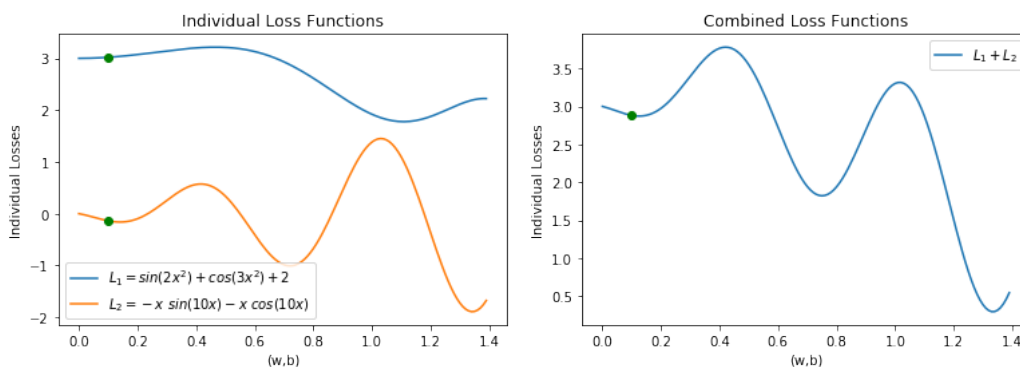


Figure 3.5: Example of possible multi-objective functions. Component and total representation.

Suppose that we wanted to find the simultaneously the minimum of the two individual loss functions, L_1 and L_2 shown in Figure 3.5 (left), and thus, we consider the global loss function composed by their addition represented in Figure 3.5 (right). The actual expressions of the functions can be drawn from the figure legends, but are irrelevant as are using them for conceptual reasons. Observe that the term L_1 is much flatter (has smaller derivatives) than the term L_2 , hence, when adding both functions together, the topology of the end result becomes mostly that of L_2 . When it comes to the term L_1 , it acts as though it is little more than a constant moving up or down the global loss, and the whole variation is dominated by L_2 . Therefore, minimizing the global loss L effectively becomes minimizing the L_2 with complete disregard to L_1 . For example, if we were to optimize the global loss starting at the green point of the figure (around 0.1), we would be moving to the right and most likely get stuck in the local minimum of the second term (around 0.7), which is the same that would happen optimizing only L_2 . This might be acceptable if we were only minimizing L_2 , since we might accept this local minimum as a good enough result. However, at around 0.7 the term L_1 still has a very clear slope that points to the right, and quite strong for L_1 standards, although it is many orders of magnitude weaker to its counterpart in L_2 . Hence, this slope for L_1 is completely ignored, meaning that we are not minimizing L_1 at all. Finally, the question is, how to balance the influence of both terms, L_1 and L_2 ?

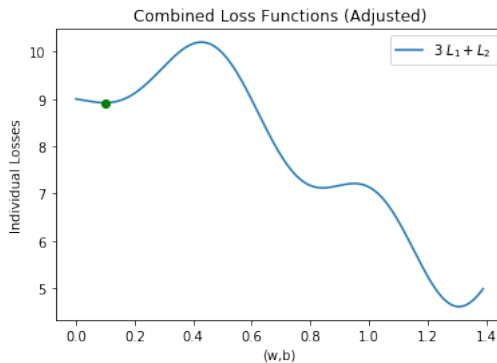


Figure 3.6: Example of possible multi-objective functions. Adjusted factors.

In Figure 3.6 we show what happens when we scale a term relatively to the other (here by a factor of 3). Note that, when multiplying a term by a factor, we also multiply their slopes and the variance it contributes to the global loss by that same factor. Hence, the final global loss function becomes a much more balanced agreement between the two terms. Here, minimizing the global loss starting at the green point and reaching the region at 0.7 does not represent such a strong sink for a gradient based approach to overcome. The region is still a good local minimum for L_2 , but the slope of L_1 which disagrees showing a strong direction of decrease for L_1 to the right, now has enough influence over the global loss to have an effect. Balancing the terms makes it so that the minimum found for global loss is at least a minimum (or close enough) for each of the individual terms L_1 and L_2 .

Going back to the problem of $\nabla_{(w,b)}L_2(w,b) \gg \nabla_{(w,b)}L_1(w,b)$ in the regime we are dealing with, i.e. the differential operator loss L_1 is much flatter than the initial/border loss L_2 , we would like to find a consistent way to balance the loss terms for any instance. The approach that we will be taking in this work to deal with this issue does not use a scale factor, as the example in Figures 3.5 and 3.6, but make use of the exact same underlying ideas.

We will still be computing the global loss as $L(w, b) = L_1(w, b) + L_2(w, b) + R(w, b)$ (2.1), and we will still be computing the each individual gradients with respect to the components $\nabla_{(w,b)}L_1(w, b)$, $\nabla_{(w,b)}L_2(w, b)$, $\nabla_{(w,b)}R(w, b)$, which is done separately since the global loss is additive. However, the difference here will be that, the gradient of the global loss will not be calculated as $\nabla_{(w,b)}L(w, b) = \nabla_{(w,b)}L_1(w, b) + \nabla_{(w,b)}L_2(w, b) + \nabla_{(w,b)}R(w, b)$ but using:

$$\begin{aligned} \nabla_{(w^{[\ell]}, b^{[\ell]})}L(w^{[\ell]}, b^{[\ell]}) \rightarrow & \\ & \frac{\nabla_{(w^{[\ell]}, b^{[\ell]})}L_1(w^{[\ell]}, b^{[\ell]})}{\|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_1(w^{[\ell]}, b^{[\ell]})\|_2^2} \cdot \min(\|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_1(w^{[\ell]}, b^{[\ell]})\|_2^2, \|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_2(w^{[\ell]}, b^{[\ell]})\|_2^2) \\ & + \frac{\nabla_{(w^{[\ell]}, b^{[\ell]})}L_2(w^{[\ell]}, b^{[\ell]})}{\|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_2(w^{[\ell]}, b^{[\ell]})\|_2^2} \cdot \min(\|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_1(w^{[\ell]}, b^{[\ell]})\|_2^2, \|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_2(w^{[\ell]}, b^{[\ell]})\|_2^2) \\ & + \lambda \cdot \frac{\nabla_{(w^{[\ell]}, b^{[\ell]})}R(w^{[\ell]}, b^{[\ell]})}{\|\nabla_{(w^{[\ell]}, b^{[\ell]})}R(w^{[\ell]}, b^{[\ell]})\|_2^2} \cdot \min(\|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_1(w^{[\ell]}, b^{[\ell]})\|_2^2, \|\nabla_{(w^{[\ell]}, b^{[\ell]})}L_2(w^{[\ell]}, b^{[\ell]})\|_2^2), \end{aligned} \tag{3.4}$$

which means that, before adding up the individual losses (and regularization) gradients, we renormalise the sub-gradient of each layer to the minimum between L_1 and L_2 . Geometrically, we are adjusting the slopes corresponding to one layer of L_1 , L_2 and R (corrected by λ), to be of the same magnitude, and thus, make each of the terms equally influential to the global loss (the regularization corrected by λ), which was the original intention.

This fixes really well the vanishing gradient problem with respect to the derivatives of the inputs, and the reason is, in some sense, analogous to why the first order methods which we called *component learning rate adapting* (Adam et al.) deal so well with the classic vanishing gradient problem. In those methods the gradients are adjusted component-wise, i.e. each parameter derivative is modified individually based on the changes for the previous vales of that same parameter. Hence, although the gradient gets deformed (we are no longer moving in maximal slope direction), it is much preferable as it allows for a much more equilibrated learning where all the parameters learn at similar rates, increasing naturally smaller parameter derivatives, when the slope in that direction is steady. Similarly, what component learning rate adapting optimizers do for parameters, (3.4) does among the global loss terms, L_1 , L_2 and R , allowing to rescale in a balanced way naturally smaller sub-gradients, which otherwise would have no contribution.

Correcting the gradient before introducing it the optimizer is much convenient than altering the global loss function $L(w, b)$ terms with coefficients (recalled we already gave some arguments when discussing weight penalties in section 2.6.2). First, we only affect the step we take during the optimization, thus $L(w, b)$ and their terms remain unaltered and consistent metric of the performance of the model throughout the whole training. Second, it is really the gradient (slope of the hyper-surface) and not the actual value of the loss (which can be naturally high or low) that conditions the optimization process.

Finally, comment there is a second multi-objective behaviours due to the loss terms effectively being the sum over a collocation of points in the domain and initial/boundary conditions. A possible fix to this is to increase the power of the norm, $|\cdot| \rightarrow \|\cdot\|_2^2 \rightarrow \|\cdot\|_n^n$, which is known to make the problem more sensible to outliers. Since we are not using noise, and the collocation points expected outputs are deterministic, as a result of satisfying the initial/boundary value problem, making the model more sensible to fit the regions where there are more discrepancies between the model and the expected outputs can be very beneficial.

3.4 Model Simulation

At last, in this section we will be numerically integrating initial/boundary value problem with the operators already introduced in Table 1.2 (except for the identity operator which already got covered in section 3.2). As stated before we will add external forces and initial/boundary conditions, so that the exact solution, $u(x)$, has a simple polynomial form and the exact loss L_{sol} (3.3) can be computed.

The general layout for all instance will consist on several distributions of layers and neurons, activation functions, optimization methods, and regularizations, that we will vary and compare among each other. In terms of the sampling, we will always approximate the solutions in the domain $\Omega = [0, 1]$ for 1D, and $\Omega = [0, 1] \times [0, 1]$ for 2D, hence, we will always draw samples of 10000 and 1000 points, for the domain and initial/boundary training sets, and validation sets, respectively, from $\sim \mathcal{U}(0, 1)$ distributions. While a very small collocation sample may lead to very poorly fitted model, there is a point where increasing the size of the sample does not improve the result of the model, which has to be located with trial and error. For most of the instances in this work it is probably around 1000 points, thus we have increased it 10 fold for safety reasons. Also implicitly applied through all the instances, there will be upper bound parameter constrains of 10^3 , gradient clipping by layer with norm 1, early stop evaluation every 500 training steps (or epochs), and the gradient adjustment (3.4).

Each of the instances will be increasing in complexity. Furthermore each instance may contain optional comparisons for some of the aspects that we have already discussed for benchmarking reasons.

3.4.1 Model 1: The 1D Divergence Operator

For the first proper initial/boundary problem we will be approximating the 1D divergence operator (derivative equals to an external force), which is one of the most simple models that we can generate. The exact instance of this ODE will be:

$$\begin{aligned} \nabla_{(x)} \cdot u(x) &= \frac{\partial u(x)}{\partial x} = 2x - 1, \\ u(0) &= 0, \end{aligned} \tag{3.5}$$

which has exact solution $u(x) = x^2 - x$. The solution can be easily obtain through the separation method by integrating both sides:

$$\begin{aligned} \int_0^x \frac{\partial u(x)}{\partial x} dx &= \int_0^x 2x - 1 dx, \\ u(x) \Big|_0^x &= x^2 - x \Big|_0^x, \\ u(x) - u(0) &= (x^2 - x) - (0^2 - 0), \\ u(x) &= (x^2 - x). \end{aligned} \tag{3.6}$$

Here the loss function we are optimizing is:

$$L(w, b) = \frac{1}{N_\Omega} \sum_{1 \leq i \leq N_\Omega} \left(\frac{\partial \hat{u}(x_i; w, b)}{\partial x} - (2x_i - 1) \right)^2 + (\hat{u}(0; w, b) - 0)^2 + R(w, b). \tag{3.7}$$

For this instance we will focus on the effect in the models of varying the activation functions. We will be training 3 model with sigmoid, hyperbolic tangent and swish activation functions respectively. All model will be trained on 3000 iterations (epochs), using a [1,5,5,1]-ANN, with no regularization, and using Adam with $\eta = 0.01$, $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Table 3.1 shows the end result losses, and Figure 3.7 plots the performance of the training (the losses for the figure have been broken into each of its components).

Activation	L	L_{sol}
Sigmoid	$1.66 \cdot 10^{-4}$	$1.71 \cdot 10^{-6}$
Tanh	$1.20 \cdot 10^{-4}$	$4.02 \cdot 10^{-7}$
Swish	$6.41 \cdot 10^{-5}$	$1.75 \cdot 10^{-6}$

Table 3.1: Results of 3 models trained for a [1,5,5,1]-ANN scheme, with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 3000 epochs. (3.5)

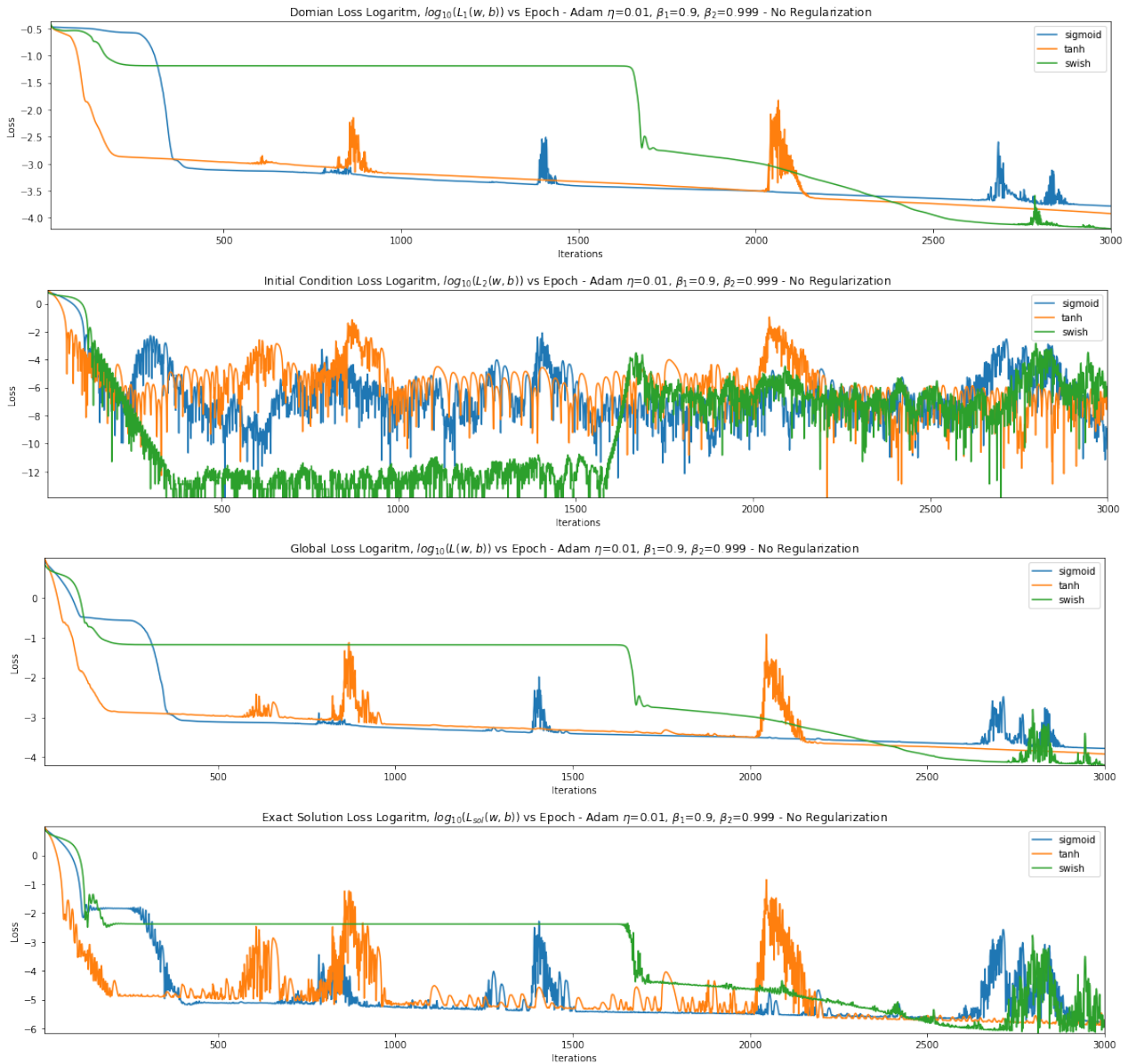


Figure 3.7: Training performance of 3 models trained for a [1,5,5,1]-ANN scheme, with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 3000 epochs. (3.5)

From the previous results we can see some interesting behaviours. First of all, the worst performing activation (the swish) with regards to the solution loss L_{sol} , was actually the best in terms of the objective loss L ; and the best performing activation (the tanh) with regards to the solution loss L_{sol} , was not the best in terms of the objective loss L . Also, we see that there is correlation between L (3rd plot of Figure 3.7) and L_{sol} (4th plot of Figure 3.7); and that L_1 dominates over L_2 , meaning L_2 is natively much smaller than L_1 . As well as this, we observe that the swish model had a much later initial decay than the other two, but the three of them start saturating at the same time. All these are expected behaviours that we have explained before.

Finally, in the next figure we plot the output of the best performing model (the one with tanh activations), against the exact solution. Note that just in 3000 epochs (2min) the match is almost perfect.

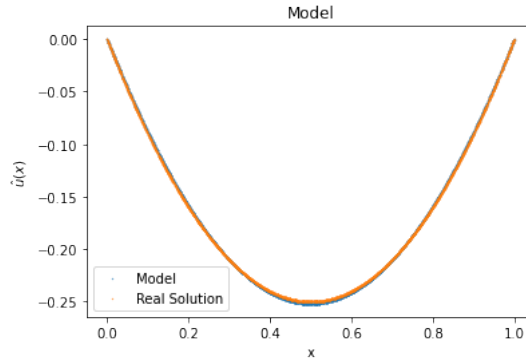


Figure 3.8: Final results. Best performing trained model (tanh) for (3.7) against the exact solution.

3.4.2 Model 2: The 2D Divergence Operator

Here we will take the previous model to the next level adding a dimension, and in doing so we will consider our first PDE. Still, this will be a very simple problem. The instance we will be considering first is:

$$\nabla_{(x,y)} \cdot u(x,y) \cdot \mathbf{1} = \frac{\partial u(x,y)}{\partial x} + \frac{\partial u(x,y)}{\partial y} = (2x-1) \cdot (y^2-y) + (x^2-x) \cdot (2y-1), \quad (3.8)$$

$$u(x,0) = 0, \quad x \in (-\infty, \infty),$$

which has exact solution $u(x,y) = (x^2-x) \cdot (y^2-y)$. The loss function for (3.8) would be:

$$L(w,b) = \frac{1}{N_\Omega} \sum_{1 \leq i \leq N_\Omega} \left(\frac{\partial \hat{u}(x_i, y_i; w, b)}{\partial x} + \frac{\partial \hat{u}(x_i, y_i; w, b)}{\partial y} - (2x_i-1)(y_i^2-y_i) - (x_i^2-x_i)(2y_i-1) \right)^2 + \frac{1}{N_\Gamma} \sum_{1 \leq i \leq N_\Gamma} (\hat{u}(x_i, 0; w, b) - 0)^2 + R(w,b). \quad (3.9)$$

However, there is an issue when using the (3.9) loss function. The border conditions are described by a curve for $x \in (-\infty, \infty)$, but effectively, we cannot draw samples from such a wide range. Since, we are limiting ourselves to approximating the solutions in the domain to $\Omega = [0, 1] \times [0, 1]$ for practical reasons, we will sample x from $(-10, 10)$ for the L_2 term.

The following Figure 3.9 shows the results of training a model under the previous assumptions (specifics in the caption). Observe that the left plot shows that, the solution approximated by the model has two separate regions, one approximating really well the exact solution, and another one that does not by a large margin. If we turn to the right plot we see that the MSE of the individual points with respect to the differential operator/external force is very even, meaning every point is equally well fitted.

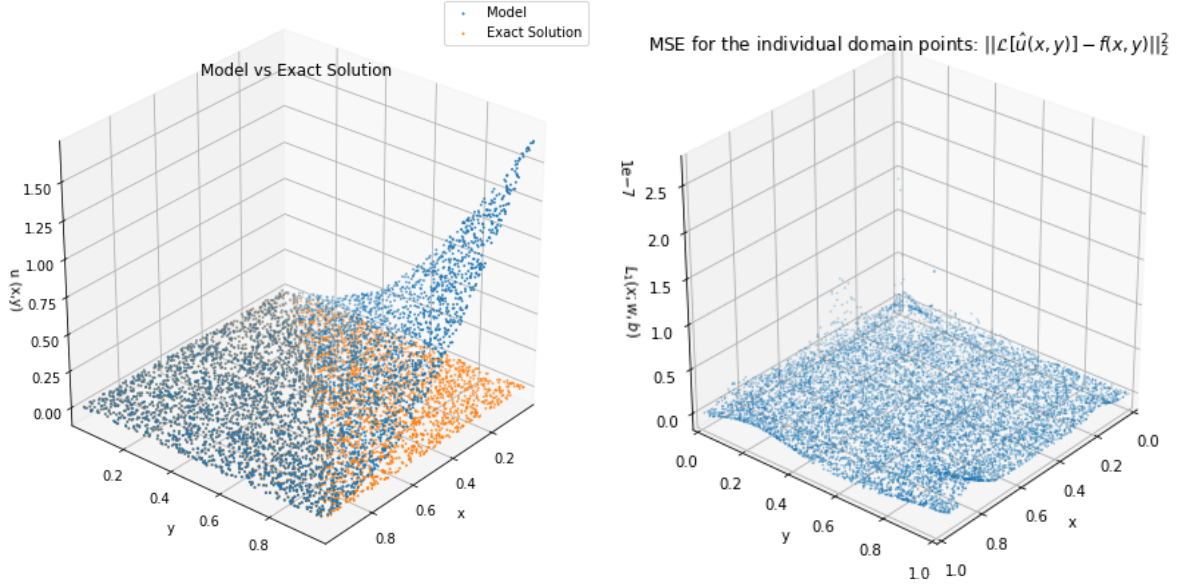


Figure 3.9: Result of a [1,10,10,1]-ANN model and tanh activations, trained with no regularization, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 12000 epochs. Left plot: model against exact solution. Right plot MSE error of the model, for each point in the domain.

This occurs because, in practice, when we draw a sample points for the border conditions, we are limiting ourselves to $x \in (-10, 10)$. Hence, for all purposes we are solving (3.8) with boundary conditions $u(x, 0) = 0$, $x \in (-10, 10)$, which are no longer Cauchy conditions and do not guarantee uniqueness. The solution we want to find is also a solution of the problem we are fitting in practice, but there are many more. In fact, what we see in Figure 3.9 is the artificial neural network overlapping two different solutions of the problem (the one closest to $y = 0$ is the one we would want). Thus, this a good example of what happens when integrating a problem which is not well-posed.

In order to fix this issue we will change the Cauchy “border” conditions, which for infinite domains would be simply an open curve, to its finite domain version, which requires the information over the border. This means, for $\Omega = [0, 1] \times [0, 1]$, changing (3.8) to:

$$\begin{aligned} \nabla_{(x,y)} \cdot u(x,y) \cdot \mathbf{1} &= \frac{\partial u(x,y)}{\partial x} + \frac{\partial u(x,y)}{\partial y} = (2x-1) \cdot (y^2-y) + (x^2-x) \cdot (2y-1), \\ u(x,0) &= 0, \quad u(x,1) = 0, \quad x \in (0,1), \\ u(0,y) &= 0, \quad u(1,y) = 0, \quad y \in (0,1), \end{aligned} \tag{3.10}$$

with its respective change in the loss function (3.9). The solution for this problem is the same as before.

In the actual experiments for (3.10) we will take more interesting features to compare than the simple activation functions of the previous instance. Here we will analyse the effects of the size of the artificial neural network and the regularization.

When choosing an artificial neural network architecture, the general rule is that deeper neural networks are able to learn more complex functions, although at a greater cost of training [54]. Furthermore, papers such as [55], focused on learning polynomials with artificial neural networks, suggest that a fully-connected network with a single hidden layer with a number of nodes equals to the degree of the polynomial, would be enough to learn a polynomial (this is an rough and imprecise extraction of what [55] states, but holds for the most part). In this work though, we have been using two hidden layers so far (an will keep using them), and a much larger number of neurons than the theoretic minimal suggests for the underlying solutions we want to approximate. The reason for doing this is to better account for the information of the derivatives during training and make use of regularization techniques, to obtain better minima.

For this instance (3.10) we will be training 6 models, all using hyperbolic tangent activations and are trained on Adam with $\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs. The models will either have a [1,10,10,1]-ANN structure or a [1,40,40,1]-ANN structure; and be trained using no regularization, the custom regularization (2.58) with $\lambda = 1$, or a Tikhonov regularization with $\lambda = 1$; which make for a total of 6 combinations. Moreover, all the models with the same architecture have been initialized with exactly the same parameters. This has been done to root out the possible effect of luck for starting at a slightly better point for the optimization, and ensure the difference is training are caused by the regularization.

In Table 3.2 we show the final results of the models, and in Figures 3.10 and 3.11 we show the performance of the training. First, we observe for these kind of problems Tikhonov regularizations do not work well and their training incurs in early stopping. For the (2.58) custom regularization we see that, in the smaller [1,10,10,1]-ANN model, the training is actually hindered and yields awful results, but used the larger [1,40,40,1]-ANN model it outperforms any other set-up. This, is due to what we have already explained in section 2.6, that regularizations clamp down on the extra degree of freedom overfitting the model. Hence, for the smaller model which is adequately parametrized, it becomes an extra condition drawing resources form the model, while for the larger model it narrows the parameters to the fit the model. Furthermore, not only the larger model with regularization outperforms the smaller one without, but if we compare their performances from Figures 3.10 and 3.11, we note that by the end of the training, the smaller model has saturated (stagnated), while the larger is still steadily decreasing (thus, have more room for improvement). This shows that is much preferable to have a larger model with regularization than simply a well adjusted one.

Architecture - Regularization Technique	L	L_{sol}
[1,10,10,1]-ANN - No Regularization	$1.04 \cdot 10^{-4}$	$7.52 \cdot 10^{-6}$
[1,10,10,1]-ANN - (2.58) Regularization with $\lambda = 0.1$	$7.18 \cdot 10^{-3}$	$2.84 \cdot 10^{-4}$
[1,10,10,1]-ANN - Tikhonov Regularization with $\lambda = 0.1$	$6.33 \cdot 10^{-4}$	$4.19 \cdot 10^{-5}$
[1,40,40,1]-ANN - No Regularization	$6.36 \cdot 10^{-4}$	$1.95 \cdot 10^{-5}$
[1,40,40,1]-ANN - (2.58) Regularization with $\lambda = 0.1$	$2.93 \cdot 10^{-4}$	$2.32 \cdot 10^{-6}$
[1,40,40,1]-ANN - Tikhonov Regularization with $\lambda = 0.1$	$3.88 \cdot 10^{-3}$	$7.91 \cdot 10^{-5}$

Table 3.2: Results of 6 models with different architectures, trained for (3.10), using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs and different regularization techniques.

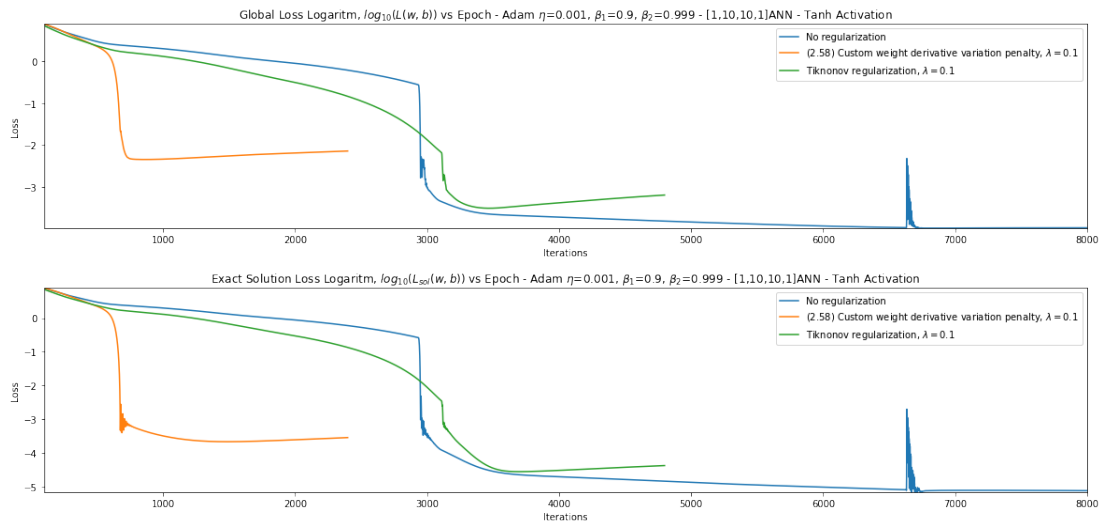


Figure 3.10: Comparison of different regularization techniques in training performance of 3 models trained for a $[1,10,10,1]$ -ANN scheme, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs. (3.10)

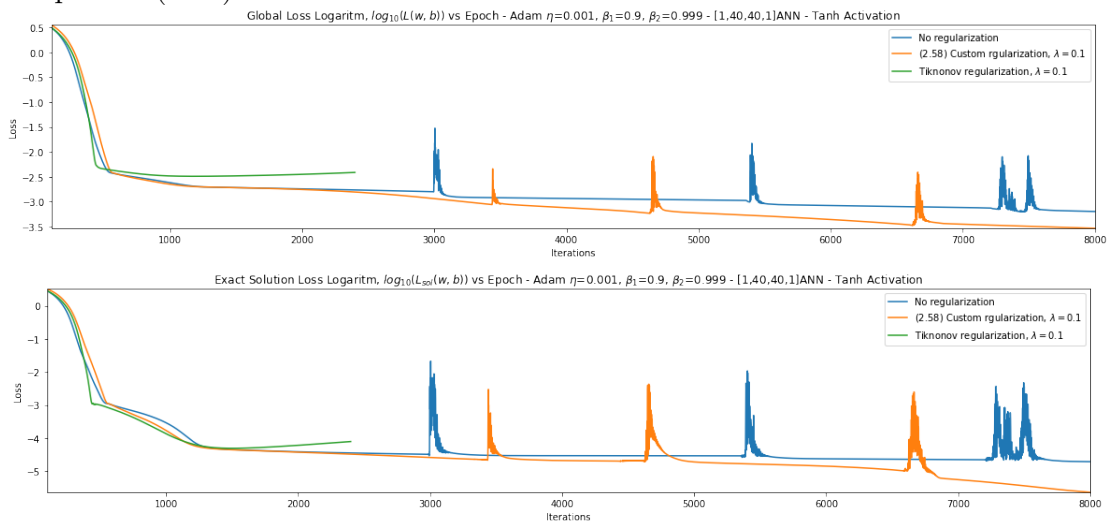


Figure 3.11: Comparison of different regularization techniques in training performance of 3 models trained for a $[1,40,40,1]$ -ANN scheme, using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, on 8000 epochs. (3.10)

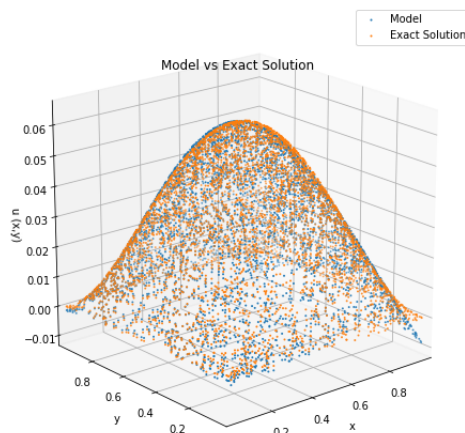


Figure 3.12: Final results of the best performing trained model ($[1,40,40,1]$ -ANN, trained with the custom regularization (2.58)) for (3.7) against the exact solution.

3.4.3 Model 3: The 2D Laplacian Operator

At this point we will complicate a bit more the differential operator by considering second order derivatives. Thus, we will consider the following boundary value problem for the Laplacian operator in 2 dimensions:

$$\Delta u(x, y) = \frac{\partial^2 u(x, y)}{\partial x^2} + \frac{\partial^2 u(x, y)}{\partial y^2} = 2 \cdot (y^2 - y) + 2 \cdot (x^2 - x),$$

$$u(\Gamma) = g_1(\Gamma) : \begin{cases} u(x, 0) = 0, & u(x, 1) = 0, & x \in (0, 1), \\ u(0, y) = 0, & u(1, y) = 0, & y \in (0, 1), \end{cases}$$

$$\frac{\partial u(\Gamma)}{\partial(x, y)} \cdot n(\Gamma) = g_2(\Gamma) : \begin{cases} \frac{\partial u(x, 0)}{\partial y} = -(x^2 - x), & \frac{\partial u(x, 1)}{\partial y} = (x^2 - x), & x \in (0, 1), \\ \frac{\partial u(0, y)}{\partial x} = -(y^2 - y), & \frac{\partial u(1, y)}{\partial x} = (y^2 - y), & y \in (0, 1), \end{cases} \quad (3.11)$$

which has exact solution $u(x, y) = (x^2 - x) \cdot (y^2 - y)$, as with the previous problem. The form of problem (3.11) in its general form, for any dimension and external force, constitutes what is called the Poisson equation, which is important throughout physics, as it is the interpretation of Gauss Law in terms of potentials.

Before training an artificial neural network to fit this model, we would like to make a brief note regarding the coding of higher order derivatives in TensorFlow. Looking at the official documentation of TensorFlow, the method given to obtain higher order derivatives in one variable is by nesting auto-differentiations calls. However, note that, TensorFlow is used in a context of training artificial neural networks, thus when auto-differentiating twice we obtain:

$$\begin{aligned} \nabla_{(x)} f(x_1, \dots, x_n) &= \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right), \\ \nabla_{(x)}^2 f(x_1, \dots, x_n) &= \left(\frac{\partial}{\partial x_1} \sum_{m=1}^n \frac{\partial f}{\partial x_m}, \dots, \frac{\partial}{\partial x_n} \sum_{m=1}^n \frac{\partial f}{\partial x_m} \right), \end{aligned} \quad (3.12)$$

which is not the Laplacian. There are two ways to overcome this issue: either use the *unstack* and *stack* functions to decouple the inputs and compute the gradients tracking only an individual variable (the option we have used in the code); or to use the *hessian* function to compute the Hessian matrix and then compute the trace, which is highly inefficient as we only require the elements in the diagonal. Without [56] where this observation is pointed out, we would not have been able to carry out this simulation.

At this point we have already experimented on all the principal options and hyper-parameter choices covered in this work, and we have studied their performance. So, from now on, we will be dropping the comparisons and limit ourselves to simply solve the next models with the best possible set-up best on what we have discussed.

The artificial neural network model trained for (3.11) has achieved a final global loss of $L = 1.23 \cdot 10^{-3}$ and final loss with respect to the solution of $L_{sol} = 4.25 \cdot 10^{-6}$. This model consisted of a [1,40,40,1]-ANN with tanh activations, trained for 6000 epochs (when early stop triggered), using Adam with $\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and the custom regularization (2.58) with $\lambda = 0.1$. The results can be seen in the following Figure 3.13.

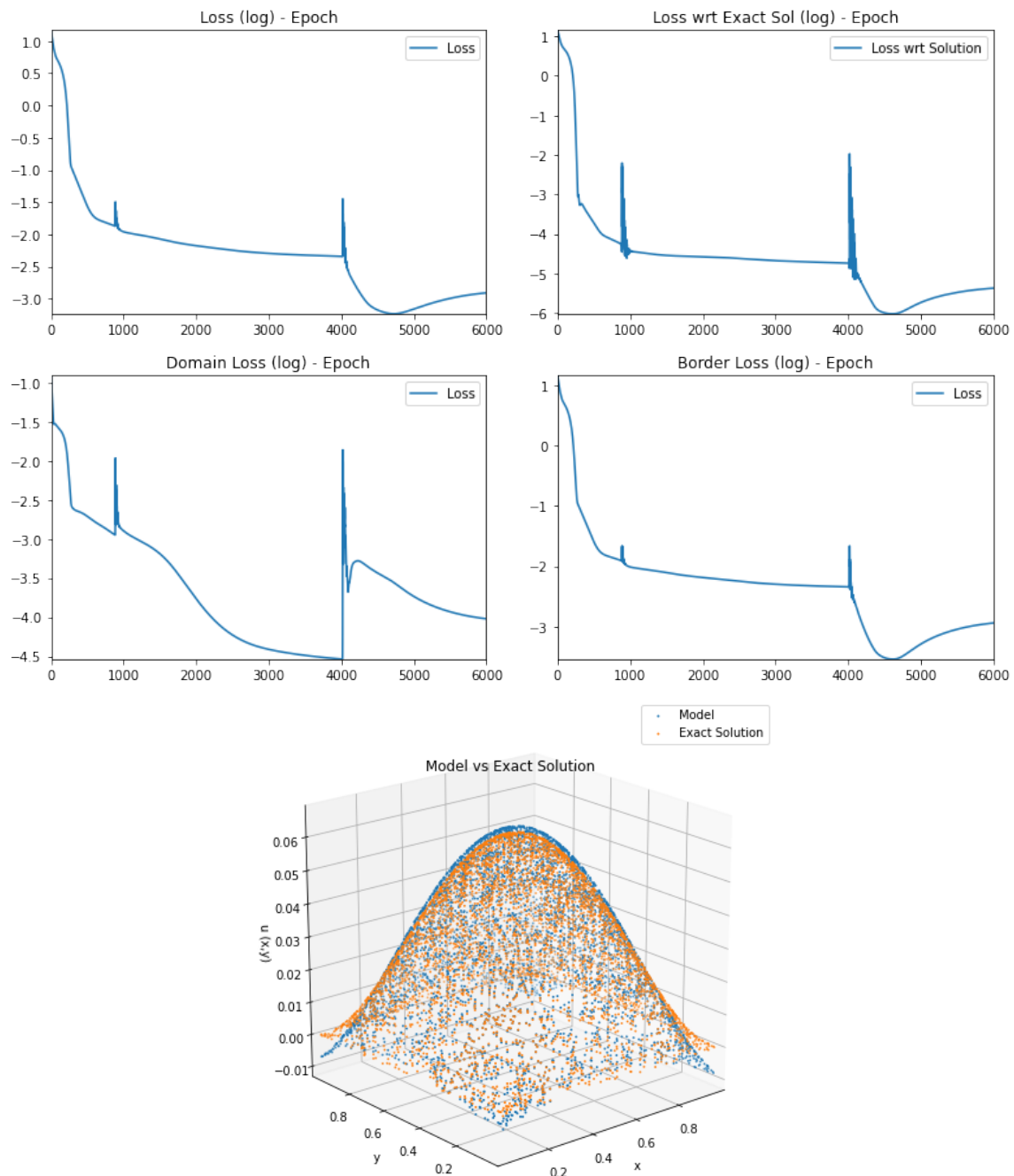


Figure 3.13: Results and performance of the model trained for (3.11).

3.4.4 Model 4: The 1D Advection Operator

For this simulation we step down from the 2D PDE cases, to go back to an ODE. The reason for this downgrade is to explain a certain issue occurring for this operator. This issue is one that happens for the final case of this section, the 2D Burgers operator, and since the advection operator we are proposing coincides with the Burgers operator in 1D, we see this as a much simpler example to introduce a discussion.

The initial value problem we want to consider is:

$$\begin{aligned} u(x) \cdot \nabla_{(x)} \cdot u(x) &= u(x) \cdot \frac{\partial u(x)}{\partial x} = 2x^3 - 3x^2 + x, \\ u(0) &= 0, \end{aligned} \tag{3.13}$$

which has exact solution $u(x) = x^2 - x$, same as the 1D divergence case. This problem look like falling under the Cauchy-Kovalevskaya conditions, so existence and uniqueness should be guaranteed. However, there is a subtlety hidden here. If we write the equation in its canonical form (isolating the higher derivative), which is required to apply the Cauchy-Kovalevskaya theorem,

$$\frac{\partial u(x)}{\partial x} = \frac{2x^3 - 3x^2 + x}{u(x)}, \tag{3.14}$$

we note that the equation is quasi-linear and its terms are analytic everywhere except for the zeroes of $u(x)$. Hence we have local existence and uniqueness almost everywhere, but since it can fail in some points, we cannot build a unique global solution using the theorem. This can be verified easily in this case, as the differential equation is separable and can be solved easily by separations of variables method:

$$\begin{aligned} \int_0^x u(x) \frac{\partial u(x)}{\partial x} dx &= \int_0^x 2x^3 - 3x^2 + x dx, \\ \frac{1}{2}(u(x))^2 \Big|_0^x &= \frac{1}{2}x^4 - x^3 + \frac{1}{2}x^2 \Big|_0^x, \\ \frac{1}{2}(u(x))^2 - 0 &= \frac{1}{2}x^4 - x^3 + \frac{1}{2}x^2 - 0, \\ u(x) &= \pm\sqrt{x^4 - 2x^3 + x^2} = \pm(x^2 - x). \end{aligned} \tag{3.15}$$

Looking at Figure 3.14 we observe that the solutions intersect (hence, are not unique) at the roots of $u(x)$.

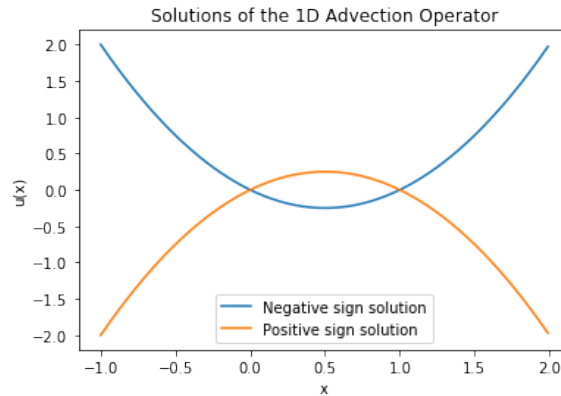


Figure 3.14: Positive and negative sign solutions of 3.13.

To fix this issue and fix a solution, it is enough to provide information about an extra derivative of one more order than the required by the Cauchy conditions. Therefore, the well-posed problem that we will consider will be:

$$\begin{aligned} u(x) \cdot \nabla_{(x)} \cdot u(x) &= u(x) \cdot \frac{\partial u(x)}{\partial x} = 2x^3 - 3x^2 + x, \\ u(0) &= 0, \quad u'(0) = -1. \end{aligned} \tag{3.16}$$

The artificial neural network model trained for (3.16) has achieved a final global loss of $L = 1.05 \cdot 10^{-3}$ and final loss with respect to the solution of $L_{sol} = 2.74 \cdot 10^{-7}$. This model consisted of a [1,20,20,1]-ANN with sigmoid activations, trained for 3000 epochs using Adam with $\eta = 0.01$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and the custom regularization (2.58) with $\lambda = 0.1$. The results can be seen in the following Figure 3.16.

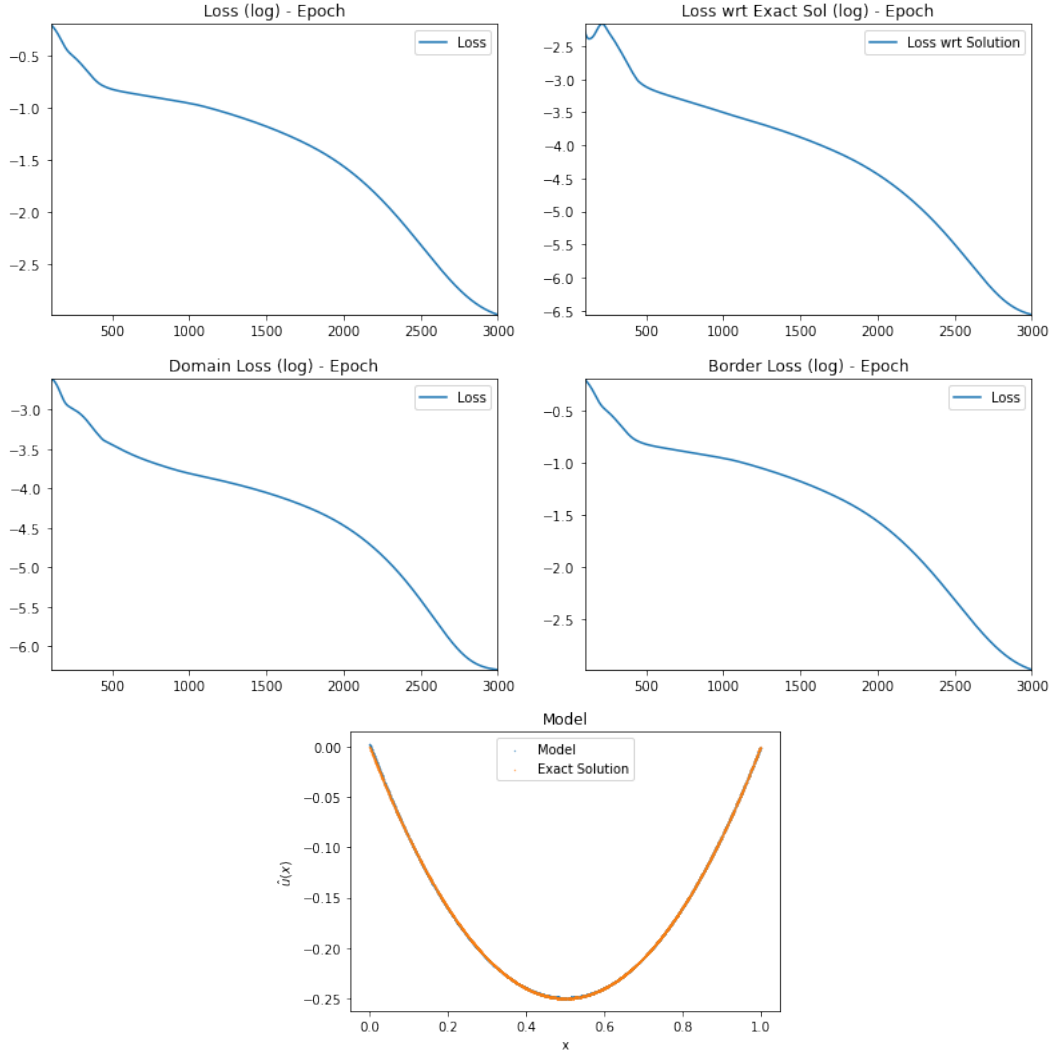


Figure 3.15: Results and performance of the model trained for (3.16).

3.4.5 Model 5: The 2D Clairaut Operator

The Clairaut operator can be seen as an upgrade to the 2D Advection case. It may not be much more complicated than what we have seen before, but it is the first PDE with non-constant coefficients that we integrate in this work. We pose its boundary problem as:

$$\begin{aligned}
 (x, y) \cdot \nabla_{(x,y)} u(x, y) &= x \cdot \frac{\partial u(x, y)}{\partial x} + y \cdot \frac{\partial u(x, y)}{\partial y} \\
 &= x \cdot (2x - 1) \cdot (y^2 - y) + (x^2 - x) \cdot y \cdot (2y - 1), \\
 u(x, 0) &= 0, \quad u(x, 1) = 0, \quad x \in (0, 1), \\
 u(0, y) &= 0, \quad u(1, y) = 0, \quad y \in (0, 1),
 \end{aligned} \tag{3.17}$$

with solution $u(x, y) = (x^2 - x) \cdot (y^2 - y)$, as always.

The artificial neural network model trained for (3.17) has achieved a final global loss of $L = 3.04 \cdot 10^{-6}$ and final loss with respect to the solution of $L_{sol} = 3.96 \cdot 10^{-6}$. This model consisted of a [1,40,40,1]-ANN with tanh activations, trained for 8000 epochs, using Adam with $\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and the custom regularization (2.58) with $\lambda = 0.1$. The results can be seen in the following Figure 3.16.

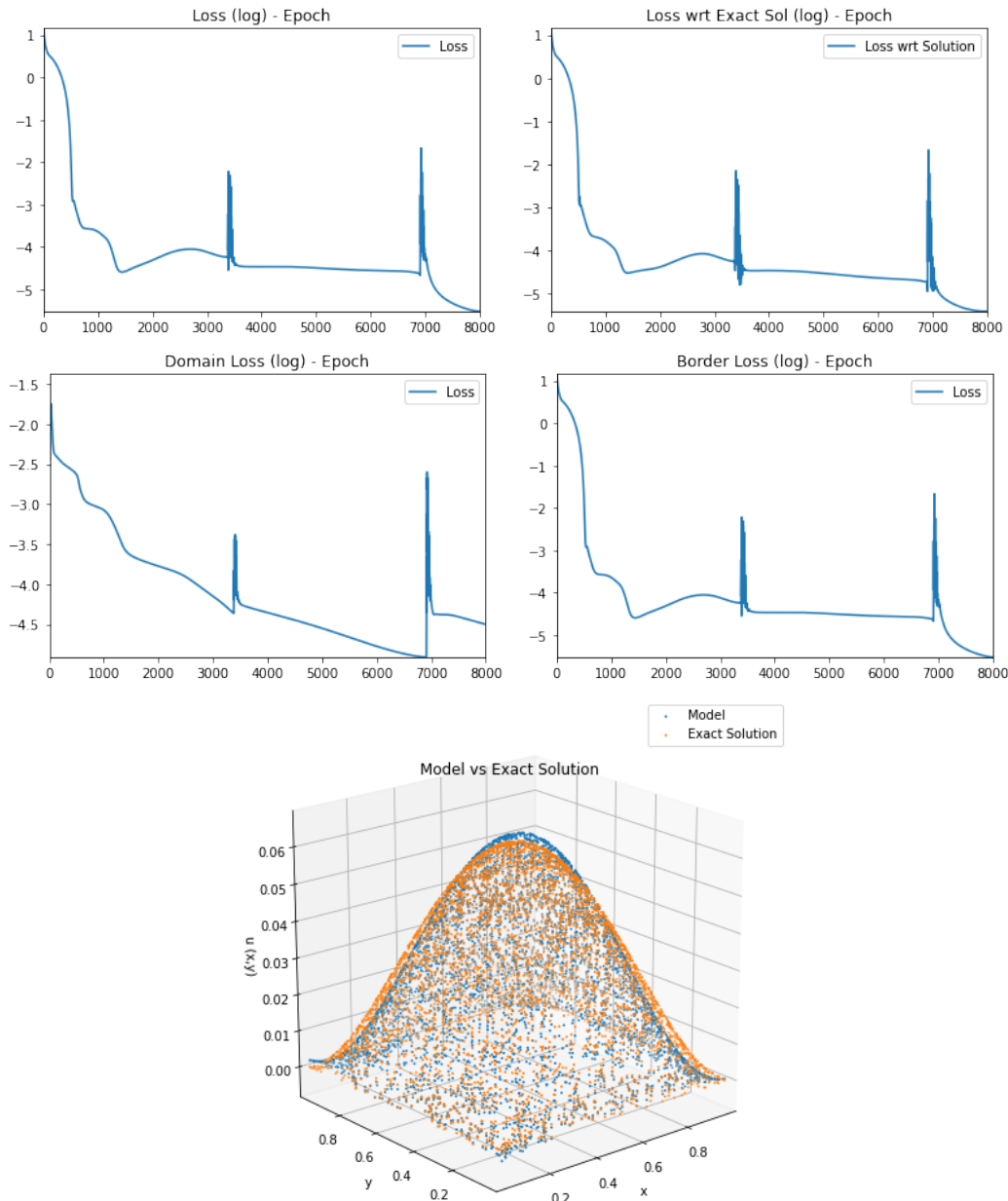


Figure 3.16: Results and performance of the model trained for (3.17).

3.4.6 Model 6: The 2D Burgers Operator

Finally, we will numerically integrate the last, and most complex, boundary problem of this work. This would be the 2D Burgers operator, and it can be regarded as the multi-dimensional case of the advection operator. While the advection operator is applied on scalar functions, the Burgers operator is applied on vector fields.

The boundary value problem requires, as explained in the 1D advection case, information on one more derivative than in the Cauchy conditions, and thus will be defined by:

$$\begin{aligned} \mathbf{u}(x, y) \cdot \nabla_{(x, y)} \mathbf{u}(x, y) &= (u_x(x, y), u_y(x, y)) \begin{pmatrix} \frac{\partial u_x(x, y)}{\partial x} & \frac{\partial u_y(x, y)}{\partial x} \\ \frac{\partial u_x(x, y)}{\partial y} & \frac{\partial u_y(x, y)}{\partial y} \end{pmatrix} \\ &= \begin{pmatrix} (x^2 - x) \cdot (y^2 - y)[(2x - 1) \cdot (y^2 - y) + (x^2 - x) \cdot (2y - 1)] \\ (x^2 - x) \cdot (y^2 - y)[(2x - 1) \cdot (y^2 - y) + (x^2 - x) \cdot (2y - 1)] \end{pmatrix} \\ \\ \mathbf{u}(\Gamma) = \mathbf{g}_1(\Gamma) : &\begin{cases} u_x(x, 0) = u_y(x, 0) = 0, & u_x(x, 1) = u_y(x, 1) = 0, & x \in (0, 1), \\ u_x(0, y) = u_y(0, y) = 0, & u_x(1, y) = u_y(1, y) = 0, & y \in (0, 1), \end{cases} \\ \\ \frac{\partial \mathbf{u}(\Gamma)}{\partial(x, y)} \cdot \mathbf{n}(\Gamma) = \mathbf{g}_2(\Gamma) : &\begin{cases} \frac{\partial u_x(x, 0)}{\partial y} = -(x^2 - x), & \frac{\partial u_x(x, 1)}{\partial y} = (x^2 - x), & x \in (0, 1), \\ \frac{\partial u_x(0, y)}{\partial x} = -(y^2 - y), & \frac{\partial u_x(1, y)}{\partial x} = (y^2 - y), & y \in (0, 1), \\ \frac{\partial u_y(x, 0)}{\partial y} = -(x^2 - x), & \frac{\partial u_y(x, 1)}{\partial y} = (x^2 - x), & x \in (0, 1), \\ \frac{\partial u_y(0, y)}{\partial x} = -(y^2 - y), & \frac{\partial u_y(1, y)}{\partial x} = (y^2 - y), & y \in (0, 1), \end{cases} \end{aligned} \quad (3.18)$$

with the vector field $\mathbf{u}(x) = ((x^2 - x) \cdot (y^2 - y), (x^2 - x) \cdot (y^2 - y))$ as its solution.

We have called the differential operator in (3.18) as the Burgers operator (although it is not the real burgers operator) because the most prominent place we can find this operator is in the inviscid Burgers equation, where it appears along an extra $\partial/\partial t$. Hence, we can think of (3.18) as the system that solves for the steady-state solutions $u_t = 0$ of an inviscid Burgers equation.

All the attempts to obtain the two components of the solution as an artificial neural network with 2 outputs resulted in failure. The models always seemed to get stuck in very bad local minima and the process was slow. We assume that the reasons for this lies in the fully-connected nature of the architecture we are using, since we were trying to adjust 2 functions on the same parameters (despite being the same in this case). Hence, we have taken the much efficient approach of approximating each component of the solution with a separate neural network, and obtained much better results. We will only show the results for the first component since the operator and external force are symmetric, therefore the behaviours for both components are the almost same.

The artificial neural network model trained for (3.18) has achieved a final global loss of $L = 6.21 \cdot 10^{-6}$ and final loss with respect to the solution of $L_{sol} = 8.75 \cdot 10^{-5}$. This model consisted of a [1,40,40,1]-ANN with tanh activations, trained for 8000 epochs, using Adam with $\eta = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$, and the custom regularization (2.58) with $\lambda = 0.1$. The results can be seen in the following Figure 3.17.

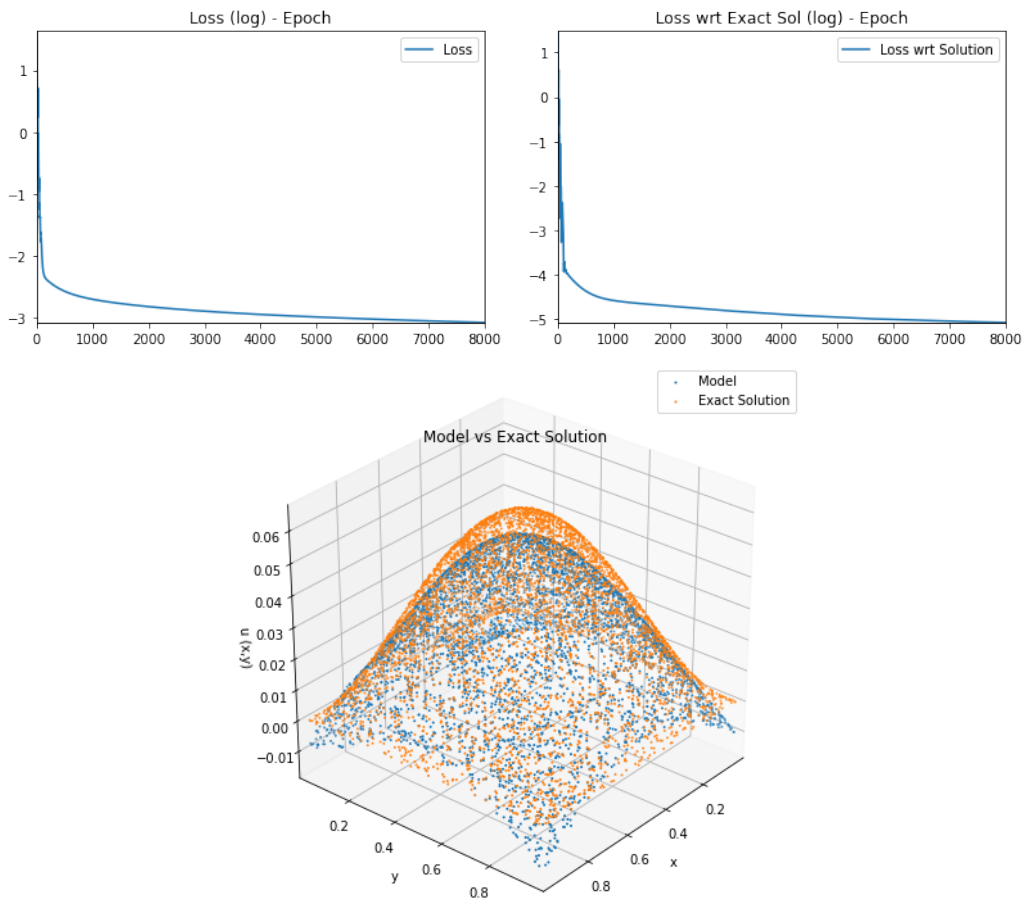


Figure 3.17: Results and performance of the model trained for (3.18).

Chapter 4

Conclusions

Regarding the method, in this work we have covered, theoretically and in practice, all the elements necessary to build a functional artificial neural network method to solve initial/boundary problem. Based on what we have seen as a whole we can conclude that this approach offers a solid option for integrating PDEs. Summarizing all of the pros and cons discussed so far:

- Pros:
 - Mesh-free algorithm: we do not need a mesh, a less costly random collocation of points approach is enough.
 - Higher dimensions of the same operator do not increase much the number of iterations necessary.
 - Good interpolation capacities for points outside of the collocation.
- Cons:
 - Unknown error behaviour: we do not know the precision we are getting and cannot use an objective stop criterion.
 - Unknown stability: we are subject to solving a non-convex optimization, which means that we incur in many potential traps such as local minima and valleys, yielding bad results.
 - Iterative method: we are required to evaluate the derivatives of the artificial neural network at each step of the training, thus the cost of an increase in complexity in the equations by having many different derivatives, accumulates exponentially through the training.

The first and the last cons can be somewhat mitigated. For the last we can try to use synthetic derivatives (i.e. draw them from a distribution or approximation). This idea is not foreign, as we have already seen in this work how the second order quasi-Newton methods (BFGS and L-BFGS) and truncated Newton methods (the Hessian-free), use approximations of the Hessian to avoid computing the whole matrix.

In the case of the error, looking back at the simulations we can observe that in the cases where we took the same operator in different dimensions, like with the divergence operator, the relation between the objective loss L and the exact solution loss L_{sol} behaved in a very similar way. This means that perhaps we could roughly use the results of a smaller problem in 1D to approximate the relation between L and L_{sol} .

To sum up, we would recommend this kind of method to deal with simple differential equations (not many different types of derivatives) in high dimensions and irregular domains. The strengths and weaknesses are close to that of Monte Carlo methods, so with further research it could become an alternative to those methods.

4.1 Author's Final Thoughts

This work contains a wide coverage of the topics general to the subject of deep learning, auto-differentiation, activation functions, optimizers, regularizations, problems like the vanishing and exploding gradient..., all of them applicable to contexts other than this project. The perspective of artificial neural networks through the view of tensors has also been very helpful.

As well as this, there has been an important part of programming involved using TensorFlow, and many bits of practical knowledge were gained through the experience of the many simulations. By experimenting on the cases we were able to identify many problems like the vanishing and exploding gradient with respect to inputs, and offer solutions such as the custom regularization (2.58) proposed in this job, and mechanism (3.4) to balance the gradients.

Finally, from a personal standpoint, I am really satisfied with this work since I have learned many things about differential equations and artificial neural networks, and I believe I have gained much expertise practical to understand and work confidently with those fields.

4.2 Further Work

As contents for further work it would be interesting to study the effect of different architectures and the relations between the objective loss and the solution loss to obtain a measure of the real error. An interesting architecture for PDEs could be the convolutional one, since it is known to preserve account for the spacial relations of its inputs.

Appendix A

Linear Algebra Formulation of 2.3.1

Here we will describe the elements and derivatives of the example model in Figure 2.6 using linear algebra notation for clarity purposes. We will limit ourselves to tensors in which the linear algebra representation is practical, i.e. up to (1,1)-tensors, or vector and matrices.

First, we will start by writing the variables and parameters (2.11) of the artificial neural network in linear algebra form:

$$\begin{aligned}
 y^{[\ell]}_{n_\ell} &\equiv \begin{pmatrix} y^{[1]}_1 \\ y^{[1]}_2 \\ y^{[3]}_3 \end{pmatrix}, & \begin{pmatrix} y^{[2]}_1 \\ y^{[2]}_2 \\ y^{[2]}_3 \\ y^{[2]}_4 \end{pmatrix}, & \begin{pmatrix} y^{[3]}_1 \\ y^{[3]}_2 \end{pmatrix}, \\
 z^{[\ell]}_{n_\ell} &\equiv \begin{pmatrix} z^{[1]}_1 \\ z^{[1]}_2 \\ z^{[3]}_3 \end{pmatrix}, & \begin{pmatrix} z^{[2]}_1 \\ z^{[2]}_2 \\ z^{[2]}_3 \\ z^{[2]}_4 \end{pmatrix}, & \begin{pmatrix} z^{[3]}_1 \\ z^{[3]}_2 \end{pmatrix}, \\
 w^{[\ell]}_{m_\ell} &\equiv \begin{pmatrix} w^{[1]}_1 & w^{[1]}_2 \\ w^{[1]}_2 & w^{[1]}_2 \\ w^{[1]}_3 & w^{[1]}_3 \end{pmatrix}, & \begin{pmatrix} w^{[2]}_1 & w^{[2]}_2 & w^{[2]}_3 \\ w^{[2]}_2 & w^{[2]}_2 & w^{[2]}_3 \\ w^{[2]}_3 & w^{[2]}_3 & w^{[2]}_3 \\ w^{[2]}_4 & w^{[2]}_4 & w^{[2]}_4 \end{pmatrix}, & \begin{pmatrix} w^{[3]}_1 & w^{[3]}_2 & w^{[3]}_3 & w^{[3]}_4 \\ w^{[3]}_2 & w^{[3]}_2 & w^{[3]}_3 & w^{[3]}_4 \end{pmatrix}, \\
 b^{[\ell]}_{n_\ell} &\equiv \begin{pmatrix} b^{[1]}_1 \\ b^{[1]}_2 \\ b^{[3]}_3 \end{pmatrix}, & \begin{pmatrix} b^{[2]}_1 \\ b^{[2]}_2 \\ b^{[2]}_3 \\ b^{[2]}_4 \end{pmatrix}, & \begin{pmatrix} b^{[3]}_1 \\ b^{[3]}_2 \end{pmatrix},
 \end{aligned} \tag{A.1}$$

The operations in each neuron, without Einstein Notation are: (2.12):

$$\begin{aligned}
 z^{[1]}_{n_1} &= b^{[1]}_{n_1} + \sum_{k=1}^2 w^{[1]k}_{n_1} \cdot x_k, & y^{[1]}_{n_1} &= a^{[1]}_{n_1} (z^{[1]}_{n_1}) \\
 z^{[2]}_{n_2} &= b^{[2]}_{n_2} + \sum_{k=1}^3 w^{[2]k}_{n_2} \cdot y^{[2]}_k, & y^{[2]}_{n_2} &= a^{[2]}_{n_2} (z^{[2]}_{n_2}) \\
 z^{[3]}_{n_3} &= b^{[3]}_{n_3} + \sum_{k=1}^4 w^{[3]k}_{n_3} \cdot y^{[3]}_k, & u_{n_3} &= y^{[3]}_{n_3} = z^{[3]}_{n_3}
 \end{aligned} \tag{A.2}$$

The total composed output (2.14) of the artificial neural network, without Einstein Notation is:

$$u_{n_3} = b^{[3]}_{n_3} + \sum_{k_1=1}^4 w^{[3]k_1}_{n_3} \cdot a^{[2]}_{k_1} \left(b^{[2]}_{k_1} + \sum_{k_2=1}^3 w^{[2]k_2}_{k_1} \cdot a^{[1]}_{k_2} \left(b^{[1]}_{k_2} + \sum_{k_3=1}^2 w^{[1]k_3}_{k_2} \cdot x_{k_3} \right) \right) \quad (\text{A.3})$$

An example of the chain rule would be the following (where underneath we have specified the dimensions of the resulting matrix):

$$\frac{\partial u_{n_3}}{\partial x_{n'_0}} = \frac{\partial y^3_{n'_3}}{\partial z^0_{n'_0}} = \underbrace{\frac{\partial y^3_{n'_3}}{\partial z^3_{n_3}}}_{2 \times 2} \cdot \underbrace{\frac{\partial z^3_{n_3}}{\partial y^2_{n_2}}}_{2 \times 4} \cdot \underbrace{\frac{\partial y^2_{n_2}}{\partial z^2_{n_2}}}_{4 \times 4} \cdot \underbrace{\frac{\partial z^2_{n_2}}{\partial y^1_{n_1}}}_{4 \times 3} \cdot \underbrace{\frac{\partial y^1_{n_1}}{\partial z^1_{n_1}}}_{3 \times 3} \cdot \underbrace{\frac{\partial z^1_{n_1}}{\partial y^0_{n_0}}}_{3 \times 2} \cdot \underbrace{\frac{\partial y^0_{n_0}}{\partial z^0_{n'_0}}}_{2 \times 2} \quad (\text{A.4})$$

Finally, the partial derivatives used in the chain rules for (2.15-2.18), in their matrix form look like:

$$\frac{\partial y^{[3]}_{n'_3}}{\partial z^{[3]}_{n_3}} = \begin{pmatrix} \frac{\partial y^{[3]}_1}{\partial z^{[3]}_1} & \frac{\partial y^{[3]}_1}{\partial z^{[3]}_2} \\ \frac{\partial y^{[3]}_2}{\partial z^{[3]}_1} & \frac{\partial y^{[3]}_2}{\partial z^{[3]}_2} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (\text{A.5})$$

$$\frac{\partial y^{[2]}_{n_2}}{\partial z^{[2]}_{n_2}} = \begin{pmatrix} \frac{\partial y^{[2]}_1}{\partial z^{[2]}_1} & \frac{\partial y^{[2]}_1}{\partial z^{[2]}_2} & \frac{\partial y^{[2]}_1}{\partial z^{[2]}_3} & \frac{\partial y^{[2]}_1}{\partial z^{[2]}_4} \\ \frac{\partial y^{[2]}_2}{\partial z^{[2]}_1} & \frac{\partial y^{[2]}_2}{\partial z^{[2]}_2} & \frac{\partial y^{[2]}_2}{\partial z^{[2]}_3} & \frac{\partial y^{[2]}_2}{\partial z^{[2]}_4} \\ \frac{\partial y^{[2]}_3}{\partial z^{[2]}_1} & \frac{\partial y^{[2]}_3}{\partial z^{[2]}_2} & \frac{\partial y^{[2]}_3}{\partial z^{[2]}_3} & \frac{\partial y^{[2]}_3}{\partial z^{[2]}_4} \\ \frac{\partial y^{[2]}_4}{\partial z^{[2]}_1} & \frac{\partial y^{[2]}_4}{\partial z^{[2]}_2} & \frac{\partial y^{[2]}_4}{\partial z^{[2]}_3} & \frac{\partial y^{[2]}_4}{\partial z^{[2]}_4} \end{pmatrix} = \begin{pmatrix} \frac{\partial a^{[2]}_1(z^2_1)}{\partial z^2_1} & 0 & 0 & 0 \\ 0 & \frac{\partial a^{[2]}_2(z^2_2)}{\partial z^2_2} & 0 & 0 \\ 0 & 0 & \frac{\partial a^{[2]}_3(z^2_3)}{\partial z^2_3} & 0 \\ 0 & 0 & 0 & \frac{\partial a^{[2]}_4(z^2_4)}{\partial z^2_4} \end{pmatrix} \quad (\text{A.6})$$

$$\frac{\partial y^{[1]}_{n_1}}{\partial z^{[1]}_{n_1}} = \begin{pmatrix} \frac{\partial y^{[1]}_1}{\partial z^{[1]}_1} & \frac{\partial y^{[1]}_1}{\partial z^{[1]}_2} & \frac{\partial y^{[1]}_1}{\partial z^{[1]}_3} \\ \frac{\partial y^{[1]}_2}{\partial z^{[1]}_1} & \frac{\partial y^{[1]}_2}{\partial z^{[1]}_2} & \frac{\partial y^{[1]}_2}{\partial z^{[1]}_3} \\ \frac{\partial y^{[1]}_3}{\partial z^{[1]}_1} & \frac{\partial y^{[1]}_3}{\partial z^{[1]}_2} & \frac{\partial y^{[1]}_3}{\partial z^{[1]}_3} \end{pmatrix} = \begin{pmatrix} \frac{\partial a^{[1]}_1(z^1_1)}{\partial z^1_1} & 0 & 0 \\ 0 & \frac{\partial a^{[1]}_2(z^1_2)}{\partial z^1_2} & 0 \\ 0 & 0 & \frac{\partial a^{[1]}_3(z^1_3)}{\partial z^1_3} \end{pmatrix} \quad (\text{A.7})$$

$$\frac{\partial y^{[0]}_{n_0}}{\partial z^{[0]}_{n'_0}} = \begin{pmatrix} \frac{\partial y^{[0]}_1}{\partial z^{[0]}_1} & \frac{\partial y^{[0]}_1}{\partial z^{[0]}_2} \\ \frac{\partial y^{[0]}_2}{\partial z^{[0]}_1} & \frac{\partial y^{[0]}_2}{\partial z^{[0]}_2} \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \quad (\text{A.8})$$

The derivatives of $z^{[\ell]}_{n_\ell}$ with respect to $y^{[\ell]}_{n_\ell}$ are simply the weight matrices in (A.1).

Appendix B

The Code

As already introduced in section 3.1, the code has been implemented using Python's version TensorFlow 2.3. The code was implemented a Google Colab notebook, hence each class was encapsulated in a cell. Next, there is a brief simplified description on what each cell/class contains:

- **imports Cell:** Imports the main libraries, which includes TensorFlow for tensor manipulation, Time to get the time stamp, Pickle to save the models, Matplotlib to plot the models, among many others. It also suppresses Warnings.
- **auxiliryPlotting Class:** Encapsulates the methods for plotting results. It contains functions to: plot the distribution of dataset collocations points, plot the output of the model along the exact solution, plot the loss vs epoch graph of the training of the model, plot the errors of individual points in the training set, and plot the loss vs epoch of multiple models in the same graph.
- **myDataSets Class:** Used to create instances of myDataSets. Each of this instances mainly generate for different options, and contain, the collocation of points for the training and validation sets.
- **problemInstance Class:** Encapsulates the methods for the specifics of the diverse instances of the initial/boundary problems. It contains functions that given the training or validation set, and the artificial neural network output and derivatives, return the values of: the differential operator, the external force, the initial/boundary conditions lhs and rhs, and the exact solution of the problem.
- **secondOrderOptimizers Class:** Used to create instances of secondOrderOptimizers implementing the BFGS and L-BFGS optimizers. Since Keras only contains first order optimizers, this custom class uses the implementation in tensorflow_probability library, which is generic, and adapts it to input artificial neural network models.
- **myLayer Class:** Overrides the keras.Layer class and it is used to create object instances of myLayer. These objects contain the parameters and composition of the neurons in an artificial neuron layer, and the *feed* method which process an input to obtain the corresponding layer output.
- **myModel Class:** Overrides the keras.Model class to create instances of myModel, which implements the artificial neural network models. These objects are based on collections of myLayer instances, and contain either, a first order optimizer instance from Keras, or second order optimizer instance from secondOrderOptimizers, which can be accessed and changed at any moment. Through the methods in these objects and given a myDataSets instance one can: obtain the model output, or train the model for a problem set-up which calls on problemInstance for its specifics. Historical information about the loss performance during training is stored in the object. Also, there are methods to save and load models in *.pickle files, for later use.

- **execution Cell:** These are the snippets of code that calls on to the previous classes to perform the experiments. One of these calls usually consist on: a call to instance a myDataSets and myModel, with some options; a call to the fit method in myModel, to train the artificial neural network; a call to one of the auxiliryPlotting methods to plot the results; and optionally, saving the model. B.8 has an example showing in comments all of the variations that can be used.

There are more functionalities implemented throughout these classes. For more details, read the comments through the code. (The code font has been reduced to preserve indentation).

B.1 imports Cell

```

1 """
2 @author: Alberto Garcia Molina
3 @latest_update: 12/10/2020
4 """
5
6 import math
7 from math import log
8 import numpy as np
9
10 import time
11 import matplotlib.pyplot as plt
12 from pylab import rcParams
13 from mpl_toolkits.mplot3d import Axes3D
14
15 import pickle
16 from google.colab import files # Only for the colab environment.
17
18 import tensorflow as tf
19 import tensorflow_probability as tfp
20 from tensorflow import keras
21 from tensorflow.keras import layers
22
23 import logging, os
24
25 # Supress Warnings.
26
27 logging.disable(logging.WARNING)
28 os.environ["TF_CPP_MIN_LOG_LEVEL"] = "3"
29
30 # Optional code to check if there is a GPU available.
31
32 #%tensorflow_version 2.x
33 #device_name = tf.test.gpu_device_name()
34 #if device_name != '/device:GPU:0':
35 # raise SystemError('GPU device not found')
36 #print('Found GPU at: {}'.format(device_name))

```

B.2 auxiliryPlotting Class

```

1 """
2 @author: Alberto Garcia Molina
3 @latest_update: 12/10/2020
4 """
5
6 class auxiliryPlotting:
7
8     #####
9     # Plots the generated sets (Only 2D).
10    #####
11    def plot_datasets (data_set):
12
13        %matplotlib inline
14        training_set, border_training_set, validation_set = data_set.get_sets()
15        _, _, _, input_dim, output_dim = data_set.get_set_dimensions()
16
17        if (input_dim == 2):
18            plt.scatter(training_set[:,0], training_set[:,1], s=0.1)
19            plt.title('Training Set')
20            plt.show()
21
22            plt.scatter(border_training_set[0][:,0], border_training_set[0][:,1], s=0.1)
23            plt.title('Border Training Set')
24            plt.show()
25
26            plt.scatter(validation_set[:,0], validation_set[:,1], s=0.1)
27            plt.title('Validation Set')
28            plt.show()
29        else:
30            print('Invalid dimensions for plot.')

```

```

31 #####
32 # Plot loss (Only for training set).
33 #####
34 def plot_loss_function (model,
35                         init_range = 0,
36                         end_range = -1,
37                         subdivide_losses = False,
38                         use_log_scale = False):
39
40
41 %matplotlib inline
42
43 #Sets the x range of the plot.
44 plot_real_sol_loss = True
45 if (end_range < 0):
46     end_range = len(model._losses)
47
48 # Plots the real loss.
49 min_loss = min(model._losses[init_range:end_range])
50 max_loss = max(model._losses[init_range:end_range])
51
52 if (use_log_scale == True):
53     min_loss = min(loss for loss in model._losses[init_range:end_range] if loss > 0)
54     plt.plot(range(init_range, end_range),
55              [log(y,10) if y !=0 else None
56               for y in model._losses[init_range:end_range]],
57              label='Loss')
58     plt.ylim(log(min_loss,10), log(max_loss,10))
59     plt.title('Loss (log) - Epoch')
60 else:
61     plt.plot(range(init_range, end_range),
62              model._losses[init_range:end_range],
63              label='Loss')
64     plt.ylim(min_loss, max_loss)
65     plt.title('Loss - Epoch')
66
67 plt.xlim(init_range, end_range)
68 plt.legend()
69 plt.show()
70 print('Minimum Loss at:', str(min_loss))
71
72 # Plots the loss wrt the real solution.
73 min_loss_wrt_solution = min(model._losses_solution[init_range:end_range])
74 max_loss_wrt_solution = max(model._losses_solution[init_range:end_range])
75
76 if (use_log_scale == True):
77     min_loss_wrt_solution = min(loss for loss in model._losses_solution[init_range:end_range] if loss > 0)
78     plt.plot(range(init_range, end_range),
79              [log(y,10) if y !=0 else None
80               for y in model._losses_solution[init_range:end_range]],
81              label='Loss wrt Solution')
82     plt.ylim(log(min_loss_wrt_solution,10), log(max_loss_wrt_solution,10))
83     plt.title('Loss wrt Exact Sol (log) - Epoch')
84 else:
85     plt.plot(range(init_range, end_range),
86              model._losses_solution[init_range:end_range],
87              label='Loss')
88     plt.ylim(min_loss_wrt_solution, max_loss_wrt_solution)
89     plt.title('Loss wrt Exact Sol - Epoch')
90 plt.xlim(init_range, end_range)
91 plt.legend()
92 plt.show()
93 print('Minimum Loss wrt Solution at:', str(min_loss_wrt_solution))
94
95 # Plots the subdivision of the loss by its components.
96 if (subdivide_losses == True):
97
98     # Domain Component
99     min_domain_loss = min(model._losses_domain[init_range:end_range])
100    max_domain_loss = max(model._losses_domain[init_range:end_range])
101
102    if (use_log_scale == True):
103        min_domain_loss = min(loss for loss in model._losses_domain[init_range:end_range] if loss > 0)
104        plt.plot(range(init_range, end_range),
105                 [log(y,10) if y !=0 else None
106                  for y in model._losses_domain[init_range:end_range]],
107                 label='Loss')
108        plt.ylim(log(min_domain_loss,10), log(max_domain_loss,10))
109        plt.title('Domain Loss (log) - Epoch')
110    else:
111        plt.plot(range(init_range, end_range),
112                 model._losses_domain[init_range:end_range],
113                 label='Loss')
114        plt.ylim(min_domain_loss, max_domain_loss)
115        plt.title('Domain Loss - Epoch')
116
117    plt.xlim(init_range, end_range)
118    plt.legend()
119    plt.show()
120    print('Minimum Domain Loss at:', str(min_domain_loss))
121
122    # Border Component
123    min_border_loss = min(model._losses_border[init_range:end_range])
124    max_border_loss = max(model._losses_border[init_range:end_range])
125
126    if (use_log_scale == True):
127        min_border_loss = min(loss for loss in model._losses_border[init_range:end_range] if loss > 0)
128        plt.plot(range(init_range, end_range),
129                 [log(y,10) if y !=0 else None
130                  for y in model._losses_border[init_range:end_range]],
131                 label='Loss')
132        plt.ylim(log(min_border_loss,10), log(max_border_loss,10))

```

```

133     plt.title('Border Loss (log) - Epoch')
134     else:
135         plt.plot(range(init_range, end_range),
136                  model._losses_border[init_range:end_range],
137                  label='Loss')
138         plt.ylim(min_border_loss, max_border_loss)
139         plt.title('Border Loss - Epoch')
140
141     plt.xlim(init_range, end_range)
142     plt.legend()
143     plt.show()
144     print('Minimum Border Loss at:', str(min_border_loss))
145
146     # Regularization Component
147     if (model._regularization != None):
148         min_reg_loss = min(model._losses_regularization[init_range:end_range])
149         max_reg_loss = max(model._losses_regularization[init_range:end_range])
150
151         if (use_log_scale == True):
152             min_reg_loss = min(loss for loss in model._losses_regularization[init_range:end_range] if loss > 0)
153             plt.plot(range(init_range, end_range),
154                      [log(y,10) if y !=0 else None
155                       for y in model._losses_regularization[init_range:end_range]],
156                      label='Loss')
157             plt.ylim(log(min_reg_loss,10), log(max_reg_loss,10))
158             plt.title('Regularization Loss (log) - Epoch')
159         else:
160             plt.plot(range(init_range, end_range),
161                      model._losses_regularization[init_range:end_range],
162                      label='Loss')
163             plt.ylim(min_reg_loss, max_reg_loss)
164             plt.title('Regularization Loss - Epoch')
165         plt.xlim(init_range, end_range)
166         plt.legend()
167         plt.show()
168         print('Minimum Regularization Loss at:', str(min_border_loss))
169
170     #####
171     # Plots the model (Uses the validation set).
172     #####
173     def plot_model (data_set,
174                    model,
175                    plot_real_sol = False):
176
177         if (model._input_dim == 1):
178             if (model._output_dim == 1):
179                 outputs, _ = model.predict(data_set._validation_set)
180                 exact_sol = problemInstance.exact_solution(inputs = data_set._validation_set,
181                                                           exact_solution = model._exact_solution,
182                                                           input_dim = 1,
183                                                           output_dim = 1)
184
185                 plt.scatter(data_set._validation_set, outputs, s=0.1, label='Model')
186                 plt.scatter(data_set._validation_set, exact_sol, s=0.1, label='Exact Solution')
187                 plt.xlabel('x')
188                 plt.ylabel(r'$\hat{u}(x)$')
189                 plt.legend()
190                 plt.title('Model')
191             else:
192                 print('Invalid dimensions for plot.')
193
194         elif (model._input_dim == 2):
195             x = data_set._validation_set[:,0]
196             y = data_set._validation_set[:,1]
197
198             if (model._output_dim == 1):
199                 outputs, _ = model.predict(data_set._validation_set)
200                 exact_sol = problemInstance.exact_solution(inputs = data_set._validation_set,
201                                                           exact_solution = model._exact_solution,
202                                                           input_dim = 2,
203                                                           output_dim = 1)
204
205                 # Modificar para los limites reales.
206                 plt.rcParams['figure.figsize'] = [8,8]
207                 fig = plt.figure()
208                 ax = plt.axes(projection='3d')
209                 ax.set_title('Model')
210                 ax.set_xlim(tf.math.reduce_min(x), tf.math.reduce_max(x))
211                 ax.set_ylim(tf.math.reduce_min(y), tf.math.reduce_max(y))
212                 ax.set_xlabel('x')
213                 ax.set_ylabel('y')
214                 ax.set_zlabel('u (x,y)')
215                 ax.scatter3D(x, y, outputs, cmap='Greens', s=1, label='Model')
216                 fig = plt.figure()
217                 ax.scatter3D(x, y, exact_sol, cmap='Greens', s=1, label='Exact Solution')
218                 fig = plt.figure()
219                 ax.set_title('Model vs Exact Solution')
220                 ax.legend()
221                 #Optional
222                 ax.view_init(20, 230)
223
224             elif (model._output_dim == 2):
225                 outputs, _ = model.predict(data_set._validation_set)
226                 exact_sol = problemInstance.exact_solution(inputs = data_set._validation_set,
227                                                           exact_solution = model._exact_solution,
228                                                           input_dim = 2,
229                                                           output_dim = 2)
230
231                 # Modificar para los limites reales.
232                 plt.rcParams['figure.figsize'] = [7,7]
233                 fig = plt.figure()
234                 ax = plt.axes(projection='3d')
235                 ax.set_title('Model')
236                 ax.set_xlim(tf.math.reduce_min(x), tf.math.reduce_max(x))

```

```

235     ax.set_ylim(tf.math.reduce_min(y), tf.math.reduce_max(y))
236     ax.set_xlabel('x')
237     ax.set_ylabel('y')
238     ax.set_zlabel('u_x(x,y)')
239     ax.scatter3D(x, y, outputs[:,0], cmap='Greens', s=0.2)
240     fig = plt.figure()
241
242     fig = plt.figure()
243     ax = plt.axes(projection='3d')
244     ax.set_title('Exact Solution')
245     ax.set_xlim(tf.math.reduce_min(x), tf.math.reduce_max(x))
246     ax.set_ylim(tf.math.reduce_min(y), tf.math.reduce_max(y))
247     ax.set_xlabel('x')
248     ax.set_ylabel('y')
249     ax.set_zlabel('u_x(x,y)')
250     ax.scatter3D(x, y, exact_sol[:,0], cmap='Greens', s=0.2)
251     fig = plt.figure()
252     plt.legend()
253
254     # Modificar para los limites reales.
255     plt.rcParams['figure.figsize'] = [7,7]
256     fig = plt.figure()
257     ax = plt.axes(projection='3d')
258     ax.set_title('Model')
259     ax.set_xlim(tf.math.reduce_min(x), tf.math.reduce_max(x))
260     ax.set_ylim(tf.math.reduce_min(y), tf.math.reduce_max(y))
261     ax.set_xlabel('x')
262     ax.set_ylabel('y')
263     ax.set_zlabel('u_y(x,y)')
264     ax.scatter3D(x, y, outputs[:,1], cmap='Greens', s=0.2)
265     fig = plt.figure()
266
267     fig = plt.figure()
268     ax = plt.axes(projection='3d')
269     ax.set_title('Exact Solution')
270     ax.set_xlim(tf.math.reduce_min(x), tf.math.reduce_max(x))
271     ax.set_ylim(tf.math.reduce_min(y), tf.math.reduce_max(y))
272     ax.set_xlabel('x')
273     ax.set_ylabel('y')
274     ax.set_zlabel('u_y(x,y)')
275     ax.scatter3D(x, y, exact_sol[:,1], cmap='Greens', s=0.2)
276     fig = plt.figure()
277     plt.legend()
278
279     else:
280         print('Invalid dimensions for plot.')
281
282     #####
283     # Plots the squared error (Prototype).
284     #####
285     def plot_error (data_set,
286                    model):
287
288         if (model._input_dim == 1):
289             if (model._output_dim == 1):
290                 outputs, _ = model.predict(data_set._validation_set)
291                 exact_sol = problemInstance.exact_solution(inputs = data_set._validation_set,
292                                                           exact_solution = model._exact_solution,
293                                                           input_dim = 1,
294                                                           output_dim = 1)
295
296                 plt.scatter(data_set._validation_set, tf.square(outputs-exact_sol), s=0.1, label='Square Error')
297                 plt.legend()
298                 plt.title('Model')
299             else:
300                 print('Invalid dimensions for plot.')
301
302         elif (model._input_dim == 2):
303             x = data_set._validation_set[:,0]
304             y = data_set._validation_set[:,1]
305
306             if (model._output_dim == 1):
307                 # Loss wrt to the operator and force
308                 domain_ind_loss = tf.reduce_sum(
309                     tf.square(
310                         problemInstance.differential_operator(
311                             inputs = data_set._training_set,
312                             outputs = model.predict(data_set._training_set,
313                                                       model._required_derivative_order)[0],
314                             outputs_derivatives = model.predict(data_set._training_set,
315                                                                    model._required_derivative_order)[1],
316                             differential_operator = model._differential_operator,
317                             input_dim = 2,
318                             output_dim = 1)
319                     - problemInstance.external_force(inputs = data_set._training_set,
320                                                       external_force = model._external_force,
321                                                       input_dim = 2,
322                                                       output_dim = 1)),
323                     axis = 1,
324                     keepdims = True)
325
326                 #border_ind_loss = tf.reduce_sum(
327                     tf.square(data_set._border_training_set[1]
328                             - model.predict(data_set._border_training_set[0])[0],
329                     model._required_derivative_order),
330                     axis = 1,
331                     keepdims = True)
332
333                 # Modificar para los limites reales.
334                 plt.rcParams['figure.figsize'] = [7,7]
335                 fig = plt.figure()
336                 ax = plt.axes(projection='3d')
337                 ax.set_xlim(0, 1)

```

```

337     ax.set_ylim(0, 1)
338     ax.set_xlabel('x')
339     ax.set_ylabel('y')
340     ax.set_zlabel(r'$L_{1}(x;w,b)$')
341     ax.scatter3D(tf.concat([data_set._training_set[:,0], data_set._border_training_set[0][:,0]], axis=0),
342                 #tf.concat([data_set._training_set[:,1], data_set._border_training_set[0][:,1]], axis=0),
343                 #tf.square(tf.concat([domain_ind_loss, border_ind_loss], axis=0)),
344                 tf.concat([data_set._training_set[:,0]], axis=0),
345                 tf.concat([data_set._training_set[:,1]], axis=0),
346                 tf.square(tf.concat([domain_ind_loss], axis=0)),
347                 cmap='Greens',
348                 s=0.2)
349     fig = plt.figure()
350     ax.set_title(r'MSE for the individual domain points:  $\|\hat{u}(x,y) - f(x,y)\|_{2}^2$ ')
351     # Optional
352     ax.view_init(30, 40)
353
354     # Loss wrt to the real sol.
355     #input_set = tf.concat([data_set._training_set[:,], data_set._border_training_set[0][:,]], axis=0)
356     input_set = tf.concat([data_set._training_set[:,], axis=0)
357     outputs, _ = model.predict(input_set)
358     exact_sol = problemInstance.exact_solution(inputs = input_set,
359                                             exact_solution = model._exact_solution,
360                                             input_dim = 2,
361                                             output_dim = 1)
362
363     # Modificar para los limites reales.
364     plt.rcParams['figure.figsize'] = [7,7]
365     fig = plt.figure()
366     ax = plt.axes(projection='3d')
367     ax.set_title('Model')
368     ax.set_xlim(0, 1)
369     ax.set_ylim(0, 1)
370     ax.set_xlabel('x')
371     ax.set_ylabel('y')
372     ax.set_zlabel('Square Error')
373     ax.scatter3D(input_set[:,0],
374                 input_set[:,1],
375                 tf.square(outputs-exact_sol),
376                 cmap='Greens',
377                 s=0.2)
378     fig = plt.figure()
379     ax.set_title('Square Error of the Real Sol')
380
381     if (model._output_dim == 2):
382         # Loss wrt to the operator and force
383         domain_ind_loss = tf.reduce_sum(
384             tf.square(
385                 problemInstance.differential_operator(
386                     inputs = data_set._training_set,
387                     outputs = model.predict(data_set._training_set,
388                                             model._required_derivative_order)[0],
389                     outputs_derivatives = model.predict(data_set._training_set,
390                                                         model._required_derivative_order)[1],
391                     differential_operator = model._differential_operator,
392                     input_dim = 2,
393                     output_dim = 2)
394                 - problemInstance.external_force(inputs = data_set._training_set,
395                                                  external_force = model._external_force,
396                                                  input_dim = 2,
397                                                  output_dim = 2)),
398             axis = 1,
399             keepdims = True)
400         border_ind_loss = tf.reduce_sum(
401             tf.square(data_set._border_training_set[1]
402                     - model.predict(data_set._border_training_set[0][0],
403                                     model._required_derivative_order),
404             axis = 1,
405             keepdims = True)
406
407         # Modificar para los limites reales.
408         plt.rcParams['figure.figsize'] = [7,7]
409         fig = plt.figure()
410         ax = plt.axes(projection='3d')
411         ax.set_title('Model')
412         ax.set_xlim(0, 1)
413         ax.set_ylim(0, 1)
414         ax.set_xlabel('x')
415         ax.set_ylabel('y')
416         ax.set_zlabel('Square Error')
417         ax.scatter3D(tf.concat([data_set._training_set[:,0], data_set._border_training_set[0][:,0]], axis=0),
418                     tf.concat([data_set._training_set[:,1], data_set._border_training_set[0][:,1]], axis=0),
419                     tf.square(tf.concat([domain_ind_loss, border_ind_loss], axis=0)),
420                     cmap='Greens',
421                     s=0.2)
422         fig = plt.figure()
423         ax.set_title('Square Error of the Loss Formula')
424
425     else:
426         print('Invalid dimensions for plot.')
427
428     def plot_loss_comparison (models,
429                             names,
430                             title,
431                             init_range = 0,
432                             end_range = -1,
433                             subdivide_losses = False,
434                             use_log_scale = False):
435
436         %matplotlib inline
437         #rcParams['figure.figsize'] = 15, 5
438         rcParams['figure.figsize'] = 20, 4

```

```

439
440 #Sets the x range of the plot.
441 if (end_range < 0):
442     end_range = 0
443     for model in models:
444         end_range_var = len(model._losses)
445         if (end_range < end_range_var):
446             end_range = end_range_var
447
448 # Plots the global loss.
449 min_loss = 1e30
450 max_loss = 0
451 for model in models:
452     min_loss_var = min(loss for loss in model._losses[init_range:end_range] if loss > 0)
453     max_loss_var = max(model._losses[init_range:end_range])
454     if (min_loss_var < min_loss):
455         min_loss = min_loss_var
456     if (max_loss_var > max_loss):
457         max_loss = max_loss_var
458
459 if (use_log_scale == True):
460     for model_ind in range(len(models)):
461         plt.plot(range(init_range, len(models[model_ind]._losses)),
462                 [log(y,10) if y !=0 else None
463                  for y in models[model_ind]._losses[init_range:len(models[model_ind]._losses)]],
464                 label = names[model_ind])
465     plt.ylim(log(min_loss,10), log(max_loss,10))
466     plt.title(r'Global Loss Logarithm,  $\log_{10}(L(w,b))$  vs Epoch - ' + title)
467 else:
468     for model_ind in range(len(models)):
469         plt.plot(range(init_range, len(models[model_ind]._losses)),
470                 models[model_ind]._losses[init_range:len(models[model_ind]._losses)],
471                 label = names[model_ind])
472     plt.ylim(min_loss, max_loss)
473     plt.title(r'Global Loss,  $L(w,b)$  vs Epoch - ' + title)
474
475
476 plt.xlim(init_range, end_range)
477 plt.xlabel('Iterations')
478 plt.ylabel('Loss')
479 plt.legend()
480 plt.show()
481
482 # Plots the loss wrt the real solution.
483 min_loss = 1e30
484 max_loss = 0
485 for model in models:
486     min_loss_var = min(loss for loss in model._losses_solution[init_range:end_range] if loss > 0)
487     max_loss_var = max(model._losses_solution[init_range:end_range])
488     if (min_loss_var < min_loss):
489         min_loss = min_loss_var
490     if (max_loss_var > max_loss):
491         max_loss = max_loss_var
492
493 if (use_log_scale == True):
494     for model_ind in range(len(models)):
495         plt.plot(range(init_range, len(models[model_ind]._losses_solution)),
496                 [log(y,10) if y !=0 else None
497                  for y in models[model_ind]._losses_solution[init_range:len(models[model_ind]._losses_solution)
498                  ]],
499                 label = names[model_ind])
500     plt.ylim(log(min_loss,10), log(max_loss,10))
501     plt.title(r'Exact Solution Loss Logarithm,  $\log_{10}(L_{\{sol\}}(w,b))$  vs Epoch - ' + title)
502 else:
503     for model_ind in range(len(models)):
504         plt.plot(range(init_range, len(models[model_ind]._losses_solution)),
505                 models[model_ind]._losses_solution[init_range:len(models[model_ind]._losses_solution)],
506                 label = names[model_ind])
507     plt.ylim(min_loss, max_loss)
508     plt.title(r'Exact Solution Loss,  $L_{\{sol\}}(w,b)$  vs Epoch - ' + title)
509
510 plt.xlim(init_range, end_range)
511 plt.xlabel('Iterations')
512 plt.ylabel('Loss')
513 plt.legend()
514 plt.show()
515
516 # Plots the border loss.
517 min_loss = 1e30
518 max_loss = 0
519 for model in models:
520     min_loss_var = min(loss for loss in model._losses_border[init_range:end_range] if loss > 0)
521     max_loss_var = max(model._losses_border[init_range:end_range])
522     if (min_loss_var < min_loss):
523         min_loss = min_loss_var
524     if (max_loss_var > max_loss):
525         max_loss = max_loss_var
526
527 if (use_log_scale == True):
528     for model_ind in range(len(models)):
529         plt.plot(range(init_range, len(models[model_ind]._losses_border)),
530                 [log(y,10) if y !=0 else None
531                  for y in models[model_ind]._losses_border[init_range:len(models[model_ind]._losses_border)]],
532                 label = names[model_ind])
533     plt.ylim(log(min_loss,10), log(max_loss,10))
534     plt.title(r'Initial Condition Loss Logarithm,  $\log_{10}(L_{\{2\}}(w,b))$  vs Epoch - ' + title)
535 else:
536     for model_ind in range(len(models)):
537         plt.plot(range(init_range, len(models[model_ind]._losses_border)),
538                 models[model_ind]._losses_border[init_range:len(models[model_ind]._losses_border)],
539                 label = names[model_ind])
540     plt.ylim(min_loss, max_loss)

```

```

540         plt.title(r'Initial Condition Loss,  $L_2(w,b)$  vs Epoch - ' + title)
541
542     plt.xlim(init_range, end_range)
543     plt.xlabel('Iterations')
544     plt.ylabel('Loss')
545     plt.legend()
546     plt.show()
547
548     # Plots the domain loss.
549     min_loss = 1e30
550     max_loss = 0
551     for model in models:
552         min_loss_var = min(loss for loss in model._losses_domain[init_range:end_range] if loss > 0)
553         max_loss_var = max(model._losses_domain[init_range:end_range])
554         if (min_loss_var < min_loss):
555             min_loss = min_loss_var
556         if (max_loss_var > max_loss):
557             max_loss = max_loss_var
558
559     if (use_log_scale == True):
560         for model_ind in range(len(models)):
561             plt.plot(range(init_range, len(models[model_ind]._losses_domain)),
562                     [log(y,10) if y !=0 else None
563                      for y in models[model_ind]._losses_domain[init_range:len(models[model_ind]._losses_domain)]],
564                     label = names[model_ind])
565     plt.ylim(log(min_loss,10), log(max_loss,10))
566     plt.title(r'Domian Loss Logaritm,  $\log_{10}(L_1(w,b))$  vs Epoch - ' + title)
567 else:
568     for model_ind in range(len(models)):
569         plt.plot(range(init_range, len(models[model_ind]._losses_domain)),
570                 models[model_ind]._losses_domain[init_range:len(models[model_ind]._losses_domain)],
571                 label = names[model_ind])
572     plt.ylim(min_loss, max_loss)
573     plt.title(r'Domian Loss,  $L_1(w,b)$  vs Epoch - ' + title)
574
575     plt.xlim(init_range, end_range)
576     plt.xlabel('Iterations')
577     plt.ylabel('Loss')
578     plt.legend()
579     plt.show()

```

B.3 myDataSets Class

```

1  """
2  @author: Alberto Garcia Molina
3  @latest_update: 12/10/2020
4  """
5
6  class myDataSets:
7
8      # Initialize myDataSets object.
9      def __init__(self,
10                 training_batch_size = 2000,
11                 border_training_batch_size = 20,
12                 validation_batch_size = 1000,
13                 input_dim = 1,
14                 method = 'uniform-hit-collocation',
15                 domain = 'hypercube-0-1',
16                 border = 'side-x_1-y_0',
17                 seed = None):
18
19         self._training_batch_size = training_batch_size
20         self._border_training_batch_size = border_training_batch_size
21         self._validation_batch_size = validation_batch_size
22         self._input_dim = input_dim
23
24         self.method = method
25         self.domain = domain
26         self.border = border
27
28         seed_1 = None
29         seed_2 = None
30         seed_3 = None
31         if (seed != None):
32             seed_1 = seed
33             seed_2 = 2*seed
34             seed_3 = 3*seed
35
36         self._training_set = self.generate_domain_set (training_batch_size, input_dim,
37                                                       method, domain, seed_1)
38
39         self._border_training_set = self.generate_border_set (border_training_batch_size,
40                                                             input_dim, method, border, seed_2)
41
42         self._validation_set = self.generate_domain_set(validation_batch_size, input_dim,
43                                                         method, domain, seed_3)
44
45     # Generates a distribution of points inside the solving domain.
46     def generate_domain_set (self,
47                             batch_size = 2000,
48                             input_dim = 1,
49                             method = 'uniform-hit-collocation',
50                             domain = 'hypercube-0-1',
51                             seed = None):
52
53         if (seed != None):
54             tf.random.set_seed(seed)

```

```

55
56     if (method == 'uniform-hit-collocation'):
57         if (domain == 'hypercube-0-1'):
58             domain_set = tf.random.uniform(shape=[batch_size, input_dim],
59                                           minval=0., maxval=1., dtype=tf.float32)
60         elif (domain == 'quarter-hypercube-0-1'):
61             domain_set = tf.random.uniform(shape=[batch_size, input_dim],
62                                           minval=0., maxval=0.5, dtype=tf.float32)
63         elif (domain == 'hypercube-0-2'):
64             domain_set = tf.random.uniform(shape=[batch_size, input_dim],
65                                           minval=0., maxval=2., dtype=tf.float32)
66
67     return domain_set
68
69 # Generates a distribution of points on the border of the solving domain.
70 def generate_border_set (self,
71                        batch_size = 2,
72                        input_dim = 1,
73                        method = 'uniform-hit-collocation',
74                        border = 'hypercube-0-1',
75                        seed = None):
76
77     if (seed != None):
78         tf.random.set_seed(seed)
79
80     if (method == 'uniform-hit-collocation'):
81         if (border == 'hypercube-0-1'):
82             if (input_dim == 1):
83                 x1 = tf.constant(0., shape=[1, input_dim], dtype=tf.float32)
84                 x2 = tf.constant(1., shape=[1, input_dim], dtype=tf.float32)
85                 border_set = tf.concat([x1, x2], axis=0)
86
87             elif (input_dim == 2):
88                 x1 = tf.random.uniform(shape=[batch_size//4],
89                                       minval=0.,
90                                       maxval=1.,
91                                       dtype=tf.float32)
92                 y1 = tf.constant(0.,
93                                 shape=[batch_size//4],
94                                 dtype=tf.float32)
95                 border_set_1 = tf.stack([x1, y1], axis=1) # y=0
96
97                 x2 = tf.random.uniform(shape=[batch_size//4],
98                                       minval=0.,
99                                       maxval=1.,
100                                      dtype=tf.float32)
101                 y2 = tf.constant(1.,
102                                 shape=[batch_size//4],
103                                 dtype=tf.float32)
104                 border_set_2 = tf.stack([x2, y2], axis=1) # y=1
105
106                 x3 = tf.constant(0.,
107                                 shape=[batch_size//4],
108                                 dtype=tf.float32)
109                 y3 = tf.random.uniform(shape=[batch_size//4],
110                                       minval=0.,
111                                       maxval=1.,
112                                       dtype=tf.float32)
113                 border_set_3 = tf.stack([x3, y3], axis=1) # x=0
114
115                 x4 = tf.constant(1.,
116                                 shape=[batch_size//4],
117                                 dtype=tf.float32)
118                 y4 = tf.random.uniform(shape=[batch_size//4],
119                                       minval=0.,
120                                       maxval=1.,
121                                       dtype=tf.float32)
122                 border_set_4 = tf.stack([x4, y4], axis=1) # x=1
123
124                 border_set = tf.concat([border_set_1, border_set_2, border_set_3, border_set_4],
125                                       axis=0)
126
127             elif (border == 'side-x_1-y_0'):
128                 if (input_dim == 1):
129                     border_set = tf.constant(0., shape=[1, input_dim], dtype=tf.float32)
130                 elif (input_dim == 2):
131                     x1 = tf.random.uniform(shape=[batch_size],
132                                           minval=-1.,
133                                           maxval=2.,
134                                           dtype=tf.float32)
135                     y1 = tf.constant(0.,
136                                     shape=[batch_size],
137                                     dtype=tf.float32)
138                     border_set = tf.stack([x1, y1], axis=1)
139
140             elif (border == 'side-x_1-y_0_expanded'):
141                 if (input_dim == 1):
142                     border_set = tf.constant(0., shape=[1, input_dim], dtype=tf.float32)
143                 elif (input_dim == 2):
144                     x1 = tf.random.uniform(shape=[batch_size],
145                                           minval=-1.,
146                                           maxval=2.,
147                                           dtype=tf.float32)
148                     y1 = tf.constant(0.,
149                                     shape=[batch_size],
150                                     dtype=tf.float32)
151                     border_set = tf.stack([x1, y1], axis=1)
152
153             elif (border == 'two_sides-x_0-y_0'):
154                 if (input_dim == 1):
155                     border_set = tf.constant(0., shape=[1, input_dim], dtype=tf.float32)
156                 elif (input_dim == 2):

```



```

157         x1 = tf.random.uniform(shape=[batch_size//2],
158                               minval=0.,
159                               maxval=1.,
160                               dtype=tf.float32)
161     y1 = tf.constant(0.,
162                    shape=[batch_size//2],
163                    dtype=tf.float32)
164
165     x2 = tf.constant(0.,
166                    shape=[batch_size//2],
167                    dtype=tf.float32)
168     y2 = tf.random.uniform(shape=[batch_size//2],
169                            minval=0.,
170                            maxval=1.,
171                            dtype=tf.float32)
172
173     border_set_1 = tf.stack([x1, y1], axis=1) # y=0
174     border_set_2 = tf.stack([x2, y2], axis=1) # x=0
175     border_set = tf.concat([border_set_1, border_set_2],
176                           axis=0)
177
178     return border_set
179
180 # Returns the sets stored in this object.
181 def get_sets(self):
182     return self._training_set, self._border_training_set, self._validation_set
183
184 # Returns the metadata of the sets stored in this object.
185 def get_set_metadata(self):
186     return self._training_batch_size, self._border_training_batch_size, self._validation_batch_size, \
187            self._input_dim, self.method, self.domain, self.border
188
189 # Drops the values which have negative loss.
190 def drop_negative_loss (data_set,
191                       model):
192
193     # TBD
194     domain_ind_dif = tf.reduce_sum(
195         problemInstance.differential_operator(
196             inputs = data_set._training_set,
197             outputs = model.predict(data_set._training_set,
198                                   model._required_derivative_order)[0],
199             outputs_derivatives = model.predict(data_set._training_set,
200                                                model._required_derivative_order)[1],
201             differential_operator = model._differential_operator,
202             input_dim = model._input_dim,
203             output_dim = model._output_dim)
204         - problemInstance.external_force(inputs = data_set._training_set,
205                                         external_force = model._external_force,
206                                         input_dim = model._input_dim,
207                                         output_dim = model._output_dim),
208         axis = 1,
209         keepdims = False)
210
211     # TBD
212     border_ind_dif = tf.reduce_sum(
213         data_set._border_training_set[1]
214         - model.predict(data_set._border_training_set[0], 0)[0],
215         axis = 1,
216         keepdims = False)
217
218     # Mask
219     filtered_training_set = tf.boolean_mask(tensor = data_set._training_set,
220                                           mask = domain_ind_dif > 0,
221                                           axis = 0)
222     filtered_border_training_set_0 = tf.boolean_mask(tensor = data_set._border_training_set[0],
223                                                    mask = border_ind_dif > 0,
224                                                    axis = 0)
225     filtered_border_training_set_1 = tf.boolean_mask(tensor = data_set._border_training_set[1],
226                                                    mask = border_ind_dif > 0,
227                                                    axis = 0)
228
229     # Replace the Dataset
230     if (filtered_training_set.shape[0] != 0):
231         data_set._training_set = filtered_training_set
232         data_set._training_batch_size = filtered_training_set.shape[0]
233     if (filtered_border_training_set_0.shape[0] != 0):
234         data_set._border_training_set[0] = filtered_border_training_set_0
235         data_set._border_training_set[1] = filtered_border_training_set_1
236         data_set._border_training_batch_size = filtered_border_training_set_0.shape[0]
237
238 # Drops the values which have negative loss.
239 def drop_best_loss (data_set,
240                   model):
241
242     # TBD
243     domain_ind_dif = tf.reduce_sum(tf.square(
244         problemInstance.differential_operator(
245             inputs = data_set._training_set,
246             outputs = model.predict(data_set._training_set,
247                                   model._required_derivative_order)[0],
248             outputs_derivatives = model.predict(data_set._training_set,
249                                                model._required_derivative_order)[1],
250             differential_operator = model._differential_operator,
251             input_dim = model._input_dim,
252             output_dim = model._output_dim)
253         - problemInstance.external_force(inputs = data_set._training_set,
254                                         external_force = model._external_force,
255                                         input_dim = model._input_dim,
256                                         output_dim = model._output_dim)),
257         axis = 1,
258         keepdims = False)

```

```

259
260 # TBD
261 border_ind_dif = tf.reduce_sum(tf.square(
262     data_set._border_training_set[1]
263     - model.predict(data_set._border_training_set[0],0)[0]),
264     axis = 1,
265     keepdims = False)
266
267 # TBD
268 domain_best = tf.math.reduce_mean(
269     domain_ind_dif)
270
271 border_best = tf.math.reduce_mean(
272     border_ind_dif)
273
274 # Mask
275 filtered_training_set = tf.boolean_mask(tensor = data_set._training_set,
276     mask = domain_ind_dif > domain_best,
277     axis = 0)
278 filtered_border_training_set_0 = tf.boolean_mask(tensor = data_set._border_training_set[0],
279     mask = border_ind_dif > border_ind_dif,
280     axis = 0)
281 filtered_border_training_set_1 = tf.boolean_mask(tensor = data_set._border_training_set[1],
282     mask = border_ind_dif > border_ind_dif,
283     axis = 0)
284
285 # Replace the Dataset
286 if (filtered_training_set.shape[0] != 0):
287     data_set._training_set = filtered_training_set
288     data_set._training_batch_size = filtered_training_set.shape[0]
289 if (filtered_border_training_set_0.shape[0] != 0):
290     data_set._border_training_set[0] = filtered_border_training_set_0
291     data_set._border_training_set[1] = filtered_border_training_set_1
292     data_set._border_training_batch_size = filtered_border_training_set_0.shape[0]

```

B.4 problemInstance Class

```

1 """
2 @author: Alberto Garcia Molina
3 @latest_update: 12/10/2020
4 """
5
6
7 class problemInstance():
8
9     #####
10    # Validates if the problem instance of a model is implemented.
11    #####
12    def instance_exists (differential_operator,
13        external_force,
14        exact_solution):
15
16        if (differential_operator == 'Constant'):
17            required_derivative_order = 0
18        elif (differential_operator == 'Divergence'):
19            required_derivative_order = 1
20        elif (differential_operator == 'Advection'):
21            required_derivative_order = 1
22        elif (differential_operator == 'Laplacian'):
23            required_derivative_order = 2
24        elif (differential_operator == 'Clairaut'):
25            required_derivative_order = 1
26        elif (differential_operator == 'Korteweg-deVries'):
27            required_derivative_order = 3
28        else:
29            raise Exception("Invalid differential operator option. Please recompile with a valid name.")
30
31        if (external_force == 'Force_Constant'
32            or external_force == 'Force_Divergence'
33            or external_force == 'Force_Laplacian'
34            or external_force == 'Force_Advection'
35            or external_force == 'Force_Clairaut'
36            or external_force == 'Force_Korteweg-deVries'):
37            pass
38        else:
39            raise Exception("Invalid external force option. Please recompile with a valid name.")
40
41        if (exact_solution == 'Sol_Polynomial_2Deg_1D_1D'
42            or exact_solution == 'Sol_Polynomial_2Deg_2D_1D'
43            or exact_solution == 'Sol_Polynomial_2Deg_2D_2D'
44            or exact_solution == 'Sol_Polynomial_4Deg_1D_1D'):
45            pass
46        else:
47            raise Exception("Invalid exact solution option. Please recompile with a valid name.")
48
49        return required_derivative_order
50
51    #####
52    # Computes the differential operator value.
53    #####
54    #@tf.function(experimental_relax_shapes = True)
55    def differential_operator (inputs,
56        outputs,
57        outputs_derivatives,
58        differential_operator,
59        input_dim,
60        output_dim):

```

```

61
62 #batch_size = inputs.shape[0]
63 # tf.print(batch_size)
64
65 if (differential_operator == 'Constant'):
66     if (input_dim == 1):
67         dif_operator_output = outputs
68     elif (input_dim == 2):
69         dif_operator_output = outputs
70     else:
71         raise Exception("Incompatible dimension (Diff Operator).")
72
73 if (differential_operator == 'Divergence'):
74     if (input_dim == 1):
75         dif_operator_output = outputs_derivatives[0]
76     elif (input_dim == 2):
77         dif_operator_output = tf.reduce_sum(outputs_derivatives[0], axis=1, keepdims=True)
78     else:
79         raise Exception("Incompatible dimension (Diff Operator).")
80
81 if (differential_operator == 'Laplacian'):
82     if (input_dim == 2):
83         dif_operator_output = tf.reduce_sum(outputs_derivatives[1],
84                                             axis=1,
85                                             keepdims=True)
86     else:
87         raise Exception("Incompatible dimension (Diff Operator).")
88
89 if (differential_operator == 'Advection'):
90     if (input_dim == 1):
91         dif_operator_output = outputs*outputs_derivatives[0]
92     elif (input_dim == 2):
93         dif_operator_output = tf.linalg.matvec(outputs_derivatives[0], outputs) # transpose_a = True
94         #outputs = tf.reshape(outputs,[batch_size,1,output_dim])
95         #outputs_derivatives[0] = tf.transpose(outputs_derivatives[0], perm=[0, 2, 1])
96         #dif_operator_output = tf.matmul(outputs, outputs_derivatives[0])
97         #dif_operator_output = tf.reshape(dif_operator_output,
98                                         # [batch_size, output_dim])
99     else:
100         raise Exception("Incompatible dimension (Diff Operator).")
101
102 if (differential_operator == 'Clairaut'):
103     if (input_dim == 2):
104         dif_operator_output = tf.reduce_sum(inputs*outputs_derivatives[0],
105                                             axis = 1,
106                                             keepdims = True)
107     else:
108         raise Exception("Incompatible dimension (Diff Operator).")
109
110 if (differential_operator == 'Korteweg-deVries'):
111     if (input_dim == 1):
112         dif_operator_output = outputs_derivatives[2]-outputs_derivatives[0]
113     else:
114         raise Exception("Incompatible dimension (Diff Operator).")
115
116 return dif_operator_output
117
118 #####
119 # Computes the external force value.
120 #####
121 #@tf.function(experimental_relax_shapes = True)
122 def external_force(inputs,
123                  external_force,
124                  input_dim,
125                  output_dim):
126
127     batch_size = inputs.shape[0]
128     # tf.print(batch_size)
129
130     if (external_force == 'Force_Constant'):
131         if (input_dim == 1 and output_dim == 1):
132             x = inputs
133             ext_force_output = x*(x-1)
134         else:
135             raise Exception("Incompatible dimension (External force).")
136
137     if (external_force == 'Force_Divergence'):
138         if (input_dim == 1 and output_dim == 1):
139             x = inputs
140             ext_force_output = 2*x-1
141         elif (input_dim == 2 and output_dim == 1):
142             x = inputs[:,0]
143             y = inputs[:,1]
144             ext_force = (2*x-1)*y*(y-1)+x*(x-1)*(2*y-1)
145             ext_force_output = tf.reshape(ext_force,
146                                         [batch_size, output_dim])
147         else:
148             raise Exception("Incompatible dimension (External force).")
149
150     if (external_force == 'Force_Laplacian'):
151         if (input_dim == 2 and output_dim == 1):
152             x = inputs[:,0]
153             y = inputs[:,1]
154             ext_force = 2*y*(y-1)+2*x*(x-1)
155             ext_force_output = tf.reshape(ext_force,
156                                         [batch_size, output_dim])
157         else:
158             raise Exception("Incompatible dimension (External force).")
159
160     if (external_force == 'Force_Advection'):
161         if (input_dim == 1 and output_dim == 1):
162             x = inputs

```

```

163     ext_force_output = x*(x-1)*(2*x-1)
164 elif (input_dim == 2 and output_dim == 2):
165     x = inputs[:,0]
166     y = inputs[:,1]
167     x_ext_force = x*(x-1)*y*(y-1)*((2*x-1)*y*(y-1)+x*(x-1)*(2*y-1))
168     y_ext_force = x*(x-1)*y*(y-1)*((2*x-1)*y*(y-1)+x*(x-1)*(2*y-1))
169     ext_force_output = tf.stack([x_ext_force, y_ext_force], axis=1)
170 else:
171     raise Exception("Incompatible dimension (External force).")
172
173 if (external_force == 'Force_Clairaut'):
174     if (input_dim == 2 and output_dim == 1):
175         x = inputs[:,0]
176         y = inputs[:,1]
177         ext_force = x*(2*x-1)*y*(y-1)+y*x*(x-1)*(2*y-1)
178         ext_force_output = tf.reshape(ext_force,
179                                     [batch_size, output_dim])
180     else:
181         raise Exception("Incompatible dimension (External force).")
182
183 if (external_force == 'Force_Korteweg-deVries'):
184     if (input_dim == 1 and output_dim == 1):
185         x = inputs
186         ext_force_output = -4*x*x*x+3*x*x+32*x-10
187     else:
188         raise Exception("Incompatible dimension (External force).")
189
190 return ext_force_output
191
192 #####
193 # Computes the rhs border functions of the boundary conditions.
194 # Returns list with [border_inputs, rhs_function_1, rhs_funtion_2, ...]
195 #####
196 def border_data_prep(border_inputs,
197                    border_type,
198                    border_batch_size,
199                    external_force,
200                    required_derivative_order,
201                    input_dim,
202                    output_dim):
203
204     border_data = None
205     if (required_derivative_order != 0):
206         border_data = [border_inputs]
207
208     if (required_derivative_order >= 1):
209
210         # Border is always zero in the perimeter.
211         if (input_dim == 1):
212             x1 = border_inputs[:,0]
213
214             if (output_dim == 1):
215                 g1 = 0.*x1
216                 border_set_outputs = tf.stack([g1], axis=1)
217
218             elif (output_dim == 2):
219                 g1 = 0.*x1
220                 g2 = 0.*x1
221                 border_set_outputs = tf.stack([g1, g2], axis=1)
222
223         elif (input_dim == 2):
224             x1 = border_inputs[:,0]
225             y1 = border_inputs[:,1]
226
227             if (output_dim == 1):
228                 g1 = 0.*x1*y1
229                 border_set_outputs = tf.stack([g1], axis=1)
230
231             elif (output_dim == 2):
232                 g1 = 0.*x1*y1
233                 g2 = 0.*x1*y1
234                 border_set_outputs = tf.stack([g1, g2], axis=1)
235
236         border_data.append(border_set_outputs)
237
238     if (required_derivative_order >= 2 or external_force == 'Force_Advection'):
239
240         if (external_force == 'Force_Laplacian' or external_force == 'Force_Advection'):
241             if (border_type == 'side-x_1-y_0'
242                 or border_type == 'side-x_1-y_0_expanded'):
243
244                 if (input_dim == 1):
245                     x1 = border_inputs[:,0]
246
247                     if (output_dim == 1):
248                         g1 = 2.*x1-1
249                         border_set_outputs = tf.stack([g1], axis=1)
250
251                     elif (input_dim == 2):
252                         x1 = border_inputs[:,0]
253                         y1 = border_inputs[:,1]
254
255                         if (output_dim == 1):
256                             g1 = (2.*y1-1)*x1*(x1-1)
257                             border_set_outputs = tf.stack([g1], axis=1)
258
259                 elif (border_type == 'two_sides-x_0-y_0'):
260                     if (input_dim == 1):
261                         x1 = border_inputs[:,0]
262
263                     if (output_dim == 1):
264                         g1 = 2.*x1-1

```

```

265         border_set_outputs = tf.stack([g1], axis=1)
266
267     elif (input_dim == 2):
268         sub_batch_size = border_batch_size//2
269         x1 = border_inputs[0*sub_batch_size:1*sub_batch_size,0]
270         y1 = border_inputs[0*sub_batch_size:1*sub_batch_size,1]
271         x2 = border_inputs[1*sub_batch_size:2*sub_batch_size,0]
272         y2 = border_inputs[1*sub_batch_size:2*sub_batch_size,1]
273
274         if (output_dim == 1):
275             g1 = x1*(x1-1)*(2.*y1-1) # First y=0
276             g2 = (2.*x2-1)*y2*(y2-1) # Second x=0
277             border_set_outputs = tf.stack([g1,g2], axis=1)
278
279     elif (border_type == 'hypercube-0-1'):
280         if (input_dim == 1):
281             x1 = border_inputs[:,0]
282
283             if (output_dim == 1):
284                 g1 = 2.*x1-1
285                 border_set_outputs = tf.stack([g1], axis=1)
286
287         elif (input_dim == 2):
288             sub_batch_size = border_batch_size//4
289             x1 = border_inputs[0*sub_batch_size:1*sub_batch_size,0]
290             y1 = border_inputs[0*sub_batch_size:1*sub_batch_size,1]
291             x2 = border_inputs[1*sub_batch_size:2*sub_batch_size,0]
292             y2 = border_inputs[1*sub_batch_size:2*sub_batch_size,1]
293             x3 = border_inputs[2*sub_batch_size:3*sub_batch_size,0]
294             y3 = border_inputs[2*sub_batch_size:3*sub_batch_size,1]
295             x4 = border_inputs[3*sub_batch_size:4*sub_batch_size,0]
296             y4 = border_inputs[3*sub_batch_size:4*sub_batch_size,1]
297
298             if (output_dim == 1):
299                 g1 = x1*(x1-1)*(2.*y1-1) # First y=0
300                 g2 = x2*(x2-1)*(2.*y2-1) # First y=1
301                 g3 = (2.*x3-1)*y3*(y3-1) # Second x=0
302                 g4 = (2.*x4-1)*y4*(y4-1) # Second x=1
303
304             border_set_outputs = tf.stack([g1,g2,g3,g4], axis=1)
305
306         border_data.append(border_set_outputs)
307
308     return border_data
309
310     #####
311     # Computes the boundary conditions (lhs).
312     #####
313     def lhs_boundary_conditions (inputs,
314                                 outputs,
315                                 outputs_derivatives,
316                                 border_type,
317                                 border_batch_size,
318                                 external_force,
319                                 required_derivative_order,
320                                 input_dim,
321                                 output_dim):
322
323         boundary_cond = []
324
325         if (required_derivative_order >= 1):
326             boundary_cond.append(outputs)
327
328         if (required_derivative_order >= 2 or external_force == 'Force_Advection'):
329             if (external_force == 'Force_Laplacian' or external_force == 'Force_Advection'):
330                 if (border_type == 'side-x_1-y_0'
331                     or border_type == 'side-x_1-y_0_expanded'):
332                     if (input_dim == 1):
333                         if (output_dim == 1):
334                             boundary_cond.append(outputs_derivatives[0])
335
336                     elif (input_dim == 2):
337                         if (output_dim == 1):
338                             boundary_cond_val = outputs_derivatives[0][:,1]
339                             boundary_cond_val_resized = tf.stack([boundary_cond_val], axis=1)
340                             boundary_cond.append(boundary_cond_val_resized)
341
342             elif (border_type == 'two_sides-x_0-y_0'):
343                 if (input_dim == 1):
344                     if (output_dim == 1):
345                         boundary_cond.append(outputs_derivatives[0])
346
347                 elif (input_dim == 2):
348                     if (output_dim == 1):
349                         sub_batch_size = border_batch_size//2
350                         boundary_cond_val_1 = outputs_derivatives[0][0*sub_batch_size:1*sub_batch_size,1] # First y=0
351                         boundary_cond_val_2 = outputs_derivatives[0][1*sub_batch_size:2*sub_batch_size,0] # Second x=0
352                         boundary_cond_val_comb = tf.stack([boundary_cond_val_1,boundary_cond_val_2], axis=1)
353                         boundary_cond.append(boundary_cond_val_comb)
354
355             elif (border_type == 'hypercube-0-1'):
356                 if (input_dim == 1):
357                     if (output_dim == 1):
358                         boundary_cond.append(outputs_derivatives[0])
359
360                 elif (input_dim == 2):
361                     if (output_dim == 1):
362                         sub_batch_size = border_batch_size//4
363                         boundary_cond_val_1 = outputs_derivatives[0][0*sub_batch_size:1*sub_batch_size,1] # y=0
364                         boundary_cond_val_2 = outputs_derivatives[0][1*sub_batch_size:2*sub_batch_size,1] # y=1
365                         boundary_cond_val_3 = outputs_derivatives[0][2*sub_batch_size:3*sub_batch_size,0] # x=0
366                         boundary_cond_val_4 = outputs_derivatives[0][3*sub_batch_size:4*sub_batch_size,0] # x=1

```

```

367         boundary_cond_val_comb = tf.stack([boundary_cond_val_1, boundary_cond_val_2, boundary_cond_val_3,
368         boundary_cond_val_4], axis=1)
369         boundary_cond.append(boundary_cond_val_comb)
370
371     return boundary_cond
372
373     #####
374     # Computes the exact solution value.
375     #####
376     #@tf.function(experimental_relax_shapes = True)
377     def exact_solution (inputs,
378         exact_solution,
379         input_dim,
380         output_dim):
381
382         batch_size = inputs.shape[0]
383         # tf.print(batch_size)
384
385         if (exact_solution == 'Sol_Polynomial_2Deg_1D_1D'):
386             if (input_dim == 1 and output_dim == 1):
387                 x = inputs
388                 exact_sol_output = x*(x-1)
389             else:
390                 raise Exception("Incompatible dimension (External force).")
391
392         if (exact_solution == 'Sol_Polynomial_4Deg_1D_1D'):
393             if (input_dim == 1 and output_dim == 1):
394                 x = inputs
395                 exact_sol_output = x*(x-1)*(x*x-4)
396             else:
397                 raise Exception("Incompatible dimension (External force).")
398
399         if (exact_solution == 'Sol_Polynomial_2Deg_2D_1D'):
400             if (input_dim == 2 and output_dim == 1):
401                 x = inputs[:,0]
402                 y = inputs[:,1]
403                 exact_sol = x*(x-1)*y*(y-1)
404                 exact_sol_output = tf.reshape(exact_sol,
405                 [batch_size, output_dim])
406             else:
407                 raise Exception("Incompatible dimension (External force).")
408
409         if (exact_solution == 'Sol_Polynomial_2Deg_2D_2D'):
410             if (input_dim == 2 and output_dim == 2):
411                 x = inputs[:,0]
412                 y = inputs[:,1]
413                 x_exact_sol = x*(x-1)*y*(y-1)
414                 y_exact_sol = x*(x-1)*y*(y-1)
415                 exact_sol_output = tf.stack([x_exact_sol, y_exact_sol],
416                 axis=1)
417             else:
418                 raise Exception("Incompatible dimension (External force).")
419
420     return exact_sol_output

```

B.5 secondOrderOptimizers Class

```

1  """
2  @author: Alberto Garcia Molina
3  @latest_update: 12/10/2020
4  """
5
6  class secondOrderOptimizers():
7
8      def __init__ (self,
9          name,
10         model):
11
12         self._name = name
13         self._model = model
14
15     def optimizer_train_data (self,
16         inputs,
17         border_data,
18         split_gradient = False,
19         display_gradient_norm = False,
20         normalize_gradient = False,
21         train_only_domain = False,
22         train_only_border = False):
23
24         self._inputs = inputs
25         self._border_data = border_data
26         self._split_gradient = split_gradient
27         self._display_gradient_norm = display_gradient_norm
28         self._normalize_gradient = normalize_gradient
29         self._train_only_domain = train_only_domain
30         self._train_only_border = train_only_border
31
32     def flatten_tensor_list (self,
33         tensor_list):
34
35         linerarized_params = tf.constant(0., shape=[1,])
36         for param_ind in range(0, len(tensor_list)//2):
37             for param_sub_ind in range(0, tensor_list[2*param_ind].shape[0]):
38                 param = tensor_list[2*param_ind][param_sub_ind]
39                 linerarized_params = tf.concat([linerarized_params, param], axis=0)
40                 param = tensor_list[2*param_ind+1]

```

```

41     linerarized_params = tf.concat([linerarized_params, param], axis=0)
42     linerarized_params = tf.reshape(linerarized_params, [1,-1])
43     return linerarized_params
44
45     def update_parameters (self,
46                             linerarized_params):
47
48         init_pos = 1
49         final_pos = 1
50         for param_ind in range(0, len(self._model._trainable_weights)//2):
51             for param_sub_ind in range(0, self._model._trainable_weights[2*param_ind].shape[0]):
52                 final_pos = init_pos + self._model._trainable_weights[2*param_ind][param_sub_ind].shape[0]
53                 self._model._trainable_weights[2*param_ind][param_sub_ind].assign(linerarized_params[0,init_pos:final_pos])
54                 init_pos = final_pos
55                 final_pos = init_pos + self._model._trainable_weights[2*param_ind+1].shape[0]
56                 self._model._trainable_weights[2*param_ind+1].assign(linerarized_params[0,init_pos:final_pos])
57                 init_pos = final_pos
58
59     def parametric_model (self,
60                             new_parameters):
61
62         self.update_parameters(new_parameters)
63
64         loss, _, _, _ = self._model.loss_function (inputs_domain = self._inputs,
65                                                     border_data = self._border_data,
66                                                     is_training = False,
67                                                     use_only_domain = self._train_only_domain,
68                                                     use_only_border = self._train_only_border)
69
70         gradient_update = self._model.back_propagation(inputs = self._inputs,
71                                                         border_data = self._border_data,
72                                                         is_training = False,
73                                                         split_gradient = True,
74                                                         display_gradient_norm = False,
75                                                         normalize_gradient = self._normalize_gradient,
76                                                         train_only_domain = self._train_only_domain,
77                                                         train_only_border = self._train_only_border)
78
79         gradient_update = self.flatten_tensor_list(gradient_update)
80         loss = tf.reshape(loss, shape=[1])
81         return loss, gradient_update
82
83     def apply_gradients (self):
84
85         initial_params = self.flatten_tensor_list(self._model._trainable_weights)
86         if (self._name == 'L-BFGS'):
87             optimization_results = tfp.optimizer.lbfgs_minimize(self.parametric_model,
88                                                                 initial_position = initial_params,
89                                                                 num_correction_pairs = 10,
90                                                                 tolerance = 1e-32,
91                                                                 max_iterations=1,
92                                                                 parallel_iterations=1,
93                                                                 stopping_condition=None,
94                                                                 max_line_search_iterations = 10)
95         elif (self._name == 'BFGS'):
96             optimization_results = tfp.optimizer.bfgs_minimize(self.parametric_model,
97                                                                 initial_position = initial_params,
98                                                                 tolerance = 1e-32,
99                                                                 max_iterations=1,
100                                                                parallel_iterations=1,
101                                                                stopping_condition=None,
102                                                                max_line_search_iterations = 10)
103
104         optimized_params = optimization_results.position
105         self.update_parameters(optimized_params)
106
107         _ = self._model.back_propagation(inputs = self._inputs,
108                                         border_data = self._border_data,
109                                         is_training = True,
110                                         split_gradient = self._split_gradient,
111                                         display_gradient_norm = self._display_gradient_norm,
112                                         normalize_gradient = self._normalize_gradient,
113                                         train_only_domain = self._train_only_domain,
114                                         train_only_border = self._train_only_border)

```

B.6 myLayer Class

```

1  """
2  @author: Alberto Garcia Molina
3  @latest_update: 12/10/2020
4  """
5
6  class myLayer(keras.layers.Layer):
7
8      #####
9      # Initialize the Layer object.
10     #####
11     def __init__(self,
12                 name = 'newLayer'):
13
14         super(myLayer, self).__init__()
15         self._name = name
16
17     #####
18     # Constructs the Layer
19     #####
20     def build(self,
21              input_dim = 2,

```

```

22         output_dim = 2,
23         activation = 'sigmoid',
24         weight_initializer = 'xavier',
25         bias_initializer = 'xavier',
26         seed = None,
27         batch_normalization = False,
28         suppress_bias = False,
29         epsilon = 1e-12):
30
31     self._input_dim = input_dim
32     self._output_dim = output_dim
33     self._activation = activation
34     self._weight_initializer = weight_initializer
35     self._bias_initializer = bias_initializer
36     self._batch_normalization = batch_normalization
37     self._has_bias = not suppress_bias
38     self._epsilon = epsilon
39
40     if (weight_initializer == 'zeros'):
41         wInit = tf.keras.initializers.Zeros()
42     elif (weight_initializer == 'ones'):
43         wInit = tf.keras.initializers.Ones()
44     elif (weight_initializer == 'normal_0_1'):
45         wInit = RandomNormal(mean=0., stddev=1., seed=seed)
46     elif (weight_initializer == 'uniform_-1_1'):
47         wInit = tf.keras.initializers.RandomUniform(minval=-1., maxval=1., seed=seed)
48     elif (weight_initializer == 'xavier'):
49         wInit = tf.keras.initializers.GlorotNormal(seed=seed)
50     elif (weight_initializer == 'he'):
51         wInit = tf.keras.initializers.he_normal(seed=seed)
52
53     if (bias_initializer == 'zeros'):
54         bInit = tf.keras.initializers.Zeros()
55     elif (bias_initializer == 'ones'):
56         bInit = tf.keras.initializers.Ones()
57     elif (bias_initializer == 'normal_0_1'):
58         bInit = RandomNormal(mean=0., stddev=1., seed=seed)
59     elif (bias_initializer == 'uniform_-1_1'):
60         bInit = tf.keras.initializers.RandomUniform(minval=-1., maxval=1., seed=seed)
61     elif (bias_initializer == 'xavier'):
62         bInit = tf.keras.initializers.GlorotNormal(seed=seed)
63     elif (bias_initializer == 'he'):
64         bInit = tf.keras.initializers.he_normal(seed=seed)
65
66     self.w = self.add_weight(
67         name = self._name + ' W',
68         shape = (self._input_dim, self._output_dim),
69         initializer = wInit,
70         trainable = True)
71     tf.cast(self.w, tf.float32)
72
73     if (self._has_bias == True):
74         self.b = self.add_weight(
75             name = self._name + ' b',
76             shape = (self._output_dim,),
77             initializer = bInit,
78             trainable = True)
79     else:
80         bInit = tf.keras.initializers.Zeros()
81         self.b = self.add_weight(
82             name = self._name + ' b',
83             shape = (self._output_dim,),
84             initializer = bInit,
85             trainable = True)
86     tf.cast(self.b, tf.float32)
87
88     #####
89     # Feeds the input into the layer.
90     #####
91     def feed(self,
92             inputs = None):
93
94         if (self._has_bias == True):
95             outputs = tf.matmul(inputs, self.w) + self.b
96         else:
97             outputs = tf.matmul(inputs, self.w)
98
99         if (self._activation == 'sigmoid'):
100             outputs = tf.nn.sigmoid(outputs)
101         if (self._activation == 'tanh'):
102             outputs = tf.keras.activations.tanh(outputs)
103         if (self._activation == 'relu'):
104             outputs = tf.nn.relu(outputs)
105         if (self._activation == 'exponential'):
106             outputs = tf.keras.activations.exponential(outputs)
107         if (self._activation == 'elu'):
108             outputs = tf.keras.activations.elu(outputs, alpha=1.0)
109         if (self._activation == 'swish'):
110             outputs = tf.keras.activations.swish(outputs)
111         if (self._activation == 'softplus'):
112             outputs = tf.nn.softplus(outputs)
113
114
115         if (self._batch_normalization == True):
116             mean, var = tf.nn.moments(outputs, axes=0, keepdims=True)
117             outputs = (outputs-mean)/(tf.math.sqrt(var + self._epsilon))
118
119     return outputs

```


B.7 myModel Class

```
1 """
2 @author: Alberto Garcia Molina
3 @latest_update: 12/10/2020
4 """
5
6 class myModel(tf.keras.Model):
7
8     #####
9     # Initializes the model instance.
10    #####
11    def __init__(self,
12                name = 'myModel'):
13
14        super(myModel, self).__init__()
15
16        # Initializes the name and flags for the model.
17        self.name = name
18        self.built = False
19        self.is_compiled = False
20        self.has_dataset = False
21
22        # Initializes the historical training variables of the model.
23        self.num_epochs_trained = 0
24        self.losses = []
25        self.losses_domain = []
26        self.losses_border = []
27        self.losses_regularization = []
28        self.losses_solution = []
29        self.losses_validation = []
30
31    #####
32    # Builds the layers of the model.
33    #####
34    def build(self,
35             input_dim = 2,
36             hidden_dim = [5,5],
37             output_dim = 2,
38             activations = 'sigmoid',
39             weight_initializers = 'xavier',
40             bias_initializers = 'xavier',
41             batch_normalization = False,
42             suppress_bias = False,
43             seed = None,
44             epsilon = 1e-12):
45
46        # Sets up the basic characteristics of the layers in the model.
47        self._input_dim = input_dim
48        self._hidden_dim = hidden_dim
49        self._output_dim = output_dim
50        self._num_hidden_layers = len(hidden_dim)-1
51        self._activations = activations
52        self._weight_initializers = weight_initializers
53        self._bias_initializers = bias_initializers
54        self._batch_normalization = batch_normalization
55        self._has_bias = not suppress_bias
56
57        self._layers = []
58
59        # Constructs the input layer.
60        layer = myLayer('Input_Layer')
61        layer.build(input_dim = self._input_dim,
62                  output_dim = self._hidden_dim[0],
63                  activation = self._activations,
64                  weight_initializer = self._weight_initializers,
65                  bias_initializer = self._bias_initializers,
66                  seed = seed,
67                  batch_normalization = self._batch_normalization,
68                  suppress_bias = suppress_bias,
69                  epsilon = epsilon)
70        self._layers.append(layer)
71        self._trainable_weights.append(layer.variables[0])
72        self._trainable_weights.append(layer.variables[1])
73
74        # Constructs the hidden layers.
75        for layer_num in range(1, self._num_hidden_layers+1):
76            layer = myLayer('Hidden_Layer_'+str(layer_num))
77            layer.build(input_dim = self._hidden_dim[layer_num-1],
78                      output_dim = self._hidden_dim[layer_num],
79                      activation = self._activations,
80                      weight_initializer = self._weight_initializers,
81                      bias_initializer = self._bias_initializers,
82                      seed = seed,
83                      batch_normalization = self._batch_normalization,
84                      suppress_bias = suppress_bias,
85                      epsilon = epsilon)
86            self._layers.append(layer)
87            self._trainable_weights.append(layer.variables[0])
88            self._trainable_weights.append(layer.variables[1])
89
90        # Constructs the output layer.
91        layer = myLayer('Output_Layer')
92        layer.build(input_dim = self._hidden_dim[-1],
93                  output_dim = self._output_dim,
94                  activation = None,
95                  weight_initializer = self._weight_initializers,
96                  bias_initializer = self._bias_initializers,
```

```

97         seed = seed,
98         batch_normalization = False,
99         suppress_bias = suppress_bias,
100        epsilon = None)
101    self._layers.append(layer)
102    self._trainable_weights.append(layer.variables[0])
103    self._trainable_weights.append(layer.variables[1])
104
105    # Raise flag if the neural network has been built successfully.
106    self.built = True
107
108    #####
109    # Sets up the optimizer.
110    #####
111    def set_up_optimizer (self,
112                        optimizer_selection,
113                        learning_rate = 1e-03,
114                        epsilon = 1e-07):
115
116        # Sets up the information of the optimizer.
117        self._optimizer_selection = optimizer_selection
118        self._learning_rate = learning_rate
119        self._epsilon = epsilon
120
121        # Adam Optimizer (1st Order)
122        if (self._optimizer_selection == 'Adam'):
123            self._optimizer1 = tf.keras.optimizers.Adam(learning_rate = self._learning_rate,
124                                                       epsilon = self._epsilon,
125                                                       amsgrad = False)
126            self._optimizer2 = tf.keras.optimizers.Adam(learning_rate = self._learning_rate,
127                                                       epsilon = self._epsilon,
128                                                       amsgrad = False)
129
130        # AMSGrad Optimizer (1st Order)
131        elif (self._optimizer_selection == 'AMSGrad'):
132            self._optimizer1 = tf.keras.optimizers.Adam(learning_rate = self._learning_rate,
133                                                       epsilon = self._epsilon,
134                                                       amsgrad = True)
135            self._optimizer2 = tf.keras.optimizers.Adam(learning_rate = self._learning_rate,
136                                                       epsilon = self._epsilon,
137                                                       amsgrad = True)
138
139        # Nadam Optimizer (1st Order)
140        elif (self._optimizer_selection == 'Nadam'):
141            self._optimizer1 = tf.keras.optimizers.Nadam(learning_rate = self._learning_rate,
142                                                       epsilon = self._epsilon)
143            self._optimizer2 = tf.keras.optimizers.Nadam(learning_rate = self._learning_rate,
144                                                       epsilon = self._epsilon)
145
146        # AdaGrad Optimizer (1st Order)
147        elif (self._optimizer_selection == 'AdaGrad'):
148            self._optimizer1 = tf.keras.optimizers.Adagrad(learning_rate = self._learning_rate,
149                                                         epsilon = self._epsilon)
150            self._optimizer2 = tf.keras.optimizers.Adagrad(learning_rate = self._learning_rate,
151                                                         epsilon = self._epsilon)
152
153        # AdaDelta Optimizer (1st Order)
154        elif (self._optimizer_selection == 'AdaDelta'):
155            self._optimizer1 = tf.keras.optimizers.Adadelta(learning_rate = self._learning_rate,
156                                                         rho = 0.95,
157                                                         epsilon = self._epsilon)
158            self._optimizer2 = tf.keras.optimizers.Adadelta(learning_rate = self._learning_rate,
159                                                         rho = 0.95,
160                                                         epsilon = self._epsilon)
161
162        # RMSProp Optimizer (1st Order)
163        elif (self._optimizer_selection == 'RMSProp'):
164            self._optimizer1 = tf.keras.optimizers.RMSprop(learning_rate=self._learning_rate,
165                                                         epsilon = self._epsilon)
166            self._optimizer2 = tf.keras.optimizers.RMSprop(learning_rate=self._learning_rate,
167                                                         epsilon = self._epsilon)
168
169        # Vanilla SDG Optimizer (1st Order)
170        elif (self._optimizer_selection == 'Vanilla_SGD'):
171            self._optimizer1 = tf.keras.optimizers.SGD(learning_rate = self._learning_rate,
172                                                         nesterov = False)
173            self._optimizer2 = tf.keras.optimizers.SGD(learning_rate = self._learning_rate,
174                                                         nesterov = False)
175
176        # SGD with Momentum Optimizer (1st Order)
177        elif (self._optimizer_selection == 'Momentum_SGD'):
178            self._optimizer1 = keras.optimizers.SGD(learning_rate = self._learning_rate,
179                                                         momentum = 0.9,
180                                                         nesterov = False)
181            self._optimizer2 = keras.optimizers.SGD(learning_rate = self._learning_rate,
182                                                         momentum = 0.9,
183                                                         nesterov = False)
184
185        # SGD with Nesterov Momentum Optimizer (1st Order)
186        elif (self._optimizer_selection == 'Nesterov_SGD'):
187            self._optimizer1 = keras.optimizers.SGD(learning_rate = self._learning_rate,
188                                                         momentum = 0.9,
189                                                         nesterov = True)
190            self._optimizer2 = keras.optimizers.SGD(learning_rate = self._learning_rate,
191                                                         momentum = 0.9,
192                                                         nesterov = True)
193
194        # BFGS Optimizer (2st Order)
195        elif (self._optimizer_selection == 'BFGS'):
196            self._optimizer = secondOrderOptimizers(name = 'BFGS',
197                                                         model = self)
198
199        # L-BFGS Optimizer (2st Order)
200        elif (self._optimizer_selection == 'L-BFGS'):
201            self._optimizer = secondOrderOptimizers(name = 'L-BFGS',
202                                                         model = self)
203
204        else:
205            self._is_compiled = False
206            raise Exception("Invalid optimizer.")
207
208    #####
209    # Builds the problem instance and training set up.

```

```

199 #####
200 def compile (self,
201             differential_operator = None,
202             external_force = None,
203             exact_solution = None,
204             optimizer_selection = None,
205             learning_rate = 1e-03,
206             epsilon = 1e-07,
207             scale_factor = 1,
208             loss_fuction = 'square_L2_error',
209             regularization = None,
210             regularization_coef = 0,
211             clip_gradient = 'global'):
212
213     # Sets up he problem solved by the model.
214     self._differential_operator = differential_operator
215     self._external_force = external_force
216     self._exact_solution = exact_solution
217
218     # Sets up the regularization and loss options.
219     self._scale_factor = scale_factor
220     self._loss_fuction = loss_fuction
221     self._regularization = regularization
222     self._regularization_coef = regularization_coef
223     self._clip_gradient = clip_gradient
224
225     # Constructs the optimizer and validates the instances.
226     if (regularization == None):
227         print('No regularization introduced, using default None')
228     self._required_derivative_order = problemInstance.\
229                                     instance_exists(differential_operator = self._differential_operator,
230                                                       external_force = self._external_force,
231                                                       exact_solution = self._exact_solution)
232     self.set_up_optimizer(optimizer_selection = optimizer_selection,
233                          learning_rate = learning_rate,
234                          epsilon = epsilon)
235
236     # Raise flag if the problem and training instance has been built successfully.
237     self._is_compiled = True
238
239 #####
240 # Feed forward of the neural network, returning also the gradient wrt inputs.
241 #####
242 def predict (self,
243             inputs,
244             return_derivative_order = 0):
245
246     if (self.built == False):
247         raise Exception("Cannot feed forward, the model is not built.")
248
249     outputs_derivatives = []
250     if (return_derivative_order in (0,1,2,3)):
251
252         # Output with 0 order derivative.
253         if (return_derivative_order == 0):
254             outputs = self._layers[0].feed(inputs)
255             for layer_ind in range(1, self._num_hidden_layers+2):
256                 outputs = self._layers[layer_ind].feed(outputs)
257             tf.debugging.check_numerics(outputs, message = 'NaN occurred in network output.')

```

```

301         outputs_2nd_der.append(outputs_2nd_der_var)
302         del tape_ord1
303         del tape_ord2
304
305     outputs_derivatives.append(tf.squeeze(tf.stack(outputs_1st_der, axis = 1), axis = 2))
306     outputs_derivatives.append(tf.squeeze(tf.stack(outputs_2nd_der, axis = 1), axis = 2))
307
308     tf.debugging.check_numerics(outputs,
309                               message = 'NaN occurred in network output.')
```

```

310     tf.debugging.check_numerics(outputs_derivatives[0],
311                               message = 'NaN occurred in network 1st derivative output.')
```

```

312     tf.debugging.check_numerics(outputs_derivatives[1],
313                               message = 'NaN occurred in network 2nd derivative output.')
```

```

314
315     # Output with 3rd order derivatives. (CORRECT FOR THE THIRD ORDER DERIVATIVE RIGHT)
316     #elif (return_derivative_order == 3):
317     #     with tf.GradientTape(persistent = False) as tape_ord3:
318     #         tape_ord3.watch(inputs)
319     #         with tf.GradientTape(persistent = False) as tape_ord2:
320     #             tape_ord2.watch(inputs)
321     #             with tf.GradientTape(persistent = False) as tape_ord1:
322     #                 tape_ord1.watch(inputs)
323     #                 outputs = self._layers[0].feed(inputs)
324     #                 for layer_ind in range(1, self._num_hidden_layers+2):
325     #                     outputs = self._layers[layer_ind].feed(outputs)
326     #                     outputs_1st_der = tape_ord1.gradient(outputs,
327     #                                                         inputs)
328     #                     outputs_derivatives.append(outputs_1st_der)
329     #                     del tape_ord1
330     #                     outputs_2nd_der = tape_ord2.gradient(outputs_1st_der,
331     #                                                         inputs)
332     #                     outputs_derivatives.append(outputs_2nd_der)
333     #                     del tape_ord2
334     #                     outputs_3rd_der = tape_ord3.gradient(outputs_2nd_der,
335     #                                                         inputs)
336     #                     outputs_derivatives.append(outputs_3rd_der)
337     #                     del tape_ord3
338     #                     tf.debugging.check_numerics(outputs,
339     #                                                 message = 'NaN occurred in network output.')
```

```

340     #                     tf.debugging.check_numerics(outputs_1st_der,
341     #                                                 message = 'NaN occurred in network 1st derivative output.')
```

```

342     #                     tf.debugging.check_numerics(outputs_2nd_der,
343     #                                                 message = 'NaN occurred in network 2nd derivative output.')
```

```

344     #                     tf.debugging.check_numerics(outputs_3rd_der,
345     #                                                 message = 'NaN occurred in network 3rd derivative output.')
```

```

346
347     else:
348         raise Exception("Invalid order for network derivative computation.")
349
350     return outputs, outputs_derivatives
351
352     #####
353     # Calculates the loss function.
354     #####
355     def loss_function (self,
356                      inputs_domain,
357                      border_data = None,
358                      is_training = False,
359                      use_only_domain = False,
360                      use_only_border = False):
361
362         # Initializes the losses variables.
363         loss_domain = tf.constant(0.)
364         loss_border = tf.constant(0.)
365         loss_regularization = tf.constant(0.)
366         loss_solution = tf.constant(0.)
367
368         # Evaluates the left-hand-side and the right-hand-side of the differential equation and solution.
369
370         # Use the gradient balancing regularization.
371         if (self._regularization != 'Gradient_Type'):
372             outputs, outputs_derivatives = self.predict(inputs = inputs_domain,
373                                                       return_derivative_order = self._required_derivative_order)
374         else:
375             with tf.GradientTape(persistent = True) as tape_reg:
376                 tape_reg.watch(self._trainable_weights)
377                 outputs, outputs_derivatives = self.predict(inputs = inputs_domain,
378                                                           return_derivative_order = self._required_derivative_order)
379                 outputs_param_der = tape_reg.gradient(outputs,
380                                                       self._trainable_weights,
381                                                       unconnected_gradients = tf.UnconnectedGradients.ZERO)
382                 outputs_derivatives_param_der = tape_reg.gradient(outputs_derivatives[0],
383                                                                 self._trainable_weights,
384                                                                 unconnected_gradients = tf.UnconnectedGradients.ZERO)
385                 for weight_ind in range(self._num_hidden_layers+2):
386                     loss_regularization += tf.reduce_mean(tf.square(outputs_derivatives_param_der[2*weight_ind]
387                                                                 - outputs_param_der[2*weight_ind]))
388                     loss_regularization += tf.reduce_mean(tf.square(outputs_derivatives_param_der[2*weight_ind+1]
389                                                                 - outputs_param_der[2*weight_ind+1]))
390                 del tape_reg
391
392         diff_op_output = problemInstance.\
393             differential_operator(inputs = inputs_domain,
394                                outputs = outputs,
395                                outputs_derivatives = outputs_derivatives,
396                                differential_operator = self._differential_operator,
397                                input_dim = self._input_dim,
398                                output_dim = self._output_dim)
399
400         ext_force_output = problemInstance.\
401             external_force(inputs = inputs_domain,
402                          external_force = self._external_force,
403                          input_dim = self._input_dim,
```

```

403             output_dim = self._output_dim)
404
405     exact_sol_output = problemInstance.\
406         exact_solution(inputs = inputs_domain,
407                        exact_solution = self._exact_solution,
408                        input_dim = self._input_dim,
409                        output_dim = self._output_dim)
410
411     # Evaluates the left-hand-side of the border conditions.
412     # Right-hand-side already computed in border_data[1:n].
413     if (border_data != None):
414         outputs_border, outputs_derivatives_border = self.predict(inputs = border_data[0],
415                                                                    return_derivative_order = self._required_derivative_order)
416
417         lhs_border = problemInstance.\
418             lhs_boundary_conditions (inputs = border_data[0],
419                                     outputs = outputs_border,
420                                     outputs_derivatives = outputs_derivatives_border,
421                                     border_type = self._border_type,
422                                     border_batch_size = self._border_training_batch_size,
423                                     external_force = self._external_force,
424                                     required_derivative_order = self._required_derivative_order,
425                                     input_dim = self._input_dim,
426                                     output_dim = self._output_dim)
427
428     # Computes the loss function for the L2 Error.
429     if (self._loss_fuction == 'L2_error'):
430         loss_domain = tf.reduce_mean(
431             tf.norm(diff_op_output-ext_force_output,
432                   ord = 'euclidean',
433                   axis = 1))
434
435         if (border_data != None):
436             for ind in range(len(border_data)-1):
437                 loss_border += tf.reduce_mean(
438                     tf.norm(lhs_border[ind] - border_data[ind+1],
439                           ord='euclidean',
440                           axis=1))
441
442             if (self._regonov == 'Tikhonov'):
443                 for weight_ind in range(self._num_hidden_layers+2):
444                     loss_regularization += tf.reduce_mean(
445                         tf.norm(self._trainable_weights[2*weight_ind],
446                               ord='euclidean',
447                               axis = 1))
448
449                 loss_regularization += tf.reduce_mean(
450                     tf.norm(self._trainable_weights[2*weight_ind+1],
451                           ord='euclidean',
452                           axis = 0))
453
454             elif (self._regonov == None
455                  or self._regonov == 'Gradient_Type'
456                  or self._regonov == 'Quadratic_Balance'):
457                 pass
458             else:
459                 print('Invalid regularization option, defaulting to none.')
460                 self._regonov = None
461
462         loss_solution = tf.reduce_mean(
463             tf.norm(outputs-exact_sol_output,
464                   ord='euclidean',
465                   axis=1))
466
467     # Computes the loss function for the Square L2 Error (MSE).
468     elif (self._loss_fuction == 'square_L2_error'):
469         loss_domain = tf.reduce_mean(
470             tf.reduce_sum(tf.square(diff_op_output-ext_force_output),
471                           axis = 1,
472                           keepdims = True))
473
474         if (border_data != None):
475             for ind in range(len(border_data)-1):
476                 loss_border += tf.reduce_mean(
477                     tf.reduce_sum(tf.square(lhs_border[ind] - border_data[ind+1]),
478                                   axis = 1,
479                                   keepdims = True))
480
481             if (self._regonov == 'Tikhonov'):
482                 for weight_ind in range(self._num_hidden_layers+2):
483                     loss_regularization += tf.reduce_mean(
484                         tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind]),
485                                       axis = 1,
486                                       keepdims = True))
487
488                 loss_regularization += tf.reduce_mean(
489                     tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind+1]),
490                                   axis = 0,
491                                   keepdims = True))
492
493             elif (self._regonov == None
494                  or self._regonov == 'Gradient_Type'
495                  or self._regonov == 'Quadratic_Balance'):
496                 pass
497             else:
498                 print('Invalid regularization option, defaulting to none.')
499                 self._regonov = None
500
501         loss_solution = tf.reduce_mean(
502             tf.reduce_sum(tf.square(outputs-exact_sol_output),
503                           axis = 1,
504                           keepdims = True))
505
506     # Computes the loss function for the Absolute Error (L1).
507     elif (self._loss_fuction == 'absolute_error'):
508         loss_domain = tf.reduce_mean(
509             tf.reduce_sum(tf.abs(diff_op_output-ext_force_output),
510                           axis = 1,
511                           keepdims = True))
512
513         if (border_data != None):
514             for ind in range(len(border_data)-1):
515                 loss_border += tf.reduce_mean(
516                     tf.reduce_sum(tf.abs(lhs_border[ind] - border_data[ind+1]),

```

```

505         axis = 1,
506         keepdims = True))
507     if (self._regularization == 'Tikhonov'):
508         for weight_ind in range(self._num_hidden_layers+2):
509             loss_regularization += tf.math.reduce_mean(
510                 tf.reduce_mean(tf.abs(self._trainable_weights[2*weight_ind]),
511                     axis = 1,
512                     keepdims = True))
513             loss_regularization += tf.math.reduce_mean(
514                 tf.reduce_mean(tf.abs(self._trainable_weights[2*weight_ind+1]),
515                     axis = 0,
516                     keepdims = True))
517
518     elif (self._regularization == None
519           or self._regularization == 'Gradient_Type'
520           or self._regularization == 'Quadratic_Balance'):
521         pass
522     else:
523         print('Invalid regularization option, defaulting to none.')
524         self._regularization = None
525     loss_solution = tf.reduce_mean(
526         tf.reduce_sum(tf.abs(outputs-exact_sol_output),
527             axis = 1,
528             keepdims = True))
529
530 # Experimental: Computes the loss (MSE) proportional to the external force.
531 elif (self._loss_fuction == 'force_proportional_error'):
532     loss_domain = tf.reduce_mean(
533         tf.square((diff_op_output-ext_force_output)
534             *(ext_force_output+1e-12)))
535
536     if (border_data != None):
537         for ind in range(len(border_data)-1):
538             loss_border += tf.reduce_mean(
539                 tf.square(lhs_border[ind] - border_data[ind+1]))
540
541     if (self._regularization == 'Tikhonov'):
542         for weight_ind in range(self._num_hidden_layers+2):
543             loss_regularization += tf.reduce_mean(
544                 tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind]),
545                     axis = 1,
546                     keepdims = True))
547             loss_regularization += tf.reduce_mean(
548                 tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind+1]),
549                     axis = 0,
550                     keepdims = True))
551
552     elif (self._regularization == None
553           or self._regularization == 'Gradient_Type'
554           or self._regularization == 'Quadratic_Balance'):
555         pass
556     else:
557         print('Invalid regularization option, defaulting to none.')
558         self._regularization = None
559     loss_solution = tf.reduce_mean(
560         tf.square(outputs-exact_sol_output))
561
562 # Computes the square of the MSE, i.e. the  $\| \cdot \|_2$  error .
563 elif (self._loss_fuction == 'square_MSE'):
564     loss_domain = tf.reduce_mean(
565         tf.reduce_sum(tf.square(tf.square(diff_op_output-ext_force_output)),
566             axis = 1,
567             keepdims = True))
568
569     if (border_data != None):
570         for ind in range(len(border_data)-1):
571             loss_border += tf.reduce_mean(
572                 tf.reduce_sum(tf.square(tf.square(lhs_border[ind] - border_data[ind+1])),
573                     axis = 1,
574                     keepdims = True))
575
576     if (self._regularization == 'Tikhonov'):
577         for weight_ind in range(self._num_hidden_layers+2):
578             loss_regularization += tf.reduce_mean(
579                 tf.reduce_mean(tf.square(tf.square(self._trainable_weights[2*weight_ind])),
580                     axis = 1,
581                     keepdims = True))
582             loss_regularization += tf.reduce_mean(
583                 tf.reduce_mean(tf.square(tf.square(self._trainable_weights[2*weight_ind+1])),
584                     axis = 0,
585                     keepdims = True))
586
587     elif (self._regularization == None
588           or self._regularization == 'Gradient_Type'
589           or self._regularization == 'Quadratic_Balance'):
590         pass
591     else:
592         print('Invalid regularization option, defaulting to none.')
593         self._regularization = None
594     loss_solution = tf.reduce_mean(
595         tf.reduce_sum(tf.square(outputs-exact_sol_output),
596             axis = 1,
597             keepdims = True))
598
599 # Experimental: Computes the loss (MSE) proportional to the square of the inputs.
600 elif (self._loss_fuction == 'input_proportional_error'):
601     loss_domain = tf.reduce_mean(
602         tf.square((diff_op_output-ext_force_output)
603             *inputs_domain*inputs_domain))
604
605     if (border_data != None):
606         for ind in range(len(border_data)-1):
607             loss_border += tf.reduce_mean(
608                 tf.square(lhs_border[ind] - border_data[ind+1]))
609
610     if (self._regularization == 'Tikhonov'):
611         for weight_ind in range(self._num_hidden_layers+2):
612             loss_regularization += tf.reduce_mean(
613                 tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind]),
614                     axis = 1,
615                     keepdims = True))

```

```

607         loss_regularization += tf.reduce_mean(
608             tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind+1]),
609                 axis = 0,
610                 keepdims = True))
611     elif (self._regularization == None
612           or self._regularization == 'Gradient_Type'
613           or self._regularization == 'Quadratic_Balance'):
614         pass
615     else:
616         print('Invalid regularization option, defaulting to none.')
617         self._regularization = None
618     loss_solution = tf.reduce_mean(
619         tf.square(outputs-exact_sol_output))
620
621 # Computes the loss (MSE) with respect to one component.
622 elif (self._loss_fuction == 'square_L2_error_1st_comp'):
623     loss_domain = tf.reduce_mean(tf.square(diff_op_output-ext_force_output)[:,:0])
624     if (border_data != None):
625         for ind in range(len(border_data)-1):
626             loss_border += tf.reduce_mean(
627                 tf.reduce_sum(tf.square(lhs_border[ind] - border_data[ind+1]),
628                     axis = 1,
629                     keepdims = True))
630     if (self._regularization == 'Tikhonov'):
631         for weight_ind in range(self._num_hidden_layers+2):
632             loss_regularization += tf.reduce_mean(
633                 tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind]),
634                     axis = 1,
635                     keepdims = True))
636             loss_regularization += tf.reduce_mean(
637                 tf.reduce_mean(tf.square(self._trainable_weights[2*weight_ind+1]),
638                     axis = 0,
639                     keepdims = True))
640     elif (self._regularization == None
641           or self._regularization == 'Gradient_Type'
642           or self._regularization == 'Quadratic_Balance'):
643         pass
644     else:
645         print('Invalid regularization option, defaulting to none.')
646         self._regularization = None
647     loss_solution = tf.reduce_mean(tf.square(outputs-exact_sol_output)[0])
648
649 # Error for invalid loss option.
650 else:
651     raise Exception("Invalid loss option. Please recompile with a valid name.")
652
653 # Experimental: Implement the quadratic loss balance regularization.
654 if (self._regularization == 'Quadratic_Balance'):
655     if (inputs_border != None and expected_outputs_border != None):
656         loss_regularization += tf.sqrt(tf.square(loss_domain-loss_border))
657
658 # Implements the train only domain or border options
659 coef_domain = 1
660 coef_border = 1
661 if (use_only_domain == True):
662     coef_border = 0.
663 if (use_only_border == True):
664     coef_domain = 0.
665
666 # Computes the total loss and checks for explosions.
667 loss = coef_domain*loss_domain + coef_border*loss_border
668 tf.debugging.check_numerics(loss_domain, message='NaN occurred in domain loss fuction.')
669 tf.debugging.check_numerics(loss_border, message='NaN occurred in border loss fuction.')
670 tf.debugging.check_numerics(loss_regularization, message='NaN occurred in regularization loss fuction.')
671 tf.debugging.check_numerics(loss, message='NaN occurred in total loss fuction.')
672
673 # If the method is set to training mode, the losses are saved on the historical training variables.
674 if (is_training == True):
675     self._losses_domain.append(loss_domain.numpy())
676     self._losses_border.append(loss_border.numpy())
677     self._losses_regularization.append(loss_regularization.numpy())
678     self._losses.append(loss.numpy())
679     self._losses_solution.append(loss_solution.numpy())
680
681     return loss, loss_domain, loss_border, loss_regularization
682
683 #####
684 # Computes the gradient wrt the weights.
685 #####
686 def back_propagation (self,
687                     inputs,
688                     border_data,
689                     is_training,
690                     split_gradient = False,
691                     display_gradient_norm = False,
692                     normalize_gradient = False,
693                     train_only_domain = False,
694                     train_only_border = False):
695
696 # Executes the back propagation
697 with tf.GradientTape(persistent = True) as tape_bp:
698     tape_bp.watch(self._trainable_weights)
699     loss, loss_domain, loss_border,\
700     loss_regularization = self.loss_function(inputs_domain = inputs,
701                                             border_data = border_data,
702                                             is_training = is_training,
703                                             use_only_domain = train_only_domain,
704                                             use_only_border = train_only_border)
705
706     total_loss_f = loss+self._regularization_coef*loss_regularization
707
708     gradient_update = tape_bp.gradient(total_loss_f,

```

```

709         self._trainable_weights,
710         unconnected_gradients = tf.UnconnectedGradients.ZERO)
711
712 # Splits the gradient wrt to each individual loss component.
713 if (split_gradient == True):
714     gradient_update_domain = tape_bp.gradient(loss_domain,
715         self._trainable_weights,
716         unconnected_gradients = tf.UnconnectedGradients.ZERO)
717     gradient_update_border = tape_bp.gradient(loss_border,
718         self._trainable_weights,
719         unconnected_gradients = tf.UnconnectedGradients.ZERO)
720     gradient_update_regularization = tape_bp.gradient(loss_regularization,
721         self._trainable_weights,
722         unconnected_gradients = tf.UnconnectedGradients.ZERO)
723
724 del tape_bp
725
726 # Avoids NAN propagation by rightfully setting them to 0.
727 gradient_update = [tf.where(tf.math.is_nan(g), tf.zeros_like(g), g)
728     for g in gradient_update]
729
730 if (split_gradient == True):
731     gradient_update_domain = [tf.where(tf.math.is_nan(g), tf.zeros_like(g), g)
732         for g in gradient_update_domain]
733     gradient_update_border = [tf.where(tf.math.is_nan(g), tf.zeros_like(g), g)
734         for g in gradient_update_border]
735     gradient_update_regularization = [tf.where(tf.math.is_nan(g), tf.zeros_like(g), g)
736         for g in gradient_update_regularization]
737
738 # Applies clipping regularization to bound the gradients.
739 if (self._clip_gradient == 'global'):
740     gradient_update = tf.clip_by_global_norm(gradient_update, 1e+1)[0]
741     gradient_update_domain = tf.clip_by_global_norm(gradient_update_domain, 1e+1)[0]
742     gradient_update_border = tf.clip_by_global_norm(gradient_update_border, 1e+1)[0]
743     gradient_update_regularization = tf.clip_by_global_norm(gradient_update_regularization, 1e+1)[0]
744 elif (self._clip_gradient == 'value'):
745     gradient_update = [tf.clip_by_value(g, clip_value_min = -1e+1, clip_value_max = 1e+1)
746         for g in gradient_update]
747     gradient_update_domain = [tf.clip_by_value(g, clip_value_min = -1e+1, clip_value_max = 1e+1)
748         for g in gradient_update_domain]
749     gradient_update_border = [tf.clip_by_value(g, clip_value_min = -1e+1, clip_value_max = 1e+1)
750         for g in gradient_update_border]
751     gradient_update_regularization = [tf.clip_by_value(g, clip_value_min = -1e+1, clip_value_max = 1e+1)
752         for g in gradient_update_regularization]
753 elif (self._clip_gradient == 'norm'):
754     gradient_update = [tf.clip_by_norm(g, 1e+1) for g in gradient_update]
755     gradient_update_domain = [tf.clip_by_norm(g, 1e+1) for g in gradient_update_domain]
756     gradient_update_border = [tf.clip_by_norm(g, 1e+1) for g in gradient_update_border]
757     gradient_update_regularization = [tf.clip_by_norm(g, 1e+1) for g in gradient_update_regularization]
758 elif (self._clip_gradient == None):
759     pass
760 else:
761     print('Invalid clipping option, defaulting to global.')
762     self._clip_gradient = 'global'
763
764 # Applies gradient normalization regularization.
765 if (normalize_gradient == True):
766     norm = tf.linalg.global_norm(gradient_update)
767     gradient_update = [g/norm for g in gradient_update]
768
769 # Rescale Gradient Regularization (Always On)
770 if (border_data != None):
771     for layer_num in range(self._num_hidden_layers+2):
772         weight_norm_domain = tf.norm(gradient_update_domain[2*layer_num],
773             ord = 'euclidean',
774             axis = 0)
775         bias_norm_domain = tf.norm(gradient_update_domain[2*layer_num+1],
776             ord = 'euclidean',
777             axis = 0)
778         weight_norm_border = tf.norm(gradient_update_border[2*layer_num],
779             ord = 'euclidean',
780             axis = 0)
781         bias_norm_border = tf.norm(gradient_update_border[2*layer_num+1],
782             ord = 'euclidean',
783             axis = 0)
784         weight_norm_regularization = tf.norm(gradient_update_regularization[2*layer_num],
785             ord = 'euclidean',
786             axis = 0)
787         bias_norm_regularization = tf.norm(gradient_update_regularization[2*layer_num+1],
788             ord = 'euclidean',
789             axis = 0)
790
791         weight_norm = tf.minimum(weight_norm_domain, weight_norm_border)
792         bias_norm = tf.minimum(bias_norm_domain, bias_norm_border)
793
794         gradient_update_domain[2*layer_num] = gradient_update_domain[2*layer_num]*weight_norm/(weight_norm_domain+1e-31)
795         gradient_update_domain[2*layer_num+1] = gradient_update_domain[2*layer_num+1]*bias_norm/(bias_norm_domain+1e-31)
796
797         gradient_update_border[2*layer_num] = gradient_update_border[2*layer_num]*weight_norm/(weight_norm_border+1e-31)
798         gradient_update_border[2*layer_num+1] = gradient_update_border[2*layer_num+1]*bias_norm/(bias_norm_border+1e-31)
799
800         gradient_update_regularization[2*layer_num] = gradient_update_regularization[2*layer_num]*weight_norm/(
801             weight_norm_regularization+1e-31)
802         gradient_update_regularization[2*layer_num+1] = gradient_update_regularization[2*layer_num+1]*weight_norm/(
803             bias_norm_regularization+1e-31)
804
805         gradient_update[2*layer_num] = (self._scale_factor*gradient_update_domain[2*layer_num]
806             + gradient_update_border[2*layer_num]
807             + self._regularization_coef*gradient_update_regularization[2*layer_num])
808         gradient_update[2*layer_num+1] = (self._scale_factor*gradient_update_domain[2*layer_num+1]
809             + gradient_update_border[2*layer_num+1]
810             + self._regularization_coef*gradient_update_regularization[2*layer_num+1])
811
812 # Displays the gradient(s) if the option is selected.

```



```

809     if (display_gradient_norm == True):
810         total_norm = tf.linalg.global_norm(gradient_update)
811         domain_norm = tf.linalg.global_norm(gradient_update_domain)
812         border_norm = tf.linalg.global_norm(gradient_update_border)
813         regularization_norm = tf.linalg.global_norm(gradient_update_regularization)
814         print(' Total Gradient Norm', str(total_norm.numpy()))
815         print(' Domain Gradient Norm', str(domain_norm.numpy()))
816         print(' Border Gradient Norm', str(border_norm.numpy()))
817         print(' Regularization Gradient Norm', str(regularization_norm.numpy()))
818     for layer_num in range(self._num_hidden_layers+2):
819         print(' ', self._layers[layer_num]._name, ' Total Weight Gradient Norm: ',
820               str(tf.norm(gradient_update[2*layer_num],
821                           ord = 'euclidean', axis = 1).numpy()))
822         print(' ', self._layers[layer_num]._name, ' Domain Weight Gradient Norm: ',
823               str(tf.norm(gradient_update_domain[2*layer_num],
824                           ord = 'euclidean', axis = 1).numpy()))
825         print(' ', self._layers[layer_num]._name, ' Border Weight Gradient Norm: ',
826               str(tf.norm(gradient_update_border[2*layer_num],
827                           ord = 'euclidean', axis = 1).numpy()))
828         print(' ', self._layers[layer_num]._name, ' Regularization Weight Gradient Norm: ',
829               str(tf.norm(gradient_update_regularization[2*layer_num],
830                           ord = 'euclidean', axis = 1).numpy()))
831         print(' ', self._layers[layer_num]._name, ' Total Bias Gradient Norm: ',
832               str(tf.norm(gradient_update[2*layer_num+1],
833                           ord = 'euclidean', axis = 0).numpy()))
834         print(' ', self._layers[layer_num]._name, ' Domain Bias Gradient Norm: ',
835               str(tf.norm(gradient_update_domain[2*layer_num+1],
836                           ord = 'euclidean', axis = 0).numpy()))
837         print(' ', self._layers[layer_num]._name, ' Border Bias Gradient Norm: ',
838               str(tf.norm(gradient_update_border[2*layer_num+1],
839                           ord = 'euclidean', axis = 0).numpy()))
840         print(' ', self._layers[layer_num]._name, ' Regularization Bias Gradient Norm: ',
841               str(tf.norm(gradient_update_regularization[2*layer_num+1],
842                           ord = 'euclidean', axis = 0).numpy()))
843
844     return gradient_update
845
846     #####
847     # Applies an optimization step.
848     #####
849     def apply_training_step (self,
850                             inputs,
851                             border_data,
852                             split_gradient = False,
853                             display_gradient_norm = False,
854                             normalize_gradient = False,
855                             train_only_domain = False,
856                             train_only_border = False):
857
858     if (self._optimizer_selection != 'L-BFGS' and self._optimizer_selection != 'BFGS'):
859         gradient_update = self.back_propagation(inputs = inputs,
860                                                 border_data = border_data,
861                                                 is_training = True,
862                                                 split_gradient = split_gradient,
863                                                 display_gradient_norm = display_gradient_norm,
864                                                 normalize_gradient = normalize_gradient,
865                                                 train_only_domain = train_only_domain,
866                                                 train_only_border = train_only_border)
867
868     if (train_only_domain == True):
869         self._optimizer1.apply_gradients(zip(gradient_update, self._trainable_weights))
870     elif (train_only_border == True):
871         self._optimizer2.apply_gradients(zip(gradient_update, self._trainable_weights))
872     else:
873         self._optimizer1.apply_gradients(zip(gradient_update, self._trainable_weights))
874
875     else:
876         self._optimizer.apply_gradients()
877
878     # Applies clipping regularization to bound the weights.
879     for layer_num in range(self._num_hidden_layers+2):
880         self._trainable_weights[2*layer_num].assign(tf.clip_by_value(self._trainable_weights[2*layer_num],
881                                                                     clip_value_min = -1e+5,
882                                                                     clip_value_max = +1e+5))
883         self._trainable_weights[2*layer_num+1].assign(tf.clip_by_value(self._trainable_weights[2*layer_num+1],
884                                                                     clip_value_min = -1e+5,
885                                                                     clip_value_max = +1e+5))
886
887     # Clip by magnitude
888     weight_tensor = self._trainable_weights[2*layer_num]
889     weight_sign = tf.math.sign(self._trainable_weights[2*layer_num])
890     clipped_weight_tensor = tf.clip_by_value(tf.abs(weight_tensor),
891                                             clip_value_min = 1e-3,
892                                             clip_value_max = 1e+2)
893     self._trainable_weights[2*layer_num].assign(weight_sign * clipped_weight_tensor)
894
895     # bias_tensor = self._trainable_weights[2*layer_num+1]
896     bias_sign = tf.math.sign(self._trainable_weights[2*layer_num+1])
897     clipped_bias_tensor = tf.clip_by_value(tf.abs(bias_tensor),
898                                           clip_value_min = 1e-3,
899                                           clip_value_max = 1e+2)
900     self._trainable_weights[2*layer_num+1].assign(bias_sign * clipped_bias_tensor)
901
902     #####
903     # Loads the training sets to train the network.
904     #####
905     def use_training_sets (self,
906                           data_set):
907
908     if (data_set == None):
909         raise Exception("No data set loaded.")
910
911     training_batch_size, border_training_batch_size, validation_batch_size, \
912     input_dim, method, domain, border = data_set.get_set_metadata()

```

```

911         if (self._input_dim != input_dim):
912             raise Exception("Data set input dimension incompatible with neural network.")
913
914         training_set, border_training_set, validation_set = data_set.get_sets()
915         print('Dataset uploaded.')
916
917         return training_batch_size, border_training_batch_size, validation_batch_size,\
918             training_set, border_training_set, validation_set,\
919             method, domain, border
920
921     #####
922     # Trains the model.
923     #####
924     def fit (self,
925             data_set = None,
926             num_epochs = 1000,
927             plot_interval = 500,
928             validation_interval = 5000,
929             stagnation_stop_tol = 1e-16,
930             split_gradient = False,
931             display_gradient_norm = False,
932             normalize_gradient = False,
933             train_only_domain = False,
934             train_only_border = False,
935             dual_training = False,
936             auto_batch_rotation = 0):
937
938         # Loads the data set.
939         training_batch_size, border_training_batch_size, validation_batch_size,\
940             training_set, border_training_set, validation_set,\
941             method, domain_type, border_type = self.use_training_sets(data_set = data_set)
942
943         self._border_type = border_type
944         self._border_training_batch_size = border_training_batch_size
945
946         border_data = problemInstance.\
947             border_data_prep(border_inputs = border_training_set,
948                             border_type = border_type,
949                             border_batch_size = self._border_training_batch_size,
950                             external_force = self._external_force,
951                             required_derivative_order = self._required_derivative_order,
952                             input_dim = self._input_dim,
953                             output_dim = self._output_dim)
954
955         # Loads data in optimizer if BFGS or L-BFGS
956         if (self._optimizer_selection == 'L-BFGS' or self._optimizer_selection == 'BFGS'):
957             self._optimizer.optimizer_train_data (inputs = training_set,
958                                                  border_data = border_data,
959                                                  split_gradient = split_gradient,
960                                                  display_gradient_norm = display_gradient_norm,
961                                                  normalize_gradient = normalize_gradient,
962                                                  train_only_domain = train_only_domain,
963                                                  train_only_border = train_only_border)
964
965         # First Iteration.
966         start_time = time.perf_counter()
967         target_total_num_epochs = self._num_epochs_trained + num_epochs - 1
968         if (self._num_epochs_trained == 0):
969             loss_validation, _, _, _ = self.loss_function (inputs_domain = validation_set,
970                                                         border_data = border_data,
971                                                         is_training = False,
972                                                         use_only_domain = False,
973                                                         use_only_border = False)
974             self._losses_validation.append(loss_validation.numpy())
975
976         self._num_epochs_trained += 1
977         self.apply_training_step (inputs = training_set,
978                                  border_data = border_data,
979                                  split_gradient = split_gradient,
980                                  display_gradient_norm = display_gradient_norm,
981                                  normalize_gradient = normalize_gradient,
982                                  train_only_domain = train_only_domain,
983                                  train_only_border = train_only_border)
984
985         print('Epoch:', str(self._num_epochs_trained), 'Training Loss:', str(self._losses[-1]),
986               'Training Loss wrt Sol:', str(self._losses_solution[-1]))
987         print('Domain Loss:', str(self._losses_domain[-1]),
988               'Border Loss:', str(self._losses_border[-1]),
989               'Regularization Loss:', str(self._losses_regularization[-1]))
990
991         # Training Loop
992         while (self._num_epochs_trained <= target_total_num_epochs):
993             self._num_epochs_trained += 1
994             if (dual_training == True):
995                 if (self._num_epochs_trained % 2 == 0):
996                     train_only_domain = True
997                     train_only_border = False
998                 else:
999                     train_only_domain = False
1000                     train_only_border = True
1001
1002             self.apply_training_step (inputs = training_set,
1003                                       border_data = border_data,
1004                                       split_gradient = split_gradient,
1005                                       display_gradient_norm = display_gradient_norm,
1006                                       normalize_gradient = normalize_gradient,
1007                                       train_only_domain = train_only_domain,
1008                                       train_only_border = train_only_border)
1009
1010             print('Epoch:', str(self._num_epochs_trained), 'Training Loss:', str(self._losses[-1]),
1011                   'Training Loss wrt Sol:', str(self._losses_solution[-1]))
1012             print('Domain Loss:', str(self._losses_domain[-1]),

```

```

1013         'Border Loss:', str(self._losses_border[-1]),
1014         'Regularization Loss:', str(self._losses_regularization[-1]))
1015
1016     # Stagnation Stop.
1017     if (abs(self._losses_domain[-1]-self._losses_domain[-2]) <= stagnation_stop_tol
1018         and abs(self._losses_border[-1]-self._losses_border[-2]) <= stagnation_stop_tol
1019         and abs(self._losses_regularization[-1]-self._losses_regularization[-2]) <= stagnation_stop_tol):
1020         break
1021
1022     # Early Stop (Validation loss evaluation).
1023     if (self._num_epochs_trained % validation_interval == 0):
1024         loss_validation, _, _ = self.loss_function (inputs_domain = validation_set,
1025                                                    border_data = border_data,
1026                                                    is_training = False,
1027                                                    use_only_domain = False,
1028                                                    use_only_border = False)
1029         self._losses_validation.append(loss_validation.numpy())
1030         print('Evaluation Loss on Epoch' , str(self._num_epochs_trained), ':',
1031               str(self._losses_validation[-1]))
1032         if (self._losses_validation[-1] > self._losses_validation[-2]):
1033             print('Early stop - Loss wrt validation set worsen.')
1034             break
1035
1036     # Autoregulates the regularization coefficient if it only improves the regularization term.
1037     if (self._losses_domain[-1] > self._losses_domain[-validation_interval]
1038         or self._losses_border[-1] > self._losses_border[-validation_interval]):
1039         print('Regularization is too strong reducing factor by /5.')
1040         self._regularization_coef = self._regularization_coef/5
1041
1042     # Rotate batches.
1043     if (auto_batch_rotation != 0):
1044         new_sample = myDataSets(training_batch_size = training_batch_size,
1045                                 border_training_batch_size = border_training_batch_size,
1046                                 validation_batch_size = validation_batch_size,
1047                                 input_dim = self._input_dim,
1048                                 method = method,
1049                                 domain = domain_type,
1050                                 border = border_type,
1051                                 seed = None)
1052
1053         training_batch_size, border_training_batch_size, validation_batch_size,\
1054         training_set, border_training_set, validation_set,\
1055         method, domain_type, border_type = self.use_training_sets(data_set = data_set)
1056
1057     # Plot loss every specified number of epochs while training.
1058
1059     if (self._num_epochs_trained % plot_interval == 0):
1060         real_start_epoch = self._num_epochs_trained - num_epochs
1061         auxiliaryPlotting.plot_loss_function(model = self,
1062                                             init_range = self._num_epochs_trained-plot_interval+1,
1063                                             end_range = -1,
1064                                             subdivide_losses = False,
1065                                             use_log_scale = False)
1066
1067     end_time = time.perf_counter()
1068     print("Execution Time:", end_time - start_time)
1069
1070     #####33
1071     # A REVISAR
1072     #####33
1073
1074     def save(self, filepath=''):
1075
1076         config_list = []
1077
1078         config_list.append(self._name)
1079         config_list.append(self._num_epochs_trained)
1080         config_list.append(self._losses_from_domain)
1081         config_list.append(self._losses_from_border)
1082         config_list.append(self._losses_from_regularization)
1083         config_list.append(self._losses)
1084         config_list.append(self._losses_wrt_solution)
1085         config_list.append(self._losses_wrt_validation)
1086
1087         config_list.append(self._input_dim)
1088         config_list.append(self._hidden_dim)
1089         config_list.append(self._output_dim)
1090         config_list.append(self._activation)
1091         config_list.append(self._weight_initializers)
1092         config_list.append(self._bias_initializers)
1093         config_list.append(self._batch_normalization)
1094
1095         config_list.append(self._trainable_weights)
1096
1097         config_list.append(self._learning_rate)
1098         config_list.append(self._scale_factor)
1099         config_list.append(self._regularization_coef)
1100
1101         config_list.append(self._differential_operator)
1102         config_list.append(self._external_force)
1103         config_list.append(self._exact_solution)
1104         config_list.append(self._loss_fuction)
1105         config_list.append(self._regularization)
1106         config_list.append(self._clip_gradient)
1107         config_list.append(self._optimizer_selection)
1108
1109         pickle_out = open(filepath + self.name + ".pickle", "wb")
1110         pickle.dump(config_list, pickle_out)
1111         pickle_out.close()
1112         files.download(self.name + ".pickle")
1113
1114     def load(filepath='', filename=''):

```

```

1115
1116 pickle_in = open(filepath + filename + ".pickle", "rb")
1117 config_list = pickle.load(pickle_in)
1118
1119 ld_name = config_list[0]
1120 model = myModel(ld_name)
1121 model._num_epochs_trained = config_list[1]
1122 model._losses_from_domain = config_list[2]
1123 model._losses_from_border = config_list[3]
1124 model._losses_from_regularization = config_list[4]
1125 model._losses = config_list[5]
1126 model._losses_wrt_solution = config_list[6]
1127 model._losses_wrt_validation = config_list[7]
1128
1129 ld_input_dim = config_list[8]
1130 ld_hidden_dim = config_list[9]
1131 ld_output_dim = config_list[10]
1132 ld_activation = config_list[11]
1133 ld_weight_initializers = config_list[12]
1134 ld_bias_initializers = config_list[13]
1135 ld_batch_normalization = config_list[14]
1136 model.build(ld_input_dim, ld_hidden_dim, ld_output_dim,
1137            ld_activation, ld_weight_initializers, ld_bias_initializers,
1138            ld_batch_normalization, None)
1139
1140 layer_copy = model._layers.copy()
1141 model._trainable_weights = config_list[15]
1142 model._layers = layer_copy
1143
1144 for layer_num in range(0, model._num_hidden_layers+2):
1145     model._layers[layer_num]._trainable_weights = model._trainable_weights[2*layer_num:2*(layer_num+1)]
1146     model._layers[layer_num].w = model._layers[layer_num]._trainable_weights[0]
1147     model._layers[layer_num].b = model._layers[layer_num]._trainable_weights[1]
1148
1149 ld_learning_rate = config_list[16]
1150 ld_scale_factor = config_list[17]
1151 ld_regularization_coef = config_list[18]
1152 ld_differential_operator = config_list[19]
1153 ld_external_force = config_list[20]
1154 ld_exact_solution = config_list[21]
1155 ld_loss_fuction = config_list[22]
1156 ld_regularization = config_list[23]
1157 ld_clip_gradient = config_list[24]
1158 ld_optimizer_selection = config_list[25]
1159
1160 model.compile(learning_rate = ld_learning_rate, scale_factor = ld_scale_factor,
1161              differential_operator = ld_differential_operator,
1162              external_force = ld_external_force,
1163              exact_solution = ld_exact_solution, loss_fuction = ld_loss_fuction,
1164              regularization = ld_regularization, regularization_coef = ld_regularization_coef,
1165              clip_gradient = ld_clip_gradient, optimizer_selection = ld_optimizer_selection)
1166
1167 return model
1168
1169 def save_weights(self, filepath=''):
1170
1171     pickle_out = open(filepath + self.name + "_weights.pickle", "wb")
1172     pickle.dump(self._trainable_weights, pickle_out)
1173     pickle_out.close()
1174
1175 def set_weights(self, filepath='', filename=''):
1176
1177     pickle_in = open(filepath + filename + ".pickle", "rb")
1178     self._trainable_weights = pickle.load(pickle_in)

```

B.8 execution Cell

```

1 """
2 @author: Alberto Garcia Molina
3 @latest_update: 12/10/2020
4 """
5
6 # Creates the training/validation set.
7 sample = myDataSets(training_batch_size = 10000,
8                     border_training_batch_size = 10000,
9                     validation_batch_size = 1000,
10                    input_dim = 1,
11                    method = 'uniform-hit-collocation',
12                    domain = 'hypercube-0-1', #'quarter-hypercube-0-1' #'hypercube-0-2'
13                    border = 'hypercube-0-1', #'side-x_1-y_0' #'side-x_1-y_0_expanded' #'two_sides-x_0-y_0'
14                    seed = 1993)
15
16 #Initializes the model.
17 model = myModel('model_name')
18 model.build(input_dim = 1,
19            hidden_dim = [3,4,3],
20            output_dim = 1,
21            activations = 'sigmoid', #'tanh' #'swish' #'exponential' #'softplus' #'relu' #'elu'
22            weight_initializers = 'xavier', #'he' #'normal_0_1' #'uniform_-1_1' #'zeros' #'ones'
23            bias_initializers = 'xavier', #'he' #'normal_0_1' #'uniform_-1_1' #'zeros' #'ones'
24            batch_normalization = False, # Actually layer normalization.
25            suppress_bias = False, # Make the neurons have no bias terms.
26            seed = 1993,
27            epsilon = 1e-12)
28 model.compile(differential_operator = 'Constant', #'Divergence' #'Advection' #'Laplacian' #'Clairaut'
29              external_force = 'Force_Constant', #'Force_Divergence' #'Force_Advection' #'Force_Laplacian' #'Force_Laplacian'
30              exact_solution = 'Sol_Polynomial_2Deg_1D_1D', #'Sol_Polynomial_2Deg_1D_1D' #'Sol_Polynomial_2Deg_2D_1D' #'

```

```

31         Sol_Polynomial_2Deg_2D_2D'
32         optimizer_selection = 'Adam', #'AMSGrad' #'Nadam' #'AdaGrad' #'AdaDelta' #'RMSProp' #'Vanilla_SGD' #'Momentum_SGD'
33         #'Nesterov_SGD' #'BFGS' #'L-BFGS'
34         learning_rate = 1e-02,
35         epsilon = 1e-31,
36         scale_factor = 1,
37         loss_fuction = 'square_L2_error', #'L2_error' #'square_L2_error' #'absolute_error' #'square_MSE' #'
38         square_L2_error_1st_comp'
39         regularization = None, #'Gradient_Type' #'Tikhonov'
40         regularization_coef = 0.1,
41         clip_gradient = 'global' #'value' 'norm'
42     )
43 # Trains the model.
44 model.fit (data_set = sample,
45           num_epochs = 15000,
46           validation_interval = 1000,
47           plot_interval = 500,
48           stagnation_stop_tol = 1e-16,
49           split_gradient = True, # Has to always be set to True
50           display_gradient_norm = False,
51           normalize_gradient = False,
52           train_only_domain = False,
53           train_only_border = False,
54           dual_training = False, # Experimental, do not use. Set to False.
55           auto_batch_rotation = False # Automatically generates new collocations with the same set-up as sample every
56           plot_interval iterations.
57 )
58 # Optional Plots:
59 auxiliaryPlotting.plot_loss_function(model = model,
60                                     init_range = 0,
61                                     end_range = -1,
62                                     subdivide_losses = True,
63                                     use_log_scale = True)
64 auxiliaryPlotting.plot_model(data_set = sample,
65                              model = model,
66                              plot_real_sol = True)
67 auxiliaryPlotting.plot_error (data_set = sample,
68                              model = model)
69 auxiliaryPlotting.plot_loss_comparison (models = [model1, model2, model3],
70                                       names = [model1_name, model3_name, model3_name],
71                                       title = ['loss of model 1', 'loss of model 2', 'loss of model 3'],
72                                       init_range = 0,
73                                       end_range = -1,
74                                       subdivide_losses = True,
75                                       use_log_scale = True)
76 # save/load options:
77
78 model.save(filepath='C//...')
79 model.load(filepath='C//...')

```

Bibliography

- [1] R. L. Bishop and S. I. Goldberg, *Tensor analysis on manifolds*. Dover, Dec. 1980.
- [2] C. Reisinger and G. Wittum, “Efficient hierarchical approximation of high-dimensional option pricing problems,” *SIAM J. Sci. Comput.*, vol. 29, pp. 440–458, 2007.
- [3] L. Grzelak, J. Witteveen, M. Suárez-Taboada, and C. Oosterlee, “The stochastic collocation Monte Carlo sampler: Highly efficient sampling from ‘expensive’ distributions,” *Quantitative Finance*, vol. 19, pp. 339–356, 2019.
- [4] J. Sirignano and K. Spiliopoulos, “DGM: A deep learning algorithm for solving partial differential equations,” *Journal of Computational Physics*, vol. 375, p. 1339–1364, Dec 2018.
- [5] W. Hughes, J. H. Merkin, and R. Sturman, “Analytic solutions of partial differential equations,” 2003/4. MATH3414, School of Mathematics, University of Leeds, lectures notes.
- [6] S. Kepley and T. Zhang, “A constructive proof of the Cauchy-Kovalevskaya theorem for ordinary differential equations.,” *arXiv: Classical Analysis and ODEs*, 2019.
- [7] D. G. Gaidashev, “CHAPTER 2 The Cauchy-Kovalevskaya Theorem.” <http://www2.math.uu.se/~gaidash/1MA216/CK.pdf>. Accessed: 2020-09-01.
- [8] A. Meade and A. Fernandez, “The numerical solution of linear ordinary differential equations by feedforward neural networks,” *Mathematical and Computer Modelling*, vol. 19, no. 12, pp. 1 – 25, 1994.
- [9] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, “Neural ordinary differential equations,” *ArXiv*, vol. abs/1806.07366, 2018.
- [10] M. Jaderberg, W. Czarnecki, S. Osindero, O. Vinyals, A. Graves, D. Silver, and K. Kavukcuoglu, “Decoupled neural interfaces using synthetic gradients,” *ArXiv*, vol. abs/1608.05343, 2017.
- [11] W. E. J. Han, and A. Jentzen, “Deep learning-based numerical methods for high-dimensional parabolic partial differential equations and backward stochastic differential equations,” *Communications in Mathematics and Statistics*, vol. 5, p. 349–380, Nov 2017.
- [12] C. Beck, E. Weinan, and A. Jentzen, “Machine learning approximation algorithms for high-dimensional fully nonlinear partial differential equations and second-order backward stochastic differential equations,” *Journal of Nonlinear Science*, vol. 29, pp. 1563–1619, Aug. 2019.
- [13] J. Berg and K. Nyström, “A unified deep artificial neural network approach to partial differential equations in complex geometries,” *Neurocomputing*, vol. 317, pp. 28–41, 2018.
- [14] T. Dockhorn, “A discussion on solving partial differential equations using neural networks,” *ArXiv*, vol. abs/1904.07200, 2019.
- [15] S. H. Kolluru, “A neural network based method to solve boundary value problems,” *ArXiv*, vol. abs/1909.11082, 2019.
- [16] I. Lagaris, A. Likas, and D. Papageorgiou, “Neural-network methods for boundary value problems with irregular boundaries,” *IEEE transactions on neural networks*, vol. 11 5, pp. 1041–9, 2000.
- [17] I. Lagaris, A. Likas, and D. Fotiadis, “Artificial neural networks for solving ordinary and partial differential equations,” *IEEE transactions on neural networks*, vol. 9 5, pp. 987–1000, 1998.
- [18] A. Al-Arabi, A. Correia, D. Naiff, G. Jardim, and Y. F. Saporito, “Solving nonlinear and high-dimensional partial differential equations via deep learning,” *arXiv: Computational Finance*, 2018.
- [19] L. Ruthotto and E. Haber, “Deep neural networks motivated by partial differential equations,” *Journal of Mathematical Imaging and Vision*, vol. 62, pp. 352–364, 2019.
- [20] P. Chaudhari, A. M. Oberman, S. Osher, S. Soatto, and G. Carlier, “Deep relaxation: partial differential equations for optimizing deep neural networks,” *Research in the Mathematical Sciences*, vol. 5, pp. 1–30, 2017.

-
- [21] M. Raissi and G. Karniadakis, “Hidden physics models: Machine learning of nonlinear partial differential equations,” *ArXiv*, vol. abs/1708.00588, 2018.
- [22] M. Raissi, P. Perdikaris, and G. Karniadakis, “Physics informed deep learning (Part I): Data-driven solutions of nonlinear partial differential equations,” *ArXiv*, vol. abs/1711.10561, 2017.
- [23] M. Raissi, P. Perdikaris, and G. Karniadakis, “Physics informed deep learning (Part II): Data-driven discovery of nonlinear partial differential equations,” *ArXiv*, vol. abs/1711.10566, 2017.
- [24] P. Simard, B. Victorri, Y. LeCun, and J. Denker, “Tangent prop - a formalism for specifying selected invariances in an adaptive network,” in *NIPS*, 1991.
- [25] W. Czarnecki, S. Osindero, M. Jaderberg, G. Swirszcz, and R. Pascanu, “Sobolev training for neural networks,” in *NIPS*, 2017.
- [26] R. Grosse, “CSC321 Lecture 10: Automatic Differentiation.” https://www.cs.toronto.edu/~rgrosse/courses/csc321_2018/slides/lec10.pdf. Accessed: 2020-09-01.
- [27] A. G. Baydin, B. Pearlmutter, A. Radul, and J. Siskind, “Automatic differentiation in machine learning: a survey,” *ArXiv*, vol. abs/1502.05767, 2017.
- [28] S. Sun, Z. Cao, H. Zhu, and J. Zhao, “A survey of optimization methods from a machine learning perspective,” *IEEE Transactions on Cybernetics*, vol. 50, pp. 3668–3681, 2020.
- [29] A. Gaier and D. R. Ha, “Weight agnostic neural networks,” in *NeurIPS*, 2019.
- [30] S. Ruder, “An overview of gradient descent optimization algorithms,” *ArXiv*, vol. abs/1609.04747, 2016.
- [31] Y. Li, C. Wei, and T. Ma, “Towards explaining the regularization effect of initial large learning rate in training neural networks,” *ArXiv*, vol. abs/1907.04595, 2019.
- [32] L. N. Smith and N. Topin, “Super-convergence: Very fast training of neural networks using large learning rates,” *ArXiv*, vol. abs/1708.07120, 2017.
- [33] P. Chaudhari, A. Choromanska, S. Soatto, Y. LeCun, C. Baldassi, C. Borgs, J. Chayes, L. Sagun, and R. Zecchina, “Entropy-SGD: Biasing gradient descent into wide valleys,” *ArXiv*, vol. abs/1611.01838, 2016.
- [34] I. Sutskever, J. Martens, G. Dahl, and G. Hinton, “On the importance of initialization and momentum in deep learning,” in *Proceedings of the 30th International Conference on Machine Learning* (S. Dasgupta and D. McAllester, eds.), vol. 28 of *Proceedings of Machine Learning Research*, (Atlanta, Georgia, USA), pp. 1139–1147, PMLR, 17–19 Jun 2013.
- [35] J. C. Duchi, E. Hazan, and Y. Singer, “Adaptive subgradient methods for online learning and stochastic optimization.,” *J. Mach. Learn. Res.*, vol. 12, pp. 2121–2159, 2011.
- [36] M. D. Zeiler, “Adadelta: An adaptive learning rate method,” *ArXiv*, vol. abs/1212.5701, 2012.
- [37] G. Hinton, “Lecture 6a. Overview of mini-batch gradient descent.” https://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf. Accessed: 2020-09-01.
- [38] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” 2014. cite arxiv:1412.6980 Comment: Published as a conference paper at the 3rd International Conference for Learning Representations, San Diego, 2015.
- [39] S. J. Reddi, S. Kale, and S. Kumar, “On the convergence of Adam and beyond.,” in *ICLR*, OpenReview.net, 2018.
- [40] T. Dozat, “Incorporating nesterov momentum into Adam,” 2016.
- [41] J. Zhang, I. Mitliagkas, and C. Ré, “Yellowfin and the art of momentum tuning,” *ArXiv*, vol. abs/1706.03471, 2019.
- [42] J. Rafati and R. F. Marcia, “Quasi-Newton optimization methods for deep learning applications.,” *CoRR*, vol. abs/1909.01994, 2019.
- [43] J. Martínez, “Practical quasi-Newton methods for solving nonlinear systems,” *Journal of Computational and Applied Mathematics*, vol. 124, pp. 97–121, 2000.
- [44] J. Nocedal and S. J. Wright, “Numerical optimization (springer series in operations research and financial engineering),” 2000.

-
- [45] R. Bollapragada, D. Mudigere, J. Nocedal, H.-J. M. Shi, and P. Tang, “A progressive batching L-BFGS method for machine learning,” in *ICML*, 2018.
- [46] J. Martens, “Deep learning via hessian-free optimization.,” in *ICML* (J. Fürnkranz and T. Joachims, eds.), pp. 735–742, Omnipress, 2010.
- [47] P. Ramachandran, B. Zoph, and Q. V. Le, “Searching for activation functions,” *ArXiv*, vol. abs/1710.05941, 2018.
- [48] S. Elfving, E. Uchibe, and K. Doya, “Sigmoid-weighted linear units for neural network function approximation in reinforcement learning,” *Neural networks : the official journal of the International Neural Network Society*, vol. 107, pp. 3–11, 2018.
- [49] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proceedings of the 27th International Conference on Machine Learning (ICML-10)* (J. Fürnkranz and T. Joachims, eds.), pp. 807–814, 2010.
- [50] X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neural networks,” vol. 9 of *Proceedings of Machine Learning Research*, (Chia Laguna Resort, Sardinia, Italy), pp. 249–256, JMLR Workshop and Conference Proceedings, 13–15 May 2010.
- [51] K. He, X. Zhang, S. Ren, and J. Sun, “Delving deep into rectifiers: Surpassing human-level performance on imagenet classification,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, December 2015.
- [52] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov, “Dropout: A simple way to prevent neural networks from overfitting,” *Journal of Machine Learning Research*, vol. 15, no. 56, pp. 1929–1958, 2014.
- [53] J. Ba, J. Kiros, and G. E. Hinton, “Layer normalization,” *ArXiv*, vol. abs/1607.06450, 2016.
- [54] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle, “Greedy layer-wise training of deep networks,” in *NIPS*, 2006.
- [55] A. Andoni, R. Panigrahy, G. Valiant, and L. Zhang, “Learning polynomials with neural networks,” in *ICML*, 2014.
- [56] Z.-Q. J. Xu, “A note of using Tensorflow to code Laplacian operator in high dimension.” <https://ins.sjtu.edu.cn/people/xuzhiqin/pub/laplaciancode.pdf>. Accessed: 2020-09-01.