

# UNIVERSITÀ DELLA CALABRIA



*Department of Mathematics and Computer Science*

---

*Master's degree in Mathematics*

## ***Neural Networks in reaction-diffusion models***

### **SUPERVISORS**

Prof. Giovanni Mascali

Prof. Gennaro Infante

### **CANDIDATE**

Riccardo Scida

M. N. 235617

---

*Academic year 2023-2024*

# Table of Contents

1. Chapter 1: Partial Differential Equations (PDEs).....	10
1.1. Introduction on PDEs.....	10
1.1.1. Definition of PDEs.....	10
1.1.2. Multi-index Notation.....	10
1.1.3. Types of PDEs.....	11
1.1.4. Definition of system of PDEs.....	11
1.1.5. Examples of PDEs.....	11
1.2. Linear differential operators.....	12
1.2.1. Definition of linear differential operators.....	12
1.2.2. Boundary conditions.....	13
1.2.3. Dirac Delta.....	13
1.2.4. Definition of Green function.....	14
1.2.5. Proposition (Green function as solution of the Dirichlet problem).....	15
1.2.6. Proposition (Building a Green function).....	15
1.3. Laplace-Poisson equations.....	16
1.3.1. Definition of Laplace and Poisson equations.....	16
1.3.2. Definition of harmonic function.....	16
1.3.3. Physical interpretation.....	16
1.3.4. Boundary conditions.....	17
1.3.5. Theorem (Green).....	17
1.3.6. n-dimensional Dirac Delta.....	17
1.3.7. Solution of the Dirichlet problem of the Laplacian.....	18
1.4. Heat equation.....	19
1.4.1. Definition of Heat equation.....	19

1.4.2. Definition of parabolic boundary.....	19
1.4.3. Proposition 1 (about the Weak Maximum principle).....	19
1.4.4. Proposition 2 (about the Weak Maximum principle).....	20
1.4.5. Theorem (Weak Maximum principle).....	21
1.4.6. Definition of subsolutions and supersolutions.....	21
1.4.7. Theorem (uniqueness of the solution).....	21
1.4.8. Theorem (Strong Maximum principle).....	22
1.4.9. Fundamental solution of the heat equation ( $n=1$ ).....	22
1.4.10. Fundamental solution of the heat equation ( $n>1$ ).....	25
1.4.11. The global Cauchy problem ( $n=1$ ).....	25
1.4.12. Theorem (Existence of Solution).....	26
1.4.13. Theorem (existence for the non-homogeneous Cauchy problem).....	27
1.4.14. Theorem (principle of global maximum).....	28
 1.5. Scalar reaction-diffusion equations.....	29
1.5.1. Introduction to scalar reaction-diffusion equations.....	29
1.5.2. Theorem (Maximum principle).....	30
1.5.3. Theorem (of comparison).....	30
1.5.4. Global existence for the Cauchy problem.....	30
1.5.5. Global existence in the case of systems.....	31
 1.6. Reaction-diffusion models for a biological species.....	32
1.6.1. Reaction rate.....	32
1.6.2. Fisher–Kolmogorov–Petrovsky–Piskunov Equation.....	33
1.6.3. Definition of monostable function.....	34
1.6.4. Proposition (reaction-diffusion equation with monostable reaction rate).....	34
1.6.5. Properties (boundedness of the solution).....	34
1.6.6. Stationary solutions.....	35
 1.7. Travelling waves.....	36
1.7.1. Definition of travelling waves.....	36

1.7.2. Theorem (reaction-diffusion equation with monostable reaction rate admits travelling waves).....	36
1.7.3. Logistic growth term.....	37
1.7.4. Asymptotic approximation.....	37
1.7.5. Stability.....	38
 1.8. Systems of reaction-diffusion equations.....	39
1.8.1. Definition of conservative system.....	39
1.8.2. Lotka – Volterra system.....	39
1.8.3. Structural instability.....	42
1.8.4. Theorem (Diffusion).....	43
1.8.5. Definition (Boundary operator).....	44
1.8.6. Definition (Invariant set).....	44
1.8.7. Definition (Almost monotonicity).....	45
1.8.8. Theorem (Comparison generalized).....	45
1.8.9. Corollary (Uniqueness).....	47
1.8.10. Corollary (Invariant set).....	47
1.8.11. Theorem (Global existence).....	48
1.8.12. Definition (Closed operator).....	48
1.8.13. Corollary (of the comparison theorem).....	48
1.8.14. Lemma (on the almost monotonicity).....	49
1.8.15. Theorem (invariant set).....	49
1.8.16. Theorem (uniqueness of the solution).....	50
 1.9. Compartmental models.....	51
1.9.1. Diffusion of rabies.....	51
1.9.2. Model improvement.....	53
1.9.3. Spatial distribution patterns of grey and red squirrels.....	57
1.9.4. Spatial pattern formation.....	59
1.9.5. Definition (Turing instability).....	61
1.9.6. Necessary and sufficient condition for Turing instability....	61
1.9.7. Theorem (necessary condition).....	63
1.9.8. Theorem (sufficient condition).....	64

1.9.9. Detailed analysis of pattern formation.....	64
1.9.10. Pattern formation in growing domains.....	68
1.9.11. Hints of pattern formation on the coat of some mammals...	68
1.9.12. Non-existence of spatial patterns in reaction-diffusion systems.....	70
1.9.13. Numerical methods.....	71
<b>2. Chapter 2: Neural Networks for solving reaction-diffusion equations.....</b>	<b>73</b>
2.1. Machine Learning.....	73
2.1.1. ML algorithms.....	73
2.1.2. Overfitting.....	74
2.2. Neural Networks.....	75
2.2.1. Definition of Neural Networks.....	75
2.2.2. Perceptrons.....	75
2.2.3. Multilayer perceptron (MLP).....	77
2.2.4. Sigmoid function.....	78
2.2.5. Hyperbolic Tangent function.....	79
2.2.6. ReLU function.....	79
2.2.7. How MLP works.....	81
2.2.8. General description of a Neural Network.....	82
2.2.9. Definition (Neural Networks as composition of functions).....	82
2.3. Neural Network learning.....	83
2.3.1. Example (Error function).....	83
2.3.2. Steepest descent method (SDM).....	83
2.3.3. Learning scheme, the toy network.....	84
2.3.4. Learning scheme, the general case.....	85
2.3.5. Backpropagation (BP).....	86
2.3.6. Algorithm (Backpropagation and SDM method).....	87

2.3.7. Definition of Loss Function.....	87
2.3.8. Gradient Descent method.....	88
2.3.9. Automatic Differentiation.....	89
 2.4. The Universal Function Approximation Theorem for Neural Networks.....	90
2.4.1. Introduction.....	90
2.4.2. Theorem (Universal function approximation theorem for Neural Networks).....	92
2.4.3. Definition of ridge function.....	92
2.4.4. Theorem (Vostrecov, Kreines).....	93
2.4.5. Proposition 1 (about density).....	94
2.4.6. Proposition 2 (about density).....	94
2.4.7. Proposition 3 (about density).....	95
2.4.8. Application.....	97
 2.5. Neural Network for Solving Differential Equations.....	98
2.5.1. A first approach.....	98
2.5.2. Neural Network for PDEs.....	100
2.5.3. A different approach.....	101
2.5.4. Mixed boundary conditions.....	103
2.5.5. Weight space.....	103
2.5.6. Monte Carlo method.....	104
 2.6. Neural Network methods for solving Reaction Diffusion models..	105
2.6.1. Introduction.....	105
2.6.2. Calculation of derivatives.....	107
2.6.3. Calculation of $\frac{\partial L^2}{\partial \omega_{ij}^{(k)}}, \frac{\partial L^2}{\partial \theta_i^{(k)}}.....$	110
2.6.4. Recursive representation.....	114
2.6.5. How it works.....	116
2.6.6. Updating weights process.....	116
2.6.7. Algorithm (Backpropagation).....	119

<b>2.7. Application to Fisher – Kolmogorov – Petrovsky - Piskunov</b>	
Equation.....	121
2.7.1. Introduction.....	121
2.7.2. Spatially homogeneous case.....	121
2.7.3. Matlab implementation.....	123
2.7.4. Matlab codes description.....	123
2.7.5. Case with spatially homogeneous initial conditions.....	128
2.7.6. Matlab implementation.....	130
2.7.7. Matlab codes description.....	130
2.7.8. Case with non-homogeneous initial conditions.....	136
2.7.9. Matlab implementation.....	137
2.7.10. Matlab codes description.....	138
<b>3. Appendix A: Application to the approximation theorem.....</b>	<b>146</b>
<b>4. Appendix B: Matlab codes in the spatially homogeneous case.....</b>	<b>155</b>
<b>5. Appendix C: Matlab codes in the general case with spatially homogeneous initial conditions.....</b>	<b>168</b>
<b>6. APPENDIX D: Matlab codes in the general case with sinusoidal initial conditions.....</b>	<b>181</b>

# **Introduction**

In recent years, the development of Neural Networks (NN), inspired by the structure and functioning of the human brain, has been a major revolution in the fields of Artificial Intelligence (AI) and Machine Learning (ML). These disciplines represent the most revolutionary tools in the contemporary scientific and technological landscape. They are concerned with building algorithms capable of performing very complex tasks, such as image recognition, sound recognition, and various autonomous decision-making processes. In doing so, Artificial Intelligence has found application in a wide range of fields, transforming the way many industries operate and improving the quality of people's lives. Some examples are:

1. in the field of health care, AI systems have been designed that can analyse medical images and clinical data so as to support doctors in diagnosing diseases such as cancer or certain heart diseases;
2. in the field of finance, AI systems capable of monitoring financial transactions have been designed to identify suspicious activities and prevent fraud;
3. in the field of transportation, some companies, such as Tesla, are using AI to develop self-driving cars that can sense their surroundings and make driving decisions, including optimizing routes and thus, for example, reducing delivery times and costs;
4. in the field of E-commerce, some platforms, such as Amazon and Netflix, use AI to analyse user preferences and recommend relevant products or content;
5. in the field of video games, AI is used to create "intelligent" and dynamic non-player characters (NPCs), making the gaming experience more realistic.

These examples are just a few of the many applications of Artificial Intelligence, and they demonstrate how this technology is becoming, more and more, a part of our daily lives and work.

In our work we will exploit a very important property of Neural Networks, that is their ability to approximate, under appropriate assumptions, any continuous function defined on  $\mathbb{R}^n$ . This led to the idea of using Neural Networks to numerically approximate the solutions of partial differential equations (PDEs). In particular, we will consider reaction-diffusion equations. PDEs occupy a central role in the creation of models describing the evolution of complex systems, as a function of space and time, such as physical, engineering, biological systems, etc., and have various applications, including fluid dynamics, wave propagation, heat diffusion, and material deformation. This combination of Neural Networks and PDEs represents a promising research frontier and opens up new possibilities for tackling problems previously considered "inaccessible" or excessively expensive to solve by traditional methods. Our goal, in this thesis, precisely, will be to analyse what has already been studied in the field of Neural Networks to solve ordinary differential equations (ODEs) and some PDEs, and then to deepen these studies also in the case of reaction-diffusion models. In particular, after developing these concepts from the theoretical point of view, we will try to build a "bridge" that will allow us to implement in Matlab algorithms to approximate, via Neural Networks, some specific models, such as the Fisher-Kolmogorov-Petrovsky-Piskunov model. The advantages and disadvantages of these tools will be discussed, with the aim of providing an overview of possible future research directions. Thus, an attempt will be made to demonstrate how Neural Networks can be an effective and innovative tool to address the study of reaction-diffusion patterns, paving the way for new applications and improvements in various fields of science.

# Chapter 1

## Partial Differential Equations (PDEs)

### 1.1 Introduction on PDEs

#### 1.1.1 Definition of PDEs

Let  $k \geq 1$  be a non-negative integer and let  $U \subseteq \mathbb{R}^n$  be an open set, a *partial differential equation* [2] of order  $k$  is an equation of the form

$$(1.1) \quad F(D^k u(x), D^{k-1} u(x), \dots, u(x), x) = 0,$$

where

- $x \in U$ ;
- $F : \mathbb{R}^{n^k} \times \mathbb{R}^{n^{k-1}} \times \dots \times \mathbb{R}^n \times \mathbb{R} \times U \rightarrow \mathbb{R}$  is a given function;
- $u : U \rightarrow \mathbb{R}$  is the unknown.

In this definition we have used the

#### 1.1.2 Multi-index Notation

A vector of the form  $\alpha = (\alpha_1, \dots, \alpha_n)$ , where each component  $\alpha_i$  is a non-negative integer, is called a *multi-index* of order

$$|\alpha| = \alpha_1 + \dots + \alpha_n.$$

Given a multi-index  $\alpha = (\alpha_1, \dots, \alpha_n)$ , define

$$D^\alpha u(x) := \frac{\partial^{|\alpha|} u(x)}{\partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n}} = \partial x_1^{\alpha_1} \dots \partial x_n^{\alpha_n} u.$$

If  $k$  is a non-negative integer, we can define with

$$D^k u(x) := \{ D^\alpha u(x) \mid |\alpha| = k \}$$

the set of all partial derivatives of order  $k$ .

### 1.1.3 Types of PDEs

A PDE is called:

1. *linear*, if it's of the form

$$\sum_{|\alpha| \leq k} a_\alpha(x) D^\alpha u = f(x).$$

This PDE is said to be *homogeneous* if  $f \equiv 0$ ;

2. *semi-linear*, if it's of the form

$$\sum_{|\alpha|=k} a_\alpha(x) D^\alpha u + a_0(D^{k-1}u, \dots, Du, u, x) = 0;$$

3. *quasilinear*, if it's of the form

$$\sum_{|\alpha|=k} a_\alpha(D^{k-1}u, \dots, Du, u, x) D^\alpha u + a_0(D^{k-1}u, \dots, Du, u, x) = 0;$$

4. *fully not linear*, if it depends non-linearly upon the highest order derivatives.

### 1.1.4 Definition of system of PDEs

Let  $k \geq 1$  be an integer and let  $U \subseteq \mathbb{R}^n$  be an open set, a *system of partial differential equations* of order  $k$  is an equation of the form

$$(1.2) \quad \mathbf{F}(D^k \mathbf{u}(x), D^{k-1} \mathbf{u}(x), \dots, \mathbf{u}(x), x) = \mathbf{0},$$

where:

- $x \in U$ ;
- $\mathbf{F}: \mathbb{R}^{mn^k} \times \mathbb{R}^{mn^{k-1}} \times \dots \times \mathbb{R}^{mn} \times \mathbb{R}^m \times U \rightarrow \mathbb{R}^m$  is a given function;
- $\mathbf{u}: U \rightarrow \mathbb{R}^m$ ,  $\mathbf{u} = (u^1, \dots, u^m)$  is the unknown.

### 1.1.5 Examples of PDEs

The following equations are linear:

- *Laplace's equation*

$$\Delta u = 0.$$

- *Eigenvalue equation*

$$-\Delta u = \lambda u.$$

- *Linear transport equation*

$$u_t + \sum_{i=1}^n b^i u_{x_i} = 0.$$

- *Liouville's equation*

$$u_t - \sum_{i=1}^n (b^i u)_{x_i} = 0.$$

- *Heat equation*

$$u_t - \Delta u = 0.$$

- *Wave equation*

$$u_{tt} - \Delta u = 0.$$

The following equations are nonlinear:

- *Eikonal equation*

$$|Du| = 1.$$

- *Nonlinear Poisson equation*

$$-\Delta u = f(u).$$

- *Scalar reaction-diffusion equation*

$$u_t - \Delta u = f(u).$$

- *Nonlinear wave equation*

$$u_{tt} - \Delta u = f(u).$$

The following are linear systems:

- *Equilibrium equations of linear elasticity*

$$\mu \Delta \mathbf{u} + (\lambda + \mu) D(\operatorname{div} \mathbf{u}) = 0.$$

- *Evolution equations of linear elasticity*

$$\mathbf{u}_{tt} - \mu \Delta \mathbf{u} - (\lambda + \mu) D(\operatorname{div} \mathbf{u}) = 0.$$

The following are nonlinear systems:

- *System of conservation laws*

$$\mathbf{u}_t + \operatorname{div} \mathbf{F}(\mathbf{u}) = 0.$$

- *reaction-diffusion system*

$$\mathbf{u}_t - \Delta(\mathbf{u}) = \mathbf{f}(\mathbf{u}).$$

## 1.2 Linear differential operators

### 1.2.1 Definition of linear differential operators

A linear *differential operator* [7] of order  $n$  is a linear operator  $A$  acting on functions  $u \in C^n$  in the ordinary way:

$$(1.3) \quad Au = \sum_{j=0}^n a_j u^{(j)},$$

where

- $u^{(j)} = \frac{d^j u}{dx^j}$ ,  $u^{(0)} = u(x)$ ,  $x \in I \subseteq \mathbb{R}$ ;
- $a_j$  are real-complex valued functions.

If  $n=2$  we have operators of the type

$$Au = au'' + bu' + cu,$$

where every function is regular on  $[0,1]$  e  $a(x) > 0$ ,  $\forall x \in [0,1]$ , where, without loss of generality, we have taken  $I = [0,1]$ .

### 1.2.2 Boundary conditions

We will study problems of the type

$$(1.4) \quad Au = f,$$

with boundary conditions. The most common are

- *Dirichlet conditions*  $\begin{cases} u(0) = a, & a, b \in \mathbb{R}; \\ u(1) = b, \end{cases}$
- *Neumann conditions*  $\begin{cases} u'(0) = c, & c, d \in \mathbb{R}; \\ u'(1) = d, \end{cases}$
- *periodic conditions*  $\begin{cases} u(0) = u(1), \\ u'(0) = u'(1), \end{cases}$
- *mixed conditions*  $\begin{cases} a_0 u(0) + b_0 u'(0) = \alpha, & \alpha, \beta \in \mathbb{R}. \\ a_1 u(1) + b_1 u'(1) = \beta, \end{cases}$

### 1.2.3 Dirac Delta

Roughly speaking, the Dirac Delta  $\delta$  is a “function” on  $\mathbb{R}$  with the following formal properties:

- $\delta(x) = 0$ ,  $\forall x \neq 0$ ;
- $\int_{-\infty}^{+\infty} \delta(x) dx = 1$ ;
- $\int_{-\infty}^{+\infty} \delta(y-x) f(y) dy = f(x)$ ;

$$\bullet \quad \int_{-\infty}^y \delta(x) dx = \begin{cases} 0, & \text{if } y < 0, \\ 1, & \text{if } y \geq 0. \end{cases}$$

Of course, precise meaning can be given to these results only within the framework of the theory of distributions, for example, see [4]. In this section, we have limited ourselves only to a heuristic discussion suitable for motivating the classical definition of

#### 1.2.4 Definition of Green function

In the framework of the classical theory, a function  $g:[0,1] \times [0,1] \rightarrow \mathbb{C}$  is said to be a *Green function* if:

- $g$  is a continuous function in the box  $[0,1] \times [0,1]$  and it is twice differentiable with respect to  $x$  in the triangles  $0 \leq x \leq y \leq 1$  and  $0 \leq y \leq x \leq 1$ ;
- $Ag=0$  in  $0 < x < y < 1$  and  $0 < y < x < 1$ ;
- $g(0, y) = g(1, y) = 0$ ;
- the “jump” of  $g_x$  across the line  $x=y$  is given by

$$g_x(y^+, y) - g_x(y^-, y) = \frac{1}{a(y)}.$$

#### Remark

Let  $A = a \frac{d^2}{dx^2} + b \frac{d}{dx} + c$  be a linear differential operator of order 2 and, without loss of generality, let's consider the homogeneous Dirichlet problem

$$(1.5) \quad \begin{cases} Au=f, \\ u(0)=u(1)=0, \end{cases}$$

with  $f : [0,1] \rightarrow \mathbb{C}$  continuous function. We are looking for solutions, if they exist, such as

$$(1.6) \quad u(x) = \int_0^1 g(x, y) f(y) dy,$$

where  $g : [0,1] \times [0,1] \rightarrow \mathbb{C}$  is called the *Green function*. Let us define, if it exists,

$$G : C([0,1]) \rightarrow D(A),$$

$$(1.7) \quad G(f(x)) = \int_0^1 g(x, y) f(y) dy,$$

as the inverse of  $A$ , such that to a function  $f$  it associates the function  $u$  which is solution of the Dirichlet problem. The function  $g(x, y)$  is also called the *kernel* of  $G$ . The Green function is the formal solution of the problem

$$(1.8) \quad \begin{cases} Ag(x, y) = \delta(x - y), \\ g(0, y) = g(1, y) = 0, \end{cases}$$

indeed, formally we have:

$$\begin{aligned} Au(x) &= \int_0^1 Ag(x, y) f(y) dy = \int_0^1 \delta(x - y) f(y) dy = f(x), \\ u(0) &= \int_0^1 g(0, y) f(y) dy = 0, \quad u(1) = \int_0^1 g(1, y) f(y) dy = 0. \end{aligned}$$

### 1.2.5 Proposition (Green function as solution of the Dirichlet problem)

Let  $A = a \frac{d^2}{dx^2} + b \frac{d}{dx} + c$  a linear differential operator with  $a, b, c \in C([0, 1])$ ,  $a \neq 0$ ,

If  $g$  is a Green function and  $f \in C([0, 1])$ , then

$$(1.9) \quad u(x) = G(f(x)) = \int_0^1 g(x, y) f(y) dy$$

is the solution of the Dirichlet problem

$$\begin{cases} Au = f, \\ u(0) = u(1) = 0. \end{cases}$$

### 1.2.6 Proposition (Building a Green function)

To construct the Green's function, we can choose two nonzero solutions  $v_1, v_2$  of the homogeneous case  $Au = 0$ , such that

$$v_1(0) = v_2(1) = 0.$$

There exists two solutions that are linearly independent if and only if the problem

$$\begin{cases} Au = 0, \\ u(0) = u(1) = 0, \end{cases}$$

has the unique solution  $u \equiv 0$ . In this case we can define a Green function as

$$(1.10) \quad g(x, y) = \begin{cases} C(y) v_1(x) v_2(y), & \text{for } 0 \leq x \leq y \leq 1, \\ C(y) v_1(y) v_2(x), & \text{for } 0 \leq y \leq x \leq 1, \end{cases}$$

where

- $C(y) = \frac{I}{a(y)W(y)}$ ;

- $W(y) = \begin{vmatrix} v_1 & v_2 \\ v_1' & v_2' \end{vmatrix} = v_1 v_2' - v_2 v_1' \neq 0$ , is said the *Wronskian* of  $v_1$  and  $v_2$ .

Let us note that

$$\begin{cases} Ag=0, \\ g(0,y)=g(1,y)=0, \forall y \in [0,1], \end{cases}$$

and moreover

$$\begin{aligned} g_x(y^+,y) - g_x(y^-,y) &= C(y)v_1(y^+)v_2'(y) - C(y)v_1'(y^-)v_2(y) \\ &= C(y)[v_1(y)v_2'(y) - v_1'(y)v_2(y)] \\ &= C(y)W(y) = \frac{I}{a(y)}. \end{aligned}$$

## 1.3 Laplace-Poisson equations

### 1.3.1 Definition of Laplace and Poisson equations

Let  $U \subseteq \mathbb{R}^n$  be a connected, open and bounded set,  $u: \bar{U} \rightarrow \mathbb{R}$  be the unknown and  $f: U \rightarrow \mathbb{R}$  be a given function. Let us define the *Laplace equation* [7] as

$$(1.11) \quad \Delta u = 0,$$

and the *Poisson equation* as

$$(1.12) \quad -\Delta u = f,$$

where the *Laplacian* of  $u$  is given by

$$\Delta u = \sum_{i=1}^n \frac{\partial^2 u}{\partial x_i^2}.$$

### 1.3.2 Definition of harmonic function

A function  $u \in C^2$  such that  $\Delta u = 0$  is called a *harmonic function*.

### 1.3.3 Physical interpretation

Laplace's equation is found in many physical contexts. Suppose that  $u$  is the density of a certain quantity in equilibrium and let  $\mathbf{F}: \bar{U} \rightarrow \mathbb{R}^n$  be a differentiable vector field, representing the flux of  $u$ , then by Gauss's Theorem:

$$(1.13) \quad \int_U \nabla \cdot \mathbf{F} d\mathbf{x} = \int_{\partial U} \mathbf{F} \cdot \mathbf{n} d\sigma.$$

We assume that the flux of  $u$  through  $\partial U$  is zero, i.e.

$$\int_{\partial U} \mathbf{F} \cdot \mathbf{n} d\sigma = 0,$$

where

- $\mathbf{n}$  is the normal versor to the surface  $\partial U$  oriented outward;
- $d\sigma$  is the area element of the surface.

hence

$$\nabla \cdot \mathbf{F} = 0, \text{ in } U.$$

In some physical examples, we can suppose that  $\mathbf{F}$  is proportional to the *gradient*  $\nabla u$ , but with opposite direction, i.e.

$$\mathbf{F} = -a \nabla u, \quad a > 0,$$

then substituting yields

$$\nabla \cdot \nabla u = \Delta u = 0.$$

### 1.3.4 Boundary conditions

For the Laplace and Poisson equations, the most common boundary conditions are

- $u = h$ , on  $\partial U$ , (*Dirichlet condition*);
- $\frac{\partial u}{\partial n} := \mathbf{n} \cdot \nabla u = h$ , on  $\partial U$ , (*Neumann condition*).

### 1.3.5 Theorem (Green)

Let's take  $u, v \in C^2(\overline{U})$ , then

$$(1.14) \quad \int_U (u \Delta v - v \Delta u) d\mathbf{x} = \int_{\partial U} \left( u \frac{\partial v}{\partial n} - v \frac{\partial u}{\partial n} \right) d\sigma.$$

### 1.3.6 n-dimensional Dirac delta

We can define the *n-dimensional Dirac delta* as a "function" with the following formal properties:

- $\delta(\mathbf{x}) = 0, \forall \mathbf{x} \neq 0;$
- $\int_{\mathbb{R}^n} \delta(\mathbf{x}) d\mathbf{x} = 1;$
- $\int_{\mathbb{R}^n} \delta(\mathbf{x} - \mathbf{y}) f(\mathbf{y}) d\mathbf{y} = f(\mathbf{x}), \forall f \in C(\mathbb{R}^n).$

The Green's function  $g(\mathbf{x}, \mathbf{y})$  associated with the Laplacian problem, with Dirichlet homogeneous boundary conditions, is solution of the problem

$$(1.15) \quad \begin{cases} -\Delta g(\mathbf{x}, \mathbf{y}) = \delta(\mathbf{x} - \mathbf{y}), & \mathbf{x} \in U, \mathbf{y} \in U, \\ g(\mathbf{x}, \mathbf{y}) = 0, & \mathbf{x} \in \partial U, \mathbf{y} \in U. \end{cases}$$

### 1.3.7 Solution of the Dirichlet problem of the Laplacian

Let's consider the problem

$$(1.16) \quad \begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}), & \mathbf{x} \in U, \\ u(\mathbf{x}) = h(\mathbf{x}), & \mathbf{x} \in \partial U, \end{cases}$$

and let  $g(\mathbf{x}, \mathbf{y})$  be the solution of the problem (1.15), then by the Green formula one gets

$$\int_U (g \Delta u - u \Delta g) d\mathbf{x} = \int_{\partial U} (g \frac{\partial u}{\partial n} - u \frac{\partial g}{\partial n}) d\sigma = \int_{\partial U} -u \frac{\partial g}{\partial n} d\sigma.$$

By substituting we obtain

$$\int_U (-g(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) + u(\mathbf{x}) \delta(\mathbf{x} - \mathbf{y})) d\mathbf{x} = \int_{\partial U} -h(\mathbf{x}) \frac{\partial g(\mathbf{x}, \mathbf{y})}{\partial n} d\sigma(\mathbf{x}),$$

and then

$$u(\mathbf{y}) = \int_U g(\mathbf{x}, \mathbf{y}) f(\mathbf{x}) d\mathbf{x} - \int_{\partial U} h(\mathbf{x}) \frac{\partial g(\mathbf{x}, \mathbf{y})}{\partial n} d\sigma(\mathbf{x}).$$

By renaming the variables and using the symmetry of  $g$ , which is due to the self-adjointness of  $-\Delta$ , see [2], we have that

$$(1.17) \quad u(\mathbf{x}) = \int_U g(\mathbf{x}, \mathbf{y}) f(\mathbf{y}) d\mathbf{y} - \int_{\partial U} h(\mathbf{y}) \frac{\partial g(\mathbf{x}, \mathbf{y})}{\partial n_y} d\sigma(\mathbf{y})$$

is solution of Dirichlet's Laplacian problem.

### Proof

For the rigorous proof see [2], for example.

## 1.4 Heat equation

### 1.4.1 Definition of Heat equation

Let  $U \subseteq \mathbb{R}^n$  be a connected, open and bounded set,  $u: \overline{U} \times [0, \bar{t}] \rightarrow \mathbb{R}$  be the unknown,  $f: U \rightarrow \mathbb{R}$  and  $g: \partial U \times [0, \bar{t}] \rightarrow \mathbb{R}$  be given functions,  $K > 0$  be the thermal diffusion. Let us define the Heat *equation* [7] as

$$(1.18) \quad \begin{cases} \frac{\partial u}{\partial t} - K \Delta u = 0, & (\mathbf{x}, t) \in U \times (0, \bar{t}), \\ u(\mathbf{x}, 0) = f(\mathbf{x}), & \mathbf{x} \in U, \\ u(\mathbf{x}, t) = g(\mathbf{x}, t), & (\mathbf{x}, t) \in \partial U \times [0, \bar{t}]. \end{cases}$$

### 1.4.2 Definition of parabolic boundary

Let  $U \subseteq \mathbb{R}^n$  be a connected, open and bounded set and let's take  $t \in [0, \bar{t}]$ , we can define the *parabolic boundary* of the cylinder  $Q_T := U \times (0, \bar{t})$  as the set

$$\partial_p Q_T := (\overline{U} \times \{t=0\}) \cup S_T,$$

where  $S_T = \partial U \times (0, \bar{t})$  is the *lateral surface* of the cylinder.

### 1.4.3 Proposition 1 (about the Weak Maximum principle)

Let  $u \in C^{2,1}(Q_T) \cap C(\overline{Q_T})$  satisfy

$$(1.19) \quad \frac{\partial u}{\partial t} - K \Delta u < 0, \quad (\mathbf{x}, t) \in U \times (0, \bar{t}), \quad K > 0.$$

If  $u \leq M$  in  $U \times \{t=0\}$  and  $u \leq M$  in  $\partial U \times [0, \bar{t}]$ , then  $u \leq M$  in  $\overline{U} \times [0, \bar{t}]$ .

#### Proof

Let's suppose by contradiction that  $(\mathbf{x}_0, t_0)$  is a maximum point, with  $\mathbf{x}_0 \in U$ ,  $t_0 \in (0, \bar{t})$ , then

$$\begin{aligned} \frac{\partial}{\partial t} u(\mathbf{x}_0, t_0) &= 0, \\ \Delta u(\mathbf{x}_0, t_0) &\leq 0, \end{aligned}$$

and this is a contradiction because one would have

$$\frac{\partial u}{\partial t} - K \Delta u \geq 0.$$

Then we have shown that

$$(1.20) \quad \max_{U \times (0, \bar{t})} u \leq \max_{\partial U \times (0, \bar{t}] \cup U \times \{0\}} u.$$

Let us suppose that there exists a maximum point  $(x_0, \bar{t})$ , with  $x_0 \in U$ , such that

$$u(x_0, \bar{t}) > M.$$

By the continuity of  $u$ , there exists  $\varepsilon > 0$ , such that  $u(x_0, \bar{t} - \varepsilon) > M$ , but  $(x_0, \bar{t} - \varepsilon) \in U \times (0, \bar{t})$ , whence it's a contradiction with (1.20). ■

#### 1.4.4 Proposition 2 (about the Weak Maximum principle)

Let us suppose that

$$(1.21) \quad \frac{\partial u}{\partial t} - K \Delta u = 0, \quad (x, t) \in U \times (0, \bar{t}), \quad K > 0,$$

with  $u \in C^{2,1}(Q_T) \cap C(\overline{Q_T})$ . If  $u \leq M$  in  $U \times \{t=0\}$  and in  $\partial U \times [0, \bar{t}]$ , then  $u \leq M$  in  $\overline{U} \times [0, \bar{t}]$ .

#### Proof

Let us define

$$v = u + \varepsilon \left( \sum_{i=1}^n x_i^2 \right),$$

and we observe that

$$\frac{\partial v}{\partial t} = \frac{\partial u}{\partial t},$$

$$\Delta v = \Delta u + 2n\varepsilon.$$

Then we obtain

$$\frac{\partial v}{\partial t} - K \Delta v = \frac{\partial u}{\partial t} - K(\Delta u + 2n\varepsilon) = -K2n\varepsilon < 0,$$

that is  $v$  satisfies (1.19). Moreover, since  $U$  is bounded, there exists  $R > 0$  such that  $U \subseteq B(0, R)$ , whence

$$v \leq M + \varepsilon R^2 \text{ in } U \times \{t=0\} \text{ and in } \partial U \times [0, \bar{t}].$$

From the previous proposition, we have that

$$v \leq M + \varepsilon R^2 \text{ in } \overline{U} \times [0, \bar{t}].$$

The thesis follows taking the limit for  $\varepsilon$  going to zero. ■

By **Proposition 1** and **Proposition 2**, we obtain the following

#### 1.4.5 Theorem (Weak Maximum principle)

Let us suppose that

$$(1.22) \quad \frac{\partial u}{\partial t} - K \Delta u \leq 0, \quad (\mathbf{x}, t) \in U \times (0, \bar{t}), \quad K > 0,$$

with  $u \in C^{2,1}(Q_T) \cap C(\overline{Q_T})$ . If  $u \leq M$  in  $U \times \{t=0\}$  and in  $\partial U \times [0, \bar{t}]$ , then  $u \leq M$  in  $\overline{U} \times [0, \bar{t}]$ , that is, the maximum is assumed on the parabolic boundary.

#### 1.4.6 Definition of subsolutions and supersolutions

The functions  $u$  such that  $\frac{\partial u}{\partial t} - K \Delta u \leq 0$  ( $\frac{\partial u}{\partial t} - K \Delta u \geq 0$ ) are called *subolutions* (*supersolutions*) of the heat equation.

#### 1.4.7 Theorem (uniqueness of the solution)

If there exists a solution for the problem

$$(1.23) \quad \begin{cases} \frac{\partial u}{\partial t} - K \Delta u = 0, & (\mathbf{x}, t) \in U \times (0, \bar{t}), \\ u(\mathbf{x}, 0) = f(\mathbf{x}), & \mathbf{x} \in U, \\ u(\mathbf{x}, t) = g(t), & (\mathbf{x}, t) \in \partial U \times [0, \bar{t}], \end{cases}$$

then this solution is unique.

#### Proof

Let's suppose by contradiction that there exist  $u_1, u_2$  both solutions of (1.23). Let's define  $u := u_1 - u_2$ , and observe that it's solution of

$$\begin{cases} \frac{\partial u}{\partial t} - K \Delta u = 0, & (\mathbf{x}, t) \in U \times (0, \bar{t}), \\ u(\mathbf{x}, 0) = 0, & \mathbf{x} \in U, \\ u(\mathbf{x}, t) = 0, & (\mathbf{x}, t) \in \partial U \times [0, \bar{t}], \end{cases}$$

that is  $u = 0$  in  $\{U \times \{t=0\}\} \cup \{\partial U \times [0, \bar{t}]\}$ . By the *Weak Maximum principle theorem*, we obtain

$$u = 0 \text{ in } \overline{U} \times [0, \bar{t}],$$

that is  $u_1 = u_2$  in  $\overline{U} \times [0, \bar{t}]$ .

■

### 1.4.8 Theorem (Strong Maximum principle)

Let us suppose that

$$\frac{\partial u}{\partial t} - K \Delta u \leq 0, \quad (\mathbf{x}, t) \in U \times (0, \bar{t}), \quad K > 0,$$

with  $u \in C^{2,1}(Q_T) \cap C(\overline{Q_T})$ . If  $u$  has maximum  $M$  in  $(\mathbf{x}_I, t_I) \in U \times (0, \bar{t}]$ , then  $u \equiv M$  in  $\overline{U} \times [0, t_I]$ .

### 1.4.9 Fundamental solution of the heat equation (n=1)

We want to study the propagation, in an infinite homogeneous metal bar placed along the real axis, of the heat produced by an instantaneous source in the origin of the axes at time  $t=0$ . We want to find the solution  $u^*$  of

$$u_t - D u_{xx} = 0, \quad t > 0.$$

The law of conservation of energy implies that

$$(1.24) \quad \rho c_v \int_{\mathbb{R}} u^*(x, t) dx = E,$$

where

- $\rho$  is the mass density of the bar;
- $c_v$  is the specific heat;
- $E$  is the total energy.

#### **STEP 1**

Let us define  $Q = \frac{E}{\rho c_v}$ , we want  $u^*$  to be an even function of  $x$  and depend also

on  $t, D, Q$ , i.e.

$$u^* = U(x, t, D, Q) \text{ and } U(x) = U(-x).$$

In order to non-dimensionalize the variables, we list the physical dimensions of each argument of  $U$ , in terms of the three dimensional units:  $\Theta$  (temperature),  $L$  (length) and  $T$  (time):

$$[x] = L, [t] = T, [D] = L^2 T^{-1}, [Q] = \Theta L,$$

or similarly

$$[x] = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, [t] = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, [D] = \begin{bmatrix} 0 \\ 2 \\ -1 \end{bmatrix}, [Q] = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}.$$

We observe that  $[t], [D], [Q]$  are linearly independent and thus form a basis for  $\mathbb{R}^3$ . It follows from these observations that

- $[u^*] = \Theta = \frac{[Q]}{L} = \frac{[Q]}{[D]^{1/2}[t]^{1/2}},$
- $[x] = L = [D]^{1/2}[t]^{1/2}.$

Let's define the non-dimensional quantities

- $\pi = \frac{u^* \sqrt{Dt}}{Q};$
- $\pi_I = \frac{x}{\sqrt{Dt}}.$

From these definitions and from that of  $u^*$ , one obtains

$$\pi = \frac{\sqrt{Dt}}{Q} U(\pi_I) = \tilde{U}(\pi_I, t, D, Q) = \tilde{U}(\pi_I),$$

since  $\pi$  and  $\tilde{U}$  are dimensionless. Then

$$(1.25) \quad u^*(x, t) = \frac{Q}{\sqrt{Dt}} \cdot \tilde{U}\left(\frac{x}{\sqrt{Dt}}\right).$$

## STEP 2

By the law of conservation of energy and making the change of variables

$y = \frac{x}{\sqrt{Dt}}$ , we have that

$$Q = \int_{\mathbb{R}} u^*(x, t) dx = \int_{\mathbb{R}} \frac{Q}{\sqrt{Dt}} \tilde{U}\left(\frac{x}{\sqrt{Dt}}\right) dx = Q \int_{\mathbb{R}} \tilde{U}(y) dy,$$

which implies

$$(1.26) \quad \int_{\mathbb{R}} \tilde{U}(y) dy = 1.$$

Calculating the derivatives  $u_t^*$  and  $u_{xx}^*$ , it is derived that

$$u_t^* - D u_{xx}^* = \frac{-Q}{t \sqrt{Dt}} \left[ \tilde{U}''(y) + \frac{1}{2} \tilde{U}'(y) + \frac{1}{2} \tilde{U}(y) \right] = 0.$$

Then  $u^*$  is solution of the problem if and only if

$$\tilde{U}''(y) + \frac{1}{2} \tilde{U}'(y) + \frac{1}{2} \tilde{U}(y) = 0,$$

that is

$$\frac{\partial}{\partial y} \left( \tilde{U}'(y) + \frac{1}{2} y \tilde{U}(y) \right) = 0,$$

then

$$\tilde{U}'(y) + \frac{1}{2}y\tilde{U}(y) = C.$$

Since  $u^*$  is an even function, we have that  $\tilde{U}'(0)=0$ . Evaluating in zero gives  $C=0$ . Thus, we obtain the equation

$$\tilde{U}'(y) + \frac{1}{2}y\tilde{U}(y) = 0,$$

which has solution

$$\tilde{U}(y) = c_0 e^{-y^2/4}.$$

From (1.26) we get

$$I = \int_{\mathbb{R}} \tilde{U}(y) dy = \int_{\mathbb{R}} c_0 e^{-y^2/4} dy = c_0 \int_{\mathbb{R}} e^{-z^2/4} 2dz = 2c_0 \sqrt{\pi},$$

which implies

$$c_0 = \frac{I}{\sqrt{4\pi}}.$$

Then we obtain

$$u^*(x, t) = \frac{Q}{\sqrt{4\pi Dt}} e^{\frac{-x^2}{4Dt}}.$$

The function

$$(1.27) \quad \Gamma_D(x, t) = \frac{1}{\sqrt{4\pi Dt}} e^{\frac{-x^2}{4Dt}}, \quad x \in \mathbb{R}, \quad t > 0,$$

is called the *fundamental solution of the heat equation* in one dimension. It is a *self-similar* solution in the sense that it keeps spatially "similar" to itself as time changes.

We observe that it possesses the following properties:

1. if  $x \neq 0$ , then  $\lim_{t \rightarrow 0^+} \Gamma_D(x, t) = \lim_{t \rightarrow 0^+} \frac{1}{\sqrt{4\pi Dt}} e^{\frac{-x^2}{4Dt}} = 0$ ;
2. if  $x = 0$ , then  $\lim_{t \rightarrow 0^+} \Gamma_D(x, t) = \lim_{t \rightarrow 0^+} \frac{1}{\sqrt{4\pi Dt}} e^{\frac{-x^2}{4Dt}} = +\infty$ ;
3.  $\lim_{t \rightarrow 0^+} \int_{\mathbb{R} \setminus [-\delta, \delta]} \Gamma_D(x, t) dx = 0, \forall \delta > 0$ .

From these properties it follows that the fundamental solution is concentrated around the origin. Furthermore, for  $t \rightarrow 0^+$ , we have that  $\Gamma_D(x, t)$  satisfies all the

properties of Dirac's Delta. The function  $\Gamma_D(x-y, t)$  can be interpreted as the solution of the unitary heat equation when the heat source is concentrated in  $y$  at  $t=0$ , that is it satisfies the initial condition

$$\Gamma_D(x-y, 0) = \delta(x-y).$$

As soon as  $t>0$ , then  $\Gamma_D(x, t)>0$  over all  $\mathbb{R}$ , that is, heat diffuses instantaneously therefore with infinite speed, which is a limit of the model of heat propagation.

#### 1.4.10 Fundamental solution of the heat equation ( $n>1$ )

Let's study the diffusion, over all  $\mathbb{R}^n$ , of heat produced by an instantaneous source placed in the origin of axes at  $t=0$ . Similarly to the  $n=1$  case, it is possible to show that the solution is

$$(1.28) \quad u^*(x, t) = \frac{Q}{(4\pi Dt)^{n/2}} e^{-\frac{|x|^2}{4Dt}}, \quad x \in \mathbb{R}^n, \quad t > 0.$$

The function

$$(1.29) \quad \Gamma_D(x, t) = \frac{I}{(4\pi Dt)^{n/2}} e^{-\frac{|x|^2}{4Dt}},$$

is called the *fundamental solution of the heat equation* in dimension  $n$ .

The function  $\Gamma_D(x-y, t)$  can be interpreted as the solution of the unitary heat equation when the heat source is concentrated in  $y$  and it's the unique solution of the problem

$$\begin{cases} u_t - D\Delta u = 0, & x \in \mathbb{R}^n, \quad t > 0, \\ u(x, 0) = \delta(x-y), & x \in \mathbb{R}^n. \end{cases}$$

#### 1.4.11 The global Cauchy problem ( $n=1$ )

##### *Homogeneous case*

Let's consider the homogeneous problem

$$\begin{cases} u_t - Du_{xx} = 0, & x \in \mathbb{R}, \quad t > 0, \\ u(x, 0) = g(x), & x \in \mathbb{R}, \end{cases}$$

where  $g$  is a given function. The idea is to determine, for example, the concentration  $u(x, t)$  of a substance whose initial concentration is  $g(x)$ .  $g(y)dy$  represents the quantity contained in the interval  $(y, y+dy)$ , and

$$\Gamma_D(x-y, t)g(y)dy$$

represents the corresponding mass concentration at point  $x$  at time  $t$ .

It's possible to prove that the solution is given by

$$(1.30) \quad u(x, t) = \int_{\mathbb{R}} \Gamma_D(x-y, t)g(y)dy = \frac{1}{\sqrt{4\pi Dt}} \int_{\mathbb{R}} g(y) e^{\frac{-(x-y)^2}{4Dt}} dy,$$

which is called the *convolution of the initial value* with the fundamental solution. Of course, these are heuristic considerations, and it must be shown that it is indeed solution of the Cauchy problem under consideration.

#### 1.4.12 Theorem (Existence of Solution)

Let  $g$  be a function with countable points of discontinuity in  $\mathbb{R}$ , such that

$$|g(x)| \leq c e^{ax^2}, \quad \forall x \in \mathbb{R},$$

with  $a, c > 0$ , and let  $u(x, t)$  be defined as in (1.30), then

1.  $u(x, t)$  is well defined and differentiable infinitely many times in  $\mathbb{R}$ ,  $\forall t < T$ , with  $T < \frac{1}{4Da}$ , and we have that

$$u_t - Du_{xx} = 0, \text{ in } \mathbb{R} \times (0, T);$$

2. if  $g$  is continuous in  $x_0$ , then

$$\lim_{(x, t) \rightarrow (x_0, 0)} u(x, t) = g(x_0), \quad t > 0;$$

3. there exist  $A, C$  constants such that

$$|u(x, t)| \leq C e^{Ax^2}, \text{ in } \mathbb{R} \times (0, T).$$

This theorem states that (1.30) is solution of the homogeneous Cauchy problem in a finite time interval.

#### *Non-homogeneous case*

Let's consider the problem

$$\begin{cases} v_t - Dv_{xx} = \delta(x-y, t-s), & s > 0, \\ v(x, 0) = 0, & x \in \mathbb{R}. \end{cases}$$

Since up to the time  $t=s>0$  nothing happens and after that time  $\delta(x-y, t-s)=0$ , we can alternatively consider the problem

$$\begin{cases} u_t - Du_{xx} = 0, & x \in \mathbb{R}, \quad t > s, \\ u(x, s) = \delta(x-y), & s \geq 0. \end{cases}$$

If  $s=0$ , the solution is

$$u(x, t) = \Gamma_D(x-y, t),$$

instead, if  $s$  is generic, the solution becomes

$$u(x, t) = \Gamma_D(x-y, t-s).$$

Now consider a source distributed in the strip  $\mathbb{R} \times (0, T)$  that produces mass density at the rate  $f(x, t)$ , then

$$f(x, t) dx dt$$

is the mass produced in the spatial region  $(x, x+dx)$  in the time interval  $(t, t+dt)$ .

If no mass is initially present, the problem becomes

$$(1.31) \quad \begin{cases} v_t - Dv_{xx} = f(x, t), & \text{in } \mathbb{R} \times (0, T), \\ v(x, 0) = 0, & x \in \mathbb{R}. \end{cases}$$

### STEP 1

We construct a family  $w(x, t, s)$  of solutions of the homogeneous Cauchy problem, in which the initial time, instead of being  $t=0$ , fixed, is a time  $0 < s < t$ , variable, and the initial datum is  $f(x, s)$ , i.e. we consider the problem:

$$\begin{cases} w_t - Dw_{xx} = 0, & x \in \mathbb{R}, \quad t > s, \\ w(x, s, s) = f(x, s), & x \in \mathbb{R}. \end{cases}$$

The solution is given by

$$w(x, t, s) = \int_{\mathbb{R}} \Gamma_D(x-y, t-s) f(y, s) dy.$$

### STEP 2

By integrating, we have

$$v(x, t) = \int_0^t w(x, t, s) ds = \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) f(y, s) dy ds,$$

which is solution of (1.31).

#### 1.4.13 Theorem (existence for the non-homogeneous Cauchy problem)

Let  $g$  be a function with numerable points of discontinuity in  $\mathbb{R}$ , such that

$$|g(x)| \leq c e^{ax^2}, \quad \forall x \in \mathbb{R},$$

with  $a, c > 0$ , and let  $f, f_t, f_x, f_{xx}$  be continuous and bounded functions in  $\mathbb{R} \times [0, \bar{t}]$ , then a solution for the problem

$$(1.32) \quad \begin{cases} z_t - Dz_{xx} = f, & \text{in } \mathbb{R} \times (0, T), \\ z(x, 0) = g, & \text{in } \mathbb{R}, \end{cases}$$

is given by

$$(1.33) \quad \begin{aligned} z(x, t) &= u(x, t) + v(x, t) \\ &= \int_{\mathbb{R}} \Gamma_D(x-y, t) g(y) dy + \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) f(y, s) dy ds. \end{aligned}$$

Furthermore, (1.33) is the only solution of (1.32)

#### 1.4.14 Theorem (principle of global maximum)

Let's consider  $z(x, t) \in C(\mathbb{R} \times [0, T])$ , with  $z_t, z_x, z_{xx} \in C(\mathbb{R} \times (0, T))$ , such that

$$(1.34) \quad \begin{cases} z_t - Dz_{xx} \leq 0, & (\geq 0) \\ z(x, t) \leq Ce^{Ax^2}, & (\geq -Ce^{Ax^2}) \end{cases}$$

with  $C, A > 0$ . Then

$$(1.35) \quad \sup_{\mathbb{R} \times [0, T]} z(x, t) \leq \sup_{\mathbb{R}} z(x, 0), \quad \left( \inf_{\mathbb{R} \times [0, T]} z(x, t) \geq \inf_{\mathbb{R}} z(x, 0) \right).$$

#### Corollary

Let  $u$  be a solution of

$$\begin{cases} u_t - Du_{xx} = 0, & \text{on } \mathbb{R} \times (0, T), \\ u(x, 0) = 0, & \text{on } \mathbb{R}, \end{cases}$$

where  $u, u_t, u_x, u_{xx}$  are continuous function. If  $|u| \leq Ce^{Ax^2}$ , with  $C, A > 0$ , then  $u \equiv 0$ .

#### Proof

From the *principle of global maximum*, we have

$$0 \leq \inf_{\mathbb{R}} u(x, 0) \leq \inf_{\mathbb{R} \times [0, T]} u(x, t) \leq \sup_{\mathbb{R} \times [0, T]} u(x, t) \leq \sup_{\mathbb{R}} u(x, 0) = 0,$$

then  $u \equiv 0$ .

■

## 1.5 Scalar reaction-diffusion equations

### 1.5.1 Introduction to scalar reaction-diffusion equations

*Reaction-diffusion models* [7] are concerned with the evolution over time of a population of individuals residing within an open spatial region  $\Omega \subset \mathbb{R}^d$ . Let  $x \in \Omega$  be a point and let  $\{O_n\}_{n \in \mathbb{N}} \subset \Omega$  be a decreasing sequence of spatial regions containing  $x$ , such that  $|O_n| \rightarrow 0$ , for  $n \rightarrow +\infty$ . Then, we can define the *density function of individuals* as

$$(1.36) \quad u(x, t) = \frac{\text{number of individuals within } O_n}{|O_n|},$$

if we take  $n$  such that  $O_n$  is small relative to  $\Omega$ , but large enough to contain a significant number of individuals. Then the *total population* present in a measurable sub-region  $O \subset \Omega$  is given by

$$(1.37) \quad \int_O u(x, t) d\mathbf{x}.$$

The population within a region may vary due to a number of factors, including:

- people entering or leaving the region (for example for a *diffusion process*);
- births or deaths caused by physical, biological, chemical reasons and so on (*reaction process*).

Historically, for example, people have always moved from highly populated regions to less populated ones, so the flow  $\mathbf{J}$  has opposite direction to that in which  $u$  grows, i.e., *Fick's law* can be assumed to apply

$$(1.38) \quad \mathbf{J}(x, t) = -D(x) \nabla u(x, t),$$

where  $D(x) \geq 0$  is said the *diffusion coefficient*.

In addition, we can denote by  $f(u, x, t)$  the *reaction rate*, which is the rate of change in population density due to biological, physical and chemical reasons.

From the above considerations, it follows

$$\frac{\partial}{\partial t} \int_O u(x, t) d\mathbf{x} = - \int_{\partial O} \mathbf{J}(x, t) \cdot \mathbf{n} d\sigma + \int_O f(u, x, t) d\mathbf{x}.$$

Applying Gauss's theorem, we obtain

$$\int_O \frac{\partial}{\partial t} u(x, t) d\mathbf{x} = \int_O \nabla \cdot (D(x) \nabla u(x, t)) d\mathbf{x} + \int_O f(u, x, t) d\mathbf{x},$$

from which we get the *reaction-diffusion equation*

$$(1.38) \quad \frac{\partial}{\partial t} u(\mathbf{x}, t) = \nabla \cdot (D(\mathbf{x}) \nabla u(\mathbf{x}, t)) + f(u, \mathbf{x}, t).$$

In the case when the spatial region is *homogeneous*, the equation becomes

$$(1.39) \quad \frac{\partial}{\partial t} u(\mathbf{x}, t) = D \Delta u(\mathbf{x}, t) + f(u, \mathbf{x}, t).$$

### 1.5.2 Theorem (maximum principle)

Let's take  $u(\mathbf{x}, t) \in C^2(Q_T) \cap C(\overline{Q_T})$  and let  $C = C(\mathbf{x}, t)$  be a bounded function on  $Q_T$ . If

$$u_t - D \Delta u \leq C u, \quad \text{in } Q_T,$$

then we have that

$$(1.40) \quad \max_{\partial_p Q_T} u \leq 0 \Rightarrow \max_{\overline{Q_T}} u \leq 0.$$

### 1.5.3 Theorem (of comparison)

Let  $u, v \in C^2(Q_T) \cap C(\overline{Q_T})$  be functions such that

$$\begin{aligned} u_t - D \Delta u &\geq f(u) \text{ in } Q_T, \\ v_t - D \Delta v &\leq f(v) \text{ in } Q_T. \end{aligned}$$

Then

$$u \geq v \text{ on } \partial_p Q_T \Rightarrow u \geq v \text{ in } \overline{Q_T}.$$

### 1.5.4 Global existence for the Cauchy problem

For simplicity, let us suppose that  $u_0 \in C(\mathbb{R})$ . Consider the spatially one-dimensional Cauchy problem

$$(1.41) \quad \begin{cases} u_t - Du_{xx} = f(u), & \text{in } \mathbb{R}, \\ u(x, 0) = u_0(x) \in L^\infty(\mathbb{R}), \end{cases}$$

where  $f$  is a globally Lipschitzian function, i.e.

$$\exists L > 0 \text{ s.t. } |f(u) - f(v)| \leq L|u - v|, \quad \forall u, v \in \mathbb{R}.$$

A solution for this problem, if it exists, must satisfy the *Duhamel's formula*:

$$(1.42) \quad u(x, t) = \int_{\mathbb{R}} \Gamma_D(x - y, t) u_0(y) dy + \int_0^t \int_{\mathbb{R}} \Gamma_D(x - y, t - s) f(u(y, s)) dy ds.$$

Let's define the operator

$$u \rightarrow A(u) = \int_{\mathbb{R}} \Gamma_D(x-y, t) u_0(y) dy + \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) f(u(y, s)) dy ds.$$

Let us look for a fixed point of  $A$ , i.e., a function  $u$  such that  $u = A(u)$ . Let's take  $T > 0$ , and define the Banach space

$$X = C(\mathbb{R} \times [0, T]) \cap L^\infty(\mathbb{R} \times [0, T]),$$

endowed with the norm

$$\|u\|_X = \sup_{\mathbb{R} \times [0, T]} |u(x, t)|.$$

Proceeding with the calculations we have

$$\begin{aligned} |A(u) - A(v)| &\leq \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) |f(u) - f(v)| dy ds \\ &\leq L \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) |u(y, s) - v(y, s)| dy ds \\ &\leq L \|u - v\|_X \int_0^t \int_{\mathbb{R}} \Gamma_D(x-y, t-s) dy ds \\ &\leq LT \|u - v\|_X. \end{aligned}$$

By choosing  $T$  such that  $LT < 1$  we have that  $A$  is a contraction, then by the *Banach-Caccioppoli Theorem* we have the local existence of the solution in the interval  $[0, T]$ . To go to the global existence of the solution, one must iterate the procedure in the interval  $[T, 2T]$  with initial datum given by  $u(x, T)$  and so on.

### Remark

The result can also be proved in the case when the function  $f$  is locally Lipschitzian, if one has a priori estimates for the solutions.

#### 1.5.5 Global existence in the case of systems

Suppose we are in the case where  $m$  chemicals interact with each other with vector source term  $\mathbf{f} = (f_1, \dots, f_m)$ . We have a vector of *unknown functions* of components  $u_i(\mathbf{x}, t)$ ,  $i=1, \dots, m$ , each endowed with its own *diffusion coefficient*  $D_i$ ,  $i=1, \dots, m$ . The Cauchy problem becomes

$$(1.43) \quad \begin{cases} \mathbf{u}_t = \mathbf{D} \Delta \mathbf{u} + \mathbf{f}(\mathbf{u}), \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \in L^\infty(\mathbb{R}^d), \end{cases}$$

where

- $\mathbf{u} \in \mathbb{R}^m$ ;
- $\mathbf{D}$  is the *diffusivity matrix*, which is positive semidefinite;
- $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$ .

If  $\mathbf{D} = \text{diag}(D_1, \dots, D_m)$ , with  $D_i \geq 0$ , and if  $f$  is globally Lipschitzian, we can rewrite the system component by component as

$$\frac{\partial u_i}{\partial t} = D_i \Delta u_i + f_i(u_1, \dots, u_m), \quad i=1, \dots, m.$$

Under these assumptions, the existence of a solution is guaranteed as in the scalar case.

## 1.6 Reaction-diffusion models for a biological species

### 1.6.1 Reaction rate

Several choices can be made for the *reaction rate* [7] in these models:

1. *Malthaus linear growth*, in which

$$(1.44) \quad f = r u,$$

where

- $r = b - d$  is called the *growth rate*;
- $b$  is called the *birth rate* (constant);
- $d$  is called *death rate* (constant).

In the spatially homogeneous case we have  $u = u(t)$  and the equation becomes

$$u_t = r u,$$

which has solution

$$u(t) = u_0 e^{rt},$$

where  $u_0$  represents the *initial population*. If  $r > 0$  there is an exponential growth, otherwise there is an exponential decrease.

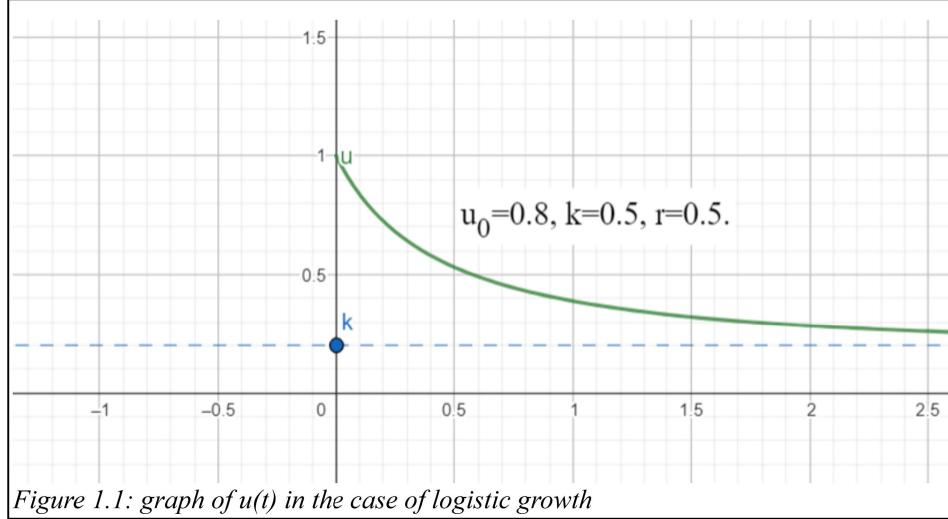
2. *logistic growth*, in which

$$(1.45) \quad f = r u \left(1 - \frac{u}{k}\right),$$

where  $k$  is the *carrying capacity* of the environment, accounting for resource limitation and competition among individuals.

In the spatially homogeneous case we have that

$$u(t) = u_0 \frac{k e^{rt}}{k + u_0(e^{rt} - 1)}.$$



Therefore, if  $r > 0$ , whatever its initial value, the population tends to reach the value of  $k$ , see Figure 1.1.

We observe that there are two equilibrium states:

- $u = 0$ , which is unstable;
- $u = k$ , which is stable and attractive.

3. *growth with an Allee effect*, which takes into account that the presence of a small population can be a disadvantage because, for example, mating may be difficult. A typical example is given by

$$(1.46) \quad f = r u \left( 1 - \frac{u}{k} \right) (u - c),$$

where  $c$  is the *constant of deficiency*.

### 1.6.2 Fisher–Kolmogorov–Petrovsky–Piskunov Equation

It is a reaction-diffusion equation in which the logistic growth rate is used as the reaction rate:

$$(1.47) \quad u_t = r u (1 - u) + D u_{xx},$$

where  $r, D > 0$ ,  $k = 1$ .

The importance of this equation is primarily didactic in that it represents the prototype of equations that admit travelling wave solutions and is useful for showing many standard techniques that are used for the analysis of models with diffusive scattering. By taking  $t^* = rt$ ,  $x^* = x\sqrt{rD}$ , and omitting the asterisks, (1.47) can be rewritten as

$$(1.48) \quad u_t = u(1-u) + u_{xx}.$$

### 1.6.3 Definition of monostable function

A reaction function  $f \in C^1(\mathbb{R})$  is said to be *monostable* if

- $f(u) > 0$ , if  $u \in (0, 1)$ ,
- $f(u) < 0$ , if  $u \notin (0, 1)$ .

#### Remark

The function  $f(u) = u(1-u)$  is monostable.

### 1.6.4 Proposition (reaction-diffusion equation with monostable reaction rate)

For a reaction-diffusion equation with monostable reaction rate, we have that every constant  $M > 1$  is supersolution. Moreover  $u \equiv 0$  and  $u \equiv 1$  are solutions.

#### Proof

Since  $M > 1$ , through a direct calculation we have

$$\frac{\partial u}{\partial t} - f(u) - \frac{\partial^2 u}{\partial x^2} = \frac{\partial M}{\partial t} - f(M) - \frac{\partial^2 M}{\partial x^2} = -f(M) > 0.$$

■

### 1.6.5 Properties (boundedness of the solution)

Let us consider the problem

$$\begin{cases} u_t - Du_{xx} = f(u), \\ u(x, 0) = u_0(x). \end{cases}$$

If  $0 \leq u_0(x) \leq M$ , with  $M > 1$ , then  $0 \leq u(x, t) \leq M$ ,  $\forall x \in \mathbb{R}, \forall t > 0$ .

### 1.6.6 Stationary solutions

We want to start studying *stationary solutions*, that is, solutions such that  $u_t=0$ , from which we get

$$u_{xx} + f(u) = 0.$$

Multiplying by  $u_x$  we get

$$u_{xx}u_x + f(u)u_x = 0.$$

We now define the potential energy function  $F$ , such that

$$\frac{\partial F}{\partial u} = f(u),$$

from which

$$\frac{\partial F}{\partial x} = \frac{\partial F}{\partial u} \cdot \frac{\partial u}{\partial x} = f(u)u_x.$$

Integrating  $u_{xx}u_x + f(u)u_x = 0$  we have

$$\int_{\mathbb{R}} u_{xx}u_x + f(u)u_x dx = \int_{\mathbb{R}} 0 dx,$$

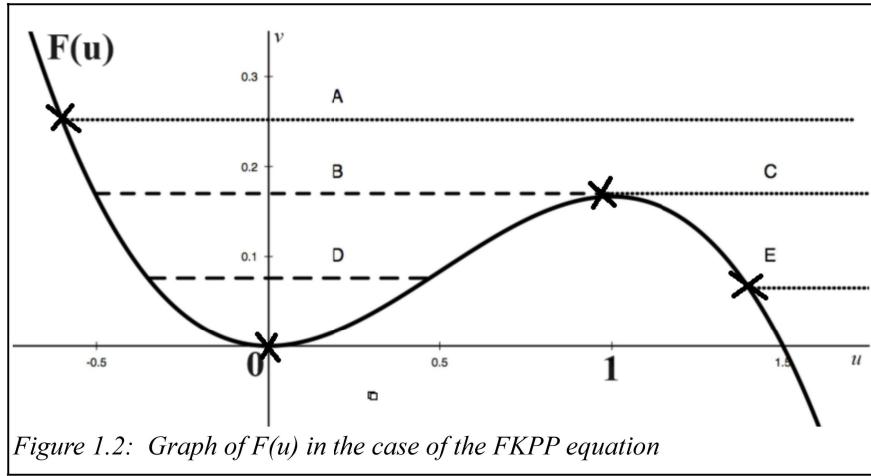
that is,  $\frac{1}{2}(u_x)^2 + F(u) = K$ , with  $K \in \mathbb{R}$ , from which we get

$$\left| \frac{\partial u}{\partial x} \right| = \sqrt{2(K - F(u))}.$$

Then

$$\pm \int_{u(x_0)}^{u(x)} \frac{1}{\sqrt{2(K - F(u))}} du = \int_{x_0}^x dx = x - x_0.$$

Drawing the graph of  $F(u)$ , for example in the case of the FKPP equation, we see that the stationary solutions are in biunivocal correspondence, up to translations and reflections, with the horizontal segments  $v = K$  in the plane  $(u, v)$ . At the values  $K_0 = F(0)$  and  $K_1 = F(1)$  we have the equilibrium states  $u = 0$  and  $u = 1$ . As we can see in Figure 1.2, the solutions corresponding to segments A, C and E are unbounded, while those corresponding to segments B and D are finite. In addition, the solutions corresponding to segments A and E have minima given by the point of intersection with the graph of the function  $F$ , while, for  $x \rightarrow \pm\infty$ , they diverge. The solution corresponding to segment C tends to 1 for  $x \rightarrow -\infty$ , and diverges for  $x \rightarrow +\infty$ , or viceversa. Segment B has negative minimum and tends to 1 for  $x \rightarrow \pm\infty$ . While D has minimum and maximum.



We remark that the only bounded non-negative solutions are  $u \equiv 0$  and  $u \equiv 1$ .

## 1.7 Travelling waves

### 1.7.1 Definition of travelling waves

A solution  $u = u(x, t)$  of an evolution equation is called a *travelling wave* [7] if there exists a function  $U$ , called the *wave profile*, and a constant  $c$  called the *wave velocity*, such that  $u(x, t) = U(x - ct)$ . We want to determine if there exist travelling waves of reaction-diffusion equations, with asymptotic conditions of the type:

$$U(+\infty) = u^+, \quad U(-\infty) = u^-, \text{ with } u^+, u^- \text{ such that } f(u^+) = f(u^-) = 0.$$

In this regard, one has the

### 1.7.2 Theorem (reaction-diffusion equation with monostable reaction rate admits travelling waves)

Let  $f$  be a monostable function such that

$$f(u) < f'(0)u, \quad u \in (0, 1).$$

Then the equation

$$u_t = u_{xx} + f(u)$$

admits travelling waves such that

$$U(-\infty) = 1, \quad U(+\infty) = 0, \quad U(z) \geq 0, \quad \forall z \in \mathbb{R},$$

if and only if

$$(1.49) \quad c \geq c_{\min} := 2 \sqrt{f'(0)}.$$

### 1.7.3 Logistic growth term

It is easy to see that the logistic growth term satisfies the required condition. In this case, in terms of dimensional quantities, we have that

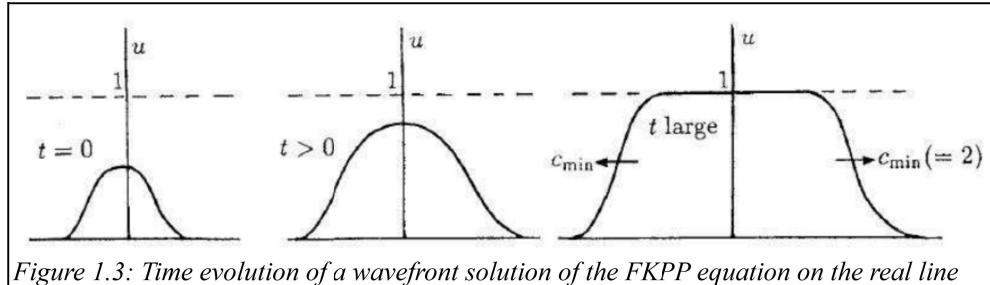
$$(1.50) \quad c_{\min} = \sqrt{r D}.$$

In addition, the conditions of  $u(x,0)$  for the FKPP solution to evolve into a travelling wave solution have been widely studied. It was shown that if  $\exists x_0, x_1 \in \mathbb{R}$ , with  $x_0 < x_1$ , such that

$$u(x,0) = u_0(x) = \begin{cases} 1, & \text{if } x \leq x_0, \\ 0, & \text{if } x \geq x_1, \end{cases}$$

with  $u_0(x)$  continuous function, then  $u(x,t)$  evolves to a travelling wave solution  $U(z)$  such that

$$(1.51) \quad z = x - 2t, \quad c_{\min} = 2.$$



### 1.7.4 Asymptotic approximation

Let us consider the reaction-diffusion equation

$$(1.52) \quad u_t = u_{xx} + u(1-u).$$

By imposing  $u = U(z)$ , with  $z = x - ct$ , we obtain

$$(1.53) \quad U'' + c U' + U(1-U) = 0.$$

By the previous theorem, it admits a monotone solution such that

$$U(-\infty) = 1, \quad U(+\infty) = 0, \quad U(z) \geq 0, \quad \forall z \in \mathbb{R}.$$

Furthermore, it can be shown that an approximation to this solution, in the smallness parameter  $\epsilon = \frac{1}{c^2}$ , is of the form

$$(1.54) \quad U(z) = \frac{1}{1+e^{z/c}} + \frac{1}{c^2} \cdot \frac{e^{z/c}}{1+e^{z/c}} \log \left[ \frac{4e^{z/c}}{(1+e^{z/c})^2} \right] + O\left(\frac{1}{c^2}\right),$$

with  $c \geq c_{min} = 2$ .

If we want to give an estimate of the slope of the wave, we must calculate

$$\max|U'(z)|,$$

which is reached at the point where  $U''=0$ , that is, for  $z=0$ . Continuing with the calculations it is possible to obtain that an estimate for the slope is given by

$$(1.55) \quad s \approx U'(0) \approx -\frac{1}{4c}, \quad c \geq 2,$$

which shows that the higher the speed, the lower the wave front steepness will be.

An estimate of the length of the wave is given by

$$(1.56) \quad L \approx \frac{1}{s} = 4c.$$

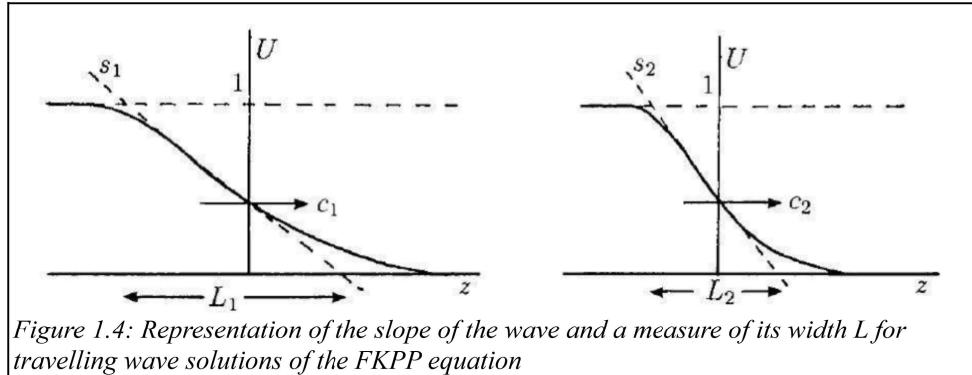


Figure 1.4: Representation of the slope of the wave and a measure of its width  $L$  for travelling wave solutions of the FKPP equation

### 1.7.5 Stability

We want to study how the travelling wave solutions behave with respect to small perturbations. Let us consider the equation

$$(1.57) \quad u_t = u_{xx} + u(1-u),$$

and we make the transformation to the system travelling with the wave

$$u(x, t) = u(z, t), \text{ with } z = x - ct,$$

with an abuse of notation. Substituting in (1.57) yields

$$(1.58) \quad u_t = cu_z + u_{zz} + u(1-u), \text{ with } c \geq 2.$$

Let  $u_c(z)$  denote the travelling wave solution  $U(z)$  of (1.53), and let us consider a small compactly supported perturbation of it

$$u(z, t) = u_c(z) + \varepsilon v(z, t), \quad 0 < \varepsilon \ll 1.$$

Substituting in (1.58) we get the equation

$$v_t = (1 - 2u_c)v + c v_z + v_{zz}.$$

It is possible to show all the solutions of the type

$$v(z, t) = g(z)e^{-\lambda t},$$

have  $\lambda$  positive, apart from the solution with  $\lambda=0$  which represents a small translation of  $u_c$ , and  $g$  is a continuous function. Therefore, for  $\lambda>0$ ,

$$\lim_{t \rightarrow +\infty} v(z, t) = g(z)e^{-\lambda t} = 0,$$

that is,  $u_c(z)$  is stable for small perturbations of the type under consideration.

## 1.8 Systems of reaction-diffusion equations

### 1.8.1 Definition of conservative system

Given  $\mathbf{u} : (0, T) \rightarrow \mathbb{R}^m$ , the ordinary system of differential operations

$$(1.59) \quad \dot{\mathbf{u}}(t) = \mathbf{f}(\mathbf{u})$$

is said to be *conservative* [7] if there exists a real valued function  $g(\mathbf{u})$  such that

$$\dot{g}(\mathbf{u}) = \frac{\partial g}{\partial u_i} \dot{u}_i = \frac{\partial g}{\partial u_i} f_i = 0$$

### 1.8.2 Lotka – Volterra system

Let us consider the prey-predator system

$$(1.60) \quad \begin{cases} \frac{du}{dt} = u(a - bv), \\ \frac{dv}{dt} = v(cu - e), \end{cases}$$

where:

- $u$  = prey population;
- $v$  = predator population;
- $a$  = preys growth rate;
- $e$  = predators decrease rate;
- $b$  = predation rate;
- $c$  = coefficient of prey-predator conversion.

We can adimensionalize the system by placing

$$\hat{t}=at, \quad \hat{u}=\frac{c}{e}u, \quad \hat{v}=\frac{b}{a}v, \quad \alpha=\frac{e}{a},$$

and we obtain

$$(1.61) \quad \begin{cases} \frac{du}{dt}=u(1-v)=f_1(u,v), \\ \frac{dv}{dt}=\alpha v(u-1)=f_2(u,v), \end{cases}$$

where, for simplicity, we have dropped the hats.

To study the equilibrium points, let us set

$$\frac{du}{dt}=\frac{dv}{dt}=0,$$

whence

$$\begin{cases} u(1-v)=0, \\ \alpha v(u-1)=0. \end{cases}$$

The equilibrium points are  $O=(0,0)$ ,  $Q=(1,1)$ .

The Jacobian matrix is

$$J=\begin{pmatrix} \frac{\partial f_1}{\partial u} & \frac{\partial f_1}{\partial v} \\ \frac{\partial f_2}{\partial u} & \frac{\partial f_2}{\partial v} \end{pmatrix}=\begin{pmatrix} 1-v & -u \\ \alpha v & \alpha(u-1) \end{pmatrix}$$

Evaluated in  $O$ , it gives

$$J(O)=\begin{pmatrix} 1 & 0 \\ 0 & -\alpha \end{pmatrix},$$

which is diagonal, then the eigenvalues are  $\lambda=-\alpha$ ,  $\lambda=1$ , and thus it is an unstable equilibrium point (*saddle point*).

Evaluated in  $Q$ , the Jacobian is

$$J(Q)=\begin{pmatrix} 0 & -1 \\ \alpha & 0 \end{pmatrix},$$

whence

$$P_\lambda(J)=|J-\lambda I|=\begin{vmatrix} -\lambda & -1 \\ \alpha & -\lambda \end{vmatrix}=\lambda^2+\alpha.$$

The eigenvalues are  $\pm i\sqrt{\alpha}$ , so we cannot apply the Hartman-Grobman Theorem, see [8].

Consider then the system

$$\begin{cases} \frac{du}{dt} = u(1-v) = f_1(u, v), \\ \frac{dv}{dt} = \alpha v(u-1) = f_2(u, v). \end{cases}$$

By formally dividing the first by the second we have

$$\frac{du}{dv} = \frac{u(1-v)}{\alpha v(u-1)},$$

from which

$$\alpha \frac{u-1}{u} du = \frac{1-v}{v} dv,$$

which when integrated gives us

$$g(u, v) := \alpha(u - \log u) + v - \log v = \text{const.}$$

We note that  $\dot{g} = 0$ , so the system is conservative.

Since

$$H(g) = \begin{pmatrix} \frac{\partial^2 g}{\partial u^2} & \frac{\partial^2 g}{\partial u \partial v} \\ \frac{\partial^2 g}{\partial v \partial u} & \frac{\partial^2 g}{\partial v^2} \end{pmatrix} = \begin{pmatrix} \frac{\alpha}{u^2} & 0 \\ 0 & \frac{1}{v^2} \end{pmatrix}$$

and  $|H(g)| = \frac{\alpha}{u^2 v^2} > 0$  then  $g$  is convex over all  $\mathbb{R}^2 \setminus \{0\}$ .

Moreover

$$\begin{cases} \frac{\partial g}{\partial u} = \alpha(1 - \frac{1}{u}) = 0 \\ \frac{\partial g}{\partial v} = (1 - \frac{1}{v}) = 0 \end{cases}$$

The system is zero only at the point  $Q = (1, 1)$ , which is a global minimum point for  $g$  and it holds

$$g(1, 1) = \alpha + 1$$

One has

$$\lim_{u \rightarrow +\infty} g(u, v) = +\infty,$$

$$\lim_{u \rightarrow 0} g(u, v) = +\infty,$$

similarly for  $v \rightarrow +\infty, 0^+$ .

It is obtained that for  $k > \alpha + 1$ , the contour lines  $\{(u, v) | g(u, v) = k\}$ , are closed regular curves containing  $(1, 1)$ . The solutions are then periodic and the period depends on  $k$ .

Furthermore, we have that, if  $T(k)$  is the period of the orbit, then

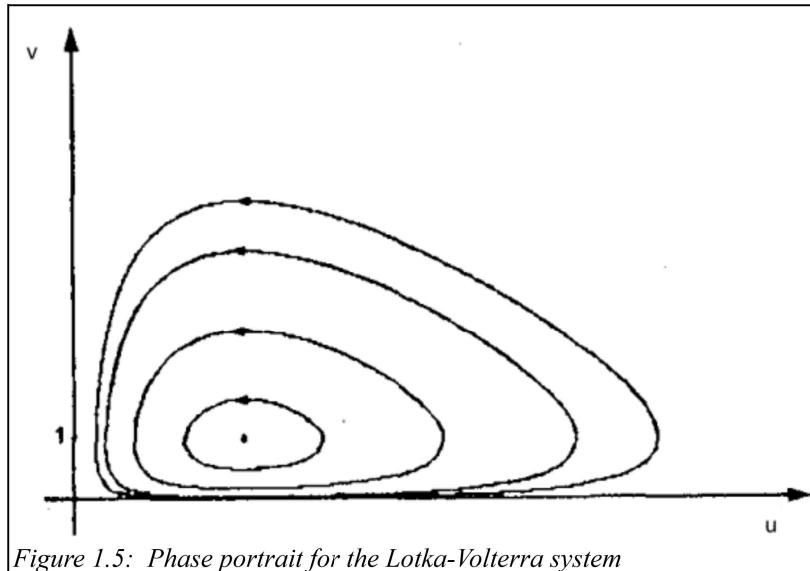
$$\int_0^{T(k)} \frac{\dot{u}}{u} dt = \ln(u(t))|_0^{T(k)} = 0,$$

$$\int_0^{T(k)} \frac{\dot{v}}{v} dt = \ln(v(t))|_0^{T(k)} = 0,$$

from which

$$0 = \int_0^T \frac{f_1}{u} dt = \int_0^T (1-v) dt = T - \int_0^T v(t) dt,$$

$$0 = \int_0^T \frac{f_2}{v} dt = \int_0^T (u-1) dt = T - \int_0^T u(t) dt.$$



### 1.8.3 Structural instability

Let us consider small perturbations of the type

$$\begin{cases} \frac{du}{dt} = u(1-v) - \varepsilon u^2, \\ \frac{dv}{dt} = \alpha v(u-1). \end{cases}$$

The equilibrium state, in which there is coexistence between the two species, is given by  $(1, 1-\varepsilon)$  which is a stable spiral. From a biological point of view, this means that a small perturbation can remove the oscillations, which is not very

realistic. This phenomenon, from a mathematical point of view, is called “structural instability”.

### 1.8.4 Theorem (Diffusion)

Consider the Lotka-Volterra system with diffusion

$$(1.62) \quad \begin{cases} u_t = u(1-v) - D\Delta u, \\ v_t = \alpha v(u-1) + D\Delta v, \\ \frac{\partial u}{\partial n} = \frac{\partial v}{\partial n} = 0, (x, t) \in \partial \Omega \times (0, T], \end{cases} \quad (\text{Neumann})$$

with  $\Omega \subseteq \mathbb{R}^n$  bounded,  $n$  external normal.

It can be shown that the system tends to a homogeneous spatial state for  $t \rightarrow +\infty$ .

### Proof

Let us consider the initial conditions

$$u(x, 0) = u_0(x), \quad v(x, 0) = v_0(x), \quad x \in \bar{\Omega}$$

and let

$$g(u, v) = \alpha(u - \log u) + v - \log v$$

be the “energy” of system as before. Obviously we have

$$\begin{aligned} g_t &= \alpha \left( u_t - \frac{u_t}{u} \right) + v_t - \frac{v_t}{v}, \\ \Delta g &= \alpha \left( \Delta u - \frac{\Delta u}{u} + \frac{|\nabla u|^2}{u^2} \right) + \Delta v - \frac{\Delta v}{v} + \frac{|\nabla v|^2}{v^2}, \end{aligned}$$

from which, by a direct calculation, one obtains

$$g_t - D\Delta g \leq 0,$$

which means that the energy is dissipated by diffusion.

The boundary conditions for  $g$  are

$$\frac{\partial g}{\partial n} = \alpha \left( \frac{\partial u}{\partial n} - \frac{1}{u} \frac{\partial u}{\partial n} \right) + \frac{\partial v}{\partial n} - \frac{1}{v} \frac{\partial v}{\partial n} = 0$$

while the initial conditions read

$$g(x, 0) = \alpha(u_0(x) - \log u_0(x)) + v_0(x) - \log v_0(x) =: g_0(x).$$

Let us consider now the total energy

$$G(t) = \int_{\Omega} g(x, t) d\mathbf{x},$$

then

$$\begin{aligned}\dot{G}(t) &= \int_{\Omega} g_t d\mathbf{x} = \int_{\Omega} D\Delta g - \alpha D \frac{|\nabla u|^2}{u^2} - D \frac{|\nabla v|^2}{v^2} d\mathbf{x} \\ &\stackrel{\text{Green}}{=} \int_{\partial\Omega} D \frac{\partial g}{\partial n} - \int_{\Omega} \alpha D \frac{|\nabla u|^2}{u^2} + D \frac{|\nabla v|^2}{v^2} d\mathbf{x} \leq 0,\end{aligned}$$

so  $G(t)$  is decreasing.

Furthermore, since  $z - \log z \geq 1$ ,  $\forall z > 0$ , we have

$$g(x, t) = \alpha(u - \log u) + v - \log v \geq \alpha + 1.$$

Therefore

$$G(t) = \int_{\Omega} g(\mathbf{x}, t) d\mathbf{x} \geq (\alpha + 1)|\Omega|,$$

from which

$$\lim_{t \rightarrow +\infty} G(t)$$

is finite. Consequently

$$\lim_{t \rightarrow +\infty} \dot{G}(t) = 0,$$

implies that

$$\nabla u \rightarrow 0, \nabla v \rightarrow 0, \text{ for } t \rightarrow +\infty.$$

So the system goes toward a spatially homogeneous situation. ■

### 1.8.5 Definition (Boundary operator)

The operator  $\mathbf{B}\mathbf{u}$  such that

$$(1.63) \quad (\mathbf{B}\mathbf{u})_i = c_i(\mathbf{x}, t)u_i(\mathbf{x}, t) + d_i(\mathbf{x}, t)\frac{\partial u_i}{\partial n_i}(\mathbf{x}, t),$$

is called the *boundary operator*, with  $c_i \geq 0, d_i \geq 0, c_i^2 + d_i^2 > 0$  and  $\frac{\partial}{\partial n_i}$  a certain

external derivative. We observe that  $\mathbf{B}$  is a diagonal operator.

### 1.8.6 Definition (Invariant set)

An *invariant set* for systems of the type

$$(1.64) \quad \mathbf{u}_t = \mathbf{f}(\mathbf{u}) + D\Delta\mathbf{u}$$

is a set  $\Sigma \subseteq \mathbb{R}^m$  such that if  $\mathbf{u}$  is solution of the system with initial values in  $\Sigma$  and boundary values in  $\mathbf{B}\mathbf{u} \in \mathcal{C}\Sigma := \{(c_1 v_1, \dots, c_n v_n)^T \mid \mathbf{v} \in \Sigma\}$ , then

$$\mathbf{u}(\mathbf{x}, t) \in \Sigma, \forall \mathbf{x} \in \Omega, \forall t \in [0, T].$$

### 1.8.7 Definition (Almost monotonicity)

- A function  $f(\mathbf{u})$  is said to be *almost monotonic increasing* if every  $f_i$  is increasing with respect to every  $u_j, j \neq i$ .
- A function  $f(\mathbf{u})$  is said to be *mixed almost monotonic* if every  $f_i$  is monotonic with respect to every  $u_j, j \neq i$ .

### Notation

In the following we will consider systems of the type

$$(P) \quad \begin{cases} N\mathbf{u} = L\mathbf{u} - f(\mathbf{u}) = \mathbf{u}_t - \mathbf{D}\Delta\mathbf{u} - f(\mathbf{u}) = 0, & \text{in } \Omega \times (0, T), \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}), & \text{on } \overline{\Omega} \times \{0\}, \\ \mathbf{B}\mathbf{u} = \mathbf{b}(\mathbf{x}, t), & \text{on } \partial\Omega \times (0, T], \end{cases}$$

where

- $\mathbf{D}$  = diagonal diffusion matrix;
- $\mathbf{B}$  = diagonal boundary operator.

### 1.8.8 Theorem (Comparison generalized)

Let  $\underline{\mathbf{u}}$  and  $\bar{\mathbf{u}}$  be a subsolution and a supersolution of the system **(P)**. Let  $f$  be globally Lipschitz and almost monotonically increasing, and let  $\Omega$  satisfy the inner sphere property.

1. The *weak comparison theorem* states that

$$(1.65) \quad \underline{\mathbf{u}} \leq \bar{\mathbf{u}} \quad \text{in } \overline{\Omega} \times [0, T],$$

where the inequalities have to be intended componentwise.

2. The *strong comparison theorem* states that

$$(1.66) \quad \underline{u}_i \leq \bar{u}_i, \forall i, \quad \text{in } \Omega \times (0, T],$$

or

$$\underline{u}_i = \bar{u}_i, \quad \text{in } \overline{\Omega} \times [0, t'] \text{ for a certain } t' \leq T.$$

### Proof

1. Let  $y$  be dependent on the parameter  $\varepsilon$ , in such that  $y \rightarrow 0$ , if  $\varepsilon \rightarrow 0$ .

Let  $\mathbf{Y} = (y, \dots, y)$  be a solution of the problem

$$\begin{cases} \dot{y} = k\|\mathbf{Y}\| = k\sqrt{n}y \\ y(0) = \varepsilon, \end{cases}$$

i.e.  $y(t) = \varepsilon^{k\sqrt{n}t}$ , where  $k$  is the Lipschitz constant of  $f$ .

Let us consider the vector

$$\mathbf{v} = \underline{\mathbf{u}} - \mathbf{Y}$$

and show that it is a subsolution of problem  $L\mathbf{u} = \mathbf{f}(\mathbf{u})$ .

We have

$$\begin{aligned} L_i v_i &\stackrel{\text{def}}{=} \frac{\partial}{\partial t} v_i - D_i \Delta v_i = \frac{\partial}{\partial t} \underline{u}_i - \dot{y} - D_i \Delta \underline{u}_i \\ &\leq f_i(\underline{\mathbf{u}}) - k \|\mathbf{Y}\| = f_i(\mathbf{v} + \mathbf{Y}) - f_i(\mathbf{v}) + f_i(\mathbf{v}) - k \|\mathbf{Y}\| \\ &\leq k \|(\mathbf{v} + \mathbf{Y} - \mathbf{v})\| + f_i(\mathbf{v}) - k \|\mathbf{Y}\| \\ &= k \|\mathbf{Y}\| + f_i(\mathbf{v}) - k \|\mathbf{Y}\| = f_i(\mathbf{v}). \end{aligned}$$

Now

$$\begin{aligned} \mathbf{v}(\mathbf{x}, 0) &= \underline{\mathbf{u}}(\mathbf{x}, 0) - \mathbf{Y}(0) \leq \underline{\mathbf{u}}_0(\mathbf{x}) - (\varepsilon, \dots, \varepsilon)^T \leq \underline{\mathbf{u}}_0(\mathbf{x}), \\ \mathbf{B}\mathbf{v} &= \mathbf{B}(\underline{\mathbf{u}} - \mathbf{Y}) = \mathbf{B}\underline{\mathbf{u}} - \mathbf{c} \mathbf{Y} \leq \mathbf{b}(\mathbf{x}, t). \end{aligned}$$

Then  $\mathbf{v} = \underline{\mathbf{u}} - \mathbf{Y}$  is a subsolution of  $L\mathbf{u} = \mathbf{f}(\mathbf{u})$ .

Let now be

$$\mathbf{w} = \bar{\mathbf{u}} - \mathbf{v} = \bar{\mathbf{u}} - \underline{\mathbf{u}} + \mathbf{Y},$$

we have

$$\begin{aligned} \mathbf{w}(\mathbf{x}, 0) &= \bar{\mathbf{u}}(\mathbf{x}, 0) - \underline{\mathbf{u}}(\mathbf{x}, 0) + \mathbf{Y}(0) \geq \mathbf{Y}(0) > 0, \\ \mathbf{B}\mathbf{w} &= \mathbf{B}\bar{\mathbf{u}} - \mathbf{B}\mathbf{v} + \mathbf{c} \mathbf{Y} \geq \mathbf{b} - \mathbf{b} + \mathbf{c} \mathbf{Y} \geq 0. \end{aligned}$$

Since  $\mathbf{w}(\mathbf{x}, 0) > 0$  then  $\exists \tau \in [0, T]$  s.t.

$$\mathbf{w}(\mathbf{x}, t) \geq 0 \quad \text{in } [0, \tau],$$

then

$$\begin{aligned} L_i w_i &= L_i \bar{u}_i - L_i \underline{u}_i + L_i y = L_i \bar{u}_i - L_i v_i \\ &\geq f_i(\bar{u}) - f_i(v) = f_i(\bar{u}_1, \dots, \bar{u}_i, \dots, \bar{u}_m) - f_i(v_1, \dots, v_i, \dots, v_m) \\ &\geq -K(\bar{u}_i - v_i) = -K w_i, \end{aligned}$$

which means that  $\mathbf{w}$  is a supersolution. Then, by the principle of the scalar maximum  $w_i > 0$ . If  $\tau \in [0, T]$  and  $\mathbf{w} > 0$  in  $[0, \tau]$  then  $\tau = T$ , i. e.

$$\mathbf{w} > 0 \quad \text{in } [0, T],$$

from which

$$\bar{\mathbf{u}} > \underline{\mathbf{u}} - \mathbf{Y},$$

and, since  $\mathbf{Y} \rightarrow 0$ , if  $\varepsilon \rightarrow 0$ , we have

$$\bar{\mathbf{u}} \geq \underline{\mathbf{u}} \quad \text{in } \overline{Q_T}.$$

2. By placing  $N_i(\mathbf{u}) = L_i u_i - f_i(\mathbf{u})$ ,  $z = \bar{\mathbf{u}} - \underline{\mathbf{u}}$ , one obtains

$$\begin{aligned}
0 &\leq N_i(\bar{\mathbf{u}}) - N_i(\underline{\mathbf{u}}) = L_i(\bar{\mathbf{u}}_i) - f_i(\bar{\mathbf{u}}_i) - L_i(\underline{\mathbf{u}}_i) + f_i(\underline{\mathbf{u}}_i) \\
&= L_i z_i - [f_i(\bar{\mathbf{u}}_1 \dots \bar{\mathbf{u}}_m) - f_i(\underline{\mathbf{u}}_1 \dots \underline{\mathbf{u}}_m)] \\
&\leq L_i z_i - [f_i(\bar{\mathbf{u}}_1 \dots \bar{\mathbf{u}}_i \dots \bar{\mathbf{u}}_m) - f_i(\underline{\mathbf{u}}_1 \dots \underline{\mathbf{u}}_i \dots \underline{\mathbf{u}}_m)] \\
&\leq L_i z_i - h_i(\bar{\mathbf{u}}_i - \underline{\mathbf{u}}_i) = L_i z_i - h_i z_i.
\end{aligned}$$

By the theorem of the strong scalar maximum it follows that the thesis holds true

■

### 1.8.9 Corollary (Uniqueness)

Under the assumptions of the comparison theorem, we have that the problem

$$(1.67) \quad \begin{cases} N\mathbf{u} = L\mathbf{u} - \mathbf{f}(\mathbf{u}) = \mathbf{u}_t - \mathbf{D}\Delta\mathbf{u} - \mathbf{f}(\mathbf{u}) = \mathbf{0}, \text{ on } \Omega \times (0, T), \\ \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}), \text{ on } \overline{\Omega} \times \{0\}, \\ \mathbf{B}\mathbf{u} = \mathbf{b}(\mathbf{x}, t), \text{ on } \partial\Omega \times (0, T), \end{cases}$$

has a unique solution.

### Proof

Suppose that  $\mathbf{v}, \mathbf{w}$  are both solutions, then by the comparison theorem we have  $\mathbf{v} \geq \mathbf{w}$  and  $\mathbf{w} \geq \mathbf{v}$ , i.e.

$$\mathbf{v} = \mathbf{w}.$$

■

### 1.8.10 Corollary (Invariant set)

Under the assumptions of the comparison theorem, let  $\mathbf{a}, \mathbf{b}$  be two vectors s.t.  $-\infty < a < b < +\infty$  with  $\mathbf{f}(\mathbf{a}) \geq \mathbf{0}$ ,  $\mathbf{f}(\mathbf{b}) \leq \mathbf{0}$ , then

$$(1.68) \quad \Sigma = \{\mathbf{u} \mid \mathbf{a} \leq \mathbf{u} \leq \mathbf{b}\}$$

is an invariant set for (1.67).

### Proof

If  $\mathbf{f}(\mathbf{a}) \geq \mathbf{0}$ , then  $N(\mathbf{a}) = -\mathbf{f}(\mathbf{a}) \leq \mathbf{0}$ , i. e.  $\mathbf{a}$  is a subsolution of (1.67).

If  $\mathbf{f}(\mathbf{b}) \leq \mathbf{0}$ , then  $N(\mathbf{b}) = -\mathbf{f}(\mathbf{b}) \geq \mathbf{0}$ , i. e.  $\mathbf{b}$  is a supersolution of (1.67).

Let  $\mathbf{u}(\mathbf{x}, t)$  be a solution of (1.67) then we have that if

$$\mathbf{a} \leq \mathbf{u}(\mathbf{x}, 0) = \mathbf{u}_0(\mathbf{x}) \leq \mathbf{b},$$

then

$$a \leq u(x, t) \leq b, \quad \forall t \in (0, T),$$

hence the thesis. ■

### 1.8.11 Theorem (Global existence)

Under the assumptions of the comparison theorem, let  $a, b$  s.t.  $-\infty < a < b < +\infty$ ,  $f(a) \geq 0$  and  $f(b) \leq 0$ , then the problem

$$(1.69) \quad \begin{cases} Nu = Lu - f(u) = u_t - D\Delta u - f(u) = 0, & \text{in } \Omega \times (0, T), \\ u(x, 0) = u_0(x), & \text{in } \overline{\Omega} \times 0, \\ Bu = b(x, t), & \text{in } \partial \Omega \times (0, T], \end{cases}$$

has a solution on  $[0, T]$ .

### Proof (Sketch)

Let  $\{v^n\}, \{w^n\}$  be defined as

$$\begin{aligned} v^0 &= a, \quad (L_i + k)v_i^n = f_i(v^{n-1}) + k v_i^{n-1}, \\ w^0 &= b, \quad (L_i + k)w_i^n = f_i(w^{n-1}) + k w_i^{n-1}. \end{aligned}$$

By induction, it can be shown that  $v^{n+1} \geq v^n$  and  $w^{n+1} \leq w^n$ .

Then we have that  $\{v^n\}, \{w^n\}$  are monotone sequences and are also bounded, so they converge; moreover, the limit is equal. We have constructed two monotone sequences that both converge to the solution of the problem, which can be calculated numerically. ■

### 1.8.12 Definition (Closed operator)

Let  $A : D(A) \subseteq H \rightarrow H$  be an operator, with  $H$  a Hilbert space. It will be called a *closed operator* if

$$\forall x_n \in D(A) \text{ s.t. } x_n \rightarrow x \text{ and } Ax_n \rightarrow y,$$

it holds that

$$x \in D(A) \text{ and } Ax_n \rightarrow y = Ax.$$

### 1.8.13 Corollary (of the comparison theorem)

Let  $f : \mathbb{R}^m \rightarrow \mathbb{R}^m$  be almost monotonically increasing and globally Lipschitz such that

$$\underline{f}(\mathbf{v}) \leq f(\mathbf{v}), \quad \forall \mathbf{v} \in \mathbb{R}^m.$$

Let  $\underline{\mathbf{u}}$  be a subsolution of **(P)**, with  $f$  substituted by  $\underline{f}$  (let us recall this problem **(P)**), and  $\bar{\mathbf{v}}$  a supersolution of **(P)**, then

$$\underline{\mathbf{u}} \leq \bar{\mathbf{v}} \quad \text{in } \overline{Q_T}.$$

Similarly, one can define  $\bar{f}$  for which there is corresponding property.

### Proof

Since  $\underline{\mathbf{u}}$  is subsolution of **(P)**, we have that

$$L(\underline{\mathbf{u}}) \leq \underline{f}(\underline{\mathbf{u}}) \leq f(\underline{\mathbf{u}}),$$

i.e.  $\underline{\mathbf{u}}$  is subsolution of **(P)**, then

$$\underline{\mathbf{u}} \leq \bar{\mathbf{v}} \quad \text{in } \overline{Q_T}.$$

■

#### 1.8.14 Lemma (on the almost monotonicity)

Let  $\underline{f}$  be defined by

$$(1.71) \quad \underline{f}_i(\mathbf{u}) := \inf_{\{\mathbf{v} | \mathbf{u} \leq \mathbf{v} \leq \bar{\mathbf{u}}, v_i = u_i\}} f_i(\mathbf{v}),$$

with  $\bar{\mathbf{u}}$  supersolution of **(P)**.

Then

$$\underline{f}(\mathbf{u}) \leq f(\mathbf{u}), \quad \forall \mathbf{u} \leq \bar{\mathbf{u}},$$

and  $\underline{f}$  is almost monotonically increasing and is uniformly Lipschitzian with the same constant as  $f$ . Analogously, it is possible to construct

$$\bar{f}(\mathbf{u}) \geq f(\mathbf{u}), \quad \forall \mathbf{u} \geq \underline{\mathbf{u}}.$$

#### 1.8.15 Theorem (invariant set)

Let  $f$  be locally Lipschitzian and  $Q_T$  satisfy the inner sphere property. Let  $\mathbf{a}, \mathbf{b}$  be such that

$$-\infty < a < b < +\infty,$$

$$f_i(\mathbf{v}) \geq 0, \quad \text{for } \mathbf{a} \leq \mathbf{v} \leq \mathbf{b}, v_i = a_i,$$

$$f_i(\mathbf{v}) \leq 0, \quad \text{for } \mathbf{a} \leq \mathbf{v} \leq \mathbf{b}, v_i = b_i.$$

Then  $\Sigma := \{\mathbf{v} | \mathbf{a} \leq \mathbf{v} \leq \mathbf{b}\}$  is invariant for **(P)**.

### 1.8.16 Theorem (uniqueness of the solution)

Under the assumptions of the previous theorem, we can consider (P) with initial conditions  $\mathbf{u}_0 \in \Sigma$  and boundary conditions  $\mathbf{B}\mathbf{u} \in \Sigma$ . Then (P) has only one global solution.

#### Example

Consider the system for two populations competing between them

$$\begin{cases} \frac{\partial u_1}{\partial t} = u_1(I - u_1 - \alpha u_2) + D_1 \Delta u_1 = f_1(u_1, u_2) + D_1 \Delta u_1 \\ \frac{\partial u_2}{\partial t} = \vartheta u_2(I - \beta u_1 - u_2) + D_2 \Delta u_2 = f_2(u_1, u_2) + D_2 \Delta u_2 \end{cases}$$

such that  $\alpha, \vartheta, \beta > 0$  and with

$$\begin{cases} \frac{\partial u_1}{\partial n} = \frac{\partial u_2}{\partial n} = 0, \\ u_1(\mathbf{x}, 0) = u_1^0(\mathbf{x}), \\ u_2(\mathbf{x}, 0) = u_2^0(\mathbf{x}). \end{cases}$$

If  $\mathbf{0} \leq (u_1^0, u_2^0) \leq \mathbf{M} = (M_1, M_2)$ , then the system has a global in time solution that satisfies

$$\mathbf{0} \leq (u_1, u_2) \leq \mathbf{b} = (b_1, b_2),$$

where  $b_1 = \max\{M_1, I\}$ ,  $b_2 = \max\{M_2, I\}$ .

#### Proof

The assumptions of theorem 1.8.13 are satisfied with  $\mathbf{a} = (0, 0)$  and  $\mathbf{b} = (b_1, b_2)$ . In fact

$$\begin{aligned} f_1(0, u_2) &= 0 \Rightarrow f_1(a_1, u_2) \geq 0, \\ f_2(u_1, 0) &= 0 \Rightarrow f_2(u_1, a_2) \geq 0, \\ f_1(b_1, u_2) &= b_1(I - b_1 - \alpha u_2) \leq 0 \Rightarrow f_1(b_1, u_2) \leq 0, \\ f_2(u_1, b_2) &= \vartheta b_2(I - b_1 - b_2) \leq 0 \Rightarrow f_2(u_1, b_2) \leq 0. \end{aligned}$$

Furthermore, we have that

$$\frac{\partial \mathbf{a}}{\partial n} \leq \frac{\partial \mathbf{u}}{\partial n} \leq \frac{\partial \mathbf{b}}{\partial n}, \text{ with } \frac{\partial \mathbf{a}}{\partial n}, \frac{\partial \mathbf{u}}{\partial n}, \frac{\partial \mathbf{b}}{\partial n} = 0, \text{ and } \mathbf{0} \leq (u_1^0, u_2^0) \leq \mathbf{M} \leq \mathbf{b}.$$

■

## 1.9 Compartmental models

### 1.9.1 Diffusion of rabies

Travelling wave solutions also occur for systems. Consider the case of the spread of rabies [7], which is a viral disease of the central nervous system that affects warm-blooded animals, including humans. There is no cure for rabies, but there are very effective vaccines. Let us consider a fox population consisting of

- $S$  = susceptible foxes;
- $I$  = infected foxes.

We get the equations

$$(1.72) \quad \begin{cases} \frac{\partial S}{\partial t} = -r S I, \\ \frac{\partial I}{\partial t} = r S I - a I + D \frac{\partial^2 I}{\partial x^2}, \end{cases}$$

where

- $r$ : infection rate;
- $a$ : mortality rate;
- $D$ : diffusion coefficient of infected foxes.

By introducing the quantities

$$S^* = \frac{S}{S_0}, \quad I^* = \frac{I}{S_0}, \quad x^* = \sqrt{\frac{r S_0}{D}} x, \quad t^* = r S_0 t, \quad m = \frac{a}{r S_0},$$

we can adimensionalize the system and obtain

$$(1.73) \quad \begin{cases} \frac{\partial S}{\partial t} = -S I, \\ \frac{\partial I}{\partial t} = S I - m I + \frac{\partial^2 I}{\partial x^2}, \end{cases}$$

where the \*'s have been omitted, for simplicity.

Looking for travelling wave solutions

$$(1.74) \quad \begin{cases} S = S(x - ct), \\ I = I(x - ct), \end{cases}$$

we obtain

$$(1.75) \quad \begin{cases} c S' = S I, \\ I'' + c I' + (S - m) I = 0. \end{cases}$$

Let us assume that before the epidemic there were no infected, that is

$$S(+\infty)=1, I(+\infty)=0.$$

It is also expected that after the passage of the wave

$$S'(-\infty)=0, I(-\infty)=0.$$

Since  $c S' = S I$ , then  $I = c \frac{S'}{S}$  and substituting in (1.75)<sub>2</sub> we get

$$I'' + c I' + c S' - m c \frac{S'}{S} = 0.$$

By integrating, it is obtained

$$(1.76) \quad I' + c I + c S - m c \log S = cost.$$

Applying the condition at  $+\infty$  gives  $cost = c$ . Applying the condition at  $-\infty$  gives

$$c S(-\infty) - c m \log(S(-\infty)) = c,$$

from which

$$(1.77) \quad \frac{S(-\infty) - 1}{\log[S(-\infty)]} = m = \frac{a}{r S_0},$$

from which we derive the fraction of survivors, as we can see in Figure 1.6.

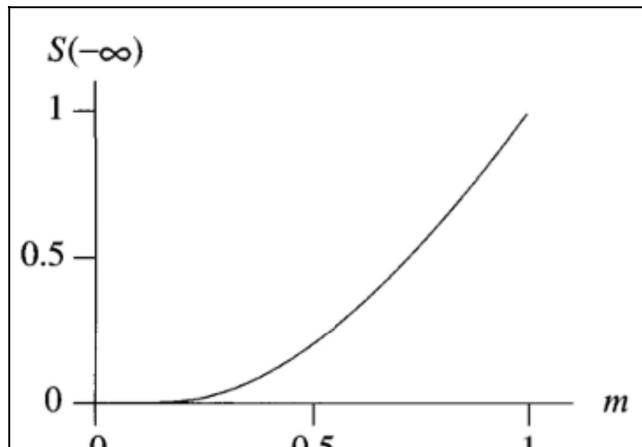


Figure 1.6: Fraction of survivors as a function of  $m$

The parameter  $m$  is a measure of severity, the smaller it is, the smaller the number of survivors. We note that if  $m > 1$ , then  $a > r S_0$ , that is, the foxes die before they can infect anyone, so the epidemic does not spread. Furthermore, the quantity  $S_0 = \frac{a}{r}$  also provides a minimum quantity of susceptibles for the epidemic to spread.

Let us now consider

$$(1.78) \quad \begin{cases} S' = \frac{1}{c} S I, \\ I' = -c(S+I) + c m \log S + c, \end{cases}$$

and look for equilibrium states  $S'=0, I'=0$ , that is

$$(1.79) \quad \begin{cases} \frac{1}{c} S I = 0, \\ -c(S+I) + c m \log S + c = 0, \end{cases}$$

which has solutions

$$P_1 = (S(-\infty), 0)^T, P_2 = (1, 0)^T.$$

The Jacobian matrix is

$$J = \begin{pmatrix} \frac{1}{c} I & \frac{1}{c} S \\ -c + \frac{mc}{S} & -c \end{pmatrix}.$$

Evaluating in  $P_1$  we have

$$J(P_1) = \begin{pmatrix} 0 & \frac{1}{c} S(-\infty) \\ -c + \frac{mc}{S(-\infty)} & -c \end{pmatrix},$$

whose eigenvalues are

$$(1.80) \quad \lambda_{1,2} = \frac{-c \pm \sqrt{c^2 + (m - S(-\infty))}}{2}, \text{ with } m > S(-\infty),$$

therefore  $P_1$  is a saddle. For  $c < 2\sqrt{1-m}$ , the point  $P_2$  is a stable focus, so there can be no heterocline orbit for which  $S$  is positive. For  $c \geq 2\sqrt{1-m}$ , the point  $P_2$  is a stable node and therefore one can have a heterocline orbit joining  $P_1$  and  $P_2$ .

Therefore

$$(1.81) \quad c = 2\sqrt{1-m}$$

is the *minimum wave velocity*.

### 1.9.2 Model improvement

The model does not take into account growth terms in the equation for  $S$ . Therefore, the following modification should be made

$$(1.82) \quad \begin{cases} \frac{\partial S}{\partial t} = -rS I + BS \left(1 - \frac{S}{S_0}\right), \\ \frac{\partial I}{\partial t} = rS I - aI + D \frac{\partial^2 I}{\partial x^2}, \end{cases}$$

where

- $B$ : growth rate;
- $S_0$ : carrying capacity.

We can adimensionalize the system, by posing

$$S^* = \frac{S}{S_0}, \quad I^* = \frac{rI}{B}, \quad x^* = \sqrt{\frac{B}{D}}x, \quad t^* = Bt,$$

from which

$$(1.83) \quad \begin{cases} \frac{\partial S}{\partial t} = -SI + S(1-S), \\ \frac{\partial I}{\partial t} = bI(S-m) + \frac{\partial^2 I}{\partial x^2}, \end{cases}$$

$$\text{where } m = \frac{a}{rS_0}, \quad b = \frac{rS_0}{B}.$$

We now search for travelling wave solutions

$$(1.84) \quad \begin{cases} S = S(z) = S(x-ct), \\ I = I(z) = I(x-ct), \end{cases}$$

We obtain the system

$$(1.85) \quad \begin{cases} cS' = SI - S(1-S), \\ cI' = -I'' - bI(S-m), \end{cases}$$

and placing  $V = I'$  we have

$$(1.86) \quad \begin{cases} S' = \frac{1}{c}S(I-1+S), \\ I' = V, \\ V' = -cV - bI(S-m). \end{cases}$$

Obviously, it is expected that

$$(S, I, V)(+\infty) = (1, 0, 0)^T.$$

We note that the homogeneous system

$$\begin{cases} \frac{1}{c}S(I-1+S)=0, \\ V=0, \\ -cV-bI(S-m)=0, \end{cases}$$

has solutions

$$\begin{cases} S=1, \\ V=0, \\ I=0, \end{cases} \text{ and } \begin{cases} I=1-m, \\ V=0, \\ S=m, \end{cases}$$

from which

$$(S, I, V)(-\infty) = (m, 1-m, 0)^T.$$

The Jacobian matrix is

$$J = \begin{pmatrix} -\frac{1-2S-I}{c}I & \frac{1}{c}S & 0 \\ 0 & 0 & 1 \\ -bI & -b(S-m) & -c \end{pmatrix}.$$

Evaluating in  $P_1 = (1, 0, 0)^T$  we have

$$J(P_1) = \begin{pmatrix} \frac{1}{c} & \frac{1}{c} & 0 \\ 0 & 0 & 1 \\ 0 & -b(1-m) & -c \end{pmatrix},$$

that has the following eigenvalues

$$(1.87) \quad \lambda_1 = \frac{1}{c}, \quad \lambda_{2,3} = \frac{-c \pm \sqrt{c^2 - 4b(1-m)}}{2}.$$

One must have

$$c \geq 2\sqrt{b(1-m)} = :c.$$

Then the stable variety has dimension two, the unstable one.

Evaluating instead in  $P_2 = (m, 1-m, 0)^T$  we have

$$J(P_2) = \begin{pmatrix} \frac{m}{c} & \frac{m}{c} & 0 \\ 0 & 0 & 1 \\ -b(1-m) & 0 & -c \end{pmatrix}.$$

The eigenvalues are therefore given by the roots of the characteristic equation

$$P_b(\lambda) = \lambda^3 + \left(c - \frac{m}{c}\right)\lambda^2 - m\lambda + \frac{b}{c}m(1-m).$$

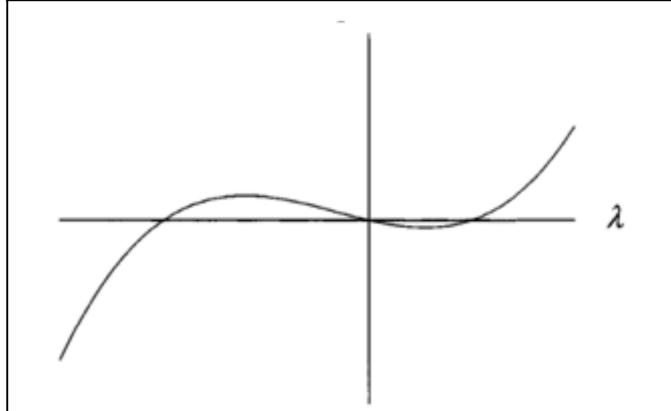


Figure 1.7: Graph of the characteristic polynomial for  $b = 0$

In the case where  $b=0$ , it has eigenvalues

$$(1.88) \quad \lambda_1 = 0, \lambda_{2,3} = \frac{-(c-m/c) \pm \sqrt{(c-m/c)^2 + 4m}}{2}.$$

We want to find  $b^*$  such that  $P_{b^*}$  has a local minimum equal to zero. Therefore, being

$$(1.89) \quad P_b'(\lambda) = 3\lambda^2 + 2\left(c - \frac{m}{c}\right)\lambda - m,$$

one obtains

$$(1.90) \quad \lambda = \frac{m/c - c \pm \sqrt{(m/c - c)^2 + 3m}}{3}.$$

The minimum is in

$$(1.91) \quad \lambda^* = \frac{m/c - c \pm \sqrt{(m/c - c)^2 + 3m}}{3},$$

and imposing  $P_b(\lambda^*)=0$ , we have

$$(1.92) \quad b^* = -\frac{c}{m(1-m)} P_0(\lambda^*).$$

The two real roots  $\lambda_{2,3}$  as  $b$  increases, will become complex with positive real part, then the point  $P_2$  has an unstable manifold of dimension 2. The trajectory joining  $P_1$  and  $P_2$  is given by the intersection of the unstable manifold of  $P_2$  and the stable manifold of  $P_1$ . The trajectory satisfies

$$\begin{aligned} \phi(z) &\rightarrow (m, 1-m, 0), \text{ for } z \rightarrow -\infty, \\ \phi(z) &\rightarrow (1, 0, 0), \text{ for } z \rightarrow +\infty. \end{aligned}$$

### 1.9.3 Spatial distribution patterns of grey and red squirrels

In the early 20th century, North American gray squirrels were introduced to various parts of Great Britain, particularly the south-east, and they spread quickly. Concurrently, native red squirrels disappeared from these localities after only a few years of coexistence. Both red and gray squirrels can breed twice a year, but red squirrels, which are smaller, rarely have more than two or three offspring per litter, while gray squirrels often have four or five pups. In addition, gray squirrels have another competitive advantage in that they have a body weight that is almost twice as high as red squirrels. In separate environments, the two species of squirrels show quite similar social organization and feeding patterns, so it is natural that placed in the same habitat they would compete for resource exploitation and are unlikely to coexist.

Let us define

- $S_1(x, t)$ : population density of grey squirrels;
- $S_2(x, t)$ : population density of red squirrels.

A model describing the competition between the two species is given by

$$(1.93) \quad \begin{cases} \frac{\partial S_1}{\partial t} = D_1 \Delta S_1 + a_1 S_1 (1 - b_1 S_1 - c_1 S_2), \\ \frac{\partial S_2}{\partial t} = D_2 \Delta S_2 + a_2 S_2 (1 - b_2 S_2 - c_2 S_1), \end{cases}$$

where

- $a_i$ : growth rates;
- $\frac{1}{b_i}$ : carrying capacity;
- $c_i$ : coefficient of competition;
- $D_i$ : diffusion coefficient.

As a first step, we switch to a dimensionless system, making a change of variables and defining new constants:

$$(1.94) \quad \begin{aligned} \theta_i &= b_i S_i, \quad i=1,2, \quad t^* = a_1 t, \quad \mathbf{x}^* = \sqrt{a_1/D} \mathbf{x}, \\ \gamma_1 &= c_1/b_2, \quad \gamma_2 = c_2/b_1, \quad k = D_2/D_1, \quad \alpha = a_2/a_1, \end{aligned}$$

we obtain the system

$$(1.95) \quad \begin{cases} \frac{\partial \theta_1}{\partial t} = \Delta \theta_1 + \theta_1(1 - \theta_1 - \gamma_1 \theta_2), \\ \frac{\partial \theta_2}{\partial t} = k \Delta \theta_2 + \alpha \theta_2(1 - \theta_2 - \gamma_2 \theta_1), \end{cases}$$

where, as usual, we have omitted the asterisks for simplicity of notation.

In the absence of diffusion, the equilibrium states are  $(0,0)$ ,  $(1,0)$ ,  $(0,1)$ . Considering diffusion, it can be shown that there is a trajectory joining  $(1,0)$  and  $(0,1)$ , which corresponds to the situation in which grey squirrels drive red squirrels to extinction. Considering the one-dimensional case, we look for wave solutions of (1.95) of the form

$$(1.96) \quad \theta_i = \theta_i(z) = \theta_i(x - ct),$$

we obtain then

$$(1.97) \quad \begin{cases} \frac{\partial^2 \theta_1}{\partial z^2} + c \frac{\partial \theta_1}{\partial z} + \theta_1(1 - \theta_1 - \gamma_1 \theta_2) = 0, \\ k \frac{\partial^2 \theta_2}{\partial z^2} + c \frac{\partial \theta_2}{\partial z} + \alpha \theta_2(1 - \theta_2 - \gamma_2 \theta_1) = 0, \end{cases}$$

with conditions

$$\begin{cases} \theta_1(-\infty) = 1, \theta_2(-\infty) = 0, \\ \theta_1(+\infty) = 0, \theta_2(+\infty) = 1, \end{cases}$$

implying that the grey squirrels drive out the red ones.

Analytical results for (1.97) can be obtained in the case  $k = \alpha = 1$ ,  $\gamma_1 + \gamma_2 = 2$ .

Adding up the two equations and setting  $\theta = \theta_1 + \theta_2$ , we obtain that

$$(1.98) \quad \frac{\partial^2 \theta}{\partial z^2} + c \frac{\partial \theta}{\partial z} + \theta(1 - \theta) = 0,$$

which is a FKPP equation, but different boundary conditions

$$\theta(\pm\infty) = 1,$$

suggesting that

$$\theta(z) = \theta_1(z) + \theta_2(z) = 1, \forall z.$$

The first equation of (1.97) can then be rewritten as

$$(1.99) \quad \frac{\partial^2 \theta_1}{\partial z^2} + c \frac{\partial \theta_1}{\partial z} + \theta_1(1 - \gamma_1)(1 - \theta_1) = 0.$$

It admits travelling wave solutions for  $\theta_1$  if

$$c \geq 2\sqrt{1-\gamma_1}.$$

Similarly, since  $\gamma_1 = 2 - \gamma_2$ , we have that

$$c \geq 2\sqrt{\gamma_2 - 1}.$$

Returning to dimensional quantities we obtain

$$(1.100) \quad c \geq c_{min} = 2\sqrt{a_1 D_1 \left(1 - \frac{c_1}{b_2}\right)}.$$

#### 1.9.4 Spatial pattern formation

Alan Turing was the first to show that, under certain conditions, “individuals” of various nature can react and spread in such a way as to produce stable distributions that are spatially non-homogeneous (patterns). Let us consider two species with densities  $u, v$  respectively, then we have the system

$$(1.101) \quad \begin{cases} \frac{\partial u}{\partial t} = D_u \Delta u + f(u, v), \\ \frac{\partial v}{\partial t} = D_v \Delta v + g(u, v), \end{cases}$$

Turing assumed that in the absence of diffusion ( $D_u = D_v = 0$ ),  $u$  and  $v$  approach over time a spatially homogeneous steady state  $(u_0, v_0)$  that is asymptotically stable for the linearised system

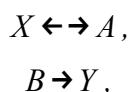
$$(1.102) \quad \begin{cases} \frac{\partial u}{\partial t} = f_u(u_0, v_0)(u - u_0) + f_v(u_0, v_0)(v - v_0), \\ \frac{\partial v}{\partial t} = g_u(u_0, v_0)(u - u_0) + g_v(u_0, v_0)(v - v_0). \end{cases}$$

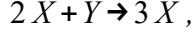
Under appropriate conditions the state  $(u_0, v_0)$  becomes unstable due to diffusion and the solution evolves towards a spatially inhomogeneous stable state. Since diffusion is usually stabilising, Turing's demonstration was surprising.

Let us consider the case where

$$(1.103) \quad \begin{cases} f(u, v) = k_1 - k_2 u + k_3 u^2 v, \\ g(u, v) = k_4 - k_3 u^2 v, \end{cases}$$

which is known as the Schnakenberg, with constants  $k_j > 0$ . This term describes a chemical reaction of the type





in which the concentration of the substances A and B is kept constant. The associated mathematical model is obtained by the law of mass action.

A second system is given by the activation-inhibition mechanism introduced by Gierer-Meinhardt in 1972, in which the reaction term is

$$(1.104) \quad \begin{cases} f(u, v) = k_1 - k_2 u + k_3 \frac{u^2}{v}, \\ g(u, v) = k_4 u^2 - k_5 v. \end{cases}$$

The model describes the evolution of the concentration,  $u$ , of a short-range autocatalytic substance, the activator, which regulates the production of the long-range antagonistic substance, the inhibitor, having concentration  $v$ .

Another system is that of Thomas, that is

$$(1.105) \quad \begin{cases} f(u, v) = k_1 - k_2 u - R(u, v), \\ g(u, v) = k_3 - k_4 v - R(u, v), \end{cases}$$

where  $R(u, v) = k_5 \frac{uv}{k_6 + k_7 u + k_8 v}$ , which describes the reaction of the substrates oxygen, of concentration  $u$ , and uric acid, of concentration  $v$ , in the presence of the enzyme urease.

For example, we can adimensionalize the system (1.101) with the reaction terms (1.103). Let  $L$  be a typical length and set

$$(1.106) \quad \begin{aligned} u^* &= u \left( \frac{k_3}{k_2} \right)^{\frac{1}{2}}, \quad v^* = v \left( \frac{k_3}{k_2} \right)^{\frac{1}{2}}, \quad t^* = D_u \frac{t}{L^2}, \quad x^* = \frac{x}{L}, \\ d &= \frac{D_v}{D_u}, \quad a = \frac{k_1}{k_2} \left( \frac{k_3}{k_2} \right)^{\frac{1}{2}}, \quad b = \frac{k_4}{k_2} \left( \frac{k_3}{k_2} \right)^{\frac{1}{2}}, \quad \gamma = \frac{L^2 k_2}{D_u}. \end{aligned}$$

Omitting the asterisks, the dimensionless equations are written as

$$(1.107) \quad \begin{cases} u_t = \Delta_u + \gamma f(u, v), \\ v_t = d \Delta_v + \gamma g(u, v), \end{cases}$$

where  $f(u, v) = a - u + u^2 v$  and  $g(u, v) = b - u^2 v$ .

### Remark

The  $\gamma$  parameter is such that

$$1. \quad \gamma^{1/2} \propto L, \text{ for } n=1;$$

2.  $\gamma^{1/2} \propto L^2$ , for  $n=2$ ;
3.  $\gamma$  represents the relative strength of the reaction term, with respect to the diffusive one;
4. increasing  $\gamma$ , is equivalent to decreasing  $d$ .

### 1.9.5 Definition (Turing instability)

A system exhibits Turing instability if a steady state of the homogeneous system without diffusion is stable for small perturbations, but is unstable with respect to spatially non-homogeneous perturbations when diffusion is present.

### 1.9.6 Necessary and sufficient condition for Turing instability

Let us consider the system

$$(1.108) \quad \begin{cases} \frac{\partial u}{\partial t} = \Delta u + \gamma f(u, v), \\ \frac{\partial v}{\partial t} = d \Delta v + \gamma g(u, v), \end{cases}$$

with conditions

$$(1.109) \quad \begin{cases} u(x, 0) = u_0(x), x \in \Omega, \\ v(x, 0) = v_0(x), x \in \Omega, \\ \frac{\partial u}{\partial n} = \frac{\partial v}{\partial n} = 0, x \in \partial \Omega. \end{cases}$$

In the stationary case without diffusion, we must have asymptotic stability. Let now be the Jacobian matrix

$$A = \begin{pmatrix} f_u & f_v \\ g_u & g_v \end{pmatrix} \Big|_{(u_0, v_0)}.$$

To have stability, we must impose

$$(1.110) \quad \begin{cases} Tr(A) = f_u(u_0, v_0) + g_v(u_0, v_0) < 0, \\ det(A) = f_u(u_0, v_0)g_v(u_0, v_0) - f_v(u_0, v_0)g_u(u_0, v_0) > 0. \end{cases}$$

In fact, the eigenvalues of  $A$  are

$$(1.111) \quad \lambda_{1,2} = \frac{Tr(A) \pm \sqrt{(Tr(A))^2 - 4 \det(A)}}{2},$$

which implies  $Tr(A) < 0, \det(A) > 0$ . Now, let us define

$$\mathbf{w} = \begin{pmatrix} u - u_0 \\ v - v_0 \end{pmatrix},$$

and consider the complete linearized system

$$(1.112) \quad \mathbf{w}_t = \mathbf{D} \Delta \mathbf{w} + \gamma \mathbf{A} \mathbf{w},$$

with

$$\mathbf{D} = \begin{pmatrix} 1 & 0 \\ 0 & d \end{pmatrix}.$$

Let  $\mathbf{W}(\mathbf{x})$  be a time-independent solution, that is such that

$$(1.113) \quad \begin{cases} \mathbf{D} \Delta \mathbf{W} + k^2 \mathbf{W} = 0, \\ \frac{\partial \mathbf{W}}{\partial n} = 0, \end{cases}$$

with  $k^2$  eigenvalue of the Laplacian.

If the domain is  $[0, a] \times [0, b]$ , it can be seen that eigenvalues and eigenvectors are

$$(1.114) \quad k^2 = \lambda_{mn} = \pi^2 \left( \frac{n^2}{a^2} + \frac{m^2}{b^2} \right),$$

$$(1.115) \quad \mathbf{W}_{m,n} = \mathbf{C}_{m,n} \cos\left(\frac{\pi n x}{a}\right) \cos\left(\frac{\pi m x}{b}\right), \quad \mathbf{C} = \begin{pmatrix} C_1 \\ C_2 \end{pmatrix}.$$

We look for solutions such as

$$(1.116) \quad \mathbf{W} = \exp(\lambda t) \mathbf{W}_k(\mathbf{x}),$$

where  $\mathbf{W}_k$  is the eigenfunction corresponding to the wave number  $k$ , and substituting in (1.112) we obtain

$$\lambda \mathbf{W}_k = \mathbf{D} \Delta \mathbf{W}_k + \gamma \mathbf{A} \mathbf{W}_k = -k^2 \mathbf{D} \mathbf{W}_k + \gamma \mathbf{A} \mathbf{W}_k,$$

which implies

$$\lambda \mathbf{W}_k + k^2 \mathbf{D} \mathbf{W}_k - \gamma \mathbf{A} \mathbf{W}_k = 0,$$

from which

$$|\lambda + k^2 \mathbf{D} - \gamma \mathbf{A}| = 0.$$

Calculating the determinant, we obtain

$$(1.117) \quad \lambda^2 + \lambda [(1+d)k^2 - \gamma(f_u + g_v)] + h(k^2) = 0,$$

where

$$(1.118) \quad h(k^2) = dk^4 - \gamma k^2 (df_u + g_v) + \gamma^2 \det(\mathbf{A}).$$

Since (1.112) is linear, its solutions will be of the form

$$(1.119) \quad \mathbf{w}(\mathbf{x}, t) = \sum_k c_k \exp(\lambda(k^2)t) \mathbf{W}_k(\mathbf{x}).$$

We have already imposed that the stationary state is stable in the absence of spatial effects, that is we have that  $\operatorname{Re}(\lambda(k^2=0))<0$ . Now, for the stationary state to be unstable for spatially non-homogeneous perturbations it must be

$$\operatorname{Re}(\lambda(k^2>0))>0, \text{ for some } k>0.$$

This can occur if and only if  $h(k^2)<0$ , that is, if and only if  $d f_u + g_v > 0$ , but since  $f_u + g_v = \operatorname{Tr}(A) < 0$ , we need  $d \neq 1$  and  $f_u \cdot g_v < 0$ .

### 1.9.7 Theorem (necessary condition)

A necessary condition to have Turing instability is that

$$(1.120) \quad D_u \neq D_v \text{ and } f_u \cdot g_v < 0.$$

#### Remark

These conditions are necessary, but not sufficient to have  $\operatorname{Re} \lambda > 0$ . To make sure that this happens, we must have that  $\min_k \{h(k)\} < 0$ .

Let us introduce  $s=k^2$ , then, by considering the parabola,

$$(1.121) \quad h(s) = d s^2 - \gamma(d f_u + g_v)s + \gamma^2 \det(A),$$

and calculating its derivative

$$h'(s) = 2d s - \gamma(d f_u + g_v) = 0,$$

we observe that the minimum is taken for

$$(1.122) \quad s_{\min} = \frac{\gamma(d f_u + g_v)}{2d},$$

and is equal to

$$(1.123) \quad h(s_{\min}) = \gamma \left[ \det(A) - \frac{(d f_u + g_v)^2}{4d} \right].$$

Therefore there must be

$$\frac{(d f_u + g_v)^2}{4d} > \det(A),$$

The bifurcation (that is, the change in behaviour) occurs for

$$(1.124) \quad \frac{(d f_u + g_v)^2}{4d} = \det(A),$$

which occurs for a certain critical coefficient  $d_c > 1$ .

Correspondingly, in (1.121), we derive the critical wave number

$$(1.125) \quad k_c = \gamma^{\frac{1}{2}} \left[ \frac{\det(\mathbf{A})}{d_c} \right]^{\frac{1}{4}}$$

For  $d > d_c$ ,  $h(k^2)$  has two zeros:

- $\underline{k^2} = \frac{\gamma}{2d} [(df_u + g_v) - \sqrt{(df_u + g_v)^2 - 4d\det(\mathbf{A})}]$ ;
- $\overline{k^2} = \frac{\gamma}{2d} [(df_u + g_v) + \sqrt{(df_u + g_v)^2 - 4d\det(\mathbf{A})}]$ .

### Remark

To have real solutions, one must have

$$(1.126) \quad (df_u + g_v)^2 - 4d(f_u g_v - f_v g_u) > 0.$$

### 1.9.8 Theorem (sufficient condition)

Conditions of the type

- $D_u \neq D_v$ ,
- $f_u \cdot g_v < 0$ ,
- $df_u + g_v > 0, f_u + g_v < 0$ ,
- $(df_u + g_v)^2 - 4d(f_u g_v - f_v g_u) > 0$ ,

are sufficient to ensure that there is an unstable range  $[\underline{k^2}, \overline{k^2}]$  for the wave numbers. That is, for  $k^2 \in [\underline{k^2}, \overline{k^2}]$ , we have that  $\operatorname{Re}(\lambda(k^2)) > 0$ .

### Consequences

For large  $t$ , we have that (1.119) becomes

$$(1.127) \quad w(x, t) \sim \sum_{k \leq k \leq \bar{k}} c_k e^{\lambda(k^2)t} W_k(x).$$

It is important to note that the exponential growth will be restrained by the non-linearity of the reaction terms and the solution will tend towards a non-homogeneous spatial equilibrium configuration.

### 1.9.9 Detailed analysis of pattern formation

Let us consider the system

$$(1.128) \quad \begin{cases} u_t = \gamma f(u, v) + u_{xx} = \gamma(a - u + u^2 v) + u_{xx}, \\ v_t = \gamma g(u, v) + d v_{xx} = \gamma(b - u^2 v) + d v_{xx}, \end{cases} \quad x \in (0, l_x),$$

in the one-dimensional case. In the absence of diffusion, the stationary state is given by

$$(1.129) \quad (u_0, v_0) = \left( a+b, \frac{b}{(a+b)^2} \right),$$

and in correspondence of it, we have

$$(1.130) \quad \begin{aligned} f_u(u_0, v_0) &= \frac{b-a}{a+b}, \quad f_v(u_0, v_0) = (a+b)^2 > 0, \\ g_u &= \frac{-2b}{a+b} < 0, \quad g_v = -(a+b)^2 < 0. \end{aligned}$$

To have an unstable range, the following properties must hold

$$(1.131) \quad \begin{cases} f_u \cdot g_v < 0 \Rightarrow -b^2 + a^2 < 0 \Rightarrow a < b, \\ f_u + g_v < 0 \Rightarrow \frac{(b-a)}{(a+b)} - (a+b)^2 < 0 \Rightarrow b - a < (a+b)^3, \\ f_u g_v - f_v g_u > 0 \Rightarrow (a+b)^2 > 0, \\ d f_u + g_v > 0 \Rightarrow d(b-a) > (a+b)^3, \\ (d f_u + g_v)^2 - 4d(f_u g_v - f_v g_u) > 0 \Rightarrow [d(b-a) + (a+b)^3]^2 > 4d(a+b)^4. \end{cases}$$

These inequalities define a domain in the parameter space  $(a, b, d)$ , known as Turing space. Consider the eigenvalue problem

$$(1.132) \quad \begin{cases} \mathbf{W}_{xx} + k^2 \mathbf{W} = 0, \\ \frac{\partial \mathbf{W}}{\partial x} \Big|_{x=0} = \frac{\partial \mathbf{W}}{\partial x} \Big|_{x=l_x} = 0, \end{cases}$$

whose solutions are

$$(1.133) \quad \mathbf{W}_n(x) = \mathbf{C}_n \cos\left(\frac{n\pi x}{l_x}\right), n = 0, 1, 2, \dots,$$

where the  $\mathbf{C}_n$  are arbitrary constant vectors. The eigenvalues are given by

$$(1.134) \quad k_n = \frac{n\pi}{l_x}.$$

Once (1.131) are satisfied and there are wave numbers included in  $[\underline{k}, \bar{k}]$ , the related eigenfunctions  $\mathbf{W}_n$  are linearly unstable. The explicit expressions of  $\underline{k}$  and  $\bar{k}$ , in the case under consideration, are

$$(1.135) \quad \begin{aligned} \underline{k}^2 &= \frac{[d(b-a) - (a+b)^3] - \sqrt{[d(b-a) - (a+b)^3]^2 - 4d(a+b)^4}}{2d(a+b)} =: \gamma q_1(a, b, d), \\ \bar{k}^2 &= \frac{[d(b-a) - (a+b)^3] + \sqrt{[d(b-a) - (a+b)^3]^2 - 4d(a+b)^4}}{2d(a+b)} =: \gamma q_2(a, b, d), \end{aligned}$$

such that in the interval  $[\underline{k}, \bar{k}]$  we have instability.

If  $\gamma$  is too small, there will be no eigenvalue in  $[\underline{k}^2, \bar{k}^2]$  and thus the stationary state is stable. From (1.127) it follows that, as time passes, the non-homogeneous spatial solution is given by

$$(1.136) \quad w(x, t) = \sum_{n_1}^{n_2} c_n \exp \left[ \lambda \left( \frac{n^2 \pi^2}{l_x^2} \right) t \right] \cos \frac{n \pi x}{l_x},$$

where  $n_1$  is the smallest value of  $n$  such that the corresponding eigenvalue is in  $[\underline{k}^2, \bar{k}^2]$  and  $n_2$  is the largest.

Let us now consider the two-dimensional case. Let  $\Omega = (0, l_x) \times (0, l_y)$ , we have

$$(1.137) \quad \begin{cases} \left( \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2} \right) \mathbf{W} + k^2 \mathbf{W} = 0, \\ \frac{\partial \mathbf{W}}{\partial x}(0, y) = \frac{\partial \mathbf{W}}{\partial x}(l_x, y) = 0, \\ \frac{\partial \mathbf{W}}{\partial y}(x, 0) = \frac{\partial \mathbf{W}}{\partial y}(x, l_y) = 0. \end{cases}$$

The eigenfunctions and the eigenvalues are

$$(1.138) \quad W_{nm} = C_{nm} \cos \frac{n \pi x}{l_x} \cos \frac{m \pi y}{l_y},$$

$$(1.139) \quad k_{nm}^2 = \pi^2 \left( \frac{n^2}{l_x^2} + \frac{m^2}{l_y^2} \right), \quad \underline{k}^2 \leq k_{nm}^2 \leq \bar{k}^2.$$

The emerging solution will be

$$(1.140) \quad \sum_{n, m} C_{nm} \exp \{ \lambda(k_{nm}^2) t \} \cos \frac{n \pi x}{l_x} \cos \frac{m \pi y}{l_y},$$

with  $n, m$  such that  $\gamma q_1(a, b, d) = \underline{k}^2 \leq k_{nm}^2 \leq \bar{k}^2 = \gamma q_2(a, b, d)$ .

Let us return to the one-dimensional case and assume that  $\gamma$  is such that in the range  $[\underline{k}, \bar{k}]$  only the mode with  $n=1$  falls, then

$$(1.141) \quad w(x, t) \approx c_1 \exp \left\{ \lambda \left( \frac{\pi^2}{l_x^2} \right) t \right\} \cos \frac{\pi x}{l_x},$$

since all other modes decay exponentially with the passage of time. Let us take  $c_1 = (\epsilon)$ , for some small  $\epsilon > 0$ , and consider

$$(1.142) \quad u(x, t) \approx u_0 + \epsilon \exp \left\{ \lambda \left( \frac{\pi^2}{l_x^2} \right) t \right\} \cos \frac{\pi x}{l_x}.$$

Clearly, if exponential growth continued over time we would have  $|u| \rightarrow +\infty$ , for  $t \rightarrow +\infty$ . However, it is hypothesizable that the solution eventually settles down to a spatial pattern which looks like the cosine function shown in figure 1.8. Similar considerations can be made if the unstable mode is the second one.

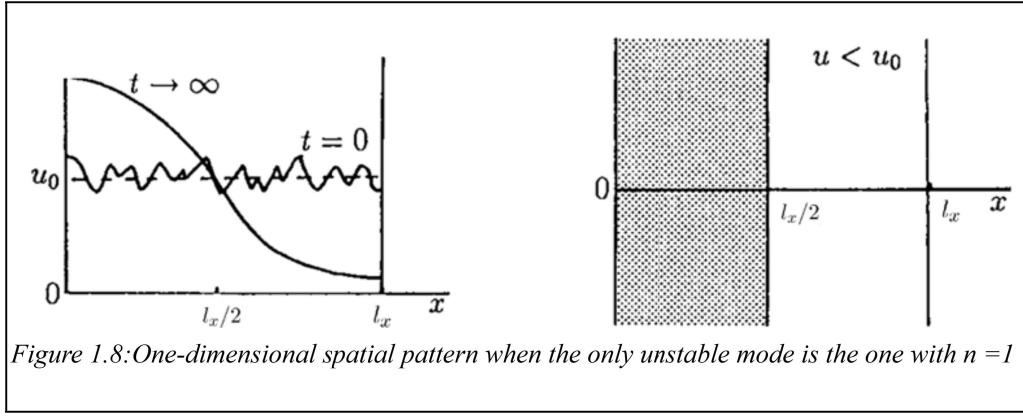


Figure 1.8: One-dimensional spatial pattern when the only unstable mode is the one with  $n = 1$

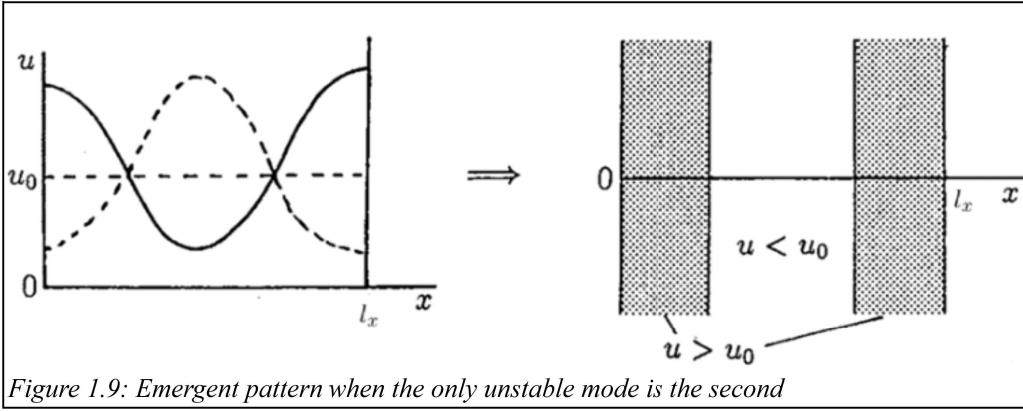


Figure 1.9: Emergent pattern when the only unstable mode is the second

The two-dimensional case, on the other hand, is much more complicated because the unstable modes depend on  $l_x$  and  $l_y$ . We have that

$$\gamma q_1(a, b, d) = \underline{k}^2 \leq k^2 = \pi^2 \left( \frac{n^2}{a^2} + \frac{m^2}{b^2} \right) \leq \bar{k}^2 = \gamma q_2(a, b, d).$$

If either dimension of the domain is too small, there is no unstable mode.

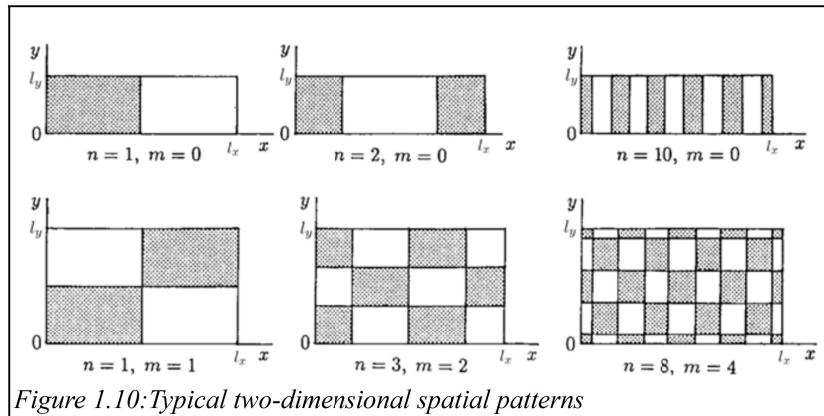


Figure 1.10: Typical two-dimensional spatial patterns

### 1.9.10 Pattern formation in growing domains

The evolution of patterns in growing domains can be quite complicated, but it has important biological implications. There are two cases:

1. as  $\gamma$  increases, the dispersion curve shifts along the wavelength axis, such that the unstable mode changes and becomes smaller (“discrete” pattern);
2. as  $\gamma$  increases, more modes are excited, that is the unstable band grows with gamma (“continuous” pattern)

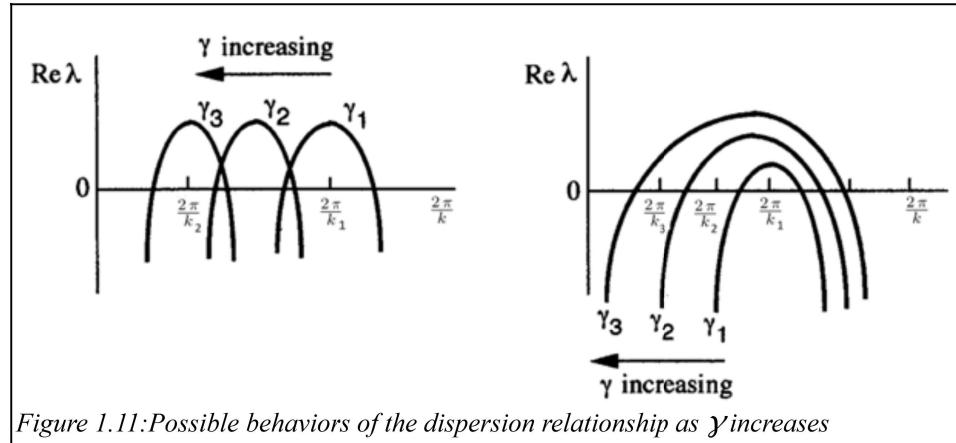


Figure 1.11: Possible behaviors of the dispersion relationship as  $\gamma$  increases

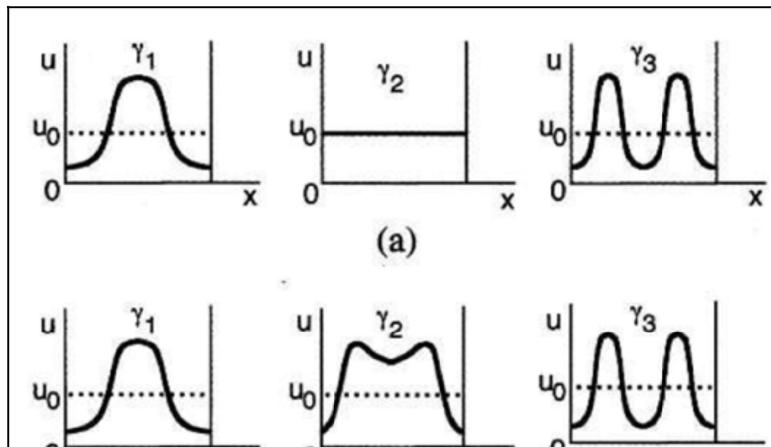


Figure 1.12: Evolution of spatial patterns for the two different types of dispersion relation behavior considered.

### 1.9.11 Hints of pattern formation on the coat of some mammals

Mammals exhibit a wide variety of patterns on their coats. According to J.D. Murray, a single mechanism may underlie the formation of the vast majority of commonly observed patterns: a diffusion reaction mechanism, characterized by diffusion instability. To create the colour patterns on animal skin, genetically

determined cells, melanoblasts, migrate to the surface of the embryo, where they transform into specialized cells, melanocytes, which are apt to produce melanin, a pigment responsible for skin, eye and hair colour. The basic assumption is that the movement and organization of melanoblasts reflects the spatial, prepattern distribution of as-yet unknown chemicals, the so-called morphogens. The formation of prepatterns occurs in the early stages of embryo development, while colour patterns form toward the end of embryo development. Because the mechanism of diffusion involved is not yet known, Murray in his studies used Thomas's, one both for illustrative purposes and because it describes truly observable kinetics with realistic parameters. As an example, one can consider the tail or paw of a feline. The spatial domain can then be thought of as rectangular, assuming that the skin covering the tail or paw is cut longitudinally and flattened. It is then natural to assume on the two edges at the lower and upper sides of the rectangle periodicity conditions for the concentrations of morphogens,  $u$  and  $v$ , and for their derivatives  $u_y$  and  $v_y$ , while considering homogeneous Neumann conditions on the other two sides:

$$(1.143) \quad \begin{cases} u(x,0,t)=u(x,l_y,t), \\ v(x,0,t)=v(x,l_y,t), \\ u_y(x,0,t)=u_y(x,l_y,t), \\ v_y(x,0,t)=v_y(x,l_y,t), \\ u_x(0,y,t)=u_x(l_x,y,t)=0, \\ v_x(0,y,t)=v_x(l_x,y,t)=0. \end{cases}$$

Proceeding with the calculations we can obtain

$$(1.144) \quad w(x,y,t) \sim \sum_{k \leq k_{nm} \leq \bar{k}} C_{nm} \exp\{\lambda_{nm} t\} \cos\left(\frac{n \pi x}{l_x}\right) \cos\left(\frac{m \pi y}{l_y}\right),$$

where the  $\lambda$ 's are solutions of

$$\lambda^2 + \lambda [k_{nm}^2(1+d) - \gamma \operatorname{Tr}(A)] + h(k_{nm}^2) = 0,$$

where

$$k_{nm}^2 = \pi^2 \left( \frac{n^2}{l_x^2} + \frac{4m^2}{l_y^2} \right).$$

We can draw the following conclusions

- unstable modes determine the colouring of the mantle;
- instability occurs for modes with  $k$  in the instability range;

- pattern formation is only possible for  $d > 1$  ;
- if  $l_y < l_x$ , the formation of transverse stripes is probable, while if  $l_x < l_y$ , non-constant spots in the y-direction are possible.

### 1.9.12 Non-existence of spatial patterns in reaction-diffusion systems

We want to show how the process of pattern formation can be destroyed if diffusion is fast enough.

In the one-dimensional case we consider

$$(1.145) \quad \begin{cases} u_t = f(u, v) + D_u u_{xx}, \\ v_t = g(u, v) + D_v v_{xx}, \end{cases}$$

with conditions

$$(1.146) \quad \begin{cases} u_x(0, t) = u_x(1, t) = v_x(0, t) = v_x(1, t) = 0, \\ u(x, 0) = u_0(x), \quad u_x(0, t) = u_x(1, t) = 0, \\ v(x, 0) = v_0(x), \quad v_x(0, t) = v_x(1, t) = 0. \end{cases}$$

The integral of energy is defined as

$$(1.147) \quad E(t) = \frac{1}{2} \int_0^1 (u_x^2 + v_x^2) dx,$$

which derived with respect to time gives us

$$\frac{dE}{dt} = \int_0^1 (u_x u_{xt} + v_x v_{xt}) dx,$$

where

$$\begin{cases} u_{xt} = u_{tx} = f_u u_x + f_v v_x + (D_u u_{xx})_x, \\ v_{xt} = v_{tx} = g_u u_x + g_v v_x + (D_v v_{xx})_x. \end{cases}$$

Substituting in the integral, we obtain

$$\begin{aligned} \frac{dE}{dt} &= \int_0^1 [u_x(D_u u_{xx})_x + u_x(f_u u_x + f_v v_x) + v_x(D_v v_{xx})_x + v_x(g_u u_x + g_v v_x)] dx \\ &= [u_x D_u u_{xx} + v_x D_v v_{xx}]_0^1 - \int_0^1 (D_u u_{xx}^2 + D_v v_{xx}^2) dx + \int_0^1 [f_x u_x^2 + g_x v_x^2 + (f_u + g_v) v_x u_x] dx. \end{aligned}$$

Let us define

$$d := \min(D_u, D_v), \quad M = \max_{u, v} (f_u^2 + f_v^2 + g_u^2 + g_v^2)^{\frac{1}{2}},$$

where the maximum is taken over all values  $u$  and  $v$  that the solutions can take, which are assumed to be bounded. Let us observe that

$$(u_x - v_x)^2 \geq 0,$$

which implies

$$u_x^2 + v_x^2 \geq 2 u_x v_x.$$

Then

$$\begin{aligned} \frac{dE}{dt} &\leq -d \int_0^1 (u_{xx}^2 + v_{xx}^2) dx + 2M \int_0^1 (u_x^2 + v_x^2) dx \\ &\leq -d \pi^2 \int_0^1 (u_x^2 + v_x^2) dx + 2M \int_0^1 (u_x^2 + v_x^2) dx = (4M - 2d\pi^2) E \leq 0, \end{aligned}$$

if  $d$  is sufficiently large, from which  $E \rightarrow 0$  for  $t \rightarrow +\infty$ . Consequently

$$(1.148) \quad u_x, v_x \rightarrow 0,$$

and this means that diffusion will cancel out any spatial inhomogeneity as time varies.

### 1.9.13 Numerical methods

We want to discretize the one-dimensional problem

$$(1.149) \quad \begin{cases} u_t = u_{xx} + f(u), \\ u_x(0, t) = u_x(1, t) = 0, \\ u(x, 0) = u^{(0)}(x). \end{cases}$$

We divide  $[0, 1]$  with  $N+1$  equispaced nodes  $x_i = i h$ ,  $h = \frac{1}{N}$ ,  $i = 0, \dots, N$ .

Approximations of the spatial derivatives can be calculated, such as:

- $u_{xx}(x_i, t) = \frac{u(x_{i+1}, t) - 2u(x_i, t) + u(x_{i-1}, t)}{h^2} + O(h^2)$ ,
- $u_x(x_0, t) = \frac{-3u(x_0, t) + 4u(x_1, t) - u(x_2, t)}{2h} + O(h^2)$ ,
- $u_x(x_N, t) = \frac{u(x_{N-2}, t) - 4u(x_{N-1}, t) + 3u(x_N, t)}{2h} + O(h^2)$ .

If the  $u_i(t)$  are appropriate approximations of the  $u(x_i, t)$ , using (4.61)-(4.62), and defining

- $\mathbf{U}(t) = (u_1(t), \dots, u_{N-1}(t))^T$ ,
- $\mathbf{F}(\mathbf{U}(t)) = (f(u_1(t)), \dots, f(u_{N-1}(t)))^T$ ,

we can discretize the problem

$$(1.150) \quad \dot{\mathbf{U}}(t) = \mathbf{A}_h \mathbf{U}(t) + \mathbf{F}(\mathbf{U}(t)),$$

$$(1.151) \quad \mathbf{U}(0) = (u_1(0), \dots, u_{N-1}(0))^T,$$

where  $\mathbf{A}_h$  is an appropriate coefficient matrix. It is called the semi-discretized problem, since it is discretized only with respect to  $x$ . The main difficulty lies in choosing an appropriate scheme to discretize the time variable. In fact, the use of the simplest scheme, the Euler forward scheme, is affected by severe limitations concerning its stability and accuracy. However, an explicit third-order Runge-Kutta method can be resorted to.

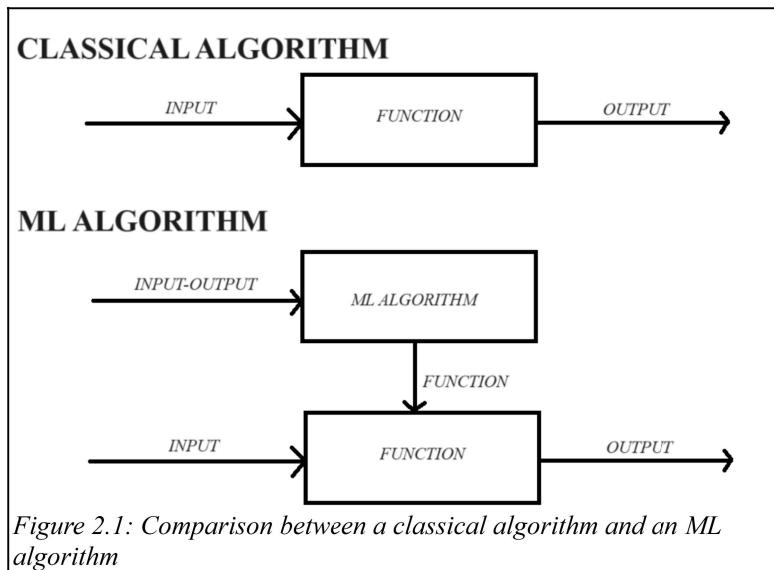
# Chapter 2

## Neural Networks for solving reaction-diffusion equations

### 2.1 Machine Learning

#### 2.1.1 ML algorithms

Sometimes it can be particularly difficult to write algorithms to solve certain types of problems, such as recognizing an image or driving a car. That's where Machine Learning (ML) algorithms [6] come in: they take advantage of a collection of examples (*training sets*), consisting of input and output, to produce a model (*target function*) for classifying new instances.



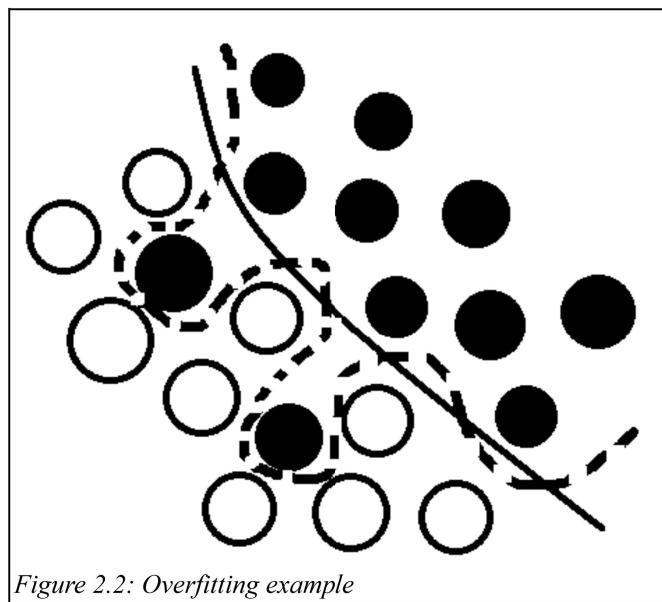
A training set  $S$  is a set of examples (*instances*), each of which is characterized by a set of features (*attributes*) and a class label. The association between instances and the corresponding class is the target function.

### 2.1.2 Overfitting

One of the biggest problems of the training sets is the presence of noise, i.e., data with the wrong class label are present and that could adversely affect the creation of the model. Another problem is the lack of representative examples, i.e., the training set has too many parameters compared to the number of instances. As a result, one can fall into the phenomenon of overfitting, that is, one finds a model that works very well on the training set used, but has very poor generalization capabilities.

#### Example

Suppose we want to classify a set of objects into two distinct classes (represented in Figure 2.2 with a black ball and a white ball, respectively).



It is possible to observe how the continuous line is constructed having a good accuracy on the training set, but without achieving 100% accuracy, in fact some black balls are considered as white balls. However, this allows us to have a better generalization ability since the constructed model is not too restrictive. Instead, the dashed line is constructed to present 100% accuracy on the training set (no black balls are classified as white balls, and vice versa). This results in a lower generalization ability, as the model turns out to be too restrictive and as a result some new instances may be misclassified. To prevent overfitting one should try to prefer simple models rather than complex ones.

## 2.2 Neural Networks

### 2.2.1 Definition of Neural Networks

Neural Networks are Machine Learning tools inspired by the way human brain works. They are composed of processing units called neurons (*perceptrons*) [1]. These neurons are organized into layers and connections, and the Neural Network is able to learn from the data through a training process. They have a wide range of applications, such as: Speech recognition, Autonomous driving, Image classification, etc..

### 2.2.2 Perceptrons

A *perceptron* is a binary classifier of the form

$$(2.1) \quad y = \sigma(Z) = \begin{cases} 0, & \text{if } Z \leq 0, \\ 1, & \text{if } Z > 0, \end{cases}$$

where:

- $Z = \sum_{i=1, \dots, n} w_i x_i + w_0$ ;
- $x_1, \dots, x_n$  are the inputs, that is, real or boolean values;
- $w_1, \dots, w_n$  are the weights;
- $w_0$  is the bias (distortion);
- $\sigma$  is the activation step function.

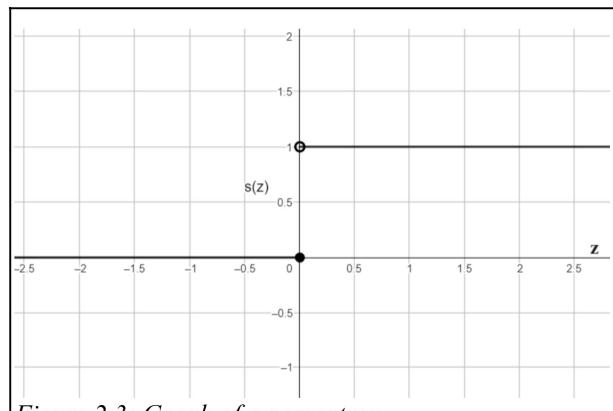
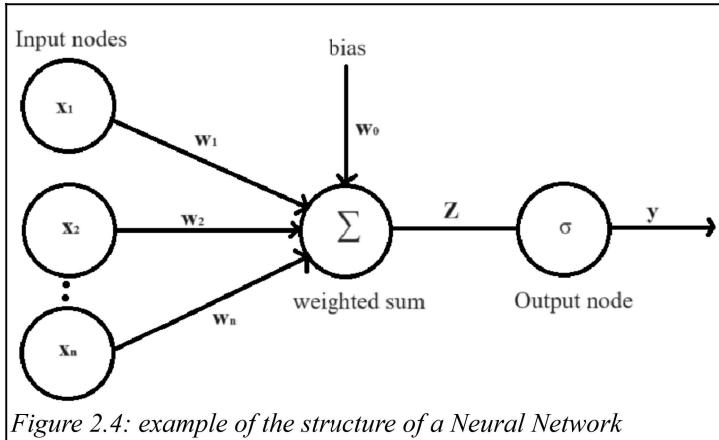


Figure 2.3: Graph of a perceptron

The "bias" is an additional parameter that represents the level of activation of a neuron when all its inputs are null. In simpler terms, bias allows the perceptron to learn even when all inputs are null. In a perceptron, each input is associated with a weight that indicates the importance of the input to the neuron's output. Bias, similarly, is an additional weight that is not associated with any specific input but still contributes to the determination of the neuron's output. Without bias, perceptron capacity would be limited to linear functions passing through the

origin. Too high or too low a bias could result in the perceptron's inability to adapt properly to the training data.



By placing  $t = -w_0$ , we can observe that

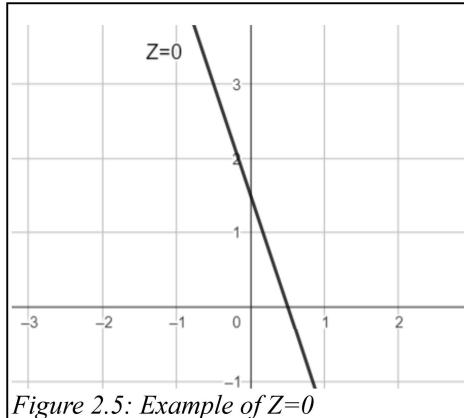
$$(2.2) \quad y = \sigma(Z) = \begin{cases} 0, & \text{if } \sum_{i=1}^n w_i x_i < t, \\ 1, & \text{if } \sum_{i=1}^n w_i x_i > t, \end{cases}$$

that is,  $t$  represents the *activation threshold*.

### Example

Let's suppose to have two input nodes  $x_1, x_2$ , then

$$Z = w_0 + w_1 x_1 + w_2 x_2.$$



We can observe that  $Z=0$  represents the *hyperplane* separating the input space into two regions:

- $P$ , if  $Z > 0$ ;
- $N$ , if  $Z < 0$ .

So the perceptron will return 1 if  $(x_1, x_2) \in P$ , otherwise it will return 0 if  $(x_1, x_2) \in N$ .

### Remark

Since our brain is composed of billions of neurons that are responsible for our logical thinking, it was natural to try to design a neuron that can represent Boolean functions such as "AND" and "OR." In fact, the perceptron is capable of computing any binary function that is *linearly separable*, that is, we can construct a perceptron from any Data Set that is linearly separable.

Unfortunately, not all Data Sets are linearly separable, consequently it is not possible to use a single perceptron to implement a generic function, but we must combine several perceptrons to form the so-called

#### 2.2.3 Multilayer perceptron (MLP)

A multilayer perceptron is an artificial Neural Network composed of multiple layers of neurons, usually organized into an input layer, one or more hidden layers, and an output layer. Each layer consists of a set of artificial neurons, each of which is connected to all the neurons in the previous and next layer. In general, for each function of the type

- $f : \{0,1\}^N \rightarrow \{0,1\}$ , boolean function,
- $f : \mathbb{R}^N \rightarrow \{0,1\}$ , real function,

there exists an MLP that implements it. As shown in Figure 2.6, by combining multiple perceptrons, it is possible to construct an MLP that returns 1 when the input is within the closed area and 0 otherwise. The greater the number of nodes, the greater the accuracy of the approximation.

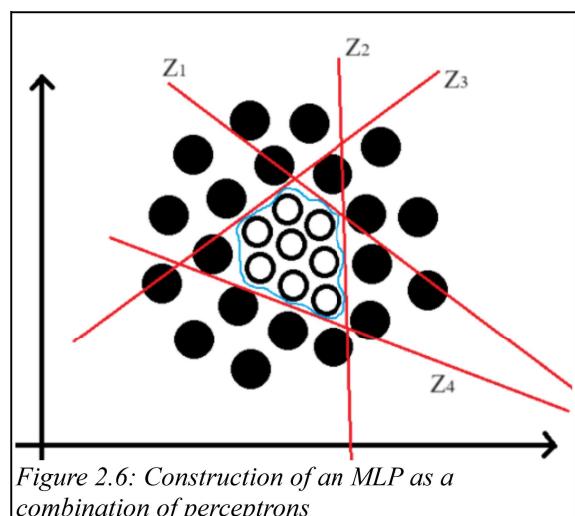


Figure 2.6: Construction of an MLP as a combination of perceptrons

### Remark

Unfortunately, the step function is a discrete function and is discontinuous at the point  $Z=0$ , so it is not derivable. This is why it is useful to introduce new activation functions that are continuous and derivable over all  $\mathbb{R}$ .

#### 2.2.4 Sigmoid function

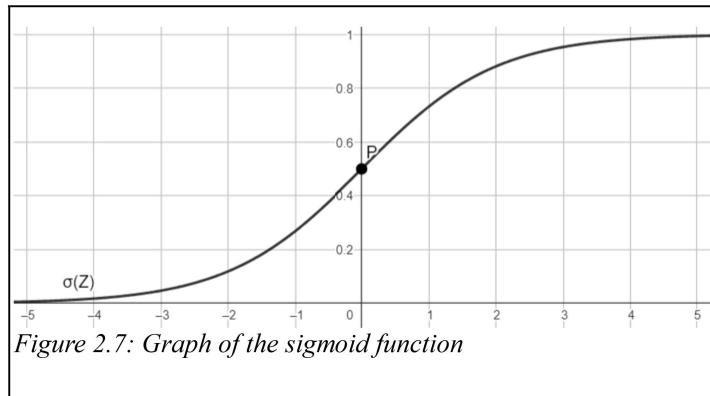
The *Sigmoid function*, also known as the logistic function, is one of the most widely used activation functions in Neural Networks, especially in the early stages of development in the field of deep learning. It's a function that produces a sigmoid curve, that is, a curve having an "S" shape. A sigmoid function is of the form

$$(2.3) \quad y = \sigma(z) = \frac{1}{1+e^{-z}},$$

where

$$z = w_0 + \sum_{i=1}^n w_i x_i.$$

A sigmoid function takes a real value and scales it between 0 and 1.



This makes the function particularly useful when you want to interpret the output as a probability. Very large negative numbers approach 0, while large positive numbers approach 1. The function returns 0.5 when the input is 0. It introduces nonlinearities into the model, allowing the network to learn complex, nonlinear representations of the data. We can observe that it is a continuous and derivable function over all  $\mathbb{R}$  and its derivative is

$$(2.4) \quad \frac{d\sigma(z)}{dz} = \sigma(z)(1-\sigma(z)).$$

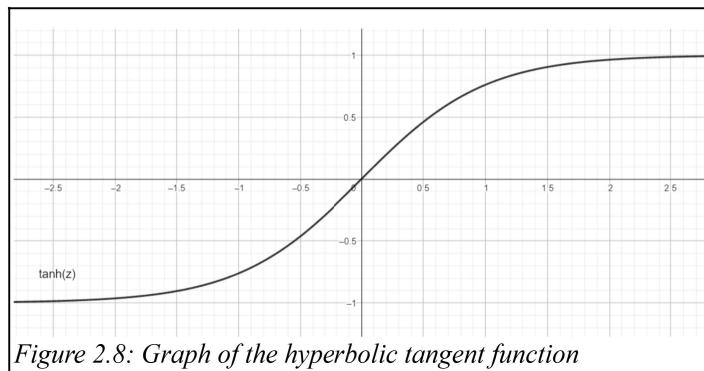
For very large or very small inputs, the derivative of the Sigmoid tends to zero, causing the vanishing gradient problem. This slows down the updating of weights in deep networks during training.

### 2.2.5 Hyperbolic Tangent function

The Tanh (*Hyperbolic Tangent*) activation function [5] is one of the most commonly used activation functions in Neural Networks. The Tanh function is defined as:

$$(2.5) \quad \text{Tanh}(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}.$$

Tanh maps any input value to a range between -1 and 1. This is advantageous compared to the Sigmoid function, whose output ranges between 0 and 1.



Unlike the Sigmoid, Tanh is zero-centered, meaning that negative and positive output values are balanced. This can improve convergence during training because the gradients are not all shifted in the same direction. As with the sigmoid function, the hyperbolic tangent also introduces nonlinearity into the model, allowing the network to learn complex, nonlinear representations of the data. Unfortunately, however, it can suffer from vanishing gradient problems (gradient becoming very small), slowing down the training of deep networks. Calculating the Tanh function requires the use of exponential functions, which are computationally more expensive compared to the basic operations used in other activation functions.

### 2.2.6 ReLU function

The ReLU (Rectified Linear Unit) activation function [5] is one of the most popular and widely used in modern Neural Networks, especially deep Neural

Networks. The ReLU function is defined as:

$$(2.6) \quad \text{ReLU}(z) = \max(0, z).$$

It's very simple to implement and computationally efficient, since it requires only one comparison and does not involve exponential or trigonometric operations. As for the previous functions, the ReLU also introduces nonlinearity into the model, allowing the network to learn nonlinear representations of the data. The ReLU sets all negative input values to zero, which introduces sparsity into the model, that is, many nodes in the Neural Network will have value zero. Sparsity allows us to perform fewer calculations because neurons with an output of zero do not contribute further to calculations in later levels. This can make the network more efficient. Unlike activation functions such as Sigmoid and Tanh, ReLU does not suffer severely from the vanishing gradient problem, allowing gradients to remain meaningful even in deep networks. This promotes faster convergence during training. The ReLU is computationally very efficient because it requires only simple operations (compare and maximum). Because the ReLU sets negative values to zero, only a fraction of the neurons are activated at any instant, reducing the risk of overfitting and improving efficiency.

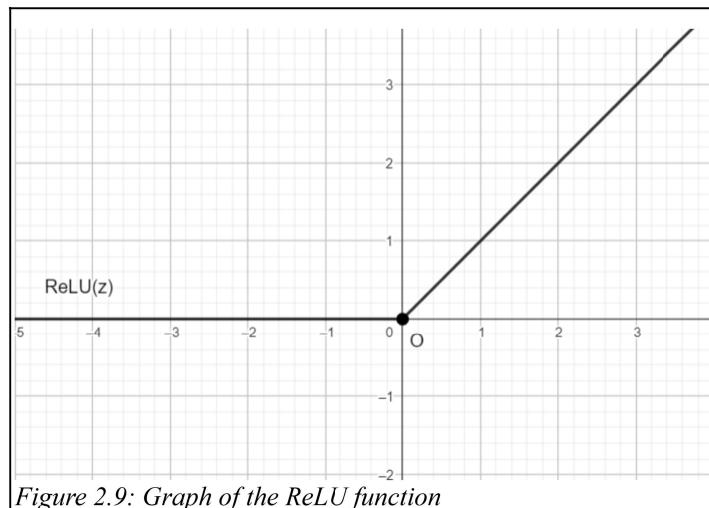
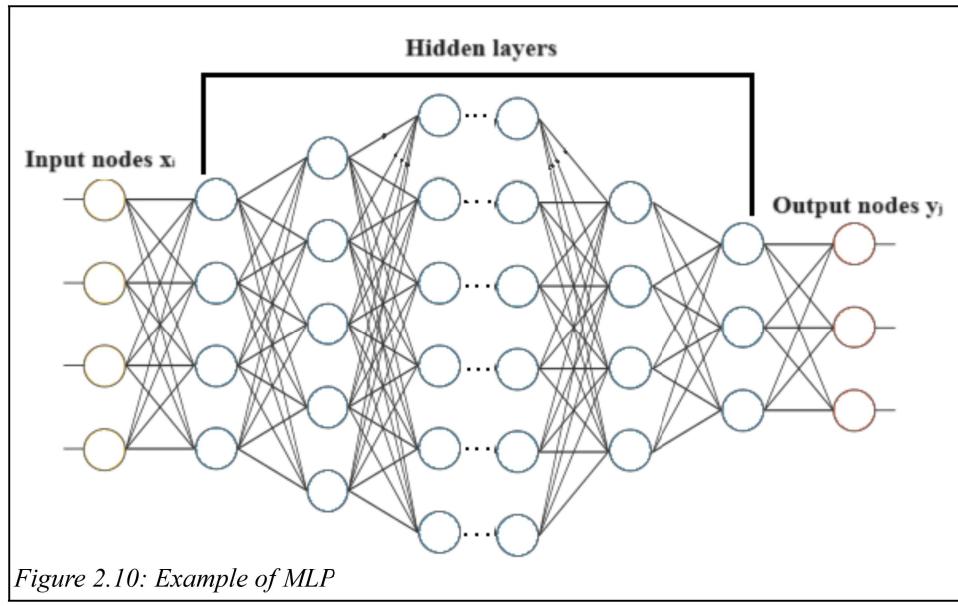


Figure 2.9: Graph of the ReLU function

A common problem with ReLU is that neurons can "die" during training. If a neuron receives negative input, its output is zero, and if this continues to happen, the neuron stops updating (its gradients become zero). This is known as the "dead neuron" problem. In some cases, very large gradients can cause training to diverge. This can be mitigated by using techniques such as gradient normalization or ReLU variants.

### 2.2.7 How MLP works

We have already mentioned that MLP consists of an input layer, one or more hidden layers, and an output layer. Each layer consists of a set of artificial neurons, each of which is connected to all the neurons in the previous layer and the next layer. An MLP uses a process called "feed-forward," which means that input is presented to the network and information flows forward through the hidden layers until it reaches the output layer. During this process, each node in each layer receives input from the nodes in the previous layer, applies the activation function, and sends the output to the nodes in the next layer.



Suppose we have a Neural Network defined by:

- $m$  input nodes  $x_1, \dots, x_m$  (sometimes, we will denote by  $\mathbf{x}^{(0)}$  or  $\mathbf{x}$  the vector of input nodes);
- $N$  hidden layers  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$ , consisting respectively of  $d_1, \dots, d_N$  nodes (in particular, we will denote by  $x_i^{(j)}$  the  $i$ -th component of the vector  $\mathbf{x}^{(j)}$ );
- $n$  output nodes  $y_1, \dots, y_n$  (sometimes, we will denote by  $\mathbf{x}^{(N+1)}$  or  $\mathbf{y}$  the vector of output nodes);
- $\omega^{(j)}$  are the weights of the arches connecting  $\mathbf{x}^{(j-1)}$  with  $\mathbf{x}^{(j)}$  (we will denote by  $\omega_{ki}^{(j)}$  the weight connecting  $x_i^{(j-1)}$  with  $x_k^{(j)}$ );
- $\theta^{(j)}$  are the biases of  $\mathbf{x}^{(j)}$  (we will denote by  $\theta_k^{(j)}$  the bias of  $x_k^{(j)}$ ).

For the  $j$ -th layer we calculate

$$x_1^{(j)} = \sigma \left( \sum_{i=0}^{d_{j-1}} \omega_{1i}^{(j)} x_i^{(j-1)} \right), \dots, x_{d_j}^{(j)} = \sigma \left( \sum_{i=0}^{d_{j-1}} \omega_{di}^{(j)} x_i^{(j-1)} \right), \quad j=1, \dots, N,$$

with  $x_0^{(j)} = 1$ ,  $j=0, \dots, N$ . Instead, for the output layer, we have

$$y_1 = \sigma \left( \sum_{i=0}^{d_N} \omega_{1i}^{(N+1)} x_i^{(N)} \right), \dots, y_n = \sigma \left( \sum_{i=0}^{d_N} \omega_{ni}^{(N+1)} x_i^{(N)} \right).$$

In the next subsection, we will rewrite everything in a compact form.

### 2.2.8 General description of a Neural Network

A Neural Network, with  $N$  hidden layers, can be seen as a map  $N_\sigma : \mathbb{R}^m \rightarrow \mathbb{R}^n$ , which depends on the parameters [3]

$$W = \{(\boldsymbol{\omega}^{(I)}, \boldsymbol{\theta}^{(I)}), \dots, (\boldsymbol{\omega}^{(N+I)}, \boldsymbol{\theta}^{(N+I)})\},$$

such that

$$\mathbf{y} = \boldsymbol{\omega}^{(N+I)} \sigma(\boldsymbol{\omega}^{(N)} \sigma(\dots \boldsymbol{\omega}^{(2)} \sigma(\boldsymbol{\omega}^{(1)} \mathbf{x} + \boldsymbol{\theta}^{(1)}) + \boldsymbol{\theta}^{(2)} \dots) + \boldsymbol{\theta}^{(N)}) + \boldsymbol{\theta}^{(N+I)},$$

where

- $\boldsymbol{\omega}^{(j)} \in \mathbb{R}^{d_j \times d_{j-1}}$ ,  $j=1, \dots, N+I$ , is the matrix of weights connecting the nodes  $\mathbf{x}^{(j-1)}$  with the nodes  $\mathbf{x}^{(j)}$ ;
- $\boldsymbol{\theta}^{(j)} \in \mathbb{R}^{d_j}$ ,  $j=1, \dots, N+I$ , is the biases vector of  $\mathbf{x}^{(j)}$ .

### 2.2.9 Definition (Neural Networks as composition of functions)

The Neural Network  $N_\sigma : \mathbb{R}^m \rightarrow \mathbb{R}^n$  can then be seen as a composition [3] of functions

$$(2.7) \quad N_\sigma = \mathbf{z}^{(k+I)} \circ \sigma \circ \mathbf{z}^{(k)} \circ \dots \circ \sigma \circ \mathbf{z}^{(2)} \circ \sigma \circ \mathbf{z}^{(1)},$$

where the  $\mathbf{z}^{(j)} : \mathbb{R}^{d_{j-1}} \rightarrow \mathbb{R}^{d_j}$ , such that  $\mathbf{z}^{(j)} \mathbf{x} = \boldsymbol{\omega}^{(j)} \mathbf{x} + \boldsymbol{\theta}^{(j)}$ , are affine maps.

## 2.3 Neural Networks learning

### 2.3.1 Example (Error function)

Let  $\mathbf{x} = \{x_1, \dots, x_m\}$  be the input,  $\mathbf{W} = \{(\boldsymbol{\omega}^{(1)}, \boldsymbol{\theta}^{(1)}), \dots, (\boldsymbol{\omega}^{(N+1)}, \boldsymbol{\theta}^{(N+1)})\}$ , be the set of weights and biases, and  $\mathbf{y} = \{y_1, \dots, y_n\}$  be the output. We can define the *quadratic error* as

$$\varepsilon(\mathbf{W}) = \frac{1}{2} \sum_{i=1}^n (y_i - t_i)^2,$$

where  $t_i$  is the target value of the output nodes  $y_i$ ,  $i = 1, \dots, n$ .

It is natural to think of a way to train the Neural Network by varying the weights and the biases in such a way that the quadratic error function is minimized.

To do this, we will use the

### 2.3.2 Steepest descent method (SDM)

The *steepest descent method* [10] is a widely used algorithm for finding the minimum of a scalar function  $f(\mathbf{x})$ , with  $\mathbf{x} = (x_1, \dots, x_m) \in \mathbb{R}^m$ . It is based on calculating the gradient of the objective function with respect to its arguments, and then deciding the direction in which to move to reach the minimum. You start with an initial point and move in the direction opposite to the gradient of the function at the current point. This means that you move along the direction in which the function is decreasing most rapidly. In mathematical terms, given a current point  $\mathbf{x}^k$ , it is possible to calculate the next point  $\mathbf{x}^{k+1}$  in the following way:

$$(2.8) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k),$$

where

- $\alpha_k$  is the *learning rate*, which can be fixed or variable;
- $\nabla f(\mathbf{x}^k)$  is the *gradient* of the function  $f$  evaluated at the point  $\mathbf{x}^k$ .

The process continues until a stop condition is reached, such as when  $|\mathbf{x}^{k+1} - \mathbf{x}^k| < \epsilon$ , with  $\epsilon$  suitably chosen, or when the value of the function reaches a

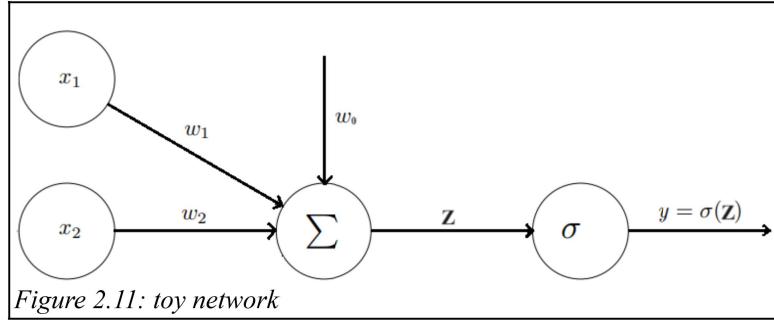
desired threshold. This method is easy to implement and is very effective for smooth functions. However, it may converge slowly in some special cases that will not be studied. Another very important task is the choice of parameters  $\alpha_k$ , indeed:

- if the steps are too small, convergence will be very slow;
- if the steps are too large, there is a risk of “jumping” beyond the minimum.

### 2.3.3 Learning scheme, the toy network

For simplicity, let us consider a perceptron consisting of 2 input nodes and a single output node.

Consider  $N$  examples of input  $x_1^k, x_2^k$ ,  $k=1, \dots, N$ , with target value



$t^k$ ,  $k=1, \dots, N$ . Since the perceptron has not yet learned, when it is applied on example  $k$ , it will produce an error  $e^k = t^k - y^k$ ,  $k=1, \dots, N$ . This error will be used to improve the values of the weights. We can observe that, in  $k$ -th step, the goal is to minimize the squared error (*squared Loss Function*)

$$E^k(\omega_0, \omega_1, \omega_2) = \frac{1}{2} \left( t^k - \sigma \left( \sum_{j=0}^2 \omega_j x_j^k \right) \right)^2.$$

The idea is to vary the weights by using the local steepest descent method in the opposite direction of the gradient of  $E$ , that is

$$\omega_j^{k+1} = \omega_j^k - \eta \frac{\partial E^k}{\partial \omega_j}, \quad j=0, 1, 2,$$

where  $\eta$ , as we have said, is the *learning rate*.

Since the derivative of the sigmoid function is  $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ , we can calculate the gradient components of  $E^k$  as

$$\frac{\partial}{\partial \omega_i} E^k(\omega_0, \omega_1, \omega_2) = -e^k x_i^k \sigma' \left( \sum_{j=0}^2 \omega_j x_j^k \right).$$

We obtain that

$$\omega_j^{k+1} = \omega_j^k + \eta \delta^k x_j^k, \quad j=0,1,2,$$

where  $\delta^k = e^k \sigma'(z^k)$ , with  $z^k = \sum_{j=0}^2 \omega_j x_j^k$ .

Due to the presence of  $\delta^k$ , this procedure is also called the *delta rule* for the learning step with the sigmoid function. A gradient sweep over all N learning examples is called an epoch. Since  $\eta$  is usually taken small (unless a reliable estimate of its value is available), a large number of epochs may be required to approximately minimise the Loss Function for all k.

### 2.3.4 Learning scheme, the general case

Consider now, as it is explained in [1], that we have a MLP consisting of  $m$  input nodes,  $n$  output nodes, and a single hidden layer consisting of  $N$  nodes. For convenience, we consider the bias as if it were an input node with value  $x_0 = 1$ . We denote by:

- $\boldsymbol{\omega}^{(I)} = (\omega_{ij}^{(I)})_{i=1,\dots,N; j=0,1,\dots,m}$  the matrix of weights and biases of the input nodes;
- $\mathbf{x} = (x_0, x_1, x_2, \dots, x_m)^T$  the vector of input nodes;
- $\mathbf{z}^{(I)} = (z_1^{(I)}, \dots, z_N^{(I)})^T$  the vector of the weighted sums of the hidden nodes.

We observe that

$$\mathbf{z}^{(I)} = \boldsymbol{\omega}^{(I)} \mathbf{x},$$

and applying the activation function  $\sigma$  we get

$$\mathbf{x}^{(I)} = (\sigma(z_1^{(I)}), \dots, \sigma(z_N^{(I)}))^T.$$

We will also write  $\mathbf{x}^{(I)} = \sigma(\mathbf{z}^{(I)}) = \sigma(\boldsymbol{\omega}^{(I)} \mathbf{x})$ , where  $\sigma$  is applied componentwise.

Now we denote by:

- $\boldsymbol{\omega}^{(2)} = (\omega_{ij}^{(2)})_{i=1,\dots,n; j=0,1,\dots,N}$  the matrix of weights and biases of the hidden layer;
- $\mathbf{z}^{(2)} = (z_1^{(2)}, \dots, z_n^{(2)})^T$  the vector of the weighted sums of the output nodes.

We observe that

$$\mathbf{z}^{(2)} = \boldsymbol{\omega}^{(2)} \mathbf{x}^{(I)} = \boldsymbol{\omega}^{(2)} \sigma(\boldsymbol{\omega}^{(I)} \mathbf{x}),$$

and applying the activation function  $\sigma$  we get

$$\mathbf{y} = \sigma(\mathbf{z}^{(2)}) = \sigma(\boldsymbol{\omega}^{(2)} \mathbf{x}^{(I)}) = \sigma(\boldsymbol{\omega}^{(2)} \sigma(\boldsymbol{\omega}^{(I)} \mathbf{x})).$$

Consider  $K$  examples of input  $(x_1, \dots, x_m)^k$ ,  $k=1, \dots, K$ , with target values  $(t_1, \dots, t_n)^k$ ,  $k=1, \dots, K$ . Now we can define a learning process for a MLP. Obviously, using what we introduced above, we can update the  $\omega^{(2)}$  weights in the following way

$$\omega_{ij}^{(2)} = \omega_{ij}^{(2)} + \eta \delta_i^{(2)} y_j^{(1)}, \quad i=1, \dots, n, \quad j=1, \dots, N,$$

where, for simplicity of notation, we have omitted to write the index  $k$ . Now we also want to improve the  $\omega^{(1)}$  weights, but unfortunately we cannot use the same rules as that for the output nodes, since no target value is available to estimate the error term of a hidden node. The idea is to propagate the error backward in the following way

$$e_j^{(1)} = \sum_{i=1}^n \omega_{ij}^{(2)} \delta_i^{(2)}, \quad j=1, \dots, N,$$

that is, in vector form we have

$$\mathbf{e}^{(1)} = \boldsymbol{\omega}^{(2)T} \boldsymbol{\delta}^{(2)}.$$

By using this new error expression, we can define

$$\delta_i^{(1)} = e_i^{(1)} \sigma'(z_i^{(1)}), \quad i=1, \dots, N.$$

Similarly to the previous step, we can formulate a delta rule for the weights  $\omega^{(1)}$ , that is

$$\omega_{ij}^{(1)} = \omega_{ij}^{(1)} + \eta \delta_i^{(1)} x_j, \quad i=1, \dots, N, \quad j=0, 1, \dots, m.$$

This idea of propagating the error and the resulting delta rule for updating the weights  $\omega^{(1)}$  is actually the result of finding the minimum of the Loss Function  $\frac{\|\mathbf{t} - \mathbf{y}\|_2^2}{2}$  with respect to the elements of  $\omega^{(1)}$ . Indeed, by calculation we derive

$$\frac{1}{2} \frac{\partial}{\partial \omega_{ij}^{(1)}} \|\mathbf{t} - \mathbf{y}\|_2^2 = - \sum_{l=1}^n \delta_l^{(2)} \omega_{li}^{(2)} \sigma'(z_i^{(1)}) x_j.$$

### 2.3.5 Backpropagation (BP)

We can conclude summarising the BP method [1] with the SDM scheme. In this algorithm, we denote with  $N_{\text{epochs}}$  the number of learning sweeps over the training data. The matrices of weights  $\omega^{(1)}$  and  $\omega^{(2)}$  are initialised such that they have full rank. For example, for  $\omega^{(1)}$  (and similarly for  $\omega^{(2)}$ ) we set  $\omega_{jj}^{(1)} = 1$ ,  $j=0, 1, \dots, m$

and  $w_{ij}^{(l)} = 0$ , if  $i \neq j$ , and this setting guarantees the injectivity of  $\omega^{(l)}$  and the surjectivity of its transpose.

### 2.3.6 Algorithm (Backpropagation and SDM method)

Input: initialise  $\omega^{(1)}$  and  $\omega^{(2)}$ , choose  $N_{\text{epochs}}$ ,  $\eta > 0$ .

Learning data: input nodes  $\mathbf{x}^k = (x_1^k, \dots, x_m^k)$ , with target values  $\mathbf{t}^k$ ,  $k = 1, \dots, K$ .

for  $l = 1$  to  $N_{\text{epochs}}$

for  $k = 1$  to  $K$

$$\mathbf{z}^{(l)} = \omega^{(l)} \mathbf{x}^k$$

$$\mathbf{x}^{(l)} = \sigma(\mathbf{z}^{(l)})$$

$$\mathbf{z}^{(2)} = \omega^{(2)} \mathbf{x}^{(1)}$$

$$\mathbf{y} = \sigma(\mathbf{z}^{(2)})$$

$$\mathbf{e} = \mathbf{t}^k - \mathbf{y}$$

$$\boldsymbol{\delta}^{(2)} = \mathbf{e} \odot \sigma'(\mathbf{z}^{(2)})$$

$$\mathbf{e}^{(1)} = \omega^{(2)^T} \boldsymbol{\delta}^{(2)}$$

$$\boldsymbol{\delta}^{(1)} = \mathbf{e}^{(1)} \odot \sigma'(\mathbf{z}^{(1)})$$

$$w_{ij}^{(l)} = w_{ij}^{(l)} + \eta \delta_i^{(l)} x_j^k, \quad i = 1, \dots, N, \quad j = 0, 1, \dots, m$$

$$w_{ij}^{(2)} = w_{ij}^{(2)} + \eta \delta_i^{(2)} x_j^{(1)}, \quad i = 1, \dots, n, \quad j = 1, \dots, N$$

end

mix the data input

end

#### Remark

The symbol  $\mathbf{a} \odot \mathbf{b}$  represents componentwise multiplication.

### 2.3.7 Definition of Loss Function

In general, a *Loss Function*  $L(u)$  [3] is a functional that quantifies how much the output of a Neural Network deviates from the actual or expected value. In other words, it measures the error between the predictions made by the model and the desired outcomes. The goal of training the model is to minimize this functional so that the predictions are as accurate as possible.

### Example (L2-Loss Function)

A commonly used Loss Function is the L<sub>2</sub>-Loss Function

$$L_2(y, f(\mathbf{x}, \boldsymbol{\omega})) = |y - f(\mathbf{x}, \boldsymbol{\omega})|^2,$$

where  $y \in \mathbb{R}$  is the exact value of the solution corresponding to the input  $\mathbf{x}$ , and  $f(\mathbf{x}, \boldsymbol{\omega}) \in \mathbb{R}$  è the approximate value calculated through the Neural Network.

Very often, however, one may run into problems that require minimizing more complex operators. To do that, we will use the

#### 2.3.8 Gradient Descent method

The *Gradient Descent method* [10] is a special version of the Steepest Descent Method in which the step is fixed, viz.

$$(2.9) \quad \mathbf{x}^{k+1} = \mathbf{x}^k - \eta \nabla f(\mathbf{x}^k).$$

In our case, let  $\boldsymbol{\omega}, \boldsymbol{\theta}$  be the parameters on which the Neural Network  $N(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta})$  depends, then we will have

$$(2.10) \quad \begin{aligned} \boldsymbol{\omega}_{new} &= \boldsymbol{\omega}_{old} - \eta \nabla_{\omega} N(\mathbf{x}, t; \boldsymbol{\omega}_{old}, \boldsymbol{\theta}), \\ \boldsymbol{\theta}_{new} &= \boldsymbol{\theta}_{old} - \eta \nabla_{\theta} N(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta}_{old}), \end{aligned}$$

where  $\eta$  is the *learning rate*.

The *Stochastic Gradient Descent method* is a variant of the Gradient method. Unlike classical Gradient Descent, which uses the entire dataset to compute gradients, SGD uses only a single sample (or a small subset of samples, called Mini-Batch Gradient Descent) to update model parameters. This leads to several advantages and disadvantages over whole-batch based methods.

Advantages:

- parameter updates are much more frequent;
- reduces computation time per update, making the algorithm faster on very large training sets;
- converges faster (if it converges at all);
- the algorithm can exit local minimum points;
- requires less memory, since the entire training set does not have to be loaded.

Disadvantages:

- convergence is not guaranteed; the algorithm may oscillate around the minimum;
- it is often necessary to reduce the learning rate to get to convergence.

To calculate the gradient of the Loss Function, we use the

### 2.3.9 Automatic Differentiation

The idea behind the backpropagation is that an algorithm can be viewed as a composition of simpler operations, and then its derivative can be derived through the chain rule [3] by deriving these simpler operations.

Let  $f$  be an algorithm that can be written as a composition of simpler functions

$$(2.11) \quad f = f_n \circ f_{n-1} \circ \dots \circ f_2 \circ f_1.$$

Then we can define the *intermediate variables*

$$\begin{aligned} w_0 &= x, \\ w_1 &= f_1(w_0), \\ w_2 &= f_2(w_1), \\ &\vdots \\ w_n &= f_n(w_{n-1}), \end{aligned}$$

from which

$$\begin{aligned} w_0' &= \frac{\partial x}{\partial x} = 1, \\ w_1' &= \frac{\partial w_1}{\partial w_0} w_0', \\ w_2' &= \frac{\partial w_2}{\partial w_1} w_1', \\ &\vdots \\ w_n' &= \frac{\partial w_n}{\partial w_{n-1}} w_{n-1}'. \end{aligned}$$

Consequently, the derivative of  $f$  respect to  $x$  can be calculated as

$$(2.12) \quad \frac{\partial f}{\partial x} = \frac{\partial w_n}{\partial w_{n-1}} \frac{\partial w_{n-1}}{\partial w_{n-2}} \dots \frac{\partial w_2}{\partial w_1} \frac{\partial w_1}{\partial w_0} = \prod_{i=1}^n \frac{\partial w_i}{\partial w_{i-1}}.$$

## 2.4 The Universal Function Approximation Theorem for Neural Networks

### 2.4.1 Introduction

Let us consider the set

$$(2.13) \quad \mathbf{M}(\sigma) = \text{span}\{\sigma(\omega \cdot \mathbf{x} - \theta) | \theta \in \mathbb{R}, \omega \in \mathbb{R}^n\}.$$

In [9] they asked the following question:

for which activation function  $\sigma$  is true that, for any  $f \in C(\mathbb{R}^n)$ , any compact subset  $K$  of  $\mathbb{R}^n$ , and any  $\epsilon > 0$ , there exists a function  $g \in \mathbf{M}(\sigma)$  such that

$$(2.14) \quad \max_{x \in K} \|N_\sigma(x) - f(x)\| < \epsilon ?$$

In other words, when do we have density of the linear space  $\mathbf{M}(\sigma)$  in the space  $C(\mathbb{R}^n)$ , in the topology of uniform convergence on compact sets? In this case we will say that  $\sigma$  has the *density property*.

Density is a necessary, but not a sufficient, condition. In fact, it theoretically represents the ability to approximate well, but it does not imply that there is a valid and efficient approximation scheme. Lack of density implies that it is impossible to approximate a large class of functions. Nevertheless, it is must be understood that density does not imply that one can approximate every function well from

$$(2.15) \quad \mathbf{M}_r(\sigma) = \left\{ \sum_{i=1}^r c_i \sigma(\omega^{(i)} \cdot \mathbf{x} - \theta^{(i)}) | c_i, \theta^{(i)} \in \mathbb{R}, \omega^{(i)} \in \mathbb{R}^n \right\},$$

for some fixed  $r$ . On the contrary, there is generally a lower bound on the degree to which one can approximate using  $\mathbf{M}_r(\sigma)$ , independent of the choice of  $\sigma$ . In the 1980s-90s many papers have been published which considered the density problem for general classes of activation functions:

- Gallant and White (1988) constructed a specific continuous, non-decreasing sigmoidal function (called *cosine squasher*) from which it was possible to obtain any trigonometric (Fourier) series.

- Irie and Miyake (1988) constructed an integral representation for any  $f \in L^1(\mathbb{R}^n)$  by using a kernel of the form  $\sigma(\omega \cdot x - \theta)$ , where  $\sigma$  was an arbitrary function in  $L^1(\mathbb{R})$ .
- Carroll and Dickinson (1989) used a discretized inverse Radon transform to approximate  $L^2$  functions with compact support in the  $L^2$  norm, using any continuous sigmoidal function as an activation function.
- The main result of Cybenko (1989) is the density property, in the uniform norm on compact sets, for any continuous sigmoidal function. (Cybenko does not demand monotonicity in his definition of sigmoidality.)
- Funahashi (1989) (independently of Cybenko (1989)) proved the density property, in the uniform norm on compact sets, for any continuous monotone sigmoidal function. He noted that, for  $\sigma$  continuous, monotone and bounded, it follows that  $\sigma(\cdot + \alpha) - \sigma(\cdot + \beta) \in L^1(\mathbb{R})$ ,  $\forall \alpha, \beta \in \mathbb{R}$ .
- Stinchcombe and White (1989) proved that  $\sigma$  has the density property for every  $\sigma \in L^1(\mathbb{R})$  with  $\int_{-\infty}^{+\infty} \sigma(t) dt \neq 0$ .
- Cotter (1990) considers different types of models and non-sigmoidal activation functions, as for example  $\sigma(t) = e^t$ , to obtain density.
- Jones (1990) shows, using ridge functions (which we will soon define), that to answer the question of density it suffices to consider only the univariate problem. He then proves that a bounded (not necessarily monotone or continuous) sigmoidal activation function suffices.
- Hornik (1991) proved density for any continuous bounded and nonconstant activation function, and also in other norms.
- Itô (1991-1992) studied the problem of density using monotone sigmoidal functions, with only weights of norm 1. He also considers conditions under which one obtains uniform convergence on all of  $\mathbb{R}^n$ .
- Mhaskar and Micchelli (1992) extended the density result to what they call *k-th degree sigmoidal functions*. They prove that if  $\sigma$  is continuous, bounded by some polynomial of degree  $k$  on all of  $\mathbb{R}$ , and

$$\lim_{t \rightarrow -\infty} \frac{\sigma(t)}{t^k} = 0, \quad \lim_{t \rightarrow +\infty} \frac{\sigma(t)}{t^k} = 1,$$

then density holds if and only if  $\sigma$  is not a polynomial.

As we have noted, a variety of techniques were used to attack a problem which many considered important and difficult. The solution to this problem, however, turns out to be surprisingly simple. Leshno, Lin, Pinkus and Schocken (1993) proved that the necessary and sufficient condition, for any continuous activation function, to have the density property is that it not be a polynomial.

#### 2.4.2 Theorem (Universal function approximation theorem for Neural Networks)

Let  $f \in C(K, \mathbb{R}^n)$  be a continuous function from a compact subset  $K \subseteq \mathbb{R}^m$  to  $\mathbb{R}^n$ , then there exists a Neural Network  $N_\sigma$  with a single hidden layer and activation function  $\sigma \in C(\mathbb{R}, \mathbb{R})$  which satisfies

$$\max_{x \in K} \|N_\sigma(x) - f(x)\| < \epsilon,$$

for some  $\epsilon > 0$ , if and only if  $\sigma$  is not a polynomial.

The proof of this theorem is rather long and we will first need to establish several lemmas and propositions.

#### Remark

If  $\sigma$  is a polynomial then density cannot possibly hold. In fact, if  $\sigma$  is a polynomial of degree  $m$ , then for every choice of  $\omega \in \mathbb{R}^n$  and  $\theta \in \mathbb{R}$ , we have that  $\sigma(\omega \cdot x - \theta)$  is a multivariate polynomial of total degree at most  $m$ , then

$$M(\sigma) = \text{span}\{\sigma(\omega \cdot x - \theta) | \theta \in \mathbb{R}, \omega \in \mathbb{R}^n\}$$

becomes the space of all polynomials of total degree  $m$  and does not span  $C(\mathbb{R}^n)$ .

#### 2.4.3 Definition of ridge function

A function  $g: \mathbb{R}^n \rightarrow \mathbb{R}$  is said to be a *ridge function*, if there exists a function  $f: \mathbb{R} \rightarrow \mathbb{R}$  and  $\omega = (\omega_1, \dots, \omega_n) \in \mathbb{R}^n$  (called *direction*), such that

$$(2.16) \quad g(x) = f(\omega \cdot x), \quad \forall x \in \mathbb{R}^n.$$

In other words, they are multivariate constant functions on the parallel hyperplanes  $\omega \cdot x = c$ ,  $c \in \mathbb{R}$ .

## Notation

We will denote by

$$(2.17) \quad R(\Omega) = \text{span}\{g(\omega \cdot x) | g \in C(\mathbb{R}), \omega \in \Omega\}$$

the span of all ridge functions defined over some set  $\Omega \subseteq \mathbb{R}^n$ .

Ridge functions are relevant in the theory of the single hidden layer perceptron model since each factor  $\sigma(\omega \cdot x - \theta)$  is a ridge function for every choice of  $\sigma$ ,  $\omega$ , and  $\theta$ .

### 2.4.4 Theorem (Vostrecov, Kreines)

The set of ridge functions

$$R(\Omega) = \text{span}\{g(\omega \cdot x) | g \in C(\mathbb{R}), \omega \in \Omega\},$$

defined over some compact set  $\Omega \subseteq \mathbb{R}^n$ , is dense in  $C(\mathbb{R}^n)$ , in the topology of uniform convergence over  $\Omega$ , if and only if there are no non-trivial homogeneous polynomials that vanishes on  $\Omega$ .

## Proof

For the proof, see [13].

## Remark

Because of the homogeneity of the directions (allowing a direction  $\omega$  is equivalent to allowing all directions  $\mu\omega$ , for every real  $\mu$ , since we vary over all  $g \in C(\mathbb{R}^n)$ ), it in fact suffices to consider directions normalized to lie on the unit ball

$$S^{n-1} = \{\mathbf{y} \mid \|\mathbf{y}\|_2 = 1\} \subset \mathbb{R}^n.$$

Then the previous theorem is equivalent to says that  $R(\Omega)$ , defined over some compact set  $\Omega \subseteq S^{n-1}$ , is dense in  $C(\mathbb{R}^n)$ , if no non-trivial homogeneous polynomial has zero set containing  $\Omega$ . In what follows we will always assume that  $\Omega \subseteq S^{n-1}$ .

## Notation

Let  $\Lambda$  be a subset of  $\mathbb{R}$ . We will denote with  $\Lambda \times \Omega$  the subset of  $\mathbb{R}^n$  given by

$$(2.18) \quad \Lambda \times \Omega = \{\lambda \omega \mid \lambda \in \Lambda, \omega \in \Omega\}.$$

### 2.4.5 Proposition 1 (about density)

Assume that  $\Lambda, \Theta$  are subsets of  $\mathbb{R}$  for which

$$(2.19) \quad N(\sigma; \Lambda, \Theta) = \text{span}\{\sigma(\lambda t - \theta) | \lambda \in \Lambda, \theta \in \Theta\},$$

is dense in  $C(\mathbb{R})$ , in the topology of uniform convergence on compact sets. Let us assume, in addition, that  $\Omega \subseteq S^{n-1}$  is such that  $R(\Omega)$  is dense in  $C(\mathbb{R}^n)$ , in the topology of uniform convergence on compact sets. Then

$$(2.20) \quad M(\sigma; \Lambda \times \Omega, \Theta) = \text{span}\{\sigma(\omega \cdot x - \theta) | \omega \in \Lambda \times \Omega, \theta \in \Theta\}$$

is dense in  $C(\mathbb{R}^n)$ , in the topology of uniform convergence on compact sets.

#### Proof

Let  $f \in C(K)$  for some compact set  $K \in \mathbb{R}^n$ . Since  $R(\Omega)$  is dense in  $C(K)$ , given  $\epsilon > 0$ , there exist  $g_i \in C(\mathbb{R})$  and  $\omega^i \in \Omega$ ,  $i = 1, \dots, r$ , such that

$$\left| f(x) - \sum_{i=1}^r g_i(\omega^i \cdot x) \right| < \frac{\epsilon}{2},$$

for all  $x \in K$ . Since  $K$  is compact, then

$$\{\omega^i \cdot x \mid x \in K\} \subseteq [\alpha_i, \beta_i],$$

for some finite interval  $[\alpha_i, \beta_i]$ ,  $i = 1, \dots, r$ .

Now, since  $N(\sigma; \Lambda, \Theta)$  is dense in  $C([\alpha_i, \beta_i])$ ,  $i = 1, \dots, r$ , there exist constants

$$c_{ij} \in \mathbb{R}, \quad \lambda_{ij} \in \Lambda, \quad \theta_{ij} \in \Theta, \quad j = 1, \dots, m_i, i = 1, \dots, r,$$

for which

$$\left| g_i(t) - \sum_{j=1}^{m_i} c_{ij} \sigma(\lambda_{ij} t - \theta_{ij}) \right| < \frac{\epsilon}{2r},$$

for all  $t \in [\alpha_i, \beta_i]$ , and  $i = 1, \dots, r$ . Then

$$\left| f(x) - \sum_{i=1}^r \sum_{j=1}^{m_i} c_{ij} \sigma(\lambda_{ij} \omega^i \cdot x - \theta_{ij}) \right| < \epsilon,$$

for all  $x \in K$ . ■

### 2.4.6 Proposition 2 (about density)

Let  $\sigma \in C^\infty(\mathbb{R})$  and assume  $\sigma$  is not polynomial. Then  $N(\sigma; \mathbb{R}, \mathbb{R})$  is dense in  $C(\mathbb{R})$ .

### Proof

It is known that, if  $\sigma \in C^\infty$  on any open interval and is not polynomial, then there exists a point  $-\theta_0$  in the interval for which  $\sigma^{(k)}(-\theta_0) \neq 0$  for all  $k = 0, 1, 2, \dots$ , see Corominas and Sunyer Balaguer (1954).

Since  $\sigma \in C^\infty(\mathbb{R})$ , and

$$\frac{[\sigma((\lambda+h)t-\theta_0)-\sigma(\lambda t-\theta_0)]}{h} \in N(\sigma; \mathbb{R}, \mathbb{R}) \text{ for all } h \neq 0,$$

it follows that

$$\left. \frac{d}{d\lambda} \sigma(\lambda t - \theta_0) \right|_{\lambda=0} = t \sigma'(-\theta_0) \in \overline{N(\sigma; \mathbb{R}, \mathbb{R})}.$$

Similarly

$$\left. \frac{d^k}{d\lambda^k} \sigma(\lambda t - \theta_0) \right|_{\lambda=0} = t^k \sigma^{(k)}(-\theta_0) \in \overline{N(\sigma; \mathbb{R}, \mathbb{R})}, \quad \forall k.$$

Since  $\sigma^{(k)}(-\theta_0) \neq 0$ ,  $k = 0, 1, 2, \dots$ , the set  $\overline{N(\sigma; \mathbb{R}, \mathbb{R})}$  contains all monomials and thus all polynomials. By the Stone-Weierstrass Theorem this implies that  $N(\sigma; \mathbb{R}, \mathbb{R})$  is dense in  $C(K)$  for every compact set  $K \subset \mathbb{R}$ . ■

Now, let us restate the previous Proposition in a more general form.

### Corollary 1

Let  $\Lambda$  be any set containing a sequence of values tending to zero, and let  $\Theta$  be any open interval. Let  $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  be such that  $\sigma \in C^\infty(\Theta)$ , and  $\sigma$  is not a polynomial on  $\Theta$ . Then  $N(\sigma; \Lambda, \Theta)$  is dense in  $C(\mathbb{R})$ .

### Corollary 2

Let

$$(2.21) \quad N_r(\sigma) = \left\{ \sum_{i=1}^r c_i \sigma(\lambda_i t - \theta_i) \mid c_i, \lambda_i, \theta_i \in \mathbb{R} \right\}.$$

If  $\Theta$  is any open interval and  $\sigma \in C^\infty(\Theta)$  is not a polynomial on  $\Theta$ , then  $\overline{N_r(\sigma)}$  contains  $\pi_{r-1}$  (the linear space of algebraic polynomials of degree at most  $r-1$ ).

#### 2.4.7 Proposition 3 (about density)

Let  $\sigma \in C(\mathbb{R})$  be not a polynomial, then  $N(\sigma; \mathbb{R}, \mathbb{R})$  is dense in  $C(\mathbb{R})$ .

## Proof

For each  $\phi \in C_0^\infty(\mathbb{R})$ , let

$$\sigma_\phi(t) = (\sigma * \phi)(t) = \int_{-\infty}^{+\infty} \sigma(t-y)\phi(y)dy,$$

be the convolution of  $\sigma$  and  $\phi$ . Since  $\sigma, \phi \in C(\mathbb{R})$  and  $\phi$  has compact support, the above integral converges for all  $t$ , and can be seen that

$$\sigma_\phi \in \overline{N(\sigma; \{I\}, \mathbb{R})} \cap C^\infty(\mathbb{R}).$$

Since

$$\sigma_\phi(\lambda t - \theta) = \int_{-\infty}^{+\infty} \sigma(\lambda t - \theta - y)\phi(y)dy,$$

for each  $\lambda \in \mathbb{R}$ , it also follows that

$$\overline{N(\sigma_\phi; \mathbb{R}, \mathbb{R})} \subset \overline{N(\sigma; \mathbb{R}, \mathbb{R})}.$$

Since  $\sigma_\phi \in C^\infty(\mathbb{R})$ , we have, from the proof of **proposition 2.4.6**, that

$$t^k \sigma_\phi^{(k)}(-\theta) \in \overline{N(\sigma_\phi; \mathbb{R}, \mathbb{R})}, \quad \forall \theta \in \mathbb{R}, \text{ and } \forall k.$$

Now, suppose by contradiction that  $N(\sigma; \mathbb{R}, \mathbb{R})$  is not dense in  $C(\mathbb{R})$ , then there exists  $k$  such that  $t^k \notin \overline{N(\sigma; \mathbb{R}, \mathbb{R})}$ . Thus  $t^k \notin \overline{N(\sigma_\phi; \mathbb{R}, \mathbb{R})}$  for each  $\phi \in C_0^\infty(\mathbb{R})$ . This implies that  $\sigma_\phi^{(k)}(-\theta) = 0$ ,  $\forall \theta \in \mathbb{R}$ , and  $\forall \phi \in C_0^\infty(\mathbb{R})$ . Then  $\sigma_\phi$  is a polynomial of degree at most  $k-1$  for each  $\phi \in C_0^\infty(\mathbb{R})$ . It is well known that there exist a sequence  $\{\phi_n\}_n \subset C_0^\infty(\mathbb{R})$  for which  $\sigma_{\phi_n}$  converges to  $\sigma$  uniformly on any compact set in  $\mathbb{R}$  (we can take, for example, the *mollifiers*, see Adams (1975) p.29). Polynomials of a fixed degree form a closed finite-dimensional linear subspace. Since  $\sigma_\phi$  is a polynomial of degree at most  $k-1$  for every  $\phi_n$ , it follows that  $\sigma$  is a polynomial of degree at most  $k-1$ . This contradicts our assumption. ■

Now, we are ready to prove the universal function approximation theorem for Neural Networks.

### Proof (Universal function approximation theorem for Neural Network)

By using the **proposition 2.4.7** and noting that if  $f$  is a polynomial, then the density property cannot hold, essentially proves the theorem. Setting the degree of

$f$  to any fixed positive integer implies that the linear span of the set of all Neural Networks with one hidden layer cannot span all of  $\mathbb{R}^n$ . ■

#### 2.4.8 Application

Let us consider the continuous function

$$f(x) = \cos(x),$$

defined over the interval  $[-\pi, \pi]$ . Our goal is to build a neural network that can approximate this function sufficiently well. For the moment, we consider this case only as an example to show the ability of Neural Networks. In the next paragraphs, more clarification will be given on how the learning process works.

Let us consider the L<sub>2</sub>-Loss Function

$$L_2(u(x; \omega, \theta), f(x)) = \int_{-\pi}^{\pi} (u(x) - f(x))^2 dx,$$

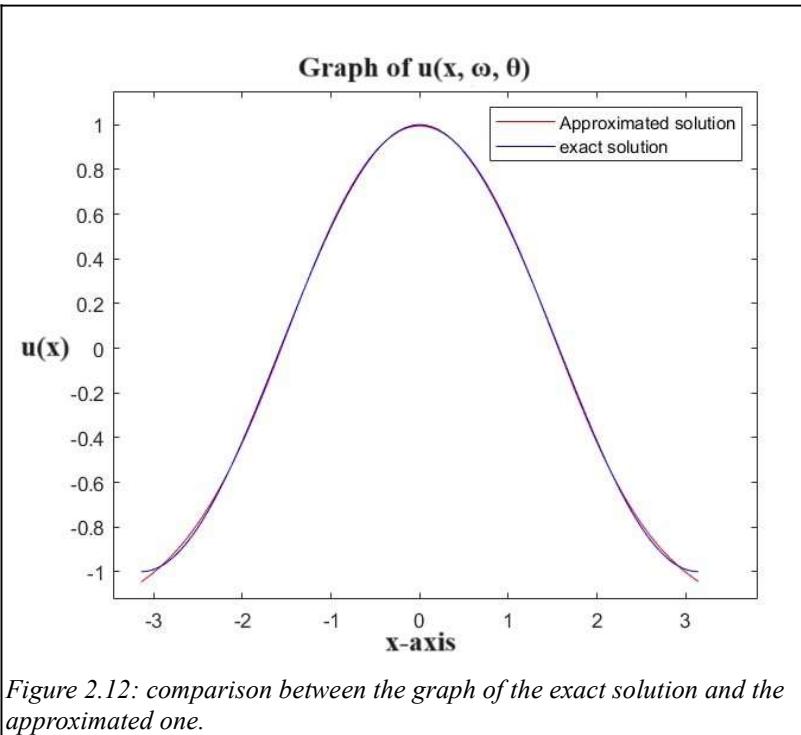
where  $u(x; \omega, \theta)$  represents the Neural Network. Let us set  $\{x_k\}_{k=1}^N$  uniformly sampled nodes on the interval  $[-\pi, \pi]$ , then the L<sub>2</sub> Loss Function can be approximated as

$$L_2(u(x; \omega, \theta), f(x)) \approx \frac{1}{N} \sum_{k=1}^N (u(x_k) - f(x_k))^2.$$

For the approximation of the function under study, we constructed a Neural Network consisting of an input node, two hidden layers each consisting of two nodes, and an output node. To train the Neural Network we set 11 uniformly sampled training nodes in the interval  $[-\pi, \pi]$ , and using the Backpropagation algorithm we obtained the following error

$$L_2(u(x; \omega, \theta), f(x)) = 9.5638 \times 10^{-4}.$$

Figure 2.12 shows the comparison between the graph of the function and that derived from the Neural Network.



For more details regarding Matlab codes, see Appendix A.

## 2.5 Neural Network for Solving Differential Equations

Our goal now is to construct and train a Neural Network so as to approximate the solution of a PDE.

### 2.5.1 A first approach

To fix ideas about how this process works, let us first consider the following one dimensional Cauchy's problem for a first-order ODE [1]:

$$(2.22) \quad \begin{cases} y'(x) = f(x, y(x)), \\ y(a) = y_a. \end{cases}$$

We can define a *trial function* of the type

$$(2.23) \quad \bar{y}(x) = y_a + (x - a)N(x; \omega, \theta),$$

where  $N(x; \omega, \theta)$ , represents the Neural Network we want to build, while  $\omega$  and  $\theta$  are the sets of weights and biases, respectively, which compose the Neural

Network. We observe that the trial function was created in such a way as to automatically satisfy the initial condition.

Suppose we construct a Neural Network with one input node, K hidden nodes and one output node. We denote by

- $\boldsymbol{\omega}^{(1)}$  the weights of the first layer;
- $\boldsymbol{\theta}$  the vector of biases;
- $\boldsymbol{\omega}^{(2)}$  the weights of the second level.

The output is linear, that is, no activation function is applied at the end. Obviously by construction we have that

$$N(x; \boldsymbol{\omega}, \boldsymbol{\theta}) = \sum_{j=1}^K \omega_j^{(2)} \sigma(\omega_j^{(1)} x + \theta_j),$$

from which, deriving with respect to  $x$ , we have

$$N'(x; \boldsymbol{\omega}, \boldsymbol{\theta}) = \sum_{j=1}^K \omega_j^{(2)} \omega_j^{(1)} \sigma'(\omega_j^{(1)} x + \theta_j).$$

By using these two expressions, it is possible to derive the derivative of the trial function, i.e.

$$\begin{aligned} \bar{y}'(x) &= N(x; \boldsymbol{\omega}, \boldsymbol{\theta}) + (x - a) N'(x; \boldsymbol{\omega}, \boldsymbol{\theta}) \\ &= \sum_{j=1}^K \omega_j^{(2)} \sigma(\omega_j^{(1)} x + \theta_j) + (x - a) \sum_{j=1}^K \omega_j^{(2)} \omega_j^{(1)} \sigma'(\omega_j^{(1)} x + \theta_j). \end{aligned}$$

*How it works:*

Suppose we have a training set  $\mathbf{x} = (x_1, \dots, x_M)$  composed of M points in the interval  $[a, b]$ . At each step, suppose we give as input the point  $x_j$ . Obviously the NN has not learned yet, so it will generate an error  $e_j = \bar{y}'(x_j) - f(x_j, \bar{y}(x_j))$ . We will use this error to improve the weights and biases via the Backpropagation algorithm, with the goal of minimizing the error. We denote by  $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_N)$  the set of test points in the interval  $[a, b]$ . We can now define our Loss Function as

$$(2.24) \quad L^2(\bar{\mathbf{x}}; \boldsymbol{\omega}, \boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^N [\bar{y}'(\bar{x}_i) - f(\bar{x}_i, \bar{y}(\bar{x}_i))]^2.$$

If now the total error is less than a fixed  $\epsilon$ , the algorithm stops, otherwise the procedure iterates until the solution is “good enough.” For more information on the results see [1], section 14.3.

### Remark

Let us observe that minimizing the Loss Function is equivalent to minimizing, for each point  $x_j$  in the training set, the quantity

$$(\bar{y}'(x_j) - f(x_j, \bar{y}(x_j)))^2.$$

Doing this, from a computational point of view, brings us a number of advantages, as we shall see later.

### 2.5.2 Neural Networks for PDEs

Let  $\mathbf{x} \in U \subset \mathbb{R}^n$ ,  $F$  be a given function and  $u: U \rightarrow \mathbb{R}$  the unknown of the second order differential equation [14]

$$(2.25) \quad F(\mathbf{x}, u(\mathbf{x}), \nabla u(\mathbf{x}), \Delta u(\mathbf{x})) = 0,$$

equipped with boundary conditions.

#### *Discretization*

First we need to discretize the domain of the function and its boundary by a grid of points  $\hat{U} = \{\mathbf{x}_i\}_i$  and  $\partial \hat{U} = \{\bar{\mathbf{x}}_j\}_j$ , respectively. Accordingly, the problem can be rewritten as a system of equations

$$(2.26) \quad F(\mathbf{x}_i, u(\mathbf{x}_i), \nabla u(\mathbf{x}_i), \Delta u(\mathbf{x}_i)) = 0, \quad \forall \mathbf{x}_i \in \hat{U},$$

with associated boundary conditions. We denote by  $\bar{u}(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})$  the test solution, where  $\boldsymbol{\omega}$  and  $\boldsymbol{\theta}$  are the sets of weights and biases, respectively, then the problem is to minimize the error

$$L^2(\hat{U}, \boldsymbol{\omega}, \boldsymbol{\theta}) = \sum_{\mathbf{x}_i \in \hat{U}} F(\mathbf{x}_i, \bar{u}(\mathbf{x}_i; \boldsymbol{\omega}, \boldsymbol{\theta}), \nabla \bar{u}(\mathbf{x}_i; \boldsymbol{\omega}, \boldsymbol{\theta}), \Delta \bar{u}(\mathbf{x}_i; \boldsymbol{\omega}, \boldsymbol{\theta}))^2,$$

subject to the constraints imposed by the boundary conditions.

#### *How to construct the test solution*

To construct the trial function  $\bar{u}(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})$ , we must keep in mind the boundary conditions to which the solution is subject, consequently we obtain

$$(2.27) \quad \bar{u}(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}) = A(\mathbf{x}) + f(\mathbf{x}, N(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})),$$

where  $A(\mathbf{x})$  is independent of  $\boldsymbol{\omega}$  and  $\boldsymbol{\theta}$ , and satisfies the boundary conditions,  $f(\mathbf{x}, N(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}))$  is constructed to satisfy the PDE, not contributing to the boundary condition that is already satisfied, while  $N(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})$  represents the

output of the Neural Network whose parameters must be adjusted so as to minimize the error.

### *Neural Networks training*

Suppose we have a training set  $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_M)$  composed of M points in the domain  $U$ . At each step, suppose we give as input the point  $\mathbf{x}_j$ . Obviously the NN has not yet learned, so it will generate an error

$$e_j = F(\mathbf{x}_j, u(\mathbf{x}_j), \nabla u(\mathbf{x}_j), \Delta u(\mathbf{x}_j)),$$

that will be used to improve the weights and biases via the Back Propagation algorithm, with the goal of minimizing the error  $L^2(\hat{U}, \omega, \theta)$  on the grid nodes. If the total error is less than a fixed  $\varepsilon$ , the algorithm stops, otherwise we iterate the process until the solution is “good enough”.

### 2.5.3 A different approach

Let  $\mathbf{x} \in \Omega \subset \mathbb{R}^n$  and  $u : U \rightarrow \mathbb{R}$  be the solution of the problem

$$(2.28) \quad \begin{cases} F(\mathbf{x}, u) = 0, & \text{in } \Omega, \\ B(\mathbf{x}, u) = 0, & \text{on } \partial\Omega, \end{cases}$$

where, for simplicity of expression, we have not specified the dependence on the derivatives. Following what was done in [12], our goal is to find the solution  $\hat{u}(\mathbf{x})$  of the problem (2.28), minimizing a given Loss Function  $L(u)$ . For this process to be valid and useful, the following properties must be satisfied:

*Property 1:* The function that minimizes  $L(u)$  must coincide with the solution of the PDE, i.e.

$$\arg \min_{u \in C^k} L(u) = \hat{u},$$

where k is the order of the equation under consideration.

Although Neural Networks are capable of approximating any function, this is not true in the case of finite Neural Networks. In particular, we cannot assume that the solution of any PDE can be expressed by a finite Neural Network. Consequently, the true minimum of L may not be reached, and so we must impose another requirement on L.

*Property 2:* For each  $\epsilon > 0$ , there exists  $\delta_\epsilon > 0$ , such that

$$L(u) - L(\hat{u}) < \delta_\epsilon \Rightarrow \|u - \hat{u}\| < \epsilon.$$

*Property 3:*  $L(u)$  has a unique global minimum.

*Property 4:* For any  $\epsilon > 0$ , there exists  $\delta_\epsilon > 0$ , such that

$$\|u - \hat{u}\| > \delta_\epsilon \Rightarrow L(u) - L(\hat{u}) > \epsilon.$$

We can now define the loss operators as

$$\begin{aligned} L_F^p(\mathbf{x}, u) &= \|F(\mathbf{x}, u)\|_p^p = \int_{\Omega} |F(\mathbf{x}, u)|^p d\mathbf{x}, \\ L_B^p(\mathbf{x}, u) &= \|B(\mathbf{x}, u)\|_p^p = \int_{\partial\Omega} |B(\mathbf{x}, u)|^p d\mathbf{x}. \end{aligned}$$

Our goal is to minimize both norms through a single Loss Function. The idea is to use the function

$$(2.29) \quad L(\mathbf{x}, u; p) = \lambda L_F(\mathbf{x}, u; p) + (1 - \lambda) L_B(\mathbf{x}, u; p),$$

where the parameter  $\lambda \in (0, 1)$  is called the *loss weight*, and means that the two operators can have different weights. A possible choice for this parameter is

$$(2.30) \quad \lambda = \frac{\|\partial\Omega\|}{\|\partial\Omega\| + \|\Omega\|}.$$

This parameter can be optimized through some methodologies, see [11].

It is possible to prove that the Loss Function (2.29) satisfies properties 1-4, that is, the following are valid:

### Theorem 1

The Loss Function (2.29) reaches its minimum when it is evaluated in the solution of the PDE.

### Theorem 2

For linear PDEs, the Loss Function (2.29), with  $p=2$ , has a unique global minimum point.

### Theorem 3

Let's consider the problem

$$\begin{cases} F(\mathbf{x}, u) = f(\mathbf{x}), \text{ in } \Omega, \\ B(\mathbf{x}, u) = b(\mathbf{x}), \text{ on } \partial\Omega, \end{cases}$$

and let  $\hat{u}$  be a solution. Then, for each  $\epsilon > 0$ , there exists  $\delta_\epsilon > 0$ , such that

$$L(u) - L(\hat{u}) < \delta_\epsilon \Rightarrow \|u - \hat{u}\| < \epsilon.$$

#### Theorem 4

Let's consider the problem

$$\begin{cases} F(\mathbf{x}, u) = f(\mathbf{x}), \text{ in } \Omega, \\ B(\mathbf{x}, u) = b(\mathbf{x}), \text{ on } \partial\Omega, \end{cases}$$

and let  $\hat{u}$  be a solution. Then, for each  $\epsilon > 0$ , there exists  $\delta_\epsilon > 0$ , such that

$$\|u - \hat{u}\| > \delta_\epsilon \Rightarrow L(u) > \epsilon.$$

#### 2.5.4 Mixed boundary conditions

Sometimes, we are in the case where the boundary  $\partial\Omega$  can be separated into several parts  $\partial\Omega_1, \dots, \partial\Omega_n$ , each of which is characterized by different boundary conditions, i.e.

$$(2.31) \quad B(\mathbf{x}, u) = \begin{cases} B_1(\mathbf{x}, u), & \mathbf{x} \in \partial\Omega_1, \\ \dots \\ B_n(\mathbf{x}, u), & \mathbf{x} \in \partial\Omega_n. \end{cases}$$

In this case, the loss operator  $L_B^p$  is defined as

$$L_B^p(\mathbf{x}, u) = \sum_i \gamma_i \int_{\partial\Omega_i} |B_i(\mathbf{x}, u)|^p d\mathbf{x},$$

where, as before, it was taken into account that different functionals may have different weights.

This theoretical framework must now be translated into practical algorithms that train Neural Networks to approximate the solution to our problem.

#### 2.5.5 Weight space

First we need to move from the space of functions to a space of weights, as it was done in [12]. To do this we introduce the parameters  $\omega, \theta$  which represent the weights and biases of the Neural Network, respectively. Through the Neural Network we define the parametrized function  $u(\mathbf{x}; \omega, \theta)$ . Using this expression, the Loss Function becomes

$$(2.32) \quad L^p(u(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})) = \lambda \int_{\Omega} |F(u(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}))|^p d\mathbf{x} + (1-\lambda) \int_{\partial\Omega} |B(u(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta}))|^p d\mathbf{x}.$$

Unfortunately, through this process, it is not assured that the properties examined above still hold, and this problem has not been solved yet. Unfortunately, there is still not much theoretical knowledge regarding the convergence of Neural Networks, and consequently we will not deal with it. Nevertheless, from an empirical point of view, this transition seems to be valid, especially for large Neural Networks. The hope is that, with the achievement of new theoretical results, it will become possible to show that, for  $N \rightarrow +\infty$ , where  $N$  is the number of degree of freedom of the network, the enunciated properties are still valid.

### 2.5.6 Monte Carlo method

The Monte Carlo method is used to approximate integrals by summations of functional evaluations, i.e.

$$(2.33) \quad \int_{\Omega} f(\mathbf{x}) d\mathbf{x} \approx \|\Omega\| \frac{1}{n} \sum_{i=1}^n f(\mathbf{x}_i),$$

where the  $\mathbf{x}_i$  are uniformly sampled points on  $\Omega$ .

Using the Monte Carlo method, the Loss Function is approximated as

$$(2.34) \quad L^p(u) \approx \lambda \|\Omega\| \frac{1}{n_F} \sum_{i=1}^{n_F} |F(u_i^F)|^p + (1-\lambda) \|\partial\Omega\| \frac{1}{n_B} \sum_{i=1}^{n_B} |B(u_i^B)|^p,$$

where, to simplify the notation, we placed

- $n_F$  is the number of uniformly sampled points on  $\Omega$ ;
- $n_B$  is the number of uniformly sampled points on  $\partial\Omega$ ;
- $u = u(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $u_i^F = u(\mathbf{x}_i^F; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $u_i^B = u(\mathbf{x}_i^B; \boldsymbol{\omega}, \boldsymbol{\theta})$ .

### Example

Let  $\Omega$  be a domain of  $\mathbb{R}^2$ , let us consider the Laplace's equation

$$(2.34) \quad \begin{cases} \Delta u(\mathbf{x}) = 0, & \text{in } \Omega, \\ u(\mathbf{x}) = G(\mathbf{x}), & \text{on } \partial\Omega. \end{cases}$$

To lead us back to a form of type (2.28), we rewrite (2.34) as

$$(2.35) \quad \begin{cases} \Delta u(\mathbf{x})=0, & \text{in } \Omega, \\ u(\mathbf{x})-G(\mathbf{x})=0, & \text{on } \partial\Omega. \end{cases}$$

The Loss Function will be

$$L^p(\mathbf{x}, u) = \lambda \int_{\Omega} |\Delta u(\mathbf{x})|^p d\mathbf{x} + (1-\lambda) \int_{\partial\Omega} |u(\mathbf{x}) - G(\mathbf{x})|^p d\mathbf{x}.$$

Let  $\Omega=[0, A] \times [0, B]$  be a square, then we can set  $\{x_i\}_{i=0}^n$  uniformly sampled points on the interval  $[0, A]$  and we can set  $\{y_j\}_{j=0}^m$  uniformly sampled points on the interval  $[0, B]$ . We thus constructed a grid of points on  $\Omega$  of the form  $\{(x_i, y_j)\}_{i=0, \dots, n; j=0, \dots, m}$ . Once we move to the space of weights, and using the Monte Carlo method, we derive

$$\begin{aligned} L^p(u(\mathbf{x}; \boldsymbol{\omega}, \boldsymbol{\theta})) &= \frac{\lambda AB}{(n+1)(m+1)} \sum_{i=0}^n \sum_{j=0}^m |\Delta(u((x_i, y_j); \boldsymbol{\omega}, \boldsymbol{\theta}))|^p + \\ &+ \frac{2(1-\lambda)(A+B)}{2(n+m)} \sum_{i=1}^{2(n+m)} |u(z_i; \boldsymbol{\omega}, \boldsymbol{\theta}) - G(z_i)|^p, \end{aligned}$$

where, to simplify the notation, we renamed the nodes on  $\partial\Omega$  with  $z_i$ ,  $i=1, \dots, 2(n+m)$ .

## 2.6 Neural Network methods for Solving Reaction-Diffusion models

### 2.6.1 Introduction

We have already mentioned that a reaction-diffusion model deals with the evolution, over time, of a population of individuals residing within an open spatial region  $\Omega \subset \mathbb{R}^d$ . Let us then consider the problem

$$(2.36) \quad \begin{cases} \frac{\partial}{\partial t} u(\mathbf{x}, t) - D \Delta u(\mathbf{x}, t) = f(u(\mathbf{x}, t)), & \text{in } \Omega \times (0, T), \\ u(\mathbf{x}, 0) = u_0(\mathbf{x}), & \text{on } \overline{\Omega} \times \{0\}, \\ u(\mathbf{x}, t) = b(\mathbf{x}, t), & \text{on } \partial\Omega \times (0, T]. \end{cases}$$

Our goal, as usual, is to find a solution  $\hat{u}(\mathbf{x}, t)$  of our problem, minimizing a given Loss Function  $L(u)$ . We can rewrite the problem as:

$$(2.37) \quad \begin{cases} \frac{\partial}{\partial t} u(\mathbf{x}, t) - D \Delta u(\mathbf{x}, t) - f(u(\mathbf{x}, t)) = 0, & \text{in } \Omega \times (0, T), \\ u(\mathbf{x}, 0) - u_0(\mathbf{x}) = 0, & \text{on } \overline{\Omega} \times \{0\}, \\ u(\mathbf{x}, t) - b(\mathbf{x}, t) = 0, & \text{on } \partial \Omega \times (0, T], \end{cases}$$

and, similarly to what was done previously, we define the Loss functionals

$$(2.38) \quad \begin{aligned} L_F^p(\mathbf{x}, t, u) &= \int_0^T dt \int_{\Omega} \left| \frac{\partial}{\partial t} u(\mathbf{x}, t) - D \Delta u(\mathbf{x}, t) - f(u(\mathbf{x}, t)) \right|^p d\mathbf{x}, \\ L_I^p(\mathbf{x}, u) &= \int_{\Omega} |u(\mathbf{x}, 0) - u_0(\mathbf{x})|^p d\mathbf{x}, \\ L_B^p(\mathbf{x}, t, u) &= \int_0^T dt \int_{\partial\Omega} |u(\mathbf{x}, t) - b(\mathbf{x}, t)|^p d\mathbf{x}. \end{aligned}$$

Obviously, our aim is to construct a single Loss Function that depends on the previous operators, then we place

$$(2.39) \quad L^p(\mathbf{x}, t, u) = \lambda_1 L_F^p(\mathbf{x}, t, u) + \lambda_2 L_I^p(\mathbf{x}, u) + \lambda_3 L_B^p(\mathbf{x}, t, u),$$

where the parameters  $\lambda_1, \lambda_2, \lambda_3$  are such that  $\lambda_1 + \lambda_2 + \lambda_3 = 1$  and they reflect the fact that the three functionals  $L_F, L_I, L_B$  can be of different importance. Now, we need to build and train a Neural Network to help us approximate the solution to our problem. In this case, both the components  $x_i$ , of the point  $\mathbf{x} \in \mathbb{R}^d$ , and time  $t$  are given as input of the Neural Network. We observe that, using a derivable activation function (e.g., logistic or tanh type), it is possible to calculate

$$\frac{\partial}{\partial t} u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta}), \quad \Delta u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta}).$$

Let  $\{t_i\}_{i=0}^m$  uniformly sampled points on the interval  $[0, T]$ ,  $\{x_k\}_{k \in K}$  be uniformly sampled points on  $\Omega$  and  $\{y_j\}_{j \in J}$  be the uniformly sampled points on  $\partial\Omega$ . Using the Monte Carlo method, we derive

$$\begin{aligned} L^p(u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta})) &\approx \frac{\lambda_1 T \|\Omega\|}{(m+1)|K|} \sum_{i=0}^m \sum_{k \in K} \left| \frac{\partial}{\partial t} u_{ki} - D \Delta u_{ki} - f(u_{ki}) \right|^p \\ &\quad + \frac{\lambda_2 \|\Omega\|}{|K|} \sum_{k \in K} |u_{k0} - u_0(x_k)|^p + \frac{\lambda_2 \|\partial\Omega\|}{|J|} \sum_{j \in J} |\bar{u}_{j0} - u_0(y_j)|^p \\ &\quad + \frac{\lambda_3 T \|\partial\Omega\|}{m|J|} \sum_{i=0}^m \sum_{j \in J} |\bar{u}_{ji} - b_{ji}|^p, \end{aligned}$$

where

- $u = u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $u_{ki} = u(x_k, t_i; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $u_{k0} = u(x_k, 0; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;

- $\overline{u}_{ji} = u(\mathbf{y}_j, t_i; \boldsymbol{\omega}, \boldsymbol{\theta});$
- $\overline{u}_{j0} = u(\mathbf{y}_j, 0; \boldsymbol{\omega}, \boldsymbol{\theta});$
- $\overline{b}_{ji} = u(\mathbf{y}_j, t_i; \boldsymbol{\omega}, \boldsymbol{\theta}).$

## 2.6.2 Calculation of derivatives

Suppose we have a Neural Network defined by:

- $m+1$  input nodes  $x_1, \dots, x_m, x_{m+1} = t;$
- $N$  hidden layers, consisting respectively of  $d_1, \dots, d_N$  nodes;
- one node output  $u;$
- $\boldsymbol{\omega}^{(I)}$  are the weights of the input nodes;
- $\boldsymbol{\omega}^{(j)}$  are the weights of the hidden layer  $j-1, j=2, \dots, N+1;$
- $\boldsymbol{\theta}^{(j)}$  are the biases of the layer  $j-1, j=1, \dots, N+1;$

For simplicity of notation we have posed  $t = x_{m+1}$  and we will denote by  $\mathbf{X}^{(0)} = (x_1, \dots, x_{m+1})^T$  the input vector.

We can schematize the operation of the Neural Network as follows

$$\begin{aligned}
& \mathbf{X}^{(0)} = (x_1, \dots, x_{m+1})^T, \\
& \mathbf{z}^{(I)} = \boldsymbol{\omega}^{(I)} \mathbf{X}^{(0)} + \boldsymbol{\theta}^{(I)}, \\
& \mathbf{X}^{(I)} = \sigma(\mathbf{z}^{(I)}), \\
& \vdots \\
& (2.41) \quad \mathbf{z}^{(k)} = \boldsymbol{\omega}^{(k)} \mathbf{X}^{(k-1)} + \boldsymbol{\theta}^{(k)}, \quad k=2, \dots, N, \\
& \mathbf{X}^{(k)} = \sigma(\mathbf{z}^{(k)}), \quad k=2, \dots, N, \\
& \vdots \\
& \mathbf{z}^{(N+1)} = \boldsymbol{\omega}^{(N+1)} \mathbf{X}^{(N)} + \boldsymbol{\theta}^{(N+1)}, \\
& u = \mathbf{z}^{(N+1)}.
\end{aligned}$$

Then, our Neural Network can be seen as

$$N_\sigma = \mathbf{z}^{(N+1)} \circ \sigma \circ \mathbf{z}^{(N)} \circ \sigma \circ \dots \circ \sigma \circ \mathbf{z}^{(2)} \circ \sigma \circ \mathbf{z}^{(1)},$$

that is  $u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta}) = \mathbf{z}^{(N+1)}(\sigma(\mathbf{z}^{(N)}(\sigma(\dots(\sigma(\mathbf{z}^{(1)}))\dots))))$ .

To compute the Loss Function, we need to calculate

$$\frac{\partial u}{\partial t}, \Delta u, f(u).$$

## Remark

Let us observe that

$$\frac{\partial \mathbf{z}_i^{(k)}}{\partial x_j^{(k-1)}} = \delta_{mj} \omega_{im}^{(k)} = \omega_{ij}^{(k)}, \quad k=1, \dots, N+1, i=1, \dots, d_k, j=1, \dots, d_{k-1}.$$

Because of what was said in the section on automatic differentiation, we get

$$\frac{\partial u}{\partial x_j} = \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \dots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(1)}} \sigma'(\mathbf{z}^{(1)}) \frac{\partial \mathbf{z}^{(1)}}{\partial x_j}, \quad j=1, \dots, m+1.$$

where  $\frac{\partial \mathbf{z}^{(I)}}{\partial x_j} = \boldsymbol{\omega}^{(I)}[:, j]$ , representing the  $j$ -th column of the weight matrix  $\boldsymbol{\omega}^{(I)}$ .

In addition, from the above observation, it follows that

$$(2.41) \quad \frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{X}^{(k-1)}} = \boldsymbol{\omega}^{(k)}, \quad k=1, \dots, N+1,$$

from which

$$(2.42) \quad \frac{\partial u}{\partial x_j} = \boldsymbol{\omega}^{(N+1)} \sigma'(\mathbf{z}^{(N)}) \odot \boldsymbol{\omega}^{(N)} \dots \boldsymbol{\omega}^{(2)} \sigma'(\mathbf{z}^{(1)}) \odot \boldsymbol{\omega}^{(1)}[:, j], \quad j=1, \dots, m+1,$$

where  $\mathbf{a} \odot \mathbf{b}$  represents the componentwise multiplication.

Now, to compute  $\Delta u$ , we need to calculate the second derivatives:

$$(2.43) \quad \begin{aligned} \frac{\partial^2 u}{\partial x_j^2} &= \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_j} \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(1)}} \sigma'(\mathbf{z}^{(1)}) \frac{\partial \mathbf{z}^{(1)}}{\partial x_j} \\ &\quad + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \frac{\partial \sigma'(\mathbf{z}^{(N-1)})}{\partial x_j} \dots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(1)}} \sigma'(\mathbf{z}^{(1)}) \frac{\partial \mathbf{z}^{(1)}}{\partial x_j} \\ &\quad + \dots + \\ &\quad + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(1)}} \frac{\partial \sigma'(\mathbf{z}^{(1)})}{\partial x_j} \frac{\partial \mathbf{z}^{(1)}}{\partial x_j}. \end{aligned}$$

Proceeding with the calculations we obtain that

$$(2.44) \quad \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_j} = \sigma''(\mathbf{z}^{(k)}) \frac{\partial \mathbf{z}^{(k)}}{\partial x_j} = \sigma''(\mathbf{z}^{(k)}) \frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{X}^{(k-1)}} \frac{\partial \mathbf{X}^{(k-1)}}{\partial x_j},$$

where, if we use the logistic activation function, we set

$$(2.45) \quad \sigma''(\mathbf{z}^{(k)}) = \sigma(\mathbf{z}^{(k)}) \odot (1 - \sigma(\mathbf{z}^{(k)})) \odot (1 - 2\sigma(\mathbf{z}^{(k)})).$$

Observing that

$$\frac{\partial \mathbf{X}^{(k)}}{\partial x_j} = \frac{\partial \sigma(\mathbf{z}^{(k)})}{\partial x_j} = \sigma'(\mathbf{z}^{(k)}) \frac{\partial \mathbf{z}^{(k)}}{\partial x_j} = \sigma'(\mathbf{z}^{(k)}) \frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{X}^{(k-1)}} \dots \sigma'(\mathbf{z}^{(1)}) \frac{\partial \mathbf{z}^{(1)}}{\partial x_j},$$

we derive

$$\frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_j} = \sigma''(\mathbf{z}^{(k)}) \frac{\partial \mathbf{z}^{(k)}}{\partial \mathbf{X}^{(k-1)}} \sigma'(\mathbf{z}^{(k-1)}) \frac{\partial \mathbf{z}^{(k-1)}}{\partial \mathbf{X}^{(k-2)}} \dots \sigma'(\mathbf{z}^{(I)}) \frac{\partial \mathbf{z}^{(I)}}{\partial x_j}.$$

### Remark

It is useful to define derivatives recursively as:

$$\begin{aligned}\frac{\partial \mathbf{X}^{(I)}}{\partial x_j} &= \sigma'(\mathbf{z}^{(I)}) \odot \boldsymbol{\omega}^{(I)}[:, j], \\ \frac{\partial \mathbf{X}^{(2)}}{\partial x_j} &= \sigma'(\mathbf{z}^{(2)}) \odot \left( \boldsymbol{\omega}^{(2)} \frac{\partial \mathbf{X}^{(I)}}{\partial x_j} \right), \\ &\vdots \\ \frac{\partial \mathbf{X}^{(k)}}{\partial x_j} &= \sigma'(\mathbf{z}^{(k)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial \mathbf{X}^{(k-1)}}{\partial x_j} \right), \\ &\vdots \\ \frac{\partial \mathbf{X}^{(N)}}{\partial x_j} &= \sigma'(\mathbf{z}^{(N)}) \odot \left( \boldsymbol{\omega}^{(N)} \frac{\partial \mathbf{X}^{(N-1)}}{\partial x_j} \right), \\ \frac{\partial u}{\partial x_j} &= \boldsymbol{\omega}^{(N+1)} \frac{\partial \mathbf{X}^{(N)}}{\partial x_j},\end{aligned}$$

where  $\sigma'(\mathbf{z}^{(k)}) = \sigma(\mathbf{z}^{(k)}) \odot (1 - \sigma(\mathbf{z}^{(k)}))$ .

This recursive method is very practical, as it allows us to “easily” derive the second derivatives as well, in fact

$$\begin{aligned}\frac{\partial^2 \mathbf{X}^{(I)}}{\partial x_j^2} &= \sigma''(\mathbf{z}^{(I)}) \odot (\boldsymbol{\omega}^{(I)}[:, j])^{\odot 2}, \\ &\vdots \\ \frac{\partial^2 \mathbf{X}^{(k)}}{\partial x_j^2} &= \sigma''(\mathbf{z}^{(k)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial \mathbf{X}^{(k-1)}}{\partial x_j} \right)^{\odot 2} + \sigma'(\mathbf{z}^{(k)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial^2 \mathbf{X}^{(k-1)}}{\partial x_j^2} \right), \\ &\vdots \\ \frac{\partial^2 \mathbf{X}^{(N)}}{\partial x_j^2} &= \sigma''(\mathbf{z}^{(N)}) \odot \left( \boldsymbol{\omega}^{(N)} \frac{\partial \mathbf{X}^{(N-1)}}{\partial x_j} \right)^{\odot 2} + \sigma'(\mathbf{z}^{(N)}) \odot \left( \boldsymbol{\omega}^{(N)} \frac{\partial^2 \mathbf{X}^{(N-1)}}{\partial x_j^2} \right), \\ \frac{\partial^2 u}{\partial x_j^2} &= \boldsymbol{\omega}^{(N+1)} \frac{\partial^2 \mathbf{X}^{(N)}}{\partial x_j^2}.\end{aligned}$$

Consequently, we obtain that

$$\Delta u = \sum_{j=1}^m \boldsymbol{\omega}^{(N+1)} \frac{\partial^2 \mathbf{X}^{(N)}}{\partial x_j^2}$$

where  $\frac{\partial^2 \mathbf{X}^{(N)}}{\partial x_i^2}$  was calculated iteratively using the procedure defined above.

### 2.6.3 Calculation of $\frac{\partial L^2}{\partial \omega_{ij}^{(k)}}$ , $\frac{\partial L^2}{\partial \theta_i^{(k)}}$

By placing

$$L_1 = (u_t - D \Delta u - f(u))^2, \quad L_2 = (u(\mathbf{x}, 0) - u_0(\mathbf{x}))^2, \quad L_3 = (u - b)^2,$$

we can define the Loss Function as

$$(2.46) \quad L^2(u, u_t, \Delta u) = \lambda_1 \int_0^T dt \int_{\Omega} L_1 \, d\mathbf{x} + \lambda_2 \int_{\Omega} L_2 \, d\mathbf{x} + \lambda_3 \int_0^T dt \int_{\partial\Omega} L_3 \, ds,$$

where

- $u = u(\mathbf{x}, t; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $u(\mathbf{x}, 0) = u(\mathbf{x}, t=0; \boldsymbol{\omega}, \boldsymbol{\theta})$ ;
- $b = b(\mathbf{x}, t)$ ;
- $u_t, \Delta u$  were calculated using the procedures shown above.

As we mentioned earlier, minimizing the Loss Function is equivalent to minimizing, for each point  $\mathbf{x}$  in the training set, the quantities

$$L_1 = (u_t - D \Delta u - f(u))^2, \quad L_2 = (u(\mathbf{x}, 0) - u_0(\mathbf{x}))^2, \quad L_3 = (u - b)^2.$$

So, studying  $\frac{\partial L^2}{\partial \omega_{ij}^{(k)}}$  is equivalent to studying:

- $\frac{\partial L_1}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_1}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u_t - D \Delta u - f(u)) \left[ \frac{\partial u_t}{\partial \mathbf{z}^{(k)}} - D \frac{\partial \Delta u}{\partial \mathbf{z}^{(k)}} - f'(u) \frac{\partial u}{\partial \mathbf{z}^{(k)}} \right] \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}}$ ;
- $\frac{\partial L_2}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_2}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u(\mathbf{x}, 0) - u_0(\mathbf{x})) \left[ \frac{\partial u}{\partial \mathbf{z}^{(k)}} \Big|_{t=0} \right] \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}}$ ;
- $\frac{\partial L_3}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_3}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u - b) \frac{\partial u}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}}$ .

Similar considerations can be made for  $\frac{\partial L^2}{\partial \theta_i^{(k)}}$ .

Let's start by calculating

$$\frac{\partial u}{\partial z^{(k)}} = \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-I)}} \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}).$$

Let us observe that

$$z_r^{(k)} = \omega_{rs}^{(k)} x_s^{(k-I)} + \theta_r^{(k)},$$

from which

$$\frac{\partial z_r^{(k)}}{\partial \omega_{ij}^{(k)}} = \delta_{ri} \delta_{sj} x_s^{(k-I)} = \delta_{ri} x_j^{(k-I)},$$

that is

$$\frac{\partial z_i^{(k)}}{\partial \omega_{ij}^{(k)}} = x_j^{(k-I)}, \text{ and } \frac{\partial z_r^{(k)}}{\partial \omega_{ij}^{(k)}} = 0, \quad \forall r \neq i,$$

which implies, in a compact form, that

$$(2.48) \quad \frac{\partial z^{(k)}}{\partial \omega_{ij}^{(k)}} = [0, \dots, 0, x_j^{(k-I)}, 0, \dots, 0]^T,$$

with  $x_j^{(k-I)}$  in the i-th position. Similarly, it is shown that

$$(2.49) \quad \frac{\partial z^{(k)}}{\partial \theta_i^{(k)}} = [0, \dots, 0, 1, 0, \dots, 0]^T,$$

with 1 in the i-th position.

Now, for all  $r=1, \dots, m+I$ , let's calculate

$$\begin{aligned} \frac{\partial u_{x_r}}{\partial z^{(k)}} &= \frac{\partial}{\partial z^{(k)}} \frac{\partial u}{\partial x_r} = \frac{\partial}{\partial x_r} \frac{\partial u}{\partial z^{(k)}} \\ &= \frac{\partial}{\partial x_r} \left( \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-I)}} \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right) \\ &= \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r} \frac{\partial z^{(N)}}{\partial X^{(N-I)}} \sigma'(\mathbf{z}^{(N-I)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\ &\quad + \dots + \\ &+ \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial z^{(s+I)}}{\partial X^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \frac{\partial z^{(s)}}{\partial X^{(s-I)}} \sigma'(\mathbf{z}^{(s-I)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\ &\quad + \dots + \\ &+ \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-I)}} \sigma'(\mathbf{z}^{(N-I)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r}, \end{aligned}$$

where

$$\begin{aligned}
\frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} &= \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} = \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \frac{\partial \mathbf{X}^{(s-1)}}{\partial x_r} \\
&= \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \frac{\partial \mathbf{z}^{(s-1)}}{\partial x_r} \\
&= \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \frac{\partial \mathbf{z}^{(s-1)}}{\partial \mathbf{X}^{(s-2)}} \sigma'(\mathbf{z}^{(s-2)}) \dots \sigma'(\mathbf{z}^{(I)}) \frac{\partial \mathbf{z}^{(I)}}{\partial x_r}.
\end{aligned}$$

### Remark

We observe that

$$\frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} = \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r},$$

then we obtain the recursive formula

$$\frac{\partial \sigma'(\mathbf{z}^{(s+1)})}{\partial x_r} = \sigma''(\mathbf{z}^{(s+1)}) \frac{\partial \mathbf{z}^{(s+1)}}{\partial x_r} = \sigma''(\mathbf{z}^{(s+1)}) \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \sigma'(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r}.$$

Now, let us consider

$$(2.50) \quad \frac{\partial \Delta u}{\partial \mathbf{z}^{(k)}} = \frac{\partial}{\partial \mathbf{z}^{(k)}} \sum_{r=1}^m \frac{\partial^2 u}{\partial x_r^2} = \sum_{r=1}^m \frac{\partial^2}{\partial x_r^2} \left( \frac{\partial u}{\partial \mathbf{z}^{(k)}} \right),$$

Then we can study

$$\begin{aligned}
&\frac{\partial^2}{\partial x_r^2} \left( \frac{\partial u}{\partial \mathbf{z}^{(k)}} \right) = \frac{\partial}{\partial x_r} \left( \frac{\partial u}{\partial \mathbf{z}^{(k)}} \right) \\
&= \frac{\partial}{\partial x_r} \left( \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r} \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right. \\
&\quad \left. + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right. \\
&\quad \left. + \dots + \right. \\
&\quad \left. + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial \mathbf{z}^{(N)}}{\partial \mathbf{X}^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} \right).
\end{aligned} \tag{2.51}$$

By using the linearity of the derivative, we can calculate, for example

$$\begin{aligned}
&\frac{\partial}{\partial x_r} \left( \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right) \\
&= \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r} \dots \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\
&\quad + \dots + \\
&\quad + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(s)})}{\partial x_r^2} \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\
&\quad + \dots + \\
&\quad + \frac{\partial \mathbf{z}^{(N+1)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial \mathbf{z}^{(s+1)}}{\partial \mathbf{X}^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \frac{\partial \mathbf{z}^{(s)}}{\partial \mathbf{X}^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r},
\end{aligned}$$

where,

$$(2.52) \quad \begin{aligned} \frac{\partial^2 \sigma'(\mathbf{z}^{(s)})}{\partial x_r^2} &= \frac{\partial}{\partial x_r} \left( \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \right) = \frac{\partial}{\partial x_r} \left( \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} \right) \\ &= \sigma'''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} + \sigma''(\mathbf{z}^{(s)}) \frac{\partial^2 \mathbf{z}^{(s)}}{\partial x_r^2}, \end{aligned}$$

with

$$(2.53) \quad \sigma'''(\mathbf{z}^{(s)}) = \sigma(\mathbf{z}^{(s)}) \odot (1 - \sigma(\mathbf{z}^{(s)})) \odot (1 - 6 \sigma(\mathbf{z}^{(s)}) \odot (1 - \sigma(\mathbf{z}^{(s)}))),$$

in the case of using the logistic activation function.

Substituting first in (2.51) and then in (2.50) gives the expression of  $\frac{\partial \Delta u}{\partial \mathbf{z}^{(k)}}$ .

### Remark

Minimizing, for each point  $\mathbf{x}$  in the training set, the quantities

$$L_1 = (u_t - D \Delta u - f(u))^2, \quad L_2 = (u(\mathbf{x}, 0) - u_0(\mathbf{x}))^2, \quad L_3 = (u - b)^2,$$

is better than minimizing the Loss Function. In fact, let us suppose that we want to minimize

$$L^2(u, u_t, \Delta u) = \lambda_1 \int_0^T dt \int_{\Omega} L_1 \ d\mathbf{x} + \lambda_2 \int_{\Omega} L_2 \ d\mathbf{x} + \lambda_3 \int_0^T dt \int_{\partial\Omega} L_3 \ ds,$$

then its derivative is

$$\frac{\partial L^2}{\partial \omega_{ij}^{(k)}} = \lambda_1 \int_0^T dt \int_{\Omega} \frac{\partial L_1}{\partial \omega_{ij}^{(k)}} \ d\mathbf{x} + \lambda_2 \int_{\Omega} \frac{\partial L_2}{\partial \omega_{ij}^{(k)}} \ d\mathbf{x} + \lambda_3 \int_0^T dt \int_{\partial\Omega} \frac{\partial L_3}{\partial \omega_{ij}^{(k)}} \ ds.$$

If we now apply the gradient method, our Neural Network would update its weights based on the total error and not from the error produced by the individual point. This implies that the Neural Network learns that a particular configuration of points generates a particular error, without understanding from which specific point the error comes.

## 2.6.4 Recursive representation

From a purely computational point of view, we want to simplify the work for the calculator by preventing it from recalculating terms that have already been calculated. We set

$$(2.54) \quad \alpha^{(k)} = \frac{\partial u}{\partial z^{(k)}}, \quad \beta_r^{(k)} = \frac{\partial u_{x_r}}{\partial z^{(k)}}, \quad r=1, \dots, m+1, \quad \gamma^{(k)} = \frac{\partial \Delta u}{\partial z^{(k)}},$$

and we want to find recursive formulas to connect

- $\alpha^{(k)}$  with  $\alpha^{(k+1)}$ ,
- $\beta_r^{(k)}$  with  $\beta_r^{(k+1)}$ ,  $r=1, \dots, m+1$ ,
- $\gamma^{(k)}$  with  $\gamma^{(k+1)}$ .

Observing that

$$\beta_r^{(k)} = \frac{\partial}{\partial x_r}(\alpha^{(k)}), \quad r=1, \dots, m+1, \quad \text{and} \quad \gamma^{(k)} = \sum_{r=1}^m \frac{\partial}{\partial x_r}(\beta_r^{(k)}),$$

and by using hidden layers of the same dimensions, proceeding with the calculations, we obtain that

$$\begin{aligned} \bullet \quad \alpha^{(k)} &= \frac{\partial u}{\partial z^{(k)}} = \frac{\partial z^{(N+1)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-1)}} \dots \sigma'(\mathbf{z}^{(k+1)}) \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\ &= \frac{\partial u}{\partial z^{(k+1)}} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) = \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}); \\ \bullet \quad \beta_r^{(k)} &= \frac{\partial}{\partial x_r}(\alpha^{(k)}) = \frac{\partial}{\partial x_r}(\alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)})) \\ &= \frac{\partial}{\partial x_r}(\alpha^{(k+1)}) \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \alpha^{(k+1)} \frac{\partial}{\partial x_r} \left( \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right) \\ &= \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r}, \quad r=1, \dots, m+1; \\ \bullet \quad \gamma^{(k)} &= \sum_{r=1}^m \frac{\partial}{\partial x_r}(\beta_r^{(k)}) = \sum_{r=1}^m \frac{\partial}{\partial x_r} \left( \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} \right) \\ &= \sum_{r=1}^m \left[ \frac{\partial}{\partial x_r}(\beta_r^{(k+1)}) \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} \right. \\ &\quad \left. + \frac{\partial}{\partial x_r}(\alpha^{(k+1)}) \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} + \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(k)})}{\partial x_r^2} \right] \\ &= \sum_{r=1}^m \left[ \frac{\partial}{\partial x_r}(\beta_r^{(k+1)}) \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} \right. \\ &\quad \left. + \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} + \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(k)})}{\partial x_r^2} \right] \\ &= \gamma^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) + 2 \sum_{r=1}^m \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} + \alpha^{(k+1)} \sum_{r=1}^m \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(k)})}{\partial x_r^2}. \end{aligned}$$

### Remark

It may also be useful to find a numerical method to calculate

- $\frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r}, \quad r=1, \dots, m+1;$
- $\frac{\partial^2 \sigma'(\mathbf{z}^{(s)})}{\partial x_r^2}, \quad r=1, \dots, m+1.$

We observe that

$$\frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} = \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} = \sigma''(\mathbf{z}^{(s)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial \sigma(\mathbf{z}^{(s-1)})}{\partial x_r} \right) = \sigma''(\mathbf{z}^{(s)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial \mathbf{X}^{(s-1)}}{\partial x_r} \right),$$

where

- $\frac{\partial \mathbf{X}^{(s-1)}}{\partial x_r} = \boldsymbol{\omega}^{(s-1)} \sigma'(\mathbf{z}^{(s-2)}) \odot \boldsymbol{\omega}^{(s-2)} \dots \boldsymbol{\omega}^{(2)} \sigma'(\mathbf{z}^{(I)}) \odot \boldsymbol{\omega}^{(I)}[:, r], \quad r=1, \dots, m+1,$
- $\sigma'(\mathbf{z}^{(k)}) = \sigma(\mathbf{z}^{(k)}) \odot (1 - \sigma(\mathbf{z}^{(k)})),$
- $\sigma''(\mathbf{z}^{(k)}) = \sigma(\mathbf{z}^{(k)}) \odot (1 - \sigma(\mathbf{z}^{(k)})) \odot (1 - 2\sigma(\mathbf{z}^{(k)})),$

in the case of the sigmoid activation function.

Moreover, we have that

$$\begin{aligned} \frac{\partial^2 \sigma'(\mathbf{z}^{(s)})}{\partial x_r^2} &= \frac{\partial}{\partial x_r} \left( \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial x_r} \right) = \frac{\partial}{\partial x_r} \left( \sigma''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} \right) = \sigma'''(\mathbf{z}^{(s)}) \frac{\partial \mathbf{z}^{(s)}}{\partial x_r} + \sigma''(\mathbf{z}^{(s)}) \frac{\partial^2 \mathbf{z}^{(s)}}{\partial x_r^2} \\ &= \sigma'''(\mathbf{z}^{(s)}) \odot \left( \boldsymbol{\omega}^{(s)} \frac{\partial \sigma(\mathbf{z}^{(s-1)})}{\partial x_r} \right) + \sigma''(\mathbf{z}^{(s)}) \left( \boldsymbol{\omega}^{(s)} \frac{\partial^2 \sigma(\mathbf{z}^{(s-1)})}{\partial x_r^2} \right) \\ &= \sigma'''(\mathbf{z}^{(s)}) \odot \left( \boldsymbol{\omega}^{(s)} \frac{\partial \mathbf{X}^{(s-1)}}{\partial x_r} \right) + \sigma''(\mathbf{z}^{(s)}) \left( \boldsymbol{\omega}^{(s)} \frac{\partial^2 \mathbf{X}^{(s-1)}}{\partial x_r^2} \right), \end{aligned}$$

where

- $\frac{\partial \mathbf{X}^{(s-1)}}{\partial x_r}$  is the same as above;
- $\frac{\partial^2 \mathbf{X}^{(s-1)}}{\partial x_r^2}$  is calculate by using the recursively procedure
  - $\frac{\partial^2 \mathbf{X}^{(k)}}{\partial x_j^2} = \sigma''(\mathbf{z}^{(k)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial \mathbf{X}^{(k-1)}}{\partial x_j} \right)^{\odot 2} + \sigma'(\mathbf{z}^{(k)}) \odot \left( \boldsymbol{\omega}^{(k)} \frac{\partial^2 \mathbf{X}^{(k-1)}}{\partial x_j^2} \right),$
  - $\frac{\partial^2 \mathbf{X}^{(I)}}{\partial x_j^2} = \sigma''(\mathbf{z}^{(I)}) \odot (\boldsymbol{\omega}^{(I)}[:, j])^{\odot 2},$
- $\sigma'(\mathbf{z}^{(k)}), \sigma''(\mathbf{z}^{(k)})$  are the same as above,
- $\sigma'''(\mathbf{z}^{(s)}) = \sigma(\mathbf{z}^{(s)}) \odot (1 - \sigma(\mathbf{z}^{(s)})) \odot (1 - 6\sigma(\mathbf{z}^{(s)}) \odot (1 - \sigma(\mathbf{z}^{(s)}))).$

## 2.6.5 How it works

The first step is to approximate  $L^2$  by the Monte Carlo method:

$$\begin{aligned} L^2 \approx & \lambda_1 \frac{T\|\Omega\|}{|I||J|} \sum_{i \in I} \sum_{j \in J} L_1(x_i, t_j) \\ & + \lambda_2 \left[ \frac{\|\Omega\|}{|I|} \sum_{i \in I} L_2(x_i, 0) + \frac{\|\partial\Omega\|}{|K|} \sum_{k \in K} L_2(\bar{x}_k, 0) \right] \\ & + \lambda_3 \frac{T\|\partial\Omega\|}{|K||J|} \sum_{k \in K} \sum_{j \in J} L_3(\bar{x}_k, t_j), \end{aligned}$$

where

- $x_i$  are uniformly sampled points on  $\Omega$ ;
- $\bar{x}_k$  are uniformly sampled points on  $\partial\Omega$ ;
- $t_j$  are uniformly sampled points on  $[0, T]$ .

Then, we need to calculate

- $\frac{\partial L_1}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_1}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u_t - D\Delta u - f(u))[\beta_{m+1}^{(k)} - D\gamma^{(k)} - f'(u)\alpha^{(k)}] \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}};$
- $\frac{\partial L_2}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_2}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u(\mathbf{x}, 0) - u_0(\mathbf{x}))[\alpha^{(k)}|_{t=0}] \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}};$
- $\frac{\partial L_3}{\partial \omega_{ij}^{(k)}} = \frac{\partial L_2}{\partial \mathbf{z}^{(k)}} \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}} = 2(u - b)[\alpha^{(k)}|_{\partial\Omega}] \frac{\partial \mathbf{z}^{(k)}}{\partial \omega_{ij}^{(k)}}.$

## 2.6.6 Updating weights process

To recapitulate, the process of updating weights is as follows:

- For each point  $(x, t)$  in the training set:
- Calculate

$$\begin{aligned} \alpha^{(N+1)} &= \frac{\partial u}{\partial \mathbf{z}^{(N+1)}} = I, \\ \beta_r^{(N+1)} &= \frac{\partial u_{x_r}}{\partial \mathbf{z}^{(N+1)}} = 0, \quad r = 1, \dots, m+1, \\ \gamma^{(N+1)} &= \frac{\partial \Delta u}{\partial \mathbf{z}^{(N+1)}} = 0, \end{aligned}$$

evaluated at the point  $(x, t)$ .

- Calculate

$$\frac{\partial L_1}{\partial \omega_{ij}^{(N+1)}}, \quad \frac{\partial L_2}{\partial \omega_{ij}^{(N+1)}}, \quad \frac{\partial L_3}{\partial \omega_{ij}^{(N+1)}},$$

evaluated at the point  $(x, t)$ .

- Update the weights  $\omega^{(N+I)}$ :

$$\omega_{ij}^{(N+I)} = \omega_{ij}^{(N+I)} - \eta \frac{\partial L_s}{\partial \omega_{ij}^{(N+I)}}.$$

Let us observe that the choice of which  $\frac{\partial L_s}{\partial \omega_{ij}^{(N+I)}}$  to calculate depends on

which point  $(x, t)$  has been fixed.

- Calculate

$$\begin{aligned}\alpha^{(N)} &= \alpha^{(N+I)} \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) = \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}), \\ \beta_r^{(N)} &= \beta_r^{(N+I)} \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) + \alpha^{(N+I)} \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r} \\ &= \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r}, \quad r=1, \dots, m+1, \\ \gamma^{(N)} &= \gamma^{(N+I)} \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \sigma'(\mathbf{z}^{(N)}) + 2 \sum_{r=1}^m \beta_r^{(N+I)} \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial x_r} \\ &\quad + \alpha^{(N+I)} \sum_{r=1}^m \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(N)})}{\partial x_r^2} \\ &= \sum_{r=1}^m \frac{\partial \mathbf{z}^{(N+I)}}{\partial \mathbf{X}^{(N)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(N)})}{\partial x_r^2},\end{aligned}$$

evaluated at the point  $(x, t)$ .

- Calculate

$$\frac{\partial L_1}{\partial \omega_{ij}^{(N)}}, \quad \frac{\partial L_2}{\partial \omega_{ij}^{(N)}}, \quad \frac{\partial L_3}{\partial \omega_{ij}^{(N)}},$$

evaluated at the point  $(x, t)$ .

- Update the weights  $\omega^{(N)}$ :

$$\omega_{ij}^{(N)} = \omega_{ij}^{(N)} - \eta \frac{\partial L_s}{\partial \omega_{ij}^{(N)}}.$$

Let us observe that the choice of which  $\frac{\partial L_s}{\partial \omega_{ij}^{(N)}}$  to calculate depends on

which point  $(x, t)$  has been fixed.

.....

- Calculate

$$\begin{aligned}\alpha^{(k)} &= \alpha^{(k+1)} \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}), \\ \beta_r^{(k)} &= \beta_r^{(k+1)} \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) + \alpha^{(k+1)} \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r}, \quad r=1, \dots, m+1, \\ \gamma^{(k)} &= \gamma^{(k+1)} \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \sigma'(\mathbf{z}^{(k)}) + 2 \sum_{r=1}^m \beta_r^{(k+1)} \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial x_r} \\ &\quad + \alpha^{(k+1)} \sum_{r=1}^m \frac{\partial \mathbf{z}^{(k+1)}}{\partial \mathbf{X}^{(k)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(k)})}{\partial x_r^2},\end{aligned}$$

evaluated at the point  $(x, t)$ .

- Calculate

$$\frac{\partial L_1}{\partial \omega_{ij}^{(k)}}, \quad \frac{\partial L_2}{\partial \omega_{ij}^{(k)}}, \quad \frac{\partial L_3}{\partial \omega_{ij}^{(k)}},$$

evaluated at the point  $(x, t)$ .

- Update the weights  $\boldsymbol{\omega}^{(k)}$ :

$$\omega_{ij}^{(k)} = \omega_{ij}^{(k)} - \eta \frac{\partial L_s}{\partial \omega_{ij}^{(k)}}.$$

Let us observe that the choice of which  $\frac{\partial L_s}{\partial \omega_{ij}^{(k)}}$  to calculate depends on

which point  $(x, t)$  has been fixed.

.....

- Calculate

$$\begin{aligned}\alpha^{(I)} &= \alpha^{(2)} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \sigma'(\mathbf{z}^{(I)}), \\ \beta_r^{(I)} &= \beta_r^{(2)} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \sigma'(\mathbf{z}^{(I)}) + \alpha^{(2)} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \frac{\partial \sigma'(\mathbf{z}^{(I)})}{\partial x_r}, \quad r=1, \dots, m+1, \\ \gamma^{(I)} &= \gamma^{(2)} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \sigma'(\mathbf{z}^{(I)}) + 2 \sum_{r=1}^m \beta_r^{(2)} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \frac{\partial \sigma'(\mathbf{z}^{(I)})}{\partial x_r} \\ &\quad + \alpha^{(2)} \sum_{r=1}^m \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{X}^{(I)}} \frac{\partial^2 \sigma'(\mathbf{z}^{(I)})}{\partial x_r^2},\end{aligned}$$

evaluated at the point  $(x, t)$ .

- Calculate

$$\frac{\partial L_1}{\partial \omega_{ij}^{(I)}}, \quad \frac{\partial L_2}{\partial \omega_{ij}^{(I)}}, \quad \frac{\partial L_3}{\partial \omega_{ij}^{(I)}},$$

evaluated at the point  $(x, t)$ .

- Update the weights  $\boldsymbol{\omega}^{(I)}$ :

$$\omega_{ij}^{(l)} = \omega_{ij}^{(l)} - \eta \frac{\partial L_s}{\partial \omega_{ij}^{(l)}}.$$

Let us observe that the choice of which  $\frac{\partial L_s}{\partial \omega_{ij}^{(l)}}$  to calculate depends on which point  $(x, t)$  has been fixed.

### 2.6.7 Algorithm (Backpropagation)

Input:  $N_{layer}$ ,  $\omega^{(1)}, \dots, \omega^{(N_{layer}+1)}$ ,  $\theta^{(1)}, \dots, \theta^{(N_{layer}+1)}$ ,  $\eta > 0$ ,  $L_{function}$ .

Learning data:  $X = \{X_k = (x_1, \dots, x_m, x_{m+1} = t)_k\}_{k=1, \dots, K}$ .

Parameters:  $error = 0.1$ ,  $errormax = 0.001$ ,  $threshold = 1$ ,  $thresholdmax = 1000$ .

while  $error > errormax$  &&  $threshold < thresholdmax$

for  $k = 1$  to  $K$

$$X_k^{(0)} = X_k;$$

for  $j = 1$  to  $N_{layer}$

$$X_k^{(j)} = \sigma(\omega^{(j)} X_k^{(j-1)} + \theta^{(j)});$$

end

$$u_k = X_k^{(N_{layer}+1)} = \omega^{(N_{layer}+1)} X_k^{(N_{layer})} + \theta^{(N_{layer}+1)}$$

$$\alpha^{(N_{layer}+1)} = 1;$$

$$\beta_r^{(N_{layer}+1)} = 0, r = 1, \dots, m+1;$$

$$\gamma^{(N_{layer}+1)} = 0;$$

$$\text{calculate } \frac{\partial L_s}{\partial z^{(N+1)}},$$

$$W^{(N_{layer}+1)} = [\theta^{(N_{layer}+1)}, \omega^{(N_{layer}+1)}];$$

$$W^{(N_{layer}+1)} = W^{(N_{layer}+1)} - \eta \frac{\partial L_s}{\partial z^{(N+1)}} [I; X^{(N_{layer})}]^T;$$

for  $k = N_{layer}$  to  $I$

$$\alpha^{(k)} = \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(z^{(k)});$$

$$\beta_r^{(k)} = \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(z^{(k)}) + \alpha^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(z^{(k)})}{\partial x_r}, \quad r = 1, \dots, m+1;$$

$$\begin{aligned} \gamma^{(k)} = & \gamma^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \sigma'(z^{(k)}) + 2 \sum_{r=1}^m \beta_r^{(k+1)} \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial \sigma'(z^{(k)})}{\partial x_r} \\ & + \alpha^{(k+1)} \sum_{r=1}^m \frac{\partial z^{(k+1)}}{\partial X^{(k)}} \frac{\partial^2 \sigma'(z^{(k)})}{\partial x_r^2}; \end{aligned}$$

```

calculate  $\frac{\partial L_s}{\partial z^{(k)}};$ 
 $\mathbf{W}^{(k)} = [\boldsymbol{\theta}^{(k)}, \boldsymbol{\omega}^{(k)}];$ 
 $\mathbf{W}^{(k)} = \mathbf{W}^{(k)} - \eta \frac{\partial L_s}{\partial z^{(k)}} [1; \mathbf{X}^{(k-1)}]^T;$ 
end
end
error =  $L^2(\mathbf{X}, \boldsymbol{\omega}, \boldsymbol{\theta});$ 
mix the input vector  $\mathbf{X} = \{\mathbf{X}_k = (x_1, \dots, x_m, x_{m+1} = t)_k\}_{k=1, \dots, K};$ 
threshold = threshold + 1;
end

```

### Remark

For the choice of training set, we have two options:

1. In each epoch, the model is trained on the same training data, but the data are shuffled for each epoch. This shuffling helps ensure that the model does not simply learn the order of the data and can generalize better.
2. If it is possible to do so, at each epoch the training set is randomly generated within the definition range of our function. This process helps even more the generalization ability of the Neural Network, as it always learns from different points.

### Remark

In the calculation

$$(2.56) \quad \mathbf{W}^{(j)} = \mathbf{W}^{(j)} - \eta \frac{\partial L_s}{\partial z^{(k)}} [1; \mathbf{X}^{(j-1)}]^T$$

we added 1 to the vector  $\mathbf{X}^{(j-1)}$  because the bias has to be taken into account.

## 2.7 Application to Fisher - Kolmogorov – Petrovsky - Piskunov Equation

### 2.7.1 Introduction

Let  $x \in \Omega = [a, b]$  and  $t \in [0, T]$ , and consider the dimensionless Fisher-Kolmogorov-Petrovsky-Piskunov equation

$$y_t - y_{xx} = y(1-y),$$

For simplicity, we start by considering spatially homogeneous initial conditions of the type

$$y(x, 0) = y_0.$$

### 2.7.2 Spatially homogeneous case

Let us consider the spatially homogeneous case

$$(2.57) \quad \begin{cases} y_t = y(1-y), & \text{in } (0, T], \\ y(0) = y_0, & \text{in } t=0, \end{cases}$$

with  $y_0 = 1.5$ ,  $T = 5$ . As we have already seen in section 1.6, the solution is of the form

$$(2.58) \quad y(t) = y_0 \frac{e^t}{1 + y_0(e^t - 1)},$$

therefore, the population tends to reach the value 1, for  $t \rightarrow +\infty$ , see Figure 1.1.

In the spatially homogeneous case, the boundary conditions lose their meaning.

Now, let us consider the trial solution

$$(2.59) \quad \bar{y}(t) = y_0 + t u(t; \omega, \theta),$$

where the initial condition is automatically satisfied. In this case, the Loss Function becomes

$$(2.60) \quad L^2(\bar{y}, \bar{y}_t) = \int_0^T [\bar{y}_t - \bar{y}(1-\bar{y})]^2 dt,$$

where

$$(2.61) \quad \bar{y}_t(t) = u(t; \omega, \theta) + t u_t(t; \omega, \theta).$$

Then, we can calculate the derivative

$$(2.62) \quad \frac{\partial}{\partial z^{(k)}} \left[ [\bar{y}_t - \bar{y}(1-\bar{y})]^2 \right] dt = 2[\bar{y}_t - \bar{y}(1-\bar{y})] \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}} - (1-2\bar{y}) \frac{\partial \bar{y}}{\partial z^{(k)}} \right],$$

where

- $\frac{\partial \bar{y}_t}{\partial z^{(k)}} = \frac{\partial u}{\partial z^{(k)}} + t \frac{\partial u_t}{\partial z^{(k)}},$
- $\frac{\partial \bar{y}}{\partial z^{(k)}} = t \frac{\partial u}{\partial z^{(k)}}.$

We have already seen that

$$\frac{\partial u}{\partial z^{(k)}} = \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-1)}} \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}),$$

and that

$$\begin{aligned} \frac{\partial u_t}{\partial z^{(k)}} &= \frac{\partial}{\partial z^{(k)}} \frac{\partial u}{\partial t} = \frac{\partial}{\partial t} \frac{\partial u}{\partial z^{(k)}} \\ &= \frac{\partial}{\partial t} \left( \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-1)}} \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \right) \\ &= \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \frac{\partial \sigma'(\mathbf{z}^{(N)})}{\partial t} \frac{\partial z^{(N)}}{\partial X^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\ &\quad + \dots + \\ &+ \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \dots \frac{\partial z^{(s+I)}}{\partial X^{(s)}} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial t} \frac{\partial z^{(s)}}{\partial X^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \sigma'(\mathbf{z}^{(k)}) \\ &\quad + \dots + \\ &+ \frac{\partial z^{(N+I)}}{\partial X^{(N)}} \sigma'(\mathbf{z}^{(N)}) \frac{\partial z^{(N)}}{\partial X^{(N-1)}} \sigma'(\mathbf{z}^{(N-1)}) \dots \frac{\partial z^{(k+I)}}{\partial X^{(k)}} \frac{\partial \sigma'(\mathbf{z}^{(k)})}{\partial t}, \end{aligned}$$

where

$$\begin{aligned} \frac{\partial \sigma'(\mathbf{z}^{(s)})}{\partial t} &= \sigma''(\mathbf{z}^{(s)}) \frac{\partial z^{(s)}}{\partial t} = \sigma''(\mathbf{z}^{(s)}) \frac{\partial z^{(s)}}{\partial X^{(s-1)}} \frac{\partial X^{(s-1)}}{\partial t} \\ &= \sigma''(\mathbf{z}^{(s)}) \frac{\partial z^{(s)}}{\partial X^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \frac{\partial z^{(s-1)}}{\partial t} \\ &= \sigma''(\mathbf{z}^{(s)}) \frac{\partial z^{(s)}}{\partial X^{(s-1)}} \sigma'(\mathbf{z}^{(s-1)}) \frac{\partial z^{(s-1)}}{\partial X^{(s-2)}} \sigma'(\mathbf{z}^{(s-2)}) \dots \sigma'(\mathbf{z}^{(1)}) \frac{\partial z^{(1)}}{\partial t}. \end{aligned}$$

## Remark

It is important to observe that we have two sets of nodes:

- $t_{train}$  is the set of uniformly sampled nodes in  $[0, T]$ , which are used to train the Neural Network;

- $\mathbf{t}_{test}$  is the set of uniformly sampled nodes in  $[0, T]$ , which are used to verify, by calculating the Loss Function, that the training is working on any point in the interval  $[0, T]$ , and not only on the training set.

### 2.7.3 Matlab implementation

For the implementation in Matlab we need to discretize the domain. Let  $\mathbf{t}_{test} = \{t_i\}_{i=1}^M$  be a set of uniformly sampled nodes on  $[0, T]$ , then the Loss Function becomes

$$(2.63) \quad L^2(\bar{y}, \bar{y}_t) \approx \frac{T}{M} \sum_{i=1}^M [\bar{y}_t(t_i) - \bar{y}(t_i)(1 - \bar{y}(t_i))]^2.$$

Let us define  $I_{L^2}(t) := [\bar{y}_t(t) - \bar{y}(t)(1 - \bar{y}(t))]^2$ .

Then, for each  $t_j \in \mathbf{t}_{train}$ , we can calculate the derivative

$$\frac{\partial I_{L^2}}{\partial \mathbf{z}^{(k)}}(t_j) = 2[\bar{y}_t(t_j) - \bar{y}(t_j)(1 - \bar{y}(t_j))] \left[ \frac{\partial \bar{y}_t}{\partial \mathbf{z}^{(k)}}(t_j) - (1 - 2\bar{y}(t_j)) \frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(t_j) \right],$$

where

- $\frac{\partial \bar{y}_t}{\partial \mathbf{z}^{(k)}}(t_j) = \frac{\partial u}{\partial \mathbf{z}^{(k)}}(t_j) + t_j \frac{\partial u_t}{\partial \mathbf{z}^{(k)}}(t_j) = \alpha^{(k)}(t_j) + t_j \beta_t^{(k)}(t_j);$
- $\frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(t_j) = t_j \frac{\partial u}{\partial \mathbf{z}^{(k)}}(t_j) = t_j \alpha^{(k)}(t_j).$

### 2.7.4 Matlab codes description

To implement a Neural Network in Matlab to approximate the solution of our problem, we proceeded in the following order:

- Construction of the **function SHCNeuralNetwork**:

This function receives as input the parameters

$$(N_{layer}, omega, theta, t),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $t$  is the input node,

and calculates the output of a Neural Network. The sigmoid function was used as the activation function.

- Construction of the ***function SHCTimeDerivative***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, t),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $t$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the derivative, with respect to  $t$ , of the Neural Network. Of course, having used the sigmoid function as the activation function, it was taken into account that

$$\sigma'(x) = \sigma(x) \odot (1 - \sigma(x)).$$

- Construction of the ***function SHCLossFunction***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, t, T, y_0),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $t$  is the input node,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $y_0$  is the initial condition,

and returns as output the value of the Loss Function.

- Construction of the ***function SHCWeightsUpdate***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, t_{test}, T, y_0, lrate)$$

where

- $N_{layer}$  is the number of hidden layers,
- $\omega$  is an array containing the weights matrices,
- $\theta$  is an array containing the biases vectors,
- $t_{test}$  is the test set of input nodes,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $y_0$  is the initial condition,
- $lrate$  is the learning rate,

and updates the weights of the Neural Network through the process described above.

### Remark

Let us observe that we have two sets of points:

- $t_{test}$ , which is given as input from outside;
- $t_{train}$ , which is randomly generated at each step.

The code is structured as follows:

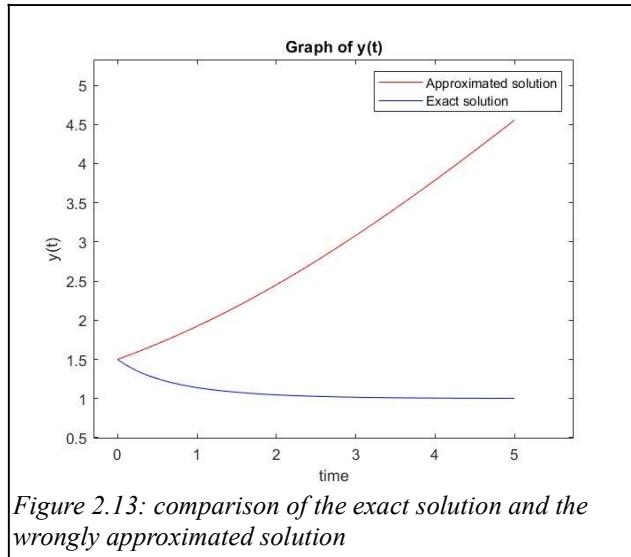
- There is a while loop that repeats until the error generated is small enough (or if a certain threshold of iterations is exceeded);
- at each step of the while loop:
  - a new training set  $t_{train}$  is generated. This is to ensure that the Neural Network always learns from different points;
  - a for loop starts on all the points in the training set. For each point,
 
$$2[\bar{y}_t(t_j) - \bar{y}(t_j)(I - \bar{y}(t_j))] \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}}(t_j) - (I - 2\bar{y}(t_j)) \frac{\partial \bar{y}}{\partial z^{(k)}}(t_j) \right],$$
 is calculated and the weights are updated;
  - at the end of the for loop a check is made on the error generated, i.e., the Loss Function  $L^2$  is calculated on the test set  $t_{test}$ .
- If the error generated is small enough, the while loop stops and outputs the value of the Loss Function and the new values of omega and theta.

Once the structure for the operation of the learning process has been built

- *function SHCNeuralNetwork,*
- *function SHCTimeDerivative,*
- *function SHCLossFunction,*
- *function SHCWeightsUpdate,*

we tried to apply it in the script *SHCtest*.

At first, having no particular criteria for choosing the parameters and weights, we have chosen them totally randomly. Unfortunately, by comparing the graph of the exact solution and the graph of the solution approximated by the Neural Network with each other, we noticed a big difference between them, see Figure 2.13. This discrepancy, is due to the fact that the output of the Neural Network is strongly influenced by the choice of the parameters and weights. To solve this drawback, we built



- ***SHCParametersOptimization:***

This script, runs for loops by varying the number of hidden layers, the dimension of the hidden layers, and the learning rate, and outputs back the "best" parameters so that, after applying the *function SHCWeightsUpdate*, the error produced by the Loss Function is as low as possible. Using this script, it was observed that a good choice for parameters is as follows:

PARAMETERS	
Number of hidden layers (N_layer)	1
Dimension of any hidden layer (dimension)	3
Learning rate (lrate)	0.001

- *script SHCWeightsOptimization:*

This script uses the parameters found through the script **SHCParameterOptimization** and generates new matrices of weights and biases so that, once the function **SHCWeightsUpdate** is applied, the error generated by the Loss Function is as low as possible.

As it was done in [12], to generate the weights and bias matrices, we used the

*Xavier initialization method (or Glorot initialization)*

It is a technique used to initialize Neural Network weights in order to improve model convergence during training. This technique is particularly useful and it is used to reduce the problems of vanishing/exploding gradients. The main idea is to keep the variance of activation values constant across the various layers of the network, preventing the outputs from becoming too small or too large. This is achieved by choosing the initial weights through a distribution with a specified variance.

For a layer with  $m$  input nodes and  $n$  output nodes, the weights

$W = \{\omega, \theta\}$  are initialized using:

- uniform distribution, i.e.

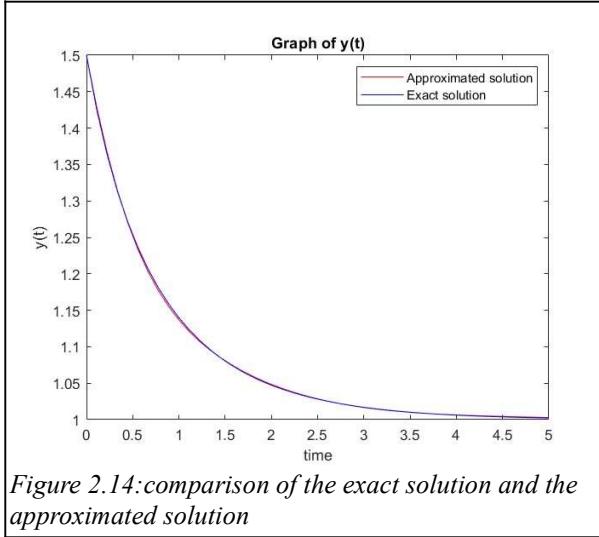
$$(2.64) \quad W \sim U\left(-\sqrt{\frac{6}{m+n}}, +\sqrt{\frac{6}{m+n}}\right);$$

- normal (Gaussian) distribution, i.e.

$$(2.65) \quad W \sim N\left(0, +\sqrt{\frac{2}{m+n}}\right).$$

Once **SHCParametersOptimization** and **SHCWeightsOptimization** were implemented, we have implemented the **SHCTest2** script using those obtained from **SHCParametersOptimization** as parameters and those obtained from **SHCWeightsOptimization** as weights and bias. We thus obtained a good

approximation of the solution, as we can see in Figure 2.14. In the interval  $[0, T]$  we have fixed  $n_{test} = 50$  nodes for error calculation.



By using the Loss Function, it was calculated that

$$Error = L^2(\bar{y}, \bar{y}_t) \approx \frac{T}{n_{test}} \sum_{i=1}^{n_{test}} [\bar{y}_t(t_i) - \bar{y}(t_i)(1 - \bar{y}(t_i))]^2 = 7.5932 \times 10^{-4},$$

then the solution found “minimizes” the Loss Function.

In addition, however, we wanted to calculate the average distance between the approximate and the real solution, obtaining

$$QuadraticError = \frac{1}{n_{test}} \sum_{i=1}^{n_{test}} |s(t_i) - u(t_i)|^2 = 0.00022.$$

For more information about the Matlab codes, see Appendix B.

### 2.7.5 General case with spatially homogeneous initial conditions

Let us now consider the dimensionless Fisher-Kolmogorov-Petrovsky-Piskunov problem with boundary conditions

$$(2.66) \quad \begin{cases} y_t - y_{xx} = y(1-y), & \text{in } (a, b) \times (0, T], \\ y(x, 0) = y_0, & \text{in } [a, b] \times \{t=0\}, \\ y(a, t) = y_a(t), y(b, t) = y_b(t), & t \in [0, T], \end{cases}$$

where, for example, we set  $a = -3$ ,  $b = 3$ ,  $T = 5$ .

We observe that  $y_a(0) = y_b(0) = y_0$ , so we choose, as an initial test, for simplicity, the boundary conditions

$$y_a(t) = y_b(t) = y_0 \frac{e^t}{1 + y_0(e^t - 1)},$$

which represents the exact solution in the spatially homogeneous case.

Now, let us consider the trial solution

$$(2.67) \quad \bar{y}(x,t) = y_a(t) \frac{b-x}{b-a} + y_b(t) \frac{x-a}{b-a} + t \frac{(b-x)(x-a)}{(b-a)^2} u(x,t; \omega, \theta),$$

and we observe that the initial condition and the boundary conditions are automatically satisfied, in fact:

- $\bar{y}(x,0) = y_a(0) \frac{b-x}{b-a} + y_b(0) \frac{x-a}{b-a} = y_0 \left[ \frac{b-x}{b-a} + \frac{x-a}{b-a} \right] = y_0;$
- $\bar{y}(a,t) = y_a(t);$
- $\bar{y}(b,t) = y_b(t).$

In this case, the Loss Function becomes

$$(2.68) \quad L^2(\bar{y}, \bar{y}_t, \bar{y}_{xx}) = \int_0^T dt \int_{\Omega} [\bar{y}_t - \bar{y}_{xx} - \bar{y}(I - \bar{y})]^2 d\mathbf{x},$$

where

- $\bar{y}_t = y_a'(t) \frac{b-x}{b-a} + y_b'(t) \frac{x-a}{b-a} + \frac{(b-x)(x-a)}{(b-a)^2} u + t \frac{(b-x)(x-a)}{(b-a)^2} u_t;$
- $\begin{aligned} \bar{y}_x &= \frac{-y_a(t)}{b-a} + \frac{y_b(t)}{b-a} + t \frac{-(x-a)+(b-x)}{(b-a)^2} u + t \frac{(b-x)(x-a)}{(b-a)^2} u_x \\ &= \frac{y_b(t) - y_a(t)}{b-a} + t \frac{-2x+b+a}{(b-a)^2} u + t \frac{(b-x)(x-a)}{(b-a)^2} u_x; \end{aligned}$
- $\begin{aligned} \bar{y}_{xx} &= -t \frac{2}{(b-a)^2} u + t \frac{-2x+b+a}{(b-a)^2} u_x + t \frac{-2x+b+a}{(b-a)^2} u_x + t \frac{(b-x)(x-a)}{(b-a)^2} u_{xx} \\ &= -t \frac{2}{(b-a)^2} u + 2t \frac{-2x+b+a}{(b-a)^2} u_x + t \frac{(b-x)(x-a)}{(b-a)^2} u_{xx}. \end{aligned}$

Let us define  $I_{L^2}(x,t) := [\bar{y}_t(x,t) - \bar{y}(x,t)(I - \bar{y}(x,t))]^2$ . Then, we can calculate the derivative

$$\frac{\partial I_{L^2}}{\partial z^{(k)}} = \frac{\partial}{\partial z^{(k)}} \left[ [\bar{y}_t - \bar{y}_{xx} - \bar{y}(I - \bar{y})]^2 \right] = 2[\bar{y}_t - \bar{y}_{xx} - \bar{y}(I - \bar{y})] \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}} - \frac{\partial \bar{y}_{xx}}{\partial z^{(k)}} - (I - 2\bar{y}) \frac{\partial \bar{y}}{\partial z^{(k)}} \right],$$

where

- $\frac{\partial \bar{y}}{\partial z^{(k)}} = t \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u}{\partial z^{(k)}},$
- $\frac{\partial \bar{y}_t}{\partial z^{(k)}} = \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u}{\partial z^{(k)}} + t \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u_t}{\partial z^{(k)}},$

- $\frac{\partial \bar{y}_{xx}}{\partial z^{(k)}} = -t \frac{2}{(b-a)^2} \frac{\partial u}{\partial z^{(k)}} + 2t \frac{-2x+b+a}{(b-a)^2} \frac{\partial u_x}{\partial z^{(k)}} + t \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u_{xx}}{\partial z^{(k)}}.$

## Reminder

Remember that  $\frac{\partial u}{\partial z^{(k)}}, \frac{\partial u_t}{\partial z^{(k)}}, \frac{\partial u_{xx}}{\partial z^{(k)}}$ , were calculated using the formulas in section 2.6.3.

### 2.7.6 Matlab implementation

For the implementation in Matlab we need to discretize the domain. Let  $t = \{t_i\}_{i=1}^M$  be the set of uniformly sampled points on  $[0, T]$ , and let  $x = \{x_j\}_{j=1}^N$  be the set of uniformly sampled points on  $[a, b]$ . Then the Loss Function becomes

$$(2.69) \quad L^2(\bar{y}, \bar{y}_t, \bar{y}_{xx}) \approx \frac{T(b-a)}{MN} \sum_{i=1}^M \sum_{j=1}^N [\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(1 - \bar{y}(x_j, t_i))]^2.$$

Then, for each  $(x_j, t_i)$  in the training set, we can calculate the derivative

$$\begin{aligned} \frac{\partial I_{L^2}}{\partial z^{(k)}}(x_j, t_i) &= 2[\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(1 - \bar{y}(x_j, t_i))] \\ &\quad \times \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}}(x_j, t_i) - \frac{\partial \bar{y}_{xx}}{\partial z^{(k)}}(x_j, t_i) - (1 - 2\bar{y}(x_j, t_i)) \frac{\partial \bar{y}}{\partial z^{(k)}}(x_j, t_i) \right], \end{aligned}$$

where

- $\frac{\partial \bar{y}}{\partial z^{(k)}}(x_j, t_i) = t_i \frac{(b-x_j)(x_j-a)}{(b-a)^2} \alpha^{(k)}(x_j, t_i);$
- $\frac{\partial \bar{y}_t}{\partial z^{(k)}}(x_j, t_i) = \frac{(b-x_j)(x_j-a)}{(b-a)^2} \alpha^{(k)}(x_j, t_i) + t_i \frac{(b-x_j)(x_j-a)}{(b-a)^2} \beta_t^{(k)}(x_j, t_i);$
- $\frac{\partial \bar{y}_{xx}}{\partial z^{(k)}}(x_j, t_i) = -t_i \frac{2}{(b-a)^2} \alpha^{(k)}(x_j, t_i) + 2t_i \frac{-2x_j+b+a}{(b-a)^2} \beta_x^{(k)}(x_j, t_i) + t_i \frac{(b-x_j)(x_j-a)}{(b-a)^2} \gamma^{(k)}(x_j, t_i).$

### 2.7.7 Matlab codes description

To implement a Neural Network in Matlab to approximate the solution of our problem, we proceeded in the following order:

- Construction of the ***function NeuralNetwork***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculate the output of a Neural Network. The sigmoid function was used as the activation function.

- Construction of the ***function TimeDerivative***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the derivative, with respect to  $t$ , of the Neural Network. Of course, having used the sigmoid function as the activation function, it was taken into account that

$$\sigma'(x) = \sigma(x) \odot (1 - \sigma(x)).$$

- Construction of the ***function SpaceDerivative***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,

- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the derivative, with respect to  $x$ , of the Neural Network.

- Construction of the ***function Laplacian***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the laplacian of the Neural Networks. Of course, having used the sigmoid function as the activation function, it was taken into account that

$$\sigma''(z^{(k)}) = \sigma(z^{(k)}) \odot (1 - \sigma(z^{(k)})) \odot (1 - 2\sigma(z^{(k)})).$$

- Construction of the ***function LossFunction***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_{test}, T, a, b, ya, yb, dya, dyb),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_{test}$  is the test set,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $a$  is the left boundary of the interval  $[a, b]$ ,
- $b$  is the right boundary of the interval  $[a, b]$ ,
- $ya$  and  $yb$  are the boundary conditions,
- $dya$  and  $dyb$  are the derivatives, respect to  $t$ , of  $ya$  and  $yb$ ,

and returns as output the value of the Loss Function, evaluated on the test set.

- Construction of the ***function WeightsUpdate***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_{test}, ya, yb, dya, dyb, a, b, T, lrate),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_{test}$  are uniformly distributed grid point on  $[a, b] \times [0, T]$ , which are used to test the error of the Neural Network,
- $ya$  and  $yb$  are the boundary conditions,
- $dya$  and  $dyb$  are the derivatives respect to  $t$ , of  $ya$  and  $yb$ ,
- $a$  is the left boundary of the interval  $[a, b]$ ,
- $b$  is the right boundary of the interval  $[a, b]$ ,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $lrate$  represents the learning rate

and updates the weights of the Neural Network using the Gradient Method.

## Remark

Let us observe that we have two sets of points:

- $X_{test}$ , which is given as input from outside;
- $X_{train}$ , which is randomly generated at each step.

The code is structured as follows:

- There is a while loop that repeats until the error generated is small enough (or if a certain threshold of iterations is exceeded);
- at each step of the while loop:
  - a new training set  $X_{train}$  is generated. This is to ensure that the Neural Network always learns from different points;
  - a for loop starts on all the points in the training set. For each point,

$$2[\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(I - \bar{y}(x_j, t_i))] \\ \times \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}}(x_j, t_i) - \frac{\partial \bar{y}_{xx}}{\partial z^{(k)}}(x_j, t_i) - (I - 2\bar{y}(x_j, t_i)) \frac{\partial \bar{y}}{\partial z^{(k)}}(x_j, t_i) \right],$$

is calculated and the weights are updated;

- at the end of the for loop a check is made on the error generated, i.e., the Loss Function  $L^2$  is calculated on the test set  $X_{test}$ .
- If the error generated is small enough, the while loop stops and outputs the value of the Loss Function and the new values of omega and theta.

Once the structure for the operation of the learning process has been built

- ***function NeuralNetwork,***
- ***function TimeDerivative,***
- ***function SpaceDerivative,***
- ***function Laplacian,***
- ***function LossFunction,***
- ***function WeightsUpdate,***

we tried to apply it in the script ***test***. After a few attempts, it was observed that a good choice for parameters is as follows:

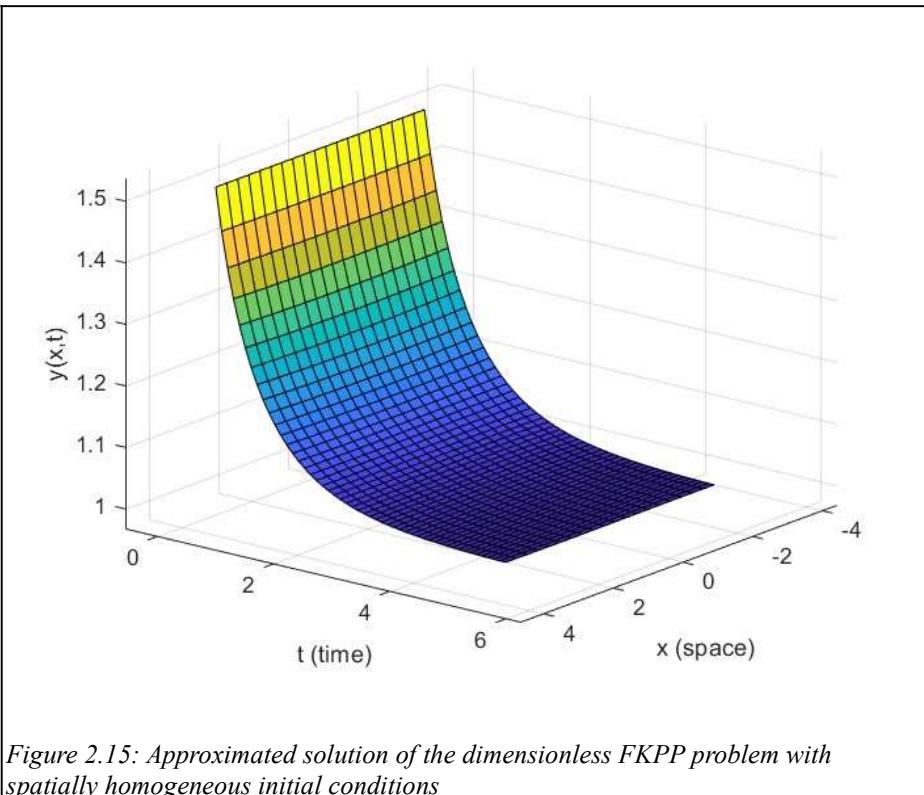
<b>PARAMETERS</b>	
Number of hidden layers (N_layer)	1
Dimension of any hidden layer (dimension)	2
Learning rate (lrate)	0.001

To generate the weights and bias matrices, we have used, as previously, the *Xavier initialization method (or Glorot initialization)*. We thus obtained a good approximation of the solution, as we can see in Figure 2.15. For the calculation of the error, we have fixed  $t_{test} = 50$  time nodes in the interval  $[0, T]$  and  $x_{test} = 20$  space nodes in the interval  $[a, b]$ .

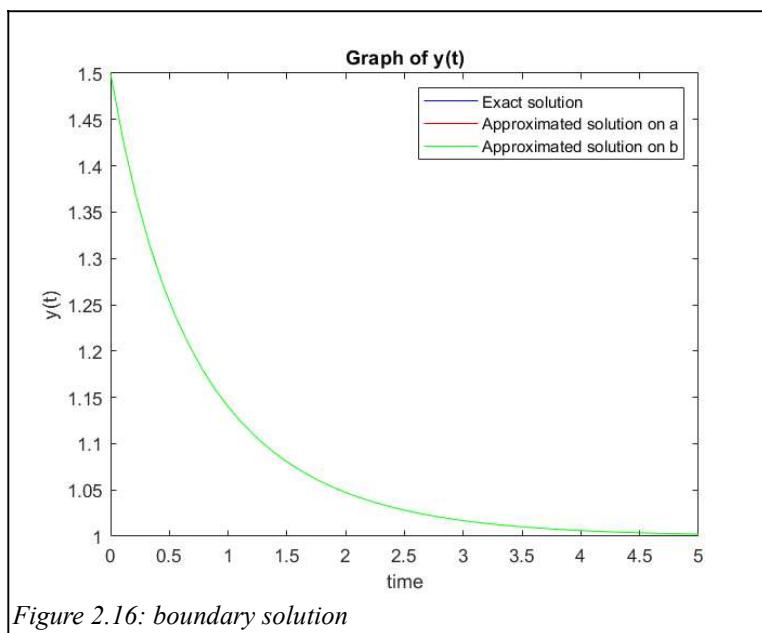
By using the Loss Function, it was calculated that

$$L^2 \approx \frac{T(b-a)}{t_{test} \times x_{test}} \sum_{i=1}^{t_{test}} \sum_{j=1}^{x_{test}} [\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(I - \bar{y}(x_j, t_i))]^2 = 3.5792 \times 10^{-4},$$

then the solution found “minimizes” the Loss Function.



In addition, we can observe that the boundary conditions turn out to be actually verified, in fact, as we can see in Figure 2.16, the graphs are "completely" overlaid.



For more information about Matlab codes, see Appendix C.

### 2.7.8 Case with non-homogeneous initial conditions

Let us now consider the dimensionless Fisher-Kolmogorov-Petrovsky-Piskunov problem with boundary conditions

$$\begin{cases} y_t - y_{xx} = y(1-y), \text{ in } (a, b) \times (0, T], \\ y(x, 0) = y_0(x), \text{ in } [a, b] \times \{t=0\}, \\ y(a, t) = y_a(t), y(b, t) = y_b(t), t \in [0, T], \end{cases}$$

where  $y_0(x) \in C([a, b])$  and  $y_a(t), y_b(t) \in C([0, T])$ . Let us observe that, in order to have continuity on the boundary, one must have

$$\begin{cases} y_a(0) = y_0(a), \\ y_b(0) = y_0(b). \end{cases}$$

Now, let us consider the trial solution

$$\bar{y}(x, t) = \frac{b-x}{b-a} y_a(t) + \frac{x-a}{b-a} y_b(t) + \frac{(b-x)(x-a)}{(b-a)^2} u(x, t; \omega, \theta),$$

and we observe that the boundary condition are automatically satisfied, in fact:

- $\bar{y}(a, t) = y_a(t);$
- $\bar{y}(b, t) = y_b(t).$

The initial condition and the reaction-diffusion equation remain to be satisfied. In this case the Loss Function becomes

$$L^2(\bar{y}, \bar{y}_t, \bar{y}_{xx}) = \int_0^T dt \int_a^b [\bar{y}_t - \bar{y}_{xx} - \bar{y}(1-\bar{y})]^2 dx + \int_a^b [\bar{y}(x, 0) - y_0(x)]^2 dx,$$

where

- $\bar{y}_t = \frac{b-x}{b-a} y_a'(t) + \frac{x-a}{b-a} y_b'(t) + \frac{(b-x)(x-a)}{(b-a)^2} u_t;$
- $\bar{y}_x = \frac{-y_a(t)}{b-a} + \frac{y_b(t)}{b-a} + \frac{-2x+b+a}{(b-a)^2} u + \frac{(b-x)(x-a)}{(b-a)^2} u_x;$
- $\bar{y}_{xx} = \frac{-2}{(b-a)^2} u + 2 \frac{-2x+b+a}{(b-a)^2} u_x + \frac{(b-x)(x-a)}{(b-a)^2} u_{xx}.$

Let us observe it is required that  $y_a(t), y_b(t) \in C^1([0, T]).$

Now, we can calculate the two derivatives

$$\frac{\partial}{\partial z^{(k)}} \left[ [\bar{y}_t - \bar{y}_{xx} - \bar{y}(1-\bar{y})]^2 \right] = 2[\bar{y}_t - \bar{y}_{xx} - \bar{y}(1-\bar{y})] \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}} - \frac{\partial \bar{y}_{xx}}{\partial z^{(k)}} - (1-2\bar{y}) \frac{\partial \bar{y}}{\partial z^{(k)}} \right],$$

$$\frac{\partial}{\partial \mathbf{z}^{(k)}} [(\bar{y}(x, 0) - y_\theta(x))^2] = 2[\bar{y}(x, 0) - y_\theta(x)] \frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(x, 0),$$

where

- $\frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}} = \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u}{\partial \mathbf{z}^{(k)}},$
- $\frac{\partial \bar{y}_t}{\partial \mathbf{z}^{(k)}} = \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u}{\partial \mathbf{z}^{(k)}} + \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u_t}{\partial \mathbf{z}^{(k)}},$
- $\frac{\partial \bar{y}_{xx}}{\partial \mathbf{z}^{(k)}} = -\frac{2}{(b-a)^2} \frac{\partial u}{\partial \mathbf{z}^{(k)}} + 2 \frac{-2x+b+a}{(b-a)^2} \frac{\partial u_x}{\partial \mathbf{z}^{(k)}} + \frac{(b-x)(x-a)}{(b-a)^2} \frac{\partial u_{xx}}{\partial \mathbf{z}^{(k)}}.$

## Reminder

Remember that  $\frac{\partial u}{\partial \mathbf{z}^{(k)}}, \frac{\partial u_t}{\partial \mathbf{z}^{(k)}}, \frac{\partial u_{xx}}{\partial \mathbf{z}^{(k)}}$ , were calculated using the formulas in section 2.6.3.

### 2.7.9 Matlab implementation

For the implementation in Matlab we need to discretize the domain. Let  $\mathbf{t} = \{t_i\}_{i=1}^M$  be the set of uniformly sampled points on  $[0, T]$ , and let  $\mathbf{x} = \{x_j\}_{j=1}^N$  be the set of uniformly sampled points on  $[a, b]$ . Then the Loss Function becomes

$$(2.69) \quad L^2(\bar{y}, \bar{y}_t, \bar{y}_{xx}) \approx \frac{T(b-a)}{MN} \sum_{i=1}^M \sum_{j=1}^N [\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(I - \bar{y}(x_j, t_i))]^2 + \frac{(b-a)}{N} \sum_{j=1}^N [\bar{y}(x_j, 0) - y_\theta(x_j)]^2,$$

Then, for each  $(x_j, t_i)$  in the training set, we can calculate the two derivatives

$$2[\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(I - \bar{y}(x_j, t_i))] \times \left[ \frac{\partial \bar{y}_t}{\partial \mathbf{z}^{(k)}}(x_j, t_i) - \frac{\partial \bar{y}_{xx}}{\partial \mathbf{z}^{(k)}}(x_j, t_i) - (I - 2\bar{y}(x_j, t_i)) \frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(x_j, t_i) \right], \\ 2[\bar{y}(x_j, 0) - y_\theta(x_j)] \frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(x_j, 0),$$

where

- $\frac{\partial \bar{y}}{\partial \mathbf{z}^{(k)}}(x_j, t_i) = \frac{(b-x_j)(x_j-a)}{(b-a)^2} \alpha^{(k)}(x_j, t_i);$
- $\frac{\partial \bar{y}_t}{\partial \mathbf{z}^{(k)}}(x_j, t_i) = \frac{(b-x_j)(x_j-a)}{(b-a)^2} \alpha^{(k)}(x_j, t_i) + \frac{(b-x_j)(x_j-a)}{(b-a)^2} \beta_t^{(k)}(x_j, t_i);$

- $$\begin{aligned} \frac{\partial \bar{y}_{xx}(x_j, t_i)}{\partial z^{(k)}} &= -\frac{2}{(b-a)^2} \alpha^{(k)}(x_j, t_i) + 2 \frac{-2x_j+b+a}{(b-a)^2} \beta_x^{(k)}(x_j, t_i) \\ &+ \frac{(b-x_j)(x_j-a)}{(b-a)^2} \gamma^{(k)}(x_j, t_i). \end{aligned}$$

For the implementation, we will set

- $[a, b] = [-3, 3];$
- $[0, T] = [0, 5];$
- $y_0(x) = \cos\left(\frac{\pi}{3}x\right) + 2.5;$
- $y_a(t) = y_b(t) = y_0(a) \frac{e^t}{1 + y_0(a)(e^t - 1)}.$

### 2.7.10 Matlab codes description

To implement a Neural Network in Matlab to approximate the solution of our problem, we proceeded in the following order:

- Construction of the ***function GCNeuralNetwork***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculate the output of the Neural Network. The sigmoid function was used as the activation function.

- Construction of the ***function GCTimeDerivative***:

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,

- $\theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the derivative, with respect to  $t$ , of the Neural Network. Of course, having used the sigmoid function as the activation function, it was taken into account that

$$\sigma'(x) = \sigma(x) \odot (1 - \sigma(x)).$$

- Construction of the ***function GCSpaceDerivative***:

This function receives as input the parameters

$$(N_{layer}, \omega, \theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $\omega$  is an array containing the weights matrices,
- $\theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the derivative, with respect to  $x$ , of the Neural Network.

- Construction of the ***function GCLaplacian***:

This function receives as input the parameters

$$(N_{layer}, \omega, \theta, X_0)$$

where

- $N_{layer}$  is the number of hidden layers,
- $\omega$  is an array containing the weights matrices,
- $\theta$  is an array containing the biases vectors,
- $X_0$  is the input node,

and calculates, through the automatic differentiation process described in Section 2.5.2, the functional evaluation of the laplacian of the Neural Networks. Of course, having used the sigmoid function as the activation function, it was taken into account that

$$\sigma''(\mathbf{z}^{(k)}) = \sigma(\mathbf{z}^{(k)}) \odot (1 - \sigma(\mathbf{z}^{(k)})) \odot (1 - 2\sigma(\mathbf{z}^{(k)})).$$

- Construction of the ***function GCLossFunction:***

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_{test}, X_{test0}, T, a, b, y_0, ya, dy),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_{test}$  is the test set on  $[a, b] \times [0, T]$ ,
- $X_{test0}$  is the test set on  $[a, b] \times \{t=0\}$ ,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $a$  is the left boundary of the interval  $[a, b]$ ,
- $b$  is the right boundary of the interval  $[a, b]$ ,
- $y_0$  is the initial condition,
- $ya$  is the boundary condition (remember that  $ya = yb$ ),
- $dy$  is the derivative, respect to  $t$ , of  $ya$  (remember that  $dy = dyb$ ),

and returns as output the value of the Loss Function, evaluated at the test sets  $X_{test}$  and  $X_{test0}$ .

- Construction of the ***function GCWeightsUpdate:***

This function receives as input the parameters

$$(N_{layer}, omega, theta, X_{test}, X_{test0}, y_0, ya, dy, a, b, T, lrate),$$

where

- $N_{layer}$  is the number of hidden layers,
- $omega$  is an array containing the weights matrices,
- $theta$  is an array containing the biases vectors,
- $X_{test}$  is the test set on  $[a, b] \times [0, T]$ ,
- $X_{test0}$  is the test set on  $[a, b] \times \{t=0\}$ ,
- $T$  is the right boundary of the interval  $[0, T]$ ,
- $a$  is the left boundary of the interval  $[a, b]$ ,
- $b$  is the right boundary of the interval  $[a, b]$ ,

- $y_0$  is the initial condition,
- $ya$  is the boundary condition (remembere that  $ya = yb$ ),
- $dya$  is the derivative, respect to  $t$ , of  $ya$  (remembere that  $dya = dyb$ ),
- $lrate$  represents the learning rate

and updates the weights of the Neural Network using the Gradient Method.

### Remark

Let us observe that we have four sets of points:

- $X_{test}, X_{test0}$  which are given as input from outside;
- $X_{train}, X_{train0}$  which are randomly generated at each step, and represent, respectively, the training set on  $[a, b] \times [0, T]$ , and the training set on  $[a, b] \times \{t=0\}$ .

The code is structured as follows:

- There is a while loop that repeats until the error generated is small enough (or if a certain threshold of iterations is exceeded);
- at each step of the while loop:
  - two new training sets  $X_{train}, X_{train0}$  are generated. This is to ensure that the Neural Network always learns from different points;
  - a for loop starts on all the points in the training sets. For each point are calculated

$$\begin{aligned} \frac{\partial L_1}{\partial z^{(k)}}(x_j, t_i) &= 2[\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(1 - \bar{y}(x_j, t_i))] \\ &\quad \times \left[ \frac{\partial \bar{y}_t}{\partial z^{(k)}}(x_j, t_i) - \frac{\partial \bar{y}_{xx}}{\partial z^{(k)}}(x_j, t_i) - (1 - 2\bar{y}(x_j, t_i)) \frac{\partial \bar{y}}{\partial z^{(k)}}(x_j, t_i) \right], \end{aligned}$$

$$\frac{\partial L^2}{\partial z^{(k)}}(x_j, 0) = 2[\bar{y}(x_j, 0) - y_0(x_j)] \frac{\partial \bar{y}}{\partial z^{(k)}}(x_j, 0),$$

depending on whether  $t \neq 0$  or  $t = 0$ , and the weights are updated;

- at the end of the for loop a check is made on the error generated, i.e., the Loss Function  $L^2$  is calculated on the test sets  $X_{test}, X_{test0}$ .
- If the error generated is small enough, the while loop stops and outputs the value of the Loss Function and the new values of omega and theta.

Once the structure for the operation of the learning process has been built

- ***function GCNeuralNetwork,***
- ***function GCTimeDerivative,***
- ***function GCSpaceDerivative,***
- ***function GCLaplacian,***
- ***function GCLossFunction,***
- ***function GCWeightsUpdate,***

we tried to apply it in the script ***test***.

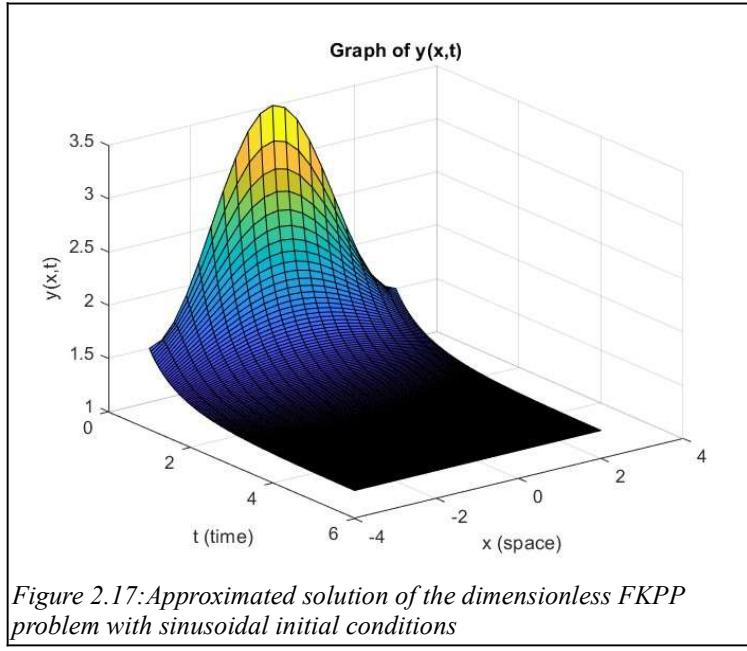
After a few attempts, it was observed that a good choice for parameters is as follows:

<b>PARAMETERS</b>	
Number of hidden layers (N_layer)	2
Dimension of any hidden layer (dimension)	19
Learning rate (lrate)	0.01

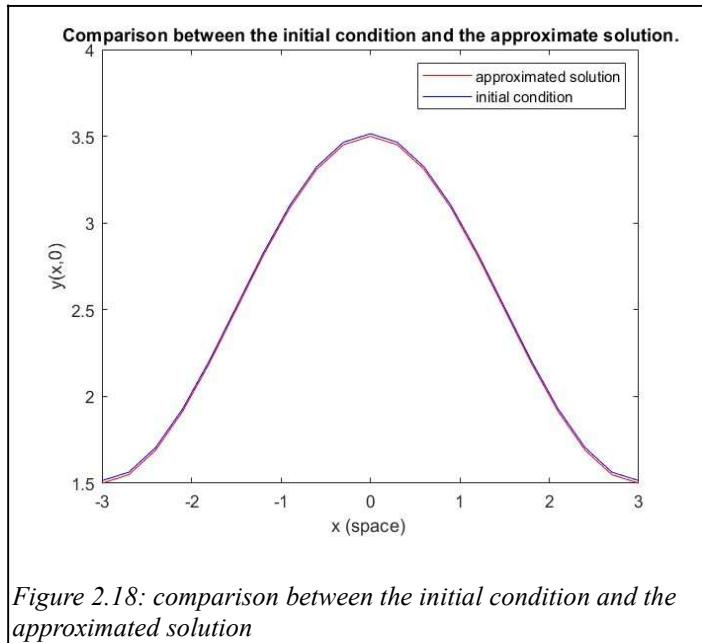
To generate the weights and bias matrices, we have used, as previously, the *Xavier initialization method (or Glorot initialization)*. We thus obtained a good approximation of the solution, as we can see in Figure 2.17. For the calculation of the error, we have fixed  $t_{test}=50$  time nodes in the interval  $[0, T]$  and  $x_{test}=x_{test,0}=20$  space nodes in the interval  $[a, b]$ . By using the Loss Function, it was calculated that

$$L^2(\bar{y}, \bar{y}_t, \bar{y}_{xx}) \approx \frac{T(b-a)}{MN} \sum_{i=1}^M \sum_{j=1}^N [\bar{y}_t(x_j, t_i) - \bar{y}_{xx}(x_j, t_i) - \bar{y}(x_j, t_i)(1 - \bar{y}(x_j, t_i))]^2 + \frac{(b-a)}{N} \sum_{j=1}^N [\bar{y}(x_j, 0) - y_0(x_j)]^2 = 6.4376 \times 10^{-3},$$

then the solution found “minimizes” the Loss Function.



In addition, we can observe that the initial condition is also verified as we can see in Figure 2.18.



For more information about Matlab codes, see Appendix D.

### Remark

In this case, we did not obtain satisfactory results from the point of view of efficiency. In fact, after numerous unsuccessful attempts, calculating the correct weights and biases took several hours, despite using a high-end PC with an Intel(R) Core(TM) i7-8750H CPU @ 2.20GHz 2.20 GHz and a 16 GB RAM.

# Conclusions

We can say that the use of Neural Networks is undoubtedly a very powerful tool for the study of solutions of partial derivative equations and, in particular, reaction-diffusion equations. However, there are significant advantages and disadvantages to consider.

Advantages:

1. model-free learning: Neural Networks, if well trained, allow us to approximate the solution of a PDE without explicit knowledge of the solution;
2. generalization ability: Neural Networks, if well-trained, allow us to find a good approximation of the solution at any point in the domain of definition, despite having trained the Neural Network only on a limited number of nodes;
3. computational speed: once trained, Neural Networks provide us with solutions very quickly;
4. high dimension: Neural Networks allow us to handle PDEs in large spaces, whereas other numerical methods, in this case, may suffer;
5. novelty: since it is still a relatively little-used tool in this area, its potential and functionality is not known exactly. Nothing prohibits us from thinking that methods can be found to further improve its effectiveness and efficiency.

Disadvantages:

1. need for training data: Neural Networks, in order to be trained, require large amounts of data. In addition, obtaining qualitatively good data can be computationally expensive;
2. training time: training Neural Networks can take a long time, especially in the case of very deep and complex Networks;
3. optimization: Neural Networks are highly dependent on the initial choice of parameters (weights and bias). This leads to the implementation of

- additional methods to find the best initial parameters so that the Neural Network can be trained properly;
4. overfitting: Neural Networks easily go into overfitting, i.e., they overfit the training data, as a result they have poor generalization ability, especially if small training sets are used or have noise;
  5. instability: the resulting solutions can sometimes be unstable and inaccurate;
  6. integration: integrating Neural Networks with other existing methods can be particularly complex and requires a lot of experience;
  7. novelty: being a relatively new tool, there are not yet many studies that allow us to understand how and why certain wrong solutions are produced.

It still remains open the problem of how to improve the algorithm to update the weights in the case where the Loss Function is the sum of several subfunctions that converge at the point of minimum with different velocities. In fact, in the case with non-homogeneous initial conditions, we have not obtained satisfactory results from the point of view of efficiency (the algorithm really takes too long to compute the correct weights, always if it finds them, because it may “explode”). We set as a future goal to continue the research in this field, trying to obtain further results, not only from a purely theoretical point of view, but also, and especially, from a practical point of view. In fact, we will focus on extending what has been done for reaction-diffusion equations, also for reaction-diffusion systems, such as, for example, the prey-predator model (the idea would be to build and train Neural Networks with two, or more, output nodes).

## APPENDIX A - APPLICATION TO THE APPROXIMATION THEOREM

### APPENDIX A - Function NeuralNetwork (1 of 1)

```
function X = NeuralNetwork (N_layer, omega, theta, x)

% sigmoid function
sigmoid = @(u) 1./(1+exp(-u));

% column vector array of node values
X=cell(1,N_layer+1);
X{1}=sigmoid(omega{1}*x+theta{1});

for i=2:N_layer
    % value of nodes in layer i
    X{i}=sigmoid(omega{i}*X{i-1}+theta{i});
end
% activation function is not applied to the last layer
X{N_layer+1}=omega{N_layer+1}*X{N_layer}+theta{N_layer+1};

end
```

## APPENDIX A - Function SpaceDerivative (1 of 1)

```
function dx = SpaceDerivative(N_layer, omega, theta, x)

% calculation of node values
X=NeuralNetwork(N_layer, omega, theta, x);

% calculation of the time derivative
dx=cell(1,N_layer+1);
dx{1}=X{1}.* (1-X{1}).*omega{1};

for i=2:N_layer
    dx{i}=omega{i}*dx{i-1};
    dx{i}=X{i}.* (1-X{i}).*dx{i};
end

dx{N_layer+1}=omega{N_layer+1}*dx{N_layer};

end
```

## APPENDIX A - Function LossFunction (1 of 1)

```
function L = LossFunction(N_layer,X_test,omega,theta,f,a,b)

%length of the test set
n=length(X_test);

% calculation of u
u=zeros(1,n);
for i=1:n
    x=X_test(i);
    X=NeuralNetwork(N_layer, omega, theta, x);
    u(i)=X{N_layer+1};
end

% calculation of the Loss Function
L=0;
for i=1:n
    x=X_test(i);
    L=L+(u(i)-f(x))^2;
end
L=L*(b-a)/n;

end
```

## APPENDIX A - Function WeightsUpdate (1 of 2)

```
function [error,thetanew,omeganew] = WeightsUpdate(N_layer,omega,theta, ...
X_test,a,b,f,lrate)

% parameters initialization
error=10;
errormax=0.0005;
threshold=1;
thresholdmax=20000;
W=cell(1,N_layer+1);
thetanew=theta;
omeganew=omega;

% weights update
while error>errormax && threshold<thresholdmax
n_train=20;
X_train=a+(b-a)*rand(1,n_train);

% calculation of u
u=zeros(1,n_train);
for i=1:n_train
    x=X_train(i);
    X=NeuralNetwork(N_layer, omega, theta, x);
    u(i)=X{N_layer+1};
end

% initialization of alfa
alfa=cell(n_train,N_layer+1);

% weights update
for i=1:n_train
    % weights update of the layer N+1
    x=X_train(i);
    X=NeuralNetwork(N_layer, omega, theta, x);
    alfa{i,N_layer+1}=1;
    dL=2*(u(i)-f(x))*[1;X{N_layer}]';
    W{N_layer+1}=[theta{N_layer+1},omega{N_layer+1}];
    W{N_layer+1}=W{N_layer+1}-lrate*dL;
    theta{N_layer+1}=W{N_layer+1}(:,1);
    omega{N_layer+1}=W{N_layer+1}(:,2:end);

    for k=N_layer:-1:2
        % weights update of the layer k
        x=X_train(i);
        X=NeuralNetwork(N_layer, omega, theta, x);
        alfa{i,k}=alfa{i,k+1}.*(omega{k+1}*(X{k}.*(1-X{k})));
        dL=2*(u(i)-f(x))*alfa{i,k}*[1;X{k-1}]';
        W{k}=[theta{k},omega{k}];
        W{k}=W{k}-lrate*dL;
        theta{k}=W{k}(:,1);
        omega{k}=W{k}(:,2:end);
    end
end

% weights update of the layer 1
```

## APPENDIX A - Function WeightsUpdate (2 of 2)

```
x=X_train(i);
X=NeuralNetwork(N_layer,omega,theta,x);
alfa{i,1}=alfa{i,2}.*(omega{2}*(X{1}.*(1-X{1})));
dL=2*(u(i)-f(x))*alfa{i,1}* [1;x]';
W{1}=[theta{1},omega{1}];
W{1}=W{1}-lrate*dL;
theta{1}=W{1}(:,1);
omega{1}=W{1}(:,2:end);

end

% error calculation
errornew=LossFunction(N_layer,X_test,omega,theta,f,a,b);

% error and weights update
if errornew<error
    error=errornew;
    omeganew=omega;
    thetanew=theta;
end

% Checking if the error is too big --->stop
if errornew>100
    threshold=thresholdmax-1;
end

% threshold update
threshold=threshold+1;

end
end
```

## APPENDIX A - Script Test1 (1 of 2)

```
% test1

N_layer=1; % number of hidden layer
dim=5; % dimension of any hidden layer
lrate=0.01; % learning rate
f=@(x) cos(x); % initial condition

a=-pi; % left bound
b=pi; %right bound

% construction of omega and theta
omega = cell(1, N_layer + 1);
theta = cell(1, N_layer + 1);

% Initialization of weights with Xavier
W = sqrt(2/(dim+2))*randn(dim,2);
theta{1} = W(:, 1);
omega{1} = W(:, 2:end);

for k = 2:N_layer
    W = sqrt(2/(2*dim))*randn(dim,dim+1);
    theta{k}=W(:,1);
    omega{k}=W(:,2:end);
end

W = sqrt(2/(1+dim))*randn(1,dim+1);
theta{N_layer+1}=W(:,1);
omega{N_layer+1}=W(:,2:end);

% Construction of the test set
ntest=100;
X_test=linspace(a,b, ntest);

% calculation of the error
error=LossFunction(N_layer,X_test,omega,theta,f,a,b);
disp("error= "+error);

% weights update
[errornew,thetanew,omeganew]=WeightsUpdate(N_layer,omega,theta, ...
                                              X_test,a,b,f,lrate);
disp("errornew= "+errornew);

% saving the new error and the new weights
if errornew<error
    error=errornew;
    for i = 1:N_layer+1
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
    omega=omeganew;
    theta=thetanew;
```

APPENDIX A - Script Test1 (2 of 2)

```
end

%calculation of u
u=zeros(1,ntest);
for i=1:ntest
    x=X_test(i);
    X=NeuralNetwork(N_layer, omega, theta, x);
    u(i)=X{N_layer+1};
end

%comparison with the exact solution
sol=zeros(1,ntest);
for i=1:ntest
    x=X_test(i);
    sol(i)=f(x);
end

plot(X_test, u, '-r','DisplayName', 'Approximated solution');
hold on;
plot(X_test, sol, '-b','DisplayName', 'exact solution');
title('Graph of y(t)');
legend;
```

APPENDIX A - Script Test2 (1 of 2)

```
% test2

N_layer=1; % number of hidden layer
dim=5; % dimension of any hidden layer
lrate=0.001; % learning rate
f=@(x) cos(x); % initial condition

a=-pi;
b=pi;

% construction of omega and theta
omega = cell(1, N_layer + 1);
theta = cell(1, N_layer + 1);

% Reading the weights and the biases
for i = 1:N_layer+1
    fileNameOmega = ['omega', num2str(i), '.xlsx'];
    omega{i} = readmatrix(fileNameOmega);
    fileNameTheta = ['theta', num2str(i), '.xlsx'];
    theta{i} = readmatrix(fileNameTheta);
end

ntest=100;
X_test=linspace(a,b, ntest);

% calculation of the error
error=LossFunction(N_layer,X_test,omega,theta,f,a,b);
disp("error= "+error);

% weights update
[errornew,thetanew,omeganew]=WeightsUpdate(N_layer,omega,theta, ...
                                              X_test,a,b,f,lrate);
disp("errornew= "+errornew);

% saving the new error and the new weights
if errornew<error
    error=errornew;
    for i = 1:N_layer+1
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
    omega=omeganew;
    theta=thetanew;
end

%calculation of u
u=zeros(1,ntest);
for i=1:ntest
    x=X_test(i);
    X=NeuralNetwork(N_layer, omega, theta, x);
    u(i)=X{N_layer+1};
```

APPENDIX A - Script Test2 (2 of 2)

```
end

%comparison with the exact solution
sol=zeros(1,ntest);
for i=1:ntest
    x=X_test(i);
    sol(i)=f(x);
end

plot(X_test, u, '-r','DisplayName', 'Approximated solution');
hold on;
plot(X_test, sol, '-b','DisplayName', 'exact solution');
title('Graph of y(t)');
legend;
```

## APPENDIX B - SPATIALLY HOMOGENEOUS CASE

### APPENDIX B - Function SHCNeuralNetwork (1 of 1)

```
function X = SHCNeuralNetwork (N_layer, omega, theta, t)

% sigmoid function
sigmoid = @(u) 1./(1+exp(-u));

% column vector array of node values
X=cell(1,N_layer+1);
X{1}=sigmoid(omega{1}*t+theta{1});

for i=2:N_layer
    % value of nodes in layer i
    X{i}=sigmoid(omega{i}*X{i-1}+theta{i});
end
% activation function is not applied to the last layer
X{N_layer+1}=omega{N_layer+1}*X{N_layer}+theta{N_layer+1};

end
```

APPENDIX B - Function SHCTimeDerivative (1 of 1)

```
function dt = SHCTimeDerivative(N_layer, omega, theta, t)

% calculation of node values
X=SHCNeuralNetwork(N_layer, omega, theta, t);

% calculation of the time derivative
dt=cell(1,N_layer+1);
dt{1}=X{1}.* (1-X{1}).*omega{1};

for i=2:N_layer
    dt{i}=omega{i}*dt{i-1};
    dt{i}=X{i}.* (1-X{i}).*dt{i};
end

dt{N_layer+1}=omega{N_layer+1}*dt{N_layer};

end
```

## APPENDIX B - Function SHCLossFunction (1 of 1)

```
function L = SHCLossFunction(N_layer,omega,theta,t,T,y0)

% calculation of the outputs
u=cell(1, length(t));
for i=1:length(t)
    X=SHCNeuralNetwork(N_layer, omega, theta, t(i));
    u{i}=X{N_layer+1};
end

% calculation of the test function
y=cell(1, length(t));
for i=1:length(t)
    y{i}=y0+t(i)*u{i};
end

% calculation of the time derivative of u
du=cell(1,length(t));
for i=1:length(t)
    dt=SHCTimeDerivative(N_layer, omega,theta, t(i));
    du{i}=dt{N_layer+1};
end

% calculation of the time derivative of y
dy=cell(1,length(t));
for i=1:length(t)
    dy{i}=u{i}+t(i)*du{i};
end

% calculation of the loss function
L=0;
for i=1:length(t)
    L=L+(dy{i}-y{i}*(1-y{i}))^2;
end
L=T/length(t)*L;

end
```

## APPENDIX B - Function SHCWeightsUpdate (1 of 3)

---

```
function [error,thetanew,omeganew]=SHCWeightsUpdate(N_layer,omega,theta, ...
t_test,T,y0,lrate)

% parameters initialization
error=1;
errormax=0.0001;
threshold=1;
thresholdmax=100000;
W=cell(1,N_layer+1);
thetanew=theta;
omeganew=omega;

% weights update
while error>errormax && threshold<thresholdmax

    % training nodes
    n_train=10;
    t_train=T*rand(1,n_train);

    % calculation of the outputs
    u=cell(1, length(t_train));
    for i=1:length(t_train)
        X=SHCNeuralNetwork(N_layer,omega,theta,t_train(i));
        u{i}=X{N_layer+1};
    end

    % calculation of the time derivatives of u
    ut=cell(1,length(t_train));
    for i=1:length(t_train)
        dt=SHCTimeDerivative(N_layer,omega,theta,t_train(i));
        ut{i}=dt{N_layer+1};
    end

    % calculation of the test function
    y=cell(1,length(t_train));
    for i=1:length(t_train)
        y{i}=y0+t_train(i)*u{i};
    end

    % calculation of the time derivatives of u
    yt=cell(1,length(t_train));
    for i=1:length(t_train)
        yt{i}=u{i}+t_train(i)*ut{i};
    end

    % initialization of alfa and beta
    alfa=cell(length(t_train),N_layer+1);
    beta=cell(length(t_train),N_layer+1);

    for i=1:length(t_train)
        % weights update of the layer N+1
        X=SHCNeuralNetwork(N_layer,omega,theta,t_train(i));
        alfa{i,N_layer+1}=1;
```

## APPENDIX B - Function SHCWeightsUpdate (2 of 2)

```

beta{i,N_layer+1}=0;
dyt=alfa{i,N_layer+1}+t_train(i)*beta{i,N_layer+1};
dy=t_train(i)*alfa{i,N_layer+1};
dL=(yt{i}-y{i}).*(1-y{i})).*(dyt-(1-2*y{i}).*dy)*[1;X{N_layer}]';
W{N_layer+1}=[theta{N_layer+1},omega{N_layer+1}];
W{N_layer+1}=W{N_layer+1}-lrate*dL;
theta{N_layer+1}=W{N_layer+1}{:,1};
omega{N_layer+1}=W{N_layer+1}{:,2:end};

for k=N_layer:-1:2
    % weights update of the layer k
    X=SHCNeuralNetwork(N_layer, omega, theta, t_train(i));
    dt=SHCTimeDerivative(N_layer, omega, theta, t_train(i));
    dsigma=X{k}.*(1-X{k}).*(1-2*X{k}).*(omega{k}*dt{k-1});
    alfa{i,k}=alfa{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}));
    beta{i,k}=beta{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}))+...
        alfa{i,k+1}.*omega{k+1}*(dsigma);
    dyt=alfa{i,k}+t_train(i)*beta{i,k};
    dy=t_train(i)*alfa{i,k};
    dL=(yt{i}-y{i})*(1-y{i}))*((dyt-(1-2*y{i})*dy)*[1;X{k-1}]';
    W{k}=[theta{k},omega{k}];
    W{k}=W{k}-lrate*dL;
    theta{k}=W{k}{:,1};
    omega{k}=W{k}{:,2:end};
end

% weights update of the layer 1
X=SHCNeuralNetwork(N_layer, omega, theta, t_train(i));
dsigma=X{1}.*(1-X{1}).*(1-2*X{1}).*omega{1};
alfa{i,1}=alfa{i,2}.*omega{2}*(X{1}.*(1-X{1}));
beta{i,1}=beta{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
    alfa{i,2}.*omega{2}*(dsigma);
dyt=alfa{i,1}+t_train(i)*beta{i,1};
dy=t_train(i)*alfa{i,1};
dL=(yt{i}-y{i})*(1-y{i}))*((dyt-(1-2*y{i})*dy)*[1;t_train(i)]';
W{1}=[theta{1},omega{1}];
W{1}=W{1}-lrate*dL;
theta{1}=W{1}{:,1};
omega{1}=W{1}{:,2:end};
end

% error calculation
errornew=SHCLossFunction(N_layer, omega, theta, t_test, T, y0)

% Checking if the error is too big
if errornew>1000
    threshold=thresholdmax-1;
end

% error and weights update
if errornew<error
    error=errornew;
    thetanew=theta;

```

APPENDIX B - Function SHCWeightsUpdate (3 of 3)

```
omeganew=omega;
end

% threshold update
threshold=threshold+1

end
end
```

APPENDIX B - Script SHC Parameters Optimization (1 of 1)

```
% SHC Parameters Optimization

T=5; % time interval [0,T]
y0=1.5; %initial condition

% training times
n_train=21;
t_train=T*rand(1,n_train);

% testing times
n_test=46;
t_test=linspace(0, T, n_test);
error=1;

for n=1:10
    for d=2:10
        for lr=0.001:0.001:0.01

            % construction of omega and theta
            omega=cell(1, n + 1);
            theta=cell(1, n + 1);

            % Initialization of weights by Xavier's method
            W=sqrt(2/(d+2))*randn(d,2);
            theta{1}=W(:,1); % bias 1
            omega{1}=W(:,2:end); % matrice 1 (spazio e tempo)
            for k = 2:n
                W=sqrt(2/(2*d))*randn(d,d+1);
                theta{k}=W(:,1);
                omega{k}=W(:,2:end);
            end
            W=sqrt(2/(1+d))*randn(1,d+1);
            theta{n+1}=W(:,1);
            omega{n+1}=W(:,2:end);

            % weights update
            [errornew,thetanew,omeganew]=SHCWeightsUpdate(n,omega,theta, ...
                t_test,T,y0,lrate);

            if errornew<error
                error=errornew; %error update
                %updating the parameters
                N_layer=n;
                dimension=d;
                lrate=lr;

            end
        end
    end
end

disp("error= "+error+", N_layer= "+N_layer+", dimension= "+...
    dimension+", lrate= "+lrate );
```

## APPENDIX B - Script SHC Weights Optimization (1 of 2)

```
% SHC Weights Optimization

T=5; % time interval [0,T]
N_layer=1; % number of hidden layer
dim=3; % dimension of any hidden layer
lrate=0.001; % learning rate
y0=1.5;

% testing times
n_test=50;
t_test=linspace(0, T, n_test);

% parameters initialization
error=0.01;
errormax=0.0005;
threshold=0;
thresholdmax=1000;
W=cell(1,N_layer+1);

while error>errormax && threshold<thresholdmax

    threshold=threshold+1;
    omega = cell(1, N_layer + 1); % array di matrici dei pesi
    theta = cell(1, N_layer + 1); % array di vettori di bias

    % Inizializzazione dei pesi con Xavier
    W = sqrt(2/(dim+2))*randn(dim,2);
    theta{1}=W(:,1); % bias 1
    omega{1}=W(:,2:end); % matrice 1 (spazio e tempo)

    for k=2:N_layer
        W=sqrt(2/ (2*dim))*randn(dim,dim+1);
        theta{k}=W(:, 1);
        omega{k}=W(:, 2:end);
    end

    W=sqrt(2/ (1+dim))*randn(1,dim+1);
    theta{N_layer+1}=W(:,1);
    omega{N_layer+1}=W(:,2:end);

    % calculation of the new error
    [errornew,thetanew,omeganew]=SHCWeightsUpdate(N_layer,omega,theta, ...
        t_test,T,y0,lrate);

    if errornew<error
        error=errornew; %error update
        for i = 1:N_layer+1
            % saving omega and theta
            fileNameOmega = ['omega', num2str(i), '.xlsx'];
            writematrix(omeganew{i}, fileNameOmega);
            fileNameTheta = ['theta', num2str(i), '.xlsx'];
            writematrix(theta{i}, fileNameTheta);
        end
    end
end
```

APPENDIX B - Script SHC Weights Optimization (2 of 2)

```
    end
end
disp("error= "+error);
```

APPENDIX B - Script Test1 (1 of 2)

```
% Test1

T=5; % time interval [0,T]
N_layer=1; % number of hidden layer
dim=3; % dimension of any hidden layer
lrate=0.001; % learning rate
y0=1.5; % initial condition

% testing times
n_test=50;
t_test=linspace(0, T, n_test);

% construction of omega and theta
omega=cell(1, N_layer + 1);
theta=cell(1, N_layer + 1);

% Initialization of weights with Xavier
W = sqrt(2/(dim+2))*randn(dim,2);
theta{1}=W(:,1);
omega{1}=W(:,2:end);

for k=2:N_layer
    W=sqrt(2/ (2*dim))*randn(dim,dim+1);
    theta{k}=W(:, 1);
    omega{k}=W(:, 2:end);
end

W=sqrt(2/(1+dim))*randn(1,dim+1);
theta{N_layer+1}=W(:,1);
omega{N_layer+1}=W(:,2:end);

%calculation of the error
error=SHCLossFunction(N_layer,omega,theta,t_test,T,y0);
disp("error= "+error);

% weights update
[errornew,thetanew,omeganew]=SHCWeightsUpdate(N_layer,omega,theta, ...
                                                 t_test,T,y0,lrate);
disp("errornew= "+errornew);

%saving the error and the weights
if errornew<error
    for i = 1:N_layer+1
        % saving omega and theta
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
end

% Comparison of the graph with the exact solution
```

APPENDIX B - Script Test1 (2 of 2)

```
y=[];
for n=1:length(t_test)
    output=SHCNeuralNetwork(N_layer,omeganew,thetanew,t_test(n));
    y=[y,y0+ t_test(n)*output{N_layer+1}];
end

solution=@(z) (y0*exp(z))/(1+y0*(exp(z)-1));
sol=[];

for n=1:length(t_test)
    sol=[sol,solution(t_test(n))];
end

err=abs(sol-y).^2;
difference=sum(err)/n_test;

plot(t_test, y, '-r','DisplayName', 'Approximated solution');
hold on;
plot(t_test, sol, '-b','DisplayName', 'Exact solution');
xlabel('time');
ylabel('y(t)');
title('Graph of y(t)');
legend;
```

---

## APPENDIX B - Script Test2 (1 of 2)

```
% Test2

T=5; % time interval [0,T]
N_layer=1; % number of hidden layer
dim=3; % dimension of any hidden layer
lrate=0.0001; % learning rate
y0=1.5; % initial condition

% testing times
n_test=50;
t_test=linspace(0, T, n_test);

% construction of omega and theta
omega=cell(1, N_layer + 1);
theta=cell(1, N_layer + 1);

%reading omega and theta
for i = 1:N_layer+1
    fileNameOmega = ['omega', num2str(i), '.xlsx'];
    omega{i} = readmatrix(fileNameOmega);
    fileNameTheta = ['theta', num2str(i), '.xlsx'];
    theta{i} = readmatrix(fileNameTheta);
end

%calculation of the error
error=SHCLossFunction(N_layer,omega,theta,t_test,T,y0);
disp("error= "+error);

% weights update
[errornew,thetanew,omeganew]=SHCWeightsUpdate(N_layer,omega,theta, ...
t_test,T,y0,lrate);
disp("errornew= "+errornew);

%saving the error and the weights
if errornew<error
    for i = 1:N_layer+1
        % saving omega and theta
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
end

% Comparison of the graph with the exact solution
y=[];
for n=1:length(t_test)
    output=SHCNeuralNetwork(N_layer,omeganew,thetanew,t_test(n));
    y=[y,y0+ t_test(n)*output{N_layer+1}];
end

solution=@(z) (y0*exp(z))/(1+y0*(exp(z)-1));
sol=[];
```

APPENDIX B - Script Test (2 of 2)

```
for n=1:length(t_test)
    sol=[sol,solution(t_test(n))];
end

plot(t_test, y, '-r','DisplayName', 'Approximated solution');
hold on;
plot(t_test, sol, '-b','DisplayName', 'Exact solution');
xlabel('time');
ylabel('y(t)');
title('Graph of y(t)');
legend;
```

## APPENDIX C - GENERAL CASE WITH SPATIALLY HOMOGENEOUS INITIAL CONDITION

### APPENDIX C - Function NeuralNetwork (1 of 1)

```
function X = NeuralNetwork (N_layer, omega, theta, X0)

% sigmoid function
sigmoid = @(u) 1./(1+exp(-u));

% column vector array of node values
X=cell(1,N_layer+1);
X{1}=sigmoid(omega{1}*X0+theta{1});

for i=2:N_layer
    % value of nodes in layer i
    X{i}=sigmoid(omega{i}*X{i-1}+theta{i});
end

% activation function is not applied to the last layer
X{N_layer+1}=omega{N_layer+1}*X{N_layer}+theta{N_layer+1};

end
```

## APPENDIX C - Function TimeDerivative (1 of 1)

```
function dt = TimeDerivative(N_layer, omega, theta, X0)

% calculation of node values
X=NeuralNetwork(N_layer, omega, theta, X0);

% calculation of the time derivative
dt=cell(1,N_layer+1);
dt{1}=X{1}.*(1-X{1}).*omega{1}(:,2);

for i=2:N_layer
    dt{i}=omega{i}*dt{i-1};
    dt{i}=X{i}.*(1-X{i}).*dt{i};
end

dt{N_layer+1}=omega{N_layer+1}*dt{N_layer};

end
```

## APPENDIX C - Function SpaceDerivative (1 of 1)

```
function dx = SpaceDerivative(N_layer, omega, theta, X0)

% calculation of node values
X=NeuralNetwork(N_layer, omega, theta, X0);

% calculation of the time derivative
dx=cell(1,N_layer+1);
dx{1}=X{1}.* (1-X{1}).*omega{1} (:,1);

for i=2:N_layer
    dx{i}=omega{i}*dx{i-1};
    dx{i}=X{i}.* (1-X{i}).*dx{i};
end
dx{N_layer+1}=omega{N_layer+1}*dx{N_layer};

end
```

## APPENDIX C - Function Laplacian (1 of 1)

```
function lap = Laplacian(N_layer, omega,theta, X0)

%calculation of the array of node values
X=NeuralNetwork(N_layer,omega,theta,X0);

%calculation of the space derivative
dx=SpaceDerivative(N_layer,omega,theta,X0);

%calculation of the laplacian
lap=cell(1,N_layer+1);
lap{1}=X{1}.* (1-X{1}).*(1-2*X{1}).*(omega{1}{:,1}.^2);

for i=2:N_layer
    lap{i}=X{i}.* (1-X{1}).*(1-2*X{i}).*((omega{i}*dx{i-1}).^2)+ ...
        X{i}.* (1-X{i}).*(omega{i}*lap{i-1});
end

lap{N_layer+1} = omega{N_layer + 1}*lap{N_layer};

end
```

## APPENDIX C - Function LossFunction (1 of 2)

```

function L = LossFunction(N_layer,omega,theta,X_test,T,a,b,ya,yb,dya,dyb)

% calculation of the outputs
u=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    X=NeuralNetwork(N_layer,omega,theta,X_test(:,i));
    u{i}=X{N_layer+1};
end

% calculation of the trial function
y=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);
    y{i}=ya(t)*(b-x)/(b-a)+yb(t)*(x-a)/(b-a)+...
        t*(b-x)*(x-a)/(b-a)^2*u{i};
end

% calculation of the time derivative of u
ut=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dt=TimeDerivative(N_layer,omega,theta,X_test(:,i));
    ut{i}=dt{N_layer+1};
end

% calculation of the space derivative of u
ux=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dx=SpaceDerivative(N_layer,omega,theta,X_test(:,i));
    ux{i}=dx{N_layer+1};
end

% calculation of the laplacian of u
uxx=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dxx=Laplacian(N_layer,omega,theta,X_test(:,i));
    uxx{i}=dxx{N_layer+1};
end

% calculation of the time derivative of y
yt=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);
    yt{i}=dya(t)*(b-x)/(b-a)+dyb(t)*(x-a)/(b-a)+...
        (b-x)*(x-a)/(b-a)^2*u{i}+t*(b-x)*(x-a)/(b-a)^2*ut{i};
end

% calculation of the laplacian of y
yxx=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);

```

## APPENDIX C - Function LossFunction (2 of 2)

```
yxx{i}=-t*(2/(b-a)^2)*u{i}+2*t*(-2*x+b+a)/(b-a)^2*ux{i}+...
t*(b-x)*(x-a)/(b-a)^2*uxx{i};
end

% calculation of the loss function
sum=0;
for i=1:length(X_test(1,:))
    sum=sum+(yt{i}-yxx{i}-y{i}*(1-y{i}))^2;
end
L=(T*(b-a)/length(X_test(1,:)))*sum;

end
```

## APPENDIX C - Function WeightsUpdate (1 of 4)

```

function [error,thetanew,omeganew] = WeightsUpdate(N_layer,omega,theta, ...
X_test,ya,yb,dya,dyb,a,b,T,lrate)

% parameters initialization
error=10;
errormax=0.0001;
threshold=1;
thresholdmax=500;
W=cell(1,N_layer+1);

% weights update
while error>errormax && threshold<thresholdmax

    % training nodes
    ntime_train=50;
    time_train=T*rand(1,ntime_train);

    nspace_train=20;
    space_train=a+(b-a)*rand(1,nspace_train);

    X_train=[];
    for i=1:ntime_train
        for j=1:nspace_train
            X_train=[X_train,[space_train(j),time_train(i)']];
        end
    end

    % calculation of the outputs
    u=cell(1, length(X_train(1,:)));
    for i=1:length(X_train(1,:))
        X=NeuralNetwork(N_layer, omega, theta, X_train(:,i));
        u{i}=X{N_layer+1};
    end

    % calculation of the trial function
    y=cell(1, length(X_train(1,:)));
    for i=1:length(X_train(1,:))
        x=X_train(1,i);
        t=X_train(2,i);
        y{i}=ya(t)*(b-x)/(b-a)+yb(t)*(x-a)/(b-a)+...
        t*(b-x)*(x-a)/(b-a)^2*u{i};
    end

    % calculation of the time derivative of u
    ut=cell(1,length(X_train(1,:)));
    for i=1:length(X_train(1,:))
        dt=TimeDerivative(N_layer, omega,theta, X_train(:,i));
        ut{i}=dt{N_layer+1};
    end

    % calculation of the space derivative of u
    ux=cell(1,length(X_train(1,:)));
    for i=1:length(X_train(1,:))

```

## APPENDIX C - Function WeightsUpdate (2 of 4)

```

dx=SpaceDerivative(N_layer, omega,theta, X_train(:,i));
ux{i}=dx{N_layer+1};
end

% calculation of the laplacian of u
uxx=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    dxx=Laplacian(N_layer, omega,theta, X_train(:,i));
    uxx{i}=dxx{N_layer+1};
end

% calculation of the time derivative of y
yt=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    x=X_train(1,i);
    t=X_train(2,i);
    yt{i}=dya(t)*(b-x)/(b-a)+dyb(t)*(x-a)/(b-a)+...
        (b-x)*(x-a)/(b-a)^2*u{i}+t*(b-x)*(x-a)/(b-a)^2*ut{i};
end

% calculation of the laplacian of y
yxx=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    x=X_train(1,i);
    t=X_train(2,i);
    yxx{i}=-t*(2/(b-a)^2)*u{i}+2*t*(-2*x+b+a)/(b-a)^2*ux{i}+...
        t*(b-x)*(x-a)/(b-a)^2*uxx{i};
end

% initialization of alfa, betax, betat and gamma
alfa=cell(length(X_train(1,:)),N_layer+1);
betax=cell(length(X_train(1,:)),N_layer+1);
betat=cell(length(X_train(1,:)),N_layer+1);
gamma=cell(length(X_train(1,:)),N_layer+1);

for i=1:length(X_train(1,:))
    % weights update of the layer N+1
    X=NeuralNetwork(N_layer,omega,theta,X_train(:,i));
    alfa{i,N_layer+1}=1;
    betax{i,N_layer+1}=0;
    betat{i,N_layer+1}=0;
    gamma{i,N_layer+1}=0;
    x=X_train(1,i);
    t=X_train(2,i);
    dy=t*(b-x)*(x-a)/(b-a)^2;
    dyt=(b-x)*(x-a)/(b-a)^2;
    dyxx=-t*2/(b-a)^2;
    dL=(yt{i}-yxx{i}-y{i}.* (1-y{i})).* ...
        (dyt-dyxx-(1-2*y{i}).*dy)*[1;X{N_layer}]';
    W{N_layer+1}=[theta{N_layer+1},omega{N_layer+1}];
    W{N_layer+1}=W{N_layer+1}-lrate*dL;
    theta{N_layer+1}=W{N_layer+1}(:,1);
end

```

## APPENDIX C - Function WeightsUpdate (3 of 4)

```

omega{N_layer+1}=W{N_layer+1} (:,2:end);

for k=N_layer:-1:2
    % weights update of the layer k
    X=NeuralNetwork(N_layer,omega,theta,X_train(:,i));
    timederiv=TimeDerivative(N_layer,omega,theta,X_train(:,i));
    spacederiv=SpaceDerivative(N_layer,omega,theta,X_train(:,i));
    lap=Laplacian(N_layer,omega,theta,X_train(:,i));
    dsigmax=X{k}.*(1-X{k}).*(1-2*X{k}).*(omega{k}*spacederiv{k-1});
    dsigmat=X{k}.*(1-X{k}).*(1-2*X{k}).*(omega{k}*timederiv{k-1});
    dsigmax2=X{k}.*(1-X{k}).*(1-6*X{k}.*(1-X{k})).*...
        (omega{k}*spacederiv{k-1}).^2+X{k}.*(1-X{k}).*...
        (1-2*X{k}).*(omega{k}*lap{k-1});
    alfa{i,k}=alfa{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}));
    betax{i,k}=betax{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}))+...
        alfa{i,k+1}.*(omega{k+1}*dsigmat);
    betat{i,k}=betat{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}))+...
        alfa{i,k+1}.*(omega{k+1}*dsigmax);
    gamma{i,k}=gamma{i,k+1}.*omega{k+1}*(X{k}.*(1-X{k}))+...
        2*betax{i,k+1}.*(omega{k+1}*dsigmax)+...
        alfa{i,k+1}.*(omega{k+1}*dsigmax2);
    x=X_train(1,i);
    t=X_train(2,i);
    dy=t*(b-x)*(x-a)/(b-a)^2*alfa{i,k};
    dyt=(b-x)*(x-a)/(b-a)^2*alfa{i,k}+...
        t*(b-x)*(x-a)/(b-a)^2*betat{i,k};
    dyxx=-t^2/(b-a)^2*alfa{i,k}+2*t*(-2*x+b+a)/(b-a)^2*betax{i,k}+...
        t*(b-x)*(x-a)/(b-a)^2*gamma{i,k};
    dL=(yt{i}-yxx{i}-y{i}.*(1-y{i})).*...
        (dyt-dyxx-(1-2*y{i}).*dy)*[1;X{k-1}]';
    W{k}=[theta{k},omega{k}];
    W{k}=W{k}-lrate*dL;
    theta{k}=W{k}(:,1);
    omega{k}=W{k}(:,2:end);
end

% weights update of the layer 1
X=NeuralNetwork(N_layer,omega,theta,X_train(:,i));
dsigmax=X{1}.*(1-X{1}).*(1-2*X{1}).*omega{1}(:,1);
dsigmat=X{1}.*(1-X{1}).*(1-2*X{1}).*omega{1}(:,2);
dsigmax2=X{1}.*(1-X{1}).*(1-6*X{1}.*(1-X{1})).*omega{1}(:,1).^2;
alfa{i,1}=alfa{i,2}.*omega{2}*(X{1}.*(1-X{1}));
betax{i,1}=betax{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
    alfa{i,2}.*(omega{2}*dsigmat);
betat{i,1}=betat{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
    alfa{i,2}.*(omega{2}*dsigmax);
gamma{i,1}=gamma{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
    2*betax{i,2}.*(omega{2}*dsigmax)+...
    alfa{i,2}.*(omega{2}*dsigmax2);
x=X_train(1,i);
t=X_train(2,i);
dy=t*(b-x)*(x-a)/(b-a)^2*alfa{i,1};
dyt=(b-x)*(x-a)/(b-a)^2*alfa{i,1}+t*(b-x)*(x-a)/(b-a)^2*betat{i,1};

```

## APPENDIX C - Function WeightsUpdate (4 of 4)

```
dyxx=-t*2/(b-a)^2*alfa{i,1}+2*t*(-2*x+b+a)/(b-a)^2*betax{i,1}+...
t*(b-x)*(x-a)/(b-a)^2*gamma{i,1};
dL=(yt{i}-yxx{i}-y{i}.*(1-y{i})).*...
(dy-t-dyxx-(1-2*y{i}).*dy)*[1;X_train(:,i)]';
W{1}=[theta{1},omega{1}];
W{1}=W{1}-lrate*dL;
theta{1}=W{1}(:,1);
omega{1}=W{1}(:,2:end);
end
% error calculation
errornew=LossFunction(N_layer,omega,theta,X_test,T,a,b,ya,yb,dy,a,dy,b);

% error update
if errornew<error
    error=errornew;
    thetanew=theta;
    omeganew=omega;
end

% threshold update
threshold=threshold+1;

end
end
```

## APPENDIX C - Script Test (1 of 3)

```
% test

T=5; % time interval [0,T]
N_layer=1; % number of hidden layer
dim=2; % dimension of any hidden layer
lrate=0.001; % learning rate
y0=1.5; % initial condition
a=-3;
b=3;

%boundary conditions
ya= @(t) (y0*exp(t))/(1+y0*(exp(t)-1));
yb=ya;
dyya= @(t) (y0*exp(t)*(1-y0))/((1+y0*(exp(t)-1))^2);
dyb=dyya;

% testing nodes
ntime_test=50;
time_test=linspace(0, T, ntime_test);

nspacetest=20;
space_test=linspace(a,b, nspacetest);

X_test=[];
for i=1:ntime_test
    for j=1:nspacetest
        X_test=[X_test, [space_test(j), time_test(i)]'];
    end
end

% construction of omega and theta
omega=cell(1, N_layer + 1);
theta=cell(1, N_layer + 1);

% Xavier initialization
W=sqrt(2/(dim + 2))*randn(dim, 3);
theta{1}=W(:, 1); % bias 1
omega{1}=W(:, 2:end); % matrice 1 (spazio e tempo)

for k=2:N_layer
    W=sqrt(2/(2*dim))*randn(dim, dim+1);
    theta{k}=W(:,1);
    omega{k}=W(:,2:end);
end

W=sqrt(2/(1+dim))*randn(1,dim+1);
theta{N_layer+1}=W(:,1);
omega{N_layer+1}=W(:,2:end);

%calculation of the error
error=LossFunction(N_layer,omega,theta,X_test,T,a,b,ya,yb,dyya,dyb);
disp("error= "+error);
```

## APPENDIX C - Script Test (2 of 3)

```
% weights update
[errornew,thetanew,omeganew]=WeightsUpdate(N_layer,omega,theta, ...
X_test,ya,yb,dya,dyb,a,b,T,lrate);
disp("errornew= "+errornew);

%saving the error and the new weights
if errornew<error
    omega=omeganew;
    theta=thetanew;
    for i = 1:N_layer+1
        % saving omega and theta
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
end

% Comparison of the graph with the exact solution on the bound

% Solution evaluated on the boundary a
bya=[];
for i=1:length(X_test(:,1))
    X=NeuralNetwork(N_layer, omega, theta, X_test(:,i));
    if X_test(1,i)==a
        x=a;
        t=X_test(2,i);
        bya=[bya, ya(t)*(b-x)/(b-a)+yb(t)*(x-a)/(b-a)+...
        t*(b-x)*(x-a)/(b-a)^2*X{N_layer+1}];
    end
end

% Solution evaluated on the boundary b
byb=[];
for i=1:length(X_test(:,1))
    X=NeuralNetwork(N_layer, omega, theta, X_test(:,i));
    if X_test(1,i)==b
        x=b;
        t=X_test(2,i);
        byb=[byb, ya(t)*(b-x)/(b-a) + yb(t)*(x-a)/(b-a)+...
        t*(b-x)*(x-a)/(b-a)^2*X{N_layer+1}];
    end
end

% exact solution on the boundary
solution=@(z) (y0*exp(z))/(1+y0*(exp(z)-1));
sol=[];

for n=1:length(time_test)
    sol=[sol,solution(time_test(n))];
end
```

APPENDIX C - Script Test (3 of 3)

---

```
%error on the boundary a
erra=abs(sol-bya);
errorA=sum(erra)/ntime_test;
disp("errorA= "+errorA);

%error on the boundary b
errb=abs(sol-byb);
errorB=sum(errb)/ntime_test;
disp("errorB= "+errorB);

%plot of the comparison with the exact solution on the bound
figure;
plot(time_test, sol, '-b','DisplayName', 'Exact solution');
hold on;
plot(time_test, bya, '-r','DisplayName', 'Approximated solution on a');
plot(time_test, byb, '-g','DisplayName', 'Approximated solution on b');
xlabel('time');
ylabel('y(t)');
title('Graph of y(t)');
legend;

%construction of the approximated solution
Y=[];
for i=1:nspace_test
    for j=1:ntime_test
        X_test=[space_test(i),time_test(j)]';
        X=NeuralNetwork(N_layer, omega, theta, X_test);
        x=X_test(1);
        t=X_test(2);
        Y(i,j)=ya(t)*(b-x)/(b-a)+yb(t)*(x-a)/(b-a)+...
            t*((b-x)*(x-a))/((b-a)^2)*X{N_layer+1};
    end
end
Y=Y';

%graph of the approximated solution
figure;
[X,T]=meshgrid(space_test, time_test);
surf(X,T,Y);
xlabel('x (space)');
ylabel('t (time)');
zlabel('y(x,t)');
```

## APPENDIX D - GENERAL CASE WITH SINUSOIDAL INITIAL CONDITION

### APPENDIX D - Function GCNeuralNetwork (1 of 1)

```
function X = GCNeuralNetwork (N_layer, omega, theta, X0)

% sigmoid function
sigmoid = @(u) 1./(1+exp(-u));

% column vector array of node values
X=cell(1,N_layer+1);
X{1}=sigmoid(omega{1}*X0+theta{1});

for i=2:N_layer
    % value of nodes in layer i
    X{i}=sigmoid(omega{i}*X{i-1}+theta{i});
end
% activation function is not applied to the last layer
X{N_layer+1}=omega{N_layer+1}*X{N_layer}+theta{N_layer+1};

end
```

APPENDIX D - Function GCTimeDerivative (1 of 1)

```
function dt = GCTimeDerivative(N_layer, omega, theta, X0)

% calculation of node values
X=GCNeuralNetwork(N_layer, omega, theta, X0);

% calculation of the time derivative
dt=cell(1,N_layer+1);
dt{1}=X{1}.* (1-X{1}).*omega{1}(:,2);

for i=2:N_layer
    dt{i}=X{i}.* (1-X{i}).*omega{i}*dt{i-1};
end
dt{N_layer+1}=omega{N_layer+1}*dt{N_layer};

end
```

#### APPENDIX D - Function GCSpaceDerivative (1 of 1)

```
function dx = GCSpaceDerivative(N_layer, omega, theta, X0)

% calculation of node values
X=GCNeuralNetwork(N_layer, omega, theta, X0);

% calculation of the time derivative
dx=cell(1,N_layer+1);
dx{1}=X{1}.* (1-X{1}).*omega{1}(:,1);

for i=2:N_layer
    dx{i}=X{i}.* (1-X{i}).*omega{i}*dx{i-1};
end
dx{N_layer+1}=omega{N_layer+1}*dx{N_layer};

end
```

## APPENDIX D - Function GCLaplacian (1 of 1)

```
function lap = GCLaplacian(N_layer, omega, theta, X0)

%calculation of the array of node values
X= GCNeuralNetwork(N_layer, omega, theta, X0);

%calculation of the derivative
dx = GCSpaceDerivative(N_layer, omega, theta, X0);

%calculation of the laplacian
lap=cell(1,N_layer+1);
lap{1}=X{1}.* (1-X{1}).*(1-2*X{1}).*(omega{1}{:,1}.^2);

for i=2:N_layer
    lap{i}=X{i}.* (1-X{1}).*(1-2*X{i}).*((omega{i}*dx{i-1}).^2)+ ...
    X{i}.* (1-X{i}).*(omega{i}*lap{i-1});
end

lap{N_layer+1} = omega{N_layer + 1}*lap{N_layer};

end
```

## APPENDIX D - Function GCLossFunction (1 of 2)

```

function L=GCLossFunction(N_layer,omega,theta,X_test,X_test0,T,a,b,y0,ya,dya)

X_test=[X_test,X_test0];

% calculation of the outputs
u=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    X=GCNeuralNetwork(N_layer,omega,theta,X_test(:,i));
    u{i}=X{N_layer+1};
end

% calculation of the trial function
y=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);
    y{i}=ya(t)+(b-x)*(x-a)/(b-a)^2*u{i};
end

% calculation of the time derivative of u
ut=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dt=GCTimeDerivative(N_layer,omega,theta,X_test(:,i));
    ut{i}=dt{N_layer+1};
end

% calculation of the space derivative of u
ux=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dx=GCSpaceDerivative(N_layer,omega,theta,X_test(:,i));
    ux{i}=dx{N_layer+1};
end

% calculation of the laplacian of u
uxx=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    dxx=GCLaplacian(N_layer,omega,theta,X_test(:,i));
    uxx{i}=dxx{N_layer+1};
end

% calculation of the time derivative of y
yt=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);
    yt{i}=dya(t)+(b-x)*(x-a)/(b-a)^2*ut{i};
end

% calculation of the laplacian of y
yxx=cell(1,length(X_test(1,:)));
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    yxx{i}=-2/(b-a)^2*u{i}+2*(-2*x+b+a)/(b-a)^2*ux{i}+...
end

```

## APPENDIX D - Function GLossFunction (2 of 2)

```
-
(b-x) * (x-a) / (b-a)^2*uxx{i};
end

% calculation of the loss function
L1=0;
L2=0;
for i=1:length(X_test(1,:))
    x=X_test(1,i);
    t=X_test(2,i);
    if t~=0
        L1=L1+(yt{i}-yxx{i}-y{i}*(1-y{i}))^2;
    else
        L2=L2+(y{i}-y0(x))^2;
    end
end
L1=0.1*T*(b-a)/length(X_test(1,:))*L1;
L2=0.9*(b-a)/length(X_test0(1,:))*L2;
L=L1+L2;

end
```

APPENDIX D - Function GCWeightsUpdate (1 of 5)

```

function [error,thetanew,omeganew]=GCWeightsUpdate(N_layer,omega,theta, ...
X_test,X_test0,y0,ya,dya,a,b,T,lrate)

% parameters initialization
error=10;
errormax=0.005;
threshold=1;
thresholdmax=10000;
W=cell(1,N_layer+1);
max=0;
errornew=1000;
omeganew=omega;
thetanew=theta;
% weights update
while error>errormax && threshold<thresholdmax
    % training nodes
    ntime_train=10;
    time_train=T*rand(1,ntime_train);

    nspace_train=10;
    space_train=a+(b-a)*rand(1,nspace_train);

    X_train=[];
    for i=1:ntime_train
        for j=1:nspace_train
            X_train=[X_train,[space_train(j),time_train(i)']];
        end
    end

    nspace_train0=200;
    space_train0=a+(b-a)*rand(1,nspace_train0);

    X_train0=[];
    for i=1:nspace_train0
        X_train0=[X_train0,[space_train0(i),0']];
    end

    X_train=[X_train,X_train0];

    X_train=X_train(:, randperm(size(X_train, 2)));

    % calculation of the outputs
    u=cell(1,length(X_train(1,:)));
    for i=1:length(X_train(1,:))
        X=GCNeuralNetwork(N_layer,omega,theta,X_train(:,i));
        u{i}=X{N_layer+1};
    end

    % calculation of the trial function
    y=cell(1,length(X_train(1,:)));
    for i=1:length(X_train(1,:))
        x=X_train(1,i);
        t=X_train(2,i);
        y{i}= ...
    end
end

```

## APPENDIX D - Function GCWeightsUpdate (2 of 5)

```

y{i}=ya(t)+(b-x)*(x-a)/(b-a)^2*u{i};
end

% calculation of the time derivative of u
ut=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    dt=GCTimeDerivative(N_layer,omega,theta,X_train(:,i));
    ut{i}=dt{N_layer+1};
end

% calculation of the space derivative of u
ux=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    dx=GCSpaceDerivative(N_layer,omega,theta,X_train(:,i));
    ux{i}=dx{N_layer+1};
end

% calculation of the laplacian of u
uxx=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    dxx=GCLaplacian(N_layer,omega,theta,X_train(:,i));
    uxx{i}=dxx{N_layer+1};
end

% calculation of the time derivative of y
yt=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    x=X_train(1,i);
    t=X_train(2,i);
    yt{i}=dyt(t)+(b-x)*(x-a)/(b-a)^2*ut{i};
end

% calculation of the laplacian of y
yxx=cell(1,length(X_train(1,:)));
for i=1:length(X_train(1,:))
    x=X_train(1,i);
    yxx{i}=-2/(b-a)^2*u{i}+2*(-2*x+b+a)/(b-a)^2*ux{i}+...
        (b-x)*(x-a)/(b-a)^2*uxx{i};
end

% initialization of alfa, betax, betat and gamma
alfa=cell(length(X_train(1,:)),N_layer+1);
alfa0=cell(nspace_train0,N_layer+1);
betax=cell(length(X_train(1,:)),N_layer+1);
betat=cell(length(X_train(1,:)),N_layer+1);
gamma=cell(length(X_train(1,:)),N_layer+1);

for i=1:length(X_train(1,:))
    % weights update of the layer N+1
    X=GCNeuralNetwork(N_layer,omega,theta,X_train(:,i));
    x=X_train(1,i);
    t=X_train(2,i);
    if t~=0

```

#### APPENDIX D - Function GCWeightsUpdate (3 of 5)

```

alfa{i,N_layer+1}=1;
betax{i,N_layer+1}=0;
betat{i,N_layer+1}=0;
gamma{i,N_layer+1}=0;
dy=(b-x)*(x-a)/(b-a)^2;
dyt=0;
dyxx=-2/(b-a)^2;
dL1=2*(yt{i}-yxx{i}-y{i}.*(1-y{i}))* (dyt-dyxx-(1-2*y{i})*dy)*[1;X ↵
{N_layer}]';
W{N_layer+1}=[theta{N_layer+1},omega{N_layer+1}];
W{N_layer+1}=W{N_layer+1}-lrate*dL1;
theta{N_layer+1}=W{N_layer+1}(:,1);
omega{N_layer+1}=W{N_layer+1}(:,2:end);
else
alfa0{i,N_layer+1}=1;
dy=(b-x)*(x-a)/(b-a)^2;
dL2=2*(y0(a)+(b-x)*(x-a)/(b-a)^2*X{N_layer+1}-y0(x))*dy*[1;X ↵
{N_layer}]';
W{N_layer+1}=[theta{N_layer+1},omega{N_layer+1}];
W{N_layer+1}=W{N_layer+1}-lrate*dL2;
theta{N_layer+1}=W{N_layer+1}(:,1);
omega{N_layer+1}=W{N_layer+1}(:,2:end);
end

for k=N_layer:-1:2
% weights update of the layer k
X=GCNeuralNetwork(N_layer,omega,theta,X_train(:,i));
x=X_train(1,i);
t=X_train(2,i);
if t~=0
timeder=GCTimeDerivative(N_layer,omega,theta,X_train(:,i));
spaceder=GCSpaceDerivative(N_layer,omega,theta,X_train(:,i));
lap=GCLaplacian(N_layer,omega,theta,X_train(:,i));
dsigmax=X{k}.* (1-X{k}).*(1-2*X{k}).*(omega{k}*spaceder{k-1});
dsigmat=X{k}.* (1-X{k}).*(1-2*X{k}).*(omega{k}*timeder{k-1});
dsigmax2=X{k}.* (1-X{k}).*(1-6*X{k}.* (1-X{k})).* ...
(omega{k}*spaceder{k-1}).^2+X{k}.* (1-X{k}).* ...
(1-2*X{k}).*(omega{k}*lap{k-1});
alfa{i,k}=alfa{i,k+1}.* (omega{k+1}*(X{k}.* (1-X{k})));
betax{i,k}=betax{i,k+1}.* (omega{k+1}*(X{k}.* (1-X{k}))) + ...
alfa{i,k+1}.* (omega{k+1}*dsigmax);
betat{i,k}=betat{i,k+1}.* (omega{k+1}*(X{k}.* (1-X{k}))) + ...
alfa{i,k+1}.* (omega{k+1}*dsigmat);
gamma{i,k}=gamma{i,k+1}.* (omega{k+1}*(X{k}.* (1-X{k}))) + ...
2*betax{i,k+1}.* (omega{k+1}*dsigmax) + ...
alfa{i,k+1}.* (omega{k+1}*dsigmax2);
dy=(b-x)*(x-a)/(b-a)^2*alfa{i,k};
dyt=(b-x)*(x-a)/(b-a)^2*betat{i,k};
dyxx=-2/(b-a)^2*alfa{i,k}+...
2*(-2*x+b+a)/(b-a)^2*betax{i,k}+...
(b-x)*(x-a)/(b-a)^2*gamma{i,k};
dL1=2*(yt{i}-yxx{i}-y{i}.*(1-y{i}))* ...
(dydt-dyxx-(1-2*y{i})*dy)*[1;X{k-1}]';

```

APPENDIX D - Function GCWeightsUpdate (4 of 5)

```

W{k}=[theta{k},omega{k}];
W{k}=W{k}-lrate*dL1;
theta{k}=W{k}(:,1);
omega{k}=W{k}(:,2:end);
else
    alfa0{i,k}=alfa0{i,k+1}.*(omega{k+1}*(X{k}.*(1-X{k})));
    dy=(b-x)*(x-a)/(b-a)^2*alfa0{i,k};
    dL2=2*(y0(a)+(b-x)*(x-a)/(b-a)^2*X{N_layer+1}-y0(x))*dy*[1;X{k-1}]';
    W{k}=[theta{k},omega{k}];
    W{k}=W{k}-lrate*dL2;
    theta{k}=W{k}(:,1);
    omega{k}=W{k}(:,2:end);
end
end
% weights update of the layer 1
X=GCNeuralNetwork(N_layer,omega,theta,X_train(:,i));
x=X_train(1,i);
t=X_train(2,i);
if t~=0
    dsigmax=X{1}.*(1-X{1}).*(1-2*X{1}).*omega{1}(:,1);
    dsigmat=X{1}.*(1-X{1}).*(1-2*X{1}).*omega{1}(:,2);
    dsigmax2=X{1}.*(1-X{1}).*(1-6*X{1}.*(1-X{1})).*omega{1}(:,1).^2;
    alfa{i,1}=alfa{i,2}.*omega{2}*(X{1}.*(1-X{1}));
    betax{i,1}=betax{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
        alfa{i,2}.*(omega{2}*dsigmax);
    betat{i,1}=betat{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
        alfa{i,2}.*(omega{2}*dsigmat);
    gamma{i,1}=gamma{i,2}.*omega{2}*(X{1}.*(1-X{1}))+...
        2*betax{i,2}.*(omega{2}*dsigmax)+...
        alfa{i,2}.*(omega{2}*dsigmax2);
    x=X_train(1,i);
    dy=(b-x)*(x-a)/(b-a)^2*alfa{i,1};
    dyt=(b-x)*(x-a)/(b-a)^2*betat{i,1};
    dyxx=-2/(b-a)^2*alfa{i,1}+...
        2*(-2*x+b+a)/(b-a)^2*betax{i,1}+...
        (b-x)*(x-a)/(b-a)^2*gamma{i,1};
    dL1=(yt{i}-yxx{i}-y{i}.*(1-y{i})).*...
        (dyt-dyxx-(1-2*y{i})*dy)*[1;X_train(:,i)]';
    W{1}=[theta{1},omega{1}];
    W{1}=W{1}-lrate*dL1;
    theta{1}=W{1}(:,1);
    omega{1}=W{1}(:,2:end);
else
    alfa0{i,1}=alfa0{i,2}.*(omega{2}*(X{1}.*(1-X{1})));
    dy=(b-x)*(x-a)/(b-a)^2*alfa0{i,1};
    dL2=2*(y0(a)+(b-x)*(x-a)/(b-a)^2*X{N_layer+1}-y0(x))*dy*[1;X_train(:,i)]';
    W{1}=[theta{1},omega{1}];
    W{1}=W{1}-lrate*dL2;
    theta{1}=W{1}(:,1);
    omega{1}=W{1}(:,2:end);
end

```

#### APPENDIX D - Function GCWeightsUpdate (5 of 5)

```
end

errorpred=errornew;
% error calculation
errornew=GCLossFunction(N_layer,omega,theta,X_test,X_test0,T,a,b,y0,ya,dya);

% error update
if errornew<error
    error=errornew;
    omeganew=omega;
    thetanew=theta;
end

if errornew>errorpred
    max=max+1;
end

if errornew>10000
    threshold=thresholdmax-1;
end

% threshold update
threshold=threshold+1;

end

end
```

## APPENDIX D - Script GC Parameters Optimization (1 of 2)

```
% GC Parameters Optimization
T=5;
N_layer=3; % number of hidden layer
dim=10; % dimension of any hidden layer
lrate=0.01; % learning rate
y0=@(x) cos(pi*x/3)+2.5; % initial condition

a=-3;
b=3;

%boundary conditions
ya= @(t) y0(a)*exp(t)/(1+y0(a)*(exp(t)-1));
dyaa= @(t) (y0(a)*exp(t)*(1-y0(a)))/((1+y0(a)*(exp(t)-1))^2);

yb= @(t) (y0(b)*exp(t))/(1+y0(b)*(exp(t)-1));
dybb= @(t) (y0(b)*exp(t)*(1-y0(b)))/((1+y0(b)*(exp(t)-1))^2);

% testing nodes
ntime_test=50;
time_test=linspace(0, T, ntime_test);

nspacetest=20;
space_test=linspace(a,b, nspacetest);

X_test=[];
for i=1:ntime_test
    for j=1:nspacetest
        X_test=[X_test,[space_test(j),time_test(i)]'];
    end
end

nspacetest0=100;
space_test0=linspace(a,b, nspacetest0);

X_test0=[];
for i=1:nspacetest0
    X_test0=[X_test0,[space_test0(i),0]'];
end

error=1;
errormax=0.005;

for n=1:4
    for d=2:20

        % construction of omega and theta
        omega = cell(1, n + 1);
        theta = cell(1, n + 1);

        % Initialization of weights with Xavier
        W =sqrt(2/(d+2))*randn(d,3);
        theta{1} = W(:, 1); % bias 1
    end
end
```

## APPENDIX D - Script GC Parameters Optimization (2 of 2)

```

omega{1} = W(:, 2:end); % matrice 1 (spazio e tempo)

for k = 2:n
    W = sqrt(2/(2*d))*randn(d,d+1);
    theta{k}=W(:,1);
    omega{k}=W(:,2:end);
end

W = sqrt(2/(1+d))*randn(1,d+1);
theta{n+1}=W(:,1);
omega{n+1}=W(:,2:end);

%Since the initial condition is an even function, we impose symmetry
for i=1:d
    omega{1}(d-i+1,1)=-omega{1}(i,1);
    omega{1}(d-i+1,2)=omega{1}(i,2);
    theta{1}(d-i+1)=theta{1}(i);
    omega{2}(:,d-i+1)=omega{2}(:,i);
end

% weights update
[errornew,thetanew,omeganew]=GCWeightsUpdate(n,omega,theta, ...
    X_test,X_test0,y0,dy,a,b,T,lrate);

%saving the error and the new weights
if errornew<error
    error=errornew;
    for i = 1:n+1
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end

    %updating the parameters
    N_layer=n;
    dim=d;
    end
end
end

disp("error= "+error+", N_layer= "+N_layer+", dimension= "+...
    "dim+", "lrate= "+lrate );

```

## APPENDIX D - Script GC Weights Optimization (1 of 2)

```
% GC Weights Optimization

T=5; % time interval [0,T]
N_layer=2; % number of hidden layer
dim=19; % dimension of any hidden layer
lrate=0.01; % learning rate
y0=@(x) cos(pi*x/3)+2.5; % initial condition

a=-3;
b=3;

ya= @(t) y0(a)*exp(t)/(1+y0(a)*(exp(t)-1));
dyya= @(t) y0(a)*exp(t)*(1-y0(a))/(1+y0(a)*(exp(t)-1))^2;

yb= @(t) (y0(b)*exp(t))/(1+y0(b)*(exp(t)-1));
dyb= @(t) (y0(b)*exp(t)*(1-y0(b)))/((1+y0(b)*(exp(t)-1))^2);

% testing nodes
ntime_test=100;
time_test=linspace(0, T, ntime_test);

nspacetest=50;
space_test=linspace(a,b, nspacetest);

X_test=[];
for i=1:ntime_test
    for j=1:nspacetest
        X_test=[X_test,[space_test(j),time_test(i)]'];
    end
end

nspacetest0=100;
space_test0=linspace(a,b, nspacetest0);

X_test0=[];
for i=1:nspacetest0
    X_test0=[X_test0,[space_test0(i),0]'];
end

% parameters initialization
error=10;
errormax=0.001;
threshold=0;
thresholdmax=1000;
W=cell(1,N_layer+1);

while error>errormax && threshold<thresholdmax

    threshold=threshold+1;

    omega = cell(1, N_layer + 1); % array di matrici dei pesi
    theta = cell(1, N_layer + 1); % array di vettori di bias
```

## APPENDIX D - Script GC Weights Optimization (2 of 2)

```
% Initialization of weights with Xavier
W = sqrt(2/(dim+2))*randn(dim,3);
theta{1} = W(:, 1);
omega{1} = W(:, 2:end);

for k = 2:N_layer
    W = sqrt(2/(2*dim))*randn(dim,dim+1);
    theta{k}=W(:,1);
    omega{k}=W(:,2:end);
end

W = sqrt(2/(1+dim))*randn(1,dim+1);
theta{N_layer+1}=W(:,1);
omega{N_layer+1}=W(:,2:end);

%Since the initial condition is an even function, we impose symmetry
for i=1:dim
    omega{1}(dim-i+1,1)=-omega{1}(i,1);
    omega{1}(dim-i+1,2)=omega{1}(i,2);
    theta{1}(dim-i+1)=theta{1}(i);
    omega{2}(:,dim-i+1)=omega{2}(:,i);
end

% weights update
[errornew,thetanew,omeganew]=GCWeightsUpdate(N_layer,omega,theta, ...
X_test,X_test0,y0,ya,dy,a,b,T,lrate);

%saving the error and the new weights
if errornew<error
    error=errornew;
    for i = 1:N_layer+1
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(thetanew{i}, fileNameTheta);
    end
end
end

disp("error= "+error);
```

---

APPENDIX D - Script GC Test (1 of 3)

```
% test
T=5;
N_layer=2; % number of hidden layer
dim=19; % dimension of any hidden layer
lrate=0.001; % learning rate
y0=@(x) cos(pi*x/3)+2.5; % initial condition

a=-3;
b=3;

ya= @(t) y0(a)*exp(t)/(1+y0(a)*(exp(t)-1));
dyu= @(t) (y0(a)*exp(t)*(1-y0(a)))/((1+y0(a)*(exp(t)-1))^2);

yb= @(t) (y0(b)*exp(t))/(1+y0(b)*(exp(t)-1));
dyb= @(t) (y0(b)*exp(t)*(1-y0(b)))/((1+y0(b)*(exp(t)-1))^2);

% testing nodes
ntime_test=50;
time_test=linspace(0, T, ntime_test);

nspacetest=20;
space_test=linspace(a,b, nspacetest);

X_test=[];
for i=1:ntime_test
    for j=1:nspacetest
        X_test=[X_test, [space_test(j), time_test(i)]'];
    end
end

nspacetest0=100;
space_test0=linspace(a,b, nspacetest0);

X_test0=[];
for i=1:nspacetest0
    X_test0=[X_test0, [space_test0(i), 0]'];
end

% construction of omega and theta
omega = cell(1, N_layer + 1); % array di matrici dei pesi
theta = cell(1, N_layer + 1); % array di vettori di bias

% Inizializzazione dei pesi con Xavier
for i = 1:N_layer+1
    fileNameOmega = ['omega', num2str(i), '.xlsx'];
    omega{i} = readmatrix(fileNameOmega);
    fileNameTheta = ['theta', num2str(i), '.xlsx'];
    theta{i} = readmatrix(fileNameTheta);
end

error=GLossFunction(N_layer,omega,theta,X_test,X_test0,T,a,b,y0,ya,dyu);
disp("error= "+ error);
```

## APPENDIX D - Script GC Test (2 of 3)

```
% weights update
[errornew,thetanew,omeganew]=GCWeightsUpdate(N_layer,omega,theta, ...
    X_test,X_test0,y0,ya,dya,a,b,T,lrate);
disp("errornew= "+errornew);
if errornew<error
    error=errornew;
    for i = 1:N_layer+1
        fileNameOmega = ['omega', num2str(i), '.xlsx'];
        writematrix(omeganew{i}, fileNameOmega);
        fileNameTheta = ['theta', num2str(i), '.xlsx'];
        writematrix(theta{i}, fileNameTheta);
    end
    omega=omeganew;
    theta=thetanew;
end

% calculation of the outputs
u=cell(1,length(X_test0(:,1)));
for i=1:length(X_test0(:,1))
    X=GCNeuralNetwork(N_layer,omega,theta,X_test0(:,i));
    u{i}=X{N_layer+1};
end

y=[];
for i=1:nspatial
    x=X_test0(1,i);
    y=[y,ya(0)+(b-x)*(x-a)/(b-a)^2*u{i}];
end

sol=[];
for i=1:nspatial
    x=X_test0(1,i);
    sol=[sol,y0(x)];
end

figure;
plot(X_test0, y, '-r','DisplayName', 'Approximated solution');
hold on;
plot(X_test0, sol, '-b','DisplayName', 'exact solution');
xlabel('x (space)');
ylabel('y(x,0)');
title('Graph of y(t)');
legend;

% Comparison of the graph with the exact solution to the bound
Y=[];
for j=1:ntime_test
    for i=1:nspatial
        xtest=[space_test(i),time_test(j)]';
        X=GCNeuralNetwork(N_layer, omega, theta, xtest);
```

APPENDIX D - Script GC Test (3 of 3)

```
x=space_test(i);
t=time_test(j);
Y(i,j)=ya(t)+(b-x)*(x-a)/(b-a)^2*X{N_layer+1};
end
end

Y=Y';
figure;
[X,T]=meshgrid(space_test, time_test);
surf(X,T,Y);
xlabel('x (space)');
ylabel('t (time)');
zlabel('u(x,t)');
```

# Bibliography

- [1] Borzi A., Modelling with Ordinary Differential Equations A Comprehensive Approach, CRC Press.
- [2] Evans C. L., Partial Differential Equations, American Mathematical Society.
- [3] Faronius H. K., Solving Partial Differential Equations With Neural Networks,  
<https://uu.diva-portal.org/smash/record.jsf?pid=diva2%3A1746454&dswid=424>.
- [4] Friedlander F. G., Joshi M., Introduction to the Theory of Distributions, Cambridge University Press.
- [5] Goodfellow I., Bengio Y., Courville A., Deep Learning, The MIT Press.
- [6] Han J., Kamber M., Pei J., Data Mining: Concepts and Techniques, Morgan Kaufmann.
- [7] Mascali G., Note di Fisica Matematica Avanzata due.
- [8] Perko L., Differential Equations and Dynamical Systems, Springer
- [9] Pinkus A., Approximation theory of the MLP model in neural networks, Cambridge University Press.
- [10] Shewchuk J. R., An Introduction to the Conjugate Gradient Method Without the Agonizing Pain,  
<https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf>.
- [11] Van der Meer R., Oosterlee C., Borovykh A., (2021), Optimally weighted loss functions for solving PDEs with Neural Networks, Journal of Computational and Applied Mathematics,  
<https://arxiv.org/abs/2002.06269>.
- [12] Van der Meer R., Solving Partial Differential Equations with Neural Networks.
- [13] Vostrecov, B.A.K.M.A., Approximation of continuous functions by superpositions of plane waves (1961),  
<http://resolver.tudelft.nl/uuid:c77e1bcc-7212-4234-af34-6586b628ab1c>.
- [14] Yadav N., Yadav A., Kumar M., An Introduction to Neural Network Methods for Differential Equations, Springer.