

UNIVERSITY OF TRIESTE

Department of Engineering and Architecture



Master's Degree in  
Computer & Electronic Engineering

Shape optimization with discrete adjoint method  
applied to RBF-FD meshless method

November 15, 2024

Candidate  
**Kevin Marzio**

Supervisor  
**Prof. Andrea De Lorenzo**

Co-supervisors  
**Dr. Riccardo Zamolo**  
**Dr. Mauro Munerato**  
**Dr. Davide Miotti**

A.Y. 2023/2024

# Abstract

Computer Aided Engineering (CAE) systems are a type of software which aggregate numerical solvers for Partial Differential Equations (PDEs), 3D CAD (Computer Aided Design) software, optimization algorithms and all the other software required for a complete assessment of the performance of a given engineering design and for its optimization. One of their most important field of application is Computational Fluid Dynamic (CFD).

Nowadays the solvers integrated in CAE systems are typically based on Finite Element Method (FEM) or Finite Volume Method (FVM): their common ground is their need for domain discretization by means of a mesh. However the process of the mesh creation has several limitations which hinder the solver flexibility and its integration with the other software packages present in CAE systems. To overcome this limitation it is possible to use different solvers that does not rely on the creation of a mesh for solving the PDEs: these are called Meshless or Meshfree Methods (MMs).

The goal of this thesis has been the implementation of an optimizer for one and three-dimensional design optimization problems based on simple models described by means of Poisson equations and solved by the Radial Basis Function-Finite Difference (RBF-FD) meshless method. To do so we studied the Automatic Differentiation (AD) and the adjoint method, and we implemented them on a gradient-based optimization algorithm.

To asses the quality of the implemented method we applied it to a generic 1D optimization problem and to a specific type of 3D optimization problem. In the latter case we also compared the results obtained by our method to those obtained with the continuous adjoint.

# Sommario

I sistemi di Computer Aided Engineering (CAE) sono un tipo di software che aggrega risolutori numerici per equazioni differenziali alle derivate parziali (PDE), software CAD (Computer Aided Design) 3D, algoritmi di ottimizzazione e tutto il resto dei software necessari per una valutazione completa delle prestazioni di un determinato progetto ingegneristico e per la sua ottimizzazione. Uno dei loro campi di applicazione più importanti è la Fluidodinamica Computazionale (CFD).

Al giorno d'oggi, i risolutori di PDE integrati nei sistemi CAE si basano tipicamente sul metodo agli elementi finiti (FEM) o sul metodo a volumi finiti (FVM): ciò che li accomuna è la necessità di discretizzare il dominio attraverso una griglia (mesh-based). Tuttavia, il processo di creazione della griglia presenta diverse limitazioni che ostacolano la flessibilità del risolutore e la sua integrazione con gli altri pacchetti software presenti all'interno dei sistemi CAE. Per superare questa limitazione è possibile utilizzare risolutori di PDE che non si basano sulla creazione di una griglia: questi sono chiamati metodi senza griglia (meshless methods, MMs).

L'obiettivo di questa tesi è stato l'implementazione di un ottimizzatore per problemi di ottimizzazione progettuale mono e tridimensionali per semplici modelli descritti tramite equazioni di Poisson e risolti con il metodo senza griglia chiamato Radial Basis Function-Finite Difference (RBF-FD). A tal fine, abbiamo studiato la Differenziazione Automatica (AD) e il metodo dell'aggiunto, implementandoli all'interno di un algoritmo di ottimizzazione basato sul gradiente.

Per valutare la qualità del metodo implementato, lo abbiamo applicato a un problema di ottimizzazione generico in 1D e ad uno specifico problema di ottimizzazione in 3D. In quest'ultimo caso, abbiamo anche confrontato i risultati ottenuti dal nostro metodo con quelli ottenuti sfruttando l'aggiunto continuo.

# Contents

|  |           |
|--|-----------|
| <b>Abstract</b>  | <b>I</b>  |
| <b>Sommario</b>  | <b>II</b> |
| <b>Introduction</b>  | <b>V</b>  |
| <b>1 Meshless methods</b>  | <b>1</b>  |
| <b>2 RBF-FD method</b>   | <b>4</b>  |
| 2.1 Scattered Data Interpolation . . . . .                                 | 4         |
| 2.2 The Mairhuber-Curtis theorem . . . . .                                 | 6         |
| 2.3 Radial Basis Functions . . . . .                                       | 7         |
| 2.4 Polynomial augmentation . . . . .                                      | 8         |
| 2.5 Problem solution . . . . .   | 10        |
| 2.6 Radial Basis Function generated Finite Differences (RBF-FD) . . .      | 11        |
| 2.6.1 Finite difference method . . . . .                                   | 12        |
| 2.6.2 Mathematical formulation . . . . .                                   | 13        |
| 2.7 Radial Basis Functions-Hermite Finite Difference (RBF-HFD) . . .       | 18        |
| 2.7.1 Hermite Interpolation . . . . .                                      | 19        |
| 2.7.2 Mathematical formulation . . . . .                                   | 21        |
| <b>3 Automatic Differentiation (AD) and adjoint method in optimization</b> | <b>24</b> |
| 3.1 Automatic Differentiation (AD) . . . . .                               | 24        |
| 3.1.1 Forward mode . . . . .   | 26        |
| 3.1.2 Reverse mode . . . . .   | 28        |
| 3.2 Design Optimization . . . . .  | 31        |
| 3.2.1 Optimization problem formulation . . . . .                           | 32        |
| 3.2.2 Models and optimization problems . . . . .                           | 34        |
| 3.2.3 Optimization Algorithms . . . . .                                    | 36        |
| 3.2.4 Design optimization with RBF-FD models . . . . .                     | 38        |

## CONTENTS

---

|          |  |           |
|----------|--|-----------|
| 3.3      | Adjoint method in RBF-FD based Design Optimization . . . . . | 39        |
| 3.3.1    | Adjoint intro . . . . .                                      | 40        |
| 3.3.2    | 1D case . . . . .  | 42        |
| 3.3.3    | 3D case . . . . .  | 45        |
| <b>4</b> | <b>Results</b>   | <b>52</b> |
| 4.1      | Poisson Equation . . . . .                                   | 52        |
| 4.2      | 1D case . . . . .  | 53        |
| 4.3      | 3D case . . . . .  | 55        |
| 4.3.1    | Results . . . . .  | 58        |
|          | <b>Conclusion</b>  | <b>61</b> |
| <b>A</b> | <b>Stereolithography .stl files</b>                          | <b>63</b> |

# Introduction

The first research in the field of Computational Fluid Dynamic (CFD) began in the early 1950s. Subsequent to this initial research, a plethora of numerical methods have been devised for solving complex fluid dynamics phenomena modeled by Partial Differential Equations (PDEs), culminating in the contemporary state-of-the-art Finite Element Method (FEM) and Finite Volume Method (FVM). Despite the differences between the two methods they both require the creation of a mesh, i.e. the domain of interest, which could be the volume of an object, needs to be divided in small interconnected elements/cells.

The mesh generation phase should not be underestimated, indeed, even if generated with the help of a computer program, it is a complex task which has several drawbacks. First, it requires a huge amount of time, especially for complex geometries like those encountered in CFD simulations, and second, in most cases, it requires the intervention of a specialized operator. These facts are particularly penalizing in the case where these mesh-based solvers need to be integrated into Computer Aided Engineering (CAE) systems.

CAE systems are a type of software whose goal is to aggregate every software tool required for a complete assessment of the performance of a given engineering design and for its optimization. Typically a set of design parameters which define a specific geometry are fed to a CAE system which uses them to generate the associated geometry. After the geometry creation, always the CAE system, discretizes it by means of a mesh generator, which may or may not require operator intervention; the mesh is subsequently fed to a mesh based solver for the necessary simulation which is then provided to a cost or objective function. From the output provided by the cost function, the CAE system, close the loop by using an optimization algorithm which modifies the design parameters, and thus the design itself, in order to minimize/maximize the cost function. In the aforementioned workflow often mesh generation is the more time-consuming task especially when it requires the intervention of an operator. Automatic mesh generator exists and may be used, but they impose constraints between computational efficiency and accuracy.

In this scenario, in order to avoid all the issues related to mesh creation, Meshless or Meshfree methods (MMs), which date back to the 90s [1], can be employed.

These methods, as the name suggest, does not require the creation of a mesh, instead they just need a set of scattered nodes in the problem's domain. Beside the fact that domain discretization using nodes is faster and cheaper since mesh and human intervention are no longer required, and the existence of automatic node generators which deal with complex-shaped domain exists [2, 3] their main advantage is related to their integration in CAE systems: nodes can be easily moved in case the domain is affected by geometrical transformations like the ones performed by the optimization algorithms. Bearing this last advantage in mind it becomes clear how CAE systems can be greatly improved by replacing the meshing step with node generation and the mesh-based solver step with a meshless solver.

In this thesis we focus on the Radial Basis Function-Finite Difference (RBF-FD) meshless methods and we show how they can be integrated in a CAE system in order to solve design optimization problems. To do so we developed a gradient-based optimization algorithm with improvements based on:

- Automatic Differentiation (AD), a technique widely used in machine learning in particular for neural network training [4]; and
- adjoint method a technique used in various disciplines as CFD and optimal control theory [5].

The implemented optimization algorithm is then used to solve one and three dimensional simple design optimization problems.

The outline of this thesis is the following:

**Chapter 1** introduces the Meshless Methods (MMs) by providing a clear definition of these methods and a brief overview of their history. Their key aspects for solving Partial Differential Equations (PDEs) are also discussed, maintaining a general overview by refraining from delving too deeply into any single method.

**Chapter 2** studies the Radial Basis Function-Finite Difference (RBF-FD) and the Radial Basis Functions-Hermite Finite Difference (RBF-HFD) which are two examples of meshless solvers. The latter is an enhanced version of the former which makes use of an Hermite interpolant (which makes use of derivatives of functions in addition to the function values alone) to approximate the solution of a generic PDE. In particular RBF-HFD turns out to be particularly useful for specific PDE boundary conditions.

**Chapter 3** initially explain the different ways derivatives of functions can be computed focusing on Automatic Differentiation and its two main implementations: forward and reverse mode.

Subsequently design optimization problems are faced emphasizing their structure and how they are solved, with a particular emphasis on problems based on RBF-FD models and on a simple gradient-based optimization algorithm.

Lastly it is explained what the adjoint method consists of and how it can be applied to 1D and 3D RBF-FD based design optimization problems. In 1D case general problems are faced, instead, in 3D case, we restricted our implementation to problems dealing with flux integrals.

**Chapter 4** shows the result of our optimization method applied to simple 1D and 3D problems modeled by means of Poisson equations. To show this it is presented the value of the objective functions which vary throughout the optimization as expected and in particular in the three dimensional case we compared the results obtained by our (discrete) adjoint method to the results obtained by the continuous adjoint.

The code developed during the implementation of the optimization algorithm discussed in this thesis has been implemented using the **Julia** [6] programming language. **Julia** was born in 2012 as an open source project hosted on GitHub in order to mix the convenience of high-level dynamic languages such as **MATLAB**, **Octave** and **R** with the performance offered by **C** and **Fortran** languages.

This programming language is mainly intended for algorithm development and data analysis in applied math, engineering and scientific computing and many libraries have already been developed for this purposes. An example of these is **Zygote** which we used to compute gradients using Automatic Differentiation (AD). Regarding the adjoint method we couldn't find any libraries that implemented it, so we had to write our own code for it. We also implemented the RBF-FD solver for the 1D case whereas for the 3D case we employed the one already developed internally at ESTECO.



# Chapter 1

## Meshless methods

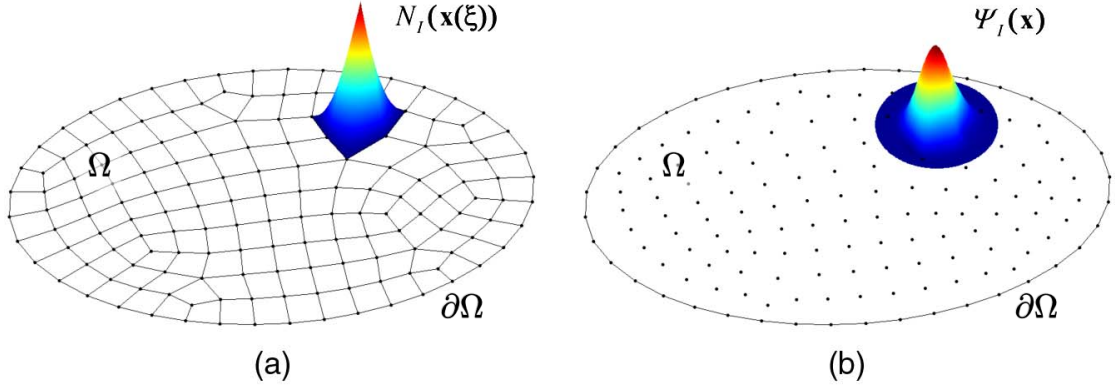
Meshless or Meshfree Methods (MMs) were developed to overcome the drawbacks of traditional mesh-based methods for the solution of Partial Differential Equations (PDE), where their true advantage is that “the approximation of unknowns in the PDE is constructed based on scattered points without mesh connectivity” [7]. The conceptual difference between Mesh Based and Meshless methods can be visualized on figure 1.1: the former patch the domain with some geometrical shapes while the latter only distributes nodes across the domain.

They appeared for the first time in 1977 with the Smooth Particle Hydrodynamics (SPH) method [8], initially used to modeling astrophysical phenomena such as exploding stars and dust clouds was later applied in solid mechanics to overcome limitations of mesh-based methods [9]. To improve accuracy, tensile instability and spatial instability of SPH, many more modern MMs have been developed: the introduction of Reproducing Kernel Particle Method (RKPM) [10] is a prime example of enhanced consistency and stability. Generalized Finite Difference (GFD) methods is another branch of numerical methods for solving PDEs that do not rely on a grid structure and many modern MMs originate from the employment of this approximation for solving PDEs.

The typical use case for MMs is the construction of an approximating field  $u^h$  for the sought solution  $u: \Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}$  of the following boundary value problem:

$$\begin{cases} \mathcal{L}u(\mathbf{x}) = f(\mathbf{x}) & \text{in } \Omega \\ \mathcal{B}u(\mathbf{x}) = g(\mathbf{x}) & \text{on } \partial\Omega \end{cases} \quad (1.0.1)$$

where  $\mathcal{L}$  is a generic linear differential operator and  $\mathcal{B}$  is a different linear operator used to enforce some Boundary Condition (BC) that does not necessarily involves partial derivatives;  $f$  and  $g$  are known functions. Usually in the Computational



**Figure 1.1:** (a) Example of reconstruction of a PDE solution, via mesh based method, for a single finite element domain; (b) approximation of the same function obtained through meshless approach. Figure taken from [7]

Fluid Mechanics (CFD) world the encountered BC are:

$$\text{Dirichlet BC:} \quad u = g \quad (1.0.2)$$

$$\text{Neumann BC:} \quad \frac{\partial u}{\partial \mathbf{n}} = g \quad (1.0.3)$$

$$\text{Robin BC:} \quad au + b \frac{\partial u}{\partial \mathbf{n}} = g \quad (1.0.4)$$

where  $\partial u / \partial \mathbf{n}$  indicate the normal derivative of  $u$  along the surface normal whereas  $a$  and  $b$  are known functions. It can also be noticed that Dirichlet and Neumann BCs are special cases of Robin BC respectively when  $b(\mathbf{x}) = 0$  and  $a(\mathbf{x}) = 0$ .

Regardless on the Meshfree approach, the following approximation for any solution  $u$  can be written:

$$u^h(\mathbf{x}) = \sum_{k=1}^N \alpha_k B_k(\mathbf{x}) \quad (1.0.5)$$

where  $B_k: \Omega \rightarrow \mathbb{R}$  are suitable basis functions and  $\alpha_k$  are the expansions coefficients that must be determined. Different choices for the basis functions  $B_k$  leads to different formulations and in literature can be found several of these, some examples are: RKPM [10], moving least square (MLS) [11], radial basis function (RBF) [12, 13] and partition of unity (PU) [14]. Furthermore solving the PDE in (1.0.1) with the approximated solution  $u^h$ , in general, yields to a non-zero error function  $\epsilon^h$  given by:

$$\epsilon^h(\mathbf{x}) = \mathcal{L}u^h(\mathbf{x}) - f(\mathbf{x}) \quad (1.0.6)$$

Once the general form of the approximated solution in (1.0.5) has been properly defined, it can be employed for discretizing the PDE, reported in (1.0.1), over a set

of  $N$  generated nodes  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  distributed in the physical domain  $\Omega \cup \partial\Omega$ . The Weighted Residual Method is used to do so: a set of test functions  $\{\Gamma_1, \dots, \Gamma_N\}$  orthogonal to  $\epsilon^h$  are used to integrate the error to zero:

$$\begin{aligned} \int_{\Omega} \Gamma_i \epsilon^h d\Omega &= \int_{\Omega} \Gamma_i (\mathcal{L}u^h(\mathbf{x}) - f(\mathbf{x})) d\Omega \\ &= \int_{\Omega} \Gamma_i \left[ \mathcal{L} \left( \sum_{k=1}^N \alpha_k B_k(\mathbf{x}) \right) - f(\mathbf{x}) \right] d\Omega = 0 \end{aligned} \tag{1.0.7}$$

where  $i = 1, \dots, N$ . From the choice of functions  $\Gamma_i$  the following two formulations are obtained [7]:

**Galerkin Meshless Methods** that find a weak solution for the PDE by using as test functions the basis functions  $B_k$ . These formulations require domain integration and special techniques to enforce boundary conditions;

**Collocation Meshless Methods** that find a strong solution for the PDE by using Dirac delta functions centered at the discretization nodes as test functions. Basically they enforce equations (1.0.1) on a finite set of nodes, and by doing so, they do not require any domain integration.

Finally, once constraints (1.0.7) are enforced, the solution of problem (1.0.1), i.e. the values  $\{u^h(\mathbf{x}_1), \dots, u^h(\mathbf{x}_N)\}$ , can be found solving the linear system resulting from PDE's discretization. We remark that the system obtained is not linear in general, but it is in this case since the PDE is linear. Nevertheless, also non-linear PDEs can be reduced to the aforementioned case thanks to a proper linearization.

# Chapter 2

## RBF-FD method

In this chapter we explain in more details the Meshless Method (MM) used within this work: the Radial Basis Function generated Finite Differences (RBF-FD) method. To do so we show how it can be used to solve the general boundary value problem defined in chapter 1. From now on, in order to be able to discretize the PDE, we consider to have a disposal a set of  $N$  distinct nodes,  $\mathcal{X}$ , defined as follow:

$$\mathcal{X} := \{ \mathbf{x}_1, \dots, \mathbf{x}_N \mid \mathbf{x}_i \in \Omega \cup \partial\Omega, i = 1, \dots, N \} \quad (2.0.1)$$

where  $\Omega \cup \partial\Omega \subset \mathbb{R}^d$  is the physical domain of the problem,  $\Omega$  and  $\partial\Omega$  indicate respectively its open subset and boundary, and  $d \in \mathbb{N}$  is its dimension.

### 2.1 Scattered Data Interpolation

Every MM, as we have already seen, on its core, is just a way to approximate the solution of a PDE and RBF-FD method is no exception; the tool used to do so is called *scattered data interpolation*. In this section we first see what scattered data interpolation is and then how it is related to PDEs solution. In the second part of the explanation we see how it is applied in general by MMs so to avoid to becoming fixated on a single implementation and losing generality; moreover this approach is still useful since it establishes all the steps that are also followed by RBF-FD. To avoid from the very beginning any kind of ambiguity we underline that scattered data interpolation is inherently connected to data fitting and not to PDEs approximation, consequently, it finds many other applications beyond the realm of MMs. Examples of its application on the field of image processing can be found in [15].

In general the interpolation problem has the following, simple, formulation. Given:

- a finite set of nodes,  $\mathcal{X} \subset \mathbb{R}^d$ , which can be the one reported in (2.0.1) and;

- a set of known real values  $u(\mathbf{x}_1), \dots, u(\mathbf{x}_N)$ , which may be obtained from a function

we want to find a continuous function  $u^h: \Omega \cup \partial\Omega \subset \mathbb{R}^d \rightarrow \mathbb{R}$  that satisfy:

$$u^h(\mathbf{x}_i) = u(\mathbf{x}_i) \quad \forall \mathbf{x}_i \in \mathcal{X} \quad (2.1.1)$$

If the locations, the nodes in  $\mathcal{X}$  where the measurements  $u(\mathbf{x}_1), \dots, u(\mathbf{x}_N)$  are taken, are placed on a uniform or regular grid we talk about interpolation, otherwise the process above is called *scattered* data interpolation. Here the main idea is to find a function  $u^h$  which is a “good” fit to the given data, where with “good” we mean a function that exactly match the given measurements at the corresponding locations.

MMs also aim to provide an approximation for an unknown function defined as a linear combination of a set of basis functions as reported in equation (1.0.5); we report here for clarity its definition:

$$u^h(x) = \sum_{k=1}^N \alpha_k B_k(\mathbf{x}) \quad (2.1.2)$$

and we remark that coefficients  $\alpha_k$  are unknown. Here is where the theory of scattered data interpolation is applied: it tells us that we are able to find the numerical values for  $\alpha_1, \dots, \alpha_N$  if we impose a number of conditions equal to the number of coefficients that we are looking for. These conditions must have a form like that shown in equation (2.1.1). Therefore if we replace the generic meshless approximation within each of the  $N$  interpolation conditions we can write the following linear system:

$$\underbrace{\begin{bmatrix} B_1(\mathbf{x}_1) & \dots & B_N(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ B_1(\mathbf{x}_N) & \dots & B_N(\mathbf{x}_N) \end{bmatrix}}_B \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_N \end{bmatrix} = \begin{bmatrix} u(\mathbf{x}_1) \\ \vdots \\ u(\mathbf{x}_N) \end{bmatrix} \quad (2.1.3)$$

that, once solved, provides us the desired coefficients that uniquely define  $u^h$  given the set of basis functions  $\{B_1, \dots, B_N\}$ . We would like to stress the fact that, however, the right hand side vector is made up of values of the unknown exact solution of problem (1.0.1).

During the solution of problem (1.0.1), via Collocation Methods, the following constraints arise instead:

$$\begin{aligned} \mathcal{L}u^h(\mathbf{x}_j) &= f(\mathbf{x}_j) & \text{if } \mathbf{x}_j \in \Omega \\ u^h(\mathbf{x}_j) &= g(\mathbf{x}_j) & \text{if } \mathbf{x}_j \in \partial\Omega \end{aligned} \quad (2.1.4)$$

and if we arrange the nodes such that the first  $N_I$  nodes belongs to  $\Omega$  and the last  $N_B$  to  $\partial\Omega$ , we can find the values of the approximated solution at the points in  $\mathcal{X}$  by solving:

$$\begin{bmatrix} c_{1,1} & \dots & c_{1,N_I} \\ \vdots & \ddots & \vdots \\ c_{N_I,1} & \dots & c_{N_I,N_I} \end{bmatrix} \begin{bmatrix} u^h(\mathbf{x}_1) \\ \vdots \\ u^h(\mathbf{x}_{N_I}) \end{bmatrix} = \mathbf{f} - \begin{bmatrix} c_{1,N_I+1} & \dots & c_{1,N_B} \\ \vdots & \ddots & \vdots \\ c_{N_I,N_I+1} & \dots & c_{N_I,N_B} \end{bmatrix} \mathbf{g} \quad (2.1.5)$$

where  $\mathbf{f} = [f(\mathbf{x}_1) \dots f(\mathbf{x}_{N_I})]^T$  and  $\mathbf{g} = [g(\mathbf{x}_{N_I+1}) \dots g(\mathbf{x}_{N_B})]^T$ , and the coefficient matrix  $\mathbf{C}$  is found as the solution of:

$$\begin{bmatrix} c_{1,1} & \dots & c_{1,N_I} \\ \vdots & \ddots & \vdots \\ c_{1,N} & \dots & c_{N_I,N} \end{bmatrix} = \mathbf{B}^{-T} \begin{bmatrix} \mathcal{L}B_1(\mathbf{x}_1) & \dots & \mathcal{L}B_1(\mathbf{x}_{N_I}) \\ \vdots & \ddots & \vdots \\ \mathcal{L}B_N(\mathbf{x}_1) & \dots & \mathcal{L}B_N(\mathbf{x}_{N_I}) \end{bmatrix} \quad (2.1.6)$$

Analyzing each row  $\mathbf{c}_i = [c_{i,1}, \dots, c_{i,N}]$  of matrix  $\mathbf{C}$  we can notice that are computed solving the following linear systems:

$$\mathbf{B}^T \mathbf{c}_i = \begin{bmatrix} \mathcal{L}B_1(\mathbf{x}_i) \\ \vdots \\ \mathcal{L}B_N(\mathbf{x}_i) \end{bmatrix} \quad i = 1, \dots, N_I \quad (2.1.7)$$

which are closely related to the ones obtained from scattered data interpolation reported in (2.1.3) due to the presence of the same matrix  $\mathbf{B}$ . We conclude by commenting that equation (2.1.6) with matrix  $\mathbf{B}$  defined as in (2.1.3) holds true only in case of Dirichlet boundary conditions, otherwise  $\mathbf{B}$  would take on a different form.

## 2.2 The Mairhuber-Curtis theorem

From the previous discussion, in particular from equation (2.1.6), it can be noticed that matrix  $\mathbf{B}$  has to be non singular in order to be able to solve the linear system associated to the boundary value problem, and this must hold for each node placement  $\mathcal{X}$  (to be read as every possible discretization of the problem domain) as long nodes are distinct. This property of  $\mathbf{B}$  turns out to be dependent on the choice of the particular set of basis functions: for example if we assume  $B_k(\mathbf{x}) \in \Pi_P^d$  and  $\{B_1(\mathbf{x}), \dots, B_N(\mathbf{x})\}$  to be a polynomial basis of the space  $\Pi_P^d$  of polynomials of degree at most  $P$  in  $\mathbb{R}^d$ , then we are not able to guarantee that  $\mathbf{B}$  is invertible for  $d > 1$ .

This issue is explained in more detail by the Mairhuber-Curtis theorem [16]; when dealing with the multidimensional case it is possible to continuously move

two nodes along a closed path  $P$ , that does not interfere with any other node in  $\mathcal{X}$ , such that they end up by interchanging their original positions without one crossing the path of the other. In the event that these two are the only nodes of  $\mathcal{X}$  that are moved,  $\mathbf{B}$  ends up with two rows exchanged leading to a change in the sign of its determinant, and, since the determinant is a continuous function, this means that there is a moment when the latter vanishes making the matrix singular.

The inconvenience arises from the fact that the set of basis functions is independent from the node position and could be solved by simply choosing a basis that is function of nodes position. By doing so we no more fall in the case of the Mairhuber-Curtis theorem since whenever we move nodes also the base itself changes and if two nodes switches their positions not only their respective rows in  $\mathbf{B}$  are switched, but also their columns, forcing the determinant not to change in sign.

## 2.3 Radial Basis Functions

Up to now, when talking about the approximated solution of the PDE,  $u^h$ , we have not specified the type of basis functions which define it. However these must be done in order to be able to find the coefficients  $\alpha_k$  in equation (2.1.2) and thus its numerical values; furthermore it would be appropriate to select a set of functions that allow the avoidance of the aforementioned Mairhuber-Curtis' theorem case. In this section we will define the ones used by the RBF-FD method: the Radial Basis Functions (RBFs).

RBFs are defined as:

$$\Phi(\mathbf{x}, \mathbf{x}_k) = \varphi(\|\mathbf{x} - \mathbf{x}_k\|_2) \quad (2.3.1)$$

where  $\mathbf{x}_k \in \mathbb{R}^d$  is a given and known point,  $\|\cdot\|_2$  is the euclidean distance and  $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ , named *basic* function, is a (univariate) function which takes the radius  $r_k = \|\mathbf{x} - \mathbf{x}_k\|_2$  as input and it is used as generator for all the (multivariate) *basis* functions associated to different  $\mathbf{x}_k$ .

Sometimes the notations  $\Phi(\mathbf{x} - \mathbf{x}_k)$  or  $\Phi_k(\mathbf{x})$  are also used instead of  $\Phi(\cdot, \mathbf{x}_k)$  where “.” act as a placeholder for the corresponding argument meaning that varies freely. However the 2-arguments notation makes clear that  $\Phi(\cdot, \cdot)$  is a real valued function defined on the space given by the Cartesian product  $\Omega \times \Omega$ . This type of functions, in accordance with the definition given in [17], are called *kernels*. In general at different basic functions  $\varphi$  are associated different kernels; some of these are reported in table 2.1.

To further clarify the RBFs name we can notice that they are called:

**Radial** since the value of each  $\Phi(\cdot, \mathbf{x}_k)$  at each point  $\mathbf{x}$  depends only on the

**Table 2.1:** Examples of basic functions where  $r$  is a real number greater than or equal to zero, and  $\epsilon$ , called shape factor, is a suitable parameter

| Name                   | $\varphi(r)$                       |
|------------------------|------------------------------------|
| Multiquadratic         | $\sqrt{1 + (\epsilon r)^2}$        |
| Inverse multiquadratic | $(\sqrt{1 + (\epsilon r)^2})^{-1}$ |
| Thin plate splines     | $r^{2l} \log l, l \in \mathbb{N}$  |
| Gaussian               | $e^{-(\epsilon r)^2}$              |
| Polyharmonics          | $r^{2l+1}, l \in \mathbb{N}$       |

distance between that point and  $\mathbf{x}_k$  through  $\|\cdot\|_2$ . Basically they satisfy radial symmetry, i.e.  $\Phi(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_j, \mathbf{x}_i)$  for any  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$ ;

**Basis** since in case of a set of  $N$  distinct nodes in  $\mathbb{R}^d$ , as can be  $\mathcal{X}$  given in (2.0.1), the set of radial functions  $\Phi_k(\mathbf{x})$  with  $k = 1, \dots, N$  associated to each point of the set, form a basis for the space of functions:

$$F_\Phi := \left\{ \sum_{k=1}^N \alpha_k \Phi_k(\mathbf{x}), \quad \alpha_k \in \mathbb{R}, \mathbf{x}_k \in \mathcal{X} \right\}$$

In case of RBF-FD method the implementation of scattered data interpolation and the solution of governing equation remain the same as explained in the previous section and leads to a symmetric matrix  $\mathbf{B}$  defined as:

$$\mathbf{B} = \begin{bmatrix} \Phi(\mathbf{x}_1, \mathbf{x}_1) & \dots & \Phi(\mathbf{x}_1, \mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \Phi(\mathbf{x}_N, \mathbf{x}_1) & \dots & \Phi(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix} \quad (2.3.2)$$

We conclude by observing that this functions are particularly convenient since they depend on nodes position through  $\mathbf{x}_k$ , thus they allow to avoid the non-invertibility of matrix  $\mathbf{B}$  in case of singular node arrangements (Mairhuber-Curtis theorem).

## 2.4 Polynomial augmentation

Setting aside the issue of the invertibility of matrix  $\mathbf{B}$  in case of particular nodes arrangement discussed in previous subsection, we should also take into account the accuracy of the interpolation that we are able to achieve, which also depends upon the type of functions that we are supposed to approximate. Indeed



RBFs approximation schemes alone are not able to exactly interpolate (i.e. with an accuracy only depending on round-off errors) constant, linear or higher degree polynomials fields. This is an issue in different important engineering applications such as modeling of constant strain in elastic bodies and steady temperature fields in differentially heated walls [18].

To overcome this limitation, a polynomial augmentation of degree  $P$  is required, leading to the overall formulation for the RBF interpolant:

$$u^h(\mathbf{x}) = \sum_{j=1}^N \alpha_j \Phi_j(\mathbf{x}) + \sum_{k=1}^M \beta_k p_k(\mathbf{x}) \quad (2.4.1)$$

where  $M = \binom{P+D}{D}$  is the number of polynomial basis functions with degree  $P \leq D$ ,  $\{p_1(\mathbf{x}) \dots p_M(\mathbf{x})\}$  is a complete polynomial basis of  $\Pi_P^D$  and  $\beta_j$  are the corresponding coefficients. An example of polynomial basis for polynomials of degree  $P = 1$  in  $2D$  has the following  $M = 3$  elements:  $p_1(x, y) = 1$ ,  $p_2(x, y) = x$ ,  $p_3(x, y) = y$ .

We must also note that using an interpolant with the introduced polynomial augmentation leads to an underdetermined system in (2.1.3). In order to obtain a square  $\mathbf{B}$  and thus having a solvable system, the following orthogonality conditions have to be imposed:

$$\sum_{i=1}^N \alpha_i p_k(\mathbf{x}_i) = 0, \quad k = 1, \dots, M \quad (2.4.2)$$

The coefficients of  $u^h$ , which are now composed not only by  $\boldsymbol{\alpha} = [\alpha_1 \dots \alpha_N]$ , but also by  $\boldsymbol{\beta} = [\beta_1 \dots \beta_M]$ , can then be found by solving the following system:

$$\underbrace{\begin{bmatrix} \mathbf{B} & \mathbf{P} \\ \mathbf{P}^T & \mathbf{0} \end{bmatrix}}_{\mathbf{M}} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{0} \end{bmatrix} \quad (2.4.3)$$

where:

$$\mathbf{P} = \begin{bmatrix} p_0(\mathbf{x}_1) & \dots & p_M(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ p_0(\mathbf{x}_N) & \dots & p_M(\mathbf{x}_N) \end{bmatrix} \quad (2.4.4)$$

$$\mathbf{u} = [u(\mathbf{x}_1) \dots u(\mathbf{x}_N)]$$

From its formulation is easy to understand that the system above is simply the composition of the system in equation (2.1.3), in the first row, with constraints (2.4.2) written in compact form, in the second row.

In practice the addition of polynomial basis to the RBF interpolant let us perfectly fit  $u(\mathbf{x})$  not only on collocation points where we have  $u^h(\mathbf{x}_i) = u(\mathbf{x}_i)$ , but

also across the rest of the domain provided that data  $u(\mathbf{x}_1), \dots, u(\mathbf{x}_N)$  come from a polynomial of total degree less than or equal to  $P$ . Nevertheless this procedure has a side effect: not all set of nodes  $\mathcal{X}$  nor all basic functions  $\varphi$  leads to a well-posed RBF interpolation with a non singular matrix  $M$ . We will discuss this limitation in more detail in the next section.

## 2.5 Problem solution

At the beginning of section 2.2 we mentioned that the system in equation (2.1.6) might not be solvable in case of a combination of non-point-dependent basis functions and particular nodes arrangements, but, up to now, we did not discuss in general in which cases the interpolation problem is solvable. In this section we will address this shortcoming.

We start by recalling that, in case of pure RBF interpolant, the system that we aim to solve has the following compact form:

$$\mathbf{B}\boldsymbol{\alpha} = \mathbf{u} \quad (2.5.1)$$

where  $\mathbf{B}$  is a symmetric matrix. A positive definite  $\mathbf{B}$  would be sufficient in order to solve the above system and this last property depends on the choice of the basic function  $\varphi$  used to define the RBFs. By definition only *strictly* positive definite [19]  $\varphi$  are associated to a positive definite  $\mathbf{B}$  and, this, restrict our choices: of those shown in table 2.1 only Inverse Multiquadratic and Gaussian functions satisfy this requirement. The same attributes of  $\varphi$  are also inherited by the associated RBFs  $\Phi_k$ .

However in previous section we have seen that, beyond the solvability issue, a polynomial augmentation of degree  $P$  is beneficial for the accuracy of the interpolant  $u^h$ . This means that the system that we have to solve, in general, is no more in the form shown in equation (2.5.1), but rather in the following:

$$\mathbf{M} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{0} \end{bmatrix} \quad (2.5.2)$$

In this case the matrix we seek to be positive definite would be  $\mathbf{M}$  that we recall being defined as:

$$\mathbf{M} = \begin{bmatrix} \mathbf{B} & \mathbf{P} \\ \mathbf{P}^T & \mathbf{0} \end{bmatrix} \quad (2.5.3)$$

which is symmetric as well. To be such the following conditions have to be met [20]:

1. basic function  $\varphi$  has to be *strictly conditionally* positive definite function of order  $P$  [19];

2. matrix  $\mathbf{P}$  has to be full-rank.

The first condition allows a greater freedom on the choice of the basic function, compared to the pure RBF interpolant: in fact, the set of strictly conditionally positive definite functions of order  $P$  is a superset of strictly positive definite functions. These functions are defined as those that require a polynomial augmentation of order at least  $P - 1$  in order to give a non singular  $\mathbf{M}$ . Strictly conditionally positive definite functions of order  $P$  are also strictly conditionally positive definite of any higher order. This means that, in case  $u^h$  include a polynomial augmentation of order 1, we can also use Multiquadratic, Thin plate splines with  $l = 0$  and Polyharmonics with  $l = 1$  (refer again to table 2.1 for their definitions) as basic functions. At this point someone might consider increasing the degree  $P$  of the polynomial augmentation up to the theoretical limit for the size of matrix  $\mathbf{P}$ , i.e.  $M = N$ , in order to increase the accuracy of the interpolant. However doing so results in ill-conditioning and singularity issues related to  $\mathbf{P}$  which in turn will have an impact on the singularity of matrix  $\mathbf{M}$ .

Therefore, since we require the non singularity of  $\mathbf{M}$ , the second condition has to be satisfied. It can be shown that to have a full-rank  $\mathbf{P}$  the set  $\mathcal{X}$ , containing the nodes respect to which we are carrying out the interpolation, must be  $P$ -unisolvent [19], where  $P$  is the degree of the polynomial augmentation. This dependency of  $\mathbf{P}$ 's rank on the node locations should not surprise since the matrix columns consist on the elements of the polynomial basis evaluated at the different points in  $\mathcal{X}$ , and they are required to be linearly independent.

Given the node generation technique used in this work, a safe rule for a stable implementation (in the sense of a well-posed interpolation problem) of the polynomial augmentation is to respect the inequality  $2M \leq N$  where  $M$  is the number of terms in the polynomial basis and  $N$  is the number of nodes in  $\mathcal{X}$ .

## 2.6 Radial Basis Function generated Finite Differences (RBF-FD)

In this section we are going to discuss in detail how Radial Basis Function generated Finite Differences (RBF-FD) meshless method is used to solve Partial Differential Equations (PDEs). Its first implementation was developed by Kansa in [12, 13]. His approach, that we denote as *global method*, due to its computational inefficiencies (that we are going to reference later on this section), was eventually dropped in preference for the one suggested by Tolstykh in [21]. We will denote the latter as *local method*. In recent years RBF-FD local method has been developed and applied with success [22, 23, 24, 25].

The goal of both methods is to approximate solutions to PDEs, i.e., to find a function (or some discrete approximation to this function) which satisfies a given relationship between various of its derivatives on some given region of space along with some boundary conditions on the boundary of this domain. In most cases an analytical solution cannot be found. What RBF-FD methods do is replacing derivatives in the differential equation by Finite Difference (FD) approximations in such a way as to obtain a large algebraic system of equations to be solved in place of the differential equation; this could be easily solved with a computer.

### 2.6.1 Finite difference method

Before tackling the approximation of PDEs solution, we first consider the more basic task of approximating the derivatives of a known function by Finite Difference (FD) formulas based only on values of the function itself at discrete points. Given  $u$ , in the simplest case a function of one variable assumed to be sufficiently smooth, we want to approximate its derivatives at a given point  $\bar{x}$  relying solely on its values at a finite number of points close to  $\bar{x}$ . In general its  $k$ -th derivative is approximated by the following FD formula:

$$\left. \frac{d^k u}{dx^k} \right|_{x=\bar{x}} = u^{(k)}(\bar{x}) \approx \sum_{i=1}^n c_i^k u(x_i) \quad (2.6.1)$$

where  $u(x_1), \dots, u(x_n)$  are the function's samples and  $c_1^k, \dots, c_n^k$ , which can be computed in different ways such as the method of undetermined coefficients or via polynomial interpolation [26], are called FD weights.

To give a concrete example we could approximate  $u'(\bar{x})$  with the following one-sided approximations:

$$D_+ u(\bar{x}) = \frac{u(\bar{x} + h) - u(\bar{x})}{h} \quad (2.6.2a)$$

$$D_- u(\bar{x}) = \frac{u(\bar{x}) - u(\bar{x} - h)}{h} \quad (2.6.2b)$$

for some value of  $h$ . This is motivated by the standard definition of the derivative as the limiting value of this expression as  $h \rightarrow 0$ . In these cases the same FD weights,  $\{1/h, -1/h\}$ , are associated to function values coming from different samples:  $\{u(\bar{x} + h), u(\bar{x})\}$  for equation (2.6.2a) and  $\{u(\bar{x}), u(\bar{x} - h)\}$  for (2.6.2b). Another possibility is to use the centered approximation:

$$D_0 u(\bar{x}) = \frac{u(\bar{x} + h) - u(\bar{x} - h)}{h} = \frac{1}{2}(D_+ u(\bar{x}) + D_- u(\bar{x})) \quad (2.6.3)$$

To derive approximations to higher order derivatives, besides the two method mentioned above, is also possible to repeatedly apply first order differences. Just

as the second order derivatives is the derivative of  $u'$ , we can view  $D^2u(\bar{x})$ , the second order derivative approximant, as being a finite difference of first differences:  $D^2u(\bar{x}) = D_+D_-u(\bar{x})$  or  $D^2u(\bar{x}) = D_-D_+u(\bar{x})$ . If we use a step size  $h/2$  in each centered approximation to the first derivative we could also define  $D^2u(\bar{x})$  as a centered difference of centered differences and obtain:

$$D^2u(\bar{x}) = \frac{1}{h} \left( \left( \frac{u(\bar{x}+h) - u(\bar{x})}{h} \right) - \left( \frac{u(\bar{x}) - u(\bar{x}-h)}{h} \right) \right) \quad (2.6.4)$$

where FD weights  $\{1/h^2, -2/h^2, 1/h^2\}$  are associated to  $\{u(\bar{x}+h), u(\bar{x}), u(\bar{x}-h)\}$ .

### 2.6.2 Mathematical formulation

After this very brief introduction to finite difference (FD), we can see how they are generalized to PDEs differential operators by RBF-FD in order to obtain easily solvable linear systems.

We start by recalling that the boundary value problem that we have to solve is defined as:

$$\begin{cases} \mathcal{L}u = f & \text{in } \Omega \\ \mathcal{B}u = g & \text{on } \partial\Omega \end{cases} \quad (2.6.5)$$

where  $\mathcal{L}$  and  $\mathcal{B}$  are linear operators and its solution  $u$  could be approximated by a function  $u^h$  defined as reported in equation (2.4.1):

$$u^h(\mathbf{x}) = \sum_{j=1}^N \alpha_j \Phi_j(\mathbf{x}) + \sum_{k=1}^M \beta_k p_k(\mathbf{x}) \quad (2.6.6)$$

which holds true all over the domain. Coefficients of the expansions are found by solving equation (2.4.3):

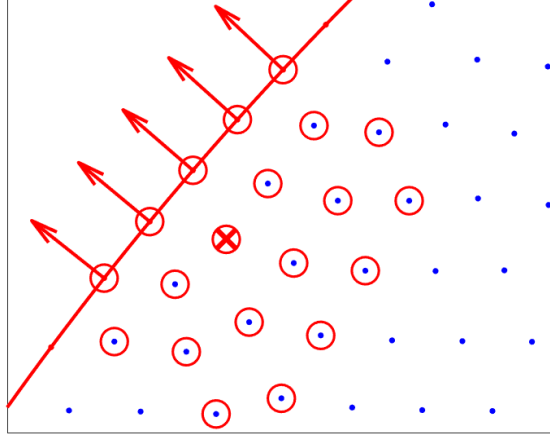
$$\underbrace{\begin{bmatrix} \mathbf{B} & \mathbf{P} \\ \mathbf{P}^T & \mathbf{0} \end{bmatrix}}_M \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{u} \\ \mathbf{0} \end{bmatrix} \quad (2.6.7)$$

In order to streamline the explanation of the procedure we split the set of nodes  $\mathcal{X}$  reported in (2.0.1) into the two following sets:

$$\mathcal{X}_I := \{ \mathbf{x}_1, \dots, \mathbf{x}_{N_I} \mid \mathbf{x}_i \in \Omega, i = 1, \dots, N_I \} \quad (2.6.8a)$$

$$\mathcal{X}_B := \{ \mathbf{x}_{N_{I+1}}, \dots, \mathbf{x}_{N_B} \mid \mathbf{x}_i \in \partial\Omega, i = N_{I+1}, \dots, N_B \} \quad (2.6.8b)$$

where  $N_I$  and  $N_B$  indicate respectively the number of nodes inside and on the boundary of the physical domain (from the definitions of the two sets it also follows that  $N_I + N_B = N$  and  $\mathcal{X}_I \cup \mathcal{X}_B = \mathcal{X}$ ).



**Figure 2.1:** Example of 2D stencil. The nodes which belong to it are marked with a red circle, its central node with a red cross. Red arrows represent the boundary normals  $\mathbf{n}$  of those nodes belonging to  $\partial\Omega$ . Figure taken from [20]

The idea that Kansa used to discretize the problem reported above is to use the theory of interpolation to approximate the differential operator  $\mathcal{L}$  in the PDE, with an operator  $\mathcal{L}^h$  represented by a matrix  $\mathbf{C}_I$ , using a FD-like method. Tolstykh's approach follows the same concept as Kansa's, but adds *stencil* as novelty.

A stencil of  $m$  nodes associated to each node  $\mathbf{x}_i \in \mathcal{X}_I$  is defined as the set  $\mathcal{X}_i = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$  formed by the  $m$  nearest neighbors of  $\mathbf{x}_i$ ; in addition  $\mathcal{X}_i$  can be interpreted as the union of  $\mathcal{X}_{i,I} = \{\mathbf{x}_1, \dots, \mathbf{x}_{m_I}\}$  and  $\mathcal{X}_{i,B} = \{\mathbf{x}_{m_I+1}, \dots, \mathbf{x}_m\}$ , respectively the set of its  $m_I$  nodes belonging to  $\Omega$  and of its other  $m_B$  nodes belonging to  $\partial\Omega$ . An example of stencil could be found in Figure 2.1.

In Tolstykh's local method the interpolation scheme is local, i.e.  $u^h$  is expanded using a basis that changes depending on the position  $\mathbf{x}$  and it is made valid only inside the stencil centered at  $\mathbf{x}$  rather than globally across the entire domain. Thus, given a point  $\mathbf{x}_i \in \mathcal{X}_I$  along its related stencil  $\mathcal{X}_i = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ , equation (2.6.6) is rewritten as:

$$u^h(\mathbf{x}_i) = \sum_{j=1}^m \alpha_j \Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{k=1}^M \beta_k p_k(\mathbf{x}_i) \quad (2.6.9)$$

In the following, we will focus solely on Tolstykh's formulation, commonly employed in practice, with a brief mention of Kansa's for comparative purposes.

Applying operator  $\mathcal{L}$  to the definition of  $u^h$  results in:

$$\begin{aligned}\mathcal{L}u^h(\mathbf{x}_i) &= \sum_{j=1}^m \alpha_j \mathcal{L}\Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{k=1}^M \beta_k \mathcal{L}p_k(\mathbf{x}_i) \\ &= [\boldsymbol{\alpha} \quad \boldsymbol{\beta}] \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix}\end{aligned}\tag{2.6.10}$$

where different vectors of coefficients  $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_m] \in \mathbb{R}^m$  and  $\boldsymbol{\beta} = [\beta_1, \dots, \beta_M] \in \mathbb{R}^M$  are associated to different  $\mathbf{x}_i$ . In order to find them we can note that, since equation (2.6.9) is still an approximated solution of the boundary value problem (at least locally), it must satisfy:

$$u^h(\mathbf{x}_i) = u(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{i,I} \tag{2.6.11a}$$

$$\mathcal{B}u^h(\mathbf{x}_i) = g(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{i,B} \tag{2.6.11b}$$

which rewritten in matrix form read as:

$$\underbrace{\begin{bmatrix} \boldsymbol{\Phi}_I & \mathbf{P}_I \\ \mathcal{B}\boldsymbol{\Phi}_B & \mathcal{B}\mathbf{P}_B \\ \mathbf{P}^T & \mathbf{0} \end{bmatrix}}_{\mathbf{M}_{BC}} \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_I \\ \mathbf{g} \\ \mathbf{0} \end{bmatrix} \tag{2.6.12}$$

where  $\mathbf{u}_I = [u(\mathbf{x}_1), \dots, u(\mathbf{x}_{m_I})]^T \in \mathbb{R}^{m_I}$  and  $\mathbf{g} = [g(\mathbf{x}_{m_I+1}), \dots, g(\mathbf{x}_m)]^T \in \mathbb{R}^{m_B}$  and the new terms in  $\mathbf{M}_{BC} \in \mathbb{R}^{(m+M) \times (m+M)}$  are defined as follow:

$$\begin{aligned}\boldsymbol{\Phi}_I &= \begin{bmatrix} \Phi(\mathbf{x}_1, \mathbf{x}_1) & \dots & \Phi(\mathbf{x}_1, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \Phi(\mathbf{x}_{m_I}, \mathbf{x}_1) & \dots & \Phi(\mathbf{x}_{m_I}, \mathbf{x}_m) \end{bmatrix} \in \mathbb{R}^{m_I \times m} \\ \mathbf{P}_I &= \begin{bmatrix} p_1(\mathbf{x}_1) & \dots & p_M(\mathbf{x}_{m_I}) \\ \vdots & \ddots & \vdots \\ p_1(\mathbf{x}_1) & \dots & p_M(\mathbf{x}_{m_I}) \end{bmatrix} \in \mathbb{R}^{m_I \times M} \\ \mathcal{B}\boldsymbol{\Phi}_B &= \begin{bmatrix} \mathcal{B}\Phi(\mathbf{x}_{m_I+1}, \mathbf{x}_1) & \dots & \mathcal{B}\Phi(\mathbf{x}_{m_I+1}, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \mathcal{B}\Phi(\mathbf{x}_m, \mathbf{x}_1) & \dots & \mathcal{B}\Phi(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix} \in \mathbb{R}^{m_B \times m} \\ \mathcal{B}\mathbf{P}_B &= \begin{bmatrix} \mathcal{B}p_1(\mathbf{x}_{m_I+1}) & \dots & \mathcal{B}p_M(\mathbf{x}_{m_I+1}) \\ \vdots & \ddots & \vdots \\ \mathcal{B}p_1(\mathbf{x}_m) & \dots & \mathcal{B}p_M(\mathbf{x}_m) \end{bmatrix} \in \mathbb{R}^{m_B \times M}\end{aligned}\tag{2.6.13}$$

When Dirichlet BC are enforced,  $\mathcal{B}$  becomes the identity operator in which case  $\mathbf{M}_{BC}$  in equation (2.6.12) coincides with the matrix  $\mathbf{M}$  of the interpolation system (2.4.3) defined in section 2.4. In this case the problem is well-posed, regardless of the shape of domain  $\Omega$  if the instructions presented in section 2.5 are adhered to. However this is no-longer true in case of Neumann or Robin BCs: we are going to tackle this issue in the next section. For the time being, to proceed with the explanation of RBF-FD, we will assume that the problem is always well-posed regardless of the type of BCs applied.

Substituting the recently determined  $\alpha$  and  $\beta$  into equation (2.6.10) we obtain:

$$\mathcal{L}u^h(\mathbf{x}_i) = [\mathbf{u}_I \quad \mathbf{g} \quad \mathbf{0}] \mathbf{M}_{BC}^{-T} \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (2.6.14)$$

We might now observe that the last two factors of this last equation are known, thus we can represent their product with a vector  $\mathbf{c}(\mathbf{x}_i) = [\mathbf{c}_I(\mathbf{x}_i), \mathbf{c}_B(\mathbf{x}_i), \mathbf{c}_p(\mathbf{x}_i)] \in \mathbb{R}^{m+M}$  that can be obtained solving the dual problem:

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (2.6.15)$$

where  $\mathbf{c}_I(\mathbf{x}_i)$ ,  $\mathbf{c}_B(\mathbf{x}_i)$  and  $\mathbf{c}_p(\mathbf{x}_i)$  denote the first  $m_I$ , the next  $m_B$  and the last  $M$  elements of  $\mathbf{c}(\mathbf{x}_i)$  respectively. Once this is done (2.6.14) can be rewritten as:

$$\begin{aligned} \mathcal{L}u^h(\mathbf{x}_i) &= \mathbf{c}_I(\mathbf{x}_i)^T \mathbf{u}_I + \mathbf{c}_B(\mathbf{x}_i)^T \mathbf{g} \\ &= \sum_{j=1}^{m_I} c_j(\mathbf{x}_i) u(\mathbf{x}_j) + \sum_{k=m_I+1}^m c_k(\mathbf{x}_i) g(\mathbf{x}_k) \end{aligned} \quad (2.6.16)$$

which is nothing else but the FD-like approximation of  $\mathcal{L}u^h$  at point  $\mathbf{x}_i$  with FD weights given by the elements of the vectors  $\mathbf{c}_I(\mathbf{x}_i)$  and  $\mathbf{c}_P(\mathbf{x}_i)$ . Finally we remark that, in this specific case, we can write equation (2.6.16) thanks to the combination of:

- the linearity of  $u^h$  respect the parameters  $\alpha_1, \dots, \alpha_m$  and  $\beta_1, \dots, \beta_M$ , as can be seen in equation (2.6.9), and;
- the linearity of operator  $\mathcal{L}$ .

The values of the approximated solution at the  $N_I$  nodes of  $\mathcal{X}$ , i.e. the solution of the PDE, can be finally found by requiring  $u^h$  to approximate the exact solution at each of these points:

$$\mathcal{L}u^h(\mathbf{x}_i) = \mathcal{L}u(\mathbf{x}_i) = f(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{i,I} \quad (2.6.17)$$



which takes the matrix form:

$$\mathbf{C}_I \begin{bmatrix} u^h(\mathbf{x}_1) \\ \vdots \\ u^h(\mathbf{x}_{N_I}) \end{bmatrix} = \begin{bmatrix} f(\mathbf{x}_1) \\ \vdots \\ f(\mathbf{x}_{N_I}) \end{bmatrix} - \mathbf{C}_B \begin{bmatrix} g(\mathbf{x}_{N_I+1}) \\ \vdots \\ g(\mathbf{x}_N) \end{bmatrix} \quad (2.6.18)$$

where rows of matrices  $\mathbf{C}_I$  and  $\mathbf{C}_B$  are formed by the elements of the vectors  $\mathbf{c}_I(\mathbf{x}_i)$  and  $\mathbf{c}_B(\mathbf{x}_i)$  found by solving equation (2.6.15)  $N_I$  times:

$$\mathbf{C}_I = \begin{bmatrix} c_1(\mathbf{x}_1) & \dots & c_{N_I}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ c_1(\mathbf{x}_{N_I}) & \dots & c_{N_I}(\mathbf{x}_{N_I}) \end{bmatrix} \in \mathbb{R}^{N_I \times N_I} \quad (2.6.19)$$

$$\mathbf{C}_B = \begin{bmatrix} c_{N_I+1}(\mathbf{x}_1) & \dots & c_N(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ c_{N_I+1}(\mathbf{x}_{N_I}) & \dots & c_N(\mathbf{x}_{N_I}) \end{bmatrix} \in \mathbb{R}^{N_I \times N_B}$$

Equation (2.6.18) can also be rewritten more compactly as

$$\mathbf{C}_I \mathbf{u} = \mathbf{f} - \mathbf{C}_B \mathbf{u}_B \quad (2.6.20)$$

What we just obtained in (2.6.18) is the algebraic system that once solved, in place of the PDE, gives the values of the PDE approximated solution  $\{u^h(\mathbf{x}_1), \dots, u^h(\mathbf{x}_{N_I})\}$ . In Kansa's global method the process is followed with the same steps as long as the global expansion of equation (2.6.6) is used to approximate the PDE solution instead of the local expansion of equation (2.6.9).

An important remark has to be done about computational efficiencies of global and local methods since linear system (2.6.15) has to be solved at any point  $\mathbf{x}_i \in \mathcal{X}_I$ , therefore requiring the inversion of  $N_I$  different  $\mathbf{M}_{BC}^T$  matrices. In case of local method each  $\mathbf{M}_{BC}$  has size  $(m + M) \times (m + M)$ ; instead in global method they have size  $(N + M) \times (N + M)$  since conditions (2.6.11) must be enforced not only in the stencil, but over the entire domain due to the global nature of the approximated solution. Remembering that the computational cost of factorizing a matrix is at least of  $\mathcal{O}(n^3)$ , where here  $n$  denotes the number of rows and columns of the matrix, then local method is computationally advantageous respect the global one.

Another difference between global and local method is that in global method derivative approximation of  $u^h$  at point  $\mathbf{x}_i$  depends on the values of  $u$  across the entire domain even if derivatives are local properties of functions: this is clearly suboptimal since leads to a matrix  $\mathbf{C}_I$ , which represent the discretized differential operator  $\mathcal{L}$ , that is full even though it should not be in principle. On the contrary,

in local method,  $\mathbf{C}_I$  is sparse since in each row  $i$  the number of non-zero entries is equal to the number of internal nodes of stencil  $\mathcal{X}_I$ .

Finally, indicating with  $\mathbf{u}^h = [u(\mathbf{x}_1), \dots, u(\mathbf{x}_N)]$  the vector of the exact solution evaluated at the RBF-FD nodes and with  $\mathbf{u} = [u^h(\mathbf{x}_1), \dots, u^h(\mathbf{x}_N)]$  the approximated solution obtained from RBF-FD method, is possible to show that the solution error  $e(N) = \frac{\|\mathbf{u}^h - \mathbf{u}\|}{\|\mathbf{u}\|}$  decrease exponentially as the number of nodes  $N$  increase in the global approach. In the local approach, conversely, the solution error decreases “only” polynomially.

The reason just listed above are those that explain the benefits of stencil introduction and that led to the choice of Tolstykh approach both in this work and within the community of researchers exploring different Meshless Methods.

What we have not covered yet is the computational cost of the stencils creation: if too expensive, it might loose all the benefits of the local approach. In fact this would be the case when using a brute force approach (which require a computational cost of  $\mathcal{O}(N^2)$  for each  $\mathbf{x}_i$ ) where all pairwise distances between nodes are computed and then sorted in order to keep just the  $m$ -nearest neighbors. To avoid this issue more efficient algorithms has to be used. An example of these is the  $k$ -d tree algorithm [27] which require  $\mathcal{O}(N \log N)$  operations for rearranging the nodes and other  $\mathcal{O}(N \log N)$  for finding a fixed number of neighbors.

Finally, regarding to the stencil, we note that the effect of its size are still under research: one of the last activities on this regard is the one of Kolar-Požun et al. [28] where they found that changing stencil’s size induce oscillations in both the solution and discretization errors for the Poisson equation.

## 2.7 Radial Basis Functions-Hermite Finite Difference (RBF-HFD)

In previous section we have overlooked a rather important fact for the implementation of the RBF-FD method: the invertibility of matrix  $\mathbf{M}_{BC}$ . In fact is always possible, given a stencil containing boundary points on which Neumann or Robin Boundary Conditions (BCs) are applied, to find normal directions associated with them that make the local matrix  $\mathbf{M}_{BC}$  in equation (2.6.12) singular [20].

Unfortunately, simply avoiding this kind of boundary condition is not an option due to their importance in most boundary value problems of engineering interest. The solution adopted in the present work to avoid this limitation is the so called Radial Basis Functions-Hermite Finite Difference (RBF-HFD). Within this section we present the formulation of the aforementioned method after explaining the Hermite interpolation framework which is at its core.

### 2.7.1 Hermite Interpolation

In the RBF-FD method, to find the approximated solution  $u^h$  valid for a given stencil  $\mathcal{X}_q = \{\mathbf{x}_1, \dots, \mathbf{x}_m\} \subseteq \mathcal{X}$ , based on a point  $\mathbf{x}_q \in \Omega$ , the following interpolation conditions are enforced:

$$u^h(\mathbf{x}_i) = u(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{q,I} \quad (2.7.1a)$$

$$\mathcal{B}u^h(\mathbf{x}_i) = g(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{q,B} \quad (2.7.1b)$$

In general the action of evaluating a function  $u$  at a point  $\mathbf{x}_i$  could be seen as the application of an operator  $\delta_{\mathbf{x}_i}$ , called point-evaluation functional, such that  $\delta_{\mathbf{x}_i}u = u(\mathbf{x}_i)$ . Furthermore functional  $\delta_{\mathbf{x}_i}$  can be combined with other operators; an example of this can be found when Neumann BC are encountered in conditions (2.7.1b):

$$\begin{aligned} \mathcal{B}u^h(\mathbf{x}_i) &= \frac{\partial}{\partial \mathbf{n}} u^h(\mathbf{x}_i) \\ &= \left( \delta_{\mathbf{x}_i} \circ \frac{\partial}{\partial \mathbf{n}} \right) u^h(\cdot) \end{aligned} \quad (2.7.2)$$

where  $\mathbf{x}_i$  belongs to  $\mathcal{X}_{q,B_N} \subseteq \mathcal{X}_{q,B}$ , the set of stencil's boundary nodes on which are applied Neumann BCs. To avoid excessive notation complexity in the following we assume that only Neumann BC are enforced on the points of the stencil that belongs to the boundary of the physical domain, i.e.  $\mathcal{X}_{q,B_N} \equiv \mathcal{X}_{q,B}$ .

Leveraging on these operators, interpolation conditions, more generally, can be defined through a set of given values  $\{u(\mathbf{x}_1), \dots, u(\mathbf{x}_m)\}$  along with a given set of functionals  $\Lambda = \{\lambda_1, \dots, \lambda_m\}$  as follow:

$$\lambda_i u^h = \lambda_i u \quad i = 1, \dots, m \quad (2.7.3)$$

These conditions are typically called Hermite interpolation conditions [19]. For RBF-FD local constraints (2.7.1) we have that  $\lambda_i = \delta_{\mathbf{x}_i}$  for those nodes  $\mathbf{x}_i$  that belong to  $\mathcal{X}_{q,I}$  and  $\lambda_i = \delta_{\mathbf{x}_i} \circ \frac{\partial}{\partial \mathbf{n}}$  for those which belong to  $\mathcal{X}_{q,B}$ .

In case of pure RBF interpolation scheme the local approximated solution  $u^h$  is given by:

$$u^h(\cdot) = \sum_{j=1}^m \alpha_j \Phi(\cdot, \mathbf{x}_j) \quad (2.7.4)$$

where “ $\cdot$ ” is a placeholder for the argument of the function meaning that it is not evaluated at any point. In this case the basis used for the interpolation is determined solely by the points of the stencil and is given by  $\{\Phi(\cdot, \mathbf{x}_1), \dots, \Phi(\cdot, \mathbf{x}_m)\} = \{\delta_{2,\mathbf{x}_1} \Phi(\cdot, \cdot), \dots, \delta_{2,\mathbf{x}_m} \Phi(\cdot, \cdot)\}$  where the subscript 2 in the point-evaluation functional indicates that it is acting on the second argument of  $\Phi(\cdot, \cdot)$ . Using interpolant (2.7.4) in case of Hermite interpolation conditions and repeating the same

scattered data interpolation procedure described in section 2.1, we would obtain a matrix  $\mathbf{B}$ :

$$\mathbf{B} = \begin{bmatrix} \lambda_1 \Phi(\mathbf{x}_1, \mathbf{x}_1) & \dots & \lambda_1 \Phi(\mathbf{x}_1, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \lambda_m \Phi(\mathbf{x}_m, \mathbf{x}_1) & \dots & \lambda_m \Phi(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix} \quad (2.7.5)$$

which is not symmetrical and does not satisfy, in general, the conditions of positive definiteness or non-singularity. This is exactly the case when RBF interpolation is applied to stencil  $\mathcal{X}_q$  where  $\mathbf{B}$  becomes:

$$\mathbf{B} = \begin{bmatrix} \delta_{1,\mathbf{x}_1} \Phi(\cdot, \mathbf{x}_1) & \dots & \delta_{1,\mathbf{x}_1} \Phi(\cdot, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \delta_{1,\mathbf{x}_{m_I}} \Phi(\cdot, \mathbf{x}_1) & \dots & \delta_{1,\mathbf{x}_{m_I}} \Phi(\cdot, \mathbf{x}_m) \\ \delta_{1,\mathbf{x}_{m_I+1}} \circ \frac{\partial}{\partial \mathbf{n}} \Phi(\cdot, \mathbf{x}_1) & \dots & \delta_{1,\mathbf{x}_{m_I+1}} \circ \frac{\partial}{\partial \mathbf{n}} \Phi(\cdot, \mathbf{x}_m) \\ \vdots & \ddots & \vdots \\ \delta_{1,\mathbf{x}_m} \circ \frac{\partial}{\partial \mathbf{n}} \Phi(\cdot, \mathbf{x}_1) & \dots & \delta_{1,\mathbf{x}_m} \circ \frac{\partial}{\partial \mathbf{n}} \Phi(\cdot, \mathbf{x}_m) \end{bmatrix} \quad (2.7.6)$$

which is not symmetric due to the rows associated to  $\delta_{1,\mathbf{x}_i} \circ \frac{\partial}{\partial \mathbf{n}}$  operators obtained from Neumann BCs. We remark once more that  $\mathbf{B}$  would have remained symmetrical in the case of Dirichlet BCs only, since their associated operator would have been just  $\delta_{1,\mathbf{x}_i}$ , the same applied on the first  $m_I$  rows.

An improvement, is done to the basic RBF interpolant in (2.7.4) in order to solve this problem: a more general interpolation basis is constructed by applying different functionals other than simple point evaluation to the same kernel  $\Phi(\cdot, \cdot)$  (where we recall that a kernel is simply a real valued function defined on a space given by the Cartesian product  $\Omega \times \Omega$ ); this leads to the so-called *Generalized Hermite Interpolant*. Suppose that the operator  $\delta_{2,\mathbf{x}_j}$  of the canonical RBF base is replaced with  $\lambda_{2,j}$ , the functional acting on the  $j$ -th Hermite condition, then the interpolation basis, in general, becomes  $\{ \lambda_{2,1} \Phi(\cdot, \cdot), \dots, \lambda_{2,m} \Phi(\cdot, \cdot) \}$ . In contrast to the classical scheme, the basis now does not depend solely on the nodes of the stencil to which  $\Phi$  is applied, but also on the conditions associated with them: this is the key feature that differentiate this approach compared to the simple RBF-FD one.

The new interpolant  $u^h$  defined according to the generalized Hermite interpolation scheme then becomes:

$$u^h(\cdot) = \sum_{j=1}^m \alpha_j \lambda_{2,j} \Phi(\cdot, \cdot) \quad (2.7.7)$$

which, once it is applied for interpolation, leads to an interpolation matrix:

$$\mathbf{B}_H = \begin{bmatrix} \lambda_{1,1}\lambda_{2,1}\Phi(\cdot, \cdot) & \dots & \lambda_{1,1}\lambda_{2,m}\Phi(\cdot, \cdot) \\ \vdots & \ddots & \vdots \\ \lambda_{1,m}\lambda_{2,1}\Phi(\cdot, \cdot) & \dots & \lambda_{1,m}\lambda_{2,m}\Phi(\cdot, \cdot) \end{bmatrix} \quad (2.7.8)$$

which is symmetric if  $\lambda_{1,j} = \lambda_{2,j}$  for  $j = 1, \dots, m$ . Because of this property, Hermite's formulation is also called *symmetric*, in contrast to the classical one which is referred to *unsymmetric*. Furthermore  $\mathbf{B}_H$  is also positive definite when the basic function  $\varphi$  used to define the kernel is strictly positive definite and operators  $\lambda_j$  are linearly independents [20]. The symmetry and positive definiteness properties previously mentioned also hold more generally when the Hermite interpolation conditions in equation (2.7.1) involves not only Neumann BCs, but also Dirichlet or Robin BCs, i.e. when  $\mathcal{X}_{q,B_N} \neq \mathcal{X}_{q,B}$ .

In addition, adoption of polynomial augmentation, as explained in section 2.4, to increase the accuracy of the RBF interpolant is still possible on condition of some changes. The augmented Hermite interpolant become:

$$u^h(\cdot) = \sum_{j=1}^m \alpha_j \lambda_{2,j} \Phi(\cdot, \cdot) + \sum_{k=1}^M \beta_k p_k(\cdot) \quad (2.7.9)$$

In order to obtain a solvable linear system to find coefficients  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$ , conditions (2.4.2) are modified as follow:

$$\sum_{i=1}^m \alpha_i \lambda_i p_k(\cdot) = 0, \quad i = 1, \dots, M \quad (2.7.10)$$

thus leading to:

$$\underbrace{\begin{bmatrix} \mathbf{B}_H & \mathbf{P}_H \\ \mathbf{P}_H^T & \mathbf{0} \end{bmatrix}}_M \begin{bmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \boldsymbol{\Lambda} \mathbf{u} \\ \mathbf{0} \end{bmatrix} \quad (2.7.11)$$

where  $\boldsymbol{\Lambda} \mathbf{u} = [\lambda_1 u, \dots, \lambda_m u]$  and  $\mathbf{P}_H$  is defined as:

$$\mathbf{P}_H = \begin{bmatrix} \lambda_1 p_1(\cdot) & \dots & \lambda_1 p_M(\cdot) \\ \vdots & \ddots & \vdots \\ \lambda_m p_1(\cdot) & \dots & \lambda_m p_M(\cdot) \end{bmatrix} \quad (2.7.12)$$

## 2.7.2 Mathematical formulation

The Hermite interpolant can be adopted also in Tolstykh's local approach with some modifications, that we explain in this subsection. We remember once again

that the problem that has to be solved is:

$$\begin{cases} \mathcal{L}u = f & \text{in } \Omega \\ \mathcal{B}u = g & \text{on } \partial\Omega \end{cases} \quad (2.7.13)$$

and for each internal node  $\mathbf{x}_i$  a stencil  $\mathcal{X}_i = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  composed of its  $m$  nearest neighbors is constructed.  $\mathcal{X}_i$  is then splitted in  $\mathcal{X}_{i,I} = \{\mathbf{x}_1, \dots, \mathbf{x}_{m_I}\}$  and  $\mathcal{X}_{i,B} = \{\mathbf{x}_{m_I+1}, \dots, \mathbf{x}_m\}$  respectively the set of stencil points that belong to the interior and boundary of the physical domain  $\Omega$ .

Interpolation conditions inside the stencil are the same as those reported at the beginning of previous subsection:

$$u^h(\mathbf{x}_i) = u(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{q,I} \quad (2.7.14a)$$

$$\mathcal{B}u^h(\mathbf{x}_i) = g(\mathbf{x}_i) \quad \text{if } \mathbf{x}_i \in \mathcal{X}_{q,B} \quad (2.7.14b)$$

At this point the interpolant or local approximated solution of the problem is defined according to the Hermite interpolation theory and takes the same form reported in equation (2.7.9). Making explicit the functionals  $\lambda_1, \dots, \lambda_m$  used in the definition of  $u^h$  and applying it at the center of the stencil  $\mathcal{X}_I$ , we can write it as:

$$u^h(\mathbf{x}_i) = \sum_{j=1}^{m_I} \alpha_j \Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{j=m_I+1}^m \alpha_j \mathcal{B}_2 \Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{k=1}^M \beta_k p_k(\mathbf{x}_i) \quad (2.7.15)$$

where  $\boldsymbol{\alpha}_I = [\alpha_1, \dots, \alpha_{m_I}]$  and  $\boldsymbol{\alpha}_B = [\alpha_{m_I+1}, \dots, \alpha_m]$  are used to distinguish the coefficients that are associated to internal and boundary nodes. Conditions (2.7.10) have to be enforced on the coefficients of the expansion in order to obtain a square coefficient matrix associated to constraints (2.7.14):

$$\sum_{j=1}^{m_I} \alpha_j p_k(\mathbf{x}_j) + \sum_{j=m_I+1}^m \alpha_j \mathcal{B} p_k(\mathbf{x}_j) = 0 \quad k = 1, \dots, M \quad (2.7.16)$$

Due to the new form of the interpolant  $u^h$ , system (2.6.12) now becomes:

$$\underbrace{\begin{bmatrix} \Phi_{I,I} & \mathcal{B}_2 \Phi_{I,B} & P_I \\ \mathcal{B}_1 \Phi_{B,I} & \mathcal{B}_1 \mathcal{B}_2 \Phi_{B,B} & \mathcal{B} P_B \\ P_I^T & \mathcal{B} P_B^T & 0 \end{bmatrix}}_{M_{BC}} \begin{bmatrix} \boldsymbol{\alpha}_I \\ \boldsymbol{\alpha}_B \\ \boldsymbol{\beta} \end{bmatrix} = \begin{bmatrix} \mathbf{u}_I \\ \mathbf{g} \\ \mathbf{0} \end{bmatrix} \quad (2.7.17)$$

Proceeding by following the same steps of the RBF-FD method we apply the linear

operator  $\mathcal{L}$  to the interpolant  $u^h$ :

$$\begin{aligned}\mathcal{L}u^h(\mathbf{x}_i) &= \sum_{j=1}^{m_I} \alpha_j \mathcal{L}_1 \Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{j=m_I+1}^m \alpha_j \mathcal{L}_1 \mathcal{B}_2 \Phi(\mathbf{x}_i, \mathbf{x}_j) + \sum_{k=1}^M \beta_k \mathcal{L}p_k(\mathbf{x}_i) \\ &= [\boldsymbol{\alpha}_I \quad \boldsymbol{\alpha}_B \quad \boldsymbol{\beta}] \begin{bmatrix} \mathcal{L}_1 \Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}_1 \mathcal{B}_2 \Phi(\mathbf{x}_i, \mathcal{X}_{i,B}) \\ \mathcal{L}p(\mathbf{x}_i) \end{bmatrix}\end{aligned}\tag{2.7.18}$$

Coefficient vectors  $\boldsymbol{\alpha}_I$ ,  $\boldsymbol{\alpha}_B$  and  $\boldsymbol{\beta}$  obtained from equation (2.7.17) are then substituted in equation (2.7.18) leading to:

$$\mathcal{L}u^h(\mathbf{x}_i) = [\mathbf{u}_I \quad \mathbf{g} \quad \mathbf{0}] \mathbf{M}_{BC}^{-T} \begin{bmatrix} \mathcal{L}_1 \Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}_1 \mathcal{B}_2 \Phi(\mathbf{x}_i, \mathcal{X}_{i,B}) \\ \mathcal{L}p(\mathbf{x}_i) \end{bmatrix}\tag{2.7.19}$$

Exactly as in RBF-FD formulation the product of the last two factors of the aforementioned equation is known and the resulting vector  $\mathbf{c}(\mathbf{x}_i) = [\mathbf{c}_I(\mathbf{x}_i), \mathbf{c}_B(\mathbf{x}_i), \mathbf{c}_p(\mathbf{x}_i)]$  contains the coefficients that make up the Finite Difference-like approximation of the differential operator  $\mathcal{L}$  in a given point  $\mathbf{x}_i$ . In this case  $\mathbf{c}(\mathbf{x}_i)$  is computed by solving:

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}_1 \Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}_1 \mathcal{B}_2 \Phi(\mathbf{x}_i, \mathcal{X}_{i,B}) \\ \mathcal{L}p(\mathbf{x}_i) \end{bmatrix}\tag{2.7.20}$$

From now on the steps to build the matrix  $\mathbf{C}_I$ , that approximate  $\mathcal{L}$ , are the same of the local RBF-FD ones. The same applies to obtaining the values of the approximate solution of the partial differential equation at the nodes inside the domain.

## Chapter 3

# Automatic Differentiation (AD) and adjoint method in optimization

In this chapter we are going to discuss Automatic Differentiation (AD) and the adjoint method which are two techniques that allows to compute efficiently the gradient of a function which depends on many parameters. We also see how they can be used to improve the efficiency of gradient-based optimization algorithms.

To do so we introduce to the possible techniques that allows to compute derivatives with a particular focus on Automatic Differentiation (AD), highlighting its advantages and its two main implementations.

Subsequently, we will provide a brief overview of design optimization problems, in particular of those related to systems which are numerically described by RBF-FD methods.

Finally, we conclude the chapter explaining in detail the application of the AD and the adjoint method to 1D and 3D RBF-FD based design optimization problems.

### 3.1 Automatic Differentiation (AD)

In general there exists different ways to compute derivatives using a computer program, these are:

**Manual Differentiation** In this case *analytical* derivatives are computed by hand and then are plugged into standard optimization procedures such as gradient descend. Of course doing so is time consuming and prone to error.

**Numerical Differentiation** In this case finite difference methods, as the ones reported in subsection [2.6.1 on page 12](#), are used to approximate the deriva-



tives *values*. Easy to implement it has the disadvantage to be inaccurate due to round-off and truncation errors [26].

**Symbolic Differentiation** This case addresses the weakness of both manual and numerical differentiation. *Analytical* derivative expressions are automatically obtained by modern computer algebra systems such as Mathematica<sup>1</sup> or SymPy<sup>2</sup>. Unfortunately, often, the outcomes are plagued with the problem of “expression swell” which means that the resulting expressions are large, complex and cryptic. Furthermore it can be applied only to models defined in closed-form.

**Automatic Differentiation** This last case refers to a family of techniques that compute derivative *values* (in contrast with symbolic differentiation) by using symbolic rules of differentiation (but keeping track of derivative values as opposed to the resulting expressions) through accumulation of values during code execution. Thanks to this, Automatic Differentiation (AD), can be applied to code involving branches, loops and recursion as opposed to symbolic differentiation [29]. The mix between symbolic and numerical differentiation gives to these methodologies an hybrid nature.

For the remaining of this section we will explain the details of Automatic Differentiation (AD), to do so we will look at its two (main) implementations:

- *forward mode*, also known as tangent linear mode;
- *reverse mode* also known as cotangent linear mode or *adjoint* mode.

To explain how each of the two modes works, we show how they are applied to a function  $\tilde{f}: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  defined as:

$$\tilde{f}(x_1, x_2) = [\tilde{f}_1(x_1, x_2) \quad \tilde{f}_2(x_1, x_2)] \quad (3.1.1)$$

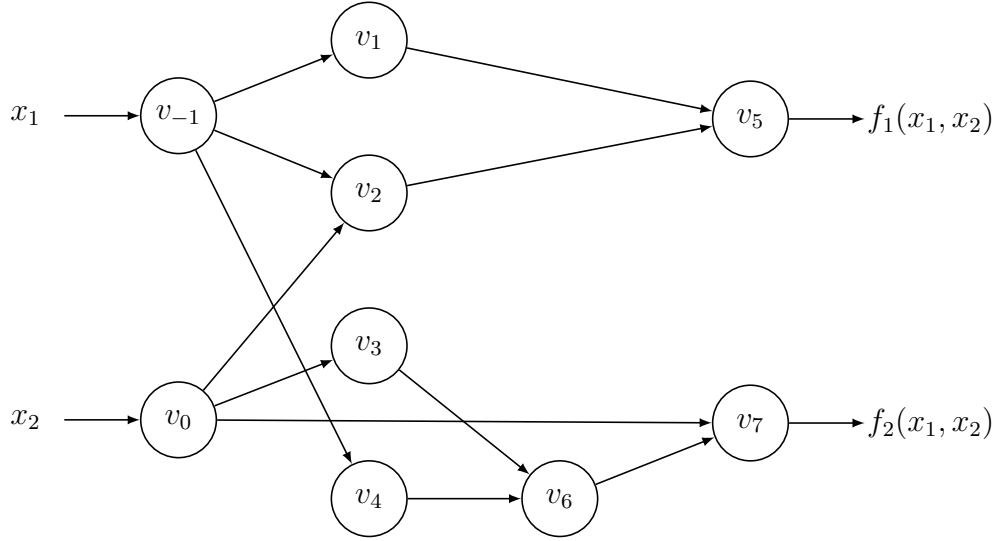
with  $\tilde{f}_1(x_1, x_2) = x_1x_2 + \cos x_1$  and  $\tilde{f}_2(x_1, x_2) = x_2^3 + \ln x_1 - x_2$ .

Before entering into the detail, we notice that every function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  can be rewritten as a computational graph. A computational graph is a direct graph whose nodes correspond to operations and each operation can feed its output into other operations; once the graph is fed with some variables each node become automatically function of those variables. The usual operations considered as nodes are: binary arithmetic operations, the unary sign switch and transcendental functions such as exponential, logarithm and trigonometric functions. Creating a computational graph for a function becomes easy by indicating with:

---

<sup>1</sup>Wolfram Research, Inc.’s proprietary software for technical computing. See <https://www.wolfram.com/mathematica/> for more informations.

<sup>2</sup>Python open source library for symbolic mathematics. The project can be found at <https://www.sympy.org/en/index.html>.



**Figure 3.1:** Computational graph of function  $\tilde{f}(x_1, x_2) = [\tilde{f}_1(x_1, x_2) \ \tilde{f}_2(x_1, x_2)]$ . Definitions of intermediate variables  $v_{-1}, \dots, v_7$  can be found in table 3.1 or table 3.2

- $v_{i-n} = x_i$ ,  $i = 1, \dots, n$  the input variables;
- $v_i$ ,  $i = 1, \dots, l$  the intermediate variable;
- $y_{m-i} = v_{l-i}$ ,  $i = m - 1, \dots, 0$  the output variables.

To better identify the relationship between the elementary operations which constitute a function it is helpful to create an *evaluation trace* for the function itself. Left-hand side of table 3.1 and figure 3.1 contains respectively the evaluation trace and the computational graph for function  $\tilde{f}$ .

Using the aforementioned representations we see that every function ultimately is a composition of elementary operations. This also means that its numerical derivatives can be computed by combining all the numerical derivatives of the constituent operations through the *chain rule*: this, is the main idea of Automatic Differentiation.

### 3.1.1 Forward mode

In order to compute the derivative of the function  $\tilde{f}$ , reported in (3.1.1), with respect to  $x_1$  we start by considering the evaluation trace on the left-hand side of table 3.1 and we associate to each intermediate variable  $v_i$  the derivative:

$$v'_i = \frac{\partial v_i}{\partial x_1}$$

CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT  
METHOD IN OPTIMIZATION

| Forward Primal Trace |                 |                | Forward Tangent (Derivative) Trace |                               |                           |
|----------------------|-----------------|----------------|------------------------------------|-------------------------------|---------------------------|
| $v_{-1}$             | $= x_1$         | $= 2$          | $v'_{-1}$                          | $= x'_1$                      | $= 1$                     |
| $v_0$                | $= x_2$         | $= 3$          | $v'_0$                             | $= x'_2$                      | $= 0$                     |
| $v_1$                | $= \cos v_{-1}$ | $= \cos 2$     | $v'_1$                             | $= -v'_{-1} \sin v_{-1}$      | $= -1 \cdot \sin 2$       |
| $v_2$                | $= v_{-1} v_0$  | $= 2 \cdot 3$  | $v'_2$                             | $= v'_{-1} v_0 + v_{-1} v'_0$ | $= 1 \cdot 3 + 2 \cdot 0$ |
| $v_3$                | $= v_0^3$       | $= 3^3$        | $v'_3$                             | $= 3v_0^2 v'_0$               | $= 3 \cdot 2^2 \cdot 0$   |
| $v_4$                | $= \ln v_{-1}$  | $= \ln 2$      | $v'_4$                             | $= v'_{-1}/v_{-1}$            | $= 1/2$                   |
| $v_5$                | $= v_2 + v_1$   | $= 6 - 0.416$  | $v'_5$                             | $= v'_2 + v'_1$               | $= 3 - 0.909$             |
| $v_6$                | $= v_3 + v_4$   | $= 27 + 0.693$ | $v'_6$                             | $= v'_3 + v'_4$               | $= 0 + 0.5$               |
| $v_7$                | $= v_6 - v_0$   | $= 27.693 - 3$ | $v'_7$                             | $= v'_6 - v'_0$               | $= 0.5 - 0$               |
| $y_1$                | $= v_5$         | $= 5.584$      | $\mathbf{y}'_1$                    | $= \mathbf{v}'_5$             | $= \mathbf{2.091}$        |
| $y_2$                | $= v_7$         | $= 24.693$     | $\mathbf{y}'_2$                    | $= \mathbf{v}'_7$             | $= \mathbf{0.5}$          |

**Table 3.1:** Forward mode AD example to evaluate the derivatives  $\frac{\partial y_1}{\partial x_1}$  and  $\frac{\partial y_2}{\partial x_1}$  of  $\tilde{f}(x_1, x_2)$  at  $[x_1, x_2] = [2, 3]$ .  $x'_1$  and  $x'_2$  are respectively set to 1 and 0 in order to derive only respect  $x_1$ . On the left is reported the forward evaluation trace, on the right the tangent one

Here the variable with respect to which we differentiate remain the same independently on  $i$ . Moving from top to bottom in the forward primal trace the corresponding tangent (derivative) trace, presented on the right-hand side of table 3.1, is generated by applying the chain rule to each encountered operation. After the primals  $v_i$  are evaluated, also the corresponding tangents  $v'_i$  are, again, from top to bottom, which means that tangent trace can be computed in parallel with the forward trace. Doing so gives us the desired derivatives in the final variables  $v'_5 = \frac{\partial y_1}{\partial x_1}$  and  $v'_7 = \frac{\partial y_2}{\partial x_1}$ .

Forward mode can also be employed to evaluate the Jacobian of a generic function  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  at a point  $\mathbf{x} = \mathbf{a}$  by performing  $n$  distinct forward passes. By setting only one variable  $x'_i = 1$  and the others to zero we obtain:

$$y'_j = \left. \frac{\partial y_j}{\partial x_i} \right|_{\mathbf{x}=\mathbf{a}} \quad j = 1, \dots, m.$$

and reiterating for  $i = 1, \dots, n$ , placing the resulting vectors side by side, we eventually obtain the desired Jacobian:

$$\mathbf{J}_f = \left[ \begin{array}{ccc} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{array} \right] \bigg|_{\mathbf{x}=\mathbf{a}}$$

Furthermore, by properly fine-tuning the values of the variables  $x'_1, \dots, x'_n$ , it is possible to compute efficiently and in a matrix-free way Jacobian-vector products:

$$\mathbf{J}_f \mathbf{r} = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \cdots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \cdots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix}$$

To accomplish this all that needs to be done is simply initializing  $\mathbf{x}' = [x'_1, \dots, x'_n]$  with  $\mathbf{r}$ ; this result is particularly important in the evaluation of directional derivatives. Before moving on, it is crucial to point out that forward mode AD:

- for a function  $f: \mathbb{R} \rightarrow \mathbb{R}^m$  allows to evaluate all its derivatives in just one forward pass, regardless of  $m$ ;
- for a scalar field  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , the evaluation of its gradient  $\nabla f = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right]$  always require  $n$  evaluations (since gradient is nothing more than a Jacobian of size  $1 \times n$ ).

This means that the here explained implementation of Automatic Differentiation (AD) has a computational cost which scales linearly with the number of function inputs. It is effective during the computation of derivative for cases  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$  where  $n \ll m$ ; for cases  $n \gg m$  *reverse mode* is more beneficial: we will discover why in the next subsections.

### 3.1.2 Reverse mode

In this case, as opposed to forward mode AD, the derivatives are propagated back from a given output: this means that the underlying mechanism is similar, yet more general, to backpropagation algorithm used in neural networks.

To do so, different variables, called adjoints, are associated to each variable  $v_i$  of the Forward Primal Trace and each of them is defined as:

$$\bar{v}_i = \frac{\partial y_j}{\partial v_i}$$

which is nothing more than the sensitivity (i.e. derivative) of the output  $y_j$  with respect to changes in  $v_i$ . It is important for the reader to note that, with this mode, the variable with respect to the derivatives are computed is not fixed as in forward mode AD. An example of reverse mode AD is presented in table 3.2 where we show how to compute the sensitivities of the output  $y_1 = \tilde{f}_1(x_1, x_2)$  with respect to the inputs  $x_1$  and  $x_2$ .

CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT  
METHOD IN OPTIMIZATION

---

| Forward Primal Trace               |  |  |           |
|------------------------------------|--|--|-----------|
| $v_{-1}$                           | $= x_1$  | $= 2$  |           |
| $v_0$                              | $= x_2$  | $= 3$  |           |
| $v_1$                              | $= \cos v_{-1}$  | $= \cos 2$   |           |
| $v_2$                              | $= v_{-1} v_0$   | $= 2 \cdot 3$  |           |
| $v_3$                              | $= v_0^3$  | $= 3^3$  |           |
| $v_4$                              | $= \ln v_{-1}$   | $= \ln 2$  |           |
| $v_5$                              | $= v_2 + v_1$  | $= 6 - 0.416$  |           |
| $v_6$                              | $= v_3 + v_4$  | $= 27 + 0.693$   |           |
| $v_7$                              | $= v_6 - v_0$  | $= 27.693 - 3$   |           |
| $y_1$                              | $= v_5$  | $= 5.584$  |           |
| $y_2$                              | $= v_7$  | $= 24.693$   |           |
| Reverse Adjoint (Derivative) Trace |  |  |           |
| $\bar{\mathbf{x}}_1$               | $= \frac{\partial y_1}{\partial \mathbf{v}_{-1}} \frac{\partial \mathbf{v}_{-1}}{\partial \mathbf{x}_1}$   | $= \bar{\mathbf{v}}_{-1} \cdot \mathbf{1}$                       | $= 2.091$ |
| $\bar{\mathbf{x}}_2$               | $= \frac{\partial y_1}{\partial \mathbf{v}_0} \frac{\partial \mathbf{v}_0}{\partial \mathbf{x}_2}$   | $= \bar{\mathbf{v}}_0 \cdot \mathbf{1}$                          | $= 2$     |
| $\bar{v}_{-1}$                     | $= \frac{\partial y_1}{\partial v_1} \frac{\partial v_1}{\partial v_{-1}} + \frac{\partial y_1}{\partial v_2} \frac{\partial v_2}{\partial v_{-1}} + \frac{\partial y_1}{\partial v_4} \frac{\partial v_4}{\partial v_{-1}}$ | $= -\bar{v}_1 \sin(v_{-1}) + \bar{v}_2 v_0 + \bar{v}_4 / v_{-1}$ | $= 2.091$ |
| $\bar{v}_0$                        | $= \frac{\partial y_1}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y_1}{\partial v_3} \frac{\partial v_3}{\partial v_0}$  | $= \bar{v}_2 v_{-1} + 3\bar{v}_3 v_0^2$                          | $= 2$     |
| $\bar{v}_1$                        | $= \frac{\partial y_1}{\partial v_5} \frac{\partial v_5}{\partial v_1}$  | $= \bar{v}_5 \cdot 1$  | $= 1$     |
| $\bar{v}_2$                        | $= \frac{\partial y_1}{\partial v_5} \frac{\partial v_5}{\partial v_2}$  | $= \bar{v}_5 \cdot 1$  | $= 1$     |
| $\bar{v}_3$                        | $= \frac{\partial y_1}{\partial v_6} \frac{\partial v_6}{\partial v_3}$  | $= \bar{v}_6 \cdot 1$  | $= 0$     |
| $\bar{v}_4$                        | $= \frac{\partial y_1}{\partial v_6} \frac{\partial v_6}{\partial v_4}$  | $= \bar{v}_6 \cdot 1$  | $= 0$     |
| $\bar{v}_6$                        | $= \frac{\partial y_1}{\partial v_7} \frac{\partial v_7}{\partial v_6}$  | $= \bar{v}_7 \cdot 1$  | $= 0$     |
| $\bar{v}_5$                        | $= \bar{y}_1$  |  | $= 1$     |
| $\bar{v}_7$                        | $= \bar{y}_2$  |  | $= 0$     |

**Table 3.2:** Reverse mode AD example with  $[y_1, y_2] = \tilde{f}(x_1, x_2)$  evaluated at  $[x_1, x_2] = [2, 3]$ . First, primal trace is evaluated (table above) and then adjoint variables are computed, from the bottom up, in a second phase (table below). The initialization  $\bar{v}_5 = \frac{\partial y_1}{\partial v_5} = \frac{\partial y_1}{\partial y_1} = \bar{y}_1 = 1$  and  $\bar{v}_7 = \frac{\partial y_1}{\partial v_7} = \frac{\partial y_1}{\partial y_2} = \bar{y}_2 = 0$  is such that it allows the sensitivities to be calculated with respect the fist output

## CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

Since typical usage of chain rule is forward derivatives propagation, reverse mode AD might appear confusing at first sight; to make it clearer we show how the chain rule can be used to back-propagate derivatives in order to compute the contribution  $\bar{v}_0$  of the change in variable  $v_0$  to the change in the output  $y_1$ . From figure 3.1 can be seen that the only way variable  $v_0$  can affect  $y_1$  is through affecting  $v_2$  and  $v_3$ , so its contribution to the change in  $y_1$  can be computed using the chain rule as follow:

$$\frac{\partial y_1}{\partial v_0} = \frac{\partial y_1}{\partial v_2} \frac{\partial v_2}{\partial v_0} + \frac{\partial y_1}{\partial v_3} \frac{\partial v_3}{\partial v_0} \quad (3.1.2)$$

which can be rewritten in

$$\begin{aligned} \frac{\partial y_1}{\partial v_0} &= \bar{v}_2 \frac{\partial v_2}{\partial v_0} + \bar{v}_3 \frac{\partial v_3}{\partial v_0} \\ &= \bar{v}_2 v_{-1} + 3\bar{v}_3 v_0^2 \end{aligned} \quad (3.1.3)$$

where quantities  $\bar{v}_2$  and  $\bar{v}_3$  are those that are computed first, from previous passages of the method, and derivatives  $\frac{\partial v_2}{\partial v_0}$ ,  $\frac{\partial v_3}{\partial v_0}$ , can be easily computed by evaluating the result of symbolic differentiation of the elementary operations on the primals  $v_i$ .

In light of the previous example, it is clear that, to compute derivatives, primals  $v_i$  must be known along their dependencies within the computational graph. To do so two phases are required:

- A *forward step* where variables  $v_i$  are populated and their dependencies recorded;
- A *reverse step* where adjoints  $\bar{v}_i$  are evaluated from outputs to inputs using values obtained in the first step.

We remark that once the backward pass is completed we get *all* the sensitivities  $\frac{\partial y_j}{\partial x_1}, \dots, \frac{\partial y_j}{\partial x_n}$  in one single pass; this means that for scalar fields  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  the computation of the gradient  $\nabla f$  requires only one application of reverse mode as opposed to the  $n$  passes required in case of forward mode. Nevertheless, this advantage comes at the cost of increased memory requirements which grows proportionally (in the worst case) to the number of operations in the evaluated function [29].

Generalizing the reasoning to functions  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ , similarly to how it was done in previous section, we can conclude that reverse mode AD can also be used to evaluate Jacobians at a point; in particular, since each pass of reverse AD allows to compute one row of the Jacobian, its usage is advantageous when  $m \ll n$ . With regard to products between Jacobian and vector, reverse AD allows to compute

efficiently vector-Jacobian products of type:

$$\mathbf{r}^T \mathbf{J}_f = \begin{bmatrix} r_1 & \dots & r_m \end{bmatrix} \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial x_1} & \dots & \frac{\partial y_m}{\partial x_n} \end{bmatrix} \quad (3.1.4)$$

by initializing  $\bar{\mathbf{y}} = [\bar{y}_1, \dots, \bar{y}_m]$  with  $\mathbf{r}$ .

## 3.2 Design Optimization

In everyday life the need for optimization arises almost naturally: birds optimize their wings' shape in real time, dogs optimize their trajectories in order to reach a specific place, and people constantly seek to improve their lives and the system that surround them. In these cases optimization is used as synonym of improvement, and from a mathematical standpoint it can be understood as the concept of: “finding the best possible solution by changing variables that can be controlled, often subject to constraints” [30].

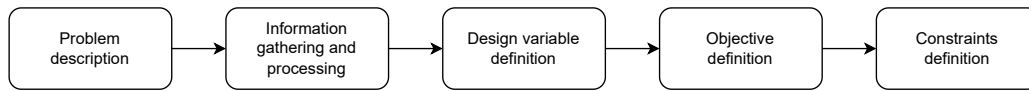
One simple way to solve general optimization problems would be a manual approach where in order to find the best possible solution of the problem a designer, the one in charge of solving the problem, has to adjust each single variable at time; however this approach has several limitation: it tends to lead suboptimal results and evaluating all possible variables might be too time-consuming. These issues become even more severe when tackling more complex problems like those found in engineering as: wing design in aerospace engineering, structural design in civil engineering, circuit design in electrical engineering and mechanism design in mechanical engineering.

Design optimization is the tool typically used to solve this last type of problems in order to accelerate the design cycle and obtain better results; it allows to automate the optimization by means of an optimizer, but in order to be used it requires a correct formulation of the problem that has to be solved, which includes:

- the design variables that are to be changed;
- the objective that has to be minimized/maximized;
- the constraints that has to be met.

A more detailed explanation of the formulation, including the points mentioned above, will be provided in the next subsections.

## CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION



**Figure 3.2:** Steps for the mathematical description of a optimization problem

### 3.2.1 Optimization problem formulation

The mathematical formulation of an optimization problem is a 5-steps procedure which is necessary in order to employ design optimization and it has the (positive) side effect of increasing the designer's insights about the problem. The steps are reported in figure 3.2 and we will explain each of this in this subsection.

At the very beginning the designer is required to write a description of the system and of the design problem, even if vague, with a list of all goals and requirements.

As a next step it has to gather as much information as possible about the problem in order to set its expectations at the right level (not too high, not too low) and to understand the constraints that the problem is subject to. We also remark that:

- in order to gain these insights about the optimization problem raw data might need to be processed and organized;
- information gathering and refinement are ongoing processes during the formulation of the problem and up to this point it is almost impossible for a designer to learn everything about the problem.

The next three steps, where mathematical definition for design variables, objective function and constraints comes into play, allows to further increase the clarity and the knowledge of the problem at hand.

Once the system's *design variables* (also called *design parameters*), which describe the system, are identified, can be represented by means of a vector:

$$\mathbf{l} = [l_1, \dots, l_N] \quad (3.2.1)$$

where it must be ensured that each design variable is independent of the others in order to allow the optimizer to perform a correct analysis of a specific design. At each vector  $\mathbf{l}$  is associated a single design and the number of variables  $N$  defines the dimensionality of the problem that has to be solved. The design variables allow also to distinguish between two different classes of problems:

- *continuous optimization problems*, where variables are allowed to vary continuously within a range:  $l_i \leq l_i \leq \bar{l}_i$  for  $i = 1, \dots, N$ ;



## CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

- *discrete optimization problems*, in case variables are restricted to assume discrete set of values regardless of whether they are real or integer.

After choosing the design variables that identify a design, we need a quantity that tell us “how good” that design is. This quantity is obtained by means of the so called *objective* or *cost function*  $J$ , which is a scalar field that assign a number to a given design variable vector  $\mathbf{l}$ : with these numbers we can therefore define an order and compare different designs. Minimizing or maximizing the objective, depending on the problem at hand, allows to find the optimal design. What matters at this stage is the choice of a function that represent the real goal of the designer otherwise regardless of the precision with which the optimal design is identified, it will not be optimal from an engineering perspective.

As final step, the formulation of *constraints* is required, these depend on the design variables and delineate the so called feasible region which is the set of allowed designs. These can be of two types:

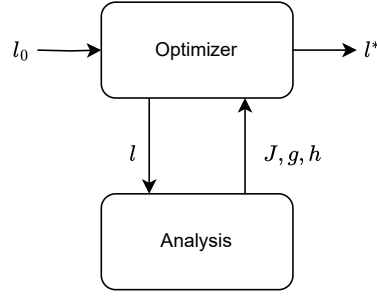
- *inequality*, indicated with  $g(\mathbf{l}) \leq 0$ . In literature is typically used the “less or equal” convention, but there is no loss of generality since a “greater or equal” constraint can be converted into “less or equal” simply by multiplying it by  $-1$ ;
- *equality*, indicated with  $g(\mathbf{l}) = 0$ . These constraints can also be seen as two separate constraints:  $g(\mathbf{l}) \leq 0$  and  $g(\mathbf{l}) \geq 0$ ;

Even if at a first glance constraints may appear limiting and superfluous, they enable the optimizer to find designs that match designer’s expectation.

Now that we have gone through the 5 steps described in this subsection a general design optimization problem can be defined as:

$$\begin{aligned}
 & \text{minimize} && J(\mathbf{l}) \\
 & \text{by varying} && \underline{l}_i \leq l_i \leq \bar{l}_i && i = 1, \dots, N \\
 & \text{subject to} && g_j(\mathbf{l}) && j = 1, \dots, n_g \\
 & && h_k(\mathbf{l}) && j = 1, \dots, n_h
 \end{aligned} \tag{3.2.2}$$

where we indicate with  $n_g$  and  $n_h$  the number of inequality and equality constraints respectively. The minimum of  $J$  is found automatically and iteratively by an optimizer starting from an initial design  $\mathbf{l}_0$ , as shown in figure 3.3. At each iteration the optimizer computes the values of the objective and constraint functions for the current design  $\mathbf{l}$  (analysis phase) and then return  $\mathbf{l}'$  which is an improvement of  $\mathbf{l}$ . The same process is then repeated using  $\mathbf{l}'$  instead of  $\mathbf{l}_0$  as input until  $\mathbf{l}'$  coincides with  $\mathbf{l}^*$  or it is deemed accurate enough, in the sense of  $J$ , by the designer.



**Figure 3.3:** Solution of an optimization problem. At the beginning the initial design  $l_0$  is provided to the optimizer, at each optimization step it produce a new design  $l$ , from the previous one, and through the analysis it evaluates the objective  $J$  and the constraints  $g, h$  in order to produce a new design  $l'$ . Once  $l'$  corresponds or is considered close enough to the optimum it is yield as the optimum design  $l^*$

### 3.2.2 Models and optimization problems

In previous subsection we have introduced the concepts of objective and constraint functions, here we cover *how* they are modeled and computed. These functions, as stated before, are evaluated in the analysis phase, and require the solution of a numerical model for a given design. A numerical model is obtained from the discretization of a mathematical model which in turn is a formal description of the physical system that undergoes the optimization.

Mathematical models are also referred as governing equations and determines the state of a physical system under specific conditions. Many of them consists of differential equations which require discretization in order to be solved, as an example, one can consider the one given in equation (1.0.1). After the discretization of such mathematical models we obtain a numerical model that can be written in residual form as:

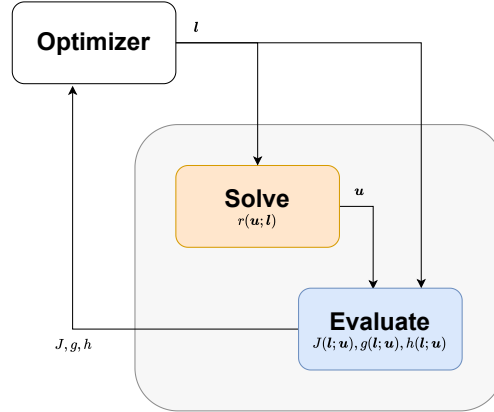
$$r_w(u_1, \dots, u_{n_r}) = 0, \quad w = 1, \dots, n_r \quad (3.2.3)$$

where  $u_1, \dots, u_{n_r}$  are the elements that compose the state vector  $\mathbf{u}$ . At this point the state  $\mathbf{u}$  can be therefore found solving equation (3.2.3) by means of a proper solver.

Often the state of a model depends on a given design  $\mathbf{l}$ , which means that this dependency enter also in equations (3.2.3). These can be rewritten as:

$$r_w(\mathbf{u}; \mathbf{l}), \quad w = 1, \dots, n_r \quad (3.2.4)$$

where the semicolon is used to underline the fact that  $\mathbf{l}$  is fixed when the equations are solved for  $\mathbf{u}$ . This implies that the objective and the constraints, which



**Figure 3.4:** Optimization scheme, where the analysis phase, shaded in gray, is broken down into the solution of the numerical model obtained from the governing equations (in red) and the evaluation of objective and constraints (in blue)

depend on the system and thus on the state  $\mathbf{u}$ , in turn depend implicitly on  $\mathbf{l}$ : this means that they are fully determined by the design variables. By making this dependency explicit we can rewrite the general design optimization problem reported in equation (3.2.2) as:

$$\begin{aligned}
 & \text{minimize} && J(\mathbf{l}; \mathbf{u}) \\
 & \text{by varying} && l_i && i = 1, \dots, N \\
 & \text{subject to} && g_j(\mathbf{l}; \mathbf{u}) \leq 0 && j = 1, \dots, n_g \\
 & && h_k(\mathbf{l}; \mathbf{u}) = 0 && k = 1, \dots, n_h \\
 & && \underline{l}_i \leq l_i \leq \bar{l}_i && i = 1, \dots, N \\
 & \text{while solving} && r_w(\mathbf{u}; \mathbf{l}) = 0 && w = 1, \dots, n_r \\
 & \text{by varying} && u_w && w = 1, \dots, n_r
 \end{aligned} \tag{3.2.5}$$

where the last two lines are solved for  $\mathbf{u}$  at each optimization step given the current design  $\mathbf{l}$ . Therefore the overall structure of an optimization problem, reported in 3.3, can be adjusted by explicitly integrating the solver into the scheme, yielding the one reported in figure 3.4.

From a different perspective the governing equations can be also considered as

equality constraints which leads to the following different formulation:

$$\begin{aligned}
 & \text{minimize} && J(\mathbf{l}, \mathbf{u}) \\
 & \text{by varying} && l_i && i = 1, \dots, N \\
 & && u_w && w = 1, \dots, n_r \\
 & \text{subject to} && g_j(\mathbf{l}; \mathbf{u}) \leq 0 && j = 1, \dots, n_g \\
 & && h_k(\mathbf{l}; \mathbf{u}) = 0 && k = 1, \dots, n_h \\
 & && l_i \leq l_i \leq \bar{l}_i && i = 1, \dots, N \\
 & && r_w(\mathbf{u}; \mathbf{l}) = 0 && w = 1, \dots, n_r
 \end{aligned} \tag{3.2.6}$$

In this case design variables  $u_1, \dots, u_{n_r}$  are considered as such since objective explicitly depends on them. In practice, however, their values can vary only nominally since practically must satisfy constraints  $r_w(\mathbf{u}; \mathbf{l}) = 0$  for  $w = 1, \dots, n_r$ . For this reason is possible to find problem formulations like the one reported above, but without the presence of the variables  $u_1, \dots, u_{n_r}$  among those that can be varied. In any case formulations (3.2.5) and (3.2.6) are equivalent and leads to the same optimal design  $\mathbf{l}^*$ .

### 3.2.3 Optimization Algorithms

So far we know that the optimizer. at each optimization step, update the current design producing a new one. In this subsection we delve in more detail how this update, performed by means of an optimization algorithm, can be produced.

A multitude of algorithms can be used to determine the aforementioned update and can be divided into diverse families, two of the most important being:

**Mathematical versus Heristic** For their proper functioning two of the essentials they need are: iterative process and evaluation criteria. The former is required in order to determine the sequence of points in the design space evaluated during the optimization, the latter in order to stop the search.

Heuristics, unlike the mathematical algorithm based on mathematical principles, are based on rule of thumb for the iterative process and the evaluation criterion. Algorithms that mix mathematical arguments with heuristic can also be defined.

**Information order** As can be also seen from the scheme reported in figure 3.4, the optimizer require the values of the objective and constraints for a fixed design: these are zeroth-order information. If the above-mentioned are the only values required by the algorithm to perform the update, then it takes the name of zeroth-order or or gradient-free algorithm.

We instead call first-order or gradient-based those algorithms which make use of gradients of one or both of the objective and constraint functions with respect to design variables (which are first-order informations). Their main advantages respect the zeroth-order algorithms are: better scalability of function evaluations with the number of design variables [30] and the possibility to easily check if a point in the design space satisfies optimality condition.

Algorithms that make use of second-order information also exists: an example are those that compute the curvature of the objective which indicate an idea of where the function might flatten out, giving more insight into the function behavior beyond simple steepness.

A final clarification about the information order is about the difference between the one provided by the user and the one actually used by the algorithm, which might differ. An example of this difference is the finite difference algorithm which estimate a first-order quantity, the gradient, starting from the function values provided by the user. In this case the gradient estimate requires additional function evaluations compared to the case where user provides the gradient value themselves.

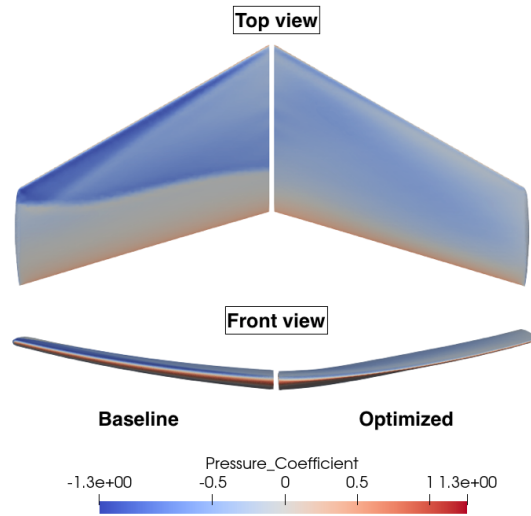
For the rest of this work we will consider only the following optimization algorithm:

$$\mathbf{l}_{k+1} = \mathbf{l}_k + \alpha \frac{dJ}{d\mathbf{l}} \quad (3.2.7)$$

where the optimizer, at each  $k$ -th steps, yields the improved design identified by  $\mathbf{l}_{k+1}$  by increasing the value of each design variable of the current design  $\mathbf{l}_k$  of a quantity proportional to the gradient of the objective respect to the design variables. This proportionality is controlled by the parameter  $\alpha$ , which is fixed at the beginning of the optimization and remains constant throughout the process.

Concerning the classification presented in this subsection we can frame (3.2.7) as an optimization algorithm that is based on a mathematical rationale with a mix of heuristics introduced by the choice of the parameter  $\alpha$ . It is also a first-order algorithm since the update require explicitly the gradient of the objective. The primary drivers behind the choice of this optimizer were its mathematical foundations, its simple implementation and its effectiveness with linear and quadratic objectives. Of course others approaches exists and may be more or less effective depending on the type of the problem that has to be faced. Some examples are:

**Brute-force** In this case all possible combinations of design parameters are evaluated. This approach could be useful in case of very simple design optimization



**Figure 3.5:** Example of design optimization: baseline and final design of ONERA M6 wing is shown. It can be seen how the relative pressure exerted by the air on the wing has been minimized. Figure taken from [31]

problems with few design variables, but results to be time-expensive for large-scale projects where parameters can be hundreds, thousands, or even more and the optimal design may not found in reasonable time-frames.

**Genetic Algorithms** These are particularly useful when objective and constraint functions are not sufficiently smooth or derivatives can not be computed with enough precision. A key trait of these algorithms is that they require many iterations for convergence, which is sensitive to the cost of the function evaluations, and scales poorly with the number of design parameters [30].

### 3.2.4 Design optimization with RBF-FD models

One of the main areas where the application of RBF-FD method bring great benefits is Computational Fluid Dynamics (CFD) where PDEs are employed to model fluid flow and heat transfer problems. In CFD, in general, many design optimization problems arise. A typical case is wing design where an example of the results that can be achieved is presented in figure 3.5.

In section 2.6 on page 11 we have seen how the the RBF-FD solver for Partial Differential Equations (PDEs), which can be employed to solve CFD related equations as well, is implemented: starting from a point cloud distributed over the physical domain, it yields a set of equations which once solved allows to find the values of the approximating field on the aforementioned nodes. In this scenario a PDE is, in every respect, what we have defined as mathematical model or governing

equations when discussing optimization problem in previous section. Moreover, for these cases, is easy to identify the numerical model as the one obtained by the RBF-FD method since it actually discretize the PDE.

Keeping in mind what is stated above a typical optimization problem which depends on the solution, found by means of an RBF-FD solver, of a model based on a PDE, can be stated as:

$$\begin{aligned}
 & \text{minimize} && J(\mathbf{l}, \mathbf{u}) \\
 & \text{by varying} && l_i && i = 1, \dots, N \\
 & \text{subject to} && g_j(\mathbf{l}; \mathbf{u}) \leq 0 && j = 1, \dots, n_g \\
 & && h_k(\mathbf{l}; \mathbf{u}) = 0 && k = 1, \dots, n_h \\
 & && l_i \leq \bar{l}_i && i = 1, \dots, N \\
 & && \mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B - \mathbf{f} = \mathbf{0}
 \end{aligned} \tag{3.2.8}$$

where in the last constraint has been used the same notation presented in chapter 2. For its solution it can be implemented the scheme reported in figure 3.4.

In the particular case when the optimizer block update the design following the rule reported in (3.2.7) the derivative of the objective  $J$  with respect to the design variables needs to be computed. In general this can be computed through the application of the chain rule as:

$$\frac{dJ^T}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{l}} + \frac{\partial J^T}{\partial \mathbf{l}} \tag{3.2.9}$$

where the notations  $d/dx$  and  $\partial/\partial x$  indicate the total and the partial derivatives respect the variable  $x$ . And if terms  $\partial J/\partial \mathbf{u}$  and  $\partial J/\partial \mathbf{l}$  can be efficiently obtained by means of reverse-mode AD, as explained in section 3.1.2 on page 28, since  $J$  is known analytically, matrix  $\partial \mathbf{u}/\partial \mathbf{l}$ , require more careful considerations, since the differentiation of the constrains  $\mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B - \mathbf{f} = \mathbf{0}$  is required for its computation. In the following section we explain the details for the computations of  $dJ/d\mathbf{l}$  using the adjoint method, which simplify the optimization of problems with RBF-FD-like constraints.

### 3.3 Adjoint method in RBF-FD based Design Optimization

In previous section we had a general overview on design optimization problems. In this section we will see how those problems based on RBF-FD models can be solved using the adjoint method in the optimization step. In particular we focus on two specific types of problems respectively in one and three-dimensional spaces.

Before that we will explain more generally how the adjoint method can be employed in the computation of  $dJ/d\mathbf{l}$ , used for the update of the design variables at each step, and which benefits it brings.

### 3.3.1 Adjoint intro

We start by considering a generic design optimization problem with the following formulation:

$$\begin{aligned} & \text{minimize} && J(\mathbf{l}, \mathbf{u}) \\ & \text{by varying} && l_i \quad i = 1, \dots, N \\ & \text{subject to} && \mathbf{A}\mathbf{u} = \mathbf{b} \end{aligned} \quad (3.3.1)$$

where the optimization algorithm used for its solution is the one that employs equation (3.2.7), which require the computation of the gradient of the objective, to update the design parameters at each step.

To directly evaluate  $dJ/d\mathbf{l}$ , in case of a particular set of design parameters defined by  $\mathbf{l}$  at a given state  $\mathbf{u}$ , following equation (3.2.9), we would compute:

$$\frac{dJ^T}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{l}} + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.3.2)$$

$\partial J/\partial \mathbf{u}$  and  $\partial J/\partial \mathbf{l}$  can be efficiently obtained by means of reverse-mode AD, as explained in the previous subsection. The matrix  $d\mathbf{u}/d\mathbf{l}$ , on the other hand, requires more careful considerations which we address here.

Differentiating the problem constraints  $\mathbf{A}\mathbf{u} = \mathbf{b}$  with respect to the parameters  $\mathbf{l}$  gives:

$$\frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} + \mathbf{A}\frac{d\mathbf{u}}{d\mathbf{l}} = \frac{d\mathbf{b}}{d\mathbf{l}} \quad (3.3.3)$$

which allows to obtain the sought term  $d\mathbf{u}/d\mathbf{l}$  as:

$$\frac{d\mathbf{u}}{d\mathbf{l}} = \mathbf{A}^{-1} \left( \frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} \right) \quad (3.3.4)$$

We notice that  $d\mathbf{u}/d\mathbf{l}$  is a matrix of size  $M \times N$  whose  $j$ -th column represent the sensitivity of the state  $\mathbf{u}$  respect to the parameter  $l_j$ . In practice it is populated column-wise by solving equation (3.3.4)  $N$  times when the derivative is taken with respect the  $j$ -th design parameter  $l_j$  rather than  $\mathbf{l}$ : this gives its  $j$ -th column. In this process  $d\mathbf{b}/d\mathbf{l}$  and  $d\mathbf{A}/d\mathbf{l}$  terms can be simply obtained using AD, if  $\mathbf{b}$  and  $\mathbf{A}$  are known analytically, and once  $d\mathbf{u}/d\mathbf{l}$  is computed it can be substituted in equation (3.2.9) to obtain the gradient which symbolically reads as:

$$\frac{dJ}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}} \left[ \mathbf{A}^{-1} \left( \frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} \right) \right] + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.3.5)$$



## CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

However the computational costs of this naive procedure scales linearly with the number of the design parameters  $N$ , since  $N$  solutions of linear systems are required in order to construct  $d\mathbf{u}/d\mathbf{l}$ . We finally remark that each of these inversions has the same computational cost of solving  $\mathbf{A}\mathbf{u} = \mathbf{b}$  for  $\mathbf{u}$ . In a scenario where CFD simulations are involved this is particularly penalizing, therefore a more efficient gradient computation is required.

To do so we can employ the *adjoint method* which simply consist of a smarter bracketing of the equation which gives the gradient. In fact is possible to rewrite (3.3.5) as:

$$\begin{aligned}\frac{dJ}{d\mathbf{l}} &= \left[ \frac{\partial J^T}{\partial \mathbf{u}} \mathbf{A}^{-1} \right] \left( \frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right) + \frac{\partial J^T}{\partial \mathbf{l}} \\ &= \boldsymbol{\lambda}^T \left( \frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right) + \frac{\partial J^T}{\partial \mathbf{l}}\end{aligned}\tag{3.3.6}$$

where the vector  $\boldsymbol{\lambda} \in \mathbb{R}^L$ , whose elements are called *adjoint variables*, can be found by solving:

$$\mathbf{A}^T \boldsymbol{\lambda} = \frac{\partial J}{\partial \mathbf{u}}\tag{3.3.7}$$

By doing so is possible to solve the system in equation (3.3.7), which is called *adjoint problem*, only *once* independently on the number  $N$  of design parameters for each gradient computation: this significantly reduces the computational burden of the whole optimization process. Bear in mind that the adjoint problem has the same size as the problem defined by the constraints  $\mathbf{A}\mathbf{u} = \mathbf{b}$  and the computational cost for their solution is the same.

The result of the whole process is then summarized by:

$$\boldsymbol{\lambda}^T \left( \frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right)\tag{3.3.8}$$

where  $\boldsymbol{\lambda}$  solves equation (3.3.7) and the Jacobian  $d\mathbf{u}/d\mathbf{l}$  is no more present. From the process outlined here, we can observe some commonalities between the adjoint method and reverse mode automatic differentiation (AD):

- both apply the chain rule to decompose the computation of the gradient;
- both propagate derivatives in reverse: reverse mode AD propagates from a single output of a function back to its inputs, the adjoint does the same from the objective back to the design variables;
- both are efficient in case of large number of inputs.

## CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

Furthermore adjoint method also share a two-steps process with the reverse mode AD:

- the forward pass, used for the computation of variables  $v_i$  in table 3.2 on page 29, now is the single solution of the system  $\mathbf{A}\mathbf{u} = \mathbf{b}$  in order to find the state  $\mathbf{u}$  used in term (3.3.8);
- the reverse pass, done in order to obtain the adjoint variables  $\bar{v}_i$  in table 3.2, is now the solution of  $\mathbf{A}^T \boldsymbol{\lambda} = \frac{\partial J}{\partial \mathbf{u}}$  (we would like to reiterate again that reverse and forward pass has the *same* computational complexity).

After this concise overview of the adjoint method applied to generic design optimization problems, we now proceed to the next section where we examine in more detail its application to those problems where the state  $\mathbf{u}$  is obtained solving a Partial Differential Equation by means of RBF-FD methods.

### 3.3.2 1D case

In this case we consider a generic cost function  $J: \mathbb{R}^{N_I \times L} \rightarrow \mathbb{R}$  which depends on the  $N_I$  values of the field  $u$  at the points located within the physical domain and on the  $L$  design parameters. The resulting optimization problem is:

$$\begin{aligned} & \text{minimize} && J(\mathbf{l}, \mathbf{u}) \\ & \text{by varying} && l_i && i = 1, \dots, N \\ & \text{subject to} && \mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B - \mathbf{f} = \mathbf{0} \end{aligned} \quad (3.3.9)$$

Thanks to this generality its sensitivities with respect to the parameters, are given by the same exact formula reported in (3.2.9):

$$\frac{dJ^T}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}_I} \frac{d\mathbf{u}_I}{d\mathbf{l}} + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.3.10)$$

On the above formula the term that needs to be computed carefully is the matrix  $d\mathbf{u}_I/d\mathbf{l}$ . The sensitivities of the state with respect to the parameters are found by differentiating equation (2.6.20) which gives:

$$\frac{d\mathbf{C}_I}{d\mathbf{l}} \mathbf{u}_I + \mathbf{C}_I \frac{d\mathbf{u}_I}{d\mathbf{l}} + \frac{d\mathbf{C}_B}{d\mathbf{l}} \mathbf{u}_B + \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} - \frac{d\mathbf{f}}{d\mathbf{l}} = \mathbf{0} \quad (3.3.11)$$

Before proceeding a note on the matrix  $\frac{d\mathbf{u}_B}{d\mathbf{l}}$  could be useful: one might think that it should not be present since values in  $\mathbf{u}_B$  are fixed from boundary conditions of problem (2.6.5) (i.e. they are not dependent on the design parameters). Nevertheless the previous statement does not hold in general since the values of  $\mathbf{u}_B$  could directly

### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

depend on the design variables (as an example one may think of Dirichlet Boundary Conditions (BCs) defined by  $u(\mathbf{x}) = g(\mathbf{x}, \mathbf{l})$ ) or indirectly depend on them through the problem geometry (e.g. Neumann or Robin BCs). Rearranging (3.3.11) we obtain the sought term:

$$\frac{d\mathbf{u}_I}{d\mathbf{l}} = \mathbf{C}_I^{-1} \left( \frac{d\mathbf{f}}{d\mathbf{l}} - \frac{d\mathbf{C}_I}{d\mathbf{l}} \mathbf{u}_I - \frac{d\mathbf{C}_B}{d\mathbf{l}} \mathbf{u}_B - \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} \right) \quad (3.3.12)$$

Then plugging the found expression for  $d\mathbf{u}_I/d\mathbf{l}$  in (3.3.10) yields:

$$\begin{aligned} \frac{dJ^T}{d\mathbf{l}} &= \frac{\partial J}{\partial \mathbf{u}_I}^T \left[ \mathbf{C}_I^{-1} \left( \frac{d\mathbf{f}}{d\mathbf{l}} - \frac{d\mathbf{C}_I}{d\mathbf{l}} \mathbf{u}_I - \frac{d\mathbf{C}_B}{d\mathbf{l}} \mathbf{u}_B - \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} \right) \right] + \frac{\partial J}{\partial \mathbf{l}}^T \\ &= \frac{\partial J}{\partial \mathbf{u}_I}^T \left[ \mathbf{C}_I^{-1} \left( \frac{d\mathbf{f}}{d\mathbf{l}} - \frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u} - \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} \right) \right] + \frac{\partial J}{\partial \mathbf{l}}^T \end{aligned} \quad (3.3.13)$$

where  $d\mathbf{C}/d\mathbf{l} \in \mathbb{R}^{N_I \times N}$  is the matrix resulting from the concatenation of the rows of  $d\mathbf{C}_I/d\mathbf{l} \in \mathbb{R}^{N_I \times N_I}$  and  $d\mathbf{C}_B/d\mathbf{l} \in \mathbb{R}^{N_B \times N_B}$ , and  $\mathbf{u} \in \mathbb{R}^N$  is the vector given by the concatenation of the elements of  $\mathbf{u}_I \in \mathbb{R}^{N_I}$  and  $\mathbf{u}_B \in \mathbb{R}^{N_B}$ . Now the adjoint method can be employed in order to avoid performing the inversion of matrix  $\mathbf{C}_I$  once for each parameter in  $\mathbf{l}$ , yielding:

$$\frac{dJ^T}{d\mathbf{l}} = \boldsymbol{\lambda}_1^T \left( \frac{d\mathbf{f}}{d\mathbf{l}} - \frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u} - \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} \right) + \frac{\partial J}{\partial \mathbf{l}}^T \quad (3.3.14)$$

where  $\boldsymbol{\lambda}_1$  is found by solving *once*:

$$\mathbf{C}_I^T \boldsymbol{\lambda}_1 = \frac{\partial J}{\partial \mathbf{u}_I} \quad (3.3.15)$$

The unknown terms involving derivatives computation on the right-hand side can be computed easily through AD except for  $\frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u}$ . This term deserves more attention since it requires the differentiation of the global matrix  $\mathbf{C}$  which is obtained by solving  $N_I$  local systems, each associated with a different stencil  $\mathcal{X}_i$ .

In order to examine it more closely we can define a matrix  $\mathbf{Q} \in \mathbb{R}^{N_I \times L}$  as follows:

$$\mathbf{Q} = \frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u} \quad (3.3.16)$$

whose elements  $q_{i,j}$  are given by:

$$q_{i,j} = \frac{d\mathbf{C}_{[i,:]} }{dl_j} \mathbf{u} \quad (3.3.17)$$

where notation  $\mathbf{C}_{[i,:]}$  is used to indicate the  $i$ -th row of the matrix  $\mathbf{C}$ . To obtain  $\frac{d\mathbf{C}}{dl}\mathbf{u}$  is thus sufficient to compute each element of  $\mathbf{Q}$  through equation (3.3.17). However, before doing so, it is worth rewriting the equation for  $q_{i,j}$  more explicitly.

Recall that elements that make up the  $i$ -th row of  $\mathbf{C}$  are found by solving the local system (2.6.15):

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (3.3.18)$$

in case of RBF-FD method, or the local system (2.7.20):

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}_1\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}_1\mathcal{B}_2\Phi(\mathbf{x}_i, \mathcal{X}_{i,B}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (3.3.19)$$

in case of RBF-HFD method. The key aspect that both cases share is that the aforementioned equations arise from the information contained in a single stencil  $\mathcal{X}_i$  and, as explained in subsection 2.6.2 on page 13, only the first  $m$  elements of  $\mathbf{c}(\mathbf{x}_i) = [\mathbf{c}_I(\mathbf{x}_i), \mathbf{c}_B(\mathbf{x}_i), \mathbf{c}_p(\mathbf{x}_i)]$  form the  $i$ -th row of  $\mathbf{C}$ . Furthermore, since  $m \ll N$ , row  $\mathbf{C}_{[i,:]}$  is sparse; which in turns implies that not all the elements of  $\mathbf{u}$  are involved in the computation of  $q_{i,j}$ .

By letting  $\tilde{\mathbf{u}}_i \in \mathbb{R}^m$  the vector composed by the components of  $\mathbf{u}$  associated to the non zero elements of  $\mathbf{C}_{[i,:]}$ , we can rewrite the elements of  $\mathbf{Q}$  reported in (3.3.17) as:

$$q_{i,j} = \frac{d\mathbf{c}(\mathbf{x}_i)^T}{dl_j} \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.20)$$

where the zero-vector concatenated to  $\tilde{\mathbf{u}}_i$  has the same length of  $\mathbf{c}_p(\mathbf{x}_i)$ . If we now differentiate equation (3.3.19) with respect to the  $j$ -th design parameter, we obtain:

$$\frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) + \mathbf{M}_{BC}^T \frac{d\mathbf{c}(\mathbf{x}_i)}{dl_j} = \frac{d\mathbf{h}}{dl_j} \quad (3.3.21)$$

where the vector  $\mathbf{h} \in \mathbb{R}^{m+M}$  is used as a shorthand for the vector present on the right-hand side of equation (3.3.19). Now from the last equation we can isolate:

$$\frac{d\mathbf{c}(\mathbf{x}_i)^T}{dl_j} = \left( \frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \mathbf{M}_{BC}^{-1} \quad (3.3.22)$$

which once plugged in (3.3.20) allows to rewrite  $q_{i,j}$  as:

$$q_{i,j} = \left[ \left( \frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \mathbf{M}_{BC}^{-1} \right] \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.23)$$

But now we can clearly see that in order to populate the matrix  $\mathbf{Q}$  we should invert  $N_I \times L$  matrices  $\mathbf{M}_{BC}$  of size  $(m + M) \times (m + M)$  since we need to compute  $\frac{d\mathbf{c}(\mathbf{x}_i)^T}{dl_j}$ . These matrices are smaller than  $\mathbf{C}_I$ , but in any case their inversion is still problematic, since both the number of parameters and the number of nodes are huge. The problem, here, is very similar to the one faced in the previous subsection where there was a linear relationship between the number of matrix inversions and the number of design variables: consequently, even here, to improve the situation it is possible to rely on the adjoint method. Modifying the parentheses of equation (3.3.23) it is possible to write:

$$\begin{aligned} q_{i,j} &= \left( \frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \left( \mathbf{M}_{BC}^{-1} \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \right) \\ &= \left( \frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \boldsymbol{\lambda}_{2,i} \end{aligned} \quad (3.3.24)$$

where the adjoint vector  $\boldsymbol{\lambda}_{2,i} \in \mathbb{R}^{m+M}$  is found as solution of:

$$\mathbf{M}_{BC} \boldsymbol{\lambda}_{2,i} = \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.25)$$

and has to be computed once for each row of  $\mathbf{Q}$  (or equivalently  $\frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u}$ ). The crucial point is that this second adjoint implementation allows to make the computational cost to obtain  $\frac{d\mathbf{C}}{d\mathbf{l}} \mathbf{u}$  independent of the number of design parameters  $L$ : the number of systems (3.3.25) that has to be solved, whose cost is similar to the inversion of  $\mathbf{M}_{BC}$ , is now always  $N_I$ .

In light of what we have observed we are now able to compute the whole gradient  $\frac{dJ^T}{d\mathbf{l}}$  with a computation effort analogous to two RBF-HFD:

1. a first plain application of RBF-HFD is required in order to compute the terms  $\mathbf{u}$ ,  $\mathbf{C}_I$  and  $\mathbf{C}_B$ ;
2. then, a comparable cost is required for the application of the presented adjoint method as shown in table 3.3.

### 3.3.3 3D case

Now, unlike in the 1D case, we consider only cost functions which describes a flux integral, i.e. functions that can be reduced into a form similar to:

$$J_c = \frac{1}{|A_c|} \int_A \frac{\partial u}{\partial \mathbf{n}} ds \quad (3.3.26)$$

# CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

**Table 3.3:** Systems that demands the computation of a solution during the application of RBF-FD and the adjoint method explained in section 3.3.2. Those of dimension  $m + M$  are solved once for each stencil in the physical domain

| Method  | Solved systems  | System dimension |
|---------|---|------------------|
| RBF-FD  | $\mathbf{C}_I \mathbf{u}_I = f - \mathbf{C}_B \mathbf{u}_B$   | $N_I$            |
|         | $\mathbf{M}_{BC}^T \mathbf{c}(\mathbf{x}_i) = \mathbf{h}$ for $i = 1, \dots, N_I$   | $m + M$          |
| Adjoint | $\mathbf{C}_I^T \boldsymbol{\lambda}_1 = \frac{\partial J}{\partial \mathbf{u}_I}$  | $N_I$            |
|         | $\mathbf{M}_{BC} \boldsymbol{\lambda}_{2,i} = \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix}$ for $i = 1, \dots, N_I$ | $m + M$          |

where  $|A_c| \in \mathbb{R}$  indicates the area of the surface  $A$  over which the integral is calculated,  $\mathbf{n}$  is the normal of the surface and  $u$  is a scalar field.

Since in our case the physical domain is described by means of `.stl` files, we will not deal with continuous quantities and the integral will therefore be approximated by the following finite sum:

$$J = \frac{1}{|A|} \sum_{j \in \mathcal{T}} \nabla \mathbf{u}(\mathbf{x}_j)^T \mathbf{a}_j \quad (3.3.27)$$

where  $|A| \in \mathbb{R}$  is still the area of the surface  $A$ , this time computed as sum of the areas of each small triangle that composes it,  $\mathcal{T}$  is the set of the indices of these triangles and each of them is identified by a single index  $j$ ;  $\mathbf{x}_j \in \mathbb{R}^3$  and  $\mathbf{a}_j \in \mathbb{R}^3$  denotes respectively the centroid and the area vector of the  $j$ -th triangle, where  $\mathbf{a}_j$  has direction given by the normal and modulus equal to the surface of the  $j$ -th triangle.  $\nabla \mathbf{u}(\mathbf{x}_j)^T \in \mathbb{R}^3$  indicates the row-vector representing the gradient of the field  $u$  computed at centroid  $\mathbf{x}_j$ . Furthermore, we recall that  $J$  is a scalar field exactly as in the case reported in previous subsection, thus  $J: \mathbb{R}^{N_I \times L} \rightarrow \mathbb{R}$ , where  $N_I$  is the number of RBF-FD nodes that lie inside the physical domain and  $L$  is the number of design parameters that we are able to control; for this specific problem design parameters  $\mathbf{l}$  are the heights of the vertices of the triangles that define surface  $A$ . Finally we denote with  $N_C$  the number of `stl` triangles that made up surface  $A$ , thus  $|\mathcal{T}| = N_C$ . Even in this case the design optimization problem obtained assumes the same form of the one presented in the 1D case, reported in equation (3.3.9).

Before proceeding we point out how  $\nabla \mathbf{u}(\mathbf{x}_j)^T$  is computed. At first one might think to obtain it using the RBF-FD method based on the  $N$  nodes scattered over the physical domain as explained in subsection 2.6.2 on page 13. However this is not

### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

the right approach, in fact by following the aforementioned procedure is possible to approximate the gradient on the RBF-FD nodes only and since the centroids of the triangles are independent from the latter another approach is required. RBF-FD method it is still a way to go since by considering the stencil  $\mathcal{X}_j$  centered on  $\mathbf{x}_j$  is possible to refer back to equation (2.6.16) and then approximate  $\nabla \mathbf{u}(\mathbf{x}_j)^T$  as:

$$\begin{aligned}\nabla \mathbf{u}(\mathbf{x}_j)^T &= \begin{bmatrix} \frac{\partial}{\partial x} u(\mathbf{x}_j) & \frac{\partial}{\partial y} u(\mathbf{x}_j) & \frac{\partial}{\partial z} u(\mathbf{x}_j) \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{c}_x^T(\mathbf{x}_j) \mathbf{u}_j & \mathbf{c}_y^T(\mathbf{x}_j) \mathbf{u}_j & \mathbf{c}_z^T(\mathbf{x}_j) \mathbf{u}_j \end{bmatrix}\end{aligned}\quad (3.3.28)$$

where  $\mathbf{u}_j \in \mathbb{R}^m$  is the vector containing the field values at the points of stencil  $\mathcal{X}_j$  while  $\mathbf{c}_x^T(\mathbf{x}_j)$ ,  $\mathbf{c}_y^T(\mathbf{x}_j)$  and  $\mathbf{c}_z^T(\mathbf{x}_j)$  are the vectors composed by the first  $m$  coefficients obtained from the solution of equation (2.6.15) respectively when the associated differential operator  $\mathcal{L}$  in the equation is  $\partial/\partial x$ ,  $\partial/\partial y$  or  $\partial/\partial z$ .

Moving forward, if we further concatenate along rows the vectors  $\nabla \mathbf{u}(\mathbf{x}_j)^T$  associated to each triangle we would obtain the matrix:

$$\nabla \mathbf{u} = \begin{bmatrix} \mathbf{c}_x^T(\mathbf{x}_1) \mathbf{u}_1 & \mathbf{c}_y^T(\mathbf{x}_1) \mathbf{u}_1 & \mathbf{c}_z^T(\mathbf{x}_1) \mathbf{u}_1 \\ \vdots & \vdots & \vdots \\ \mathbf{c}_x^T(\mathbf{x}_{N_C}) \mathbf{u}_{N_C} & \mathbf{c}_y^T(\mathbf{x}_{N_C}) \mathbf{u}_{N_C} & \mathbf{c}_z^T(\mathbf{x}_{N_C}) \mathbf{u}_{N_C} \end{bmatrix}\quad (3.3.29)$$

which belongs to  $\mathbb{R}^{N_C \times 3}$  and can be rewritten more compactly as:

$$\nabla \mathbf{u} = [\mathbf{C}_x \mathbf{u} \quad \mathbf{C}_y \mathbf{u} \quad \mathbf{C}_z \mathbf{u}]\quad (3.3.30)$$

where  $\mathbf{u} \in \mathbb{R}^N$  is the vector composed by the field values on the  $N$  RBF-FD nodes and all the matrices  $\mathbf{C}_x$ ,  $\mathbf{C}_y$  and  $\mathbf{C}_z$  are sparse rectangular matrices with  $N_C$  rows and  $N$  columns; the  $i$ -th row of matrix  $\mathbf{C}_x$  is formed by the elements of the vector  $\mathbf{c}_x^T(\mathbf{x}_i)$  present in equation (3.3.29) and in similar way the  $i$ -th rows of matrices  $\mathbf{C}_y$  and  $\mathbf{C}_z$  are also structured, where the elements of vectors  $\mathbf{c}_y^T(\mathbf{x}_i)$  and  $\mathbf{c}_z^T(\mathbf{x}_i)$  are used instead of those of vector  $\mathbf{c}_x^T(\mathbf{x}_i)$ .

By defining the following three vectors:

$$\mathbf{A}_x = \begin{bmatrix} a_{1,x} \\ \vdots \\ a_{N_C,x} \end{bmatrix} \quad \mathbf{A}_y = \begin{bmatrix} a_{1,y} \\ \vdots \\ a_{N_C,y} \end{bmatrix} \quad \mathbf{A}_z = \begin{bmatrix} a_{1,z} \\ \vdots \\ a_{N_C,z} \end{bmatrix}\quad (3.3.31)$$

where the subscripts  $x$ ,  $y$ ,  $z$  accompanying  $a_1, \dots, a_{N_C}$  are used to indicate the elements in the first, second and third position of the vectors  $\mathbf{a}_1, \dots, \mathbf{a}_{N_C}$ , is then possible to rewrite the cost function  $J$  as:

$$J = \frac{1}{|A|} (\mathbf{A}_x^T \mathbf{C}_x \mathbf{u} + \mathbf{A}_y^T \mathbf{C}_y \mathbf{u} + \mathbf{A}_z^T \mathbf{C}_z \mathbf{u})\quad (3.3.32)$$

### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

Additionally, defining the vector  $\mathbf{k}^T = \mathbf{A}_x^T \mathbf{C}_x + \mathbf{A}_y^T \mathbf{C}_y + \mathbf{A}_z^T \mathbf{C}_z$ , which belongs to  $\mathbb{R}^N$ , equation (3.3.32) can be rewritten in an even more compact manner as:

$$J = \frac{1}{|A|} \mathbf{k}^T \mathbf{u} \quad (3.3.33)$$

Now that we have rewritten the cost function in a more convenient way we can compute its gradient. To do so we have to differentiate equation (3.3.33) with respect to each design parameter; focusing on the  $i$ -th parameter we obtain the following derivative:

$$\frac{\partial J}{\partial l_i} = \left( \frac{\partial}{\partial l_i} |A|^{-1} \right) \mathbf{k}^T \mathbf{u} + \frac{1}{|A|} \left( \frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u} + \mathbf{k}^T \frac{\partial \mathbf{u}}{\partial l_i} \right) \quad (3.3.34)$$

which once a distinction between the field values corresponding to internal and boundary nodes, denoted respectively as  $\mathbf{u}_I \in \mathbb{R}^{N_I}$  and  $\mathbf{u}_B \in \mathbb{R}^{N_B}$ , is made, can be rewritten as:

$$\frac{\partial J}{\partial l_i} = \left( \frac{\partial}{\partial l_i} |A|^{-1} \right) \mathbf{k}^T \mathbf{u} + \frac{1}{|A|} \left( \frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u} + \mathbf{k}_I^T \frac{\partial \mathbf{u}_I}{\partial l_i} + \mathbf{k}_B^T \frac{\partial \mathbf{u}_B}{\partial l_i} \right) \quad (3.3.35)$$

where  $\mathbf{k}_I^T$  and  $\mathbf{k}_B^T$  are the vectors made up of the elements of  $\mathbf{k}^T$  associated respectively to  $\mathbf{u}_I$  and  $\mathbf{u}_B$ .

The term  $\partial \mathbf{u}_I / \partial l_i$  can be found in a similar way as done in the previous subsection in 1D case: differentiating the RBF-FD constraint reported in equation (2.6.20). Doing so we obtain again:

$$\frac{d\mathbf{C}_I}{dl} \mathbf{u}_I + \mathbf{C}_I \frac{d\mathbf{u}_I}{dl} + \frac{d\mathbf{C}_B}{dl} \mathbf{u}_B + \mathbf{C}_B \frac{d\mathbf{u}_B}{dl} = \frac{d\mathbf{f}}{dl} \quad (3.3.36)$$

which once replaced in equation (3.3.34) gives:

$$\begin{aligned} \frac{\partial J}{\partial l_i} = & \left( \frac{\partial}{\partial l_i} |A|^{-1} \right) \mathbf{k}^T \mathbf{u} \\ & + \frac{1}{|A|} \left[ \frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u} + \mathbf{k}_B^T \frac{\partial \mathbf{u}_B}{\partial l_i} - \mathbf{k}_I^T \mathbf{C}_I^{-1} \left( \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u} + \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i} \right) \right] \end{aligned} \quad (3.3.37)$$

where matrix  $\partial \mathbf{C} / \partial l_i$  is obtained by stacking vertically the matrices  $\partial \mathbf{C}_I / \partial l_i$  and  $\partial \mathbf{C}_B / \partial l_i$  whereas vector  $\mathbf{u}$  is given by doing the same but with vectors  $\mathbf{u}_I$  and  $\mathbf{u}_B$ . Then is possible to apply the adjoint method to compute the product  $\mathbf{k}_I^T \mathbf{C}_I^{-1}$  which makes equation (3.3.37):

$$\begin{aligned} \frac{\partial J}{\partial l_i} = & \left( \frac{\partial}{\partial l_i} |A|^{-1} \right) \mathbf{k}^T \mathbf{u} \\ & + \frac{1}{|A|} \left[ \frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u} + \mathbf{k}_B^T \frac{\partial \mathbf{u}_B}{\partial l_i} - \lambda_1^T \left( \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u} + \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i} \right) \right] \end{aligned} \quad (3.3.38)$$



### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

where  $\boldsymbol{\lambda}_1$  is found by solving:

$$\mathbf{C}_I^T \boldsymbol{\lambda}_1 = \mathbf{k}_I \quad (3.3.39)$$

and remain the same even when the derivative is computed respect a parameter different than  $l_i$ .

Now what remains to be computed are the 3 unknown terms that appear within the square brackets of equation (3.3.38). These are:  $\frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u}$ ,  $\mathbf{k}_B^T \frac{\partial \mathbf{u}_B}{\partial l_i}$  and  $\left( \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u} + \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i} \right)$ .

With regard to the term  $\frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u}$ , we can take advantage of the equivalence between equations (3.3.32) and (3.3.33) to write:

$$\begin{aligned} \frac{\partial \mathbf{k}^T}{\partial l_i} \mathbf{u} &= \frac{\partial}{\partial l_i} (\mathbf{A}_x^T \mathbf{C}_x + \mathbf{A}_y^T \mathbf{C}_y + \mathbf{A}_z^T \mathbf{C}_z) \mathbf{u} \\ &= \underbrace{\left( \frac{\partial \mathbf{A}_x^T}{\partial l_i} \mathbf{C}_x + \frac{\partial \mathbf{A}_y^T}{\partial l_i} \mathbf{C}_y + \frac{\partial \mathbf{A}_z^T}{\partial l_i} \mathbf{C}_z \right)}_{s_1} \mathbf{u} \\ &\quad + \underbrace{\left( \mathbf{A}_x^T \frac{\partial \mathbf{C}_x}{\partial l_i} + \mathbf{A}_y^T \frac{\partial \mathbf{C}_y}{\partial l_i} + \mathbf{A}_z^T \frac{\partial \mathbf{C}_z}{\partial l_i} \right)}_{s_2} \mathbf{u} \end{aligned} \quad (3.3.40)$$

Therefore we only need to find the values of  $s_1$  and  $s_2$ . In the  $s_1$  case, matrices  $\mathbf{C}_x$ ,  $\mathbf{C}_y$ ,  $\mathbf{C}_z$  and the vector  $\mathbf{u}$  are already known from the application of the initial global RBF-FD method, whereas the row vectors  $\partial \mathbf{A}_x^T / \partial l_i$ ,  $\partial \mathbf{A}_y^T / \partial l_i$  and  $\partial \mathbf{A}_z^T / \partial l_i$  can be computed via automatic differentiation since they contain only geometrical information. Alternatively,  $s_1$  can be computed by accumulating the contributions  $s_{1,j}$  of each sensitivity of the cost function with respect to each design parameter (i.e. triangle), which is given by

$$s_{1,j} = \nabla \mathbf{u}(\mathbf{x}_j)^T \frac{\partial \mathbf{a}_j}{\partial l_i} \quad (3.3.41)$$

where  $\partial \mathbf{a}_j / \partial l_i$  can be computed by concatenating the elements in position  $j$  of vectors  $\partial \mathbf{A}_x^T / \partial l_i$ ,  $\partial \mathbf{A}_y^T / \partial l_i$  and  $\partial \mathbf{A}_z^T / \partial l_i$ .

On the other hand,  $s_2$ , can be computed by applying the adjoint method to the local RBF-FD systems. Following a similar approach to what was done for  $s_1$ , the contribution given by a single triangle to  $s_2$  is given by:

$$s_{2,j} = a_{j,x} \frac{\partial \mathbf{c}_x^T(\mathbf{x}_j)}{\partial l_i} \mathbf{u}_j + a_{j,y} \frac{\partial \mathbf{c}_y^T(\mathbf{x}_j)}{\partial l_i} \mathbf{u}_j + a_{j,z} \frac{\partial \mathbf{c}_z^T(\mathbf{x}_j)}{\partial l_i} \mathbf{u}_j \quad (3.3.42)$$

### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

Now we recall that  $\mathbf{c}_x^T(\mathbf{x}_j)$  is found as solution of the following local RBF-FD system:

$$\mathbf{M}_{BC,j}^T \mathbf{c}_x(\mathbf{x}_j) = \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_j, \mathcal{X}_{j,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_j) \end{bmatrix} \quad (3.3.43)$$

which means that  $\partial \mathbf{c}_x(\mathbf{x}_j) / \partial l_i$  can be found differentiating the equation above with respect to the  $i$ -th design parameter, resulting in:

$$\frac{\partial \mathbf{c}_x(\mathbf{x}_j)}{\partial l_i} = \mathbf{M}_{BC,j}^{-T} \left( \frac{\partial \mathbf{h}_{j,x}}{\partial l_i} - \frac{\partial \mathbf{M}_{BC,j}^T}{\partial l_i} \mathbf{c}_x(\mathbf{x}_j) \right) \quad (3.3.44)$$

where we have used the notation  $\mathbf{h}_{j,x}$  to indicate the right-hand side of equation (3.3.43). Substituting what we have just found into the equation (3.3.42), rearranging the result and then applying the adjoint method finally yields:

$$\begin{aligned} s_{2,j} = & \left( a_{j,x} \frac{\partial \mathbf{h}_{j,x}^T}{\partial l_i} + a_{j,y} \frac{\partial \mathbf{h}_{j,y}^T}{\partial l_i} + a_{j,z} \frac{\partial \mathbf{h}_{j,z}^T}{\partial l_i} \right) \boldsymbol{\lambda}_{2,j} \\ & - \left( a_{j,x} \mathbf{c}_x^T(\mathbf{x}_j) + a_{j,y} \mathbf{c}_y^T(\mathbf{x}_j) + a_{j,z} \mathbf{c}_z^T(\mathbf{x}_j) \right) \frac{\partial \mathbf{M}_{BC,j}}{\partial l_i} \boldsymbol{\lambda}_{2,j} \end{aligned} \quad (3.3.45)$$

where  $\boldsymbol{\lambda}_{2,j} \in \mathbb{R}^{m+M}$  is found by solving the adjoint system:

$$\mathbf{M}_{BC,j} \boldsymbol{\lambda}_{2,j} = \begin{bmatrix} \mathbf{u}_j \\ \mathbf{0} \end{bmatrix} \quad (3.3.46)$$

once, irrespective of the number  $L$  of design parameters.

The term  $\mathbf{k}_B^T \frac{\partial \mathbf{u}_B}{\partial l_i}$  does not require any particular care except for the computation of  $\frac{\partial \mathbf{u}_B}{\partial l_i}$  whose elements are found by applying the RBF-FD procedure. On the other hand,  $\mathbf{k}_B^T$ , can be found from its definition since  $\mathbf{A}_x$ ,  $\mathbf{A}_y$ ,  $\mathbf{A}_z$  and  $\mathbf{C}_x$ ,  $\mathbf{C}_y$ ,  $\mathbf{C}_z$  are known.

The last term that remain to be computed,  $\boldsymbol{\lambda}_1^T \left( \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u} + \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i} \right)$ , can be divided as follow:

$$\boldsymbol{\lambda}_1^T \left( \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u} + \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i} \right) = \underbrace{\boldsymbol{\lambda}_1^T \frac{\partial \mathbf{C}}{\partial l_i} \mathbf{u}}_{t_1} + \underbrace{\boldsymbol{\lambda}_1^T \mathbf{C}_B \frac{\partial \mathbf{u}_B}{\partial l_i}}_{t_2} \quad (3.3.47)$$

$t_2$  is composed by already known terms:  $\mathbf{C}_B$  is found from the initial global RBF-FD application while the other two terms that comprise it have already been calculated during the previous steps; we do not dwell too much on it. A term that deserve more attention is  $t_1$  since it requires to differentiate matrix  $\mathbf{C}$  which

### CHAPTER 3. AUTOMATIC DIFFERENTIATION (AD) AND ADJOINT METHOD IN OPTIMIZATION

---

is built row by row analyzing the stencil associated to each RBF-FD node. The contribution  $t_{1,k}$  given by the  $k$ -th RBF-FD node to  $t_1$  is found as:

$$t_{1,k} = \lambda_{1,k} \frac{\partial \mathbf{c}^T(\mathbf{x}_K)}{\partial l_i} \mathbf{u}_k \quad (3.3.48)$$

This makes clearer the relationship with the adjoint method applied in 1D case: the last two terms in equation (3.3.48) made up the element  $q_{k,i}$  of the matrix  $\mathbf{Q}$  defined in equation (3.3.16) of subsection 3.3.2. Thus repeating the passages reported in equations (3.3.18 - 3.3.23) we found that  $t_{1,k}$  can be rewritten as:

$$t_{1,k} = \lambda_{1,k} \left( \frac{\partial \mathbf{h}}{\partial l_i} - \frac{\partial \mathbf{M}_{BC,k}^T}{\partial l_i} \mathbf{c}(\mathbf{x}_k) \right)^T \mathbf{M}_{BC,k}^{-1} \mathbf{u}_k \quad (3.3.49)$$

and subsequently, by means of the adjoint vector  $\boldsymbol{\lambda}_{2,k}$  as:

$$t_{1,k} = \lambda_{1,k} \left( \frac{\partial \mathbf{h}}{\partial l_i} - \frac{\partial \mathbf{M}_{BC,k}^T}{\partial l_i} \mathbf{c}(\mathbf{x}_k) \right)^T \boldsymbol{\lambda}_{2,k} \quad (3.3.50)$$

where  $\boldsymbol{\lambda}_{2,k}$  is found by solving the system in (3.3.46) based on  $k$ -th, instead of the  $j$ -th, stencil.

What we have just shown is the application of the adjoint method to the optimization of a 3D problem based on an RBF-FD model, but the reasoning, with small changes, can be also applied similarly to problems which are based on RBF-HFD models as has been done in the 1D case.

# Chapter 4

## Results

In this chapter we will show the results obtained from the solution of two design optimization problems related to:

- the length of a simple 1D segment;
- the shape of a face of a prism.

These has been formulated as optimization problems and the adjoint method was used to compute the gradient of the corresponding cost functions throughout the optimization.

However, since both optimization problems are built on top of physical systems described by means of Poisson equations, we will first briefly describe this equation in a dedicated section before delving deeper into them.

### 4.1 Poisson Equation

Given a region  $\Omega \subset \mathbb{R}^d$ , with  $d \in \mathbb{N}$  and bounded by the boundary  $\partial\Omega$ , the generic heat equation with internal heat generation [32], at steady state, is described by the Poisson Equation which is defined by the following Partial Differential Equation (PDE):

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}) & \text{in } \Omega \\ \mathcal{B}u(\mathbf{x}) = g(\mathbf{x}) & \text{on } \partial\Omega \end{cases} \quad (4.1.1)$$

where  $\Delta$  indicates the Laplacian operator,  $\mathcal{B}$  is the (possibly differential) generic linear operator that enforce the boundary conditions (BCs) and  $f$  and  $g$  are known functions. The Poisson equation is a well-known example of a PDE that can be effectively addressed using meshless techniques, such as the RBF-FD approach detailed in this work.

We decided to evaluate the results of the adjoint method applied to RBF-FD method on the aforementioned equation due to both its wide range of applications in different areas (e.g. electrostatic, chemistry, gravitation and others) and the simplicity of its analytical manipulation.

## 4.2 1D case

In this case we have  $d = 1$ . The physical domain is simply a segment and its boundary consists on its two endpoints that we indicate with  $a$  and  $b$ , both belonging to  $\mathbb{R}$ . Therefore we have  $\Omega = ]a, b[$  and  $\partial\Omega = \{a, b\}$ . For the remaining portion of the boundary value problem we define  $f(x) = -\omega^2 \sin(x)$  and we impose the following homogeneous BCs:  $u(a) = 0$  and  $u(b) = 0$ . Fixing the values  $\omega = 1$ ,  $a = 0$  and  $b = \pi$  allows to rewrite the Poisson equation reported in (4.1.1) as

$$\begin{cases} -\Delta u(\mathbf{x}) = -\sin(\mathbf{x}) & \text{in } ]a, b[ \\ u(a) = 0 \\ u(b) = 0 \end{cases} \quad (4.2.1)$$

In order to find an approximated solution to the problem (4.2.1) we employ the RBF-FD method parametrized as follow:

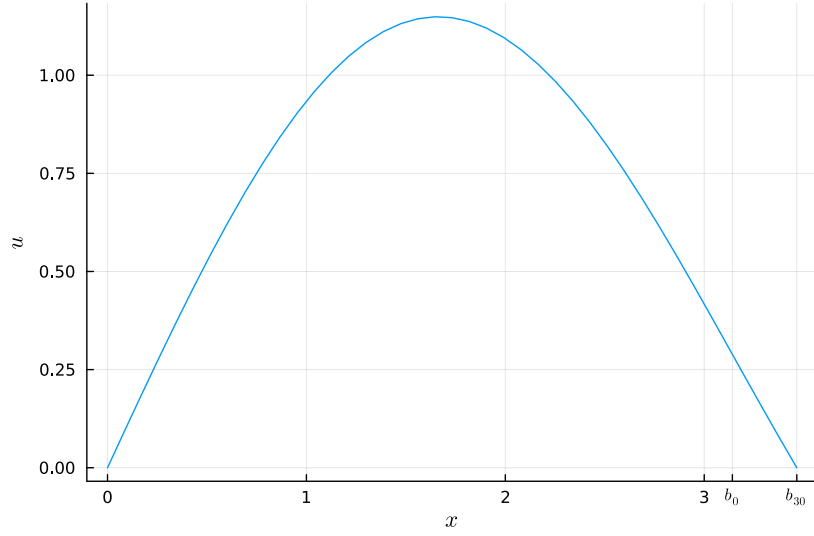
- $N = 41$  evenly spaced nodes are used for domain discretization:  $N_I = 39$  are placed in  $\Omega$  and the other  $N_B = 2$  are placed respectively in  $a$  and  $b$ ;
- $\varphi(r) = r^3$  Polyharmonic function is used as basic function to define the RBF interpolant;
- $P = 2$  is the degree of its polynomial augmentation;
- $m = 7$  is the number of nodes which constitute a single stencil;

As done previously we indicate with  $\mathbf{u}_I$  the vector obtained by solving the system (2.6.20) derived from RBF-FD discretization. The discretized version of equation (4.2.1) will be the constraints during the optimization

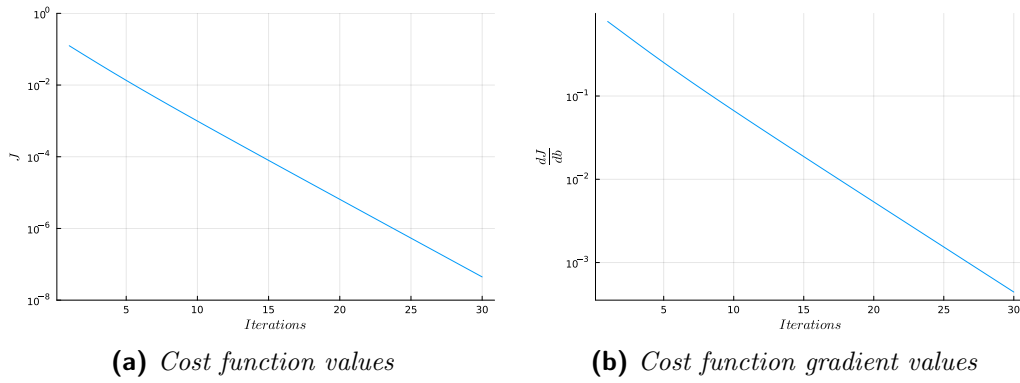
The design optimization problem, related to segment  $[a, b]$ , that we want to solve is the minimization of the following cost function:

$$J(\mathbf{u}_I, b) = \frac{1}{2} \frac{b}{N-1} \left( \sum_{i=2}^{40} u(x_i) - E \right)^2 \quad (4.2.2)$$

where  $E \in \mathbb{R}$  is a parameter set by us to a value of 2.5 and  $l = b$  is the single design variable. This particular configuration of the design variable means that in order to meet the objective we are allowed to vary the length of the segment  $[a, b]$ .



**Figure 4.1:** discretized approximated solution obtained solving (4.2.1) via RBF-FD at the end of the design optimization.  $b_0$  and  $b_{30}$  indicate respectively the initial and the final position of the  $b$  endpoint



**Figure 4.2:** Trajectories of the cost function and its gradient throughout the optimization process. The ordinates are on a logarithmic scale

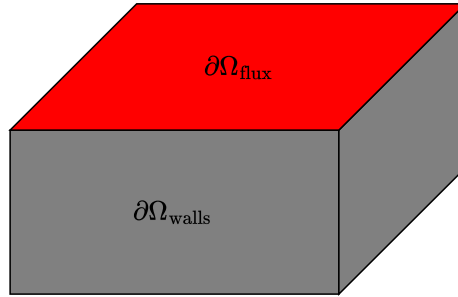
The results after 30 steps in the optimization procedure are shown in figure 4.1, while in figure 4.2 are shown the dynamic of the cost function and of its gradient respect the design variable  $b$ .

### 4.3 3D case

Now the physical domain  $\Omega$  consists of a square-based prism with a side length of 1 and a height of 0.5, its faces are indicated with  $\partial\Omega$ . In particular  $\partial\Omega$  is further partitioned in two subsets:

- $\partial\Omega_{\text{flux}}$  which is formed by the north face of the prism;
- $\partial\Omega_{\text{walls}}$  which includes the remaining faces.

Figure 4.3 show a representation of the domain.



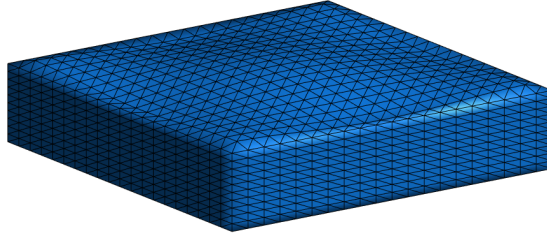
**Figure 4.3:** Physical domain used for the 3D Poisson equation. Boundaries  $\partial\Omega_{\text{flux}}$  and  $\partial\Omega_{\text{walls}}$  are highlighted in red and gray respectively

The Poisson equation associated to the geometry is defined as:

$$\begin{cases} \Delta u = -\omega^2 f & \text{in } \Omega \\ \frac{\partial u}{\partial \mathbf{n}} = \frac{\partial u}{\partial \mathbf{n}} & \text{in } \partial\Omega_{\text{walls}} \\ u = f & \text{in } \partial\Omega_{\text{flux}} \end{cases} \quad (4.3.1)$$

where  $\mathbf{n}$  is the surface normal at the boundary and  $f$  is defined as  $f(\mathbf{x}) = \sin(\omega x) \sin(\omega y) \sin(\omega z)$ , with  $\mathbf{x} \in \Omega \cup \partial\Omega$  being a point in space with coordinates  $x$ ,  $y$  and  $z$ , and  $\omega = 2\pi$ .

The problem under consideration allows for various physical interpretations: if we consider  $u$  as the temperature and  $-\omega^2 f$  as a generic heat source, the aforementioned problem constrains the temperature of the north face to be constant



**Figure 4.4:** Representation of the `.stl` file containing the prism used in the 3D design optimization problem

and to follow a profile defined by  $f$  (Dirichlet BC), and impose a prescribed heat flux on the other faces (Neumann BC); by solving it we obtain the temperature values within the domain.

On top of the boundary value problem (4.3.1) a design optimization problem is subsequently formulated: we want to maximize the flux through the north wall per unit area by modifying the shape of the wall itself. This is the problem that we are going to solve in this section. The associated cost function can be defined in a similar way as was done in subsection 3.3.3 on page 45:

$$J_c = \frac{1}{|\partial\Omega_{\text{flux}}|} \int_{\partial\Omega_{\text{flux}}} \frac{\partial u}{\partial \mathbf{n}} ds \quad (4.3.2)$$

where  $|\partial\Omega_{\text{flux}}| \in \mathbb{R}$  is the area of the surface defined by  $\partial\Omega_{\text{flux}}$ .

In this case too a physical interpretation of the problem is possible: if we suppose we are dealing with a heat conduction problem, like in the example given earlier in this section, then problem (4.3.5) is asking to find the best shape for the north face in order to maximize the average heat flux density with the exterior of the prism.

However, in our case the geometry is not described by means of continuous quantities and it is separately defined in a `.stl` file which characterizes the surface of the prism using a triangular mesh as it is shown in figure 4.4.

This means that in order to modify the shape of the north face of the cube we can act on the height of its triangles control points (i.e. their vertices). Thus it is useful to define the vector of the design variables as  $\mathbf{l} = [l_1 \dots l_L]$  where  $L \in \mathbb{N}$  is the number of triangle vertices which belongs to  $\partial\Omega_{\text{flux}}$  and the  $i$ -th variable  $l_i \in \mathbb{R}$  control the height  $z_i$  of the  $i$ -th vertex belonging to  $\partial\Omega_{\text{flux}}$ . Defining the design variables in this way allows to define a finite-dimensional optimization problem which requires the definition of a cost function based on the variables  $\mathbf{l}$  in order to be solved. We defined this cost function, which aims to approximate the one



presented in equation (4.3.2), as:

$$J = \frac{1}{|A|} \sum_{j \in \mathcal{T}} \nabla \mathbf{u}(\mathbf{x}_j)^T \mathbf{a}_j \quad (4.3.3)$$

where  $|A|$  is the area of  $\partial\Omega_{\text{flux}}$  computed as the sum of the areas of the triangles that belong to it,  $\mathcal{T}$  is the set of the indices used to identify those triangles,  $\nabla \mathbf{u}(\mathbf{x}_j)^T$  is the gradient vector of the field  $u$  evaluated at  $\mathbf{x}_j$ , centroid of the  $j$ -th triangle, and  $\mathbf{a}_j$  is the area vector of the  $j$ -th triangle which has direction given by its normal and modulus given by its surface. Additionally the boundary value problem (4.3.1), in general, must be discretized to be solved, doing so via RBF-HFD method yields the following constraint for the optimization:

$$\mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B = \mathbf{f} \quad (4.3.4)$$

where  $\mathbf{u}_I$  and  $\mathbf{u}_B$  are the vectors composed of the values of the field  $u$  at the internal and boundary nodes, respectively,  $\mathbf{C}_I$  and  $\mathbf{C}_B$  are the matrices representing the contribution of internal and boundary RBF-HFD nodes, respectively, and  $\mathbf{f}$  is the vector composed of the values of  $f$  at the internal nodes. Therefore our overall design optimization problem reads as follow:

$$\begin{aligned} & \text{maximize} && J(\mathbf{u}, \mathbf{l}) \\ & \text{by varying} && l_i \quad i = 1, \dots, L \\ & \text{subject to} && \mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B - \mathbf{f} = \mathbf{0} \end{aligned} \quad (4.3.5)$$

To solve problem (4.3.5) we have used a fixed number of 50 optimization steps and  $\mathbf{l}$  is updated using the following rule:

$$\mathbf{l}_{k+1} = \mathbf{l}_k + \alpha \frac{\partial J}{\partial \mathbf{l}} \quad (4.3.6)$$

where  $\alpha$  is fixed heuristically, by trial and error, to 0.1 and the vector  $\partial J / \partial \mathbf{l}$  is computed as explained in subsection 3.3.3. Constraint (4.3.4), on the other hand, is obtained employing the RBF-HFD method parametrized as follow:

- the  $N$  RBF-FD nodes scattered over the domain are generated according to the technique presented in [18] which yields isotropic node distributions where equal node spacing is obtained along all directions. Unlike the 1D case, thanks to this technique, the number of nodes scattered across the cube varies as the physical domain changes throughout the optimization;
- for the RBF basis we utilized the multiquadric function  $\varphi(r) = \sqrt{1 + (\epsilon r)^2}$

- the degree of the polynomial augmentation is set to  $P = 3$
- the number of nodes which constitute a single stencil is forced to  $m = 50$ .

The method also requires the solution of the sparse linear system reported in equation (2.6.18) for which we used the `bicgstabl()` function present in the julia package called `IterativeSolvers.jl`<sup>1</sup> setting `GMRES=2` and the left preconditioner `Pl` set to `ilu_LHS`; this function is an implementation of the iterative solver found in [33] to which we refer for further details along with the `IterativeSolvers.jl` package documentation. As for the left preconditioner `ilu_LHS`, it is the incomplete LU factorization of matrix  $\mathbf{C}_I$  obtained via the `ilu()` function of the `IncompleteLU.jl`<sup>2</sup> julia package with parameter  $\tau$  set to 6. Its implementation is based on [34]. When `ilu_LHS` is passed as parameter to `bicgstabl()`, it is used to pre-multiply on the left both sides of equation (2.6.18) so as to accelerate the convergence to the solution of the iterative solver.

### 4.3.1 Results

The results of the optimization, where the gradient at each step is computed employing the discreet adjoint method, are reported in figure 4.5 where the values of the cost functions at each step are shown. In order to be sure of the correctness of the results we also compared them to the results obtained from the solution of the continuous optimization problem, which require the maximization of (4.3.2) with constraints defined by the Poisson equation (4.3.1).

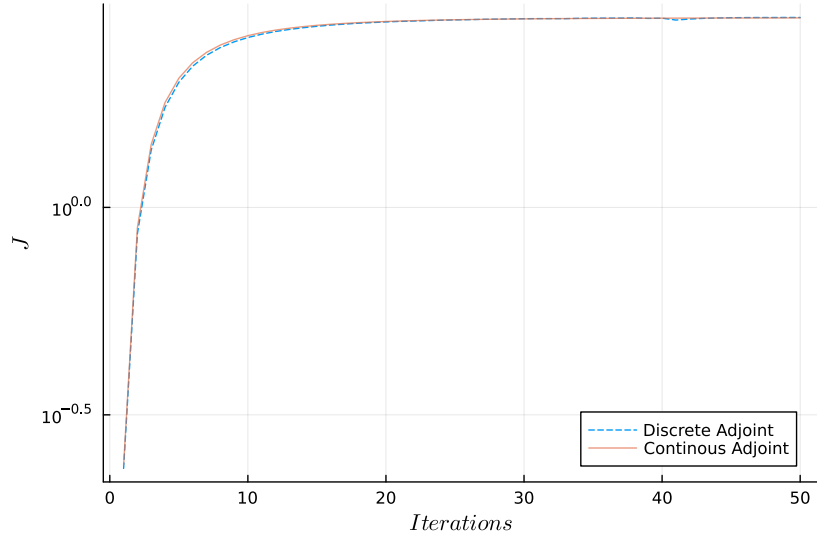
Its solution is found by transforming the boundary integral in equation (3.3.26) in a volume integral using the divergence theorem. In the thermal analogy this is equivalent to state that the sum of total heat fluxes at the boundaries of the domain equals the total internal heat generation, i.e.  $-\omega^2 f$ . This step from heat fluxes on the boundary to internally generated heat is what allows to express the objective (3.3.26) independently on the state variables  $\mathbf{u}$ , making it depend only on the design variable  $\mathbf{l}$ .

This simplified objective function can thus be used to provide a simple and reliable reference solution to the discrete optimization problem reported in (4.3.5). We also remark that this simplified approach, which is called continuous adjoint, does not decrease the importance in the study of the RBF-FD model and discrete

---

<sup>1</sup>Julia open source library which provides iterative algorithms for solving linear systems, eigen-systems, and singular value problems. The project can be found at <https://iterativesolvers.julia.linearalgebra.org/stable/>.

<sup>2</sup>Julia open source library which implements the Crout version of the incomplete LU factorization of a sparse matrix. The project can be found at <https://github.com/haampie/IncompleteLU.jl>.

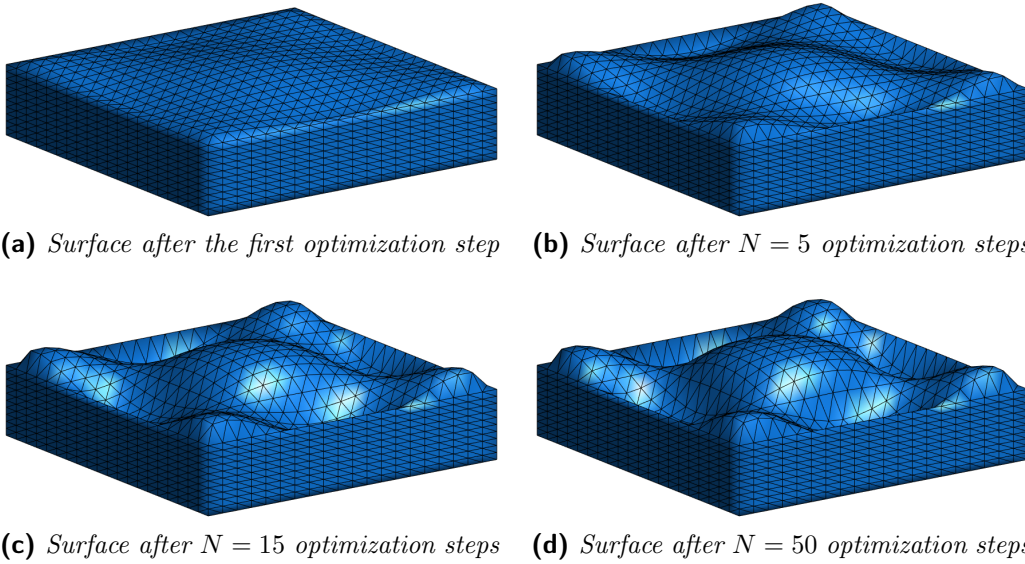


**Figure 4.5:** Trajectories of the cost functions associated to the optimization problems obtained considering a prism defined in a continuous domain (orange, solved via continuous adjoint) and the same prism defined by means of a `.stl` file (blue, solved via discrete adjoint)

adjoint (the adjoint method covered in section 3.3 on page 39) since it can be applied to a very narrow class of problems such as the one under consideration.

What can be observed from figure 4.5 is a very precise agreement between the results obtained with the discrete and continuous cases which proves the validity of the discrete adjoint method for solving design optimization problem with simple gradient-based optimization algorithms similar to the one reported in equation (3.2.7).

The deformations of surface  $\partial\Omega_{\text{flux}}$  at 4 different times of the optimization are instead shown in figure 4.6. In this case it can be noted that the convergence to the optimal design for the prism was very rapid since the northern surface had already converged to a shape similar to the final one after only 15 iterations.



**Figure 4.6:** Deformation of the north face of the prism at different optimization times

# Conclusion

The present study aims to investigate design optimization problems based on physical models whose behavior can be described by means of Poisson equation.

For the solution of the Poisson equations we decided to leverage on meshless methods, in particular RBF-FD and RBF-HFD methods, which are quickest approaches respect the standard Finite Element Methods (FEMs) or Finite Volume Methods (FVMs) and, even more importantly, does not ask the creation of a mesh which require a well trained operator in order to be created.

We addressed the possibility to introduce:

- Automatic Differentiation (AD), a family of techniques for evaluating derivatives of numeric functions expressed as computer programs which enable efficient and accurate computations;
- adjoint method, a well established technique used in Computational Fluid Dynamics (CFD) design optimization based on FEMs and FVMs.

in simple 1D and 3D problems defined by us. After an introduction to the RBF-FD and RBF-HFD methods we analyzed the structure of design optimization problems. In particular we have shown how AD and adjoint method can be used for gradient-based optimization algorithms which require the computation of the gradient of the cost function.

The mathematical steps shown are then subsequently implemented using the `Julia` programming language where we have used the `Zygote` library to implement the part related to the automatic differentiation. For the adjoint method, instead, the code has been written completely by us without using any pre-existent package.

The implemented algorithms are then initially tested on a simple 1D problem to test their correctness leading to the correct minimization of the defined objective function. After this case we also faced a 3D optimization problem which, although simple, allowed for a physical interpretation related to thermal CFD. In this case we shown the correctness of our approach which has been compared with the continuous adjoint providing the same results.

Beyond the good results achieved with the optimization problems faced on this work, further studies can be conducted regarding the topics discussed in this thesis:

## CONCLUSION

---

- test AD and adjoint method based on more complex models different from Poisson equations with more complex 3D domains;
- generalizing the proposed approach to more general objective functions in the 3D case;
- modify the implemented code so that it can be executed using multiple threads at once, thus reducing the computational time;
- explore the possibility of using automatic differentiation, through the functions provided by `Zygote`, directly for the calculation of the objective function gradient.

# Appendix A

## Stereolithography **.stl** files

**.stl** is a file format which dates back to 1987 when it has been created by *3D Systems*<sup>1</sup> for its stereolithography (from which the file extension originates) printing technology for commercial 3D printers. Eventually it became the 3D printing and rapid prototyping industry's defacto standard. In general **.stl** files are also referred to as standard triangle language or standard tessellation language files.

The **.stl** format is used to describe an unstructured triangulated surface using a list of triangles. Each of these triangles is described separately by its own normal and vertices using a Cartesian coordinate system [35]. Their main purpose is to describe the surface geometry of a 3D object, thus in this files, differently from CAD files, no information about colors or textures is present. They also do not contain any scale information and can be specified both in ASCII and binary representation.

An example of of a **.stl** file in ASCII format is:

```
solid topFace
  facet normal -0.000000e+00  0.000000e+00  1.000000e+00
    outer loop
      vertex  0.000000e+00 -0.000000e+00  5.000000e-01
      vertex  5.000001e-02 -5.000001e-02  5.000000e-01
      vertex  5.000001e-02 -0.000000e+00  5.000000e-01
    endloop
  endfacet
...
endsolid topFace
```

---

<sup>1</sup>American company that engineers, manufactures and sells 3D printers, 3D printing materials, 3D printed parts and application engineering services. For further information we refer to <https://www.3dsystems.com/>

```

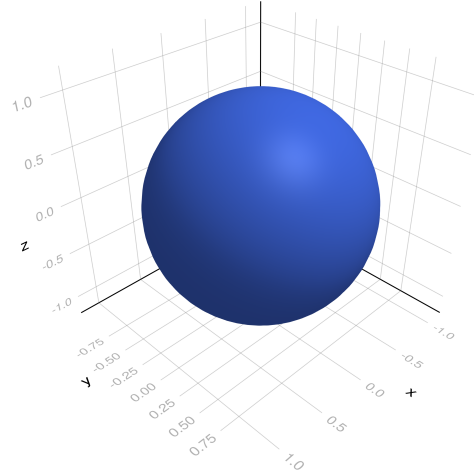
solid fixedFaces
  facet normal 0.000000e+00 1.000000e+00 0.000000e+00
    outer loop
      vertex 0.000000e+00 -0.000000e+00 0.000000e+00
      vertex 5.000001e-02 -0.000000e+00 5.000000e-02
      vertex 5.000001e-02 -0.000000e+00 0.000000e+00
    endloop
  endfacet
...
endsolid fixedFaces

```

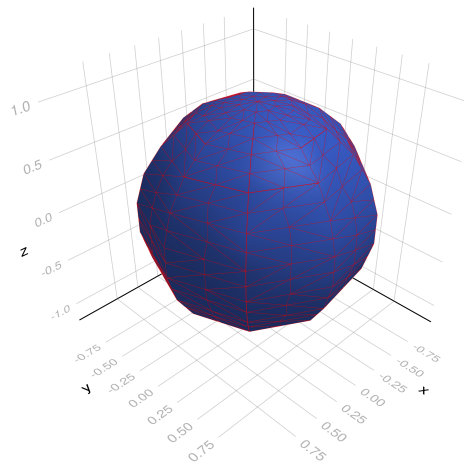
This is the `.stl` file that we have used to represent the geometry used in the 3D design optimization problem in section 4.3 on page 55. In the same file showed above are stored two different surfaces named respectively `topFace` and `fixedFaces`; each triangle of a surface is described by means of a `facet` (we reported only one triangle per surface for brevity) which specify its normal and the coordinates of its vertices arranged counterclockwise. Those files can also be stored in binary format, which is particularly convenient for large files.

We finally remark that, in general, `.stl` files contain only an *approximation* of geometry surfaces. An example of this situation can be found in figure A.1 on the next page, specifically in figure A.1b it can be seen how the red lines of the triangle edges are not capable to perfectly reconstruct the smoothness of the original surface (figure A.1a). Augmenting the number of triangles of the mesh improve the approximation, at cost of larger `.stl` files, but will not be possible to perfectly reconstruct curved surfaces.





(a) *Original sphere*



(b) *Sphere .stl representation*

**Figure A.1:** Example of approximating the surface of a unit-radius sphere centered at the origin using an .stl file

# Bibliography

- [1] S. R. Idelsohn and E. Oñate. “To mesh or not to mesh. That is the question...” In: *Computer Methods in Applied Mechanics and Engineering* (2005).
- [2] O. Vlasniuk et al. “Fast high-dimensional node generation with variable density”. In: *Computers & Mathematics with Applications* (2018).
- [3] R. Zamolo and E. Nobile. “Two algorithms for fast 2D node generation: Application to RBF meshless discretization of diffusion problems and image halftoning”. In: *Computers & Mathematics with Applications* (2018).
- [4] D. Rumelhart, G. Hinton, and R. Williams. “Learning representations by back-propagating errors”. In: *Nature* (1986).
- [5] M. B. Giles and N. A. Pierce. “An Introduction to the Adjoint Approach to Design”. In: *Flow, Turbulence and Combustion* (2000).
- [6] J. Bezanson et al. *Julia: A Fast Dynamic Language for Technical Computing*. 2012. arXiv: [1209.5145 \[cs.PL\]](https://arxiv.org/abs/1209.5145). URL: <https://arxiv.org/abs/1209.5145>.
- [7] Michael Hillman Jiun-Shyan Chen and Sheng-Wei Chi. “Meshfree Methods: Progress Made after 20 Years”. In: *Journal of Engineering Mechanics* (2017).
- [8] T. Belytschko et al. “Meshless methods: An overview and recent developments”. In: *Computer Methods in Applied Mechanics and Engineering* (1996).
- [9] W. Benz and E. Asphaug. “Simulations of brittle solids using smooth particle hydrodynamics”. In: *Computer Physics Communications* (1995).
- [10] S. Jun W. K. Liu and Y. F. Zhang. “Reproducing kernel particle methods”. In: *International Journal for Numerical Methods in Fluids* (1996).
- [11] P. Lancaster and K. Salkauskas. “Surfaces Generated by Moving Least Squares Methods”. In: *Mathematics of Computation* (1981).
- [12] E. J. Kansa. “Multiquadrics—A scattered data approximation scheme with applications to computational fluid-dynamics—I surface approximations and partial derivative estimates”. In: *Computers & Mathematics with Applications* (1990).

## BIBLIOGRAPHY

---

- [13] E. J. Kansa. “Multiquadrics—A scattered data approximation scheme with applications to computational fluid-dynamics—II solutions to parabolic, hyperbolic and elliptic partial differential equations”. In: *Computers & Mathematics with Applications* (1990).
- [14] M. A. Schweitzer. *Partition of Unity Method*. Lecture Notes in Computational Science and Engineering. 2003. URL: [https://doi.org/10.1007/978-3-642-59325-3\\_2](https://doi.org/10.1007/978-3-642-59325-3_2).
- [15] I. Amidror. “Scattered data interpolation methods for electronic imaging systems: a survey”. In: *Journal of Electronic Imaging* (2002).
- [16] J. C. Mairhuber. “On Haar’s Theorem Concerning Chebychev Approximation Problems Having Unique Solutions”. In: *Proceedings of the American Mathematical Society* (1956).
- [17] R. Schaback and H. Wendland. “Kernel Techniques: From Machine Learning to Meshless Methods”. In: *Acta Numerica* (2006).
- [18] R. Zamolo. “Radial Basis Function-Finite Difference Meshless Methods for CFD Problems”. PhD thesis. Università degli studi di Trieste, 2018.
- [19] G. Fasshauer. *Meshfree Approximation Methods with Matlab*. World Scientific Publishing Company, 2007.
- [20] D. Miotti. “Stable meshless methods for 3D CFD simulation on complex geometries based on Radial Basis Functions”. PhD thesis. Università degli studi di Trieste, 2024.
- [21] A. Tolstykh. *On using RBF-based differencing formulas for unstructured and mixed structured-unstructured grid calculations*. 2000. URL: [https://scholar.google.com/scholar\\_lookup?title=On+using+RBF+based+differencing+formulas+for+unstructured+and+mixed+structured%E2%80%9393unstructured+grid+calculations&author=Tolstykh+A.+I.&publication+year=2000&journal=Proc.+16th+IMACS+World+Congress&volume=228&pages=4606-4624](https://scholar.google.com/scholar_lookup?title=On+using+RBF+based+differencing+formulas+for+unstructured+and+mixed+structured%E2%80%9393unstructured+grid+calculations&author=Tolstykh+A.+I.&publication+year=2000&journal=Proc.+16th+IMACS+World+Congress&volume=228&pages=4606-4624).
- [22] L. A. Bueno, E. A. Divo, and A. J. Kassab. “A coupled localized RBF meshless/DRBEM formulation for accurate modeling of incompressible fluid flows”. In: *International Journal of Computational Methods and Experimental Measurements* (2017).
- [23] L. A. Bueno, E. A. Divo, and A. J. Kassab. “Multi-scale cardiovascular flow analysis by an integrated meshless-lumped parameter model”. In: *International Journal of Computational Methods and Experimental Measurements* (2018).
- [24] G. Kosec and J. Slak. “Rbf-fd based dynamic thermal rating of overhead power lines”. In: *Advances in Fluid Mechanics XII* (2018).

## BIBLIOGRAPHY

---

- [25] G. Kosec and J. Slak. “Radial basis function-generated finite differences solution of natural convection problem in 3D”. In: *AIP Conference Proceedings* 2293 (2020).
- [26] R. J. LeVeque. *Finite Difference Methods for Ordinary and Partial Differential Equations: Steady-State and Time-dependent Problems*. Society for Industrial and Applied Mathematics, 2007.
- [27] J. L. Bentley. “Multidimensional Binary Search Trees Used for Associative Searching”. In: *Communications of the ACM* (1975).
- [28] A. Kolar-Požun et al. “Oscillatory behaviour of the RBF-FD approximation accuracy under increasing stencil size”. In: *Computational Science – ICCS 2023* (2023).
- [29] A. G. Baydin, B. A. Pearlmutter, and A. A. Radul. “Automatic Differentiation in Machine Learning: a Survey”. In: *Journal of Machine Learning Research* (2018).
- [30] Joaquim R. R. A. Martins and Andrew Ning. *Engineering Design Optimization*. Cambridge University Press, 2020.
- [31] R. Bombardieri et al. “Aerostructural Wing Shape Optimization assisted by Algorithmic Differentiation”. In: *Structural and Multidisciplinary Optimization* (2020).
- [32] H. Brezis. *Functional Analysis, Sobolev Spaces and Partial Differential Equations*. 2011. URL: <https://doi.org/10.1007/978-0-387-70914-7>.
- [33] G. L. G. Sleijpen and D. R. Fokkema. “BiCGstab(l) for linear equations involving unsymmetric matrices with complex spectrum”. In: *Electronic Transactions on Numerical Analysis* (2000).
- [34] N. Li, Y. Saad, and E. Chow. “Crout versions of ILU factorization with pivoting for sparse symmetric matrices”. In: *Electronic Transactions on Numerical Analysis* (2005).
- [35] Wikipedia contributors. *STL (file format)*. Online; accessed 30-October-2024. 2024. URL: [https://en.wikipedia.org/wiki/STL\\_\(file\\_format\)](https://en.wikipedia.org/wiki/STL_(file_format)).
- [36] Michel Goossens, Frank Mittelbach, and Alexander Samarin. *The L<sup>A</sup>T<sub>E</sub>X Companion*. Reading, Massachusetts: Addison-Wesley, 1993.
- [37] B. Fornberg and N. Flyer. *A Primer on Radial Basis Functions with Applications to the Geosciences*. 2015.
- [38] B. Fornberg and N. Flyer. “Solving PDEs with radial basis functions”. In: *Acta Numerica* (2015).

## BIBLIOGRAPHY

---

- [39] G. R. Liu and Y. T. GU. *An Introduction to Meshfree Methods and Their Programming*. Cambridge University Press, 2005.
- [40] D. Miotti, R. Zamolo, and E. Nobile. “A Fully Meshless Approach to the Numerical Simulation of Heat Conduction Problems over Arbitrary 3D Geometries”. In: *Energies* (2021).
- [41] R. Zamolo, D. Miotti, and E. Nobile. “Accurate stabilization techniques for RBF-FD meshless discretizations with neumann boundary conditions”. In: *SSRN Electronic Journal* (2022).
- [42] The MathWorks Inc. *What Is Design Optimization?* 2024. URL: <https://it.mathworks.com/discovery/design-optimization.html>.
- [43] L. Bacer. “Fully meshless approaches for the numerical solution of partial differential equations: radial basis function finite difference method and physics-informed neural network”. MA thesis. Università degli studi di Trieste, 2023.
- [44] K. C. Jeyanthi. “Efficient node generation over complex 3D geometries for meshless RBF-FD method”. MA thesis. Università degli studi di Trieste, 2022.
- [45] Wikipedia contributors. *Finite difference method*. Online; accessed 21-February-2024. 2024. URL: [https://en.wikipedia.org/wiki/Finite\\_difference\\_method](https://en.wikipedia.org/wiki/Finite_difference_method).