

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

July 31, 2024

Candidate
Kevin Marzio

Supervisor
Prof. Andrea De Lorenzo

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

A.Y. 2023/2024

*Lorem ipsum
dolor sit amet*

- Cicero

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Contents

Abstract	I
Introduction	IV
1 Meshless methods	1
1.1 Motivation of case study	3
2 RBF-FD method	5
2.1 Scattered Data Interpolation	5
2.2 The Mairhuber-Curtis theorem	7
2.3 Radial Basis Functions	8
2.4 Polynomial augmentation	9
2.5 Problem solution	11
2.6 RBF-FD	12
2.6.1 Finite difference method	13
2.6.2 Formulation	14
2.7 RBF-HFD	19
2.7.1 Hermite Interpolation	19
2.7.2 Formulation	22
3 Adjoint method	25
3.1 Automatic Differentiation	25
3.1.1 Forward mode	27
3.1.2 Reverse mode	28
3.2 Application to design optimization	30
3.3 Application to RBF-FD	35
3.3.1 1D	36
3.3.2 3D	39
4 Results	41

CONTENTS

Conclusion	42
4.1 Lorem Ipsum	42
4.1.1 Dolor sit amet	43
A Sit Amet	44

Introduction

Spiegazione/storia mesh-based.

Perché siamo arrivati ai meshless e differenze in generale

Perché diff. automatica e aggiunto nel meshless

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” [1].

They appeared for the first time in 1977 with the Smooth Particle Hydrodynamics (SPH) method [2], 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 [3]. 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) [4] 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

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 $\frac{\partial u}{\partial \mathbf{n}}$ indicate the normal derivative of u 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 α_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 [4], moving least square (MLS) [5], radial basis function (RBF) [6, 7] and partition of unity (PU) [8]. Furthermore solving the PDE in (1.0.1) with the approximated solution u^h , in general, yields to a non-zero error function ϵ^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 ϵ^h are used to integrate the error to zero:

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

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

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 equaitons (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.

1.1 Motivation of case study

In this section we motivate the specific case that will be studied in the following chapters: the generic heat equation with internal heat generation [9] at steady state given by

$$\begin{cases} -\Delta u(\mathbf{x}) = f(\mathbf{x}) & \text{in } \Omega \\ u(\mathbf{x}) = g(\mathbf{x}) & \text{on } \partial\Omega \end{cases} \quad (1.1.1)$$

The reason behind this choice are 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. Equation (1.1.1) is one of the many PDEs that can be solved using the meshless approaches.

In particular for this thesis we will:

- consider domains $\Omega \subset \mathbb{R}^d$ with $d = 1, 3$;
- use RBFs augmented with polynomial terms as basis functions B_k to approximate the solution u in a similar way to what has been done in [10];
- formulate the problem in its strong form, so we will use the collocation technique in combination with the Weighted Residual Method.

thus we will apply the RBF-generated Finite Differences (RBF-FD) [11, 12] to solve (1.1.1) in one and three dimensional physics domain.

Mono dimensional physics domain are used to gain confidence with the techniques explained while three dimensional domain are used for their application in real case problems as those faced at Esteco. 2D cases, instead, have been skipped in the analysis since they would only be used to bridge the 2 previous scenarios and therefore would have had no practical utility.

RBFs are used for scattered data interpolation both for their physical foundation[overview 20, Hardy 1990] and for their profitable use in applications like meteorology, turbulence analysis and neural network [1]

Finally, collocation technique is employed because thanks to its ability to discretize the Boundary Value problem (1.1.1) expressed in strong form leads to a real fully meshless approach [13].

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 on page 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, Ω 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 [14].

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

- a finite set of nodes, $\mathcal{X} \subset \mathbb{R}^d$, which could 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 α_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 Ω 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 Π_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 [15]; 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 α_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 [16], are called *kernels*. In general at different basic functions φ 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 ϵ , 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 [17].

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 Π_P^d and β_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:

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

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 φ 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. This means that the above system require a positive definite \mathbf{B} in order to be solved and this property depends on the choice of the basic function φ used to define the RBFs. By definition only *strictly* positive definite [18] φ are associated to a positive definite \mathbf{B} and, this, restrict our choices: of those shown in table 2.1 on page 9 only the Inverse Multiquadratic satisfy this requirement. The same attributes of φ are also inherited by the associated RBFs Φ_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 that has 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 [19]:

1. basic function φ has to be *strictly conditionally* positive definite function of order P [18];

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

Condition 1 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 as these are strictly conditionally definite of order 1. 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, but doing so results in ill-conditioning and singularity issues related to \mathbf{P} that, since non-singularity of \mathbf{P} is a necessary condition for the non singularity of \mathbf{M} , also affect \mathbf{M} .

This is why condition 2 is also required. 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 [18], 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 RBF-FD

In this section we are going to discuss in detail how RBF-FD meshless method is used to solve Partial Differential Equations (PDEs). Its first implementation was developed by Kansa in [6, 7]. 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 [20]. We will denote the latter as *local method*. In recent years RBF-FD local method has been developed and applied with success [21, 22, 23, 24].

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 along the edges of this domain. In most cases an analytic formula for the solution can not be found. What RBF-FD methods do is replacing derivatives in the differential equation by Finite Difference 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 [25], 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:

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

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 Formulation

Now that we have gained a certain level of familiarity with Finite Difference (FD) methods 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. Of course $N_I + N_B = N$ and $\mathcal{X}_I \cup \mathcal{X}_B = \mathcal{X}$.

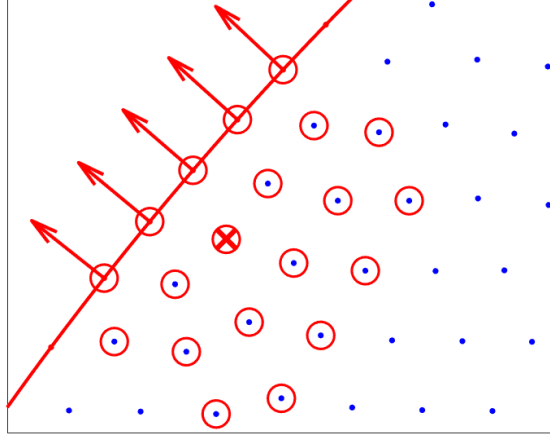


Figure 2.1: Example of 2D stencil. The nodes which belong to it are marked with a red circle, its central node with a red cross. Red arrows represent the boundary normals \mathbf{n} of those nodes belonging to $\partial\Omega$. Figure taken from [19]

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 is associated to each node $\mathbf{x}_i \in \mathcal{X}_I$ and it 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 Ω 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 on page 9. In this case the problem is well-posed, regardless of the shape of domain Ω if the instructions presented in 2.5 on page 11 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 α and β 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 equation (2.6.14) 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 $\mathbf{c}_p(\mathbf{x}_i)$. In general the approximation is not exact, however it is in this specific case thanks to the combination of:

- the linearity of u^h respect the parameters $\alpha_1, \dots, \alpha_m$ and β_1, \dots, β_M , as can be seen in equation (2.6.9), and;
- the linearity of operator \mathcal{L} .

We stress that it can not be properly called FD since not only values of u^h (remember that $u(\mathbf{x}_j) = u^h(\mathbf{x}_j)$ thanks to conditions (2.6.11a)), but also of g are used for the approximation.

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:

$$\begin{aligned} \mathbf{C}_I &= \begin{bmatrix} c_1(\mathbf{x}_1) & \dots & c_{N_I}(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ c_1(\mathbf{x}_{N_I}) & \dots & c_{N_I}(\mathbf{x}_{N_I}) \end{bmatrix} \in \mathbb{R}^{N_I \times N_I} \\ \mathbf{C}_B &= \begin{bmatrix} c_{N_I+1}(\mathbf{x}_1) & \dots & c_N(\mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ c_{N_I+1}(\mathbf{x}_{N_I}) & \dots & c_N(\mathbf{x}_{N_I}) \end{bmatrix} \in \mathbb{R}^{N_I \times N_B} \end{aligned} \quad (2.6.19)$$

What we just obtained in (2.6.18) is the algebraic system that once solved, in place of the PDE, using an iterative program, 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 equation (2.6.6) is used to approximate the PDE solution instead of the interpolant in (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 inverting a matrix is at least of $\mathcal{O}(n^3)$, where here n denotes the number of rows and columns of the matrix to be inverted, is clear how local method is computationally advantageous respect the global one. Furthermore 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 . The two reason just listed are the ones that explain the benefits of stencil introduction and that led to the abandonment of Kansa's method in favor of the one proposed by Tolstykh.

What we have not covered yet is the computational cost of the stencils creation: if too costly, it might negate 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 [26] 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. Regarding the effects of stencil size

Finally, in addition 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. [27] where they found that changing stencil's size induce oscillations in both the solution and discretization errors for the Poisson equation.

2.7 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 BCs are applied, to find normal directions associated with them that make the local matrix \mathbf{M}_{BC} in equation (2.6.12) singular [19].

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

During 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 condition, more generally, can be seen as a set of given values $\{u(\mathbf{x}_1), \dots, u(\mathbf{x}_m)\}$ along with a given set of functionals $\Lambda = \{\lambda_1, \dots, \lambda_m\}$ such that:

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

Conditions (2.7.3) are typically called Hermite interpolation conditions [18]. 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 on page 5, 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 δ_{1,\mathbf{x}_i} , the same applied on the first m_I rows.

An improvement, alternative to polynomial augmentation, 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 $\Phi(\cdot, \cdot)$; this leads to the so-called *Generalized Hermite Interpolant*. Suppose that the operator δ_{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 Φ 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 *unsymmetryc*. Furthermore \mathbf{B}_H is also positive definite when

the basic function φ used to define the kernel is strictly positive definite and operators λ_j are linearly independents [19]. Symmetry and positive definitiveness properties that were just referred to are met also in case the Hermite interpolation conditions are those of the local RBF-FD method reported in equation (2.7.1).

In addition, adoption of polynomial augmentation, as explained in section 2.4 on page 9, 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 α and β , 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:

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

where $\Lambda \mathbf{u} = [\lambda_1 u, \dots, \lambda_m u]$ 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 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 Ω .

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), that previously had to be resolved in order to find the coefficients, 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 :

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

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:

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

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:

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

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

Adjoint method

3.1 Automatic Differentiation

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 13, are used to approximate the derivatives *values*. Easy to implement it has the disadvantage to be inaccurate due to round-off and truncation errors [25];

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 or SymPy. 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, can be applied to code involving branches, loops and recursion as opposed to symbolic

differentiation [28]. 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 showcase 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 $\tilde{f}_1(x_1, x_2) = x_1 x_2 + \cos x_1$ 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:

- $v_{i-n} = x_i, i = 1, \dots, n$ 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 creating 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 AD.

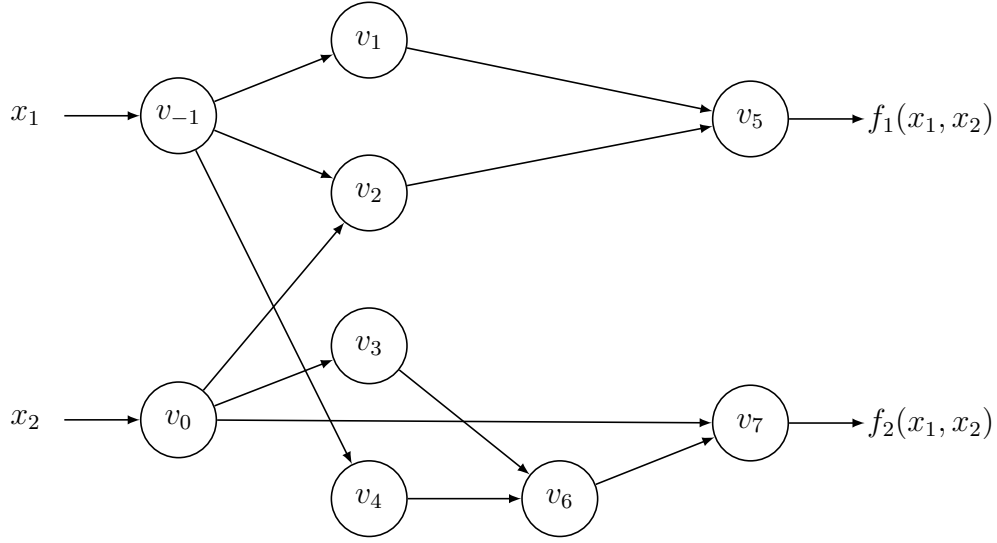


Figure 3.1: Computational graph of function $\tilde{f}(x_1, x_2) = [\tilde{f}_1(x_1, x_2) \ \tilde{f}_2(x_1, x_2)]$. Definitions of intermediate variables v_{-1}, \dots, v_7 can be found in table 3.1 or table 3.2

3.1.1 Forward mode

In order to compute the derivative of the function \tilde{f} , reported in (3.1.1), 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:

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

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. This 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 is possible computing efficiently and in a matrix-free way Jacobian-vector products:

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

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 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 following subsection.

3.1.2 Reverse mode

In this case, as opposed to forward mode AD, the derivatives are propagated back from a given output. A demonstration of how it works is presented in table 3.2 where we compute the sensitivities of output $y_1 = \tilde{f}_1(x_1, x_2)$ of \tilde{f} respect the inputs x_1 and x_2 . Different variables, called adjoints, are associated to each variable v_i of the Forward Primal Trace and each of them is defined as:

$$\bar{v}_i = \frac{\partial y_1}{\partial v_i}$$

An adjoint variable is nothing more than the sensitivity of the output y_1 with respect to changes in v_i . It is important for the reader to note that, with this mode,

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

the variable with respect to the derivatives are computed is not fixed as in forward mode AD.

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-propagating 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 known from previous passages 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 ∇f require 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 [28].

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 Application to design optimization

Design optimization is the process of finding the best design parameters $\mathbf{l} = [l_1, \dots, l_N]$ that satisfy project requirements [29]. Requirements and objectives, which are usually multiple and conflicting, are typically expressed in the form of a scalar function $J: \mathbb{R}^M \times \mathbb{R}^N \rightarrow \mathbb{R}$ named *cost* or *objective function*. J depends on the state $\mathbf{u} \in \mathbb{R}^M$ of the problem at hand and the design parameters $\mathbf{l} \in \mathbb{R}^N$; \mathbf{u} is generally found as solution of a system of linear equations $\mathbf{A}\mathbf{u} = \mathbf{b}$, which arise from the imposition of the constraints of the problem, where $\mathbf{A} \in \mathbb{R}^{M \times M}$ and $\mathbf{b} \in \mathbb{R}^M$ depends on some way on \mathbf{l} . In figure 3.2 can be found a scheme of the whole process.

This type of process is of interest since it is encountered in the most diverse fields, from computer science, where for example the proper weights of a neural

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

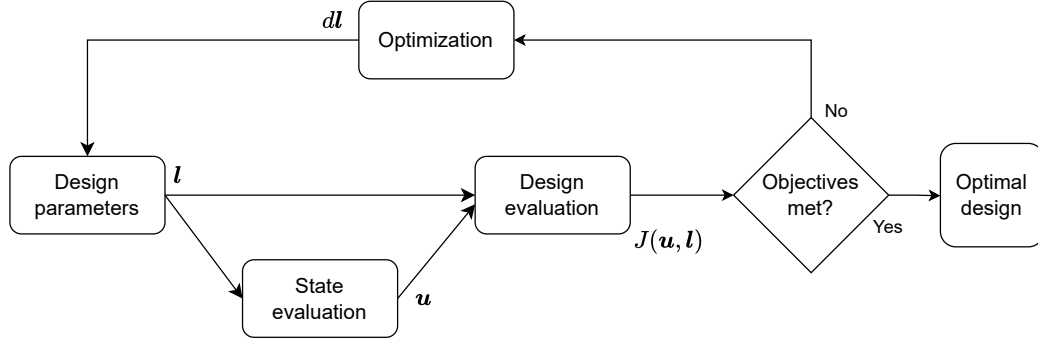


Figure 3.2: Scheme of a design optimization process. From an initial design defined by a set of parameters \mathbf{l} the associated state \mathbf{u} is derived. The state in conjunction with the parameters are then used to evaluate the design through a function J . If the objectives are met the design is kept otherwise the whole procedure is repeated with a new set of parameters obtained by modifying the previous ones by a quantity $d\mathbf{l}$

network has to be found in order to minimize the error between its predicted and the desired output respect a given set of inputs, to mechanical engineering where one may want to modify the shape of an heat sink to maximize the amount of the heat transferred to the environment.

To find the best parameters for a given project a variety of iterative approaches can be pursued during the optimization phase:

Manual approach where each parameter l_1, \dots, l_N is adjusted one at a time. Unfortunately doing so tends to leads to suboptimal results;

Brute-force approach where all possible combinations of design parameters are evaluated. However, this is time-expensive and for large-scale projects, where parameters can be hundreds, thousands, or even more, optimal parameters may not be found in reasonable time-frames;

Gradient approach where the gradient $\frac{dJ}{d\mathbf{l}}$ permits both to update all the parameters l_1, \dots, l_N at the same time and to disregard most of the parameter values that are sub-optimal by considering only those suggested by its direction.

For the effectiveness of the last mentioned approach, is crucial to ensure an efficient computation of the gradient $\frac{dJ}{d\mathbf{l}}$ which must require as few steps as possible and be efficient even for problems involving huge number of design parameters. To meet these requirements Automatic Differentiation (AD) in its reverse-mode declination is necessary.

Given its advantages in design optimization we have learned that using the gradient is a good course of action. To directly evaluate $\frac{dJ}{d\mathbf{l}}$, in case of a particular set of design parameters defined by \mathbf{l} at a given state \mathbf{u} , we would compute:

$$\frac{dJ^T}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{l}} + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.2.1)$$

where notations $\frac{d}{dx}$ and $\frac{\partial}{\partial x}$ indicate the total and the partial derivatives respect the variable x respectively. Of the required elements $\frac{\partial J}{\partial \mathbf{u}}$ and $\frac{\partial J}{\partial \mathbf{l}}$ can be efficiently obtained by means of reverse-mode AD, as explained in subsection 3.1.2 on page 28, since J is known analytically. The matrix $\frac{d\mathbf{u}}{d\mathbf{l}}$, on the other hand, require more careful considerations.

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

$$\frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} + \mathbf{A}\frac{d\mathbf{u}}{d\mathbf{l}} = \frac{d\mathbf{b}}{d\mathbf{l}} \quad (3.2.2)$$

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

$$\frac{d\mathbf{u}}{d\mathbf{l}} = \mathbf{A}^{-1} \left(\frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} \right) \quad (3.2.3)$$

We notice that $\frac{d\mathbf{u}}{d\mathbf{l}}$ is a matrix of size $M \times N$ whose j -th column represent the sensitivity of the state \mathbf{u} respect parameter l_j . In practice it is populated column-wise by solving equation (3.2.3) 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 $\frac{d\mathbf{b}}{d\mathbf{l}}$ and $\frac{d\mathbf{A}}{d\mathbf{l}}$ terms can be simply obtained using AD, if \mathbf{b} and \mathbf{A} are known analytically, and once $\frac{d\mathbf{u}}{d\mathbf{l}}$ is computed it can be substituted in equation (3.2.1) to obtain the gradient which symbolically reads as:

$$\frac{dJ}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}} \left[\mathbf{A}^{-1} \left(\frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}}\mathbf{u} \right) \right] + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.2.4)$$

However the computational costs of this naive procedure scales linearly with the number of the design parameters N since N inversions of matrix \mathbf{A} are required for a single gradient computation. We finally remark that each of these inversions has the same computational cost of solving $\mathbf{A}\mathbf{u} = \mathbf{b}$ for \mathbf{u} and if the solution represent the output of a particularly time consuming computation, as a a CFD simulation is an example of, this is particularly penalizing. A more efficient gradient computation is therefore 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 (3.2.4)

as:

$$\begin{aligned}\frac{dJ}{d\mathbf{l}} &= \left[\frac{\partial J^T}{\partial \mathbf{u}} \mathbf{A}^{-1} \right] \left(\frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right) + \frac{\partial J^T}{\partial \mathbf{l}} \\ &= \boldsymbol{\lambda}^T \left(\frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right) + \frac{\partial J^T}{\partial \mathbf{l}}\end{aligned}\tag{3.2.5}$$

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

$$\mathbf{A}^T \boldsymbol{\lambda} = \frac{\partial J}{\partial \mathbf{u}}\tag{3.2.6}$$

By doing so is possible to solve the system in equation (3.2.6), 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.

At this point it is worth noting that the adjoint method explained above is just a particular application of reverse mode Automatic Differentiation (AD). To make the connections between the two methods clearer we note that the first computation of the objective function's gradient, presented in equation (3.2.1), contains the following vector-Jacobian product:

$$\frac{\partial J^T}{\partial \mathbf{u}} \frac{d\mathbf{u}}{d\mathbf{l}}\tag{3.2.7}$$

where $\frac{\partial J^T}{\partial \mathbf{u}} \in \mathbb{R}^M$ is the vector that multiply $\frac{d\mathbf{u}}{d\mathbf{l}} \in \mathbb{R}^{M \times N}$, the Jacobian of the state \mathbf{u} respect the design parameters \mathbf{l} . From subsection 3.1.2 on page 28 we known that AD allows to compute the resulting vector without explicitly computing the Jacobian $\frac{d\mathbf{u}}{d\mathbf{l}}$. But this is precisely what the adjoint method does through equations (3.2.1 - 3.2.6) by first applying the chain rule to the leftmost term in the gradient formula and then parenthesizing the result in order to make the number of inversions of the matrix \mathbf{A} independent on the number of parameters. The result of the whole process is then summarized by:

$$\boldsymbol{\lambda}^T \left(\frac{d\mathbf{b}}{d\mathbf{l}} - \frac{d\mathbf{A}}{d\mathbf{l}} \mathbf{u} \right)\tag{3.2.8}$$

where $\boldsymbol{\lambda}$ solves equation (3.2.6) and the Jacobian $\frac{d\mathbf{u}}{d\mathbf{l}}$ is no more present. Furthermore the two-steps process of the reverses mode AD for the gradient computation is still there even in the adjoint method:

- the forward pass, used for the computation of variables v_i in table 3.2 on page 31, 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 (3.2.8);
- the reverse pass, done in order to obtain the adjoint variables \bar{v}_i in table 3.2, is now equivalent to 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).

Now that we have acquired a certain grade of familiarity with the adjoint method applied to generic design optimization problems we can see how it is applied to a specif subset of deign optimization problems: those which state \mathbf{u} is obtained solving a Partial Differential Equation using RBF-FD methods.

3.3 Application to RBF-FD

In section 2.6 on page 12 we have seen how the RBF-FD solver for Partial Differential Equations (PDEs) is implemented. Computational Fluid Dynamics (CFD), where PDEs involve fluid flows and heat transfers, is one of the main areas of its application. In all these cases design optimization problems arises almost naturally: regarding heat transfers we have already mentioned the heath sink case where heath exchanged with the environment should be maximized, as to fluid flows a typical example is to find the best shape for the wings of an aircraft in order to maximize its lift and/or minimize its drag.

Once a generic description for the design optimization problem has been written it has to be translated into a mathematical statement in order to solve it through an optimization algorithm. To do so the following items are used for a correct formulation:

- *Design variables* represented by a vector $\mathbf{l} = [l_1, \dots, l_L] \in \mathbb{R}^L$. We consider only continuous variables since they control the shape objects which are allowed to vary continuously in the space;
- *Objective function* indicated with $J: \mathbb{R}^{N_I} \times \mathbb{R}^L \rightarrow \mathbb{R}$. It is a scalar function that yields a quantity used to determines if one geometry is better than another. It is an explicit function of the state \mathbf{u}_I of the problem, which is obtained from the RBF-FD solver, and the design parameters \mathbf{l} . Its choice is crucial: it has to represent the real goal of the whole optimization. A poor choice would lead to a mathematical optimum which does not match with the engineering optimum;
- *Constraints* represented by the RBF-FD governing equations: these are reported in equation (2.6.18) in subsection 2.6.2 on page 14. We reiterate

them here, in a more compact form, since they will be used extensively throughout the rest of this section:

$$\mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B = \mathbf{f} \quad (3.3.1)$$

where we have used the following notation: $\mathbf{u}_I = [u(\mathbf{x}_1), \dots, u(\mathbf{x}_{N_I})]$, $\mathbf{u}_B = [g(\mathbf{x}_{N_I+1}), \dots, g(\mathbf{x}_N)]$ and $\mathbf{f} = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_{N_I})]$, recalling that N is the number of RBF-FD nodes scattered throughout the physical domain $\Omega \cup \partial\Omega$ of the system which is represented by the volume of the object. The first N_I nodes are contained within it, the remaining ones are placed on its surface.

Therefore the resulting optimization problem can be written as:

$$\begin{aligned} & \text{minimize} && J(\mathbf{u}_I, \mathbf{l}) \\ & \text{by varying} && u_i && i = 1, \dots, N_I \\ & && l_j && j = 1, \dots, L \\ & \text{subject to} && \mathbf{C}_I \mathbf{u}_I + \mathbf{C}_B \mathbf{u}_B = \mathbf{f} \end{aligned} \quad (3.3.2)$$

which optimum is found through the gradient approach explained section 3.2 on page 30. In general with a similar formulation can be solved also problems which require the maximization of the objective function since $\max[J(\mathbf{u}, \mathbf{l})] = -\min[-J(\mathbf{u}, \mathbf{l})]$. Finally we remark that in practice the variables that can be modified at each optimization cycle are only the design variables \mathbf{l} since the result of CFD simulation \mathbf{u}_I is fully determined from the solution of constraints (3.3.1);

In the two following subsections we will explain how the gradient $\frac{dJ}{d\mathbf{l}}$ can be computed efficiently by means of the adjoint method for the optimization problem (3.3.2) in case of one and three-dimensional objects.

3.3.1 1D

In this case we consider a general cost function $J: \mathbb{R}^{N_I \times L} \rightarrow \mathbb{R}$. Thanks to this generality its sensitivities respect the design parameters are given by the same exact formula used in general design optimization problems:

$$\frac{dJ^T}{d\mathbf{l}} = \frac{\partial J^T}{\partial \mathbf{u}_I} \frac{d\mathbf{u}_I}{d\mathbf{l}} + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.3.3)$$

Once more, what is now lacking is the term $\frac{d\mathbf{u}_I}{d\mathbf{l}}$. The sensitivities of the state respect the parameters can be computed by differentiating equation (3.3.1) which gives:

$$\frac{d\mathbf{C}_I}{d\mathbf{l}} \mathbf{u}_I + \mathbf{C}_I \frac{d\mathbf{u}_I}{d\mathbf{l}} + \frac{d\mathbf{C}_B}{d\mathbf{l}} \mathbf{u}_B + \mathbf{C}_B \frac{d\mathbf{u}_B}{d\mathbf{l}} = \frac{d\mathbf{f}}{d\mathbf{l}} \quad (3.3.4)$$

Before proceeding a note on the matrix $\frac{d\mathbf{u}_B}{dl}$ could be useful: one might think that it should not be present since values in \mathbf{u}_B are fixed from boundary conditions of the problem 2.6.5 on page 14 (i.e. they are not dependent on the design parameters). The previous statements holds true only for Dirichlet BCs, but not in general since in case of Neumann and Robin BCs the values of \mathbf{u}_B depends also on their normal derivatives respect the surface $\partial\Omega$ whose shape depends on \mathbf{l} . This simply means that only the rows of $\frac{d\mathbf{u}_B}{dl}$ associated to Dirichlet points will make up of zeros elements.

Once $\frac{d\mathbf{u}_I}{dl}$ is isolated from the previous equation it can be plugged in (3.3.4) which in turns become:

$$\begin{aligned} \frac{dJ^T}{dl} &= \frac{\partial J}{\partial \mathbf{u}_I}^T \left[\mathbf{C}_I^{-1} \left(\frac{d\mathbf{f}}{dl} - \frac{d\mathbf{C}_I}{dl} \mathbf{u}_I - \frac{d\mathbf{C}_B}{dl} \mathbf{u}_B - \mathbf{C}_B \frac{d\mathbf{u}_B}{dl} \right) \right] + \frac{\partial J^T}{\partial \mathbf{l}} \\ &= \frac{\partial J}{\partial \mathbf{u}_I}^T \left[\mathbf{C}_I^{-1} \left(\frac{d\mathbf{f}}{dl} - \frac{d\mathbf{C}}{dl} \mathbf{u} - \mathbf{C}_B \frac{d\mathbf{u}_B}{dl} \right) \right] + \frac{\partial J^T}{\partial \mathbf{l}} \end{aligned} \quad (3.3.5)$$

where $\frac{d\mathbf{C}}{dl} \in \mathbb{R}^{N_I \times N}$ is the matrix resulting from the concatenation of the rows of $\frac{d\mathbf{C}_I}{dl} \in \mathbb{R}^{N_I \times N_I}$ and $\frac{d\mathbf{C}_B}{dl} \in \mathbb{R}^{N_B \times N_B}$, 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:

$$\frac{dJ^T}{dl} = \boldsymbol{\lambda}_1^T \left(\frac{d\mathbf{f}}{dl} - \frac{d\mathbf{C}}{dl} \mathbf{u} - \mathbf{C}_B \frac{d\mathbf{u}_B}{dl} \right) + \frac{\partial J^T}{\partial \mathbf{l}} \quad (3.3.6)$$

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

$$\mathbf{C}_I^T \boldsymbol{\lambda}_1 = \frac{\partial J}{\partial \mathbf{u}_I} \quad (3.3.7)$$

The unknown terms involving derivatives computation on the right-hand side can be computed easily through AD except for $\frac{d\mathbf{C}}{dl} \mathbf{u}$. This terms deserves more attentions since 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:

$$\mathbf{Q} = \frac{d\mathbf{C}}{dl} \mathbf{u} \quad (3.3.8)$$

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

$$q_{i,j} = \frac{d\mathbf{C}_{[i,:]} }{dl_j} \mathbf{u} \quad (3.3.9)$$

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 (3.3.9). 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):

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (3.3.10)$$

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

$$\mathbf{M}_{BC}^T \begin{bmatrix} \mathbf{c}_I(\mathbf{x}_i) \\ \mathbf{c}_B(\mathbf{x}_i) \\ \mathbf{c}_p(\mathbf{x}_i) \end{bmatrix} = \begin{bmatrix} \mathcal{L}_1\Phi(\mathbf{x}_i, \mathcal{X}_{i,I}) \\ \mathcal{L}_1\mathcal{B}_2\Phi(\mathbf{x}_i, \mathcal{X}_{i,B}) \\ \mathcal{L}\mathbf{p}(\mathbf{x}_i) \end{bmatrix} \quad (3.3.11)$$

in case of RBF-HFD method. The key aspect that both cases share is that the aforementioned equations arises from the information contained in a single stencil \mathcal{X}_i and, as explained in subsection 2.6.2 on page 14, 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 (3.3.9) as:

$$q_{i,j} = \frac{\mathbf{c}(\mathbf{x}_i)^T}{dl_j} \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.12)$$

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 (3.3.11) respect the j -th design parameter, we obtain:

$$\frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) + \mathbf{M}_{BC}^T \frac{d\mathbf{c}(\mathbf{x}_i)}{dl_j} = \frac{d\mathbf{h}}{dl_j} \quad (3.3.13)$$

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 (3.3.11). Now from the last equation we can isolate:

$$\frac{d\mathbf{c}(\mathbf{x}_i)^T}{dl_j} = \left(\frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \mathbf{M}_{BC}^{-1} \quad (3.3.14)$$

which once plugged in (3.3.12) allows to rewrite $q_{i,j}$ as:

$$q_{i,j} = \left[\left(\frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \mathbf{M}_{BC}^{-1} \right] \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.15)$$

But now we can clearly see that in order to populate the matrix \mathbf{Q} we should invert $N_I L$ 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 usually both the number of parameters and the number of nodes within the physical domain are huge. The problem, here, is very similar to the one related to design optimization in section 3.2 on page 30 where there was a linear relationship between the number of matrix inversions and the number of design parameters: consequently, even here, to improve the situation it is possible to rely on the adjoint method. Modifying the parentheses of equation (3.3.15) it is possible to write:

$$\begin{aligned} q_{i,j} &= \left(\frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \left(\mathbf{M}_{BC}^{-1} \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \right) \\ &= \left(\frac{d\mathbf{h}}{dl_j} - \frac{d\mathbf{M}_{BC}^T}{dl_j} \mathbf{c}(\mathbf{x}_i) \right)^T \boldsymbol{\lambda}_{2,i} \end{aligned} \quad (3.3.16)$$

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

$$\mathbf{M}_{BC} \boldsymbol{\lambda}_{2,i} = \begin{bmatrix} \tilde{\mathbf{u}}_i \\ \mathbf{0} \end{bmatrix} \quad (3.3.17)$$

and has to be computed once for each row of \mathbf{Q} (or equivalently $\frac{d\mathbf{C}}{d\mathbf{t}} \mathbf{u}$). The crucial point is that this second adjoint implementation allows to make the computational cost to obtain $\frac{d\mathbf{C}}{d\mathbf{t}} \mathbf{u}$ independent of the number of design parameter L : the number of systems (3.3.17) 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 been observed we are now able to compute the whole gradient $\frac{dJ^T}{d\mathbf{t}}$ with a computation effort analogous to two RBF-HFD:

1. a first vanilla 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.

3.3.2 3D

Table 3.3: Systems that demands the computation of a solution during the application of RBF-FD and the adjoint method explained in section 3.3.1. 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 $i = 1, \dots, N_I$	$m + M$

Chapter 4

Results

Conclusion

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

4.1 Lorem Ipsum

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tellus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

CONCLUSION

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

4.1.1 Dolor sit amet

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

Sed commodo posuere pede. Mauris ut est. Ut quis purus. Sed ac odio. Sed vehicula hendrerit sem. Duis non odio. Morbi ut dui. Sed accumsan risus eget odio. In hac habitasse platea dictumst. Pellentesque non elit. Fusce sed justo eu urna porta tincidunt. Mauris felis odio, sollicitudin sed, volutpat a, ornare ac, erat. Morbi quis dolor. Donec pellentesque, erat ac sagittis semper, nunc dui lobortis purus, quis congue purus metus ultricies tellus. Proin et quam. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos hymenaeos. Praesent sapien turpis, fermentum vel, eleifend faucibus, vehicula eu, lacus.

Appendix A

Sit Amet

Bibliography

- [1] Michael Hillman Jiun-Shyan Chen and Sheng-Wei Chi. “Meshfree Methods: Progress Made after 20 Years”. In: *Journal of Engineering Mechanics* (2017).
- [2] T. Belytschko et al. “Meshless methods: An overview and recent developments”. In: *Computer Methods in Applied Mechanics and Engineering* (1996).
- [3] W. Benz and E. Asphaug. “Simulations of brittle solids using smooth particle hydrodynamics”. In: *Computer Physics Communications* (1995).
- [4] S. Jun W. K. Liu and Y. F. Zhang. “Reproducing kernel particle methods”. In: *International Journal for Numerical Methods in Fluids* (1996).
- [5] P. Lancaster and K. Salkauskas. “Surfaces Generated by Moving Least Squares Methods”. In: *Mathematics of Computation* (1981).
- [6] 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).
- [7] 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).
- [8] 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.
- [9] H. Brezis. *Functional Analysis, Sobolev Spaces and Partial Differential Equations*. 2011. URL: <https://doi.org/10.1007/978-0-387-70914-7>.
- [10] G. R Liu and Y. T. GU. *An Introduction to Meshfree Methods and Their Programming*. Cambridge University Press, 2005.
- [11] B. Fornberg and N. Flyer. *A Primer on Radial Basis Functions with Applications to the Geosciences*. 2015.

BIBLIOGRAPHY

- [12] B. Fornberg and N. Flyer. “Solving PDEs with radial basis functions”. In: *Acta Numerica* (2015).
- [13] 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).
- [14] I. Amidror. “Scattered data interpolation methods for electronic imaging systems: a survey”. In: *Journal of Electronic Imaging* (2002).
- [15] J. C. Mairhuber. “On Haar’s Theorem Concerning Chebychev Approximation Problems Having Unique Solutions”. In: *Proceedings of the American Mathematical Society* (1956).
- [16] R. Schaback and H. Wendland. “Kernel Techniques: From Machine Learning to Meshless Methods”. In: *Acta Numerica* (2006).
- [17] R. Zamolo. “Radial Basis Function-Finite Difference Meshless Methods for CFD Problems”. PhD thesis. Università degli studi di Trieste, 2018.
- [18] G. Fasshauer. *Meshfree Approximation Methods with Matlab*. World Scientific Publishing Company, 2007.
- [19] 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.
- [20] 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%93unstructured+grid+calculations&author=Tolstykh+A.+I.&publication+year=2000&journal=Proc.+16th+IMACS+World+Congress&volume=228&pages=4606-4624.
- [21] 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).
- [22] 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).
- [23] G. Kosec and J. Slak. “Rbf-fd based dynamic thermal rating of overhead power lines”. In: *Advances in Fluid Mechanics XII* (2018).

BIBLIOGRAPHY

- [24] G. Kosec and J. Slak. “Radial basis function-generated finite differences solution of natural convection problem in 3D”. In: *AIP Conference Proceedings* 2293 (2020).
- [25] 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.
- [26] J. L. Bentley. “Multidimensional Binary Search Trees Used for Associative Searching”. In: *Communications of the ACM* (1975).
- [27] A. Kolar-Požun et al. “Oscillatory behaviour of the RBF-FD approximation accuracy under increasing stencil size”. In: *Computational Science – ICCS 2023* (2023).
- [28] A. G. Baydin, B. A. Pearlmutter, and A. A. Radul. “Automatic Differentiation in Machine Learning: a Survey”. In: *Journal of Machine Learning Research* (2018).
- [29] The MathWorks Inc. *What Is Design Optimization?* 2024. URL: <https://it.mathworks.com/discovery/design-optimization.html>.
- [30] Michel Goossens, Frank Mittelbach, and Alexander Samarin. *The L^AT_EX Companion*. Reading, Massachusetts: Addison-Wesley, 1993.
- [31] R. Zamolo, D. Miotti, and E. Nobile. “Accurate stabilization techniques for RBF-FD meshless discretizations with neumann boundary conditions”. In: *SSRN Electronic Journal* (2022).
- [32] D. Rumelhart, G. Hinton, and R. Williams. “Learning representations by back-propagating errors”. In: *Nature* (1986).
- [33] 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.
- [34] K. C. Jeyanthi. “Efficient node generation over complex 3D geometries for meshless RBF-FD method”. MA thesis. Università degli studi di Trieste, 2022.
- [35] Wikipedia contributors. *Finite difference method*. Online; accessed 21-February-2024. 2024. URL: https://en.wikipedia.org/wiki/Finite_difference_method#Further_reading.

Ei fu. Siccome immobile,
Dato il mortal sospiro,
Stette la spoglia immemore
Orba di tanto spiro

Alessandro Manzoni – *Il Cinque Maggio*